

# IEEE Communications MAGAZINE

www.comsoc.org

- *Software Defined Radio — 20 Years Later*
- *Software Defined 5G Networks for Anything as a Service*
- *Radio Communications*
- *Network Testing*



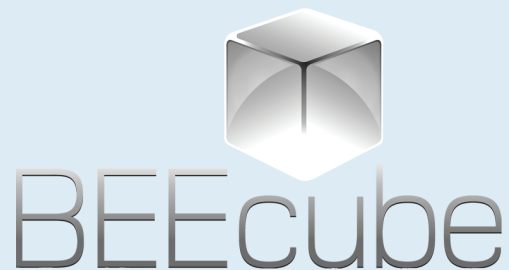
**IEEE**



**IEEE  
COMMUNICATIONS  
SOCIETY**

A Publication of the IEEE Communications Society

THANKS OUR CORPORATE SUPPORTERS



# IEEE Communications MAGAZINE

www.comsoc.org

- *Software Defined Radio — 20 Years Later*
- *Software Defined 5G Networks for Anything as a Service*
- *Radio Communications*
- *Network Testing*



**IEEE**



**IEEE  
COMMUNICATIONS  
SOCIETY**

A Publication of the IEEE Communications Society

# PAM-4 insights don't schedule meetings.

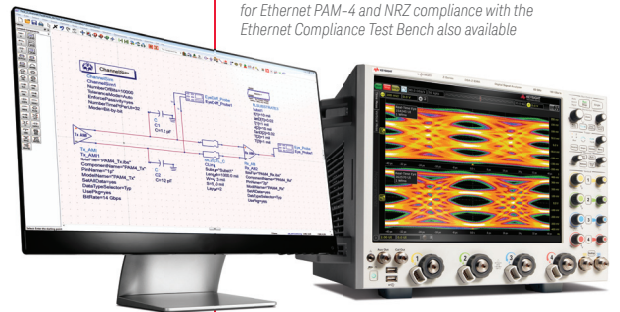
They come when they're good and ready.

Some call them Eureka moments. Others call them epiphanies. We call them insights, the precise moments when you know you've found great answers. As the networking industry considers transitioning to more complex signaling, we can help you achieve insights to meet the technical challenges of PAM-4 that lie ahead. From simulating new designs to characterizing inputs, outputs and connectors, we have the software, hardware and measurement expertise you need to succeed.

## HARDWARE + SOFTWARE + PEOPLE = PAM-4 INSIGHTS



Keysight Advanced Design System bundle for signal integrity  
Simulation-measurement correlation and workflow for Ethernet PAM-4 and NRZ compliance with the Ethernet Compliance Test Bench also available



Keysight Infiniium Z-Series oscilloscopes  
Compliance solutions available for current and emerging PAM-4/Ethernet standards

## PEOPLE

- Member representatives in test working groups including IEEE, OIF-CEI, and Fibre Channel Industry Association
- Applications engineers in more than 100 countries around the world
- Nearly 1,000 patents granted or pending

Download our app note **PAM-4 Design Challenges and the Implications on Test** at [www.keysight.com/find/PAM-4-insight](http://www.keysight.com/find/PAM-4-insight)



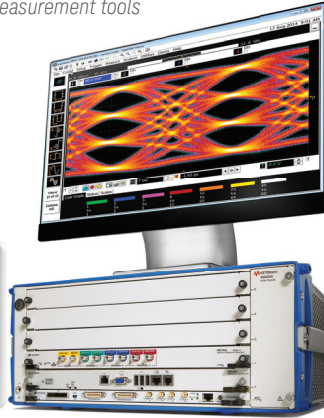
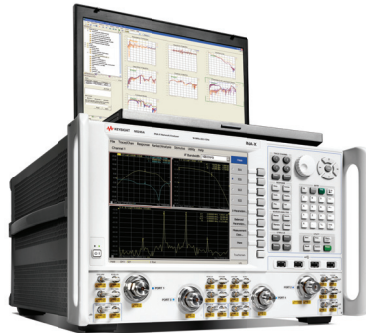
USA: 800 829 4444 CAN: 877 894 4414

© Keysight Technologies, Inc. 2015

## HARDWARE + SOFTWARE

- Instruments designed for testing PAM-4 from simulation to compliance
- Advanced Design System software for simulation-measurement correlation and workflow
- More than 4,000 electronic measurement tools

Keysight 86100D Infiniium DCA-X wide-bandwidth oscilloscope  
Compliance solutions for emerging optical and electrical PAM-4/Ethernet standards



Keysight M8195A 65-GSa/s arbitrary waveform generator  
Flexible PAM-4 pattern generation for 400G Ethernet and beyond



Keysight N5245A PNA-X microwave network analyzer with N1930B physical-layer test system software  
Gigabit Ethernet interconnect and channel test solutions

Keysight J-BERT M8020A high-performance BERT  
The most integrated solution for 100G Ethernet input testing



 **KEYSIGHT**  
TECHNOLOGIES

Unlocking Measurement Insights

#### Director of Magazines

Steve Gorshe, PMC-Sierra, Inc (USA)

#### Editor-in-Chief

Osman S. Gebizlioglu, Huawei Tech. Co., Ltd. (USA)

#### Associate Editor-in-Chief

Zoran Zvonar, MediaTek (USA)

#### Senior Technical Editors

Nim Cheung, ASTRI (China)

Nelson Fonseca, State Univ. of Campinas (Brazil)

Steve Gorshe, PMC-Sierra, Inc (USA)

Sean Moore, Centripetal Networks (USA)

Peter T. S. Yum, The Chinese U. Hong Kong (China)

#### Technical Editors

Sonia Aissa, Univ. of Quebec (Canada)

Mohammed Atiqzaman, Univ. of Oklahoma (USA)

Guillermo Atkin, Illinois Institute of Technology (USA)

Mischa Dohler, King's College London (UK)

Frank Effenberger, Huawei Technologies Co., Ltd. (USA)

Tarek El-Bawab, Jackson State University (USA)

Xiaoming Fu, Univ. of Goettingen (Germany)

Stefano Galli, ASSIA, Inc. (USA)

Admela Jukan, Tech. Univ. Carolo-Wilhelmina zu

Braunschweig (Germany)

Vimal Kumar Khanna, mCalibre Technologies (India)

Myung J. Lee, City Univ. of New York (USA)

Yoichi Maeda, TTC (Japan)

Nader F. Mir, San Jose State Univ. (USA)

Seshradi Mohan, University of Arkansas (USA)

Mohamed Moustafa, Egyptian Russian Univ. (Egypt)

Tom Oh, Rochester Institute of Tech. (USA)

Glenn Parsons, Ericsson Canada (Canada)

Joel Rodrigues, Univ. of Beira Interior (Portugal)

Jungwoo Ryo, The Penn. State Univ.-Altoona (USA)

Antonio Sánchez Esguevillas, Telefonica (Spain)

Mostafa Hashem Sherif, AT&T (USA)

Tom Starr, AT&T (USA)

Ravi Subrahmanyam, InVisage (USA)

Danny Tsang, Hong Kong U. of Sci. & Tech. (China)

Hsiao-Chun Wu, Louisiana State University (USA)

Alexander M. Wyglinski, Worcester Poly. Institute (USA)

Jun Zheng, Nat'l. Mobile Commun. Research Lab (China)

#### Series Editors

##### *Ad Hoc and Sensor Networks*

Edoardo Biagioni, U. of Hawaii, Manoa (USA)

Silvia Giordano, Univ. of App. Sci. (Switzerland)

##### *Automotive Networking and Applications*

Wai Chen, Telcordia Technologies, Inc (USA)

Luca Delgrossi, Mercedes-Benz R&D N.A. (USA)

Timo Kosch, BMW Group (Germany)

Tadao Saito, Toyota Information Technology Center (Japan)

##### *Consumer Communications and Networking*

Ali Begen, Cisco (Canada)

Mario Kolberg, University of Sterling (UK)

Madjid Merabti, Liverpool John Moores U. (UK)

##### *Design & Implementation*

Vijay K. Gurbani, Bell Labs/Alcatel Lucent (USA)

Salvatore Loreto, Ericsson Research (Finland)

Ravi Subrahmanyam, Invisage (USA)

##### *Green Communications and Computing Networks*

Daniel C. Kilper, Univ. of Arizona (USA)

John Thompson, Univ. of Edinburgh (UK)

Jinsong Wu, Alcatel-Lucent (China)

Honggang Zhang, Zhejiang Univ. (China)

##### *Integrated Circuits for Communications*

Charles Chien, CreoNex Systems (USA)

Zhiwei Xu, HRL Laboratories (USA)

##### *Network and Service Management*

George Pavlou, U. College London (UK)

Juergen Schoenwaelder, Jacobs University (Germany)

##### *Networking Testing*

Ying-Dar Lin, National Chiao Tung University (Taiwan)

Erica Johnson, University of New Hampshire (USA)

##### *Optical Communications*

Osman Gebizlioglu, Huawei Technologies (USA)

Vijay Jain, Sterlite Network Limited (India)

##### *Radio Communications*

Thomas Alexander, Ixia Inc. (USA)

Amitabh Mishra, Johns Hopkins Univ. (USA)

#### Columns

##### *Book Reviews*

Piotr Cholda, AGH U. of Sci. & Tech. (Poland)

#### Publications Staff

Joseph Milizzo, Assistant Publisher

Susan Lange, Online Production Manager

Jennifer Porcello, Production Specialist

Catherine Kemelmacher, Associate Editor

# IEEE Communications MAGAZINE

SEPTEMBER 2015, Vol. 53, No. 9

www.comsoc.org/commag

- 6 **PRESIDENT'S PAGE**
- 10 **CONFERENCE PREVIEW/IEEE ICC 2016**
- 12 **BOOK REVIEWS**
- 14 **CONFERENCE CALENDAR**
- 17 **GLOBAL COMMUNICATIONS NEWSLETTER**
- 152 **ADVERTISERS' INDEX**

## SOFTWARE DEFINED RADIO — 20 YEARS LATER: PART 1

GUEST EDITORS: JOSEPH MITOLA, PRESTON MARSHALL, KWANG-CHENG CHEN, MARKUS MUECK, AND ZORAN ZVONAR

- 22 **GUEST EDITORIAL**
- 24 **SOFTWARE RADIO: A CATALYST FOR WIRELESS INNOVATION**  
CHRISTOPHE MOY AND JACQUES PALICOT
- 31 **THE SOFTWARE COMMUNICATIONS ARCHITECTURE: TWO DECADES OF SOFTWARE RADIO TECHNOLOGY INNOVATION**  
CLAUDE BELISLE, VINCE KOVARIK, LEE PUCKER, AND MARK TURNER
- 38 **THE ETSI STANDARD ARCHITECTURE, RELATED INTERFACES, AND RECONFIGURATION PROCESS FOR RECONFIGURABLE MOBILE DEVICES**  
YONG JIN, KYUNGHON KIM, DONGHYUN KUM, SEUNGWON CHOI, AND VLADIMIR IVANOV
- 48 **SECURING PHYSICAL-LAYER COMMUNICATIONS FOR COGNITIVE RADIO NETWORKS**  
YULONG ZOU, JIA ZHU, LIUQING YANG, YING-CHANG LIANG, AND YU-DONG YAO
- 56 **PROTOTYPING REAL-TIME FULL DUPLEX RADIOS**  
MINKEUN CHUNG, MIN SOO SIM, JAEWEON KIM, DONG KU KIM, AND CHAN-BYOUNG CHAE
- 64 **A LOW-COST DESKTOP SOFTWARE DEFINED RADIO DESIGN ENVIRONMENT USING MATLAB, SIMULINK, AND THE RTL-SDR**  
ROBERT W. STEWART, LOUISE CROCKETT, DALE ATKINSON, KENNETH BARLEE, DAVID CRAWFORD, IAIN CHALMERS, MIKE MCLERNON, AND ETHEM SOZER

## SOFTWARE DEFINED 5G NETWORKS FOR ANYTHING AS A SERVICE

GUEST EDITORS: DAVID SOLDANI, BERNARD BARANI, RAHIM TAFAZOLLI, ANTONIO MANZALINI, AND CHIH-LIN I

- 72 **GUEST EDITORIAL**
- 74 **NON-ORTHOGONAL MULTIPLE ACCESS FOR 5G: SOLUTIONS, CHALLENGES, OPPORTUNITIES, AND FUTURE RESEARCH TRENDS**  
LINGLONG DAI, BICHAJ WANG, YIFEI YUAN, SHUANGFENG HAN, CHIH-LIN I, AND ZHAOCHENG WANG
- 82 **RETHINK FRONTHAUL FOR SOFT RAN**  
CHIH-LIN I, YANNAN YUAN, JINRI HUANG, SHUIA MA, CHUNFENG CUI, AND RAN DUAN
- 90 **BASEBAND UNIT CLOUD INTERCONNECTION ENABLED BY FLEXIBLE GRID OPTICAL NETWORKS WITH SOFTWARE DEFINED ELASTICITY**  
JIawei ZHANG, YUEFENG JI, JIE ZHANG, RENTAO GU, YONGLI ZHAO, SIMING LIU, KUN XU, MEI SONG, HAN LI, AND XINBO WANG
- 100 **NETWORK CODED SOFTWARE DEFINED NETWORKING: ENABLING 5G TRANSMISSION AND STORAGE NETWORKS**  
JONAS HANSEN, DANIEL E. LUCANI, JEPPE KRIGSLUND, MURIEL MÉDARD, AND FRANK H. P. FITZEK

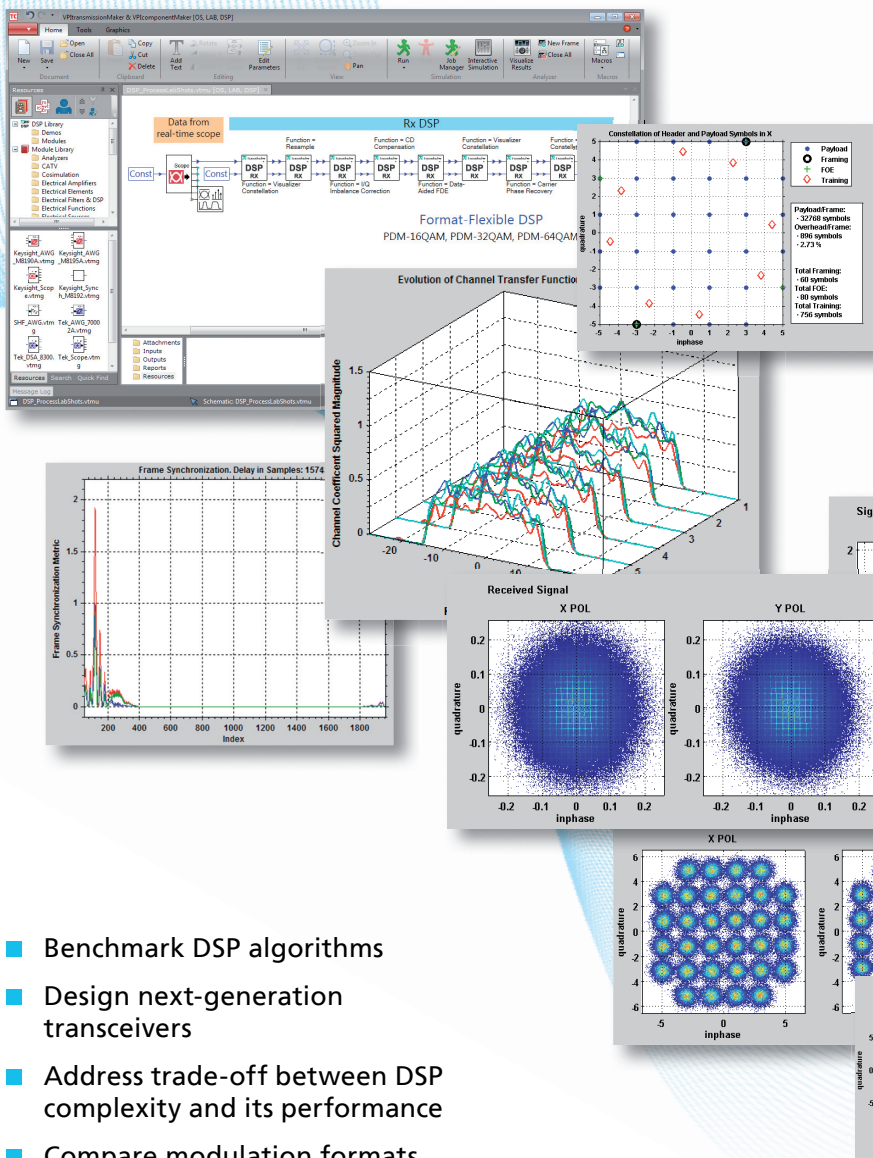




# Fraunhofer

## Heinrich Hertz Institute

# Ready-to-use DSP-Library for optical system simulations and experiments



The DSP-Library for coherent optical systems is available as pluggable toolkit for VPItransmissionMaker™ Optical Systems and VPIlabExpert™. It provides an extensive collection of lab-proven DSP algorithms designed to speed up your development of 100G, 400G and Terabit applications.

- Benchmark DSP algorithms
- Design next-generation transceivers
- Address trade-off between DSP complexity and its performance
- Compare modulation formats
- System performance analysis
- Define component requirements

In Cooperation with:



Further Information:

<http://www.vpiphotonics.com/Tools/OpticalSystems/Toolkits/>

**2015 IEEE Communications Society  
Elected Officers**

Sergio Benedetto, *President*  
Harvey A. Freeman, *President-Elect*  
Khaled Ben Letaief, *VP-Technical Activities*  
Hikmet Sari, *VP-Conferences*  
Stefano Bregni, *VP-Member Relations*  
Sarah Kate Wilson, *VP-Publications*  
Robert S. Fish, *VP-Standards Activities*

**Members-at-Large**

Class of 2015  
Nirwan Ansari, Stefano Bregni  
Hans-Martin Foisel, David G. Michelson

Class of 2016  
Sonia Aissa, Hsiao Hwa Chen  
Nei Kato, Xuemin Shen

Class of 2017  
Gerhard Fettweis, Araceli García Gómez  
Steve Gorshe, James Hong

**2015 IEEE Officers**

Howard E. Michel, *President*  
Barry L. Shoop, *President-Elect*  
Parviz Famouri, *Secretary*  
Jerry L. Hudgins, *Treasurer*  
J. Roberto B. de Marca, *Past-President*  
E. James Prendergast, *Executive Director*  
Harvey A. Freeman, *Director, Division III*

**IEEE COMMUNICATIONS MAGAZINE** (ISSN 0163-6804) is published monthly by The Institute of Electrical and Electronics Engineers, Inc. Headquarters address: IEEE, 3 Park Avenue, 17th Floor, New York, NY 10016-5997, USA; tel: +1 (212) 705-8900; <http://www.comsoc.org/commag>. Responsibility for the contents rests upon authors of signed articles and not the IEEE or its members. Unless otherwise specified, the IEEE neither endorses nor sanctions any positions or actions espoused in *IEEE Communications Magazine*.

**ANNUAL SUBSCRIPTION:** \$27 per year print subscription. \$16 per year digital subscription. Non-member print subscription: \$400. Single copy price is \$25.

**EDITORIAL CORRESPONDENCE:** Address to: Editor-in-Chief, Osman S. Gebizlioglu, Huawei Technologies, 400 Crossing Blvd., 2nd Floor, Bridgewater, NJ 08807, USA; tel: +1 (908) 541-3591, e-mail: [Osman.Gebizlioglu@huawei.com](mailto:Osman.Gebizlioglu@huawei.com).

**COPYRIGHT AND REPRINT PERMISSIONS:** Abstracting is permitted with credit to the source. Libraries are permitted to photocopy beyond the limits of U.S. Copyright law for private use of patrons: those post-1977 articles that carry a code on the bottom of the first page provided the per copy fee indicated in the code is paid through the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923. For other copying, reprint, or republication permission, write to Director, Publishing Services, at IEEE Headquarters. All rights reserved. Copyright © 2015 by The Institute of Electrical and Electronics Engineers, Inc.

**POSTMASTER:** Send address changes to *IEEE Communications Magazine*, IEEE, 445 Hoes Lane, Piscataway, NJ 08855-1331. GST Registration No. 125634188. Printed in USA. Periodicals postage paid at New York, NY and at additional mailing offices. Canadian Post International Publications Mail (Canadian Distribution) Sales Agreement No. 40030962. Return undeliverable Canadian addresses to: Frontier, PO Box 1051, 1031 Helena Street, Fort Eire, ON L2A 6C7.

**SUBSCRIPTIONS:** Orders, address changes—IEEE Service Center, 445 Hoes Lane, Piscataway, NJ 08855-1331, USA; tel: +1 (732) 981-0060; e-mail: [address.change@ieee.org](mailto:address.change@ieee.org).

**ADVERTISING:** Advertising is accepted at the discretion of the publisher. Address correspondence to: Advertising Manager, *IEEE Communications Magazine*, 3 Park Avenue, 17th Floor, New York, NY 10016.

**SUBMISSIONS:** The magazine welcomes tutorial or survey articles that span the breadth of communications. Submissions will normally be approximately 4500 words, with few mathematical formulas, accompanied by up to six figures and/or tables, with up to 10 carefully selected references. Electronic submissions are preferred, and should be submitted through Manuscript Central: <http://mc.manuscriptcentral.com/commag-ieee>. Submission instructions can be found at the following: <http://www.comsoc.org/commag/paper-submission-guidelines>. For further information contact Zoran Zvonar, Associate Editor-in-Chief ([zoran.zvonar@mediatek.com](mailto:zoran.zvonar@mediatek.com)). All submissions will be peer reviewed.



**108 SOFTWARE DEFINED SERVICE MIGRATION THROUGH LEGACY SERVICE INTEGRATION INTO 4G NETWORKS AND FUTURE EVOLUTIONS**

YEUNWOONG KYUNG, TRI M. NGUYEN, KIWON HONG, JONGKWAN PARK, AND JINWOO PARK

**NETWORK TESTING**

SERIES EDITORS: YING-DAR LIN AND ERICA JOHNSON

**116 SERIES EDITORIAL**

**118 TESTING THE CAPACITY OF OFF-THE-SHELF SYSTEMS TO STORE 10GbE TRAFFIC**

VICTOR MORENO, JAVIER RAMOS, JOSÉ LUIS GARCÍA-DORADO, IVAN GONZALEZ, FRANCISCO J. GOMEZ-ARRIBAS, AND JAVIER ARACIL

**126 PLATFORM FOR BENCHMARKING OF RF-BASED INDOOR LOCALIZATION SOLUTIONS**

TOM VAN HAUTE, ELI DE POORTER, FILIP LEMIC, VLADO HANDZISKI, NIKLAS WIRSTRÖM, THIEMO VOIGT, ADAM WOLISZ, AND INGRID MOERMAN

**134 HOW FAR IS FACEBOOK FROM ME? FACEBOOK NETWORK INFRASTRUCTURE ANALYSIS**

REZA FARAHBAKHS, ANGEL CUEVAS, ANTONIO M. ORTIZ, XIAO HAN, AND NOEL CRESPI

**RADIO COMMUNICATIONS: COMPONENTS, SYSTEMS, AND NETWORKS**

SERIES EDITORS: AMITABH MISHRA AND TOM ALEXANDER

**144 SERIES EDITORIAL**

**145 IEEE 802.11AH: SUB-1-GHz LICENSE-EXEMPT OPERATION FOR THE INTERNET OF THINGS**

MINYOUNG PARK

**CURRENTLY SCHEDULED TOPIC**

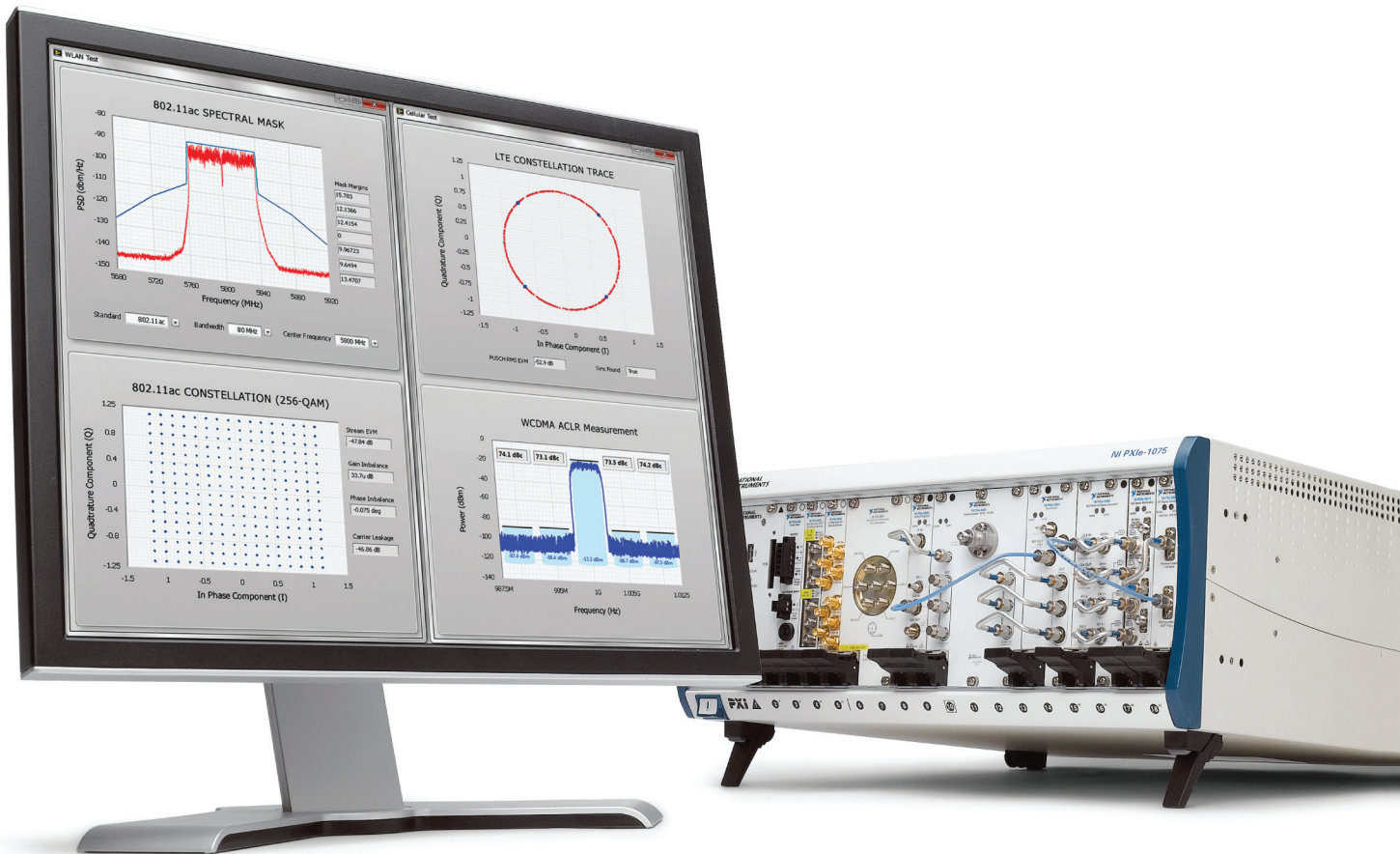
	PUBLICATION DATE	MANUSCRIPT DUE DATE
SEMANTICS FOR ANYTHING-AS-A-SERVICE	MARCH 2016	SEPTEMBER 15, 2015
CRITICAL COMMUNICATIONS AND PUBLIC SAFETY NETWORKS	APRIL 2016	OCTOBER 1, 2015
WIRELESS COMMUNICATIONS, NETWORKING, AND POSITIONING WITH UNMANNED AERIAL VEHICLES	MAY 2016	NOVEMBER 1, 2015
BIO-INSPIRED CYBER SECURITY FOR COMMUNICATIONS AND NETWORKING	JUNE 2016	NOVEMBER 1, 2015
LTE EVOLUTION	JUNE 2016	NOVEMBER 30, 2015
WIRELESS TECHNOLOGIES FOR DEVELOPMENT	JULY 2016	DECEMBER 1, 2015
RECENT ADVANCES IN GREEN INDUSTRIAL NETWORKING	OCTOBER 2016	DECEMBER 15, 2015
COMMUNICATIONS, CACHING, AND COMPUTING FOR CONTENT-CENTRIC MOBILE NETWORKS	AUGUST 2016	JANUARY 1, 2016
SOCIAL AND MOBILE SOLUTIONS IN AD HOC AND SENSOR NETWORKING	JULY 2016	JANUARY 11, 2016

[www.comsoc.org/commag/call-for-papers](http://www.comsoc.org/commag/call-for-papers)



# Redefining RF and Microwave Instrumentation

with open software and modular hardware



Achieve speed, accuracy, and flexibility in your RF and microwave test applications by combining National Instruments open software and modular hardware. Unlike rigid traditional instruments that quickly become obsolete by advancing technology, the system design software of NI LabVIEW coupled with NI PXI hardware puts the latest advances in PC buses, processors, and FPGAs at your fingertips.

## WIRELESS TECHNOLOGIES

National Instruments supports a broad range of wireless standards including:

802.11a/b/g/n/ac	LTE
CDMA2000/EV-DO	GSM/EDGE
WCDMA/HSPA/HSPA+	Bluetooth

>> Learn more at [ni.com/redefine](http://ni.com/redefine)

800 813 5078

© 2012 National Instruments. All rights reserved. LabVIEW, National Instruments, NI, and ni.com are trademarks of National Instruments. Other product and company names listed are trademarks or trade names of their respective companies. 05532



## ACHIEVEMENTS IN TECHNICAL ACTIVITIES AND CONFERENCES

The President Pages from September to December 2015 will be devoted to a description of the activities and related achievements of the leadership of the IEEE Communication Society during my term as ComSoc President (2014-2015). The first page, September 2015, is coauthored by Khaled Ben Letaief, Hikmet Sari, and myself, and summarizes the activities in the area of Technical Activities and Conferences.

Dr. Letaief received the B.S. degree with distinction, and M.S. and Ph.D. degrees in electrical engineering from Purdue University in West Lafayette, Indiana, USA. From 1990 to 1993 he was a faculty member at the University of Melbourne, Australia. He has been with the Hong Kong University of Science & Technology (HKUST) since 1993, where he has held numerous administrative positions, including the Head of the Electronic and Computer Engineering Department, Director of the Center for Wireless IC Design, Director of the Huawei Innovation Laboratory, and Director of the Hong Kong Telecom Institute of Information Technology.

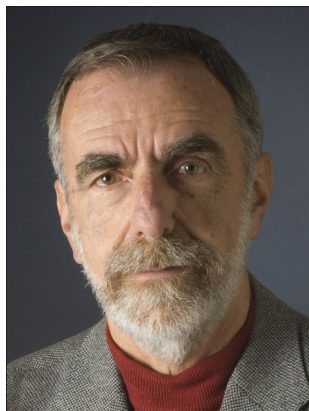
He is currently Chair Professor and Dean of the HKUST School of Engineering. He is also an internationally recognized leader in wireless communications and networks. He has served as a consultant for different organizations, including Huawei, ASTRI, ZTE, Nortel, PricewaterhouseCoopers, and Motorola. He is the founding Editor-in-Chief of the *IEEE Transactions on Wireless Communications* and has served on the editorial board of other prestigious journals. He has also been involved in organizing a number of flagship international conferences and events.

He is the recipient of many other distinguished awards and honors, including the 2007 IEEE Communications Society Joseph LoCicero Publications Exemplary Award, 2009 IEEE Marconi Prize Award in Wireless Communications, 2010 Purdue University Outstanding Electrical and Computer Engineer Award, 2011 IEEE Communications Society Harold Sobol Award, 2011 IEEE Wireless Communications Technical Committee Recognition Award, and 12 IEEE Best Paper Awards.

Dr. Letaief is a long time volunteer with dedicated service to professional societies, and in particular IEEE, where he has served in many leadership positions. These include Treasurer of the IEEE Communications Society, Vice-President for Conferences of the IEEE Communications Society, Chair of the IEEE Committee on Wireless Communications, and elected member of the IEEE Product Services and Publications Board.

Dr. Letaief is a Fellow of IEEE and a Fellow of HKIE. He is currently serving as the IEEE Communications Society Vice-President for Technical Activities, member of the IEEE Fellow Evaluation Committee, and member of the IEEE TAB Periodicals Committee.

Hikmet Sari is currently a Professor and Head of the



SERGIO BENEDETTO



KHALED BEN LETAIEF



HIKMET SARI

Telecommunications Department at SUP-ELEC, near Paris, and also Chief Scientist of Sequans Communications. Prior to moving to these positions, he held various research and management positions at Philips, SAT (SAGEM Group), Alcatel, Pacific Broadband Communications, and Juniper Networks. He has served as an Editor of the *IEEE Transactions on Communications*, Guest Editor of the *European Transactions on Telecommunications*, Guest Editor of *IEEE JSAC*, Associate Editor of *IEEE Communications Letters*, Chair of the Communication Theory Symposium of ICC 2002, Technical Program Chair of ICC 2004, Vice General Chair of ICC 2006, General Chair of PIMRC 2010, General Chair of WCNC 2012, Chair of the GITC Committee in 2010–2011, Distinguished Lecturer of the IEEE Communications Society (2001–2006), Member of the IEEE Fellow Evaluation Committee (2002–2007), and Member of the Awards Committee (2005–2007). His distinctions include the IEEE Fellow Grade and the Andre Blondel Medal in 1995, the Edwin H. Armstrong Award in 2003, the Harold Sobol Award in 2012, and election to the European Academy and to the Science Academy of Turkey in 2012.

## TECHNICAL ACTIVITIES

Over the past two years, a significant amount of activities were initiated and executed within the technical activities portfolio. Such activities are managed and led by the Vice President for Technical Activities (VPTA), Dr. Khaled B. Letaief, who has been assisted by a number of senior volunteers and staff.

In the following, we shall summarize some of the key initiatives, as well as describe some of the major accomplishments that have been achieved over the past two years.

**Positioning IEEE Communications Society to Target New and Emerging Technologies:** This was one of the critical goals set out at the beginning of our term. Under the leadership of Zhisheng Niu, Chair of the Emerging Technologies Committee, volunteers were encouraged to form emerging technology subcommittees in developing technologies and areas, within the society's field of interest. This is critical if we wish to position the IEEE Communications Society to target new and emerging technologies. We are very pleased to witness over the past two years the incubation of numerous committees, including Social Networks; Big Data Processing, Analytics, and Networking; Tactile Internet; Innovation and Standards Information and Communications Technologies; and Quantum Communications and Information Technology.

**Education and Training:** Two of the IEEE Communications Society's major products and services have been our esteemed publications and conferences. In order to evolve and stay ahead, it was decided to build a third key pillar on Education and Training, which will eventually represent the third leg in

the IEEE Communications Society leadership, revenues policy, and future development. In this regard, a new IEEE Communications Society Continuing Professional Education and Training (CPET) program has been thoroughly discussed among the IEEE Communications Society leadership volunteers and within the IEEE Communications Society's Board of Governors. Such a program will continue to be developed over the next few years. It will expand our existing training and education programs and platforms, and also promote the creation of additional new IEEE Communications Society courses and training to build credibility as a preferred provider of knowledge in the areas of information and communications technologies. In particular, such a program will include both online and face-to-face education. To achieve this, significant efforts are under development. These include 1) the creation of an IEEE Communications Society education portal, which will provide an easy, intuitive, personalized, and user-customizable web-interface for facilitating access to information and services related to the IEEE Communications Society's Education and Training offerings; and 2) the development of new online courses that are attractive to our members and meet the needs of the community at large, and especially of professional engineers.

**IEEE Communications Society Summer School:** Young members and especially students are the future of our society and as such, it is critical that we provide special membership development opportunities for them. To help achieve this goal, in July 2015 we launched the first ever IEEE Communications Society Student Summer School in Trento, Italy under the leadership of Fabrizio Granelli. The IEEE Communications Society summer school is designed for graduate students studying communications and related areas. It consists of lectures by international experts and includes poster presentations by participating graduate students. The program covers fundamental, advanced, and hot topics in communications with invited speakers including Andrea Goldsmith, Nelson Fonseca, Giuseppe Bianchi, Lajos Hanzo, and Antonio Capone. The event was very successful and the plan is to continue this activity with the goal of providing high-quality courses while engaging our local chapters and linking our distinguished lecturers to relevant membership development activities.

**Expanding our Awards Program:** To further recognize contributions that advance the fields of interest to the IEEE Communications Society as well as further recognize individual contributions and distinction, two new IEEE Communications Society awards have been created. The first award is the IEEE Communications Society Education Award, which recognizes distinguished and significant contributions to education within the society's technical scope. The second award is the IEEE Communications Society Young Author Best Paper Award, which honors the author(s) of an especially meritorious paper dealing with a subject related to the society's technical scope and who, upon the date of submission of the paper, is younger than 30 years of age.

**Updating and Augmenting the Policies and Procedures:** To guarantee proper operation and management of the IEEE Communications Society technical activities, a thorough review of our policies and procedures have been conducted and appropriate actions were taken to augment and update them.

**Technical Committees Review:** In order to evolve and continue to be at the heart of communications technology development in the world, it is critical that our technical committees stay relevant so as to be able to play a significant and leadership role. To further enhance the role of our technical committees, a review was conducted to evaluate the technical committees' mission and structure and recommend how they can be better positioned to help the IEEE Communications Soci-

ety achieve its mission and grow our values to members from academia and industry. The review was conducted by a task force appointed by VPTA. In particular, a "portfolio analysis" of the current structure and operations, which covers all fields from the technical activities portfolio, has been initiated by the task force under the leadership of Sherman Shen.

**ABET Accreditation of Telecommunications Engineering:** After several years of efforts, primarily led by Tarek El-Bawab with the help of several volunteers in the IEEE Communications Society, "Telecommunications Engineering" has been recognized as a formal engineering course of study. This recognition came from ABET, the accrediting body for academic programs in applied science, computing, engineering, and technology, and brought to a fruitful conclusion the IEEE Communications Society's efforts.

As can be seen, there have been many accomplishments within the Technical Activities portfolio in the past two years. We have only mentioned a few and there are many others that have been omitted due to limited space. Without any doubt, Technical Activities represent the heart and soul of the IEEE Communications Society, and the achievements over the past two years would not have happened without the hard work and dedication of the members of our Technical Communities and education and service-related committees, as well as the technical leadership team, which includes Sherman Shen, Technical Activities Vice-Chair and Secretary; Michele Zorzi, Director of Education and Training; Lajos Hanzo, Chair of the Awards Committee; Zhisheng Niu, Chair of the Emerging Technologies Committee; Kin Leung, Chair of the Fellow Committee; and Steve Weinstein, Chair of the History Committee. These volunteers have worked hard and brought new ideas and enthusiasm that helped us bring to fruition the activities listed above.

### CONFERENCES

With the Vice President-Conferences, Hikmet Sari, the Conferences Leadership Team during 2014-2015 included the Director of Conference Operations, Tarek El-Bawab; Director of Conference Development, Nelson Fonseca; Director of Conference Publications, Chengshan Xiao; GIMS Chair, Paul Hartmann; and GITS Chair, Mike Devetsikiotis. The latter two focused on our society's flagship conferences, ICC and GLOBECOM. At the beginning of our term, this team had a number of ambitious goals and many of them have been achieved. We will briefly review here our accomplishments as well as the initiatives that still have not been completed.

**Conference Operations:** First, the Steering Committee Charters, which describe the operation of the conference steering committees, as well as the composition and the terms of office of their chair and members, are now in place for all of our fully sponsored portfolio conferences.

Another major accomplishment in conference operations concerns technically co-sponsored conferences. Our society receives 80 to 100 technical co-sponsoring (TCS) requests per year, and the process of analyzing, accepting, or declining TCS applications has been significantly improved. We now have a rigorous process and an integrated TCS policy document, which was published in October 2014. It is posted on the community site: [http://www.comsoc.org/files/Conferences/comsoc\\_TCS\\_review\\_process\\_v3.1.pdf](http://www.comsoc.org/files/Conferences/comsoc_TCS_review_process_v3.1.pdf).

This policy document is the first of its kind, not only within ComSoc, but also within IEEE and its other societies. Since the publication of its original version, the TCS policy document has gone through several improvements, including the terms of our partnerships with sister societies in the area of conferences, making our sponsorship of their events conditional on their fulfillment of the same requirements as those

applied to non-ComSoc conferences. We have modified our TCS charging policy to comply with the new IEEE policy that will be in place in 2016. At this moment, we are in Phase I in which charging is uniform across all TCS requests. Phase II, which consists of charging based on IEEE Xplore revenues from TCS conferences, is subject to receiving from the IEEE the full data for all our conferences. This will require additional work to decide the best way to proceed.

A third action in conference operations was the undertaking of our Conferences Financial Accountability (CFA) analysis. The results of this analysis will be useful for cost cutting and more efficient use of resources, and it will also give precise guidelines to the next leadership team if a revision is decided for conference budgeting.

A fourth accomplishment is the revision of The Conflict of Interest (COI) Policy, which was voted by the Conferences Council in June 2014. The new revision allows the conference chairs and leaders with access to the paper review tool to submit a limited number of papers, but those papers are not eligible for paper awards.

**Conference Development:** In order to keep its leadership, ComSoc continuously needs to develop new conferences to address new fields and topics and also to strengthen its presence in different regions. Recently, two regional conferences, LATIN-COM and IEEE BlackSeaCom, became ComSoc portfolio conferences. Following these moves, we have made five major accomplishments in conference development. The first is our involvement in the IEEE-IEEMA INTELEC conference, whose 2015 edition was held in Mumbai, India, on 22–24 January 2015. Together with two other IEEE Societies, ComSoc is a financial co-sponsor of this conference, which primarily addresses industry.

The second accomplishment was the timely launch of the new IEEE Conference on Network Function Virtualization and Software Defined Networks (IEEE NFV-SDN), which will be held on 18–21 November 2015 in San Francisco. This conference is technically co-sponsored by the ETSI, which is a major player in the conference topics. NFV and SDN are two topics of strategic importance, which decouple network functions from the underlying physical infrastructure and open up new avenues to deployment of smart cities, homes, businesses, and other.

The third and fourth accomplishments are in the emerging field of Internet of Things (IoT). One of those is our financial participation in NetSoft 2015, which was held in London on 13–17 April 2015. This conference was initiated by the IEEE Future Directions Committee (FDC) and has other financial co-sponsors, but ComSoc has the biggest share. The other initiative in IoT is our involvement in the IEEE World Forum on Internet of Things (WF-IoT), whose 2015 edition will be held in Milan on 14–16 December 2015.

Our fifth accomplishment in conference development is the signing of an MoU to be a financial co-sponsor of the European Conference on Networks and Communications (EuCNC), which is the showcase of European projects funded by the European Commission. This conference has a very strong industry presence and a large expo. EuCNC 2015 in Paris attracted 550 participants. The 2016 edition will be in Athens, Greece and the 2017 edition in Oulu, Finland.

**Conference Publications:** A major challenge in conference publications is the total time it takes from when papers are submitted to when they appear in IEEE Xplore. This includes a pre-conference time (the time from paper submission to conference dates) and a post-conference time (the time from the conference dates to inclusion of the conference papers in IEEE Xplore). The post-conference time, which was short until a few years ago, suddenly increased to five to six months when the “no-show” checks were introduced. As a conse-

quence, ICC 2012, ICC 2013, and GLOBECOM 2013 suffered from substantial post-conference delays.

A major accomplishment during our term was reducing the post conference delay to four to five weeks. This was achieved for ICC 2014, GLOBECOM 2014, and ICC 2015. A new process was introduced that includes pre-conference and post-conference meetings involving the Director of Conference Publications, the TPC Chair and Publications Chair of the conference, ComSoc Staff, and the paper processing vendor staff.

Another important factor is the pre-conference time, which has been almost constant for decades despite the introduction of an electronic paper submission and review process. For ICC and GLOBECOM, the pre-conference time has been approximately nine month since the 1980s. Starting with GLOBECOM 2015, the paper submission deadline has been shifted in such a way as to reduce the pre-conference time to seven to seven and a half months. ICC 2016, GLOBECOM 2016, and ICC 2017 are following this path.

**Flagship Conferences:** Since GLOBECOM 2002 in Taipei, GLOBECOM was restricted to be held in US cities, and similarly, since ICC 2003 in Anchorage, ICC was always outside the US. This rule had the drawback of excluding cities in the northern part of the US and also many cities in EMEA, Asia-Pacific, and other regions, which are ideal locations for GLOBECOM. The GIMS defined and implemented a new rotation plan ensuring that each region (The Americas, EMEA, and Asia-Pacific) will get an ICC or a GLOBECOM every 18 months. According to this plan, the first GLOBECOM outside the US since 2002 will be held in Singapore in 2017. Next, ICC 2018 will be in Kansas City, GLOBECOM 2018 in Abu Dhabi, and ICC 2019 in Shanghai. Next, GLOBECOM 2019 will be in the Americas.

A major change in the technical content of ICC and GLOBECOM, implemented by the GITC, concerns the status of Green Communications, which was a Track in the Symposium on Selected Areas in Communications (SAC). Starting with ICC 2016, Green Communications will be a stand-alone symposium, and the total number of symposia in the technical program including the SAC Symposium will be 13.

Another accomplishment was the streamlining of the conference review and reporting process prior to the face-to-face GIMS and GITC meetings held during ICCs and GLOBECOMs. This has improved the depth of review performed and allowed GITC and GIMS to spot issues earlier in the conference development process and take corrective actions. This process has allowed future conferences and GITC/GIMS to spend more time discussing relevant issues during the face-to-face meetings.

**Actions Still to be Completed:** On the financial side, the KPI analysis of our portfolio conferences is still to be completed. In order to reduce costs and improve financial accountability, some conferences may be merged or co-located, and some may be discontinued. Further work is needed to come to meaningful conclusions.

On the technical side, the INFOCOM 2015 experience with optimized paper assignment and double-blind reviewing is still to be shared and tested, primarily by the GITC for our flagship conferences, but also by our other major conferences. Finding ways to further improve the efficiency, relevance, and fairness of the paper review process must remain a priority for our successors.

In concluding, we believe that our team has made substantial accomplishments since the beginning of our term in January 2014. We would like to thank all team members for their hard work and dedication. We also extend our thanks to the ComSoc Staff dedicated to conferences under the leadership of Bruce Worthman for their invaluable support.

# With this team on your bench, the sky's the limit.

Visit us at  
EuMW in Paris,  
booth 113

## Signal generation and analysis for demanding requirements

When working at the cutting edge of technology, you shouldn't waste your time with inferior tools. Rely on measuring instruments evolved in the spirit of innovation and based on industry-leading expertise. Instruments like the R&S®SMW200A vector signal generator and the R&S®FSW signal and spectrum analyzer. Each is at the crest of today's possibilities. As a team, they open up new horizons.

See for yourself at [www.rohde-schwarz.com/ad/highend](http://www.rohde-schwarz.com/ad/highend)



**ROHDE & SCHWARZ**

## IEEE ICC 2016 TO EXPLORE “COMMUNICATIONS FOR ALL THINGS” 23 – 27 MAY IN KUALA LUMPUR, MALAYSIA

CALL FOR PAPERS ENDS 16 OCTOBER 2015 FOR PREMIER GLOBAL CONFERENCE  
HIGHLIGHTING ENTIRE COMMUNICATIONS SPECTRUM RANGING FROM IOT TO ENERGY HARVESTING

The IEEE International Conference on Communications (ICC 2016) will hold its next premier international event May 23–27 in Kuala Lumpur, a rising hub of economic and social innovations as well as the cultural, financial and economic centre of Malaysia. Served by modern infrastructure and surrounded by lush parks, mega-sized shopping malls and iconic colonial architecture, IEEE ICC 2016 is expected to host thousands of global communications experts attending more than 1,500 presentations at the globally-benchmarked Kuala Lumpur Convention Centre located close to some large ICT industries and research facilities.

“Ever since Malaysia launched the Multimedia Super Corridor in 1996 as a utopian hub for multimedia advancements, scores of global corporations have come ashore and thousands of start-ups have sprung up to reap the cyber benefits offered through the country’s robust high-speed amenities and intensive R&D activities,” says Borhanuddin M. Ali, IEEE ICC 2016 Executive Co-Chair. “Coupled with the ongoing growth of government communications initiatives and rich infrastructure programs, many new and innovative products and services have emerged over the past few years alone. This has made the city into a rising



cultural and economic power that ranks among the world’s top cities by Foreign Policy’s Global Cities Index.

“The Malaysian capital also combines a buzzing digital economy and ever-rising opportunities for knowledge workers that embody a vibrant mix of Malay, Chinese, Indian and other smaller indigenous population from Borneo Island. With the tagline ‘Malaysia Truly Asia’ this is certainly one of the most exotic locales ever serving as the backdrop for an IEEE ICC event. In fact, every year organizers try to outdo what’s come before. This is no different. We are sure our scientific and learning programs will blend with this sultry, inviting and constantly-evolving technological force in Southeast Asia to offer one of the most rewarding and colorful programs ever produced by our organization.”

Dedicated to the research, advancements and implementations of the next wave of wired and wireline technologies, IEEE ICC 2016 will host hundreds of technical presentations, panels, keynotes, and forums exploring key industry topics ranging from 5G and IoT to mobile cloud computing and green ICT. All interested professionals are invited to visit <http://www.ieee-icc.org/2016> for ongoing conference updates and detailed “Call for Papers” information. The deadline for original technical paper submissions is October 16, 2015.

IEEE ICC 2016 will begin Monday, May 23 with the first of two full days of tutorials and workshops highlighting topics such as cooperative wireless system design, M2M communications, next generation IoT, network coding practices, small cell and 5G networking and forensic computing. It will then proceed over the next three days with more than 2,000 professionals, scientists, academics and government officials attending sessions highlighting the latest research and business policies surrounding communications advancements worldwide.

Those interested in either attending or presenting at IEEE ICC 2016 are invited to visit <http://www.ieee-icc.org/2016> for conference details and specific Call for Papers submission information. The conference’s comprehensive symposia program will be highlighted by the discovery and introduction of next stage technologies in the areas of:

- Ad-Hoc and Sensor Network
- Communication and Information System Security
- Communications QoS, Reliability and Modeling
- Cognitive Radio and Networks
- Communications Software, Services and Multimedia Applications
- Communication Theory
- Green Communications Systems and Networks
- Mobile and Wireless Networks
- Next Generation Networking and Internet
- Optical Networks and Systems
- Signal Processing for Communications
- Wireless Communications

In addition to these fields of research, Selected Areas of Communication will also be devoted to subjects like Access Systems and Networks; Cloud Communications and Networking; Communications for the Smart Grid; Data Storage; E-health; Internet of Things; Satellite and Space Communications; and Social Networking.

For ongoing updates on IEEE ICC 2016, please visit <http://www.ieee-icc.org/2016>. All website visitors are also invited to network with colleagues and peers, share their professional experiences through the conference’s Facebook, LinkedIn and Twitter pages.

## 2016-2017 IEEE-USA Government Fellowships



**Congressional Fellowships**

Seeking U.S. IEEE members interested in spending a year working for a Member of Congress or congressional committee.



**Engineering & Diplomacy Fellowship**

Seeking U.S. IEEE members interested in spending a year serving as a technical adviser at the U.S. State Department.



**USAID Fellowship**

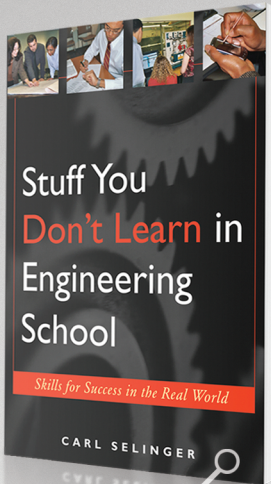
Seeking U.S. IEEE members who are interested in serving as advisors to the U.S. government as a USAID Engineering & International Development Fellow.

The application deadline for 2016-2017 Fellowships is 15 January 2016.

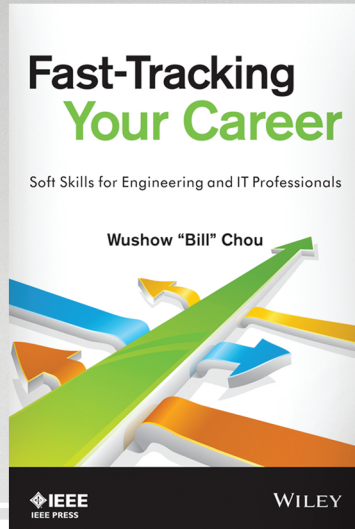
For eligibility requirements and application information, go to [www.ieeeusa.org/policy/govfel](http://www.ieeeusa.org/policy/govfel) or contact Erica Wissolik by emailing [e.wissolik@ieee.org](mailto:e.wissolik@ieee.org) or by calling +1 202 530 8347.



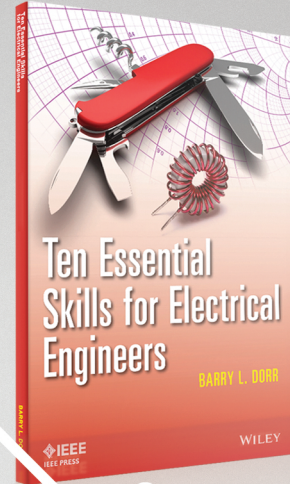

# Advance Your Career With the Experts at Wiley-IEEE Press



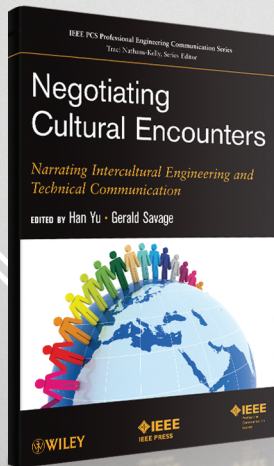
ISBN: 978-0-471-65576-3



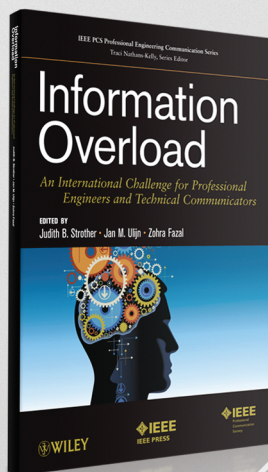
ISBN: 978-1-118-52178-6



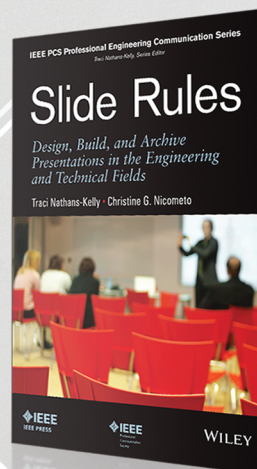
ISBN: 978-1-118-52742-9



ISBN: 978-1-118-06161-9



ISBN: 978-1-118-23013-8



ISBN: 978-1-118-00296-4



ISBN: 978-0-471-77616-1

**ARE YOU AN EXPERT IN YOUR FIELD?  
TO LEARN MORE ABOUT PUBLISHING WITH WILEY-IEEE VISIT:**

[HTTP://WWW.IEEE.ORG/PUBLICATIONS\\_STANDARDS/PUBLICATIONS/PRESS/AUTHORS/BOOKS.HTML](http://www.ieee.org/publications_standards/publications/press/authors/books.html)



VISIT: [WWW.WILEY/IEEE.COM](http://WWW.WILEY/IEEE.COM)

**WILEY**

## TELECOMMUNICATION NETWORK ECONOMICS: FROM THEORY TO APPLICATIONS

BY PATRICK MAILLÉ AND BRUNO TUFFIN

CAMBRIDGE UNIVERSITY PRESS, 2014, ISBN 978-1-107-03275-0, HARDCOVER, 291 PAGES

REVIEWER: JAN DERKACZ

Recently, more and more methods inspired by procedures developed by business departments are being applied in network planning and management. A network engineer or researcher should be acquainted with them, at least at the basic level, and the book by Maillé and Tuffin can help in this learning effort. The book is a comprehensive overview of crucial issues in the economics of telecommunication networks and services that are valid today.

The introductory section presents a general background for later sections. The evolution of telecommunications and the associated economics models are outlined. The main actors of the telecommunication market are described, including: end users, access network service providers, transit providers, content and service providers, as well as regulatory bodies. The next chapter describes the mathematical foundations for models used in the book. The key criteria used for quantification of economic performance are introduced. The general framework of game theory is provided, together with different types of games relevant to the interactions of various stakeholders seeking to optimize their own objectives. The authors also take into account the mechanism design problem, which can be a tool for a system planner dealing with self-interested players. Auctions are analyzed as a specific type of mechanisms used for the allocation of resources to be shared between different actors (operators or users).

The next sections analyze the economics at the level of access service providers as well as content and application providers. This analysis is performed in the historical and evolutionary context. The practice of Internet and telecommunications charging, as well as research activities in this area, are described. The authors draw on their extensive scientific experience and provide their own ideas concerning the pricing models.

Finally, interactions between content/application providers and access service providers in a competitive environment are addressed. The elements of a value chain typical for the current telecommunications market are speci-

fied. Incentives and results of vertical integration of services are discussed. In particular, network neutrality issues arising with the vertical integration are analyzed. The arguments of proponents and opponents of net neutrality are presented.

The book gives a broad perspective on the relations between different actors in telecommunication networks and services accompanied by models, which can be used for choosing the possibly most efficient strategies. As declared by the authors, they aim at “providing tools for a better understanding of the telecommunications ecosystem and better decision-making” for all involved actors, and they succeed in this goal. Exhaustive references for further development of these areas are given in the text.

The book maintains a balance between theory and practice, therefore it can be used either as a textbook for academic courses or a handbook for practitioners. However, the mathematical models presented in the work are not trivial and may require some effort when going into detail, especially for readers who are not familiar with these models.

## MOBILE CLOUDS: EXPLOITING DISTRIBUTED RESOURCES IN WIRELESS, MOBILE AND SOCIAL NETWORKS

BY FRANK H. P. FITZEK AND MARCOS D. KATZ

WILEY, 2014, ISBN 978-0-470-97389-9, HARDCOVER, 220 PAGES

REVIEWER: ROBERT CHODOREK

A mobile cloud is not simply a use of cloud computing via mobile/wireless devices. Mobile cloud technology assumes that mobile/wireless devices (such as smartphones) can actively form mobile clouds. These devices can share their resources, such as connectivity, radio, or energy, in order to achieve a new level of quality of various services (e.g. networking, computing, and social cooperation). The book by Fitzek and Katz presents a very good overview of the promising emerging technology of mobile clouds. The book also contains details about technologies used in mobile clouds and explains how to efficiently use these technologies.

The book is divided into six parts. Part I, comprising Chapters 1 through 3, presents an introduction and gives a solid background on mobile cloud technology. Chapter 1 briefly describes the evolution of mobile and wireless networking and demonstrates the need to build mobile clouds. Chapter 2 is an

introduction to mobile clouds. This chapter provides the definition of mobile clouds as a cooperative arrangement of dynamically connected nodes opportunistically sharing resources. The resources available in the mobile clouds are analyzed and the types of cooperation that have an effect on the way these resources are shared and used are identified. The authors also analyze the benefits of using mobile clouds. Chapter 3 is devoted to the analysis of the new services that can be built in mobile clouds. The new services can utilize shared resources of mobile and wireless devices in the cloud.

Part II is devoted to various technologies necessary to build mobile clouds and provide methods for their effective use. The main communication technologies that are used in a mobile cloud, such as cellular technology and short-range communication (IEEE 802.11, Bluetooth), are described in Chapter 4. The next chapter shows both inter-flow and intra-flow network coding. The chapter presents the most popular methods of coding to improve the efficiency of communications and improve reliability. Part II ends with Chapter 6, which is devoted to explaining cloud formation, operation, and maintenance.

The social aspects of mobile clouds are presented in Part III. To explain the various types of cooperation, an important issue in mobile clouds, Chapter 7 shows examples of cooperation found in nature. Chapter 8 shows diverse forms of cooperation in mobile clouds and the related social factors.

Part IV (Chapter 9) is devoted to green communication. In this chapter the benefits of using mobile clouds are shown. They embrace improved energy efficiency of mobile devices leading to lower power consumption. In Part V (Chapter 10) many examples of applications that use mobile clouds in the mobile cloud testbed are presented. The last part (Chapter 11) gives an overview of the current development related to mobile clouds.

The book is full of insights for researchers, developing engineers, students, and IT professionals. It contains a wide bibliography related to already implemented solutions and solutions being studied in scientific research. The authors present the theory that is skillfully illustrated by many practical examples of mobile clouds, including information about testbeds. The topics presented in the book allow us to understand mobile clouds and learn about the benefits for users, developers, and IT practitioners.



# MILCOM2015

LEVERAGING TECHNOLOGY – THE JOINT IMPERATIVE

OCTOBER 26–28, 2015 • TAMPA, FLORIDA

Register now and join your colleagues in government, military, academia and industry at MILCOM 2015! Experience an in-depth technical program with paper presentations, panel discussions, tutorials, and technology exhibits at the state-of-the-art Tampa Convention Center. *New this year*, one tutorial included with each paid conference registration! For more information – including registration details, technical program outline, and schedule of events – visit [www.milcom.org](http://www.milcom.org).

COHOSTED BY  
AFCEA AND IEEE COMMUNICATIONS SOCIETY



# CONFERENCE CALENDAR

Updated on the Communications Society's Web Site  
[www.comsoc.org/conferences](http://www.comsoc.org/conferences)

**2015**

**OCTOBER**

*LANOMS 2015 — Latin American Network Operations and Management Symposium, 1–3 Oct.*

Joao Pessoa, Brazil

<http://www.lanoms.org/2015/>

**IEEE CLOUDNET 2015 — 4th IEEE Int'l. Conference on Cloud Networking, 5–7 Oct.**

Niagara Falls, Canada

<http://www.ieee-cloudnet.org/>

*RNDM 2015 — 7th Int'l. Workshop on Reliable Networks Design and Modeling, 5–7 Oct.*

Munich, Germany

<http://www.rndm.pl/2015/>

*WMNC 2015 — 8th IFIP Wireless and Mobile Networking Conference, 5–7 Oct.*

Munich, Germany

<http://www.wmnc2015.com/>

**ATC 2015 — Int'l. Conference on Advanced Technologies for Communications, 14–16 Oct.**

Ho Chi Minh, Vietnam

<http://www.rev-conf.org/>

*APCC 2015 — 21st Asia-Pacific Conference on Communications, 14–16 Oct.*

Kyoto, Japan

<http://www.apcc2015.ieice.org/>

**IEEE HEALTHCOM 2015, 17th IEEE Int'l. Conference on e-Health Networking, Application & Services, 14–17 Oct.**

Boston, MA

<http://www.ieee-healthcom.org/index.html>

*WCSP 2015 — Int'l. Conference on Wireless Communications & Signal Processing, 15–17 Oct.*

Nanjing, China

<http://www.ic-wcsp.org/>

**MILCOM 2015 — Military Communications Conference, 26–28 Oct.**

Tampa, FL

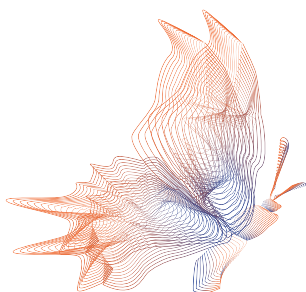
<http://events.jspargo.com/milcom15/public/enter.aspx>

–Communications Society portfolio events appear in bold colored print.

–Communications Society technically co-sponsored conferences appear in black italic print.

–Sister Society conferences appear in plain black print.

–Individuals with information about upcoming conferences, Calls for Papers, meeting announcements, and meeting reports should send this information to: IEEE Communications Society, 3 Park Avenue, 17th Floor, New York, NY 10016; e-mail: [p.oneill@comsoc.org](mailto:p.oneill@comsoc.org); fax: + (212) 705-8996. Items submitted for publication will be included on a space-available basis.



 PACIFIC TELECOMMUNICATIONS COUNCIL

**PTC'16**  
**Reimagining Telecoms**

17-20 January 2016 | Honolulu, Hawaii  
Hilton Hawaiian Village® Waikiki Beach Resort

Register for PTC'16 by 31 October and Save!

Register  
today  
at [ptc.org/ptc16](http://ptc.org/ptc16)

Keynotes by Industry Leaders  
Executive Insight Roundtables  
Industry Briefings  
Workshops  
Topical Sessions  
Social Events for Networking  
Awards Ceremony  
PTC Pavilion with Exhibitors & Meeting Tables  
Hosted Suites by Corporate Partners  
and more...



## CONFERENCE CALENDAR

**IOT 2015** — 5th Int'l. Conference on the Internet of Things, 26–28 Oct.  
Seoul, Korea  
<http://www.iot-conference.org/iot2015/>

**CNSM 2015** — 11th Int'l. Conference on Standards for Communications and Networking, 26–30 Oct.  
Barcelona, Spain  
<http://www.cnsm-conf.org/2015/>

**IEEE ICSOS 2015** — IEEE Int'l. Conference on Space Optical Systems and Applications, 27–28 Oct.  
New Orleans, LA  
<http://icsos2015.nict.go.jp/>

**IEEE CSCN 2015** — IEEE Conference on Standards for Communications and Networking, 28–30 Oct.  
Tokyo, Japan  
<http://www.ieee-cscn.org/>

**GIIS 2015** — Global Information Infrastructure and Networking Symposium, 28–30 October  
Guadalajara, Mexico.  
<http://www.giis-conf.org/>

### NOVEMBER

**IEEE/CIC ICC 2015** — IEEE/CIC Int'l. Conference on Communications in China, 2–4 Nov.  
Shenzhen, China  
<http://www.ieee-icc.org/2015/>

**IEEE COMCAS 2015** — IEEE Int'l. Conference on Microwaves, Communications, Antennas and Electronic Systems, 2–4 Nov.  
Tel Aviv, Israel  
<http://www.comcas.org/>

**IEEE SmartGridComm 2015** — 6th IEEE Int'l. Conference on Smart Grid Communications, 2–5 Nov.  
Miami, FL  
<http://sgc2015.ieee-smartgridcomm.org/>

**IEEE LATINCOM 2015** — IEEE Latin American Conference on Communications, 4–6 Nov.  
Arequipa, Peru  
<http://www.ieee-comsoc-latincom.org/2015/>

**IEEE OnlineGreenComm 2015** — IEEE Online Conference on Green Communications, 10–12 Nov.  
Virtual  
<http://www.ieee-onlinegreencomm.org/2015/>

**AINL-ISMW FRUCT 2015** — Artificial Intelligence and Natural Language & Information Extraction, Social Media and Web Search FRUCT Conference, 9–14 Nov.  
St. Petersburg, Russia  
<http://fruct.org/node/364339>

**IEEE NFV-SDN 2015** — IEEE Conference on Network Function Virtualization and Software Defined Networks, 18–21 Nov.  
San Francisco, CA  
<http://www.ieee-nfv-sdn.org/>

*“If what you want is RF Power, high performance, reliability, and customization, then we are a No Brainer”*



Choosing the right RF power amplifier is critical. But, thanks to AR Modular RF, it's an easy choice.

Our RF power amplifiers give you exactly the power and frequency you need.

With power up to 5kW; and frequency bands from 200 kHz to 6 GHz.

They also deliver the performance and the dependability required for any job. When everything depends on an amplifier that performs without fail, time after time, you can count on AR Modular RF. These amplifiers are compact and rack-mountable; and versatile enough to power all kinds of units, for easy field interchangeability.

For military tactical radios, wireless communication systems, homeland defense systems, high-tech medical equipment, sonar systems, and so much more, your best source for RF power amplifiers is AR Modular RF.

To get the power you need, the features you want, and the performance you demand, visit us at [www.arworld.us](http://www.arworld.us) or call us at 425-485-9000.



**modular rf**

Other **ar** divisions: rf/microwave instrumentation • receiver systems • ar europe

Copyright © 2015 AR. The orange stripe on AR products is Reg. U.S. Pat. & TM. Off.

# CONFERENCE CALENDAR

## DECEMBER

**NETGAMES 2015** — *Int'l. Workshop on Network and Systems Support for Games, 3–4 Dec.*

Zagreb, Croatia  
<http://netgames2015.fer.hr/>

**IEEE GLOBECOM 2015** — **IEEE Global Communications Conference 2015, 6–10 Dec.**

San Diego, CA  
<http://globecom2015.ieee-globecom.org/>

**ITU-K 2015** — *ITU Kaleidoscope: Trust in the Information Society, 9–11 Dec.*

Barcelona, Spain  
<http://www.itu.int/en/ITU-T/academia/kaleidoscope/2015/Pages/default.aspx>

**WF-IOT 2015** — **IEEE World Forum on Internet of Things, 14–16 Dec.**

Milan, Italy  
<http://www.ieee-wf-iot.org/>

**ICSPCS 2015** — *Int'l. Conference Signal Processing and Communication Systems, 14–16 Dec.*

Cairns, Australia.  
[http://www.dspsc-witisp.com/icspcs\\_2015/index.html](http://www.dspsc-witisp.com/icspcs_2015/index.html)

**IEEE ANTS 2015** — **IEEE Int'l. Conference on Advanced Networks and Telecommunications Systems, 15–18 Dec.**

Kolkata, India  
<http://www.ieee-comsoc-ants.org/>

**IEEE VNC 2015** — **IEEE Vehicular Networking Conference, 16–18 Dec.**

Kyoto, Japan  
<http://www.iitmk.ac.in/coconet2015/index.html>

**COCONET 2015** — *Int'l. Conference on Computing and Network Communications, 16–19 Dec.*

Trivandrum, India  
<http://www.iitmk.ac.in/coconet2015/index.html>

## 2016

## JANUARY

**COMSNETS 2016** — *8th Int'l. Conference on Communication Systems & Networks, 5–9 Jan.*

Bangalore, India  
<http://www.comsnets.org/index.html>

**IEEE CCNC 2016** — **IEEE Consumer Communications and Networking Conference, 8–11 Jan.**

Las Vegas, NV  
<http://ccnc2016.ieee-ccnc.org/>

**WONS 2016** — *12th Annual Conference on Wireless On-Demand Network Systems and Services, 20–22 Jan.*

Cortina d'Ampezzo, Italy  
<http://2016.wons-conference.org/>

## MARCH

**ICBDSC 2016** — *3rd MEC Int'l. Conference on Big Data and Smart City, 15–16 Mar.*

Muscat, Oman  
<http://www.mec.edu.om/conf2016/index.html>

**OFC 2016** — **Optical Fiber Conference, 20–24 Mar.**

Anaheim, CA  
<http://www.ofconference.org/en-us/home/>

**IEEE CogSIMA 2016** — **IEEE Int'l. Multi-Disciplinary Conference on Cognitive Methods in Situation Awareness and Decision Support, 21–25 Mar.**

San Diego, CA  
<http://www.cogsima2016.org/>

**WD 2016** — *Wireless Days 2016, 23–25 Mar.*

Toulouse, France  
<http://wd2015.sciencesconf.org/>

**IEEE ISPLC 2016** — **2016 IEEE Int'l. Symposium on Power Line Communications and Its Applications, 29 Mar.–1 Apr.**

Bottrop, Germany.  
<http://www.ieee-isplc.org/>

## APRIL

**IEEE WCNC 2016** — **IEEE Wireless Communications and Networking Conference, 3–6 Apr.**

Doha, Qatar  
<http://wcnc2016.ieee-wcnc.org/>

**IEEE INFOCOM 2016** — **IEEE Int'l. Conference on Computer Communications, 10–15 April**

San Francisco, CA  
<http://infocom2016.ieee-infocom.org/>

**WTS 2016** — *Wireless Telecommunications Symposium, 18–20 Apr.*

London, U.K.  
<http://www.cpp.edu/~wtisi/>

**IEEE/IFIP NOMS 2016** — **IEEE/IFIP Network Operations and Management Symposium, 25–29 Apr.**

Istanbul, Turkey  
<http://noms2016.ieee-noms.org/>

## JUNE

**IEEE BlackSeaCom 2016** — **4th Int'l. Black Sea Conference on Communications and Networking, 6–9 June**

Varna, Bulgaria  
<http://www.ieee-blackseacom.org/>

**EUCNC 2016** — **European Conference on Networks and Communications, 27–30 June**

Athens, Greece  
<http://eucnc.eu/>

## JULY

**TEMU 2016** — *Int'l. Conference on Telecommunications and Multimedia, 25–27 July*

Heraklion, Greece  
<http://www.temu.gr/>

## AUGUST

**EUSIPCO 2016**, 29 Aug.–2 Sept.

Budapest, Hungary  
<http://www.eusipco2016.org/>

### OMBUDSMAN

COMSOC BYLAWS ARTICLE 3.8.10

The Ombudsman shall be the first point of contact for reporting a dispute or complaint related to Society activities and/or volunteers. The Ombudsman will investigate, provide direction to the appropriate IEEE resources if necessary, and/or otherwise help settle these disputes at an appropriate level within the Society...

IEEE Communications Society  
Ombudsman  
c/o Executive Director  
3 Park Avenue  
17 Floor  
New York, NY 10017, USA

[ombudsman@comsoc.org](mailto:ombudsman@comsoc.org)  
[www.comsoc.org](http://www.comsoc.org) "About Us"  
(bottom of page)



September 2015

ISSN 2374-1082

## CONFERENCE REPORT

## The 14th IEEE International Conference on Communication Systems (ICCS), 19-21 November 2014, Macau, China

By Koichi Adachi, ICCS 2014 Publicity Chair

The 14th IEEE International Conference on Communication Systems (ICCS 2014) was successfully held in Macau, China, from 19–21 November 2014. It was co-organized by the IEEE Communications Society Singapore Chapter, the IEEE Vehicular Technology Society Singapore Chapter, and the IEEE Singapore Section, and technically co-sponsored by the IEEE Communications Society and the IEEE Macau Section. Started in Singapore in 1988, the IEEE ICCS has been a biennial conference providing a platform to researchers from both academia and industry to exchange the state-of-the-art communications technologies over two decades.

The organizing committee of IEEE ICCS 2014 includes academics and researchers from both Singapore and Macau. The key members consist of General chair, Prof. L. Wong (NUS, Singapore); General co-chairs, Dr. Y.-C. Liang (I2R, Singapore) and Dr. Benjamin Yue (CEM, Macau); TPC chair, Prof. T. J. Lim (NUS, Singapore); and TPC co-chairs, Prof. J. Q. Li (MUST, Macau) and Dr. S. Sun (I2R, Singapore).

IEEE ICCS 2014 showcased a technical program consisting of 24 technical sessions, covering many exciting aspects of wireless communications, optical communications, devices, and new emerging technologies. Among the 24 technical sessions, 11 were special sessions, organized by distinguished researchers in their respective areas, on current hot topics in communications and networking, featuring well known speakers sharing their latest discoveries. The topics covered by the special sessions and the regular sessions include machine-to-machine/device-to-device (M2M/D2D) communications, cognitive cellular networks, 5G wireless communications, green wireless communications, signal processing for wireless and optical communications, network economics, cyber security, and so on. A total of 151 submissions to the open call were received from more than 20 countries and territories, with 69 papers accepted after a rigorous and challenging technical review process. Each paper was peer reviewed by at least three independent reviewers.



Conference banquet and award ceremony.



Opening addresses from General Co-chairs by Dr. Ying-Chang Liang (left) and Dr. Benjamin Yue (right).



Keynote speeches by Mr. Seizo Onoe (left), Prof. Kwang-Cheng Chen (center), Prof. Shinji Shimojo (right).

IEEE ICCS 2014 featured three keynote talks. The first keynote talk, entitled "Eternal Evolution Toward 5G and Beyond", was delivered by Mr. Seizo Onoe, the Chief Technology Officer (CTO) of NTT DoCoMo, Japan, on the first day of ICCS 2014. The second keynote talk, titled "Networking and Computing for Internet of Things", was brought to the audience by Prof. Kwang-Cheng Chen from National Taiwan University, Taiwan on the second day, 20 November 2014. The third keynote talk, given on the third day by Professor Shinji Shimojo from Osaka University, Japan, focused on future Internet, titled "Future Internet: Managing Innovation and Testbed".

The conference banquet was held on the second day at the restaurant in the conference venue to promote the social interaction among the attendees. The award ceremony was also held during the conference banquet. The Best Paper Award was given to "Matrix Optimization Problems for MIMO Systems with Matrix Monotone Objective Functions", authored by Chengwen Xing (Beijing Institute of Technology and University of Hong Kong, P.R. China), Shaodan Ma (University of Macau, P.R. China), and Yiqing Zhou (Chinese Academy of Science, P.R. China). The Best Student Paper Award was given to "SOAR: Strategy-Proof Auction Mechanisms for Distributed Cloud Bandwidth Reservation", authored by Yang Gui, Zhenzhe Zheng, Fan Wu, Xiaofeng Gao, and Guihai Chen (all from Shanghai Jiaotong University).

This is the second time that ICCS was held outside Singapore, and it attracted more than 120 participants from the Asia Pacific Region as well as from Europe and North America. IEEE ICCS 2014 was a very successful event, thanks to the active participation and involvement from the authors, TPC members, and attendees, and strong support from colleagues from Macau.

The IEEE Communications Society Singapore Chapter and IEEE Vehicular Technology Society Singapore Chapter are now planning to host the next IEEE ICCS in November 2016.

## The First IEEE Communications Society Summer School, Trento, Italy, July 6-9, 2015

By Fabrizio Granelli, DISI – University of Trento, Italy

Membership services are of paramount importance for IEEE ComSoc, and especially involving younger members, they represents a big challenge. With the Education Board, we addressed young engineers with several small outreach activities targeted to high schools and universities, mainly in conjunction with other Chapters' activities or major ComSoc flagship events (IEEE ICC, Globecom, etc.).

The idea to organize a ComSoc Summer School for our Ph.D. student members was originally conceived by IEEE ComSoc President Prof. Sergio Benedetto as well as Prof. Khaled Letaief, Vice-President of Technical Activities, Prof. Stefano Bregni, Vice-President of Member Relations, and Prof. Michele Zorzi, Director of Education.

I'm honored that my university in Trento was chosen for the first edition of this event, as I believed that the city of Trento, the educational, scientific, financial, and political center of Northeastern Italy, could effectively offer the ideal launching point for the first IEEE Communications Society (ComSoc) Summer School. Indeed, Trento is recognized throughout the country for its high quality of life, prosperous business opportunities, and advanced research centers, and it hosts one of the youngest and most modern universities in Italy. The University of Trento was founded in the 1960s, and its Department of Information Engineering and Computer Science is the second highest ranked ICT department in Italy.

Held July 6-9, the first IEEE ComSoc Summer School provided participants with top-level lectures on hot topics in communications as well as an exciting and unprecedented program, consisting of several site visits and an interactive poster session.

We did our best to select the most promising students from more than 100 applicants during the Spring of 2015, aiming at a good balance between young talents at the beginning of their Ph.D. studies as well as more mature researchers. The result was a heterogeneous class of 43 highly motivated participants from all over the world.

The event started on July 6, with the introductory talks given by the Rector of the University of Trento, Prof. Paolo Collini; Head of the Department of Information Engineering and Com-



Opening session: introduction by the Rector of the University of Trento, Paolo Collini. From left to right: Stefano Bregni (IEEE ComSoc VP Member Relations), Gian Pietro Picco (Head of the Dept. of Information Engineering and Computer Science, DIEC, Univ. of Trento), Paolo Collini, Fabrizio Granelli (Head of the IEEE ComSoc Education Board Training Working Group, DIEC, Univ. of Trento).



Poster session.

puter Science (host and co-founder of the initiative), Prof. Gian Pietro Picco; the ComSoc VP of Member Relations, Prof. Stefano Bregni; and myself as the Head of the ComSoc Education Board Training Working Group and local organizer of the event.

The technical program started after a short break with an interactive lecture on "Collaborative Near-Capacity Wireless System Design" by Prof. Lajos Hanzo, University of Southampton, UK.

*(Continued on Newsletter page 4)*



Group photo of students and professors participating to the Summer School.

## Orange: A Romanian Success Story

### An Interview by Nicolae Oaca

Jean Francois Fallacher has been Orange Romania's CEO since July 1, 2011. He is also the president of the BoD. Orange Romania has been a market leader for more than ten years, while in the last four years, under his leadership, its leading market position was consolidated.

Q: Please explain what is behind your success.

JFF: In 2004 Orange became the leader in the Romanian mobile market by offering customers last-generation technologies and reliable services, based on an extensive and performant network.

Many things have happened since I came to Romania: HD Voice, 4G and 4G+, and a multi-screen TV service. Despite the very competitive Romanian market, we remained the customers' first choice by listening to them and investing in both infrastructure and innovative products and services.

In order to stay on top, we pay a lot of attention to the professionals we recruit for our team, people keen on listening to customers' voices and looking to the latest trends. We have been making massive investments in our 4G, 3G, and 2G infrastructures, and we are focusing on value in the market in this digital world.

Q: How about the use of mobile technologies in Romania?

JFF: In recent years, we've invested in a very performant 4G network, already covering 95 percent of the urban population and 67 percent of the entire population. We've also made investments to improve our national 3G+ network. We have by far the largest 4G coverage in the country, mainly in cities, where we saw a greater demand for high mobile Internet speeds. At the same time, we are currently running an extensive network refresh plan to improve our 2G and 3G coverage in cities and rural areas across the country, as well as the mobile data experience for our customers.

Our ambition for the next five years, set through our strategy Essentials2020, is to provide everyone with an unmatched customer experience and allow customers to be connected to what is essential to them. One out of two subscribers buying a subsidized phone from our retail network chooses a smartphone. They need good coverage to keep in touch with their families and friends through voice and SMS, and our investments are going toward them also.

Q: What is Orange's 4G approach for Romania?

JFF: In 2012 and 2013 we massively invested in 4G network preparation and modernization in the cities. In 2013 alone, the level of investments, mainly directed to 4G, was the highest since the beginning of the economic crisis, 30 percent higher than in 2009. Last year we started a new chapter, investing particularly to extend 4G and to offer performant and reliable services regardless of the network, and we also launched the first 4G+ network in Romania, in six cities.

Customers can also enjoy a very rich portfolio of 4G devices for all needs and budgets, 4G services in all our subscriptions, and also a wide range of services like Orange TV Go, Orange Cloud, or Deezer. Orange PrePay customers were the first prepaid users on the market to have access to 4G and the only ones for a long period of time.

The results so far show us that we've made the correct decisions. In the second quarter of 2015, our Internet traffic increased 112 percent compared with 2014, while the number of 4G users reached 600,000. According to the latest ANCOM official statistics, Orange ranked first in the speed tests run in 2015 in the Netograf

app, offering its customers the best download speeds on mobile for indoor locations at 34.15 Mbps, and for outdoor locations at 52.81 Mbps.

Q: How about the Romanian market?

JFF: Compared to other countries where I have been, the quality of the telecom infrastructure in Romania is extraordinary. This is an interesting place for telecom, very competitive and very dynamic. Romania offers high quality Internet services, ranking fifth worldwide on the current Net Index for data transfer speeds, while having among the lowest prices in Europe because of the low wages. Last years' European regulations, as well as local legislation changes, have led to some challenges that affected the entire market.

In the four years I have been at Orange, I have watched the country change in a good way. I came here just after the economic crisis, in the middle of a price war, which seemed normal due to the context. Romanians are open to technology and to what is new. Price remains very important, but the overall experience with an operator has also become a significant factor in the purchase decision, encouraged by the country's recovering economy.

Q: ... and 2015?

JFF: In the entire telecom market, the first quarter's performance was still affected by the interconnection tariff regulations. Starting with the second quarter these effects are becoming weaker, while some special taxes were reduced at the beginning of the year. Therefore, we expect a positive evolution of the market.

2015 started off on the right foot for us. For the first time since the start of the economic crisis in 2009, our half-year revenues increased this year by 1.2 percent to €453.9mn, Romania being a front-runner in our group.

In terms of investments, this year we plan to invest more than €110m in Romania, with plans to invest close to a half billion euros by 2020. Besides a network refresh and 4G consolidation, our investment plan also includes services we have recently launched. For example, Orange TV, which we launched as a multi-screen TV service and the platform with the biggest number of HD channels in Romania, has already gained more than 221,000 subscribers in approximately two years, while the app for smartphones and tablets has been downloaded or upgraded more than 1.3 million times. This year Romania was the first country in our group to launch the Orange TV stick, a portable device transforming any HDMI TV into a smart TV. With this device, customers are able to watch on their TV screen movies and shows from the mobile app, even when they are not at home.

Q: What are the rationales behind the Orange Educational Program (OEP)?

JFF: OEP is one of our oldest social responsibility programs and a good example of how CSR strategy is included in Orange's core business. Besides scholarships for the best students in telecommunications, the program offers participants the chance to complement their theoretical knowledge with practical information offered by Orange employees, as well as get to know the company and make themselves known in the company through internships. Twenty percent of the engineers, experts, and managers working in our technology department, including Stefan Slavnicu, our Chief Technology Officer, are OEP graduates. Moreover, it is the best example of a sustainable program. We are investing in preparing future professionals who will create the technologies of the future. They will be the people who will have this market "in their hands".



Jean Francois Fallacher

## SUMMER SCHOOL/Continued from page 2

talk started by addressing the limitations of MIMO's reliance on co-located array-elements, and then showed how single-antenna-aided cooperative mobiles can circumvent these limitations through the formation of MIMO's with distributed elements, a concept also referred to as Virtual Antenna Arrays (VAA). Then the speech focused on amplify-forward and decode-forward protocols as well as iterative channel coding schemes. Finally, EXIT-chart-aided designs were introduced for creating near-capacity solutions in addition to a range of future research directions and open problems. The session was closed by a comprehensive discussion with the audience about the relationship between the presented subjects and each participant's Ph.D. topic.

An informal dinner at a local pub completed the day, giving the attendees an opportunity to get to know each other better and have discussions in a relaxed environment.

The following day the school hosted Prof. Giuseppe Bianchi from the University of Rome – Tor Vergata, Italy, who delivered a talk entitled "From Dumb to Smarter Switches in Software Defined Networks: Towards a Stateful Data Plane." The seminar was motivated by a crucial shortcoming in today's software defined network architectures, namely the need to mandate the execution of all control tasks to a remote controller, and the relevant emerging concerns in terms of latency and signaling overhead. After a brief overview of software defined networking principles as well as a review of OpenFlow, the talk focused on recent approaches (recent OpenFlow extensions, Reconfigurable Match Tables, Protocol Oblivious Forwarding, etc.) devised to improve data plane programmability on the fast path, i.e. directly inside the switches themselves. In addition, it introduced Open-State, a very recent proposal devised to permit platform-agnostic programmability of stateful tasks.

Then, on Tuesday afternoon a two-hour poster session was held in the open spaces of the Department of Information Engineering and Computer Science. Each participant was instructed to prepare a poster presentation of his/her ongoing and future research plans in order to foster discussion among participants as well as with school instructors and local professors.

The topics were extremely diverse, from signal processing and resource allocation in LTE to green communications, from cooperation to streaming services.

Participants had one free night for sightseeing or additional networking, before the following seminar, which was held on Wednesday morning. Prof. Andrea Goldsmith from Stanford University, USA, presented a four-hour lecture on "The Next Wave in Wireless Communications." She provided the audience with an extensive discussion of multiuser systems, a survey of current cellular systems and a review of future expectations. The speech then focused on ad hoc wireless networks and on secondary user access in the framework of cognitive radio networks.

Finally, Prof. Goldsmith introduced the issues related to object connectivity (sensors, battery-limited devices) to the Internet, opening a discussion of green wireless networks and the ways to save or harvest power for communications. The session ended with additional hot topic examples and transversal areas of application of communications theory (e.g. neural science).

The Summer School included specific sessions enabling participants to understand the actual problems and technology in the field of communications. To this aim, practical sessions were held on Tuesday and Wednesday and included visits to the datacenter of the University of Trento and to the local network provider, Trentino Network and its Network Operation Center.

Prof. Nelson L.S. da Fonseca of the State University of Campinas, SP, Brazil, closed the event with a seminar on "Networking for Big Data". The seminar started with a discussion of the Big Data ecosystem, perspectives on society and science and processing capabilities. Prof. Fonseca introduced network virtualization as a fundamental pillar to networking support for the Big Data area, as well as research opportunities not only in networking for Big Data but also in Big Data for networking. Participants were provided with lab exercises to practice on their way back home.

Summarizing, the Summer School was a challenging and exciting event. Considering the feedback I collected from the attendees and the speakers, I'm very proud of the outcome, and I believe it could pave the way to becoming an integral and permanent component of the yearly events organized by the IEEE Communications Society for its members.

For additional information on the IEEE ComSoc Summer School, please visit <http://www.comsoc.org/summer-school>.

### OMBUDSMAN

#### COMSOC BYLAWS ARTICLE 3.8.10

The Ombudsman shall be the first point of contact for reporting a dispute or complaint related to Society activities and/or volunteers. The Ombudsman will investigate, provide direction to the appropriate IEEE resources if necessary, and/or otherwise help settle these disputes at an appropriate level within the Society....."

#### IEEE Communications Society Ombudsman

c/o Executive Director

3 Park Avenue

17 Floor

New York, NY 10017, USA

[ombudsman@comsoc.org](mailto:ombudsman@comsoc.org)

[www.comsoc.org](http://www.comsoc.org) "About Us" (bottom of page)

**GLOBAL COMMUNICATIONS NEWSLETTER**

**STEFANO BREGNI**  
Editor  
Politecnico di Milano – Dept. of Electronics and Information  
Piazza Leonardo da Vinci 32, 20133 MILANO MI, Italy  
Tel: +39-02-2399.3503 – Fax: +39-02-2399.3413  
Email: [bregni@elet.polimi.it](mailto:bregni@elet.polimi.it), [s.bregni@ieee.org](mailto:s.bregni@ieee.org)

**IEEE COMMUNICATIONS SOCIETY**  
STEFANO BREGNI, VICE-PRESIDENT MEMBER RELATIONS  
PEDRO AGUILERA, DIRECTOR OF LA REGION  
MERRILY HARTMANN, DIRECTOR OF NA REGION  
HANNA BOGUCKA, DIRECTOR OF EAME REGION  
WANJUN LIAO, DIRECTOR OF AP REGION  
CURTIS SILLER, DIRECTOR OF SISTER AND RELATED SOCIETIES

**REGIONAL CORRESPONDENTS WHO CONTRIBUTED TO THIS ISSUE**  
NICOLAE OACA ([nicolae\\_oaca@yahoo.com](mailto:nicolae_oaca@yahoo.com))  
EWELL TAN, SINGAPORE ([ewell.tan@ieee.org](mailto:ewell.tan@ieee.org))



A publication of the  
IEEE Communications Society

[www.comsoc.org/gcn](http://www.comsoc.org/gcn)  
ISSN 2374-1082



# SEPTEMBER 2015



## Low Power, Short Range RF Mesh Network Communications & Emerging Residential Applications

This tutorial focuses on the development of emerging, interoperable narrowband wireless communications designed to enhance the home experience and everything from residential energy management to home security.

IEEE ComSoc content sponsored by:



## Error Vector Magnitude measurements fit for 5G

Error vector magnitude, EVM, measurements have been the mainstay of modulation performance analysis for more than twenty years. Each new technology has defined a specific measurement to suit the characteristics of the physical layer signal. The interest in signals for 5G that are much wider bandwidth, operating at much higher frequencies means it's time to draw a comparison between the different waveforms and the impact on the measurement of EVM.

This presentation reviews what an EVM measurement is and what it can tell us about the device being measured. A combination of real life and simulated examples are used, with single and multi-carrier waveforms having bandwidths of 20 MHz – 2 GHz, to demonstrate the impact of a variety of signal impairments, including broadband noise and phase noise. The examples will show how to make measurements that give the expected, and consistent results.

Sponsor content provided by:



Limited Time Only at >> [ww.comsoc.org/freetutorials](http://ww.comsoc.org/freetutorials)



For this and other sponsor opportunities, please contact Mark David // 732-465-6473 // [m.david@ieee.org](mailto:m.david@ieee.org)

## SOFTWARE DEFINED RADIO — 20 YEARS LATER: PART 1



Joseph Mitola



Preston Marshall



Kwang-Cheng Chen



Markus Mueck



Zoran Zvonar

It has been two decades since the publication of the seminal tutorial “The Software Radio Architecture” in *IEEE Communications Magazine*. The evolution of Software Defined Radio (SDR) systems had profound impact both in military and multi-standard commercial space. The intent of the feature topic is to capture key elements of the evolution of the enabling technologies as well as SDR solutions for flexible and reconfigurable wireless platforms.

During the first twenty years of software radio, the concept has evolved from a non-linear and somewhat academic projection of how radio engineering could benefit over time from the migration of analog to digital conversion (ADC) and digital to analog conversion (DAC) from baseband, where it was then, to radio frequency (RF), which now is termed the ideal software radio (or Mitola radio) which radio systems engineering is approaching, at least in some high end market niches. Most of us have adopted Stephen Blust’s concept of SDR as the crucial near term value proposition on the way to the ideal software radio which remains a useful abstraction. Good portion of radio systems today employ Intermediate Frequency (IF) ADCs and DACs with digital IF conversion to baseband, multi-carrier, multi-frequency and heterogeneous infrastructure even in fourth generation commercial wireless leading the way to incredibly multi-band shared spectrum and extremely heterogeneous infrastructure of the next generation of commercial wireless.

Cognition emerged a decade later and still means different things to different people. Unfortunately software-defined radio also brings with it the malware menaces of low cost commercial off the shelf (COTS) hardware and software from global supply chains. Open source appears to be an important and growing software radio niche with GNU Radio and the ever present USRP leading the way, especially for academic research and start-up companies where we find much global innovation. Twenty years ago, Europe was clearly in the lead globally with the success of GSM. Over time an exploding base of products, systems, and radio research in China, Singapore, Taiwan, and South Korea created a shift of the center of gravity of software

radio from the US and EU twenty years ago to more global leadership over the next twenty years. Thus, the globalization of software radio technology while challenging provides unprecedented opportunities for all of us. Perhaps the next twenty years of software radio will see the emergence of truly hackproof radio systems and networks so that our children and grandchildren can experience the benefits of software radio technology with fewer of the disadvantages of the today’s Wild Wild West of wired and wireless Cyberspace. Looking ahead, many exciting opportunities and breakthroughs certainly are in store.

Articles selected for this series focus on the key issues and emerging concepts of contemporaneous software radio technology. This Feature Topic is composed of two parts: the second part is to appear in January 2016 issue of this magazine.

Part 1 provides overview of trends in SDR over last two decades. It opens with the overview paper “Software Radio, a Catalyst for Wireless Innovation” by Moy and Palicot, followed by “The Software Communications Architecture (SCA): Two Decades of Software Radio Technology Innovation” by Lee *et al.* that provides an overview of Software Communication Architecture that has been deployed in more than 400,000 radios worldwide. “ETSI-Standard Architecture, Related Interfaces, and Reconfiguration Process for Reconfigurable Mobile Devices” by Jin *et al.* introduces the ETSI-standard architecture for a reconfigurable Mobile Devices which will bring SDR technology to mainstream cellular networks. Continuing, “Securing Physical-Layer Communications for Cognitive Radio Networks” by Zou *et al.* addresses increasingly important topic of physical layer security. “Prototyping Real-Time Full Duplex Radios” by Chung *et al.* describes the SDR approach to prototyping 5G system and quantifying capacity gains. In a similar spirit “A Low Cost Desktop Software Defined Radio Design Environment using MATLAB, Simulink and the RTL-SDR” by Stewart *et al.* describes the flexible approach for prototyping that can aid lab and digital communication curriculum.

## ACKNOWLEDGEMENTS

The Guest Editors would like to thank the large number of colleagues who significantly contributed to this feature Topic, including the authors, reviewers and *IEEE Communications Magazine* editorial team.

## BIOGRAPHIES

JOSEPH MITOLA III (joseph.mitola@federatedwireless.com) is President and CEO, Hackproof Technologies Inc. Previously he held positions of technology leadership with Stevens Institute of Technology, U.S. DoD, The MITRE Corporation; ITT, Harris, E-Systems and the Executive Office of the President of the United States. B.S. E.E. from Northeastern University; M.S.E. from The Johns Hopkins University; Licentiate and Doctorate in Teleinformatics from KTH, The Royal Institute of Technology, Stockholm, Sweden.

PRESTON MARSHALL (pres@google.com) has over 30 year background in communications systems. Currently he is with Google Inc. Previously he was Deputy Director at USC Information Science Institute, before that with the Defense Advanced Research Projects Agency (DARPA) Strategic Technology Office, and served as a Program Manager for the DARPA Wireless and Cognitive Radio and Network programs. He holds a B.S.E.E and M.S. Information Science from Lehigh University, and Ph.D. from Trinity College, Dublin.

KWANG-CHENG CHEN (ckc@ntu.edu.tw) is the Distinguished Professor and Associate Dean in the College of Electrical Engineering and Computer Science, National Taiwan University. He has contributed essential technology to IEEE 802, Bluetooth (adaptive frequency hopping as the first widely used cognitive radio technology), and LTE and LTE-A wireless standards. He received 2011 IEEE COMSOC WTC Recognition Award, 2014 IEEE Jack Neubauer Memorial Award and 2014 IEEE COMSOC AP Outstanding Paper Award.

MARKUS MUECK (markus.dominik.mueck@intel.com) received the Dipl.-Ing. and ing. dipl. degrees from the University of Stuttgart, Germany and the Ecole Nationale Supérieure des Télécommunications (ENST), Paris, France respectively in 1999. In 2006, he received the Doctorate degree of ENST in Communications. He is with Intel Deutschland GmbH, Munich, Germany and he acts as ETSI Board Member supported by INTEL, as general Chairman of ETSI RRS Technical Body (Software Radio and Cognitive Radio Standardization).

ZORAN ZVONAR (zoran.zvonar@mediatek.com) is the Senior Director and MediaTek Fellow, leading advanced technology development. Previously he was with Analog Devices (ADI) and has been recipient of the company's highest technical honor of ADI Fellow. He has been elected IEEE Fellow in 2010. He received the Dipl. Ing. and the M.S. degree from the School of Electrical Engineering, University of Belgrade, and the Ph.D. degree in Electrical Engineering from the Northeastern University, Boston.

## CALL FOR PAPERS IEEE COMMUNICATIONS MAGAZINE SEMANTICS FOR ANYTHING-AS-A-SERVICE

### BACKGROUND

Services (including anything-as-a-service) are the buzz in the industry. Networks are morphing to utilize new technologies like Network Functions Virtualization and Software-Defined Networking that are changing the way Services are ordered, configured and monitored. To support the evolving infrastructure, new network and service management platforms need to support standard mechanisms for communication within and across administrative domains. In order to support on-demand, dynamic, configuration and monitoring, both common application programming interfaces (APIs) and a common language that has agreed semantics are required. Standards bodies are using information and data modeling to describe the abstract representations and the detailed structured data needed by the orchestrators and controllers in the ecosystem.

This Feature Topic addresses the standards industries usage and advancements in the area of Information and Data Modeling that support the semantics needed for End-to-End Service Management. Comparing and contrasting the top-down vs. bottom-up approach to API (Application Programming Interface) development is also invited. Solicited topics include (but are not limited to):

- Information modeling
- Data modeling
- Transforming information models to data models
- Service development lifecycle aspects
- End-to-End service management frameworks
- Model driven development
- Modeling tools
- Landscape of YANG models
- Survey of modeling work from industry groups
- Advances needed in network management protocols
- Interaction of Open Source and Traditional Industry Fora and Standards Development Organizations

### SUBMISSIONS

Articles should be tutorial in nature and written in a style comprehensible to readers outside the specialty of the article. Authors must follow the IEEE Communications Magazine's guidelines for preparation of the manuscript. Complete guidelines for prospective authors can be found at <http://www.comsoc.org/commag/paper-submission-guidelines>. It is very important to note that the IEEE Communications Magazine strongly limits mathematical content, and the number of figures and tables. Paper length should not exceed 4,500 words. All articles to be considered for publication must be submitted through the IEEE Manuscript Central site (<http://mc.manuscriptcentral.com/commag-ieee>) by the deadline. Submit articles to the "March 2016/Semantics for Anything as a Service" category.

### SCHEDULE FOR SUBMISSIONS

- Manuscript Submission Due: September 15, 2015
- Decision Notification: November 15, 2015
- Final Manuscript Due: January 1, 2016
- Publication Date: March 2016

### GUEST EDITORS

Scott Mansfield  
Ericsson Inc.  
scott.mansfield@ericsson.com

Hing-Kam Lam  
Alcatel-Lucent  
kam.lam@alcatel-lucent.com

Nigel Davis  
Ciena  
ndavis@ciena.com

Yuji Tochio  
Fujitsu  
tochio@jp.fujitsu.com

# Software Radio: A Catalyst for Wireless Innovation

Christophe Moy and Jacques Palicot

## ABSTRACT

In this article, we deal with software radio in an original manner, emphasizing the fact that software radio is a major evolution of radio technologies and a convergence of different pre-existing fields. Joseph Mitola deserves credit for formalizing and conceptualizing this evolution. We do not define software radio. Rather, we look at software radio from different perspectives. New equipments include more and more software radio “capacities,” that is, they can be reconfigured thanks to software. Certain former analog functionalities are now performed digitally. We think that the timeframe during which software radio is advantageous, compared to a velcro<sup>1</sup> approach for multi-standard terminals, is probably over. However, software radio has taken concrete form in the military field. We show that software radio is used de facto in all equipment and products because of the technological progress to which it has led. Software radio’s legacy indeed is that it has been and still is a catalyst for wireless innovation. We conclude by showing that software radio will be used to its full potential when it is the support technology for other more complex fields such as cognitive radio.

## INTRODUCTION

In the same way that automotive engines have rested for more than 100 years on the same principle (explosion) and mechanical parts (pistons, valves, crankshaft, etc.), radio has been resting on the same architecture (superheterodyne) and electronic parts (amplifiers, mixers, filters, etc.). Some technological evolutions, such as the transistor, have of course enabled miniaturization of the systems and increasing their performance, but the same basis of discreet electronic components has been kept. Software radio [1] is the most recent major technological change in the field. Indeed, the amazing evolution of digital technologies during the 1980s enabled this well established paradigm to be challenged, and since the beginning of the 1990s,<sup>2</sup> certain baseband processing techniques started to be performed and gathered within digital integrated circuits. This development led to software radio, the main ideas of which are that circuit processing is

a general-purpose processor, and that some of the processing that was previously performed analogically by means of several discreet components can now be performed by a sole processor.

First of all, it is interesting to put software radio in its historical context in order to fully understand the interests that led to its advent, at both the design and utilization phases of equipment. Software radio is a convergence of different technological fields. Hence, each concerned scientific community has appropriated the concept and uses this naming to illustrate its work. This fact entails as many standpoints as there are concerned communities: this is dealt with in the third section. The objective reality of software radio, due to the technological developments it has generated, are highlighted after that. The software radio vs. velcro approach is discussed in the final section, in which we also present our vision for the future of software radio.

## THE FOUNDATIONS OF SOFTWARE RADIO

### GENESIS OF SOFTWARE RADIO

In the 1980s and '90s, the U.S. Department of Defense’s (DoD’s) long-term prospects aimed to anticipate what would happen thanks to more futuristic studies taking place over the next 10- to 20-year timeframe. The opportunity to carry out such more conceptual research, which would be less subject to constraints and practicalities, led Joseph Mitola to define the founding concepts of software radio<sup>3</sup> [1, 2].

One can consider that software radio, or at least its main concepts, would have existed without the work of Mitola. In fact, in the 1990s, the following notions were very trendy:

- Convergence of standards and networks
- Reconfigurability of equipment
- Connectivity at any place, any time, and with any standard
- Distributed intelligence in networks and terminals

It is therefore important to highlight the impact of Mitola’s work, which contributed to homogenize all those preexisting notions in a single and formalized concept named software radio [3, 4].

Software radio is a technological break-

The authors are with CentraleSupélec/IETR

<sup>1</sup> A dedicated hardwired module per standard.

<sup>2</sup> Since the beginning of the 1970s, from certain points of view, in military research.

<sup>3</sup> “Software radio” was first coined in the 1980s, in the company E-systems (renamed Raytheon since) where Joseph Mitola worked.

through, which is the fruit of a major development and a convergence of different preexisting fields. Many companies or research teams work on the topic of software radio one way or the other, without knowing it. Software radio brings together a wide spectrum of techniques, each of them representing a specific subject, but without necessarily being integrated in software radio.

However, the scientific community has quickly supported the idea of software radio as illustrated by the curve in Fig. 1 (evolution of software-radio-related articles since 1995, curves obtained with Google Scholar). The first *software defined radio* (SDR) articles were published when the SDR Forum<sup>4</sup> was set up in 1996. The number of *software defined radio* labeled publications is similar to that of *software radio* labeled publications in 2003–2004. This shows that the field then becomes mainly a concrete problem to be solved, even if it is necessary to allow gaps with the original principles (*software radio*) in order to find practical solutions (*software defined radio*). A slight decrease in the number of articles has been observed in the past two years. Unfortunately, this does not reflect the fact that the subject has been solved. However, as software radio previously encompassed certain topics, some of them are now dealt with under new headings. Cognitive radio is one of those. This is actually probably correlated with the boom of publications dealing with *cognitive radio* in the last few years. However, one can see that cognitive radio has its own existence, which is independent of software radio, and has experienced huge growth since 2005. A new expression has become popular, and we can observe an exponential growth of publications on this topic: software defined networking (SDN). Will it absorb cognitive radio in the long term?

### ADVANTAGES OF SOFTWARE RADIO

There is no real consensus on the exact meaning of the term software radio. We discuss this in the next section. Numerous classifications exist, which experienced more or less success when they were published. Interested readers can refer to [4] to deepen their knowledge of those classifications. Software radio can be summarized by the following idea: it symbolizes the transition from the design of transmission/reception radio devices made with numerous dedicated electronic components as had been done for 100 years, to the design of such devices made with general-purpose parts (we consider here processing units rather than electronic components). The processing unit, which is typically a processor, enables several different operations to be performed. The latter are typically functions using computer language, that is, software, hence the name software radio. Therefore, one can easily understand the benefits that can be derived from it. They are summarized in Table 1.

### DIFFERENT PERSPECTIVES ON SOFTWARE RADIO

As we have seen, software radio is a convergence of different technological fields. Sometimes the term software radio may provoke some misunderstanding. Hence, each concerned com-

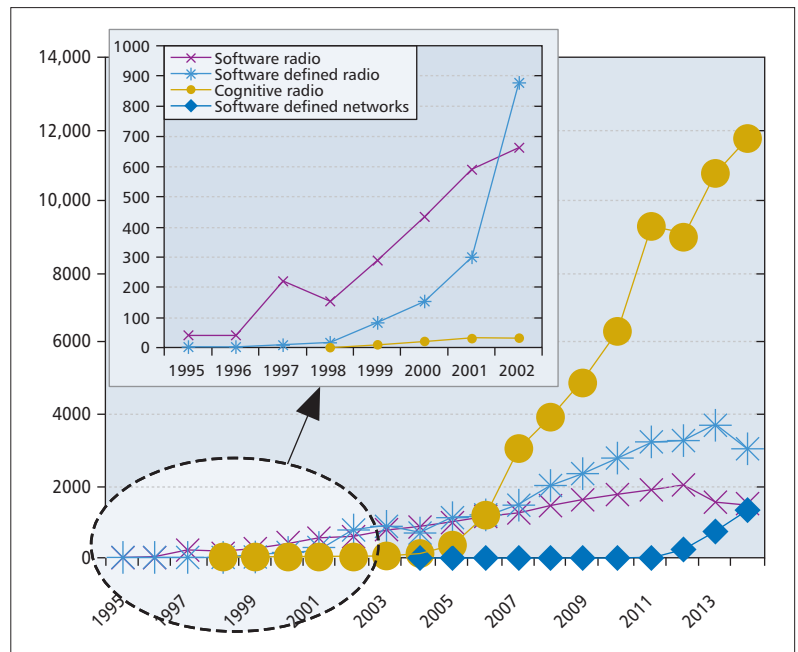


Figure 1. Evolution of the number of software-radio-related articles published since 1995 (based on Google Scholar).

munity has appropriated the concept and uses its own naming to illustrate its work. This means that there are as many standpoints as there are concerned communities. Consequently, instead of defining software radio, we propose to detail some of those in the following subsections.

We have already considered software radio under the prism of a sole general processor that executes radio processing. However, we see below that certain stakeholders think they achieve a software radio design when processing is performed digitally, on a configurable application-specific integrated circuit (ASIC), even without the intention of changing the parameters during operation. Other ones consider the frontier between software and traditional radio from a purely architectural point of view, as a question of RF architecture, as mentioned below. The place of digitization in this architecture is then a distinguishing factor for the different types or levels of software radio. Software radio can thus refer to RF digitization in some cases, or intermediate frequency (IF) digitization in other cases. In the military field, a software radio is a unit of equipment based on a specific software architecture, as explained below. In order to introduce this section, we sum up in Table 2 some of the point of views we do not detail further.

### DIGITAL PROCESSING ARCHITECTURE

One of the basic principles of software radio is to perform processing by means of a general-purpose processor (GPP). At the end of the 1990s, Vanu Bose [7] used that technique to make GSM base stations benefiting from the acceleration speed of the processors with time (Moore's Law), for example, executing in each processor generation a greater number of simultaneous communications on one processor, without changing a single line of code. But several limits oppose the generalization of this approach

<sup>4</sup> Renamed Wireless Innovation Forum in December 2009 (<http://www.wirelessinnovation.org>).

to mobile communicating devices: price, power consumption, and processing units' performance. In fact, a combination of the different processing unit categories of Table 3 provides the best compromise. It is the heterogeneous architecture concept, that is, an architecture comprising digital signal processors (DSPs), field programmable gate arrays (FPGAs), and GPPs. This is the trend followed by the manufacturers of the so-called system on chip (SoC) or network on chip (NoC) components, which include several units of a different nature within a unique chip. Actually, SoC providers, like Qualcomm and Samsung, currently dominate the markets of mobile-phone-dedicated circuits [4].

From the standpoint of the designers of RF transmitters/receivers, software radio is a question of where digitization is performed: RF, IF, or baseband. A pure software radio should be able to demodulate at reception (modulate at transmission) all communication standards thanks to software that can be used on any platform. A traditional radio performs in an analog fashion, by means of dedicated hardware, all the functions of the RF transceiver (channel selection, suppression of interference, amplification, and baseband transposition). Unlike the latter, a software radio would sample the broadband RF signal directly after fil-

<b>Benefits of software radio at design phase</b>	<b>Very first advantage that was targeted at the beginning of software radio in the 1990s</b>
Transition from analog to digital	<p>Digital instead of analog:</p> <ul style="list-style-type: none"> <li>• Electronic design automation (EDA) tools that have tremendously progressed in number and quality for the digital domain these last 20 years</li> <li>• Digital compensation of analog defects (dirty RF) [5]</li> <li>• Cheaper digital solutions than analog</li> </ul> <p>Flexibility of the digital world:</p> <ul style="list-style-type: none"> <li>• Reprogrammability or reconfigurability</li> <li>• Require adequate hardware at runtime</li> <li>• Need adequate software deployment [6]</li> </ul> <p>Digital way that opens a new era merging IT and telecommunications:</p> <ul style="list-style-type: none"> <li>• Convergence of radio and computer science</li> <li>• At both the radio equipment and network levels (SDN)</li> </ul>
Application and platform separation	<p>Earlier design of the hardware platform in the development cycle:</p> <ul style="list-style-type: none"> <li>• Delay resulting from needed hardware adjustments is not so painful</li> <li>• New versions can be made within the specified time limit</li> </ul> <p>Reuse of pre-existing hardware components:</p> <ul style="list-style-type: none"> <li>• Processing units instead of dedicated hardware components</li> <li>• New components instability is avoided</li> <li>• Standard instead of tailor made</li> </ul> <p>(Software) application can be improved until the very last minute before its delivery:</p> <ul style="list-style-type: none"> <li>• Last moment errors can be fixed at reduced cost</li> <li>• Application changes do not entail any hardware adjustments</li> <li>• Transferring design difficulties in the programming phase of the application</li> <li>• More comfort and better robustness</li> </ul> <p>Reutilization of the same platform for different products:</p> <ul style="list-style-type: none"> <li>• Several radio applications in a unique multi-standard product (unlike a velcro design, which comprises as many circuits as there are radio applications)</li> <li>• Several products whose operations are specific to different software versions</li> </ul> <p>Third parties new market:</p> <ul style="list-style-type: none"> <li>• Radio application waveforms</li> <li>• Radio design environments and software facilities</li> </ul>
<b>Benefits of software radio in terms of functionality</b>	<b>This is the main and true novelty of software radio, according to us [6], due to the possibility it offers to benefit from added flexibility during the lifetime of the product, that is, once it has been manufactured and is on the market. In that respect, software radio still remains a driver of wireless innovation since this capacity is far from being fully used in most radio systems.</b>
For mobile phone providers, wireless operators and users	<p>New technical means:</p> <ul style="list-style-type: none"> <li>• Over the air reconfiguration (OTAR)</li> <li>• Download new complete software (e.g., a radio application or waveform)</li> <li>• Correct defects through patch</li> </ul> <p>New services:</p> <ul style="list-style-type: none"> <li>• Adapt all processing thanks to a simple change of software</li> <li>• Offer an ubiquitous connexion, that is, which is able to demodulate at reception (and respectively to modulate at transmission) all communication standards with the same equipment</li> <li>• Minimize the impact on clients of maintenance operation thanks to OTAR</li> <li>• Extend product life duration thanks to software update can concern an improvement of the radio equipment, or even the download of a new standard that did not exist when the product was sold, such as for satellite industry and remote radio access point for operations and maintenance (O&amp;M).</li> </ul>

**Table 1.** Software radio benefits.

tering and low noise amplification. Thereafter, the digital processing unit performs the following operations: frequency transposition, channel selection, and demodulation, thanks to software processing.

The software radio paradigm cannot be implemented for a holistic radio, supporting all radio applications. Thus, lots of works have focused on architectures dedicated to a subset of the radio domain, and then featuring more relaxed constraints. The term software defined radio is then used, since the digitized band is reduced. In the literature, there are mainly three types of SDR from the RF point of view, which are described in [4, Ch. 6]: direct conversion, low-IF, and undersampling. Many consumer products include such software radio capacities, implementing direct RF digitization in a given band, or already digitally perform IF processing.

Since one bottleneck of the ideal architecture is conversion (analog-to-digital and digital-to-analog), the conversion community has heavily invested in research to solve the problem raised by software radio. Even if the problem has not yet been fully resolved, lots of progress has been made, and a software defined radio can now be built (e.g., in low to intermediate frequency). We find, for example, products called “Wideband SDR Solution,” which include on a single chip all the wideband conversion functions (12 bits, 4 Gsamples/s on the ADC12J4000 from Texas Instruments) as well as the filtering and frequency transposition functions.

It must be noted that another bottleneck of the ideal architecture is the required computation power, which is directly linked to the sampling frequency, in which case we now have to deal with the considerations mentioned previously in the section on digital processing architecture. Requirements in terms of computation power and digitization speed also raise the power consumption issue, which involves software radio approaches. Moreover, other functions of the ideal architecture may prove to be problematic; for instance, the large bandwidth antennas and the power amplification of a broadband signal featuring a multiplex of modulated carriers.

### SOFTWARE RADIO: SDR AND SCA

The military software radio community always uses the term software defined radio (SDR). It is important to remember that this community was the founder of the field, since Mitola’s works on the subject were funded by the U.S. DoD. This is why this is mainly an American phenomenon. Besides, the DoD has decided to opt for the software radio approach in order to facilitate the interoperability of the communication means of some of the American armed forces (Navy, Air Force, and Army): this is a fundamental characteristic of what software radio can offer. Using software radio to find solutions to interoperability issues makes even more sense within the NATO framework where the forces of several countries have to intercommunicate. This need is also present in the mobile communications of the civil defense forces (PMR<sup>5</sup>) where SDR is sometimes also considered in a way that is similar to the military approach [9]. It is interesting to know that, beyond the purely military aspects, another event has largely contributed to increase

Point of view	
Signal processing	<ul style="list-style-type: none"> <li>• Digital processing of formerly analogically processed functions</li> <li>• Correction of the remaining analog functions’ defects</li> <li>• Advanced signal processing algorithms in general</li> </ul>
Computer science	<ul style="list-style-type: none"> <li>• Software architecture standards</li> <li>• Downloading protocols (cf. the OTAR example given earlier)</li> <li>• High-level software that enable to model electronic engineering applications</li> <li>• Real-time OS, hardware/software distribution tools, etc.</li> </ul>
Application-specific software radio	<ul style="list-style-type: none"> <li>• Limited to a specific application, it is technologically possible to achieve a software radio realization in the case of applications featuring small bandwidth, low carrier frequencies, consumption that is not limited by mobility, and medium cost-related constraints</li> </ul>
Amateur radio	<ul style="list-style-type: none"> <li>• Amateur radio community is the first one to use software radio on a daily basis</li> <li>• Low-cost low-complexity software radio platforms and applications</li> <li>• USB keys for DVB-T demodulators can cover ranges of frequency from a few tens of megahertz up to 2 GHz for a few tens of dollars (R820T tuner chip)</li> <li>• They can be piloted by a free software environment run on a laptop for different radio amateur applications, as amateur radio satellite with Gpredict, for instance, as aeronautic data with ADS-B# from more generic SDR#, and so on.</li> </ul>

**Table 2.** Different points of view for addressing software radio. DVB-T: digital video broadcast — terrestrial; ADS-B: automatic dependent surveillance broadcast.

the DoD funding of software radio: 9/11. Indeed, communication gaps on that day led to losses that could have been avoided through more efficient coordination between police officers and firefighters. To a lesser extent, the loss of civil communications also meant increased confusion and contributed hindrance of emergency aid. Software radio appears as a solution that could enable quick redeployment of the interoperable communication means in case of a disaster destroying all the infrastructures of a region [9].

SDR is often equivalent to software communication architecture (SCA) in the military field. SCA is a software architecture abstracting the underlying hardware (HW) platform that runs radio applications or waveforms. This enables radio waveforms to be developed independent of the HW platform, via standardized interfaces (application programming interface, API). SCA is a standard originating from studies financed by the DoD via the Joint Tactical Radio Systems (JTRS) since 1997, with the aim of specifying a software radio’s development standard [10]. Any manufacturer of a communication system or subsystem provider that wishes to respond to the DoD’s calls for tender must comply with the SCA standard. In theory, the SCA must enable the systems’ integrators to use any hardware or software done by one or more providers, in an efficient and open manner. Keeping in mind the secrecy imposed on the making of military systems, it is thus easy to understand the impor-

<sup>5</sup> Professional Mobile Radio, also named Private Mobile Radio.

tance of the challenge, beyond any purely technical considerations that represent other challenges. SCA was specified in the middle of the '90s on the basis of keywords such as “portability,” “heterogeneity,” “multi-standard,” and “interoperability.” The technical solutions chosen at the time did not exist in the embedded real-time electronics field, and the answers provided by computer science were then chosen, since they were the only existing ones — Common Object Request Broker Architecture (CORBA), Interface Description Language (IDL), Posix, and so on.

The European military is also adopting “SCA-inspired” solutions for communications systems. A European adaptation of SCA has been developed with the framework of the European Defense Agency’s (EDA) European Secure Software Defined Radio (ESSOR)<sup>6</sup> project and is now entering European military products. The French Ministry of Defense, for instance, announced on April 17, 2012, the launch of the strategic CONTACT program, featuring a budget exceeding €1 billion. The CONTACT<sup>7</sup> program aims at providing, as of 2018, the French armed forces with new generation tactical radio sets, resting on an innovative “software radio” technology.

Some civil SDR standardization efforts should also be pointed out, such as the European Telecommunications Standards Institute Reconfigurable Radio System (ETSI RRS), which aims to specify reconfigurable radio architectures in [11], and an attempt in IEEE P1900.3 on conformance evaluation for SDR software modules,

which has been disbanded. However, civil standardization is more active on cognitive radio aspects in both ETSI and IEEE.

## SOFTWARE RADIO AS A CATALYST FOR TECHNOLOGICAL ADVANCES

Software radio has generated intense research activity for each of the previously presented standpoints. It has brought together numerous studies and collaborative projects. This has resulted in many technological advances, and many research subjects remain to be studied in the future. Let us highlight a few technological advances for each previous standpoint.

### DIGITAL PROCESSING ARCHITECTURE

The development of heterogeneous architectures, comprising DSPs, FPGAs, and GPPs, has been strongly boosted by software radio. The move toward ever more efficient digital architectures has enabled manufacturers to propose software radio base stations. The NoC or SoC architectures that are installed on chips and include programmable or configurable units are commonly used in the most recent mobile devices. While architectures were moving toward heterogeneity, another fundamental evolution took place at the same time: making hardware more and more “flexible” so that it can be used similarly to software. This led, for instance, to the partial reconfiguration concept for FPGAs. The company Xilinx, which is the main provider of FPGAs that can cater for this capacity, has developed and reserved the partial reconfiguration [8] market as the sole software radio stakeholders for several years in order to launch this market. A last new trend toward a coarse-grained processing strategy may be arising in the software radio field: the manycore era [12]. Manycore, or chip multi-processor (CMP), which can have several thousand processing cores, applies the FPGA paradigm at the CPU level, taking advantage of a “processing array” instead of a “gate array,” coming back to software radio basics in processing homogeneity. A further gap is also reached in this field as software radio itself is foreseen inside the chip in order to improve inter-core communications [13].

### RF ARCHITECTURE

The direct conversion architecture has become achievable since its main defect, which resulted from the mix between the local oscillator and the RF signal (cf. [4, Ch. 6]), has now been corrected by manufacturing companies. The analog-to-digital and digital-to-analog conversion community has heavily invested in research to solve the issue raised by software radio. It is now easy to find broadband converters, featuring 12 bits and several gigahertz sampling frequencies. Software radio also requires broadband or multi-band antennas. Numerous advances have been achieved in this field too. Software radio also imposes new and important constraints on the amplification stages. A large band of frequencies, which feature a great number of modulated carriers, must indeed be amplified. To that purpose, software radio has once again generated numerous studies in order to solve the problems linked to wideband amplification, such as new

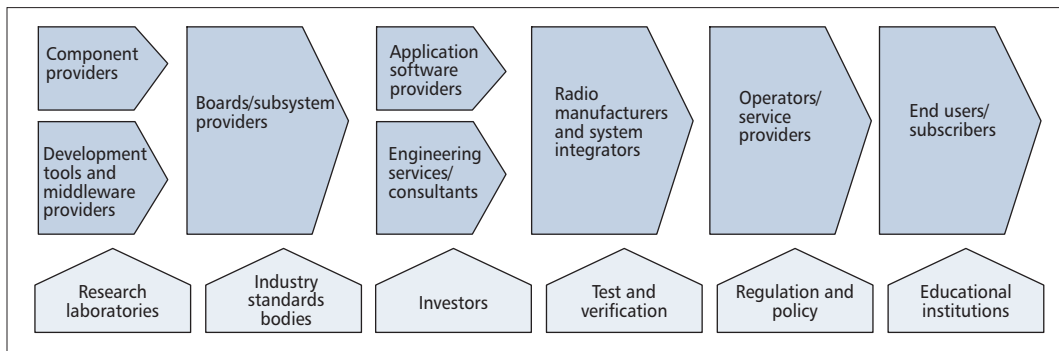
<sup>6</sup> <http://www.eda.europa.eu/otheractivities/sdr/essor>

<sup>7</sup> <http://mil-embedded.com/articles/conquering-radio-market-challenges/>

Processing unit	Software radio requirements in terms of computation power and power consumption, so as to execute real time modulation/demodulation, are so important that it is often necessary to use several processing technologies instead of GPP only:
GPP: general purpose processor	<ul style="list-style-type: none"> <li>• Adequate for complex computation, data storage management, and connectivity, with a high degree of flexibility, but very high power consumption</li> <li>• Particularly usual in base stations</li> </ul>
DSP: digital signal processor	<ul style="list-style-type: none"> <li>• Specialized for signal processing computation, with a high degree of flexibility and low power consumption</li> <li>• Particularly suited for terminals</li> </ul>
FPGA: field programmable gate array	<ul style="list-style-type: none"> <li>• Parallel intensive computation, with higher power consumption than a DSP, and with a degree of flexibility at design time only (usual FPGA use), or a high degree of flexibility at runtime (using ultra-fast dynamic partial reconfiguration [8]),</li> <li>• Particularly suited for base stations</li> </ul>
ASIC: application-specific integrated circuit	<ul style="list-style-type: none"> <li>• Hyper intensive computation, ultra-low power consumption and for a large number of parts produced, with a very low degree of flexibility</li> <li>• Particularly suited for terminals for power consumption considerations, and also in base stations for ultra-high computation performance</li> </ul>

**Table 3.** Strengths and weaknesses of processing units from a software radio perspective.





**Figure 2.** WinnF's view of the software radio value chain.

peak-to-average-power ratio (PAPR) issues related to the mix of multi-carrier/multi-standard signals [14]. IF sampling has become popular in today's products, which implies performing certain digital functions such as modulation/demodulation, selection of the channel of interest, and adaptation of the sampling frequency. These functions operate thanks to signal processing algorithms [4].

### SOFTWARE ARCHITECTURE

We saw previously that in this field, the main result is the SCA definition, which was devised originally for the military field. The needs in terms of embedded real-time OS generated by software radio have contributed significantly to push embedded OS progress. The optimum hardware/software distribution on the different processing units (hardware or software co-design) implies a high-level model process for both the application and software, as this was mentioned earlier. Those working in amateur radio have developed many free software radio environments that can be used with a few dollars of hardware (DVB-T dongles).

### OTHER POINTS OF VIEW

In the signal processing field, the new "dirty RF" concept originates directly from software radio work [5]. This concept will generate a large amount of research in coming years. As they focused on a specific application, which occupies a small bandwidth supported by low frequencies, and used software radio concepts, numerous manufacturers have proposed some so-called software radio products for their application. All the latter considerations prove that software radio features multiple realities, which depend on the different standpoints of the different scientific communities. Figure 2 illustrates the Wireless Innovation Forum (WinnF's) view of the software radio value chain. It can therefore be deduced that software radio has entered, through all these technological advances, into products and equipment.

## SOFTWARE RADIO AND THE FUTURE

It is already possible to make a software radio in a laboratory, without manufacturing and power consumption constraints, for specific applications. But the concept of a universal radio that could cater for a wide range of telecommunication applications, in particular for daily life (WiFi, multi-band 2G, 3G, 4G, TV broadcast, FM, Blue-

tooth, GPS, etc.) is not commercially viable in a terminal. This is corroborated by the fact that all current mobile phones include these radio links thanks to an equal number of dedicated circuits or sub-circuits (perhaps several in one single multi-standard chip). This is known as the velcro approach: a component (e.g., a scratch) is added as another radio functionality is added. This approach is, without question, very efficient to move from one standard to another. However, it is not upgradeable, since the processing of a new standard, which was not originally planned in the making of the circuit, requires the making of a new circuit and the sale of a new device.

In the base station, on the contrary, due to lower constraints on power consumption and cost, software radio is in common use for its interesting features in terms of wideband digitization and upgrade capabilities. Many manufacturers propose some products at different stages of the software radio's value chain. Those different stages of the value chain fit in with the different standpoints mentioned previously. Thus, waveform application software and execution platforms, software and middleware providers, integrators, and other stakeholders can also be found in that chain.

We thought that, since the wireless connection means are multiplying, a threshold beyond which it would not be viable to have many velcro circuits would be reached, even in terminals. Nevertheless, the continuing progress in the integration and microelectronics field will certainly offer the possibility to include even more velcro standards. This approach offers sufficient capacity for the current flexibility needs, whereas those of software radio are well above those needs. It can be inferred that, with regard to the multi-standard technology aspect, the software radio's window of opportunity is probably over for terminals. However, when all its capacities are necessary to offer applications using cognitive radio, there will be a new window of opportunity and relevance for software radio inside terminals. Software radio will no longer be considered as a simple alternative to the current technologies, but as a truly new service solution. In a vertical handover context, software radio will facilitate the transition between standards during communication. In that respect, as we wrote in [4], software radio becomes a support technology for cognitive radio.

We think that software radio, as a whole, will not become a standard as such, since software radio's capacities are introduced through a "nat-

*While architectures were moving toward a heterogeneous nature, another fundamental evolution took place at the same time: it consisted in making hardware more and more "flexible" so that it can be used similarly to software. This led, for instance, to the partial reconfiguration concept for FPGAs.*

Instead of combining software pieces coming from a radio library, as with GNU Radio, one can build one's own mobile phones, piece by piece, from COTS26 batteries, screens, camera, WiFi module, and so on, and reconfigure it at the hardware level. Even if not reduced further to a software sub-set, this is a completely "software radio" philosophy!

ural" evolution of products and equipment. Rather, we think applications using software radio will be standardized. Examples of the latter include dynamic spectrum access [15] and other cognitive radio [4] applications. This one, considered under very limited aspects, is beginning to be integrated into different standards. Indeed, the functions that enable the waveform to be adapted to the environment, which can be found in standards such as WiMAX and Long Term Evolution (LTE), fall under cognitive radio concepts. However, cognitive radio can be fully operational only if comprehensive software radio capacities are available.

But maybe the most valuable contribution of software radio is a philosophical one. Software radio initiated the "open era" of the radio world. The first impact came at the research level with GNU Radio,<sup>8</sup> which aimed at introducing the GNU paradigm into radio domain. The GNU paradigm was established by the Free Software Foundation and aims to give the right for users to use software pieces for free (associated with certain duties). This has been a great success for academic research on software radio and cognitive radio. However, GNU Radio moved beyond the pure research field and has also been adopted by the amateur radio community, and is also entering companies' prototyping and development strategies. New current trends confirm this philosophical impact, and its propagation into the wireless commercial market. "Open Basestation"<sup>9</sup> is no longer a concept but is becoming a reality, driven by small cell and green approaches, toward a low-cost and "programmable" infrastructure. Most recently, the ARA<sup>10</sup> project from Google (and previously Motorola) has introduced the concept of "make your own mobile phone." Instead of combining software pieces coming from a radio library, as with GNU Radio, one can build one's own mobile phone, piece by piece, from consumer off the shelf (COTS) batteries, screens, camera, WiFi modules, and so on, and reconfigure it at the hardware level. No longer reduced to a software subset, this is a completely "software radio" philosophy!

## CONCLUSION

Software radio was an initiative and has generated intense research activity over the last 20 years. Many technological advances resulted from this, in particular for some components (analog-to-digital converters, NoC, power amplifiers, etc.) used today in numerous products and equipment types. Software radio is an objective reality from this point of view. With regard to multi-standard terminals, the time window during which software radio is advantageous in comparison to a velcro approach is probably over. The software radio stakeholders were not able to propose competitive software-radio-like products, in comparison with velcro products such as smartphones. Results are better on the commercial base station side where power consumption and cost constraints of software radio are easier to support than in terminals. In the military field, software radio has materialized via the SCA at the international level, and via the ESSOR programs at the European level. This success is due to the

important fundings allocated to software radio in the military field. The operational needs of the military have proven that software radio is the necessary technology to respond to multi-standard communications' problems. This was not the case for consumer goods because of the lesser requirements of this market, for which telecommunication stakeholders do not need to use the full potential of the said capacities. Software radio's capacities in terms of flexibility and signal processing are well above the current needs of the market. Software radio will only be able to express its full potential when all software radio's "capacities" are necessary in the future to offer complex services and applications, based on cognitive radio, for instance. In that respect, software radio still holds a bright future.

## REFERENCES

- [1] J. Mitola, "The Software Radio Architecture," *IEEE Commun. Mag.*, May 1995, pp. 26–38.
- [2] J. Mitola, "The Software Radio," *Proc. IEEE Nat'l. Telesystems Conf.*, 1992.
- [3] J. H. Reed and B. D. Woerner, *Software Radio: A Modern Approach to Radio Engineering*, Prentice Hall, 2002.
- [4] J. Palicot, Ed., *Radio Engineering: From Software Radio to Cognitive Radio*, Wiley, 2011.
- [5] G. Fettweis et al., "Dirty RF, A New Paradigm," *Proc. IEEE PIMRC '05*, Berlin, Germany, 11–14 Sept. 2005.
- [6] A. Kountouris, C. Moy, and L. Rambaud, "Reconfigurability: A Key Property in Software Radio Systems," *Proc. First Karlsruhe Wksp. Software Radios*, Karlsruhe, Germany, 29–30 Mar. 2000.
- [7] V. Bose, *Design and Implementation of Software Radios Using General Purpose Processors*, MIT Ph.D. thesis, June 1999.
- [8] J. Delorme et al., "New OPBHWICAP Interface for Real-Time Partial Reconfiguration of FPGA," *Proc. Int'l. Conf. ReConfigurable Computing and FPGAs*, Cancun, Mexico, 9–11 Dec. 2009.
- [9] G. Baldini et al., "The EULER Project: Application of Software Defined Radio in Joint Security Operations," *IEEE Commun. Mag.*, vol. 49, no. 10, Oct. 2011, pp. 55–62.
- [10] "Software Communications Architecture Specification," *JTRS Standards*, v. 2.2.2, 15 May 2006.
- [11] ETSI RRS EN 303 095, "Radio Reconfiguration Related Architecture for Mobile Devices," in progress.
- [12] S. Borkar and A. A. Chien, "The Future of Microprocessors," *Commun. ACM*, vol. 54, no. 5, 2011, pp. 67–77.
- [13] A. Briere et al., "A Dynamically Reconfigurable RF NOC for Manycore," *Proc. 25th Great Lakes Symp. VLSI*, Pittsburgh, PA, 20–22 May 2015.
- [14] J. Palicot et al., "Frequency Domain Interpretation of Power Ratio Metric for Cognitive Radio Systems," *IEE Commun. J.*, 2008, 2 (2009-06-01), pp. 783–93.
- [15] Q. Zhao and B. Sadler, "A Survey of Dynamic Spectrum Access," *IEEE Signal Processing Mag.*, vol. 24, no. 3, May 2007, pp. 79–89.

## BIOGRAPHIES

CHRISTOPHE MOY (Eng. 1995, M.Sc. 1995, Ph.D. 1999, HDR 2008) is a professor with CentraleSupélec (created in January 2015), formerly Supélec, and worked before in Mitsubishi Electric European Research Lab on digital communications from 1999 to 2005. He is also now the co-head of the Communications Department of the CNRS lab IETR (Institute of Electronics and Telecommunications of Rennes — UMR 6164). His research has been focused for 15 years on software radio and cognitive radio. He has been involved in many European and French collaborative projects. He has written more than 20 book chapters, 20 journal papers, and 100 conference papers on software radio and cognitive radio.

JACQUES PALICOT received, in 1983, his Ph.D. degree in signal processing from the University of Rennes. He is now a professor with CentraleSupélec, formerly Supélec. He is head of the Signal, Communications and Embedded Electronics (SCEE) research team. His research focuses on adaptive signal processing, software radio, cognitive radio, and green rRadio. He was Editor of the book *Radio Engineering: From Software Radio to Cognitive Radio* (Wiley, 2011).

<sup>8</sup> <http://gnuradio.org/>

<sup>9</sup> IEEE ComSocCTN Special Issue on "Towards The Open Basestation", February 2015, <http://www.com-soc.org/ctn/ieee-comsoc-ctn-special-issue-towards-open-basestation>

<sup>10</sup> <http://www.projectara.com/mdk>

# The Software Communications Architecture: Two Decades of Software Radio Technology Innovation

*Claude Belisle, Vince Kovarik, Lee Pucker, and Mark Turner*

## ABSTRACT

What began more than 20 years ago as a U.S. Department of Defense project aimed at implementing a radio system that could be reprogrammed on the fly to support multiple waveforms has evolved to become a widely adopted, versatile, and industry-changing architecture. The resulting Software Communication Architecture has now been deployed in more than 400,000 radios worldwide. This article explores the evolution of the SCA starting with its origins back in the SPEAKEasy Program to the upcoming release of SCA 4.1.

## INTRODUCTION

The Software Communications Architecture (SCA) is an implementation-independent architectural framework that specifies a standardized infrastructure for a software defined radio (SDR). Initially developed and published by the U.S. Department of Defense (DoD), the SCA is maintained by the Joint Tactical Networking Center (JTNC) in collaboration with various industry partners and organizations, such as the Wireless Innovation Forum [1]. The specification has significantly influenced the evolution of the SDR domain, and its concepts have been used within multiple industries, products, and countries worldwide as depicted in Fig. 1 [2].

Advances in digital processor technology, increases in analog-to-digital sampling rates, and other technological developments have enabled the continuing growth of complex signal processing in the digital domain. This increase in digital processing has appreciably altered the architecture and design of radio systems. Recent generations of SDRs evolved to become highly software-intensive complex systems facilitating further advancement of communications capabilities. SDRs have enabled more cost-effective radio platform life cycles by providing for the update and addition of system functions and features without requiring hardware modifications. Prior to the establishment of the SCA as an open standard, SDRs were developed using pro-

prietary software architectures that tightly coupled hardware platforms and waveform applications in a manner that was unique to each SDR manufacturer. The SCA has built on the capabilities of these preceding generations of SDRs, moving today's radios substantial steps further forward by leveraging large-scale commercial software industry investments in technology and promoting open standardization. The SCA specification and associated technologies facilitate broad software reuse and application portability across SDR platforms, while enabling achievement of the key industry-wide objectives:

- Enhanced interoperability between SDRs and across entire radio communications systems, especially important for mission-critical communications
- Reduction of the time and cost required to develop and deploy SDRs and associated systems, including the incremental rollout of new SDR and radio system features and functions

The SCA provides a set of rules and constraints that define the interactions between software applications (i.e., waveforms) and radio hardware platforms, leveraging an object oriented (OO) software paradigm and employing component-based development (CBD) technologies. CBD technologies are sometimes referred to as the "industrial revolution" of software, fostering the advent of interchangeable software parts built to predefined specifications [3]. With CBD technologies, software components can be thought of as software integrated circuits with a set of defined functionality, performance, and input/output. Components can be assembled together to create entire applications, such as waveform applications for an SDR. The SCA specification also defines a core set of open system interfaces and profiles that provide for the configuration, assembly, deployment, and management of components, which ultimately comprise the software waveform applications. The components of these software waveform applications can be distributed across various SDR hardware processing elements in a manner determined by the particular radio developers that

---

*Claude Belisle is with Nordiasoft.*

*Vince Kovarik is with PrismTech.*

*Lee Pucker is with the Wireless Innovation Forum.*

*Mark Turner is with Vanguard Wireless.*

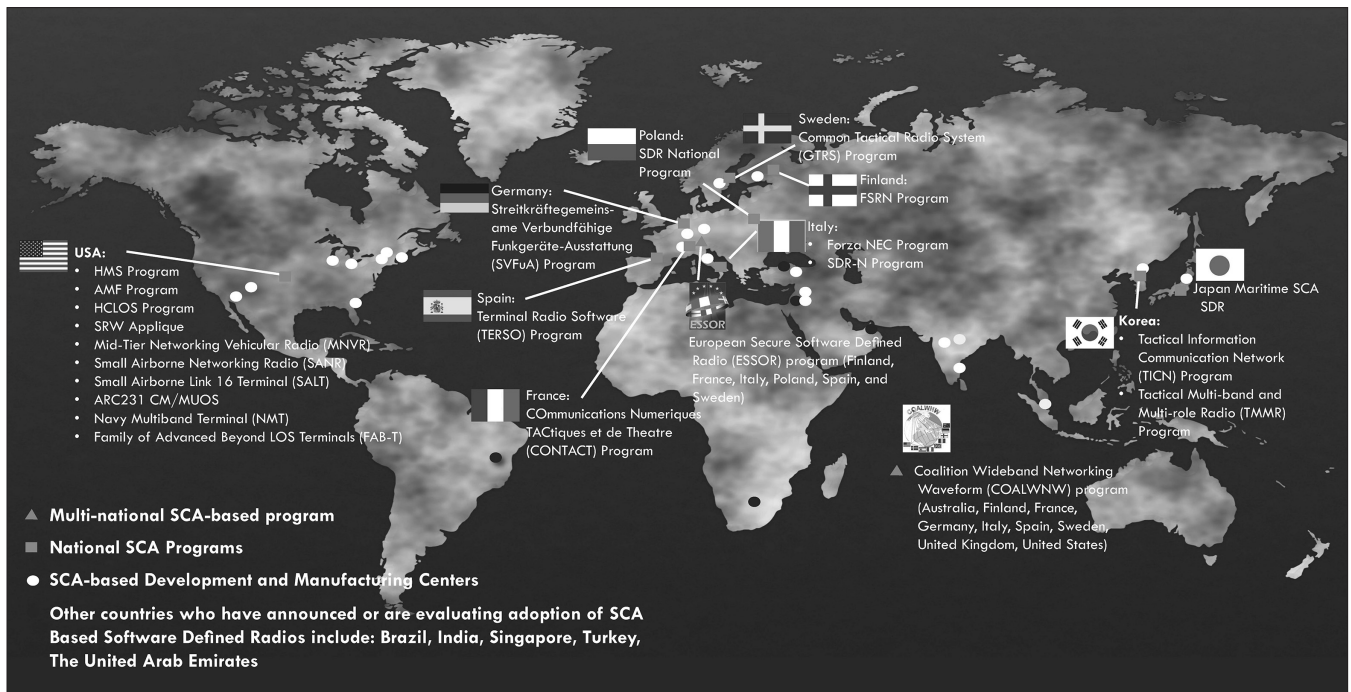


Figure 1. SCA utilization worldwide.

support the overall SDR requirements, radio hardware platform capabilities, and design, in conjunction with the SDR software design and configuration.

## ORIGINS OF THE SCA

Early development of the SCA began under the U.S. DoD SPEAKeasy project. SPEAKeasy began in 1991 with a goal of implementing an SDR system capable of being reprogrammed while in operation to support multiple air interface standards, referred to as waveforms [4]. This capability was very attractive to the military community as a means to reduce the long-term operations and maintenance costs associated with hardware-based radio systems. The SPEAKeasy platform consisted of an array of TMS320C40 processors from Texas Instruments as well as other special-purpose processors mapped onto an industry standard backplane. The platform also included a separate transceiver subsystem and programmable information security (INFOSEC) module. In moving beyond SPEAKeasy, funding was provided for some proof of concept systems such as the Joint Combat Information Terminal (JCIT) and Digital Modular Radio (DMR). These systems further established the feasibility of implementing and deploying an SDR.

Based on the experience and lessons learned in creating the SPEAKeasy and follow-on systems, it became clear that a common radio management system was required to deploy, configure, and manage the signal processing software and other components of the waveform. This management system was developed as a custom software module for each of the early projects, which limited the reuse of application code across platforms. The push toward a common application framework and operating environment to future-proof the system and allow

multiple teams to develop waveforms for the platform in parallel started in the SPEAKeasy Phase II program [5], which began in 1995.

Elements of this software architecture were introduced to the Modular Multifunction Information Transfer Systems (MMITS) Forum, which was created at the request of the U.S. Air Force as an industry association focused on advancing the development of software radio technology. In 1997, the MMITS Forum published its Technical Report 1.0, defining “Architecture and Elements of Software Radio Systems” [6, 7]. In 1998, the MMITS Forum rebranded as the Software Defined Radio (SDR) Forum, and continued to mature this architecture through its Mobile Working Group, resulting in the “Software Radio Architecture” (SRA) published in Technical Report 2.1 as an industry standard in 1999 [8].

The SRA was designed to support the functional interfaces outlined in Fig. 2. Key requirements for the architecture included scalability, allowing the architecture to be implemented across as wide a range of radio platforms as possible, and upgradability, to allow new capabilities and waveforms to be added to a radio without replacing the underlying hardware. To achieve these goals, the SRA adopted an OO core framework consisting of a domain manager, file manager, resource manager, and devices necessary to set up, tear down, and control waveform applications. The Common Object Request Broker Architecture (CORBA) was selected as the middleware layer or “software bus” in this architecture, and a common operating environment was also defined based the POSIX specification. Applications utilizing this core framework are created based on the Object Management Group’s component model, and were instantiated on the radio platform through the use of an “application factory.”

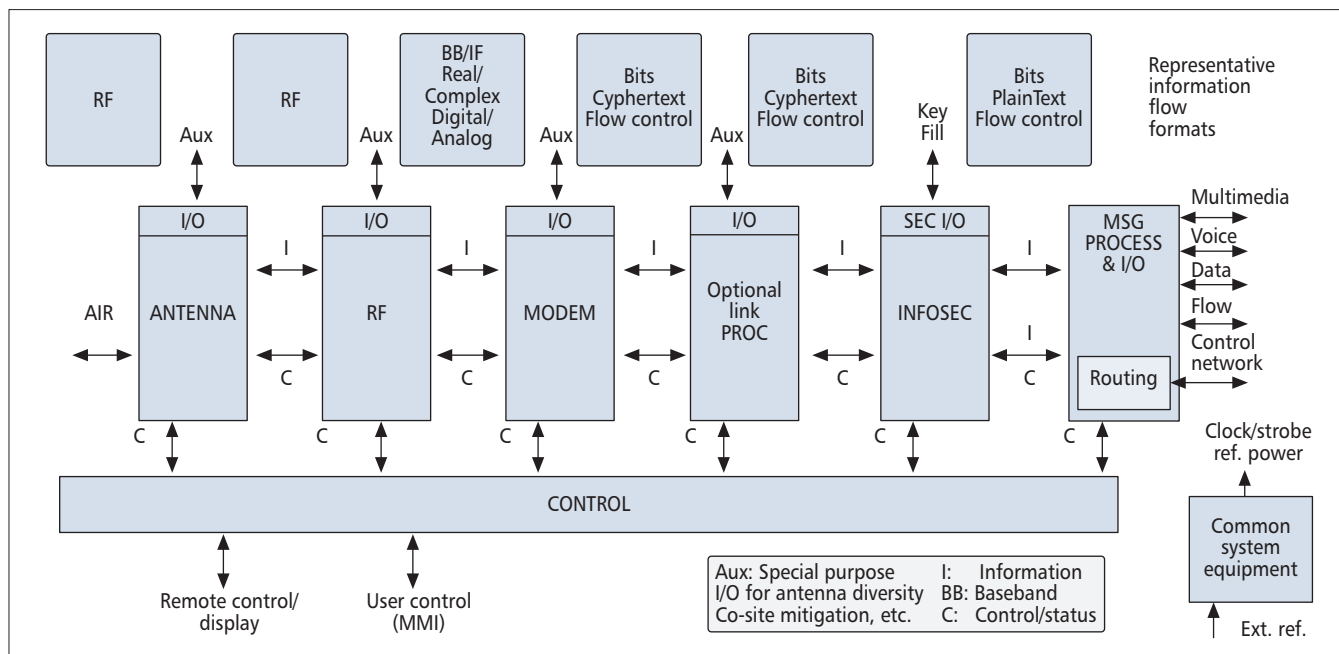


Figure 2. SDR Forum functional interface diagram.

## JOINT TACTICAL RADIO SYSTEM

In the late 1990s the Joint Tactical Radio System (JTRS) Joint Program Office (JPO) was formed to develop a new family of software-based reconfigurable radio systems. One of the first activities was to define a common software infrastructure that would be applied to this new family of radio systems. Initial project funding was provided to address the development of a radio framework, and the Modular Software Radio Research Consortium (MSRC) was formed to develop a common radio framework. The concepts in the SRA were further matured by the members of the MSRC working in cooperation with the Forum's Mobile Working Group, and this resulted in what became the JTRS Software Communications Architecture (JTRS SCA)[9].

### EARLY SCA VERSIONS

The timeline in Fig. 3 illustrates several key milestones in the evolution of the SCA specification. In the late 1990s there were several preliminary versions prior to v. 1.0, which was released in early 2000. Although portions of the SCA were implemented as prototypes, the specification was not sufficiently mature and complete to implement a complete software framework for radio applications. Version 1.0 was the first version considered to be sufficiently complete for a full prototype, and several prototypes were implemented by members of the MSRC.

An incremental release, v. 1.1, was used for a proof of concept demonstration that the SCA could be used to configure and manage an existing radio system. One of these demonstrations, performed in 2000 as part of step 2B of the JTRS program, was a collaborative demonstration by Harris Corporation and Exigent. In this demonstration, a Harris Corporation tactical radio was interfaced to a PC running the SCA

core framework (CF) developed by Exigent. It was successfully demonstrated that deployment and configuration of the waveform on the Harris radio system could be performed and managed by the Exigent SCA CF running on the PC.

The demonstration identified several shortcomings in the specification with regard to the control interfaces and the specification and configuration of devices. Another iteration of review and modification of the specification was performed, resulting in v. 2.0. Another round of preliminary implementations resulted in further incremental iterations, and in November 2011 v. 2.2 was released. Version 2.2 was generally considered to be sufficiently mature to use as the basis for the development and deployment of a tactical SDR system.

With the release of SCA 2.2, the U.S. government initiated the procurement process for the first set of SCA-compliant radio systems to be developed. In June 2002, the first major program to apply the SCA was awarded. The Cluster 1 program, later renamed the Ground Mobile Radio (GMR) program, was the inaugural project using v. 2.2 of the SCA. Other JTRS Cluster programs were later awarded. In April 2004, almost three years after the 2.2 specification, v. 2.2.1 was released. This version corrected several errors in the 2.2 specification, and incorporated several clarifications and enhancements.

### SCA 3.0

At the same time, issues with waveform portability were being raised through the ongoing JTRS Cluster programs. The basic problem was that the code developed for a general-purpose processor (GPP) was reasonably portable between platforms. However, the code developed for a digital signal processor (DSP) and field programmable gate array (FPGA) generally remained specific to the particular processor and architecture of the radio.

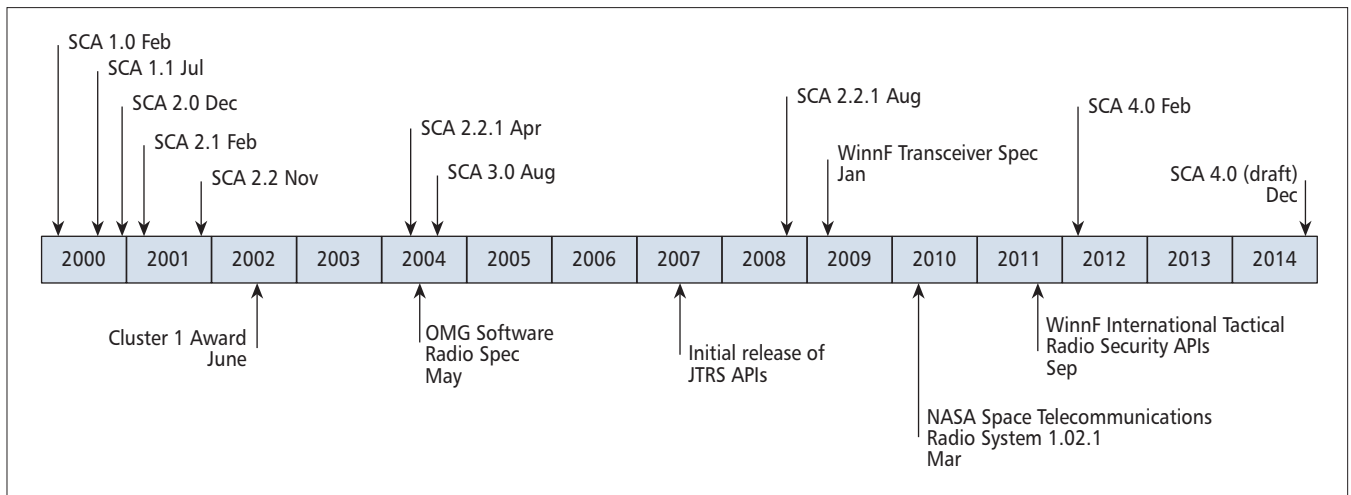


Figure 3. SCA specification timeline.

The DSP and FPGA portability issue came to a head in late 2004, resulting in several special workshops called by the JPO to address the topic. The subsequent result was the SCA 3.0 specification in 2004. This version of the SCA changed few of the core requirements describing the SCA. It did, however, define additional constraints on DSP software related to what system calls could be used by DSP code, proposed a set of waveform components and a high-level data transport design (called HAL-C), and also included an antenna application programming interface (API) section. The general reaction in the community was that the specification required additional work, and, while the concepts and approaches were potentially useful, more detail and analysis were required to achieve a set of descriptions that could be implemented efficiently.

#### JTRS APIs AND SCA 2.2.2

In late 2005 and early 2006, the JPO was reorganized as a Joint Program Executive Office (JTRS JPEO), and undertook a task to revisit some of the waveform portability issues raised by the JTRS program and evaluate alternatives to address the problem. The consensus was that there were insufficient specifications for common APIs to the underlying hardware elements. This resulted in the formation of an API working group to analyze and specify a set of interfaces for the radio system hardware. The initial work drew heavily from the Cluster 1 program and resulted in the specification of several key APIs for common hardware elements. With the release of SCA 2.2.2 in 2006 and the ongoing API work, SCA 3.0 was deprecated and shown on the JPEO website as “not supported.” SCA 2.2.2, as the number implies, is an incremental version from the 2.2.1 version of the SCA.

#### SDR FORUM CONTRIBUTIONS

With the initial ramp up of the JTRS program, the MSRC was dissolved. Most of the MSRC participants were also founding members of the SDR Forum, which consequently acted as a catalyst for the SCA and SCA related activi-

ties. An early activity of the Forum was sponsorship of a project to implement an open source implementation of the SCA. This project resulted in the SCA reference implementation (SCARI) project developed by the Communications Research Centre (CRC). SCARI-OPEN is an open source implementation of the SCA in the Java language and is still available for download [10]. Another project was to develop a reference implementation of the FM3TR waveform. The FM3TR waveform is a multi-mode waveform supporting voice and data. This was developed in the mid-2000s and is available to the community for experimentation.

In 2004, member representatives from Harris Corporation and CRC gave an interoperability demonstration of two independently developed implementations of the SCA specification. This demonstration underscored the power of having a common set of interfaces and definitions across multiple radio systems.

Numerous other contributions were made through the Forum supporting the advancement of the SCA worldwide:

- API Position Paper (SDRF-03-R-0005-V1.0.0)
- Submission to JTRS JPO from SDR Forum regarding DSP and FPGA portability standardization effort (SDRF-04-R-0003-V1.0.0)
- SDRF Change Proposals and Comments on JTRS SCA 3.0 Specialized Hardware Supplement (SDRF-05-R-0001-V1.0.0)
- Comments on Software Communications Architecture Specification Version 2.2.2 (SDRF-06-R-0012-V1.0.0)
- Endorsement of JTRS SCA 2.2.2 (SDRF-08-R-0006-V1.0.0)
- Test and Certification Guide for SDRs based on SCA — Part 1: SCA (SDRF-08-P-0007-V1.0.0)

#### SCA VARIANTS

From 2000 to 2009, a number of variations of the SCA were introduced, the most significant of which follow.

## OMG SOFTWARE-BASED COMMUNICATIONS STANDARD

One of the original goals of the MSRC and government sponsors was to evolve the SCA into an industry standard rather than a military-only specification, with a long-term goal to deprecate the government SCA specification. Many organizations participating in the MSRC and JTRS Program were already members of the Object Management Group (OMG), so the U.S. government sponsored some of these organizations to facilitate the introduction, adoption, and release of the SCA as an OMG standard. However, as the SCA went through the OMG standards process, the impact of multiple organizations and individual perspectives resulted in a specification that diverged from the SCA specification. When the OMG specification was finally released in 2004 [11], the JTRS program was already heavily invested in its radio acquisition programs, so consequently, the specification was never formally adopted by the SCA community. However, a number of the concepts and interfaces of the OMG specifications later found their way into both the SCA 4.0 and Space Telecommunications Radio System (STRS) specifications.

### SPACE TELECOMMUNICATIONS RADIO SYSTEM

In the mid-2000s, NASA kicked off an effort to explore the feasibility of deploying software radio systems in space. NASA became a member of the Wireless Innovation Forum and initiated the STRS project. STRS drew on concepts and capabilities of the OMG specification and the SCA. The STRS specification [12] interface definitions and functionality show significant similarity in interface signatures and logical behaviors. STRS was implemented on radio systems provided by General Dynamics, NASA Jet Propulsion Lab, and Harris Corporation. These radio systems were integrated onto a pallet with a common avionics control system. The pallet was flown to the International Space Station on a shuttle mission and installed on the space station. It has since been in use for SDR experiments in the space environment.

### THE ESSOR ARCHITECTURE

The European Secure Software Defined Radio (ESSOR) program was launched in 2009 under the umbrella of the European Defence Agency (EDA) and sponsored by the governments of Finland, France, Italy, Poland, Spain, and Sweden. The program was awarded by the Organisation Conjointe de Coopération en Matière d'Armement (OCCAR) to the dedicated joint venture Alliance for ESSOR (a4ESSOR S.A.S.) in charge of managing the industrial consortium composed of the following respective national champions: Elektrobit, Indra, Radmor, Saab, Selex Elsag, and Thales Communications and Security. The main scope of this project was to provide the normative referential required for development and production of software radios in Europe through the use of:

- The ESSOR architecture of SDR for military purposes
- A military high data rate waveform (HDR WF) compliant with such architecture.

The ESSOR architecture is an SDR architecture relying on the already published JTRS SCA and APIs. The ESSOR architecture is a complete and consistent secure SDR architecture addressing the European military radio communications market, and fostering waveform portability among heterogeneous SDR platforms.

### THE SVFuA PROGRAM

Also in 2009, Germany launched an SDR program known as SVFuA (Armed Forces Joint Network Radio Equipment in German). The SVFuA program includes the development of a multi-channel SDR platform together with WF development and porting onto the platform using a national variant of the SCA and APIs.

### WIRELESS INNOVATION FORUM SPECIFICATIONS

The members of the Wireless Innovation Forum have also been active in the development of multiple specifications extending the SCA including the International Radio Security Services (IRSS) API, transceiver facility specification, PIM IDL specification, and lightweight and ultra-lightweight AEP specification.

## SCA NEXT AND SCA 4.1: DRIVING TOWARD MARKET HARMONIZATION

The proliferation of SCA variants required SDR manufacturers and technology vendors to fork their development efforts in supporting the international community. The result was a gradual increase in cost that was the natural result of a loss of economies of scale. Recognizing this issue, the members of the SDR Forum, now rebranded as the Wireless Innovation Forum, formed the Coordinating Committee on International SCA Standards in 2010, with a mandate to support “the harmonization of SCA based standards at the international level for the mutual benefits of all stakeholders.” Early work of this committee defined a “Coordination Model for International SCA Standards” (Document WINNF-10-R-0018) and endorsement of a three-level model for SCA standardization:

1. Areas of open standards with unlimited public access, which are managed in the best interest of all stakeholders (U.S. DoD, EDA, radio providers, waveform applications developers and others) by an independent international organization like the Wireless Innovation Forum
2. Areas of limited distribution for sensitive multi-national needs, such as for coalition interoperability, which are managed by a relevant multi-national body such as NATO or EDA
3. Areas of highly restricted access for specific national and sovereign interests, which are managed by national bodies

In parallel with this activity, the Joint Tactical Networking Center (JTNC) began work on the SCA Next initiative to develop the next major version of the SCA: SCA 4.0. Building on the lessons learned from the deployment of more than 400,000 SCA-enabled radios worldwide, the goal of this initiative was to improve performance of the SCA and expand its applicability in

*The proliferation of “SCA Variants” was requiring SDR manufacturers and technology vendors to fork their development efforts in supporting the international community. The result was a gradual increase in cost that is the natural result of a loss of economies of scale.*

*To complete the radio development framework, enabling a truly interoperable standard, it is now important that a set of APIs, on which interoperable platforms and applications can be developed and agreed upon. A number of those APIs have already been published but many are still to be defined.*

the SDR market. Key features added in this version of the specification included:

- Removing dependencies on CORBA by providing a baseline specification that was technology-neutral with appendices providing the necessary definition of the implementation technologies, thereby allowing radio and waveform providers to select the most appropriate technology for a particular radio and waveform combination
- Adding profiles, called units of functionality, allowing the SCA to be tailored for an individual radio such that only those features needed by the radio and waveforms were implemented, thus reducing overhead
- Reducing the processing overhead, and thus reducing the time required to load and start a new waveform

Elements of the ESSOR architecture and other international variants were introduced into SCA 4.0 through recommendations made by the Wireless Innovation Forum [13]. SCA 4.0 was formally released on 28 February 2012. After some experimental development by several organizations, a number of items were identified that required further refinement in order to better meet the objectives of the new specification. Thus, the SCA 4.1 initiative was kicked off with the SCA 4.1 Working Group of the Wireless Innovation Forum serving to coordinate across all the stakeholders and as a liaison with the JTNC. Finalizing SCA 4.1 is well underway, with six recommendations and two new specifications provided by the Wireless Innovation Forum to the JTNC for evaluation to be included in the revised standard. The SCA 4.1 draft was released on 26 January 2015, and the final approved specification is anticipated to be released later in the year. Although ongoing radio programs remain based on SCA 2.2.2, it is anticipated that future programs and internal projects will evolve toward the SCA 4.x version.

## WHAT'S NEXT?

The SCA has now proven its value as a framework of reference for the design and development of military tactical radios. It has been deployed and is in use on the battlefield by soldiers worldwide. With its component-based design approach, it has considerably changed the way radios are built, enabling a higher degree of design flexibility and leading to manufacturing cost reduction. From an original U.S. DoD vision of a standard military radio development architecture, the SCA, with v. 4.1, has now moved toward becoming an international specification, thanks in part to the Wireless Innovation Forum, which spearheaded a process of gathering contributions from worldwide government agencies and industry to improve on the earlier SCA 2.2.2 version.

The future of the SCA will reside in the capability of the community to expand the reach of the SCA to break the adoption barriers that new technologies face. On the technology adoption curve, within the tactical military radio market, the SCA has no doubt passed the “early adoption” stage and is ramping up into the “early majority” phase. To ensure that it continues on

this rising path, a few more things are being put in place by the community to ensure wider adoption and sustainability of the technology.

## THE DEVELOPMENT OF INTERNATIONALLY ADOPTED APIS

To complete the radio development framework enabling a truly interoperable standard, it is now important to have a set of APIs on which interoperable platforms and applications can be developed and agreed. A number of those APIs have already been published [14], but many are still to be defined, such as the transceiver and security APIs. Here again, it will be important that those APIs receive international recognition and be mandated by military acquisition programs. The Wireless Innovation Forum is playing a vital role here.

### EXPANSION INTO OTHER MILITARY MARKETS

The military tactical radio domain has been the leader in the development and adoption of this technology. By design, however, the SCA is not specific to military radio environments. Its component-based design (CBD) approach to software development can certainly be applied to many other domains, particularly those based on complex heterogeneous and distributed embedded systems that other commercial CBD approaches are not efficient at supporting. By simply modifying the domain APIs, many other military systems could benefit from the component-based standardized process brought by this technology. Electronic warfare systems, radar, avionics, signal intelligence, and armament systems are very complex systems nowadays, heavily dependent on signal processing and heterogeneous embedded computing for their operations. Advocacy to military organizations worldwide, outside of the radio domain, now needs to be done to expand the adoption of this technology. This is especially true as cognitive RF systems are beginning to be studied. Standardization on a proven technology should certainly be easier than in the early days of the U.S. JTRS radio program, when the leaders were truly visionary to bank on such a new revolutionary technology.

It will be important, especially within the military family, to ensure global adoption of this standard to avoid costly duplication of effort in parallel programs that would result in similar but non-compatible architectures.

### EXPANSION TO COMMERCIAL MARKETS

The commercial sector is also an important market where the SCA could be of great benefit, maybe not so much in the consumer domain but certainly in the industrial domain, such as backbone telecommunications systems, test and instrumentation, robotics, transportation, and more. In Japan, for example, part of the train telecommunications system is built on the SCA specification [15]. Moving the SCA to those sectors will give it additional visibility for it to become a standard for embedded system software development and deployment environments.

The larger the adoption cloud becomes, the stronger the development ecosystem will be, providing manufacturers with more choices and



even greater integrated software tools, hardware platforms and signal processing components to be used by manufacturers to facilitate the development of embedded systems, lowering development time and cost.

The Wireless Innovation Forum, as an international organization, will continue to have a significant role in bringing together the technical forces to continue to evolve the specification, including the definition of standard APIs. It will also play a vital advocacy role to promote the use of the specification into domains other than military radios.

## REFERENCES

- [1] JTRS Standards, "Software Communications Architecture Specification," v. 4.0, Feb. 2012.
- [2] J. Bard and V. Kovarik, Jr., *Software Defined Radio: The Software Communications Architecture*, Wiley, Apr. 2007.
- [3] M. R. Turner "Software Defined Radio Solutions — New Technology and JTRS Push the Envelope," SDR '02 Tech. Conf., Nov. 2002.
- [4] R. J. Lackey and D. W. Upmal, "Speakeasy: The Military Software Radio," *IEEE Commun. Mag.*, vol. 33, no. 5, May 1995, pp. 56–61.
- [5] P. G. Cook and W. Bonser, "Architectural Overview of the SPEAKeasy System," *IEEE JSAC*, vol. 17, no. 4, Apr. 1999, pp. 650–61.
- [6] <http://groups.winnforum.org/d/do/2991>
- [7] <http://groups.winnforum.org/d/do/2992>
- [8] <http://groups.winnforum.org/d/do/1725>
- [9] <http://groups.winnforum.org/SCA-History>
- [10] [http://groups.winnforum.org/reference\\_implementation](http://groups.winnforum.org/reference_implementation)
- [11] OMG, "PIM and PSM for Software Radio Components," May 2004.
- [12] NASA, Space Telecommunications Radio System, Revisions 1.02.1, STRS-AR-00002, Mar. 2010.
- [13] <http://groups.winnforum.org/Recommendations>
- [14] <http://jtnc.mil/sca>
- [15] [http://www.nec.com/en/press/201302/global\\_20130212\\_02.html](http://www.nec.com/en/press/201302/global_20130212_02.html)

## BIOGRAPHIES

CLAUDE BELISLE is chief executive officer of Nordiasoft. He graduated in 1983 from the Royal Military College of Canada in engineering physics and obtained, in 1985, an M.Sc. in physics with a specialty in optical systems. Over his 28-year career in research and development at Defence R&D Canada and Communications Research Centre, he has been involved in a wide range of technologies, from microwave-photonics, communications networks, software defined radios, and satellite communications.

VINCE KOVARIK has more than 30 years in systems and software development. He has been involved in software defined radio technology development for the last decade. As a member of the Modular Software Radio Consortium, he worked on the initial version of the Software Communications Architecture. He was the product manager for the Domain Management ToolKit (dmTK), a commercial off-the-shelf product that implemented the SCA, first available in 2000.

LEE PUCKER ([lee.pucker@wirelessinnovation.org](mailto:lee.pucker@wirelessinnovation.org)) is chief executive officer of the Wireless Innovation Forum. He has more than 25 years of experience in the development, management, marketing, and production of embedded signal processing and advanced wireless systems in organizations that include spectrum signal processing, ITT industries, and Computer Sciences Corporation. He is the named inventor on multiple patents and holds a Bachelor of Science degree in electrical engineering from the University of Illinois, a Master of Science degree from Johns Hopkins University and is Pragmatic Marketing Certified.

MARK TURNER has more than 33 years of experience in the software engineering and radio communications industries. He has been responsible for the design and engineering of multiple military radio and high-grade security products for U.S. domestic and international markets, including development and delivery of solutions for SCA-based requirements. He has been the author and presenter of more than 40 technical papers and presentations on embedded software development, software defined radio, cognitive radio, the SCA, and programmable security.

*The Wireless Innovation Forum, as an international organization, will continue to have a significant role in bringing together the technical forces to continue to evolve the specification, including the definition of standard APIs. It will also play a vital advocacy role to promote the use of the specification into domains other than military radios.*

# The ETSI Standard Architecture, Related Interfaces, and Reconfiguration Process for Reconfigurable Mobile Devices

Yong Jin, Kyunghoon Kim, Donghyun Kum, Seungwon Choi, and Vladimir Ivanov

## ABSTRACT

This article introduces the ETSI standard architecture for a reconfigurable MD, the configuration of which is determined by the downloaded radio application code. The ultimate goal of standardizing the MD architecture is to resolve the problem of portability between the radio application codes and hardware platforms of MDs. First, this article introduces key components of the standard architecture. Then the related interfaces that allow interaction of the key components with one another are shown. In addition, this article also shows how the radio application codes can be distributed from a public domain to target MDs when the codes are given in platform-specific executable code, platform-independent source code, or platform-independent intermediate representation with a consideration of both static and dynamic linking of functional blocks of each radio application code. Finally, in order to verify the feasibility of the proposed standard architecture and related interfaces, this article shows sample procedures of data transfer communication service performed in an MD adopting the standard architecture.

## INTRODUCTION

Software defined radio (SDR)-related research for both developing new technologies and standardizing the technologies has been performed for more than a decade since the early 2000s [1–5]. Software communication architecture (SCA) is one of the main results of the SDR research, which was performed under the Joint Tactical Radio System (JTRS) project [3]. Although SCA provides many advantages in designing communication systems with the functionality of reconfigurability, it is still stuck with many inherent problems, one of which lies with its computational burden, which is too heavy to be applied to commercial handsets.

The key point that this article addresses is the problem of portability between the radio application code and hardware platform of mobile

devices (MDs). In order to develop a commercially feasible technology to decouple the software (i.e., radio application code) and hardware (i.e., the platform of an MD), Working Group2 of the Technical Committee on Reconfigurable Radio System (TC RRS) of the European Telecommunications Standards Institute (ETSI) has been performing research on radio system architecture and related interfaces.

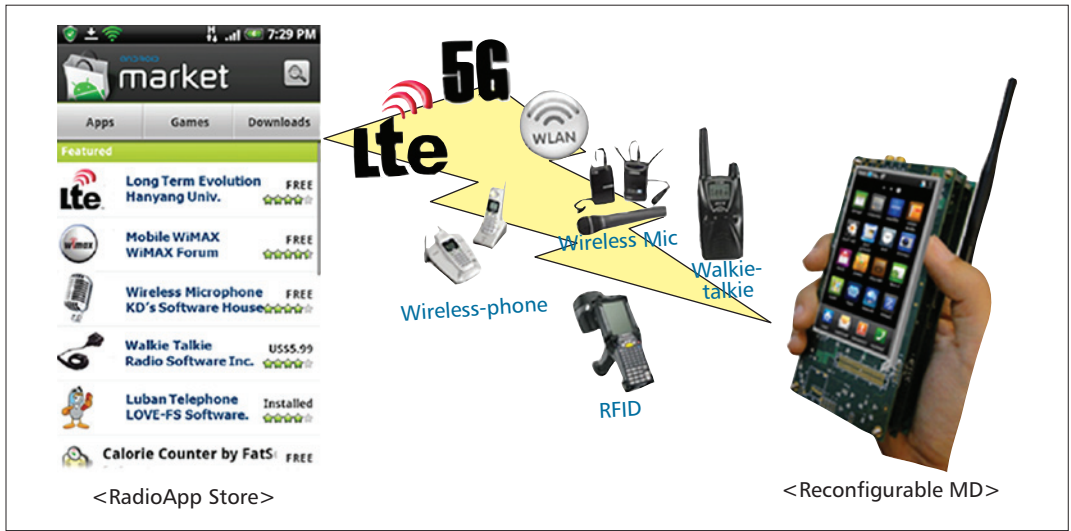
As shown in Fig. 1, the reconfigurable MD is a generic concept based on technologies such as SDR and cognitive radio (CR) with a system that exploits the capabilities of reconfigurable radio and networks for self-adaptation to a dynamically changing environment with the aim of improving the supply chain, equipment, and spectrum utilization [1]. This article proposes software architecture of a reconfigurable MD together with interfaces that interact among the components defined in the proposed architecture. This architecture has been set up as a standard architecture by TC RRS. With the standard architecture and related interfaces, a single radio application code can be ported on every different kind of hardware platform under an assumption that the hardware platform is compliant with the architecture and interfaces while the radio application code is programmed with the related interfaces. Note that the ETSI standard architecture introduced in this article does not require the operating system (OS) to be standardized.

## ETSI RRS SOLUTIONS FOR MD RECONFIGURATION

ETSI's TC RRS has been developing a standard architecture for a reconfigurable MD since 2011. The key motivation to develop the standard architecture is to resolve the problem of portability between radio application code and the hardware platform of a reconfigurable MD. Radio application code is considered to be a specific software code executable on an abstract machine of the reconfigurable MD, which provides all the necessary computing functionalities

Yong Jin, Kyunghoon Kim, Donghyun Kum, and Seungwon Choi are with Hanyang University.

Vladimir Ivanov is with LG Electronics.



**Figure 1.** Conceptual view of MD reconfiguration with a radio application downloaded from the RadioApp Store [1, Fig.1].

required in a given radio [6]. With standard interfaces defined in accordance with the standard architecture, the problem of portability can be overcome under the assumption that the hardware platform is compliant with the standard architecture and interfaces, while the radio application codes are programmed with the standard interfaces.

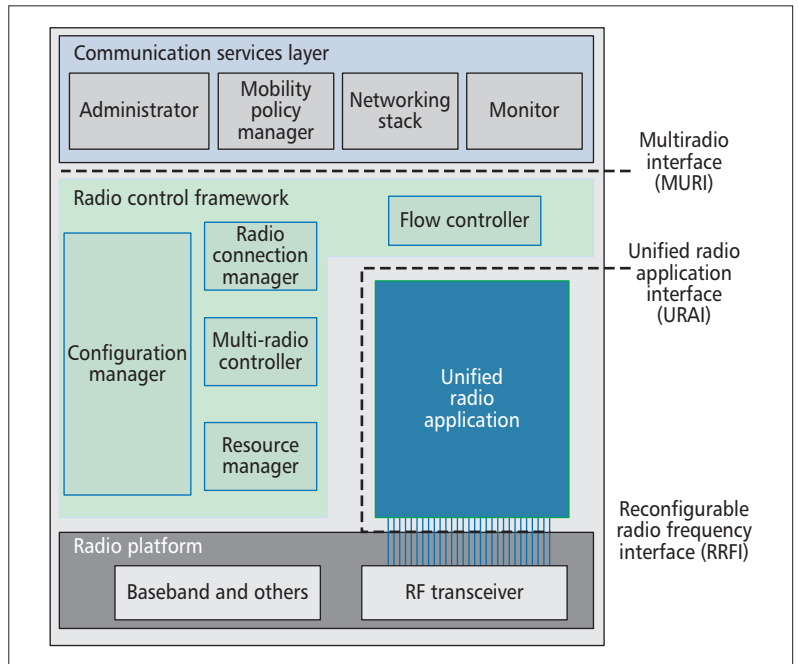
As a Technical Specification (TS) for the reconfigurable MD architecture was published in Working Group2 of TC RRS of ETSI in January 2013 [6], it is now in the procedure of becoming a corresponding European Standard (EN). With the standard architecture and related interfaces, the dependency between the software and hardware can actually vanish. In this section, the reconfigurable MD architecture and related interfaces are described.

**STANDARD ARCHITECTURE FOR MD RECONFIGURATION**

The fundamental architecture of a reconfigurable MD consists of a unified radio application (URA), communication services layer, radio control framework, and radio platform as illustrated in Fig. 2 [7].

As defined in [7], a radio application is software that enforces the generation of the transmit RF signals or the decoding of the receive RF signal. Radio application includes all the modem functionalities required in layers 1, 2, and 3. Since all radio applications exhibit common behavior from a reconfigurable MD’s perspective, those radio applications are called URAs. The URA is a key component for reconfigurability. Details about the URA are given separately in the next section.

The communication services layer supports both generic applications such as Internet access and specific applications related to multiradio applications. As shown in Fig. 2, there are four entities included in the communication services layer. Each of these four entities has different responsibilities, as follows.



**Figure 2.** Reconfigurable MD architecture components.

**Administrator:** Requests (un)installation of a URA, and creates or deletes URA instances. This typically includes the provision of information about the spectral and computational requirements for each URA, status of each URA, and so on.

**MPM:** Mobility policy manager; monitors the radio environments and MD capabilities, requests (de)activation of the URA, and provides information about the URA list. It also makes a selection among different radio access technologies (RATs), and discovers peer communication equipment and arrangement of associations.

**Networking Stack:** Sends and receives user data.

*Under the communication services layer, another important part of the reconfigurable MD is the radio control framework, which provides a generic environment for the execution of URAs, and a uniform way of accessing the functionality of the communication services layer and an individual URA.*

**Monitor:** Presents context information including received signal strength indication, packet error rate, precoding matrix indicator, and so on to a user upon the user's request and transfers context information from the URA to the proper destination entity(-ies) in MD(s) through the context information interface [8].

Under the communication services layer, another important part of the reconfigurable MD is the radio control framework, which provides a generic environment for the execution of URAs, and a uniform way of accessing the functionality of the communication services layer and an individual URA. As shown in Fig. 2, the radio control framework consists of five entities with responsibilities summarized below.

**Configuration Manager:** (Un)installs, creates/deletes instances of the URA, and manages the access to the radio parameters of the URA.

**Radio Connection Manager:** Activates/deactivates the URA according to user requests, and manages user data flows.

**Flow Controller:** Sends and receives user data packets and controls the flow of signaling packets.

**Multiradio Controller:** Schedules the requests for radio resources issued by concurrently executing URAs, and detects and manages the interoperability problems among the concurrently executed URAs.

**Resource Manager:** Manages computational resources in order to share them among simultaneously active URAs, and guarantees their real-time execution.

The radio platform typically consists of programmable hardware(s), dedicated hardware accelerator(s), RF transceiver, and antenna(s). Note that the radio platform included in the reconfigurable MD architecture shown in Fig. 2 is not a part of the ETSI standard.

#### RELATED INTERFACES FOR MD RECONFIGURATION

The above described four components, that is, URA, communication services layer, radio control framework, and radio platform, are interconnected through interfaces as follows:

- Multi-radio interface (MURI) for interconnecting the communication services layer and radio control framework
- Unified radio application interface (URAI) for interconnecting URA and radio control framework
- Reconfigurable RF interface (RRFI) for interconnecting URA and RF transceiver

MURI provides a uniform way for the reconfigurable MD to access all the URAs by interacting between the radio control framework and communication services layer as shown in Fig. 2. For that functionality, MURI supports the following three kinds of services [9]: administrative services, access control services, and data flow services.

URAI is an interface between the URA and the radio control framework, as shown in Fig. 2, to harmonize the behavior of the URA toward the operating system of a reconfigurable MD. URAI is a bidirectional service interface, where both provided and used services are visible.

URAI supports five kinds of services [10], that is, radio application management services, user data flow services, multiradio control services, resource management services, and parameter administration services.

RRFI is an interface between URA and RF transceiver as shown in Fig. 2, with which the reconfigurable MD can manage the RF transceiver regardless of URA. The key functionality of RRFI is to allocate the spectral resources simultaneously or sequentially for each of multiple URAs within the spectral boundary that is physically supported by a given MD. RRFI supports five kinds of services [11]: spectrum control services, power control services, antenna management services, transmit (Tx)/receive (Rx) chain control services, and radio virtual machine (RVM) protection services.

In addition to the three interfaces introduced in this subsection — MURI, URAI, and RRFI — in order to support the MD reconfigurability of multiple URAs, another interface regarding the independent and uniform production of radio applications as software entities should be defined as a programming interface according to the standard architecture shown in Fig. 2 [12]. This programming interface, which is denoted as radio programming interface in this article, is both a radio software development time concept as well as a runtime interface between radio software entities and the hardware platform of a reconfigurable MD. This interface, as it is a programming interface, also needs to include a uniform radio programming model that combines required runtime dynamism with real-time guarantees and efficiency. The programming model needs to be platform-independent and allow multiple radio compilers to be used for generating runtime radio packages for different platforms from the same source program. Additional aspects to be taken into account in the radio programming interface are virtualization of hardware peripherals of the reconfigurable MD such as reconfigurable RF devices. The radio programming interface is being developed in WG2 of TC RRS with special concern given to the granularity of functional blocks required for radio applications.

#### ARCHITECTURE REFERENCE MODEL FOR MULTIRADIO APPLICATIONS

In this subsection, a reference model of reconfigurable MD architecture for multiple URAs is introduced. The reference model is based on the standard architecture shown in Fig. 2.

Figure 3 exemplifies a reconfigurable MD architecture reference model for multiradio applications. Note that the standardization of TC RRS is limited to the four entities of the communication services layer and the five entities of the radio control framework, mentioned earlier, and related interfaces among those entities and the URA. As shown in Fig. 3, the reconfigurable MD architecture can be implemented with an application processor and a radio computer. In the example of Fig. 3, the red dotted part belongs to either the radio computer or the

application processor depending on the specific implementation.

In the example of Fig. 3, the operation of the application processor is performed by a given OS, which is preferably performed on non-real-time bases, whereas the radio computer's operation is performed by another OS, which should support real-time operations of the URA. The OS of the radio computer is referred to as the radio OS (ROS) in this article. Any appropriate OS empowered by a radio control framework can be an ROS [7]. In other words, the reconfigurable MD architecture is independent of OS or radio platform because the radio control framework accumulates all functionality related to MD reconfiguration.

The application processor in Fig. 3 includes the following components:

- A driver, which has the purpose of activating the hardware devices (camera, speaker, display, etc.) on a given MD
- A non-real-time OS for execution of the administrator, MPM, networking stack, and monitor, which are part of the communication services layer as previously described
- A radio controller (RC) in the radio application for sending context information to the monitor and Tx/Rx data to/from the networking stack

Note that the driver included in the application processor is to activate general hardware components such as camera, speaker, and display. In particular, this means that the driver is completely irrelevant to the MD reconfiguration.

The radio computer includes the following components:

- An ROS, a real-time OS, for executing functional blocks of the URA
- A radio platform driver, which is a hardware driver for the ROS to interact with the radio platform

From Fig. 3, it can also be observed that the five entities of the radio control framework introduced earlier in this section are classified into two groups. One group relates to real-time execution and the other group to non-real-time execution, as shown in Fig. 3. Which entities of the radio control framework relate to real-time and non-real-time execution can be determined by each vendor.

## PROPOSAL OF UNIFIED RADIO APPLICATIONS FOR MD RECONFIGURATION

As mentioned earlier, the URA is the key component for the reconfigurable MD to support reconfigurability, because the MD implements each new RAT through the procedure of download, installation, and activation of the URA without changing the hardware platform. In recent decades, the concept of MD reconfiguration through software download has been actively researched with a focus on reconfigurable architecture and framework [13, 14]. In this subsection, the distribution of radio application codes and the operational structure of the URA are explained in detail.

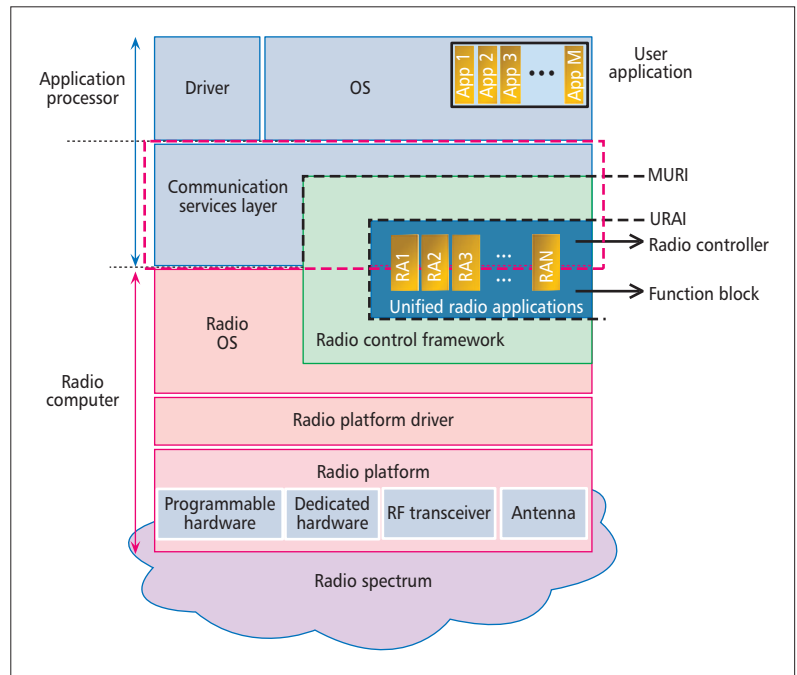


Figure 3. Reconfigurable MD architecture reference model for multiradio applications.

### DISTRIBUTION AND INSTALLATION OF RADIO APPLICATION CODES

In this subsection, the procedure of distribution and installation of radio application codes on a target reconfigurable MD is presented. During the design time, the radio application provider will generate a radio application package (RAP) that includes metadata (e.g., for pipeline configuration) and radio application codes. The radio application codes consist of two parts: one is configuration codes (Configcodes), and the other is RC codes. The former, that is, Configcodes, determines the configuration of an MD because they include functional blocks representing RAT(s) to be implemented in a given MD. The latter (i.e., RC codes), which are for transferring the context information as described in the previous section, are compiled to be executed in a given application processor before they are included in the RAP, because the context information might in general be executed in a non-real-time environment. Note that Configcodes can be distributed to the target MD in the form of either platform-specific executable code, platform-independent source code, or platform-independent intermediate representation (IR).

During the installation time, the RAP will be downloaded from a radio application store (RadioApp Store) and installed in the target reconfigurable MD. Note that the RC codes are installed in the application processor for operations that do not have to be executed in real-time processing such as context information processing, while the functional block codes are installed in the radio computer to be processed in real time. The functional blocks consist of standard functional blocks (SFBs) and user defined functional blocks (UDFBs).

When Configcodes are provided in a platform-independent IR, all the processes during installation and run-time are exactly the same as in the case of the platform-independent source code except that the radio application source codes, which include the UDFB codes, are front-end compiled during the design.

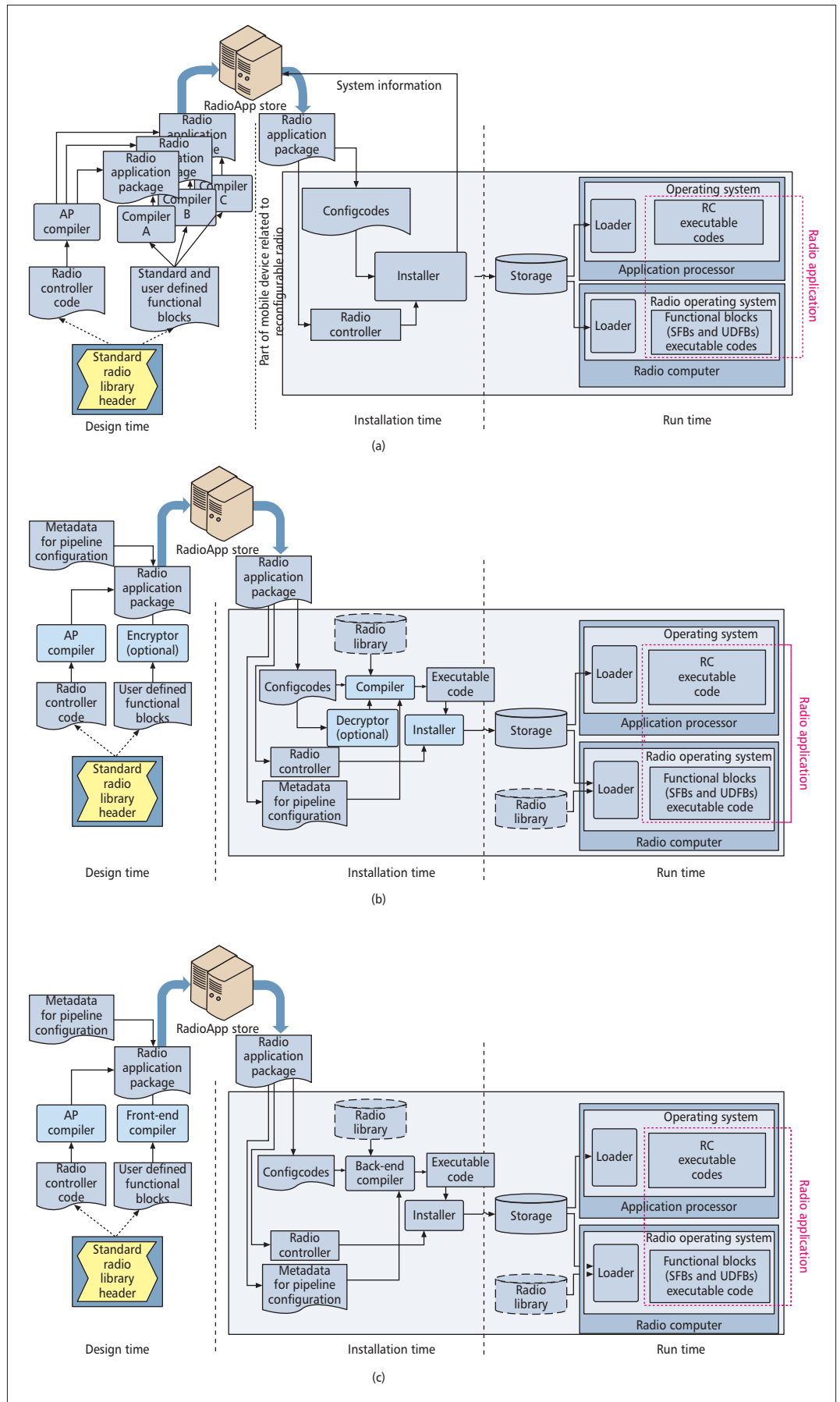


Figure 4. Conceptual diagram of distributing RAP: a) platform-specific executable code; b) platform-independent source code; c) platform-independent IR.

The former, SFB, is to be defined as a standard set of functional blocks, whereas the latter, UDFB, can be defined by radio application providers. Functional blocks that are needed in many RATs, such as forward error correction, fast Fourier transform, (de)interleaver, turbo coding, multiple input multiple output functioning, and beamforming, could be typical candidates for SFBs.

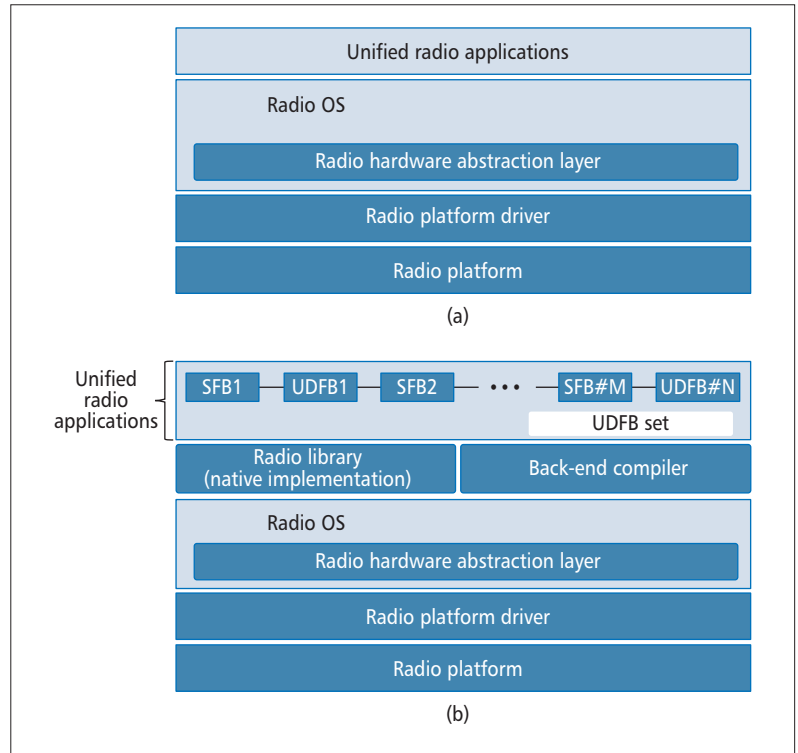
Figure 4a illustrates a block diagram corresponding to the case of distributing Configcodes that are executable in a given reconfigurable MD. When the Configcodes are executable, the functional blocks, including both SFBs and UDFBs, are executed on the radio computer of a given MD. They are compiled for each target platform during the design time to generate the corresponding Configcodes. This means that UDFB and SFB codes are compiled in accordance with a given radio computer before they are included in the RAP during the design. After compilation, the Configcodes, including both UDFB and SFB codes, are installed and loaded into a reconfigurable MD to be operated on the ROS.

Figure 4b illustrates a block diagram corresponding to the case of distributing Configcodes in the form of platform-independent source code. In the case of static linking, since the linking of UDFBs with SFBs is performed during installation time, the radio library in which all the SFBs are stored is utilized during installation time. In the case of dynamic linking, on the other hand, since the linking of UDFBs with SFBs is performed during runtime, the radio library in which all the SFBs are stored is utilized during runtime. When Configcodes are provided in a platform-independent source code, the radio application codes include the RC and UDFB codes only. As for the SFBs, the metadata provides information for efficient compilation. Only the function calls of the SFBs that are needed to execute the target URAs are contained in Configcodes. Configcodes consisting of the UDFBs are compiled (e.g., in a reconfigurable MD or in the cloud) during installation. The native implementation of SFBs is done before runtime and is contained in the radio library. As shown in Fig. 4b, the radio application source code might optionally be encrypted. If the radio application source code was originally encrypted, the corresponding Configcodes should be decrypted before the compilation during installation.

Figure 4c illustrates a block diagram corresponding to the case of distributing Configcodes in the form of a platform-independent IR. When Configcodes are provided in a platform-independent IR, all the processes during installation and runtime are exactly the same as in the case of the platform-independent source code except that the radio application source codes, which include the UDFB codes, are front-end compiled during the design.

### OPERATIONAL STRUCTURE OF URA

In the previous subsection, we have shown how the radio application codes should be processed during design and installation. In this subsection,



**Figure 5.** Two kinds of operational structure of URA: a) when URA Configcodes are executable on a target platform; b) when URA Configcodes are source codes or IR to be compiled.

the operational structure of URA in runtime is presented. The following two cases are considered:

- Configcodes are executable on a given MD.
- Configcodes are source codes or IR to be compiled at a given MD.

The first case is illustrated in Fig. 5a. Here, the SFBs and UDFBs needed to implement a given set of URAs are already bound in the executable Configcodes.

The second case is illustrated in Fig. 5b. In this case, the UDFBs needed to implement a given set of URAs are included in Configcodes and compiled in the back-end compiler. Note that the native implementation of the radio library is prepared in a given MD separately in this case because the radio library native implementation cannot be contained in Configcodes. As mentioned earlier, the function calls of SFBs are provided in the metadata. Generally, the native implementation of a radio library is provided by the radio computer vendor because the radio library includes SFBs that are implemented on the radio computer. These SFBs can be implemented without using dedicated hardware accelerator.

In the 2 cases explained above a radio hardware abstraction layer (HAL) includes hardware abstraction for SFB implementation using a hardware accelerator. This means that whenever any SFB to be implemented using hardware accelerator is called in a given radio application code, they are implemented directly on a corresponding hardware accelerator via the radio HAL.

## GENERIC PROCEDURES FOR MD RECONFIGURATION

Referring to the standard architecture of the reconfigurable MD shown in Fig. 2, several procedures for implementing the functionalities of the reconfigurable MD have been described in [7]. In this section, in order to verify the feasibility of the standard architecture together with the related interfaces, we show an operational procedure of a reconfigurable MD that adopts the standard architecture. Generic procedures introduced in this section consist of the following three steps:

- Procedures for installing URAs
- Procedures for activating installed URAs
- Procedures for creating data flow

*URA installation procedures* are shown in Fig. 6a. The procedures can be summarized as follows.

The administrator in the communication services layer sends a *DownloadRAPReq* signal including the RAP identification (ID) to RadioApp Store.

The administrator receives a *DownloadRAPCnf* signal including the RAP ID and RAP from the RadioApp Store.

The administrator sends an *InstallRReq* signal including the RAP ID to the configuration manager in the radio control framework to request URA installation.

The configuration manager first performs the radio application code certification procedure in order to verify its compatibility, authentication, and other required aspects.

The configuration manager sends an *InstallRReq* signal including the RAP ID to the file manager of the ROS to perform installation of the URA.

The file manager performs installation of the URA and transfers an *InstallRACnf* signal including the URA ID to the configuration manager, which transfers the *InstallRACnf* signal including the URA ID to the administrator.

If the downloaded radio application code is given in an IR, the configuration manager first sends a *CompileReq* signal including the RAP ID to the back-end compiler. Upon completion of

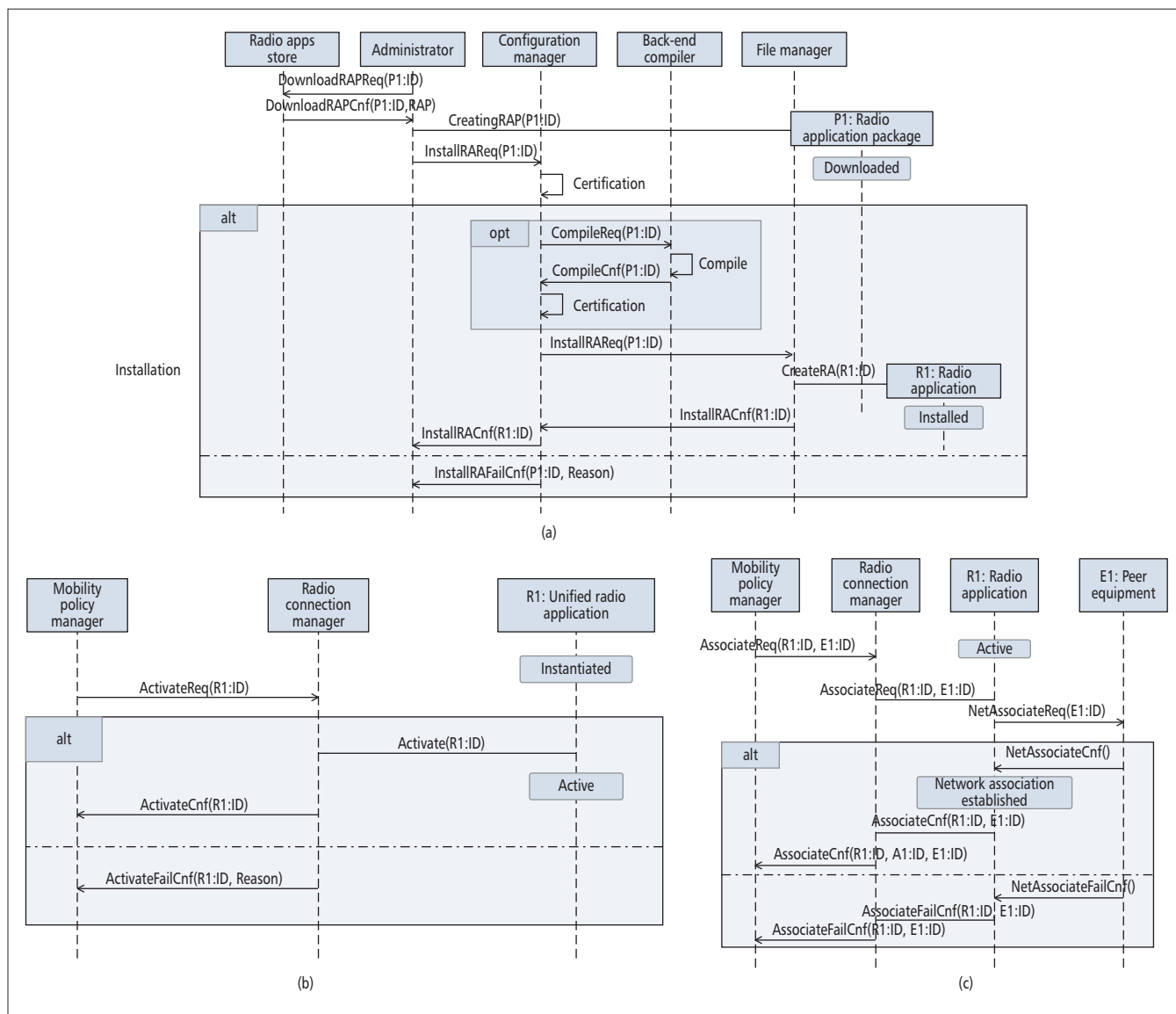


Figure 6. URA installation procedures; b) URA activation procedures; c) data flow creation procedures.



back-end compilation, the back-end compiler transfers a *CompileCnf* signal including the RAP ID to the configuration manager, which performs the certification of the back-end compiled radio application code. Only after the radio application code certification procedure is successfully completed can the URA installation take place.

In the case of installation failure, the configuration manager reports to the administrator the failure of URA installation using an *Install-RAFailCnf* signal including the RAP ID and failure reason.

*URA activation procedures* are shown in Fig. 6b. The procedure can be summarized as follows.

The MPM in the communication services layer transfers an *ActivateReq* signal including the URA ID to the radio connection manager in the radio control framework.

Upon request from the radio connection manager, the ROS activates the designated URA.

After the ROS completes the activation of the URA, the radio connection manager sends back to the MPM an *ActivateCnf* signal.

If URA activation fails, the radio connection manager reports the failure to the MPM by transferring the failed URA ID and failure reason in an *ActivateFailCnf* signal.

*Data flow creation procedures* are shown in Fig. 6c. The procedure can be summarized as follows.

The MPM in the communication services layer transfers an *AssociateReq* signal including the URA ID and peer equipment ID to the radio connection manager, where the peer equipment might be an access point of a wireless local area network, IP access node(s) (gateway general packet radio service support node, etc.) in cellular networks, or Bluetooth headsets, digital radio/television broadcasting station(s), global positioning system satellite(s), and so on.

Upon request from the radio connection manager for the ROS to create a network association, the ROS transfers the *AssociateReq* signal from the radio connection manager to the URA. Then the URA transfers the ID of corresponding peer equipment using a *NetAssociateReq* signal.

Upon completion of the network association creation, peer equipment transfers a *NetAssociateCnf* signal to the URA. Then the ROS transfers an *AssociateCnf* signal to the radio connection manager, which in turn transfers it to the MPM.

In the case of a network association failure, peer equipment transfers a *NetAssociateFailCnf* signal to the URA. Then the ROS transfers an *AssociateFailCnf* signal to the radio connection manager, which in turn transfers it to the MPM.

The three steps described above, that is, procedures for installing URAs, for activating the installed URAs, and for creating data flow, represent explicit procedures for a reconfigurable MD with the standard architecture shown in Fig. 2 to implement a communication service with a desired set of URAs. From those procedures, we have confirmed the feasibility of the standard architecture and related interfaces by observing that a desired set of radio application codes are first downloaded from RadioApp Store, then

installed and activated in the target MD; and finally, corresponding data flow is performed.

## CONCLUSION

This article introduces a standard MD architecture that has been proposed by ETSI TC-RRS for multi-radio reconfigurability. Functionalities of the key components included in the standard architecture, such as the communication services layer, radio control framework, and URA, together with the mutual interactions among those key components have been introduced in this article. This article also presents how the radio application codes can be distributed from a public domain, that is, the RadioApp Store, to the target MD. Finally, by showing procedures of downloading new radio application codes, installing and activating the downloaded codes, and performing data transfer with the desired RATs, we have verified the feasibility of the standard architecture and related interfaces. The advantages of the TC RRS standard can be summarized in the following three viewpoints. From the viewpoint of manufacturers, due to software and hardware reusability of MDs compliant with the TC-RRS standard, manufacturers can speed up new device development and decrease development cost. From the viewpoint of software developers, this standard can help them to extend business (i.e., software for reconfigurable radio). From the viewpoint of network operators, this standard can help them to improve quality of service by selection of the best communication scenario and optimize radio resource usage. Large-scale integrated projects in Europe such as SANDRA [5], WiSHFUL, and SOLDER have expressed keen interest in the standard architecture and interfaces of TC RRS for reconfigurable systems. Another important aspect of the multi-radio reconfigurability provided by the standard architecture is that a new market of radio application codes will be created, which will in turn bring about the RadioApp Store, and it will proliferate rapidly as the ETSI standard is widely deployed.

## ACKNOWLEDGMENT

This research was supported by the MSIP (Ministry of Science, ICT&Future Planning), Korea, under the ITRC (Information Technology Research Center) support program (NIPA-2014-H0301-14-1017) supervised by the NIPA (National IT Industry Promotion Agency).

## REFERENCES

- [1] M. Mueck *et al.*, "Future of Wireless Communication: Radioapps and Related Security and Radio Computer Framework," *IEEE Wireless Commun.*, vol. 19, Aug. 2012, pp. 9–16.
- [2] M. Mueck *et al.*, "ETSI Reconfigurable Radio Systems: Status and Future Directions on Software Defined Radio and Cognitive Radio Standards," *IEEE Commun. Mag.*, vol. 48, Sept. 2010, pp. 78–86.
- [3] Wireless Innovation Forum, "Software Communication Architecture 4.0," Feb. 2012.
- [4] R. C. Reinhart *et al.*, "Open Architecture Standard for NASA's Software-Defined Space Telecommunications Radio Systems," *Proc. IEEE*, vol. 95, Oct. 2007, pp. 1986–93.
- [5] Y. Cheng *et al.*, "Technology Demonstrator of a Novel Software Defined Radio-Based Aeronautical Communications System," *Science, Measurement & Technology, IET*, vol. 8, Nov. 2014, pp. 370–79.

Another important aspect of the multi-radio reconfigurability provided by the standard architecture is that a new market of radio application codes will be created, which will in turn bring about the RadioApp Store, and it will proliferate rapidly as the ETSI standard is widely deployed.

- [6] ETSI TS 103 095 V1.1.1, "Reconfigurable Radio Systems (RRS); Radio Reconfiguration Related Architecture for Mobile Devices," Jan. 2013.
- [7] ETSI EN 303 095 V1.2.1: "Reconfigurable Radio Systems (RRS); Radio Reconfiguration Related Architecture for Mobile Devices," Feb. 2015.
- [8] ETSI TR 102 944 V1.1.1, "Reconfigurable Radio Systems (RRS); Use Cases for Baseband Interfaces for Unified Radio Applications of Mobile Device," July 2011.
- [9] ETSI TS 103 146-1 V1.1.1, "Reconfigurable Radio Systems (RRS); Mobile Device Information Models and Protocols; Part 1: Multiradio Interface (MURI)," Nov. 2013.
- [10] ETSI TR 102 680 V1.1.1, "Reconfigurable Radio Systems (RRS); SDR Reference Architecture for Mobile Device," Mar. 2009.
- [11] ETSI TS 103 146-2 V1.1.1, "Reconfigurable Radio Systems (RRS); Mobile Device Information Models and Protocols; Part 2: Reconfigurable Radio Frequency Interface (RRFI)," Mar. 2015.
- [12] ETSI TR 102 839 V1.1.1, "Reconfigurable Radio Systems (RRS); Multiradio Interface for Software Defined Radio (SDR) Mobile Device Architecture and Services," Apr. 2011.
- [13] B. Bing *et al.*, "A Fast and Secure Framework for Over-the-Air Wireless Software Download Using Reconfigurable Mobile Device," *IEEE Commun. Mag.*, vol. 44, no. 6, June 2006, pp. 58–63.
- [14] T. Ulversoy, "Software Defined Radio: Challenges and Opportunities," *IEEE Commun. Surveys & Tutorials*, vol. 12, no. 4, 4th qtr. 2010, pp. 531–50.

## BIOGRAPHIES

YONG JIN received his B.S. degree from Yanbian University, Yanji, China, in 2008, and his M.S. degree from Kwandong University, Gangneung, Korea, in 2011. He is currently working toward his Ph.D. degree in the Department of Electronics and Computer Engineering at Hanyang University, Seoul, Korea. His current research focuses on signal processing techniques for software defined radio systems and cognitive radio.

KYUNGHOO KIM received his B.S. and M.S. degrees in electrical and computer engineering from Hanyang University, Seoul, Korea, in 2011 and 2013, respectively. Since 2013, he has been working toward a Ph.D. at the Communication Signal Processing (CSP) Laboratory of Hanyang Univer-

sity. His research interests include multi-antenna systems, multi-user MIMO technologies, and SDR technologies.

DONGHYUN KUM received his B.S. and M.S. from Hanyang University in electronic communication engineering in 2011 and 2013, respectively. Since 2013, he has been a Ph.D. candidate in electronic and computer engineering at Hanyang University. His research interests include MU-MIMO and SDR technology.

SEUNGWON CHOI [M91] (choi@ieee.org) received his B.S. degree from Hanyang University and his M.S. degree from Seoul National University, Korea, in 1980 and 1982, respectively, both in electronics engineering, and an M.S. degree (computer engineering) in 1985 and a Ph.D. degree (electrical engineering) in 1988, both from Syracuse University, New York. From 1988 to 1989 he was with the Department of Electrical and Computer Engineering at Syracuse University as an assistant professor. In 1989 he joined the Electronics and Telecommunications Research Institute, Daejeon, Korea. From 1990 to 1992 he was with the Communications Research Laboratory, Tokyo, Japan, as a Science and Technology Agency fellow, developing adaptive antenna array systems and adaptive equalizing filters. He joined Hanyang University in 1992 as an assistant professor. He is a professor in the School of Electrical and Computer Engineering at Hanyang University. His research interests include digital communications and adaptive signal processing with a recent focus on the implementation of smart antenna systems for both mobile communication systems and wireless data systems.

VLADIMIR IVANOV earned his M.S. degree in mathematics from Novosibirsk State University, Russia, in 1979 and his Ph.D. degree in computer science from the Military Communication Academy, St. Petersburg, Russia, in 1989. He was a senior scientist in the DSP Research Center, State University of Telecommunication in St. Petersburg. He has held leading engineering positions in high-tech companies in Israel and the United States. He joined Intel as co-director of the Communication Technology Lab in Russia in 2003. He represented Intel in the IEEE P1900.4 WG (protocols for cognitive radio) in 2008 where he was a chair of the System Architecture subgroup. In December 2014, he joined LG Electronics. His research interests include parallel and distributed computing, embedded systems, tool-aided electronic design, and mathematical foundations of electronic design.



## IEEE ICC 2016 CALL FOR PAPERS AND PROPOSALS

The 2016 IEEE International Conference on Communications (ICC) will be held from 23-27 May 2016 at Kuala Lumpur Convention Center, Malaysia, conveniently located in the middle of Southeast Asia, the region home to many of the world's largest ICT industries and research labs. Themed "Communications for All Things," this flagship conference of IEEE Communications Society will feature a comprehensive Technical Program including 13 Symposia and a number of Tutorials and Workshops. IEEE ICC 2016 will also include an attractive Industry Forum & Exhibition Program featuring keynote speakers, business and industry panels, and vendor exhibits.

### TECHNICAL SYMPOSIA

We invite you to submit original technical papers in the following areas:

#### Symposium on Selected Areas in Communications

##### - Access Systems and Networks

Ahmed E. Kamal, Iowa State University, USA

##### - Cloud Communications and Networking

Dzmitry Kliazovich, University of Luxembourg, Luxembourg

##### - Communications for the Smart Grid

Lutz Lampe, University of British Columbia, Canada

##### - Data Storage

Edward Au, Huawei Technologies, Canada

##### - E-Health

Joel Rodrigues, University of Beira Interior, Portugal

##### - Internet of Things

Antonio Skarmeta, University of Murcia, Spain

##### - Satellite and Space Communications

Song Guo, University of Aizu, Japan

##### - Social Networking

Pan Hui, HKUST, Hong Kong

#### Ad-Hoc and Sensor Networks

Abdelhakim Hafid, University of Montreal, Canada  
Cheng Li, Memorial University of Newfoundland, Canada  
Pascal Lorenz, University of Haute-Alsace, France

#### Communication and Information System Security

Kejie Lu, University of Puerto Rico, Mayaguez, Puerto Rico  
Yu Cheng, Illinois Institute of Technology, USA

#### Communications QoS, Reliability and Modelling

Kohei Shiimoto, NTT, Japan  
Christos Verikoukis, CTTC, Spain  
Charalabos Skianis, Aegean University, Greece

#### Cognitive Radio and Networks

Norman C. Beaulieu, BUPT, China  
Linyang Song, Peking University, China

#### Communications Software, Services and Multimedia Applications

Shingo Ata, Osaka City University, Japan  
Fen Hou, University of Macau, China

#### Communication Theory

Marios Kountouris, Supelec, France  
Marco Chiani, University of Bologna, Italy  
Xu (Judy) Zhu, University of Liverpool, UK

#### Green Communications Systems and Networks

Sumei Sun, Institute for Infocomm Research, Singapore  
Anura Jayasumana, Colorado State University, USA

#### Mobile and Wireless Networks

Adlen Ksentini, University of Rennes, France  
Mohammed Atiquzzaman, University of Oklahoma, USA  
Jalel Ben-Othman, University of Paris 13, France

#### Next Generation Networking and Internet

Rami Langar, University of Paris 6, France  
Shiwen Mao, Auburn University, USA  
Abdelhamid Mellouk, University of Paris-Est, France

#### Optical Networks and Systems

Walter Cerroni, University of Bologna, Italy  
Krishna Sivalingam, IIT Madras, India

#### Signal Processing for Communications

Hsiao-Chun Wu, Louisiana State University, USA  
Shaodan Ma, University of Macau, China  
Tomohiko Taniguchi, Fujitsu Labs, Japan

#### Wireless Communications

Xiaohu Ge, Huazong University of Science and Technology, China  
Dimitrie Popescu, Old Dominion University, USA  
Hossam Hassanein, Queen's University, Canada  
Rui Zhang, National University of Singapore

### INDUSTRIAL FORUM AND EXHIBITION PROGRAM

IEEE ICC 2016 will feature several prominent keynote speakers, major business and technology forums, and a large number of vendor exhibits. Submit your proposals to the IF&E Chair.  
Khaled B. Letaief (eekhaled@ee.ust.hk)

### TUTORIALS

Proposals are invited for half- or full-day tutorials in all communication and networking topics. For enquiries, please contact Tutorial Program Co-Chairs.  
Mike Devetsikiotis (mdevets@ncsu.edu)  
Koichi Asatani (asatani@ieee.org)

### WORKSHOPS

Proposals are invited for half- or full-day workshops in all communication and networking topics. For enquiries, please contact Workshop Program Co-Chairs.  
Tarek El-Bawab (telbawab@ieee.org)  
Fabrizio Granelli (granelli@disi.unitn.it)

### ORGANIZING COMMITTEE

#### General Chair

**Dato' Sri Jamaludin Ibrahim**  
CEO, Axiata Group, Malaysia

#### Executive Co-Chairs

**Hikmet Sari**  
Supelec, France  
**Borhanuddin Mohd Ali**  
Universiti Putra, Malaysia

#### Technical Program Co-Chairs

**Stefano Bregni**  
Politecnico di Milano, Italy  
**Nelson Fonseca**  
State University of Campinas, Brazil

#### Technical Program Vice-Chair

**Jiang Linda Xie**  
University of North Carolina,  
Charlotte, USA

#### Industry Forums & Exhibition Chair

**Khaled B. Letaief**  
Hong Kong University of Science  
and Technology, Hong Kong

#### Tutorial Program Co-Chairs

**Mike Devetsikiotis**  
North Carolina State University, USA  
**Koichi Asatani**  
Kogakuin University, Japan

#### Workshop Program Co-Chairs

**Tarek El-Bawab**  
Jackson State University, USA  
**Fabrizio Granelli**  
University of Trento, Italy

#### Conference Operations Chair

**Hafizal Mohamad**  
MIMOS Berhad, Malaysia

#### Advisory Executive Vice-Chair

**Datuk Hod Parman**  
Past Communication Commission  
General Director, Malaysia

#### Exhibition Chair

**Nordin Ramli**  
MIMOS Berhad, Malaysia

### IMPORTANT DATES

Paper Submissions:  
**16 October 2015**

Tutorial Proposals:  
**13 November 2015**

IF&E Proposals:  
**13 November 2015**

Workshop Proposals:  
**17 July 2015**

Paper Acceptance Notification:  
**29 January 2016**

Camera-Ready Papers:  
**29 February 2016**

# Securing Physical-Layer Communications for Cognitive Radio Networks

Yulong Zou, Jia Zhu, Liuqing Yang, Ying-Chang Liang, and Yu-Dong Yao

## ABSTRACT

This article investigates the physical-layer security of CR networks, which are vulnerable to various newly arising attacks targeting the weaknesses of CR communications and networking. We first review a range of physical-layer attacks in CR networks, including primary user emulation, sensing falsification, intelligence compromise, jamming, and eavesdropping attacks. Then we focus on the physical-layer security of CR networks against eavesdropping and examine the secrecy performance of cognitive communications in terms of secrecy outage probability. We further consider the use of relays for improving CR security against eavesdropping and propose an opportunistic relaying scheme, where a relay node that makes CR communications most resistant to eavesdropping is chosen to participate in assisting the transmission from a cognitive source to its destination. It is illustrated that the physical-layer secrecy of CR communications relying on opportunistic relaying can be significantly improved by increasing the number of relays, showing the security benefit of exploiting relay nodes. Finally, we present some open challenges in the field of relay-assisted physical-layer security for CR networks.

## INTRODUCTION

Cognitive radio (CR) [1, 2] has emerged as an intelligent radio communications system that is capable of learning its surrounding context and reconfiguring its operating parameters adapted to the time-varying environment. As an enabling technology for spectrum sharing, CR allows an unlicensed user, also called a cognitive user (CU), to sense the RF environment for detecting whether or not spectrum bands licensed to primary users (PUs) are occupied by PUs [3]. If a licensed band is detected to be unoccupied by PUs, meaning that a spectrum hole is identified, the CU changes its communications parameters for the sake of transmitting over the detected spectrum hole. Until now, extensive efforts have been devoted to the research and development

of CR spectrum sharing systems from different aspects in terms of spectrum sensing, spectrum shaping, spectrum access, and spectrum management [4, 5].

As aforementioned, the physical layer of CR networks is supposed to have the ability to sense and learn about its surrounding RF environment. This, however, is also a critical weakness that can be exploited by an adversary to launch malicious activities [6]. For example, the adversary can emit an interfering signal with the intention to modify the actual RF environment, leading legitimate CUs to be misled, compromised, and malfunctioning. Also, due to the broadcast nature of radio propagation, any network node within a CU's transmit coverage can overhear the CU's confidential communications and may illegally interpret the confidential information. Therefore, the highly dynamic and open nature of the CR physical layer makes cognitive communications extremely vulnerable to various malicious activities resulting from both the internal and external attacks.

Recently, the physical-layer security of CR networks has attracted increasing research attention [7]. Considerable studies have been conducted to protect CR communications against primary user emulation attack (PUEAs) and denial of service (DoS) attacks. Specifically, a PUEA intends to emulate a PU and transmits a radio signal with the PU's characteristics over a licensed band, leading the band to be falsely detected as occupied by the PU and denying access to legitimate CUs [8]. In contrast, a DoS attacker emits a radio signal (not necessarily with the same characteristics as the PU's signal) to interfere with the signal reception at legitimate CUs to disrupt CR communications services [9], which is also known as a jammer. It needs to be pointed out that both the PUEA and jammer transmit active signals, which may be detected by legitimate CUs so that certain prevention strategies can be adopted.

In addition to the active PUEA and jammer, cognitive transmission is also vulnerable to an eavesdropper, which is a passive attacker and becomes undetectable since the eavesdropper

Yulong Zou and Jia Zhu are with Nanjing University of Posts and Telecommunications.

Liuqing Yang is with Colorado State University.

Ying-Chang Liang is with the Institute for Information Research (I2R) and also with the University of Electronic Science and Technology of China.

Yu-Dong Yao is with Stevens Institute of Technology.

just overhears and interprets the CR transmission without transmitting any active signals. Generally, cryptographic techniques relying on secret keys are adopted to protect the transmission confidentiality against eavesdropping, which, however, introduces additional system complexity resulting from the secret key management. Moreover, the secret key distribution relies on a trusted infrastructure, which may be unavailable and even compromised in some cases. To this end, physical-layer security is now emerging as a promising paradigm by exploiting physical characteristics of wireless channels to achieve perfect secrecy against eavesdropping in an information-theoretic sense [10]. This also has great potential to address the security of CR communications against eavesdropping.

In this article, we are motivated to examine the security of physical-layer communications for CR networks. We first present an in-depth overview of CR physical-layer attacks, including the PUEA, sensing falsification, intelligence compromise, jamming, and eavesdropping attacks. Next, we examine CR physical-layer security in the face of an eavesdropper and show that increasing the transmit power is not always beneficial in terms of defending against eavesdropping. Then we propose the employment of opportunistic relaying for protecting the security of CR communications, which is shown to be an effective means, especially with an increased number of relays. Finally, we present a range of open challenging issues, followed by our concluding remarks.

## PHYSICAL-LAYER ATTACKS IN CR NETWORKS

In this section, we focus on discussing physical-layer attacks in CR networks. As shown in Fig. 1, a CR cycle comprises three typical stages: observation, reasoning, and action. Although these three cognitive stages enable a CU to learn its surrounding RF environment and adapt its transmission parameters to any changes in the environment, they are vulnerable to various attacks and introduce additional security threats. Table 1 summarizes various physical-layer attacks in the observation, reasoning, and action phases, including the PUEA, sensing falsification, intelligence compromise, jamming, and eavesdropping attacks, which are detailed in the following.

### PUEA

PUEA refers to an attacker that emulates a PU by transmitting radio signals with the same characteristics as the PU, which prevents legitimate CUs distinguishing the real PU's signal from the PUEA's faked one. In order to defend against PUEA, a so-called transmitter verification scheme was proposed in [8] by exploiting the location information to verify whether a signal is transmitted from a PU or not. It was assumed in [8] that the PU and PUEA are spatially separated, and, moreover, the PU's location is known. However, the location information of a PU may be unavailable in some cases. As a consequence, an authentication approach could be employed

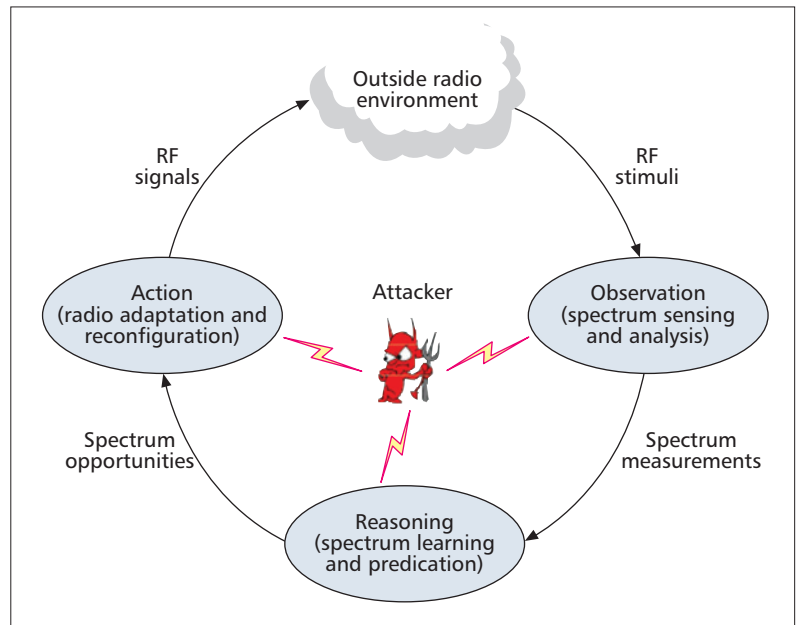


Figure 1. Illustration of a typical CR cycle.

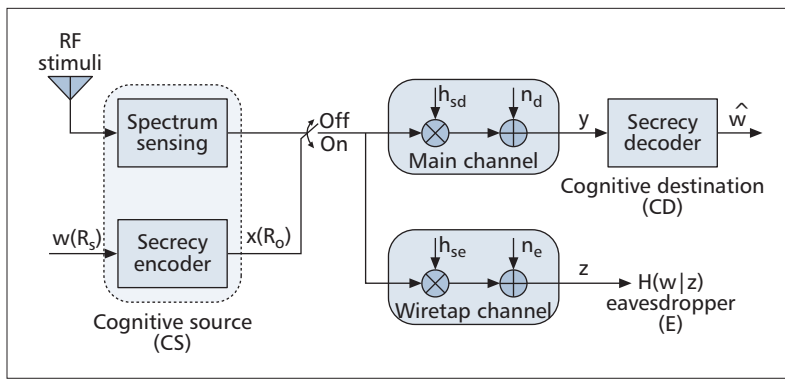
CR cycle	Attack types	Characteristics and features
Observation	PUEA	Emulating a primary user (PU) and emitting radio signals with the same characteristics as the PU
	Sensing falsification	Falsifying spectrum sensing results for the sake of intentionally misleading a cognitive user
Reasoning	Intelligence compromise	Inserting malware to maliciously alter learning and reasoning algorithms
Action	Jamming	Disrupting legitimate cognitive transmissions by emitting radio interference
	Eavesdropping	Intercepting confidential information transmissions between cognitive users

Table 1. Summarization of physical-layer attacks in different stages of the CR cycle.

to differentiate the legitimate PU from PUEA. To be specific, the legitimate PU is registered, with identity information, for example, the media access control (MAC) address, is pre-stored and authenticated. In contrast, the PUEA is typically not registered, and its identity remains unknown to legitimate users.

### SENSING FALSIFICATION

A sensing falsification attacker intends to falsify the spectrum observation and inject its fabricated results to CR networks for the sake of intentionally misleading legitimate CUs. Typically, sensing falsification attackers are sparsely distributed and are only a small fraction of the total network nodes. Thus, majority voting is an effective means to mitigate the adverse impact of fab-



**Figure 2.** A channel model for secrecy-coding-based cognitive radio communications.

ricated observation results on the spectrum sensing performance. As an alternative, a data-cleansing-based robust spectrum sensing approach was proposed in [11], where the sparsity of the falsification attack is exploited to effectively filter out the abnormal sensing data. It was shown that data-cleansing-based robust spectrum sensing significantly outperforms conventional spectrum sensing methods in terms of improving the detection probability and false alarm probability in the presence of falsified sensing data.

#### INTELLIGENCE COMPROMISE

The intelligence compromise is a legitimate CU compromised by an adversary, which maliciously inserts malware into the legitimate CU for the sake of altering its learning and reasoning algorithms, resulting in a negative impact on the node intelligence. An intelligence compromise attacker would inflict damage on spectrum learning and prediction, which may even paralyze the whole CR network. The intelligence compromise may be just a legitimate CU that is captured and enslaved by the adversary, which is thus considered as an inside attacker. Since the intelligence compromised legitimate CU infected by malware still has valid identity, it is difficult to detect and identify the presence of an intelligence compromise attacker. To this end, the automatic code patch is a promising paradigm to protect legitimate CUs against the intelligence compromise, enabling a legitimate CU to be periodically updated. If the code patch fails, it indicates that the legitimate CU may be compromised by an adversary.

#### JAMMING

A jamming attacker (also known as a jammer) attempts to emit a radio signal for interfering with the desired communications between legitimate CUs. As shown in Fig. 1, after identifying an available spectrum opportunity in the observation and reasoning stages, a legitimate CU would be scheduled to transmit its signal to its intended destination over the detected spectrum hole. Due to the broadcast nature of radio propagation, a jammer can easily disrupt the legitimate transmissions between CUs by sending a radio interference with sufficiently high power. If a jammer is present to interfere with the cognitive transmission, the received signal strength (RSS) and bit error rate (BER) experienced at

the desired destination would significantly increase, which can thus be considered as appropriate indicators for detecting the jamming attack. For example, an unusually high RSS (or an excessive BER) may indicate the presence of a jammer. Additionally, spread spectrum is considered as an effective means of defending against jamming attacks. The main spread spectrum techniques include frequency hopping spread spectrum (FHSS) and direct-sequence spread spectrum (DSSS).

#### EAVESDROPPING

An eavesdropping attacker intercepts the confidential information transmissions of legitimate CUs. The broadcast nature of wireless propagation makes cognitive transmissions vulnerable to eavesdropping attacks. When a legitimate CU transmits its data over a detected spectrum hole, any network node within the CU's transmit coverage is capable of overhearing and tapping the CU's transmission. Presently, the cryptography is adopted to protect the communications confidentiality against eavesdropping. The success of cryptography typically relies on a trusted infrastructure, which, however, may be compromised and become untrustworthy [12]. To this end, information-theoretic security emerges for CR transmissions by exploiting physical characteristics of wireless channels, referred to as physical-layer security [7], which is discussed in detail in what follows.

### PHYSICAL-LAYER SECURITY OF COGNITIVE RADIO COMMUNICATIONS

This section presents the physical-layer security of cognitive transmissions from a cognitive source (CS) to its cognitive destination (CD) in the presence of an eavesdropper. As shown in Fig. 2, a CS first performs spectrum sensing to detect whether or not a spectrum band is occupied by a primary source (PS) transmitting to its primary destination (PD). If a PS is detected to be actively transmitting, a CS is not allowed to access the spectrum band in order to avoid interfering with the reception of the PS's signal. If the PS is detected to be inactive, and thus an available spectrum hole is identified, the CS would transmit its data to the CD over the detected spectrum hole. For notational convenience, let  $P_0$  represent the probability that the spectrum band becomes unoccupied by the PS. Additionally, the probability of detection of the presence of PS is denoted by  $P_d$ , while  $P_f$  is the probability of false alarm of the presence of the PS.

Once a spectrum hole is detected, the CS switches on transmission of its confidential data to the CD, which may also be overheard by an eavesdropper (E) due to the broadcast nature of radio propagation. It is proved in [10, 15] that when the main channel (from CS to CD) has better conditions than the wiretap channel (from CS to E), physical-layer security can achieve perfect secrecy against eavesdropping. The *secrecy capacity* is shown as the difference between the

capacity of the main channel and that of the wiretap channel, which is the maximum rate at which CS can reliably and securely transmit to CD. In order to achieve the secrecy capacity, various secrecy codes (e.g., polar code and lattice code) are devised for practical wireless systems. As shown in Fig. 2, a secrecy encoder (e.g., polar code) encapsulates the CS's confidential data  $w$  (with a secrecy rate of  $R_s$ ) into an overall codeword  $x$  (with an increased rate of  $R_o$ ). The rate increase  $R_i = R_o - R_s$  represents extra redundancy, which is the cost of providing additional secrecy against eavesdropping. As shown in [12], if the rate cost  $R_i$  is higher than the capacity of the wiretap channel, perfect secrecy can be achieved, that is, the CS's data transmission is completely secure. Otherwise, the eavesdropper would succeed in intercepting the CS's transmission; a secrecy outage event happens in this case.

Next, the CS transmits its codeword  $x$  to the CD at a power of  $P_s$ , which is scaled with a wireless fading  $h_{sd}$  of the main channel and deteriorated by an additive white Gaussian noise (AWGN)  $n_d$ . Meanwhile, the codeword transmission is also overheard by E over the wiretap channel, where a wireless fading  $h_{se}$  and an AWGN  $n_e$  are encountered. Throughout this article, both the main channel and wiretap channel are independent of each other and modeled as Rayleigh fading, implying that  $|h_{sd}|^2$  and  $|h_{se}|^2$  are independent exponential random variables (RVs) with respective means of  $\sigma_{sd}^2$  and  $\sigma_{se}^2$ . Moreover, the AWGNs received at the CD and E are assumed to have zero mean and a variance of  $N_0$ . It is worth mentioning that misdetection of the presence of a PS may happen due to background noise, which would cause mutual interference between the primary and cognitive users. To limit the mutual interference level, IEEE 802.22 requires  $P_d > 0.9$  and  $P_f < 0.1$  [2], which is used throughout this article. The transmit power of the PS is represented by  $P_p$ . In addition, fading magnitudes of the wireless channels from PS to CD and E are denoted by  $|h_{pd}|^2$  and  $|h_{pe}|^2$ , respectively, which are independent exponential RVs with respective means of  $\sigma_{pd}^2$  and  $\sigma_{pe}^2$ .

In order for the CS to be able to achieve an ergodic capacity of the main channel, the codeword rate  $R_o$  is set to  $C_{sd}$ , which represents an instantaneous capacity of the CS-CD channel. Similarly, an instantaneous capacity of the wiretap channel (from CS to E) is denoted by  $C_{se}$ . As discussed above, a secrecy outage event occurs when the wiretap channel capacity becomes higher than the rate cost  $R_i$ . It needs to be pointed out that the CS starts transmitting its data only when a spectrum hole is detected. Hence, the probability of occurrence of a secrecy outage event (called secrecy outage probability) is calculated under the condition that the spectrum band is detected to be unoccupied by the PS. Hence, the secrecy outage probability of CS-CD transmissions is given by

$$P_{sout} = \Pr(C_{se} \geq R_i | \hat{H}_0) = \Pr(C_{sd} - C_{se} < R_s | \hat{H}_0),$$

where  $\hat{H}_0$  means that the spectrum band is detected idle. In Fig. 3, we show the secrecy outage

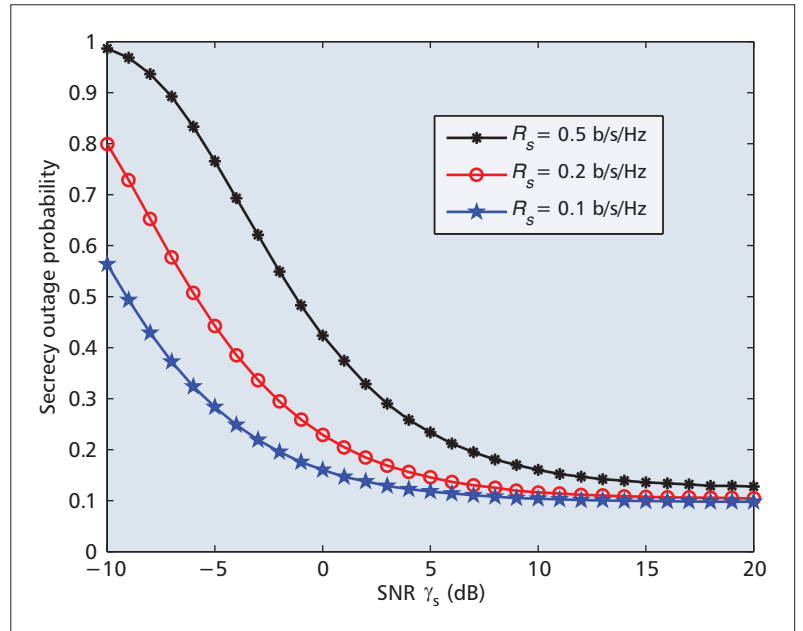
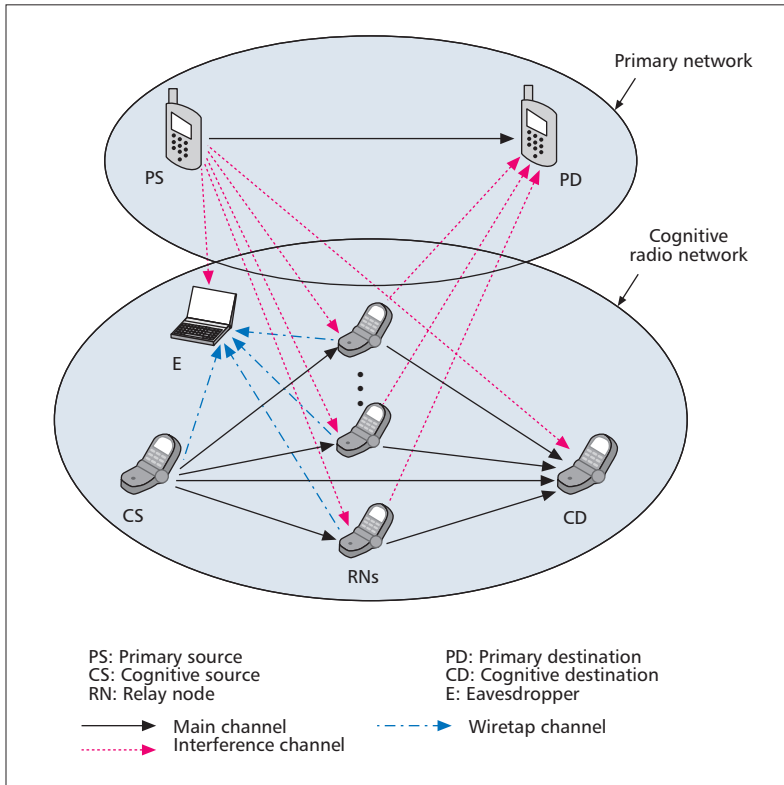


Figure 3. Secrecy outage probability vs. SNR  $\gamma_s$  for different secrecy rates.

probability vs. signal-to-noise ratio (SNR)  $\gamma_s = P_s/N_0$  of cognitive radio communications for different secrecy rates with  $P_0 = 0.8$ ,  $\gamma_p = P_p/N_0 = 5$  dB,  $\sigma_{sd}^2 = 1$ ,  $\sigma_{pd}^2 = \sigma_{pe}^2 = 0.2$ , and  $\sigma_{se}^2 = 0.1$ . It must be pointed out that the primary and secondary users are spatially separated in two different wireless networks; thus, a channel gain between two heterogeneous users from different wireless networks (e.g.,  $\sigma_{pd}^2$ ) is set to be smaller than that between two homogeneous users from the same network (e.g.,  $\sigma_{sd}^2$ ) [5], [14]. Moreover, following the physical-layer security literature [7, 10, 15], the wiretap channel is typically assumed to be a degraded version of the main channel, and thus the gain of wiretap channel  $\sigma_{se}^2$  is considered to be less than that of the main channel  $\sigma_{sd}^2$ .

As shown in Fig. 3, as the secrecy rate increases from  $R_s = 0.1$  b/s/Hz to 0.5 b/s/Hz, the secrecy outage probability of CR communications increases accordingly. This means that the physical-layer security degrades with an increased rate, showing a trade-off between security and throughput. One can also see from Fig. 3 that as the SNR  $\gamma_s$  increases, the secrecy outage probability initially decreases and finally converges to a constant value. This implies that a secrecy outage floor happens in a high SNR region, which cannot be improved by increasing the transmit power. This is because although increasing the transmit power can improve the RSS at the legitimate CD, an enhanced signal version is also received at the eavesdropper, which leads to the fact that no secrecy outage improvement is achieved with an increasing transmit power, that is, a secrecy outage floor occurs in a high SNR region. We are thus motivated to explore how the secrecy outage floor can be reduced by using, for example, opportunistic relaying, as discussed in the following section.



**Figure 4.** A cognitive relay network consists of one CS, one CD and  $N$  RNs in the presence of an E.

## OPPORTUNISTIC RELAYING FOR ENHANCING PHYSICAL-LAYER SECURITY

In this section, we examine the employment of opportunistic relaying for the enhancement of physical-layer security in CR networks. As shown in Fig. 4,  $N$  relay nodes (RNs) are assumed to be available for assisting the transmission from CS to CD, where the amplify-and-forward (AF) protocol is considered when RNs retransmit the CS's data to the CD. To be specific, when a spectrum hole is detected, a CS first transmits its signal  $x$  to a CD, which can be overheard by E and  $N$  RNs. In the opportunistic relaying, only a single RN will be chosen among the  $N$  RNs to forward an amplified version of its received signal using a scaling factor (without any sort of decoding), which is also overheard by E for interception purposes. In this way, both the CD and E can receive two copies of the CS's signal, which are transmitted from the CS and the selected RN, respectively. For simplicity, the selection diversity combining (SDC) method is considered for both the CD and E, meaning that a received signal with higher SNR is adopted for decoding the CS' signal.

Given  $N$  RNs available in the CR networks shown in Fig. 4, opportunistic relaying chooses the "best" RN to participate in forwarding the CS's transmission to the CD, aiming to maximize the cognitive physical-layer security against eavesdropping. Without loss of generality, we consider that  $RN_i$  is selected among  $N$  RNs, which first performs a coherent reception of the

CS's signal and then forwards its received signal with a scaling factor for normalization. Due to the broadcast nature of radio propagation, both CD and E can receive the  $RN_i$ 's signal retransmission and the corresponding signal-to-interference-and-noise ratio (SINR) at CD given by

$$\text{SINR}_d^i = \frac{|h_{si}|^2 |h_{id}|^2 \gamma_s}{|h_{id}|^2 (|h_{pi}|^2 \alpha \gamma_p + 1) + |h_{si}|^2 (|h_{pd}|^2 \alpha \gamma_p + 1)}, \quad (2)$$

where  $h_{si}$ ,  $h_{id}$ ,  $h_{pi}$ , and  $h_{pd}$  represent the CS- $RN_i$ ,  $RN_i$ -CD, PS- $RN_i$ , and PS-CD channels, respectively. Moreover, the parameter  $\alpha$  is given by 0 when the spectrum band is idle (i.e., no primary signal is transmitted from the PS). By contrast, if the band is currently occupied by the PS,  $\alpha$  is set to 1. Meanwhile, the SINR received at E, denoted by  $\text{SINR}_e^i$ , can be obtained similarly by replacing  $h_{id}$  and  $h_{pd}$  in Eq. 2 with  $h_{ie}$  and  $h_{pe}$ , which represent the  $RN_i$ -E and PS-E channels, respectively. In practice, obtaining the eavesdropper's channel state information (CSI) is impossible, since E is passive and typically keeps silent in CR networks. Motivated by this observation, an RN that maximizes the CD's received SINR (i.e.,  $\text{SINR}_d^i$ ) is generally selected to forward its received signal, yielding the best RN selection criterion as

$$\text{Best RN} = \arg \max_{i \in \mathcal{R}} \text{SINR}_d^i, \quad (3)$$

where  $\mathcal{R}$  denotes the set of  $N$  RNs and  $\text{SINR}_d^i$  is given by Eq. 2. It can be observed from Eq. 3 that the CSIs of the CS- $RN_i$ ,  $RN_i$ -CD, PS- $RN_i$ , and PS-CD channels are required to carry out the relay selection without needing the eavesdropper's CSI knowledge. Moreover, when  $\alpha$  is set to 0, the relay selection criterion as given by Eq. 3 degrades to the conventional so-called harmonic mean selection [13]. This is because  $\alpha = 0$  implies no mutual interference occurring between the primary and secondary users; thus, the cognitive transmission in this case becomes the same as in the conventional wireless communications scenario. From Eq. 3, the capacity achieved at CD, denoted by  $C_d$ , can be determined by using the SDC to combine the two received signals from the "best" RN and CS, respectively. Also, the wiretap channel capacity achieved at E, denoted by  $C_e$ , can be similarly obtained. Like Eq. 1, the secrecy outage probability of the opportunistic relaying scheme can be obtained by calculating the probability that the difference between  $C_d$  and  $C_e$  falls below the secrecy rate  $R_s$ . Additionally, all the CS-CD, CS- $RN_i$ ,  $RN_i$ -CD, PS-CD, PS- $RN_i$ , PS-E, CS-E, and  $RN_i$ -E channels are modeled as independent Rayleigh fading with respective variances of  $\sigma_{sd}^2$ ,  $\sigma_{si}^2$ ,  $\sigma_{id}^2$ ,  $\sigma_{pd}^2$ ,  $\sigma_{pi}^2$ ,  $\sigma_{pe}^2$ ,  $\sigma_{se}^2$ , and  $\sigma_{ie}^2$ .

In Fig. 5, we show the secrecy outage probability vs. SNR  $\gamma_s$  of the direct transmission (i.e., CS directly transmits to CD without using RNs) and the opportunistic relaying for different number of RNs  $N$  with  $R_s = 0.1$  b/s/Hz,  $P_0 = 0.8$ ,  $\gamma_p = 5$  dB,  $\sigma_{sd}^2 = \sigma_{si}^2 = \sigma_{id}^2 = 1$ ,  $\sigma_{pd}^2 = \sigma_{pi}^2 = \sigma_{pe}^2 = 0.2$ , and  $\sigma_{se}^2 = \sigma_{ie}^2 = 0.1$ . As shown in Fig. 5, for all cases



of  $N = 2, 4,$  and  $6,$  the secrecy outage probability of the opportunistic relaying is even worse than that of the direct transmission in a low SNR region (e.g.,  $\gamma_s < -6$  dB. This is because in the opportunistic relaying scheme, one half of a time slot is wasted by the chosen “best” RN to retransmit the CS’ signal to the CD, resulting in a certain loss of secrecy capacity. It is pointed out that although the SDC method is considered at the CD for combining its received signals from the CS and the “best” RN, the capacity of the CS-CD channel (i.e.,  $C_{sd}$ ) is also scaled by half in the opportunistic relaying scheme, since the CS transmits only in the first half time slot and remains silent in the second half slot, which is occupied by the “best” RN to retransmit the CS’s signal. One can observe from Fig. 5 that as the SNR continues increasing, opportunistic relaying becomes better than direct transmission in terms of the secrecy outage probability, showing the performance benefit achieved by the proposed opportunistic relaying.

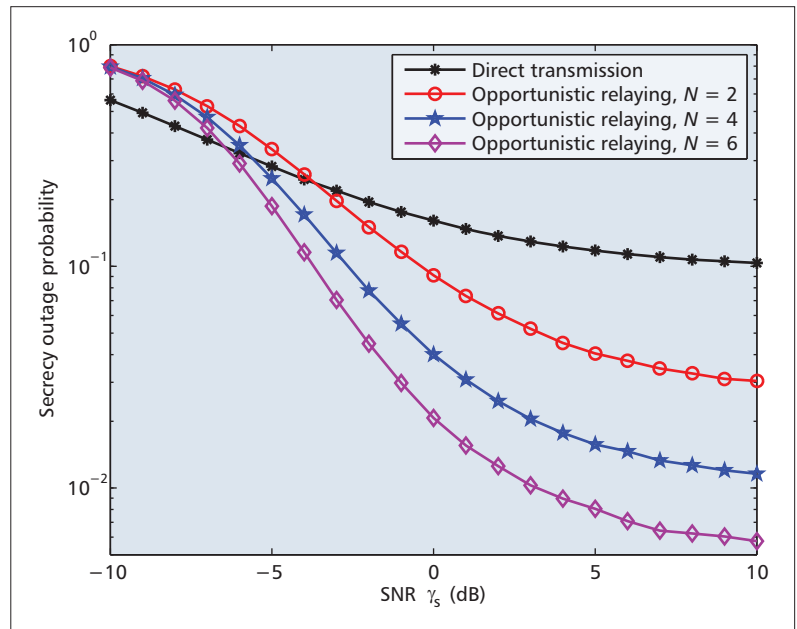
Figure 5 also shows that with a sufficiently high SNR, the direct transmission and opportunistic relaying schemes converge to their respective secrecy outage floors. Moreover, the secrecy outage floor of the opportunistic relaying is lower than that of the direct transmission. As shown in Fig. 5, as the number of RNs increases from  $N = 2$  to  $6,$  the secrecy outage floor of the opportunistic relaying is significantly reduced, showing the physical-layer security advantage of exploiting RNs. This is due to the fact that with an increasing number of RNs, it is more likely to choose an RN that can succeed in defending against eavesdropping, thus leading to a reduced secrecy outage floor. Although the opportunistic relaying scheme can effectively protect the wireless transmissions against eavesdropping, it introduces additional system complexity due to the distributed relay management and synchronization. To be specific, multiple RNs are distributed spatially in CR networks, which need to be effectively managed and synchronized for the sake of performing opportunistic relay selection. Additionally, in the opportunistic relaying scheme, the CD needs to combine its received signals from the “best” RN and CS, which comes at the cost of extra computational complexity for signal combining.

## OPEN CHALLENGES AND FUTURE WORK

This section presents some future directions in the research field of cognitive relay security. Although opportunistic relaying is shown to enhance cognitive communications security, many challenging issues remain open at the time of writing.

### JOINT RELAY-AND-JAMMER SELECTION

When a CS transmits its signal to a CD in the presence of an eavesdropper, a partner node can either be employed as a relay to assist the CS’ transmission for enhancing the signal quality received at the CD or act as a jammer to emit artificial noise to contaminate the eavesdropper’s signal reception. It is unclear whether it is beneficial to employ the node as a relay (or jammer) in terms of defending CR communications against



**Figure 5.** Secrecy outage probability comparison between the direct transmission and opportunistic relaying for different numbers of RNs  $N.$

eavesdropping. Additionally, given multiple available partner nodes, some nodes may be selected to assist the CS-CD transmission, while others may be used as jammers to generate artificial noise to interfere with the eavesdropper. This is called joint relay-and-jammer selection, which can be considered as a means of improving cognitive communications security against eavesdropping. Although there are some existing efforts devoted to joint relay-and-jammer selection, they are limited to single-relay and single-jammer selection in non-CR networks. It is of interest to explore a more general framework of multi-relay and multi-jammer selection in CR networks.

### UNTRUSTED RELAY DETECTION AND PREVENTION

As discussed above, the physical-layer security of CR communications is significantly improved by using opportunistic relaying in terms of secrecy outage probability. Although the employment of relays is capable of enhancing the security of cognitive communications against eavesdropping, the relays by themselves may not be trusted and attempt to tap CR communications. For example, if a relay is captured and compromised by an adversary, it becomes untrusted and launches malicious activities (e.g., eavesdropping) in CR networks. The secrecy performance of CR communications in the face of untrusted relays remains uncertain, which may be considered for future work. Also, it is of high importance to explore the detection and prevention of untrusted relays in CR networks.

### FIELD EXPERIMENT FOR OPPORTUNISTIC RELAYING

IEEE 802.22 is the first worldwide standard designed for a CR-based wireless regional area network (WRAN), which enables unlicensed

Although opportunistic relaying is shown to enhance the security of cognitive communications in terms of secrecy outage probability, its security benefit is only proved theoretically based on some simplified assumptions (e.g. perfect CSI knowledge is assumed).

devices to operate in white spaces of the TV broadcast spectrum without causing harmful interference to incumbent users including TV users and wireless microphones. It is necessary to carry out field experiments for testing the effectiveness of opportunistic relaying in real IEEE 802.22 WRANs in the presence of various attacks. Although opportunistic relaying is shown to enhance the security of cognitive communications in terms of secrecy outage probability, its security benefit is only proved theoretically based on some simplified assumptions (e.g., perfect CSI knowledge is assumed). It will be of great interest to investigate whether opportunistic relaying is still effective in real WRAN environments in terms of defending against CR attacks.

## CONCLUSION

In this article, we first present a comprehensive review on physical-layer attacks in CR networks, including the PUEA, sensing falsification, intelligence compromise, jamming, and eavesdropping attacks. The physical-layer security of CR communications in the presence of an eavesdropper is then examined in terms of secrecy outage probability. It is shown that as the transmit power increases, the secrecy outage probability of cognitive communications initially decreases and finally converges to a fixed value, showing that a secrecy outage floor occurs in high SNR regions. In order to improve the physical-layer security of cognitive communications, we consider the use of relays to assist the cognitive communications and propose an opportunistic relaying scheme. Numerical results show that upon increasing the number of relays, opportunistic relaying can significantly reduce the secrecy outage floor of cognitive communications. Additionally, we point out some open challenges in the research field of exploiting relays for the physical-layer security of CR networks.

## ACKNOWLEDGMENT

This work was supported by the National Natural Science Foundation of China (Grant Nos. 61401223 and 61522109), the Natural Science Foundation of Jiangsu Province (Grant No. BK20140887), and the Key Project of Natural Science Research of Higher Education Institutions of Jiangsu Province (No. 15KJA510003).

## REFERENCES

- [1] J. Mitola and G. Q. Maguire, "Cognitive Radio: Making Software Radios More Personal," *IEEE Personal Commun.*, vol. 6, no. 4, Aug. 1999, pp. 13–18.
- [2] Y. Zou, Y.-D. Yao, and B. Zheng, "Diversity-Multiplexing Trade-Off in Selective Cooperation for Cognitive Radio," *IEEE Trans. Commun.*, vol. 60, no. 9, Sept. 2012, pp. 2467–81.
- [3] Y.-C. Liang et al., "Cognitive Radio Networking and Communications: An Overview," *IEEE Trans. Vehic. Tech.*, vol. 60, no. 7, Sept. 2011, pp. 3386–3407.

- [4] S. Haykin, D. J. Thomson, and J. H. Reed, "Spectrum Sensing for Cognitive Radio," *Proc. IEEE*, vol. 97, no. 5, May 2009, pp. 849–77.
- [5] Y. Zou, Y.-D. Yao, and B. Zheng, "Cooperative Relay Techniques for Cognitive Radio Systems: Spectrum Sensing and Secondary User Transmissions," *IEEE Commun. Mag.*, vol. 50, no. 4, Apr. 2012, pp. 98–103.
- [6] G. Baldini et al., "Security Aspects in Software Defined Radio and Cognitive Radio Networks: A Survey and a Way Ahead," *IEEE Commun. Surveys & Tutorials*, vol. 14, no. 2, May 2012, pp. 355–79.
- [7] Y. Zou, X. Li, and Y.-C. Liang, "Secrecy Outage and Diversity Analysis of Cognitive Radio Systems," *IEEE JSAC*, vol. 32, no. 11, Nov. 2014, pp. 2222–36.
- [8] R. Chen and J.-M. Park, "Ensuring Trustworthy Spectrum Sensing in Cognitive Radio Networks," *Proc. 2006 IEEE Workshop Net. Tech. Software Defined Radio Net.*, Reston, VA, Sept. 2006, pp. 110–19.
- [9] Y. Tan, S. Sengupta, and K.P. Subbalakshmi, "Analysis of Coordinated Denial-Of-Service Attacks In IEEE 802.22 Networks," *IEEE JSAC*, vol. 29, no. 42, Apr. 2011, pp. 890–90.
- [10] F. Oggier and B. Hassibi, "The Secrecy Capacity of the MIMO Wiretap Channel," *IEEE Trans. Info. Theory*, vol. 57, no. 8, Oct. 2007, pp. 4961–72.
- [11] G. Ding et al., "Robust Spectrum Sensing with Crowd Sensors," *IEEE Trans. Commun.*, vol. 62, no. 9, Sept. 2014, pp. 3129–43.
- [12] X. Tang et al., "On the Throughput of Secure Hybrid-ARQ Protocols for Gaussian Block Fading Channels," *IEEE Trans. Info. Theory*, vol. 55, no. 4, Apr. 2009, pp. 1575–91.
- [13] E. Ardestanizadeh et al., "Wiretap Channel with Secure Rate-Limited Feedback," *IEEE Trans. Info. Theory*, vol. 55, no. 12, Sept. 2009, pp. 5353–61.
- [14] T. Yucek and H. Arslan, "A Survey of Spectrum Sensing Algorithms for Cognitive Radio Applications," *IEEE Commun. Surv. & Tut.*, vol. 11, no. 1, Mar. 2009, pp. 116–30.
- [15] A. Bletsas et al., "A Simple Cooperative Diversity Method Based on Network Path Selection," *IEEE JSAC*, vol. 24, no. 3, Mar. 2006, pp. 659–72.

## BIOGRAPHIES

YULONG ZOU [SM'13] (yulong.zou@njupt.edu.cn) is a professor at Nanjing University of Posts and Telecommunications (NUPT), China. He received his B.Eng. degree in information engineering from NUPT in July 2006, his first Ph.D. degree in electrical engineering from Stevens Institute of Technology, New Jersey, in May 2012, and his second Ph.D. degree in signal and information processing from NUPT in July 2012.

JIA ZHU (jjazhu@njupt.edu.cn) is an associate professor at NUPT. She received her B.Eng. degree in computer science and technology from Hohai University, Nanjing, China, in July 2005, and her Ph.D. degree in signal and information processing from NUPT in April 2010.

LIUQING YANG [F'15] (lqyang@engr.colostate.edu) is a professor in the Electrical and Computer Engineering Department, Colorado State University. She received her B.Eng. degree in electrical power engineering from Huazhong University of Science and Technology, China, in 1994, and her M.S. and Ph.D. degrees in electrical and computer engineering from the University of Minnesota in 2002 and 2004, respectively.

YING-CHANG LIANG [F'11] (ycliang@i2r.a-star.edu.sg) is a principal scientist in the Institute for Infocomm Research (I2R), Agency for Science, Technology & Research (A\*STAR), Singapore. He was an adjunct staff member of the National University of Singapore and Nanyang Technological University from 2004 – 2009.

YU-DONG YAO [F'11] (yyao@stevens.edu) is a professor at the Electrical and Computer Engineering Department, Stevens Institute of Technology. He received his B.Eng. and M.Eng. degrees from NUPT in 1982 and 1985, respectively, and his Ph.D. degree from Southeast University, Nanjing, China, in 1988, all in electrical engineering.

Now...

# 2 Ways to Access the IEEE Member Digital Library

With two great options designed to meet the needs—and budget—of every member, the IEEE Member Digital Library provides full-text access to any IEEE journal article or conference paper in the IEEE *Xplore*® digital library.

Simply choose the subscription that's right for you:

## IEEE Member Digital Library

Designed for the power researcher who needs a more robust plan. Access all the IEEE content you need to explore ideas and develop better technology.

- 25 article downloads every month

## IEEE Member Digital Library Basic

Created for members who want to stay up-to-date with current research. Access IEEE content and rollover unused downloads for 12 months.

- 3 new article downloads every month

Get the latest technology research.

**Try the IEEE Member Digital Library—FREE!**

[www.ieee.org/go/trymdl](http://www.ieee.org/go/trymdl)



IEEE Member Digital Library is an exclusive subscription available only to active IEEE members.

# Prototyping Real-Time Full Duplex Radios

MinKeun Chung, Min Soo Sim, Jaeweon Kim, Dong Ku Kim, and Chan-Byoung Chae

## ABSTRACT

In this article, we present a real-time full duplex radio system for 5G wireless networks. Full duplex radios are capable of opening new possibilities in contexts of high traffic demand where there are limited radio resources. A critical issue, however, in implementing full duplex radios in real wireless environments is being able to cancel self-interference. To overcome the self-interference challenge, we prototype our design on a software-defined radio platform. This design combines a dual-polarization antenna-based analog part with a digital self-interference canceler that operates in real time. Prototype test results confirm that the proposed full duplex system achieves about 1.9 times higher throughput than a half duplex system. This article concludes with a discussion of implementation challenges that remain for researchers seeking the most viable solution for full duplex communications.

## INTRODUCTION

### NEW BREAKTHROUGH: FULL DUPLEX RADIOS

How much does it cost to purchase wireless spectrum? In a wireless spectrum auction in January 2015, the Federal Communications Commission (FCC) raised, for 65 MHz bandwidth, a record-breaking \$44.9 billion. This illustrates how valuable the wireless spectrum has become; it offers more high-speed connectivity and satisfies more user demand for data within a limited wireless spectrum. The FCC is considering releasing more spectrum for wireless broadband usage. For the endless surge in wireless data traffic, however, this cannot be the ultimate solution.

Mobile devices with advanced wireless network capabilities, such as smartphones and tablets, are becoming ubiquitous, and keeping pace with their growth is the ever increasing demand for bandwidth. Global mobile data traffic will increase nearly 10-fold between 2014 and 2019. In that time, mobile data traffic is expected to grow at a compounded annual growth rate of 57 percent, reaching 24.3 exabytes per month by 2019 [1]. These trends could create a *spectrum crunch* as the frequencies used to carry this traffic are exhausted.

Although the laws of physics prohibit the pro-

duction of more spectrum, there is a lot of potential for *aggressive expansion* in scarce resource, that is, boosting spectral efficiency using novel technologies. A candidate for creating a new breakthrough to alleviate the spectrum crunch is full duplex. It theoretically doubles spectral efficiency, making it worth billions of dollars. Full duplex thus holds tremendous potential to carry out the solutions needed in the future evolution of wireless systems.

### KEY CHALLENGE: SELF-INTERFERENCE

Since Guglielmo Marconi developed the wireless telegraph in 1895, the bane of wireless networks has been self-interference. It is the presence of self-interference that represents the key challenge to implementing full duplex wireless systems. Self-interference is the phenomenon where, through the coupling of transceivers in a wireless network, a signal is transmitted from a transmitter to its own receiver while that receiver is attempting to receive a signal sent by another device. It compels the fundamental assumption that a wireless network has to be operated in half duplex mode on the same channel. For example, Long Term Evolution (LTE) frequency-division duplex (FDD) today is operated so that downlink and uplink transmission take place in two different frequency bands. In other words, the existence of self-interference cuts in half the amount of resources available, such as time and frequency, for wireless communications. For this reason, it is essential to manage self-interference to achieve the highest throughput performance with limited radio resources.

### THE BEGINNING OF AGGRESSIVE EXPANSION: SDR PLATFORM-BASED PROTOTYPING

Up to this point, researchers have mostly depended on software simulations to test their theories that exploit simplified channel models, for example, additive white Gaussian noise (AWGN) and Rayleigh fading. In real-world wireless systems, however, impairments occur that are often overlooked in simulations, such as amplifier nonlinearity, gain/phase offset, I/Q imbalance, quantization effects, and timing jitter. Such impairments make prototyping imperative if the feasibility and commercial viability of any

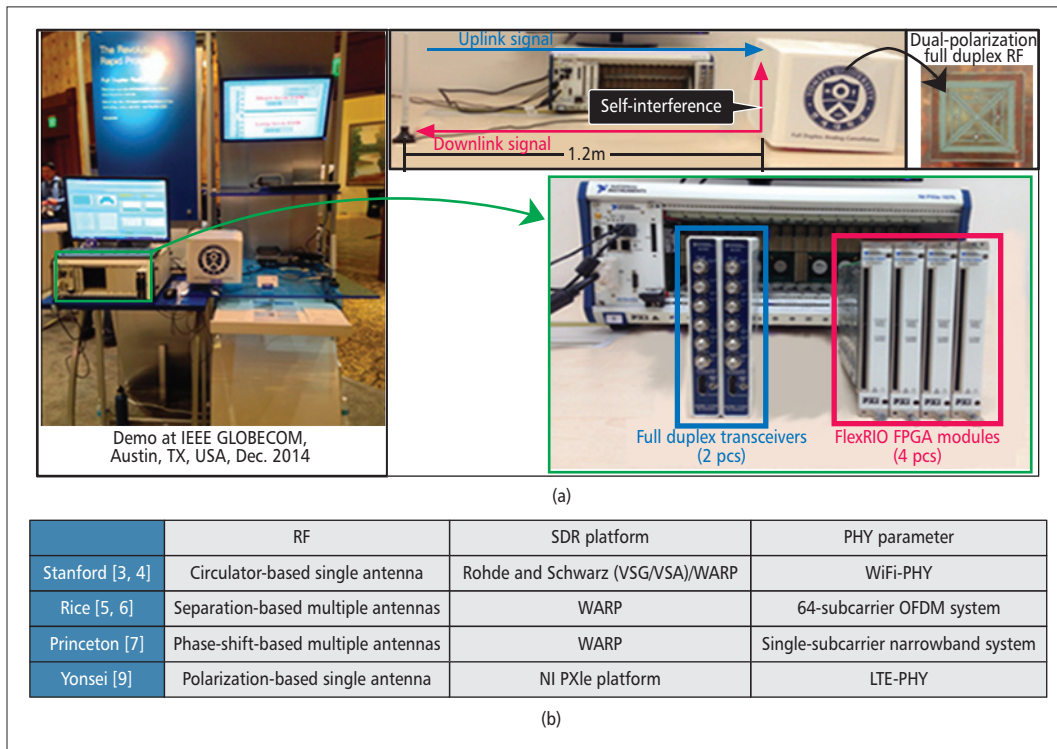
MinKeun Chung, Min Soo Sim, Dong Ku Kim, and Chan-Byoung Chae are with Yonsei University.

C.-B. Chae and D. Kim are co-corresponding authors.

Jaeweon Kim is with National Instruments.

This work was in part supported by the the MISP under the “IT Consilience Creative Program” (IITP- 2015-R0346-15-1008) and by the ICT R&D Program of MSIP/IITP [B0126-15-1012].

*It is the presence of self-interference that represents the key challenge to implementing full duplex wireless systems. Self-interference is the phenomenon where, through the coupling of transceivers in a wireless network, a signal is transmitted from a transmitter to its own receiver while that receiver is attempting to receive a signal sent by another device.*



**Figure 1.** a) Real-time full duplex radio experiment setup in our laboratory/the exhibition hall at IEEE GLOBECOM in Austin, Texas, December 2014; b) the categorized comparison of implementation characteristics by each research group. In the SDR platform, VSG and VSA denote vector signal generator and analyzer, respectively.

new wireless standard or technology are to be validated.

For next generation wireless research, a viable prototyping option has emerged known as software-defined radio (SDR) [2]. SDR enables researchers to rapidly prototype a system. Researchers at Stanford [3, 4], Rice [5, 6], and Princeton [7] have implemented various testbeds to build in-band full duplex radios using combined RF antennas and SDR platform [8]. As shown in Fig. 1a, a real-time full duplex LTE system was also demonstrated at IEEE GLOBECOM in Austin, Texas in December 2014.<sup>1</sup> The categorized comparison of implementation characteristics by each research group is summarized in Fig. 1b. The two sections that follow elaborate on how to solve the key challenge and implement real-time full duplex radios.

## PROTOTYPE SETTINGS: SYSTEM SPECIFICATIONS AND HARDWARE ARCHITECTURE

The demonstrated full duplex prototype [9] is based on the LTE downlink standard [10] with the following system specifications: a transmission bandwidth of 20 MHz, 30.72 MHz sampling rate, 15 kHz subcarrier spacing, 2048 fast Fourier transform (FFT) size, and variable 4/16/64 quadrature amplitude modulation (QAM). The prototype is implemented, as shown in Fig. 1a, using LabVIEW system design software and the state-of-the-art PXIe SDR platform, where two

full duplex nodes consist of the following four main components.

**Dual-Polarization Full Duplex RF Antenna:** This is a dual-polarization slot antenna with high cross-polarization discrimination (XPD) in all directions [11].

**PXIe-8133:** This is a real-time (RT) controller equipped with a 1.73 GHz quad-core Intel Core i7-820 processor and 8 GB of dual-channel 1333 MHz DDR3 random access memory (RAM) [12].

**NI 5791R:** A 100 MHz bandwidth baseband transceiver module is used equipped with a dual 130 MS/s analog-to-digital converter (ADC) with 14-bit accuracy and a dual 130 MS/s digital-to-analog converter (DAC) with 16-bit accuracy [13].

**PXIe-7965R:** This is a field-programmable gate array (FPGA) module equipped with a Virtex-5 SX95T FPGA optimized for digital signal processing, 512 MB of onboard RAM, and 16 direct memory access (DMA) channels for high-speed data streaming at more than 800 MB/s [14].

In addition, all these modules, except for the analog cancellation part including a dual-polarization full duplex RF antenna, sit in the NI PXIe-1075 chassis. The chassis plays a role in data aggregation with both FPGA processors and an RT controller for real-time signal processing. As explained above, for transmitting and receiving simultaneously, the NI 5791R transceiver includes both transmit (Tx) and receive (Rx) ports connected to the DAC and ADC, respectively.

As can be seen in Fig. 1a, we constructed a link for full duplex radios in an exhibition hall (a severe channel environment) where a great

<sup>1</sup> The demo video is available at <http://www.cbchae.org/>.

For analog self-interference cancellation, we introduce a novel RF antenna. Our approach is based on a dual-polarization antenna with a high XPD characteristic. XPD is defined as the ratio of the co-polarized average received power to the cross-polarized average received power.

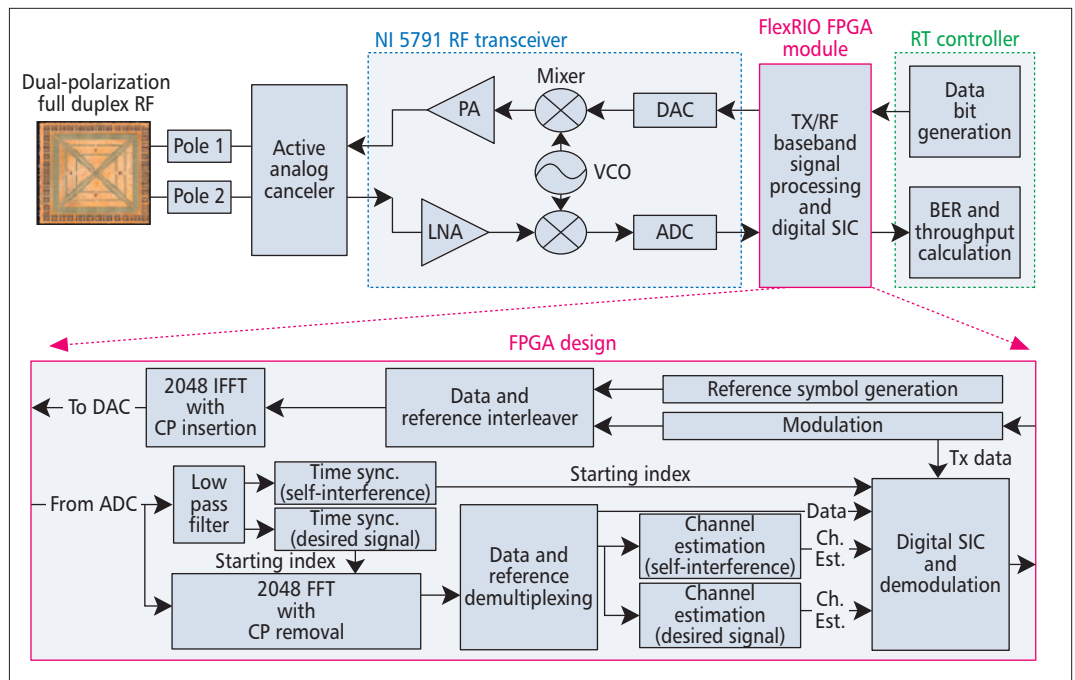


Figure 2. Block diagram of the proposed full duplex radio architecture.

crowd of people was present, as well as in an indoor open space environment. The distance between full duplex communicating nodes was about 1.2 m. Note that in fact much longer ranges are possible. In this demo/experiment, one transceiver is connected with a dual-polarization full duplex RF antenna (in the white box in Fig. 1a), and the other is connected with an omni-antenna for simplicity. In other words, the transceiver connected with the omni-antenna only transmits an uplink signal. We then observe results at the transceiver equipped with a full duplex RF antenna, where both the Tx and Rx ports are connected.

## PROPOSED FULL DUPLEX SYSTEM

In this section, we elaborate, in processing order, on our design blocks for the real-time full duplex LTE system, from transmission to reception and self-interference cancellation. The block diagram of our full duplex radio architecture is illustrated in Fig. 2.

### TRANSMISSION

As illustrated in Fig. 3a, we follow the frame structure of the LTE downlink with a frame duration of 10 ms for transmission. Each frame is divided into 20 slots, each being 0.5 ms in duration. Each slot contains 6 orthogonal frequency-division multiplexing (OFDM) symbols with 512 cyclic prefix (CP) length (extended mode). The data bit is generated on the PXIe-8133 RT controller. After the modulation block, the data symbols are interleaved with reference symbols (RS) stored in a lookup table. An array of interleaved symbols is padded with zeros to form an array of 2048 samples. The 2048 samples are passed through a 2048-point inverse FFT (IFFT) block transforming the frequency domain samples into the time domain. The 2048 IFFT with 512 CP

insertion blocks is executed on the PXIe-7965R FPGA module. To operate the discrete Fourier transform (DFT), it uses a Xilinx fast Fourier transform intellectual property (IP) core.

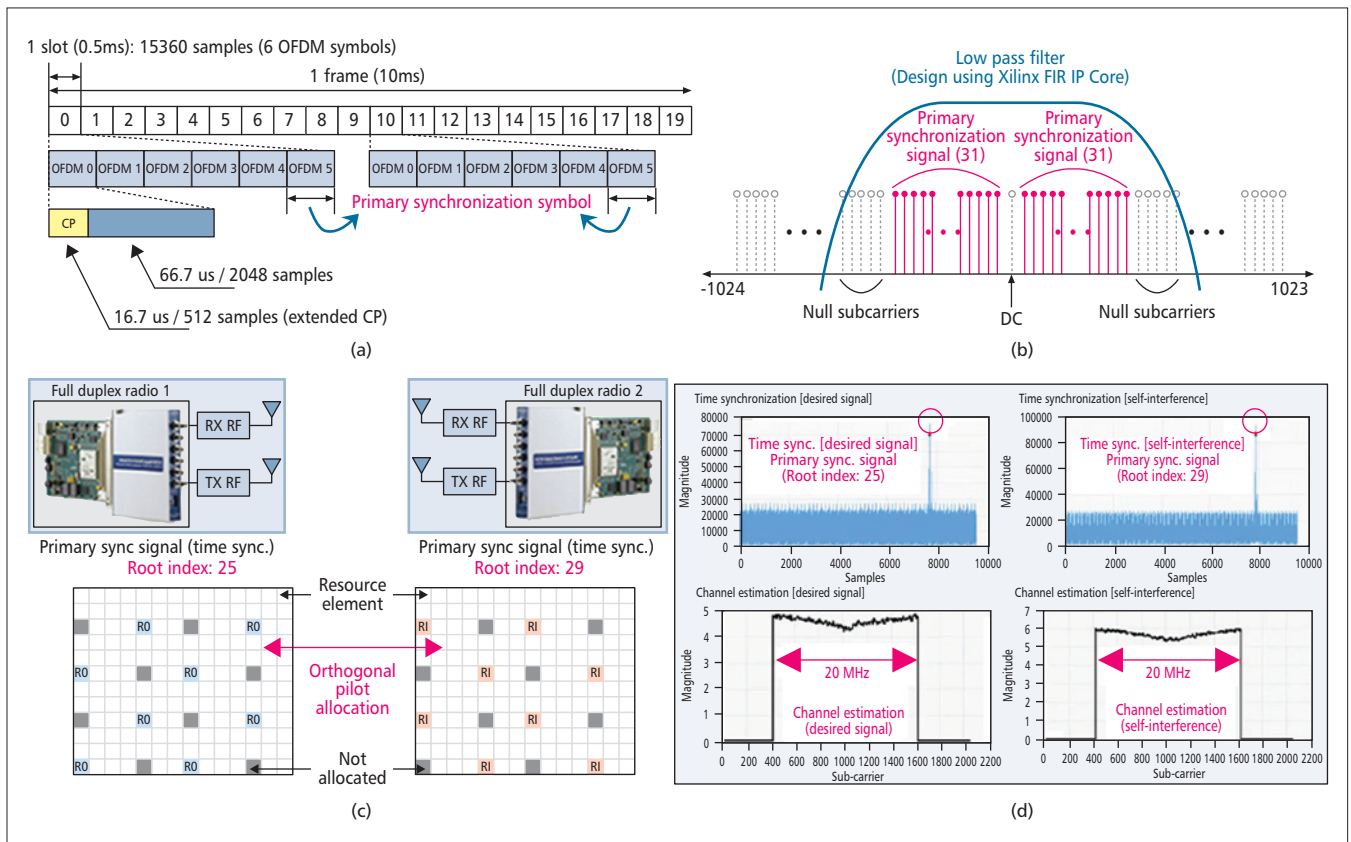
### ANALOG SELF-INTERFERENCE CANCELLATION

Conventional approaches to deal with self-interference as passive analog cancellation are:

- Isolation between Tx and Rx signals [3, 4]
- Antenna separation between the Tx and Rx antennas [5, 6]
- Signal inversion with a  $\pi$ -phase shifter [7]

Although these strategies have been extensively studied and adapted to full duplex radios as a good solution, we focus on a simpler, more compact strategy that provides outstanding self-interference cancellation performance. Furthermore, we discuss a solution that provides more robustness to environmental effects, such as the Doppler effect and multi-path, to achieve stable analog cancellation performance.

For analog self-interference cancellation, we introduce a novel RF antenna. Our approach is based on a dual-polarization antenna with a high XPD characteristic. XPD is defined as the ratio of the co-polarized average received power to the cross-polarized average received power. In other words, it represents the ability to maintain radiated or received polarization purity between horizontally and vertically polarized signals. As shown in Fig. 1a, the proposed RF unit is a compact antenna with two poles. One pole is used as a radiated Tx output; the other is used as a received Rx input in a full duplex radio. XPD is an important characteristic, particularly in full duplex systems, where cross-talk between Tx and Rx ports can curb the system's throughput performance. Since XPD has a relationship with inter-port isolation, the dual-polarization antenna with high XPD is, in full duplex systems, an excellent solution. We find that the dual-polar-



**Figure 3.** a) Frame structure of our prototype; b) arrangement of primary synchronization symbol in frequency domain; c) reference symbol patterns for a full duplex link; d) measurement of synchronization and channel estimation in the front panel of our prototype.

ization antenna itself achieves 42 dB of isolation. Active analog cancellation provides an additional cancellation gain up to 18 dB by tuning the attenuation, phase shift, and delay parameters, that is, a total of 60 dB by analog cancellation.

### DIGITAL SELF-INTERFERENCE CANCELLATION

The goal of digital self-interference cancellation is to suppress, after canceling self-interference in the analog domain, any residual self-interference. Digital self-interference cancellation consists of rebuilding self-interference and subtracting it from the received signal. A key parameter to consider in the real-time digital self-interference canceler is the guaranteed throughput of digital data, within a given time, between transmitting and receiving streams in a node. Unlike unidirectional communications that uses a radio only for transmission or reception, a digital self-interference canceler of a full duplex node requires, without a bottleneck, high-speed computation/data throughput between transmit and receive elements. To perform rebuilding self-interference and subtracting it from the received signal in real time, we use handshake protocols, shift registers, shared registers with scheduled access, and dedicated first-in-first-out (FIFO) buffers.

In [15], the authors implemented a full duplex solution that needs no additional synchronization or channel estimation for self-interference in the digital domain. Our prototype, however, focuses on an independent system that operates in real time to maximize cancellation perfor-

mance in the analog/digital domain, respectively. Furthermore, additional synchronization and estimation is exploited to reduce the complexity in our digital self-interference canceler.

At the moment of decoding the desired symbol, it is critical to know the perfect timing between self-interference and the received symbol in full duplex mode. Thus, key issues include designing synchronization and channel estimation strategies for residual self-interference as well as for a desired link. We produce a process for implementing a digital self-interference canceler from synchronization and channel estimation. In order to operate a real-time digital self-interference canceler with high performance, we focus on FPGA implementation using the LabVIEW system design software and PXIe SDR platform.

**Synchronization:** Synchronization is one of the key blocks in real-time full duplex radios. Under synchronization for full duplex, there are two operations: synchronization for decoding the desired symbol and rebuilding self-interference. In the synchronization block for decoding the desired symbol, we estimate time offset by random propagation delays and sampling clock offsets between two full duplex radios. In the synchronization block for rebuilding self-interference, we estimate the time offset between the Tx port and Rx port of a full duplex radio.

To facilitate timing synchronization, the LTE downlink standard specifies a primary synchro-

To facilitate timing synchronization, the LTE downlink standard specifies a primary synchronization signal. Accordingly, the receiver can successfully perform timing synchronization in half duplex mode. Note, however, that we need to keep performing synchronization for the self-interference signal as well as for the desired signal.

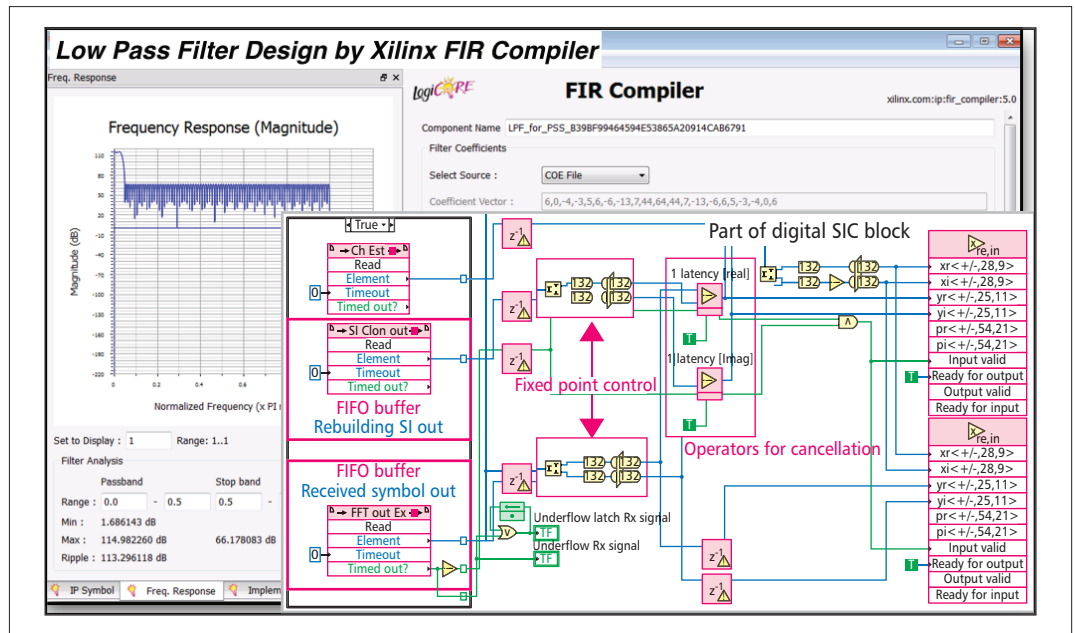


Figure 4. A part of source code for a digital self-interference canceller that operates in real time on the software-defined radio platform.

nization signal (PSS). Accordingly, the receiver can successfully perform timing synchronization in half duplex mode. Note, however, that we need to keep performing synchronization for the self-interference signal as well as for the desired signal. Thus, we use a property where the Zadoff-Chu (ZC) sequence with a different root index is orthogonal to each other. The PSS is modulated by a ZC sequence given as,  $P[k] = e^{-j(\pi/N)uk(K+1)}$ ,  $-31 \leq k \leq -1$ , and  $P[k] = e^{-j(\pi/N)u(k+1)(k+2)}$ ,  $1 \leq k \leq 31$ , where  $k$  is the subcarrier index,  $u$  is the root index, and  $N$  is the sequence length ( $N = 63$ ). We use a different root index relatively prime to  $N$  for the PSS of each full duplex radio, that is,  $u_1 = 25$ ,  $u_2 = 29$ . These symbols are located on the 62 subcarriers, symmetrically arranged around the DC carrier in the last OFDM symbol of the first and 11th slots of each frame, as shown in Figs. 3a and 3b. As the duration of a frame is 10 ms, the PSS is transmitted after every 5 ms time interval or once per half-frame.

To calculate the correlation between the ideal sequence and the estimated PSS signal, it is necessary to extract the PSS subcarrier from the received signal. For this reason, we design a low-pass filter (LPF) using Xilinx's finite impulse response (FIR) IP core, as shown in Fig. 3b. The designed LPF has a cutoff frequency of 1.4 MHz, a stop-band attenuation of 50 dB, and a pass-band ripple of 0.1 dB. After the received signal samples are passed through the LPF, each synchronization block for decoding the desired symbol and rebuilding the self-interference is executed to calculate, independently, the correlation with its own PSS. As a result, a maximum peak is detected at the sample index of the first sample of the OFDM symbol following the PSS symbol, as illustrated in Fig. 3d. A starting index of the desired signal is delivered into FFT block, and a starting index of self-interference signal is delivered into the digital cancellation block.

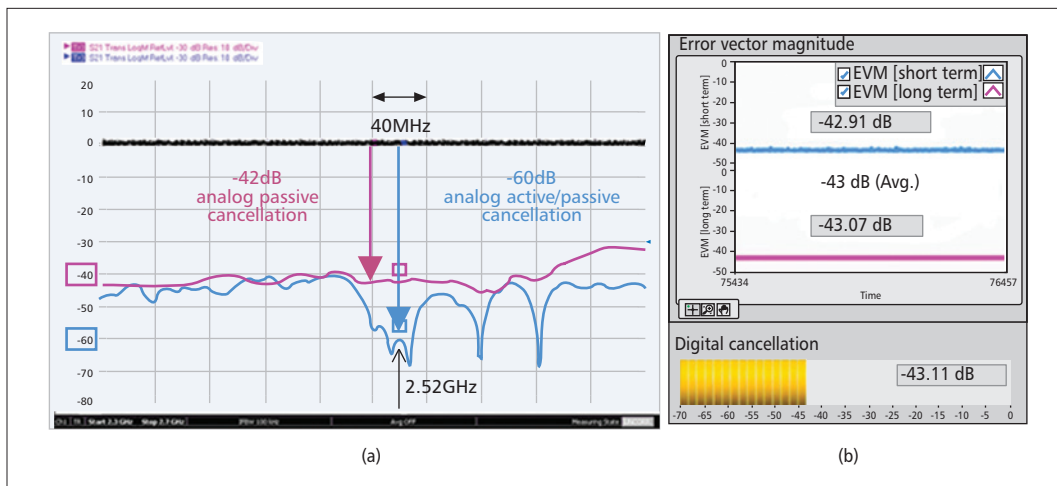
**Channel Estimation:** Channel estimation also

has, for full duplex, two operations: estimations for the channels between two full duplex nodes, and the channel between the Tx and Rx ports in its own full duplex node. The former operates to decode the desired symbol after digital self-interference cancellation, while the latter operates to rebuild the self-interference using the known Tx data. In order to handle the two operations simultaneously, we design RS patterns that are orthogonal between two full duplex nodes, exploiting the pattern of cell-specific reference signals for multiple antenna ports. The RS patterns are shown in Fig. 3c.

Both channel estimation blocks have two steps in common: RS extraction and interpolation in order. After the received samples are passed through the FFT block, the RS subcarriers of each channel are extracted from an OFDM symbol in a data and reference demultiplexing block. A channel coefficient of each RS subcarrier is calculated using original RSs stored in block memory. To estimate the channel coefficients of RS subcarriers, a least-square method is exploited. After passing through the data and reference demultiplexing block, the channel estimates of the RS subcarriers are split into two groups. One is for the channel estimation between nodes; the other is for the channel estimation between antenna ports. In each channel estimation block, we implemented a linear interpolator using Xilinx's FIR IP core. The linear interpolator in each block estimates the channel coefficients of data subcarriers as well as RS subcarriers. In Fig. 3d, the bottom left and right figures are screen shots of the instantaneous channel estimation result between nodes and ports, respectively, in the frequency domain.

**Digital Cancellation:** In most research on full duplex implementation, digital self-interference cancellation is performed in the time domain. To carry out digital cancellation in the time domain,





**Figure 5.** Measurement results of a) analog self-interference cancellation; b) digital self-interference cancellation measurement.

an additional IFFT block is needed for rebuilding self-interference. As mentioned above, in the real-time digital self-interference canceler, it is critical to operate the process for rebuilding self-interference and canceling it out without a bottleneck. For lower computational complexity and faster rebuilding of self-interference, we execute it in the frequency domain. One might argue that if the SI is not overlapped coherently with the received signal, it would be difficult to cancel the SI in the frequency domain. However, if the difference in the received time of the SI and the desired signal is less than the CP length, which is common in practice, there is no problem in canceling the SI after FFT. We will also investigate this issue in our future work.

Digital cancellation utilizes the baseband samples of the transmitted signal to rebuild self-interference in the digital domain and subtracts them from the received samples. Note that we know the baseband samples of the transmitted signal from its own node. Self-interference can be rebuilt in the digital domain using the baseband samples of the transmitted signal and the channel estimates between the ports of its own node. As mentioned above, we should know which self-interference (a sample index) is mixed in the received sample at the moment of decoding the desired symbol. Accordingly, we include a counter in the digital cancellation block. As soon as the starting index of the self-interference signal arrives in the digital cancellation block from the synchronization block to rebuild self-interference, the counter operates to choose a rebuilt digital sample for subtraction processing. After digital cancellation, a zero-forcing channel equalizer operates to decode the desired symbols. Illustrated in Fig. 4 is a part of the source code for the digital self-interference canceler operated in the frequency domain.

## PROTOTYPE TEST RESULTS

Using the real-time full duplex LTE prototype described in the previous section, we measure the level of analog and digital self-interference and calculate the bit error rate (BER) and throughput performance. In this prototype, the

carrier frequency is the 2.52 GHz in LTE bands. As shown in Fig. 5a, we find that the dual-polarization antenna provides about 42 dB isolation from our experiments, that is, the self-interference leaking to the Rx port is reduced by about 42 dB. Moreover, by tuning the attenuation, phase shift, and delay parameters, we achieve 60 dB of analog self-interference cancellation with analog active cancellation. In the digital domain, we calculate error vector magnitude (EVM) for self-interference to measure the average level of digital self-interference cancellation. As a result, we achieve 43 dB of self-interference cancellation in the digital domain, as shown in Fig. 5b.

In order to compare the throughput improvement, we also implemented the LTE-FDD prototype. Our design is based on the LTE downlink standard with system specifications that include a transmission bandwidth of 10 MHz, 15.36 MHz sampling rate, 15 kHz subcarrier spacing, 1024 FFT size, 256 CP length, and variable 4/16/64 QAM. Figure 6a shows the constellation, taken during an over-the-air test of the full duplex communications link. One full duplex radio transmits a 4-QAM modulated signal as the downlink and receives a 64-QAM modulated signal as the uplink. As a result, the goal of this full duplex radio is to decode 64-QAM, the desired symbol after perfectly canceling out the 4-QAM symbol as self-interference. In Fig. 6a, the left constellation shows that, with only analog cancellation, self-interference is not perfectly cancelled out, while the right constellation shows that with both analog and digital cancellation self-interference is perfectly cancelled out.

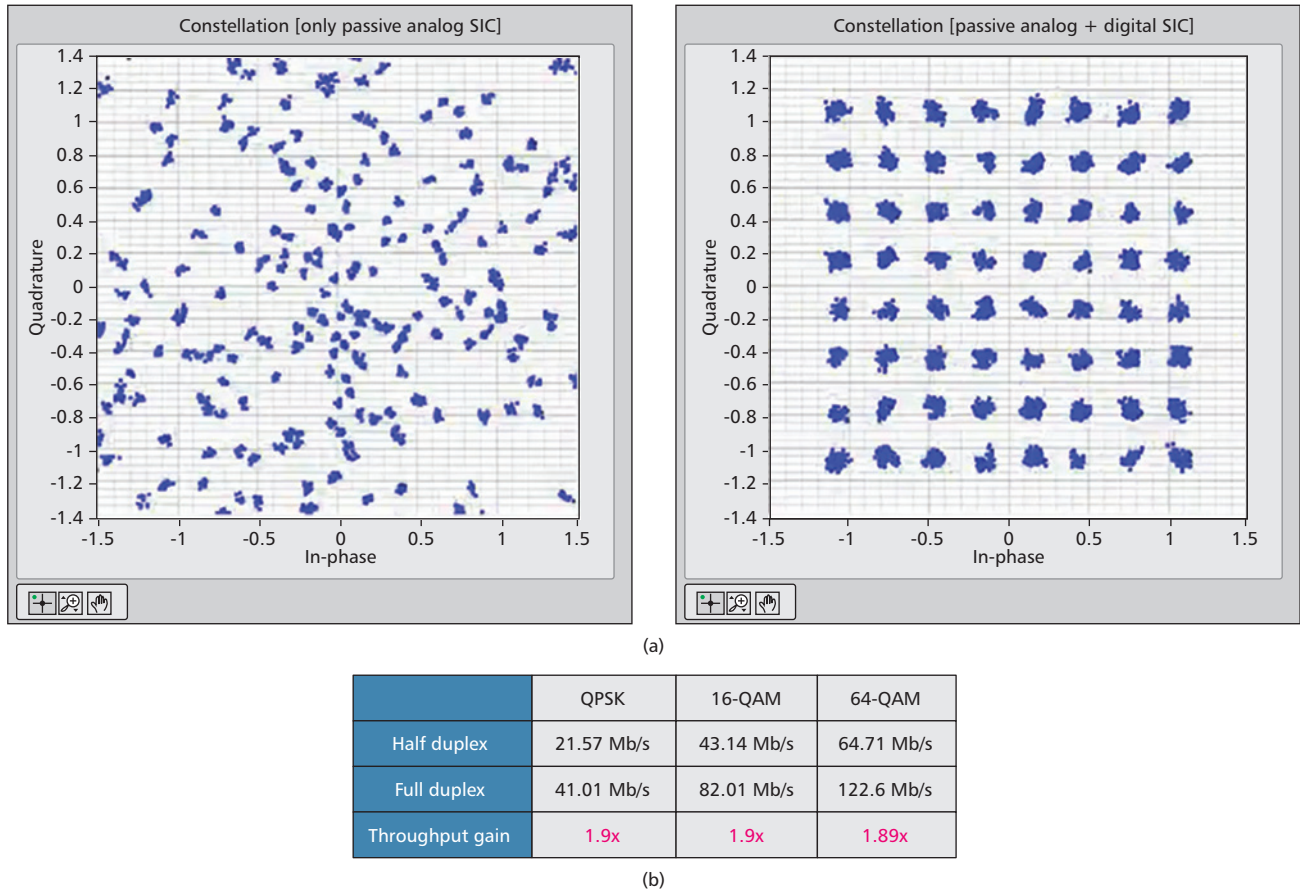
As shown in Fig. 6b, our prototype provides high throughput performance in real time. It delivers a throughput increase of  $1.9\times$  on 4- and 16-QAM, and  $1.89\times$  on 64-QAM compared to the conventional half duplex mode.<sup>2</sup>

## RESEARCH CHALLENGES

Notwithstanding our focus on designing more practical full duplex radios, several research challenges remain before the most viable solution for next generation communication systems is achieved.

*In the digital domain, we calculate error vector magnitude for self-interference to measure the average level of digital self-interference cancellation. As a result, we achieve 43 dB of self-interference cancellation in the digital domain.*

<sup>2</sup> Half duplex means out-of-band full duplex, that is, FDD throughout this article.



**Figure 6.** a) Constellation result with only analog passive self-interference cancellation (left) constellation result with analog and digital self-interference cancellation (right); b) throughput performance for each constellation.

### HARDWARE IMPAIRMENTS

The performance of a full duplex system depends heavily on hardware impairments: amplifier non-linearity, gain/phase offset, I/Q imbalance, quantization effects, and timing jitter. For example, a nonlinearly amplified OFDM signal incurs inter-modulation distortion (IMD), which is the amplitude modulation of signals containing two or more different frequencies in a system. IMD raises the noise floor and causes inter-carrier interference, which induces performance degradation of the full duplex system. Because most analog components in the system have the non-linearity property, the cancellation of all nonlinear components as well as linear components is a significant burden on a real-time system. Thus, preprocessing strategies such as predistortion for reducing hardware impairments represent an interesting research topic.

### JOINT PHY/MAC PROTOTYPING

Most implementations of full duplex radios have mainly focused on the physical layer design, which enables bidirectional communications between a single pair link. There are apparent limitations in translating the performance gains obtained from the demonstration of a single pair link into network performance. Transmissions in full duplex mode create potential interference outside the full duplex link. This calls for the

prototyping of media access control (MAC) layer protocols, including discovering and exploiting full duplex opportunities in a distributed manner. Another interesting area for future work is a joint PHY/MAC approach to prototyping.

### FULL DUPLEX SYSTEM WITH OFDM AND SC-FDMA

Since single-carrier frequency-division multiple access (SC-FDMA) has a peak-to-average power ratio (PAPR) lower than that of OFDMA, it is used for the uplink multiple access scheme in the LTE of cellular systems. Most implementation studies of full duplex, however, only deal with OFDM frame structures. There are many potential challenges in asymmetric uplink/downlink frame structures in LTE.

### COMPARISON WITH LTE-TDD

LTE-time-division duplexing (TDD) generally has been known to have many benefits, such as low latency, spectrum flexibility, uplink/downlink flexibility, and lower cost per bit. To discuss the various performance characteristics such as latency, throughput, power consumption, and flexibility between TDD and full duplex, we believe that it is worth comparing a full duplex prototype with a comparable LTE-TDD prototype.

## NOVEL SOLUTION FOR RF/ANALOG CANCELLATION

RF/analog cancellation plays a critical part in attenuating high-powered self-interference sufficiently such that Rx saturation and dynamic range are not an issue when operating digital cancellation. We showed that the analog cancellation based on dual polarization would be a good option. However, it does have two main weaknesses. First, it struggles to perform active analog cancellation in real time. While high XPD makes the passive analog cancellation level high, to estimate the coefficients for active cancellation is difficult. Second, channel reciprocity that can simplify the link overhead may not be assumed for the polarization antennas. Novel solutions for these issues will be an interesting research topic.

## CONCLUSION

Full duplex radio technologies could be a major contributor to increasing spectrum efficiency in areas of explosive traffic demand where there are limited radio resources. To validate the feasibility and commercial viability of any new wireless standard or technology like full duplex radio, SDR-based prototyping is imperative. We have prototyped a design that combines dual-polarization full duplex RF and a digital self-interference canceler that operates in real time on an SDR platform. We have focused on a more practical prototype that exhibits outstanding self-interference cancellation performance. The main portion of this article is dedicated to presenting the design, implementation, and evaluation of a real-time full duplex LTE system, a candidate for next generation wireless communication systems. We expect our prototype design to provide worthwhile insights into developing the most viable solution for future wireless communication systems.

## REFERENCES

- [1] C. V. Forecast, "Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2014–2019," Cisco Public Info., 2015.
- [2] J. Mitola, "The Software Radio Architecture," *IEEE Commun. Mag.*, vol. 33, no. 5, May 1995, pp. 26–38.
- [3] D. Bharadia, E. McMillin, and S. Katti, "Full Duplex Radios," *Proc. ACM SIGCOMM*, 2013, pp. 375–86.
- [4] S. Hong *et al.*, "Applications of Self-Interference Cancellation in 5G and Beyond," *IEEE Commun. Mag.*, vol. 52, no. 2, Feb. 2014, pp. 114–21.
- [5] M. Duarte, C. Dick, and A. Sabharwal, "Experiment-Driven Characterization of Full-Duplex Wireless Systems," *IEEE Trans. Wireless Commun.*, vol. 11, no. 12, Dec. 2012, pp. 4296–4307.
- [6] M. Duarte *et al.*, "Design and Characterization of a Full-Duplex Multiantenna System for WiFi Networks," *IEEE Trans. Vehic. Tech.*, vol. 63, no. 3, March 2014, pp. 1160–77.
- [7] E. Aryafar *et al.*, "MIDU: Enabling MIMO Full Duplex," *Proc. ACM MobiCom*, 2012, pp. 257–68.
- [8] D. Kim, H. Lee, and D. Hong, "A Survey of In-Band Full-Duplex Transmission: From the Perspective of PHY and MAC Layers," *IEEE Commun. Surveys and Tutorials*, 2015.
- [9] C.-B. Chae, "Double the Capacity: A Real-Time Full Duplex LTE System", <https://www.youtube.com/watch?v=AS11IQEQzDE>, 2014.
- [10] S. Sesia, I. Toufik, and M. Baker, *LTE: The UMTS Long Term Evolution*, Wiley, 2009.

- [11] T. Oh, Y. Lim, C.-B. Chae, and Y. Lee, "Dual-Polarization Slot Antenna with High Cross Polarization Discrimination for Indoor Small-Cell MIMO Systems," *IEEE Ant. Wireless Prop. Lett.*, vol. 14, Feb. 2014, pp. 374–77.
- [12] NI PXIe-8133 User Manual, <http://www.ni.com/pdf/manuals/372870d.pdf>, 2012.
- [13] NI 5791R User Manual and Specifications, <http://www.ni.com/pdf/manuals/373845c.pdf>, 2013.
- [14] NI FlexRIO FPGA Modules Data Sheet, <http://www.ni.com/datasheet/pdf/en/ds-366>, 2014.
- [15] Y.-S. Choi and H. Shirani-Mehr, "Simultaneous Transmission and Reception: Algorithm, Design and System Level Performance," *IEEE Trans. Wireless Commun.*, vol. 12, no. 12, Dec. 2013, pp. 5992–6010.

## BIOGRAPHIES

MINKEUN CHUNG [S'11] received his B.S. degree from the School of Electrical and Electronic Engineering at Yonsei University, Korea, in 2010. He is now working toward his Ph.D. degree under the joint supervision of Prof. D. K. Kim and Prof. C.-B. Chae. He did his graduate internship with the Advanced Wireless Research Team at National Instruments in Austin, Texas in 2013 and 2015. His research interests include the design and implementation of architectures for next-generation wireless communication systems.

MIN SOO SIM received his B.S. degree from the School of Integrated Technology at Yonsei University in 2014. He is now with the same school and university working toward his Ph.D. degree. His research interest includes emerging technologies for 5G communications.

JAWEON KIM [M'11] received his B.S. and M.S. in electrical engineering from Korea Advanced Institute of Science and Technology (KAIST) in 1994 and 1996, respectively, and his Ph.D. in electrical and computer engineering from the University of Texas at Austin in 2011. From 1996 to 2002, he was with SK Telecom, Seoul, Korea, where he worked on 2G and 3G CDMA systems and their applications. During his Ph.D. program, he held a member of technical staff position at Bandspeed, Inc., Austin, Texas from 2006 to 2008 and a senior member of technical staff position at MediaExcel, Inc., Austin, from 2008 to 2011. Currently he is with National Instruments, Austin, as a senior wireless platform architect. His current research interests include 5G wireless communications, digital signal processing, and prototyping.

DONG KU KIM [SM'15] received his B.S. from Korea Aerospace University in 1983, and his M.S. and Ph.D. from the University of Southern California, Los Angeles, in 1985 and 1992, respectively. He worked on CDMA systems in the cellular infrastructure group of Motorola at Fort Worth, Texas, in 1992. He has been a professor in the School of Electrical and Electronic Engineering, Yonsei University, since 1994, and was the principal investigator professor of the Qualcomm Yonsei Joint Research Program from 1999 to 2010. Currently, he is a vice president for academic research affairs of the Korean Institute of Communications and Information Systems. He received the Minister Award for the Distinguished Service for ICT R&D from the Ministry of Information, Science, and Future Planning in 2013, and the Award of Excellence in leadership of 100 Leading Core Technologies for Korea 2020 from the National Academy of Engineering of Korea.

CHAN-BYOUNG CHAE [SM'12] is an associate professor in the School of Integrated Technology, Yonsei University. Before joining Yonsei University, he was with Bell Labs, Alcatel-Lucent, Murray Hill, New Jersey, as a member of technical staff, and Harvard University, Cambridge, Massachusetts, as a postdoctoral research fellow. He received his Ph.D. degree in electrical and computer engineering from the University of Texas at Austin in 2008. He was the recipient/co-recipient of the IEEE INFOCOM Best Demo Award (2015), the IEIE/IEEE Joint Award for Young IT Engineer of the Year (2014), the KICS Haedong Young Scholar Award (2013), the *IEEE Signal Processing Magazine* Best Paper Award (2013), the IEEE ComSoc AP Outstanding Young Researcher Award (2012), the IEEE Dan. E. Noble Fellowship Award (2008), and two Gold Prizes (1st) in the 14th/19th Humantech Paper Contest. He currently serves as an Editor for *IEEE Transactions on Wireless Communications*, the *IEEE/KICS Journal on Communications Networks*, and *IEEE Transactions on Molecular, Biological, and Multi-scale Communications*.

Full duplex radio technologies could be a major contributor to increasing spectrum efficiency in areas of explosive traffic demand where there are limited radio resources. To validate the feasibility and commercial viability of any new wireless standard or technology like full duplex radio, SDR-based prototyping is imperative.

# A Low-Cost Desktop Software Defined Radio Design Environment Using MATLAB, Simulink, and the RTL-SDR

*Robert W. Stewart, Louise Crockett, Dale Atkinson, Kenneth Barlee, David Crawford, Iain Chalmers, Mike McLernon, and Ethem Sozer*

## ABSTRACT

In the last five years, the availability of powerful DSP and communications design software, and the emergence of relatively affordable devices that receive and digitize RF signals, has brought SDR to the desktops of many communications engineers. However, the more recent availability of very low cost SDR devices such as the RTL-SDR, costing less than \$20, has brought SDR to the home desktops of undergraduate and graduate students, as well as professional engineers and the maker communities. Since the release of the various open source drivers for the RTL-SDR, many in the digital communications community have used this device to scan the RF spectrum and digitize I/Q signals that are being transmitted in the range 25 MHz to 1.75 GHz. This wide operating range enables the sampling of frequency bands containing signals such as FM radio, ISM signals, GSM, 3G and LTE mobile radio, GPS, and so on. In this article we will describe the opportunity and operation of the RTL-SDR, and the development of a hands-on, open-courseware for SDR. These educational materials can be integrated into core curriculum undergraduate and graduate courses, and will greatly enhance the teaching of DSP and communications theory, principles, and applications. The lab and teaching materials have recently been used in senior (fourth year undergraduate) courses and are available as open course materials for all to access, use, and evolve.

## INTRODUCTION

In this article we present our experience of developing university and CPD (continuous professional development) materials for teaching SDR in the form of DSP-enabled-radio systems. The availability of SDR receivers such as the RTL-SDR, along with hardware support software drivers, means that we now have devices that are very low cost and can interface directly with MATLAB and Simulink software, allowing users

to develop real software defined radio systems from the desktop. The RTL-SDR plug-in device (which comes with a simple but useable omnidirectional antenna) currently costs less than \$20 (twenty dollars) to buy, and can be powered and connected via a USB port to Windows, Linux, and or Mac desktop computers. Students can acquire the MATLAB and Simulink student version ([http://www.mathworks.com/academia/student\\_version/](http://www.mathworks.com/academia/student_version/)), along with the relevant DSP and Communications System Toolboxes, for around \$100, and after the installation of appropriate support drivers, can be up and running with a complete SDR design environment. As we will summarize in this article, the opportunities for using this for education and learning are immense, and a whole new generation of engineers will see more and more RF and communications design being done as a coding task. The full set of SDR open-course educational materials referred to in this article are available to download (<http://www.desktopSDR.com>) as a 670 page workbook with more than 120 hands-on examples [1].

In this article we will first outline what the RTL-SDR is, “where” and how it evolved, and then introduce our open-source teaching and support materials for learning SDR from a DSP-enabled-radio perspective. The SDR design environment and open-course materials described herein have the potential to be used in classes ranging from EE freshman (first year bachelor) university environments, for courses featuring applications such as spectral viewing and first experiences in radio, all the way to EE senior (fourth year bachelor) or masters level classes teaching, for example, the challenging aspects of QPSK receivers with synchronization, and other digital communications systems [2, 3]. Real practical experience in these communications applications and theory can be achieved using a low cost SDR receiver that students and home users can keep in their pocket, and connect via USB to run SDR algorithms directly on the desktop of their laptop device, and all having spent less than \$20 on the RTL-SDR hardware!

*Robert W. Stewart, Louise Crockett, Dale Atkinson, Kenneth Barlee, David Crawford, and Iain Chalmers are with University of Strathclyde.*

*Mike McLernon and Ethem Sozer are with MathWorks Inc.*

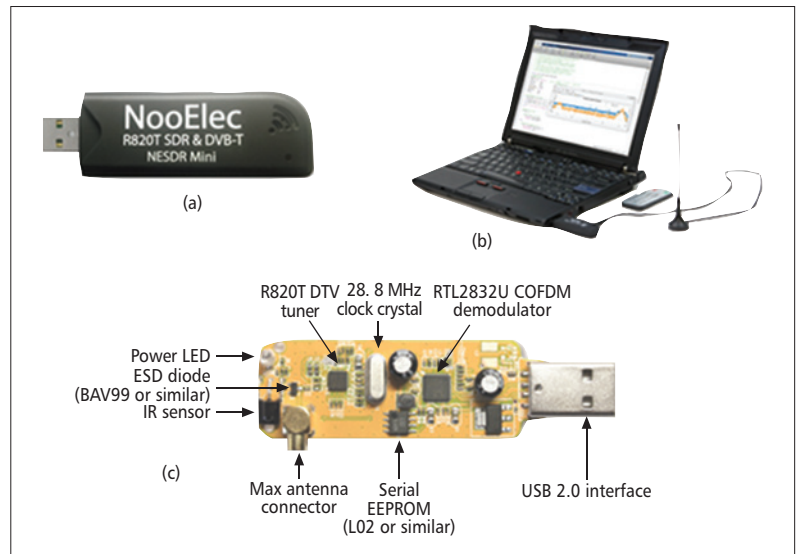
SDR is now in the domain of being “home-work” where students can work at home using their RTL-SDR and MATLAB and Simulink software. There is no longer a requirement for expensive spectrum analyzer hardware, and no requirement for signal generators; all that is needed is the ability to buy a USB RTL-SDR stick device on-line, download drivers, and then develop the appropriate SDR receivers and systems to find signals of interest that can be spectrally viewed, analyzed, and even decoded.

## THE ADVENT OF DESKTOP SDR

Over the last 20 years the prospect of software (defined) radios has been greeted with enthusiasm by the DSP (digital signal processing), digital communications, and radio/RF sectors [4, 5]. In some ways the term software defined radio (SDR) has perhaps diverged in recent years to have different meanings for different engineering groups. Many in the DSP community considered that, by virtue of very high speed ADCs (analog to digital converters) and DACs (digital to analog converters), SDR was in effect the engineering of “DSP-enabled-radio systems,” that is, where analog oscillators would be replaced with digital numerically controlled oscillators, analog filtering with digital filter chains, and phase locking components with digital phase locked loops (PLLs), and so on. Whereas in other communications system engineering domains, SDR actually refers to middleware, which is the software that could define the radio and provide the framework for the deployment of software objects over networks and between devices in the radio hardware [6, 7]. (SDR middleware would ultimately be the term to describe the control and design of high power computing platforms that would allow radio standards and waveforms to be switched in and out and downloaded on the fly, as pursued in applications such as JTRS (joint tactical radio service) from 1998 to 2011 [7, 8].)

In both closely related interpretations of SDR, i.e. DSP-enabled-radio and middleware, its concept and promise were easy to understand, but the hardware and software that was required 20 years ago was far beyond the then-affordable state of the art. But of course Moore’s law never fails (or hasn’t yet). Hence the reality and ease of access to SDR technologies is definitely here, both for the DSP-enabled-radio and the middleware groupings.

In the last five years or so, SDR in the DSP-enabled-radio category has been achievable in the lab at a reasonable cost (less than \$1500) for FPGA-enabled hardware with ADC and DAC units typically sampling at rates of a few hundred MHz, and front end radio cards that worked at up to 6 GHz, such as the ubiquitous USRP series from Ettus Research (<http://www.ettus.com/product/category/USRP-Bus-Series>). Similarly, various software platforms allow users to code and configure these devices using FPGA design environments that are often driven from DSP and communications development tools. SDR hardware products such as the URSP have been widely used to stream samples of down-converted RF signals to the desktop, where they are input



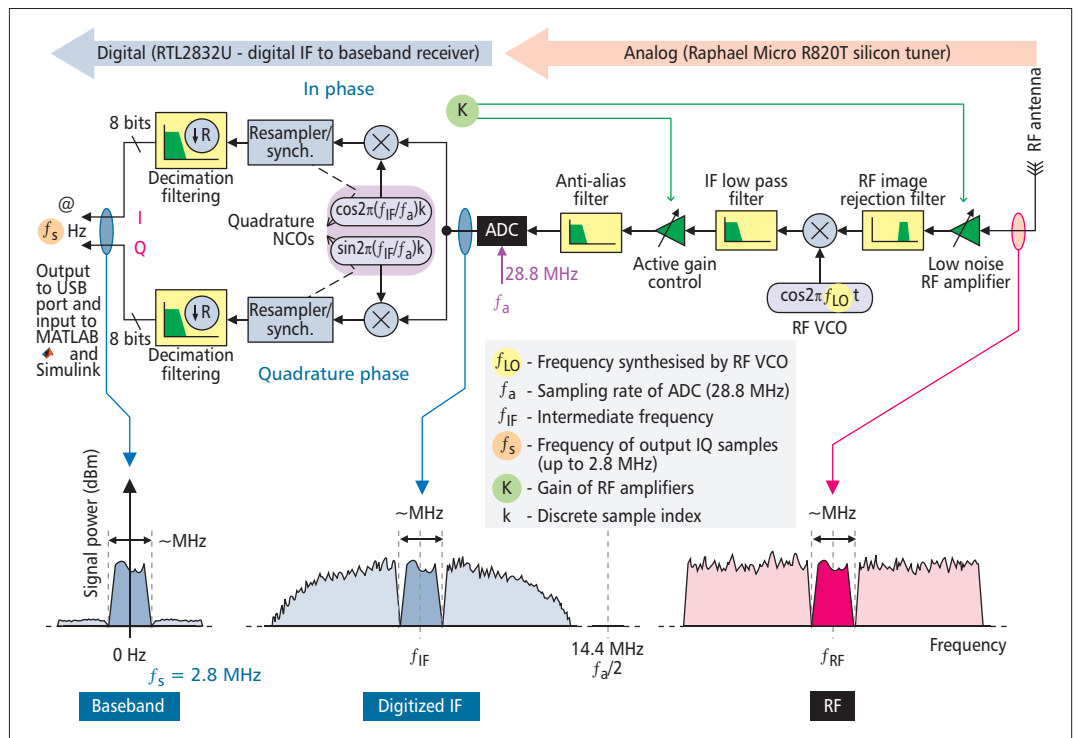
**Figure 1.** a) The RTL-SDR Mini USB device; b) a typical RTL-SDR receiver setup on a laptop running MATLAB and Simulink using a simple omni-antenna (comes with the RTL-SDR); c) the main internal components of the RTL-SDR.

to software such as MATLAB and Simulink for real-time processing or recorded off-line for off-line use. Additionally, where drivers were available and the desktop was a high speed computing platform, then real-time DSP software algorithms could be implemented (in floating point on the desktop processor) and receivers and transmitters implemented. Software defined radio is now established in a number of institutions as part of the curriculum, and in the May 2014 special SDR education feature topic in *IEEE Communications Magazine*, a number of authors reviewed their successful experiences introducing USRP-based SDR into their integrated course curricula and laboratory sessions for EE students [9–11]. However, the advent of the RTL-SDR device brings the affordability of a device down to a level lower than a textbook, and many of these successful courses can now also consider using the RTL-SDR as part of their laboratories, or as stand-alone learning assignments for students to do at home [12].

## WHAT IS THE RTL-SDR?

As shown in Fig. 1, the RTL-SDR is a small, compact, and easy-to-use USB stick device that is capable of receiving RF radio signals (“RTL” is actually not an acronym, but derives from the Realtek RTL2832U chip on which the device is based). Originally these devices were designed for use as DVB-T (digital video broadcast–terrestrial) receivers and featured custom-designed, tunable RF front end chips (e.g., the Rafael Micro R820T and the Elonics E4000) that allowed consumers to receive and watch UHF broadcast TV on their computers. In other words, these receivers were *not* originally designed or conceived to be used as generic programmable SDRs. The uptake of these devices as SDR receivers results from the efforts of a number of independent engineers and developers in the SDR community, who discovered their

The advent of the RTL-SDR device brings the affordability of a device down to a level lower than a textbook, and many of these successful courses can now also consider using the RTL-SDR as part of their laboratories, or as stand-alone learning assignments for students to do at home.



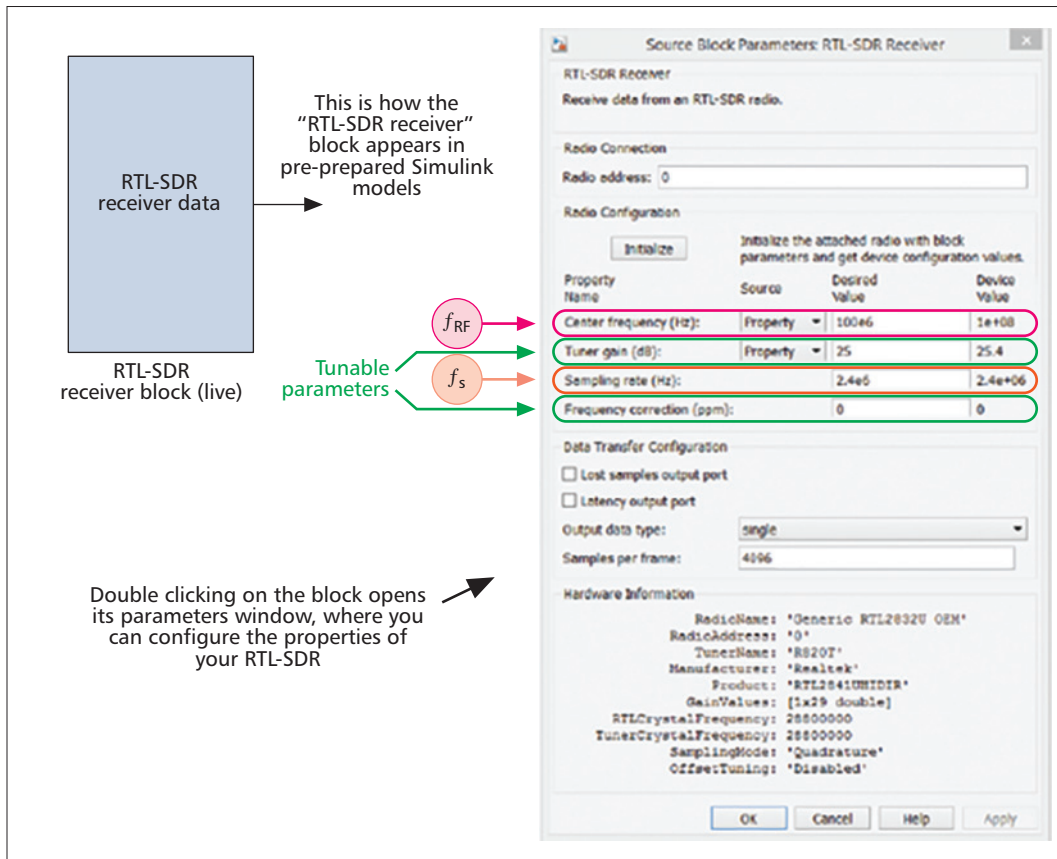
**Figure 2.** The main components of the RTL-SDR USB device. The main MATLAB and Simulink parameters of  $f_{RF}$ ,  $f_{fs}$ , and  $K$  can be set to configure the device (Fig. 3).

programmability for SDR. Specifically, it was found that the devices could be placed in a “test mode,” which meant that the RTL2832U chip (Fig. 1c) bypassed the DVB decoding stage and produced raw, 8-bit I/Q data samples. Further, it was also found to be possible to program the center frequency of the RF chip over a working range of approximately 25 MHz to 1.75 GHz, and have some control over the data sampling rate. Soon after this discovery, the name RTL-SDR was coined, which referred to the fact that the RTL (Realtek) based DVB receivers could be used as SDRs. With such a wide front end tuning frequency range, many different applications using various modulation schemes ranging from AM and FM, to ISM, GSM, LTE, and GPS applications, have become signals that we can attempt to capture. The noise floors, signal resolutions, and frequency accuracy of these devices is not optimal in some frequency bands, nor sufficient for some applications. However this investigation is all part of learning what the RTL-SDR offers, and it does work very well and successfully receives in many frequency bands and for a variety of applications.

As a quick review of the history of the RTL-SDR’s emergence, its origins were evident in some 2012 forum posts by a Linux developer on the V4L GMANE forum, stating that “radio sniffs” were possible using an RTL-based DVB device. It was discovered that when the device was tuned to receive FM and DAB radio stations, it was programmed into a different mode and that raw, modulated data samples could be transferred to the computer and demodulation performed in software (<http://comments.gmane.org/gmane.linux.drivers.video-input-infrastructure/44461>). Seventeen seconds of data originating from a

Finnish radio station was captured and posted online, along with a query asking if anyone could work out how to demodulate it manually. This was accomplished only 36 hours later, after some collaborative effort. In the original post the last line is the optimistic statement, “I smell a very cheap poor man’s software defined radio here :)”! This discovery led to further investigation of the RTL-SDR’s USB protocol. The commands transmitted when tuning to a radio station were captured, and used to force the device to stay in this special mode continuously. It turned out to be a test mode, and when the RTL2832U was in this mode, it output 8-bit unsigned samples of baseband I/Q data, rather than decoded DVB signals as per its designed operation. Work reported at the open source website Osmocom included reports from developers who had produced an independent SDR device called “OSMO-SDR,” and had experience in writing software that was able to program the DTV tuners used with the RTL2832U. After examining the Windows drivers provided by Realtek, they devised how to program the tuner via the demodulator, and the drivers for the RTL-SDR were released to the open-source community (<http://sdr.osmocom.org/trac/wiki/rtl-sdr>).

RTL-SDR, as we now know it, came onto the market in early 2013, and various devices and software kits became available, produced by a number of companies and developers around the world. Judging by the communities on the web, the RTL-based DVB-T devices appear to be more popular as SDR receivers than they were for their original intended purpose of digital TV reception! NooElec is one company with worldwide distribution of these devices (<http://www.nooelec.com/store/sdr.html>). Based on their use of the R820T tuner, the NooElec RTL-SDR



**Figure 3.** The RTL-SDR block in Simulink and the configuring parameters.

*RTL-SDR, as we now know it, came onto the market in early 2013, and various devices and software kits became available, produced by a number of companies and developers around the world. Judging by the communities on the web, the RTL-based DVB-T devices appear to be more popular as SDR receivers than they were for their original intended purpose of digital TV reception!*

devices are capable of reliably sampling the frequency spectrum at a rate up to 2.8 MHz, and receiving signals in the RF frequency range 25 MHz to 1.75 GHz.

MathWorks released a hardware support package for the RTL-SDR in early 2014 (<http://www.mathworks.com/hardware-support/rtl-sdr.html>) which enables both MATLAB and Simulink to interface with and control the RTL-SDR. With this support package, baseband samples output from the RTL-SDR device are supplied into the software environment, enabling users to implement any kind of DSP receiver or spectrum sensing system they desire as either a Simulink model or MATLAB code. I/Q data can be locally recorded to data disk files for later processing, or if processing power allows on the desktop computer, live demodulation and decoding can be performed.

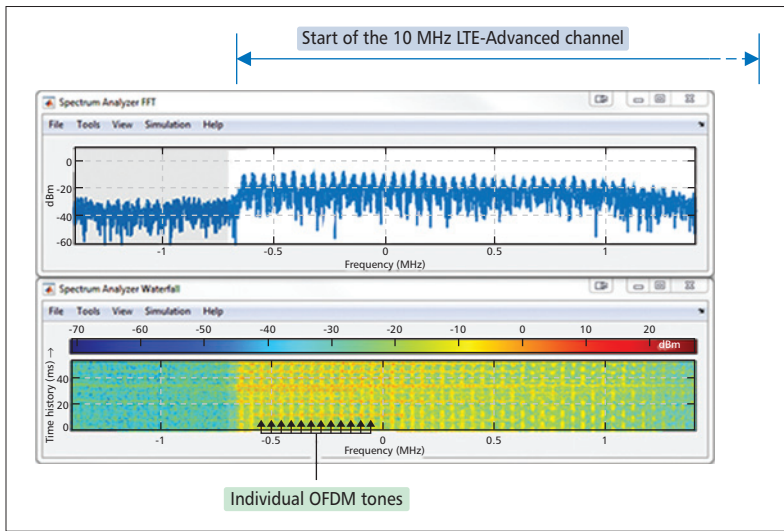
Figure 2 shows a signal processing flow diagram of the main stages that are carried out on the RTL-SDR. RF signals entering the R820T tuner (on the right hand side) are downconverted to a low-IF (intermediate frequency) using a voltage controlled oscillator (VCO). This VCO is programmable, and is controlled by the RTL2832U over an I<sup>2</sup>C interface. After an active gain control (AGC) stage, the IF signal then needs to be brought down to baseband. The classical method of doing this is to pass the IF signal through an anti-alias filter, sample the output with an ADC, and then demodulate to baseband using quadrature NCOs (numerically controlled oscillators, i.e., a sine and a cosine

oscillating at the IF frequency). Finally (on the left hand side of the diagram) the I/Q 8 bit samples are ready to stream to MATLAB or Simulink running on the desktop.

## A SOFTWARE DEFINED RADIO DESIGN ENVIRONMENT USING THE RTL-SDR

With the capability to tune over the range of 25 MHz to 1.7 GHz, the RTL-SDR can be used to investigate, view the spectra of, and receive and decode a wide range of radio signals transmitted for various applications using different modulation methods. The actual signals available to a user will of course depend on their geographical location and the surrounding radio environment. To provide an example, from our location in central Glasgow (Scotland), we can receive, view, and variously analyze and decode a selection of RF signals including:

- FM radio stations 87.5 MHz to 108 MHz
- Aeronautical 108 MHz to 117 MHz
- Meteorological ~117 MHz
- Fixed mobile 140 MHz to 150 MHz
- Special events 174 MHz to 217 MHz
- Fixed mobile (space to Earth) 267 MHz to 272 MHz
- Fixed mobile (earth to space) 213 MHz to 315 MHz
- ISM band (short range) ~433 MHz



**Figure 4.** Using the RTL-SDR and a real-time Simulink spectrum analyzer and 2D waterfall plot to view part of a 10 MHz 4G LTE signal spectrum in the 800 MHz band, clearly showing the OFDM carriers.

- Emergency services      450 MHz to 470 MHz
- UHF TV broadcasting    470 MHz to 790 MHz
- 4G LTE and GSM bands   800 MHz to 900 MHz
- Short range devices     863 MHz to 870 MHz
- GPS systems              1227 MHz to 1575 MHz

## THE RTL-SDR LABORATORY ENVIRONMENT

In Simulink, the RTL-SDR interface support takes the form of a library block, which represents both a source for the Simulink model, and a location to set parameters supplied to the RTL-SDR hardware device. As highlighted in Fig. 3, in the Simulink dialog window three main parameters are used to configure the device: the RF center frequency,  $f_{RF}$ ; tuner gain parameters,  $K$ ; and the baseband sampling frequency,  $f_s$ . In addition, a frequency correction parameter can be used to correct for offsets due to component tolerances and frequency drift which may affect the device. (In MATLAB these parameters can be set by initializing and configuring an RTL-SDR System object.)

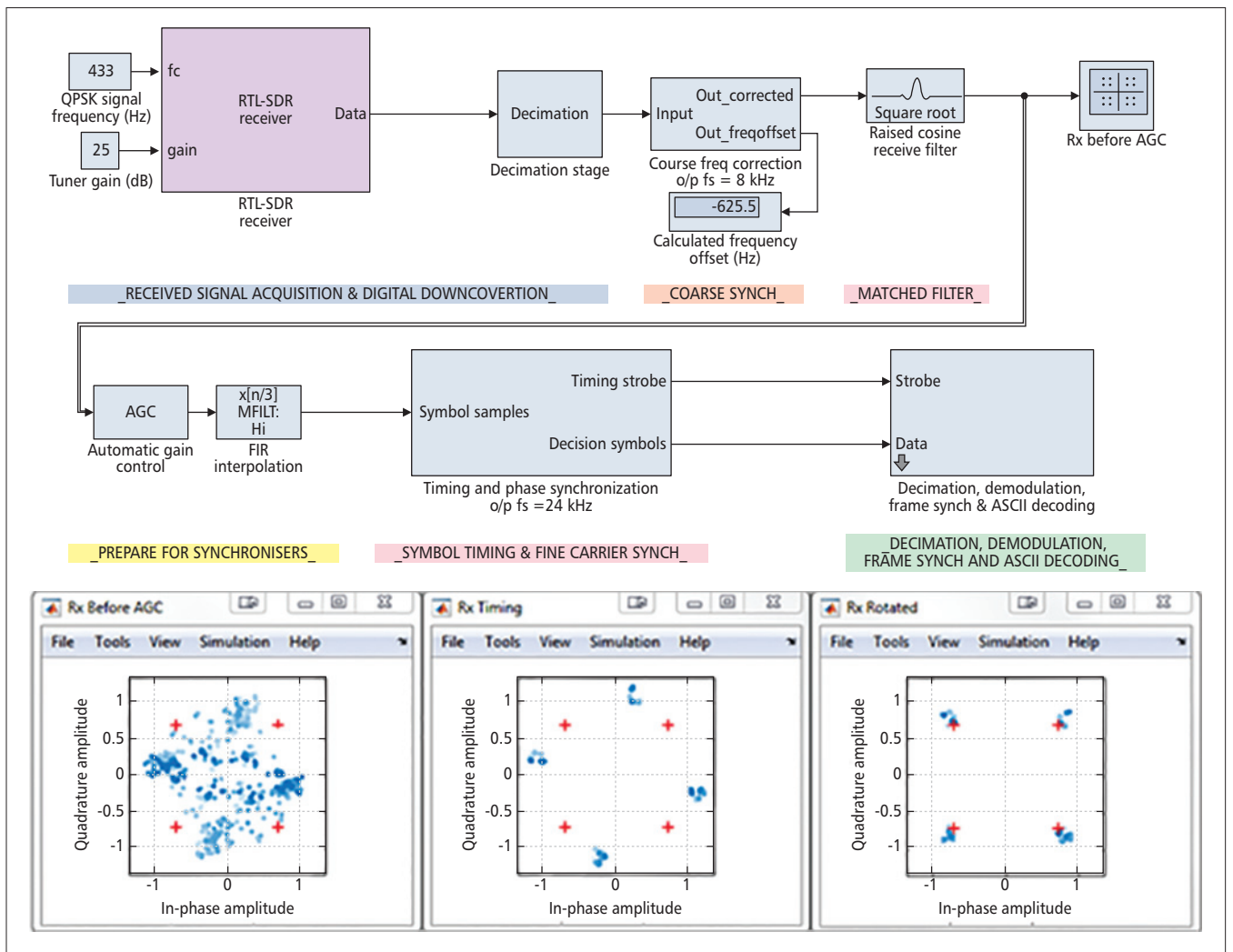
One of the first opportunities for SDR education with the RTL-SDR is simple spectrum viewing: finding and observing the frequency spectra (and in some cases the time domain representation) of some of the RF signals in the applications being broadcast around you. With knowledge (or guess work!) about available FM radio signals in the vicinity, an easy first example to run is to parameterize the RTL-SDR receiver interface block such that the device tunes to and demodulates a certain portion of the FM radio spectrum, and supplies the resulting baseband samples into the Simulink model for spectral viewing and or further processing to demodulate the signal [1]. The I/Q (complex) baseband sampling rate of the RTL-SDR has a recommended maximum of 2.8 MHz (the actual maximum is 3.2 MHz, although data loss occurs at this rate), and samples have an 8 bit resolution. Hence the maxi-

um bandwidth can be considered to be 2.8 MHz. While this bandwidth is insufficient for viewing, for example, a 5 MHz band of UMTS 3G spectra, it is more than adequate for capturing and viewing a number of other signal types, including FM signals (bandwidth = 200 kHz), keyfob signaling centered at 433 MHz, exploring the 200 kHz wide GSM channels, and so on. And while it is not possible to see the full spectrum of, for example, a 10 MHz wide LTE signal, you can easily scan over the wider band in 2 MHz sections, and observe the guard bands and spectrum edges incrementally. The output of one of the spectral viewing LTE examples from the workbook is shown in Fig. 4.

To progress to more advanced and challenging examples in the teaching lab situation, we often need signals that are locally generated and controlled. Before transmitting RF signals, however, one must be very sure that transmission of a given signal power in a particular frequency band is legal, otherwise there is a danger of being an unwelcome jammer! Recognizing our desire for a low-cost teaching and learning setup, we can generate signals locally using devices such as FM transmitters that can be plugged into smartphone headphone sockets (these devices cost less than \$20 and low-power versions are legal in many regions), or by acquiring single-chip devices such as the RT4 433 MHz device and building simple AM transmitter circuits [1].

To begin to teach more advanced digital communications using the RTL-SDR, we need to be able to generate appropriate RF signals in the lab, such as QPSK and other QAM transmissions. At the receive side, students can then design QPSK and QAM SDR receivers in MATLAB and/or Simulink, featuring numerically controlled oscillators, phase locked loops, frame synchronizers, digital receive filter chains, and other design elements. An example of such a design is shown in Fig. 5. To generate suitable signals in the laboratory we can use a programmable, transmit-capable SDR device such as the USRP, or Zynq SDR platform (featuring a Xilinx FPGA and Analog Devices FMComms card) to the class environment. In our *Information, Transmission and Security* seniors class (fourth year Bachelor) at the University of Strathclyde, our final laboratory challenge session in the Winter/Spring 2015 semester was to decode and receive a multiplex of signals consisting of two AM, two FM, and two QPSK data channels. This multiplex was transmitted in a 2 MHz band on 602 MHz (the University of Strathclyde has a UK Government Ofcom UHF white space test licence at this frequency, and hence can legally use this in the lab for test purposes). If radio transmission over the air is not practical, perhaps due to local environmental or legal concerns, then an alternative is to use a cable and MCX connectors to make an RF cable connection between the transmitter and the RTL-SDR receiver, in place of the free-space wireless channel. Figure 6 shows one of the Strathclyde students at work in a seniors (fourth year Bachelor) lab, with just a PC, software, and the RTL-SDR and simple antenna supplied with the device.





**Figure 5.** Design of a real-time RTL-SDR system receiving a QPSK signal transmitted in the laboratory from a USRP at 433 MHz and implemented from first principles in Simulink, featuring receive filters, carrier and timing synchronization, and decimation stages.

## A HANDS-ON SDR COMMUNICATIONS WORKBOOK

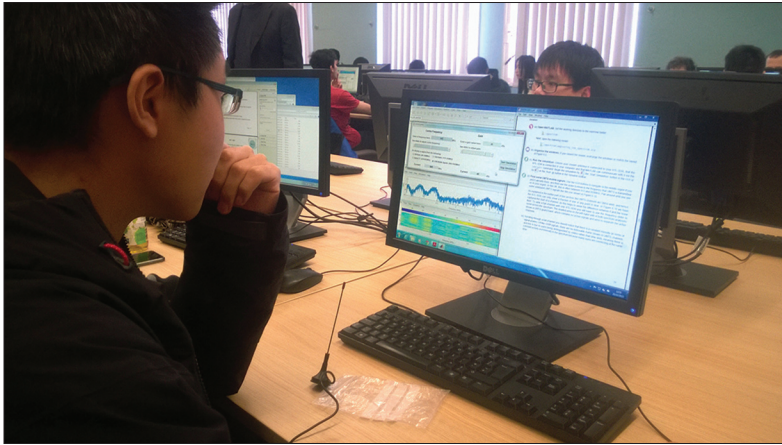
In the context of the curriculum requirements of our DSP and digital communications courses at the University of Strathclyde, the RTL-SDR has created the opportunity to invigorate our teaching and learning with real-world signals, reception, and examples. In response, we have jointly developed a complete workbook for EE students that allows them to experience and explore the radio spectrum, and to design, test, and implement radio receivers. A large selection of reference designs is provided with the workbook. MATLAB is now the de facto technical computing environment in many schools, including Strathclyde, and student familiarity with the software provides an excellent platform for developing a complete curriculum. Nevertheless, the course aims to appeal to all levels of prior experience, and includes sections for those new to the various tools and themes covered.

The open-source course materials and teaching and learning examples are all available on-

line (<http://www.desktopSDR.com>) and feature 650 pages of practical exercises (starting from first principles), descriptions and theory, and more than 120 MATLAB and Simulink example files [1]. The materials are openly available, and also likely to be of use to practicing professional engineers, the maker community, or perhaps amateur radio enthusiasts looking to learn more about SDR real-time implementation. Add-on cards are available for the RTL-SDR to upconvert short wave to frequency ranges where the RTL-SDR functions, i.e. 25 MHz and above.

Overall, the objectives of the workbook and materials are to:

- Convey the fundamental concepts and applications of SDR systems, from the RF, IF, and baseband stages of DSP enabled radio algorithms.
- Encourage an intuitive understanding of the RF spectrum, by demonstrating how to tune across the spectrum range of 25 MHz to 1.7 GHz, and to capture and view different signals in I/Q format, recognize modulation schemes, and plot live RF spectra on screen.



**Figure 6.** A software defined radio laboratory session at the University of Strathclyde to receive and decode an in-class transmitted signal. Note that the only hardware is the PC, and the RTL-SDR. Spectrum analyzer functionality is all provided by the same RTL-SDR device.

- Provide an appreciation of the different communications systems and standards in use, and the bands of RF frequencies they use, ranging from FM radio, to GSM, to ISM band and LTE.
- Demonstrate the fundamentals of the analog modulation schemes of AM and FM radio, and be able to construct real-time digital receivers for both AM and FM analog signals based around digitized I/Q SDR receiver algorithms.
- Review aspects of DSP digital receiver design (filters, demodulators, decimators, NCOs), and implement practical digital receivers from first principles.
- Consider the requirements for tuning, setting offset frequencies, carrier synchronization and phase locking, and symbol and data timing, and demonstrate how to design and implement these components as part of an SDR receiver.
- Show how to generate and transmit RF signals (using low-cost FM transmitters, USRP SDR hardware, custom designs, etc.) to build simple signaling layers and design PHY implementations to send data, music, images, and control information.

## CONCLUSIONS AND SOME NEXT STEPS

The open course materials discussed in this article create the opportunity to build and experiment with SDR techniques. Of course, with the our \$20 RTL-SDR, we do not have precise control over ADC rates, programmable RF subsystems, nor controllable wideband antennas, and hence need to deal with high noise floors and frequency drift. However, this can be turned into part of the learning experience, e.g. finding the frequency offset of a particular RTL-SDR is one of the early exercises in the workbook [4]. Also, this first course on SDR implementation is working with single-channel antenna systems. Low-cost multichannel SDR is not so far away however. In fact with the current drivers for

MATLAB and Simulink, we can host multiple RTL-SDRs (there are examples using two RTL-SDRs in a Simulink design in Chapter 3 of [1]). Therefore, multiple input desktop opportunities for students and maker communities is here. We can also expect more low-cost and accessible SDR transmitters to become available, creating more opportunities and exciting prospects for the lab, of course realizing that wherever devices are adopted, we need to be aware of the available (legal) frequency bands that we can and cannot use.

We can conclude by stating that the RTL-SDR is an excellent first SDR device that can form an IF digital radio and a front end for floating or even fixed point implementations of digital demodulators, receivers, and decoders using MATLAB and Simulink to bring SDR opportunities to the desktop. Finally, it is perhaps interesting to note that the RTL-SDR device is also currently *trending* and defining SDR in the “consumer” marketplace. In July 2015 a search on <http://www.amazon.com> for the term “Software Defined Radio” listed the RTL-SDR as the top hit, with the next three hits also being products related to the RTL-SDR! Stay tuned. The wireless (SDR-enabled) revolution is just beginning!

## REFERENCES

- [1] R. W. Stewart *et al.*, *Software Defined Radio using the MATLAB & Simulink and the RTL-SDR*, Strathclyde Academic Media, 2015. ISBN-13: 978-0-9929787-1-6.
- [2] f. harris, *Multirate Signal Processing for Communication Systems*, Prentice Hall, 2004.
- [3] M. Rice, *Digital Communications: A Discrete Time Approach*, Prentice Hall, 2008.
- [4] J. Mitola, “The Software Radio Architecture,” *IEEE Commun. Mag.*, vol. 33, no. 5, May 2015, pp. 26–38.
- [5] J. Mitola *et al.*, “Guest Editorial on Software Radios,” *IEEE JSAC*, vol. 17, no. 4, April 1999, pp 509–12.
- [6] W. Tuttlebee, (Ed.), *Software Defined Radio: Enabling Technologies*, John Wiley, 2002, ISBN 0-470-84318-7.
- [7] E. Grayver, *Implementing Software Defined Radio*, Springer, 2012, ISBN-13: 978-1441993311.
- [8] L. Goeller and D. Tate, “A Technical Review of Software Defined Radios: Vision, Reality, and Current Status,” *Proc. Military Communications Conference (MILCOM)*, 2014 IEEE, 6–8 Oct. 2014, pp. 1466–70.
- [9] S. G. Biln *et al.*, “Software-Defined Radio: A New Paradigm for Integrated Curriculum Delivery,” *IEEE Commun. Mag.*, vol. 52, no. 5, May 2014, pp. 184–93.
- [10] El-Hajjar *et al.*, “Demonstrating the Practical Challenges of Wireless Communications using USRP,” *IEEE Commun. Mag.*, vol. 52, no. 5, May 2014, pp. 194–201.
- [11] M. Petrova *et al.*, “System-Oriented Communications Engineering Curriculum: Teaching Design Concepts with SDR Platforms,” *IEEE Commun. Mag.*, vol. 52, no. 5, May 2014, pp. 202–09.
- [12] M. B. Sruthi *et al.*, “Low Cost Digital Transceiver Design for Software Defined Radio using RTL-SDR,” *Proc. Int’l Multi-Conf. Automation, Computing, Communication, Control and Compressed Sensing (iMac4s)*, 22–23 March 2013, pp. 852–55.

## BIOGRAPHIES

BOB STEWART is the MathWorks professor of signal processing at the University of Strathclyde, and is also currently the head of the Department of Electronic and Electrical Engineering. He also manages a research group working on DSP, FPGAs, white space radio, and low-cost SDR implementation. He has a bachelors and Ph.D. from the University of Strathclyde.

LOUISE CROCKETT is the Xilinx lecturer in FPGAs and programmable logic at the University of Strathclyde in Glasgow. She is also the principal author of *The Zynq Book*, published in 2014, and has a core research interest in

---

FPGA systems design for software defined radio and DSP systems. She has a masters and Ph.D. degree from the University of Strathclyde.

DALE ATKINSON is a Ph.D. student at the University of Strathclyde, working on SDR receiver systems for low-cost implementation. He received a bachelor's degree from the University of Strathclyde in 2014.

KENNETH BARLEE is a Ph.D. student at the University of Strathclyde, working on novel DSP enabled radio algorithms and implementation for software defined radio. He received a bachelor's degree from the University of Strathclyde in 2014.

DAVID CRAWFORD is the manager of the Centre for Wireless White Space at the University of Strathclyde, and also lectures on digital signal processing, using the RTL-SDR in the laboratory sessions. He was the managing director at EPSON Semiconductor before joining Strathclyde in 2010. He has four degrees from the University of Strathclyde: a bachelor's, masters, MBA, and a Ph.D.

IAIN CHALMERS is a Ph.D. student at the University of Strathclyde, working on wireless white space radio architectures. He received a master's degree from the University of Strathclyde in 2012. Previously he studied at California State University during his master's degree.

MIKE McLERNON is a development manager at MathWorks Inc. in Natick, MA. He leads software development activity on communications software products, with a particular interest in filter receiver design, SDR, standards-based modeling, and channel modeling. He has a master's degree from Rensselaer Polytechnic Institute, and a bachelor's from the University of Virginia.

ETHEM SOZER is a principal software engineer at MathWorks Inc. in Natick, MA. He specializes in software development for signal processing and communications toolboxes to support SDR. Previously he was a research engineer at Massachusetts Institute of Technology, where he developed underwater acoustic communication hardware and software platforms. He has bachelor's and master's degrees from Middle East Technical University, and a Ph.D. from Northeastern University.

## SOFTWARE DEFINED 5G NETWORKS FOR ANYTHING AS A SERVICE



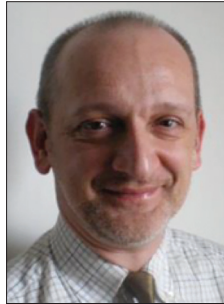
David Soldani



Bernard Barani



Rahim Tafazolli



Antonio Manzalini



Chih-Lin I

The advanced fifth generation (5G) infrastructure is expected to become the “nervous system” of the digital society, digital economy, and silver economy. New service paradigms such as “immersive experience” and “anything as a service” (XaaS) everywhere are envisioned as among the primary drivers for global adoption and market uptake of new 5G technology components. Above all, 5G networks will support mission-critical machine communications and massive machine type of traffic. As a result, the key performance metrics that 5G is expected to improve are in terms of, but not limited to, *latency, reliability, capacity, and spectrum and network agility*. This calls for a complete rethinking of all functional domains, including access stratum (AS), non-access stratum (NAS), and transport network layer (TNL), in terms of protocols and procedures. New emerging technologies, such as software defined networking (SDN), network functions virtualization (NFV), mobile edge computing (MEC), and high-performance computing (HPC), provide momentum for new design principles toward software (service) defined 5G networks, targeting a software defined air interface (SDAI) for available bands (spectrum); sliced “networks on demand” for multiple industries, especially for vertical markets (new architecture); and flexibility and spectral efficiency for mobile broadband and machine type communications (new air interface).

The main challenge is to enable dynamic slicing of the networks, end to end, to make specific physical network infrastructures capable of supporting a much wider range of requirements with the ability to create dedicated networks that business customers can control, for instance, through self-service enterprise portals. A core element of the solution is the capability of flexibly adapting the infrastructure (including the air interface and supporting networks in the fronthaul, backhaul, aggregation, and core) to different “vertical” application requirements. The goal is that, 10 years from now, communication and information technologies will be integrated in common high performing platforms and provide a malleable service defined 5G

infrastructure, with seamless integration of heterogeneous wired and wireless capabilities, and powering business business solutions, while offering multi-tenant technical and commercial control.

In *IEEE Communications Magazine*, this timely Feature Topic brings together key contributions of researchers from industry and academia, which address the above challenging issues and sheds light on some fundamental technology aspects of the advanced 5G network infrastructure.

In response to the Call for Papers, a large number of submissions were received. The submissions underwent a rigorous review process, following which only five outstanding contributions were selected for publication. The five articles are in the fields of multiple access, front/backhauling and networking technologies, and architectures for 5G. These articles are expected to stimulate new ideas and contributions within the research and innovation community, in addition to providing readers with relevant background information and feasible solutions to the main technical design issues of future 5G networks.

The first article in this Feature Topic, “Non-Orthogonal Multiple Access for 5G: Solutions, Challenges, Opportunities, and Future Research Trends,” is by Linglong Dai, Bichai Wang, Yifei Yuan, Shuangfeng Han, Chih-Lin I, and Zhaocheng Wang. The authors present one potential technology for 5G, denoted as non-orthogonal multiple access (NOMA), and the concept of software defined multiple access (SoDeMA), which enables adaptive configuration of available multiple access schemes to support diverse services and applications, and where any of the multiple access schemes may be provisioned based on stakeholders’ needs.

The second article, “Rethink Fronthaul for Soft RAN,” by Chih-Lin I, Yannan Yuan, Jinri Huang, Shijia Ma, Chunfeng Cui, and Ran Duan, proposes and analyzes a new Ethernet-based fronthaul (FH) interface, called next-generation fronthaul interface (NGFI), which better sup-

ports some of the key 5G technologies, such as cloud RAN, UDN, and large-scale antenna systems.

In the third article, “Baseband Unit Cloud Interconnection Enabled by Flexible Grid Optical Networks with Software Defined Elasticity,” Jiawei Zhang, Yuefeng Ji, Jie Zhang, Rentao Gu, Yongli Zhao, Siming Liu, Kun Xu, Mei Song, Han Li, and Xinbo Wang propose a software defined centralized control plane to coordinate heterogeneous resources from baseband units (BBUs), radio, and the optical domain, and improve elasticity and intelligence of the cloud-based radio access network (C-RAN).

The fourth article, “Network Coded Software Defined Networking: Enabling 5G Transmission and Storage Networks,” by Jeppe Krigslund, Jonas Hansen, Daniel E. Lucani, Frank H. P. Fitzek, and Muriel Médard, argues for the use of SDN along with network coding (NC) functionalities to enable more 5G network services.

The fifth article, “Software-Defined Service Migration through Legacy Service Integration into 4G Networks and Future Evolutions,” by Yeunwoong Kyung, Tri M. Nguyen, Kiwon Hong, Jongkwan Park, and Jinwoo Park, closes this Feature Topic with a proposal on how to design a next generation network architecture that integrates legacy network services into 4G networks through SDN and NFV software-based characteristics.

In closing, we would like to thank all the stakeholders who have made this Feature Topic possible and hope it meets readers’ expectations, for whom this Feature Topic on the advanced 5G infrastructure has been prepared.

## BIOGRAPHIES

DAVID SOLDANI (david.soldani@huawei.com) received a M.Sc. degree magna cum laude approbatur in electronic engineering from the University of Florence, Italy, in 1994; and a D.Sc. degree in technology with distinction from Aalto University, Finland, in 2006. He is currently vice president of Strategic Research and Innovation in Huawei. Prior to that, he was a scientific director at Nokia and Nokia Siemens Networks.

BERNARD BARANI graduated from the École Nationale Supérieure des Télécommunications de Bretagne in 1982. He served as a communications engineer in industry and with the European Space Agency on advanced satcom programs. He is currently deputy head of the Network Technologies unit in the European Commission, responsible for the definition of research related to future networks (5G) and systems based on connected objects (IoT) in support of EU industry competitiveness.

RAHIM TAFAZOLLI is director of the Institute for Communication Systems (ICS) and the 5G Innovation Centre (5GIC), University of Surrey. He has published more than 500 research papers in refereed journals and international conferences, and as an invited keynote speaker. He is the Editor of two books, including *Technologies for Wireless Future*. He established the first in the world innovation center on 5G in 2012 and launched the first International Forum on IoT.

ANTONIO MANZALINI received his M.Sc. degree in electronic engineering from the Politecnico of Turin. In 1990 he joined Telecom Italia, starting with activities on optical transport networks. He has been involved in several standardization bodies and European projects. His achievements have been published in more than 100 publications. He is Chair of the IEEE Initiative on SDN. He is currently senior manager at the Strategy an Innovation Department (Future Centre) of Telecom Italia.

CHIH-LIN I received her Ph.D. degree in electrical engineering from Stanford University. As China Mobile chief scientist of wireless technologies, she launched 5G technologies R&D in 2011, and is also leading C-RAN, green communications, and applications initiatives. She received the *IEEE Transactions on Communications* Stephen Rice Best Paper Award, and has been an IEEE ComSoc Board Member, ComSoc M&C Board Chair, and WCNC Steering Committee Founding Chair. She is on the Executive Board of Green-Touch, ETSI/NFV NOC, WWRF Steering Board, and SAB of Singapore NRF.

# Non-Orthogonal Multiple Access for 5G: Solutions, Challenges, Opportunities, and Future Research Trends

Linglong Dai, Bichai Wang, Yifei Yuan, Shuangfeng Han, Chih-Lin I, and Zhaocheng Wang

## ABSTRACT

The increasing demand of mobile Internet and the Internet of Things poses challenging requirements for 5G wireless communications, such as high spectral efficiency and massive connectivity. In this article, a promising technology, non-orthogonal multiple access (NOMA), is discussed, which can address some of these challenges for 5G. Different from conventional orthogonal multiple access technologies, NOMA can accommodate much more users via non-orthogonal resource allocation. We divide existing dominant NOMA schemes into two categories: power-domain multiplexing and code-domain multiplexing, and the corresponding schemes include power-domain NOMA, multiple access with low-density spreading, sparse code multiple access, multi-user shared access, pattern division multiple access, and so on. We discuss their principles, key features, and pros/cons, and then provide a comprehensive comparison of these solutions from the perspective of spectral efficiency, system performance, receiver complexity, and so on. In addition, challenges, opportunities, and future research trends for NOMA design are highlighted to provide some insight on the potential future work for researchers in this field. Finally, to leverage different multiple access schemes including both conventional OMA and new NOMA, we propose the concept of software defined multiple access (SoDeMA), which enables adaptive configuration of available multiple access schemes to support diverse services and applications in future 5G networks.

## INTRODUCTION

In the history of wireless communications from the first generation (1G) to 4G, the multiple access scheme has been the key technology to distinguish different wireless systems. It is well known that frequency-division multiple access (FDMA) for 1G, time-division multiple access (TDMA) mostly for 2G, code-division multiple

access (CDMA) for 3G, and orthogonal frequency-division multiple access (OFDMA) for 4G are primarily orthogonal multiple access (OMA) schemes. In these conventional multiple access schemes, different users are allocated to orthogonal resources in either the time, frequency, or code domain in order to avoid or alleviate inter-user interference. In this way, multiplexing gain can be achieved with reasonable complexity.

However, the fast growth of mobile Internet has propelled 1000-fold data traffic increase by 2020 for 5G. Hence, the spectral efficiency becomes one of the key challenges to handle such explosive data traffic. Moreover, due to the rapid development of the Internet of Things (IoT), 5G needs to support massive connectivity of users and/or devices to meet the demand for low latency, low-cost devices, and diverse service types. To satisfy these requirements, enhanced technologies are necessary. So far, some potential candidates have been proposed to address challenges of 5G, such as massive MIMO, millimeter wave communications, ultra dense network, and non-orthogonal multiple access (NOMA) [1]. In this article, we focus on NOMA, which is highly expected to increase system throughput and accommodate massive connectivity. Note that Third Generation Partnership Project (3GPP) Long Term Evolution (LTE) Rel-13 is doing ongoing studies toward NOMA in the form of multi-user superposition transmission (MUST). NOMA allows multiple users to share time and frequency resources in the same spatial layer via power domain or code domain multiplexing. Recently, several NOMA schemes have attracted lots of attention, and we can generally divide them into two categories,<sup>1</sup> that is, power domain multiplexing [2–4] and code domain multiplexing, including multiple access with low-density spreading (LDS) [5, 6], sparse code multiple access (SCMA) [7], multi-user shared access (MUSA) [8], and so on. Some other multiple access schemes such as pattern-division multiple access (PDMA) and bit division multiplexing (BDM) [9] are also proposed. Key features and advantages of NOMA are discussed

Linglong Dai, Bichai Wang, and Zhaocheng Wang are with Tsinghua University.

Yifei Yuan is with ZTE Corporation.

Shuangfeng Han and Chih-Lin I are with China Mobile Research Institute.

This work was supported in part by the International Science & Technology Cooperation Program of China (Grant No. 2015DFG12760), the National Natural Science Foundation of China (Grant Nos. 61571270 and 61201185), and the Beijing Natural Science Foundation (Grant No. 4142027)..

<sup>1</sup> Note that “NOMA” is also used by NTT DoCoMo to refer to NOMA via power domain multiplexing.

later. The design principles, key features, advantages and disadvantages of existing dominant NOMA schemes are discussed and compared. More importantly, although NOMA can provide attractive advantages, some challenging problems should be solved, such as advanced transmitter design and the trade-off between performance and receiver complexity. Thus, opportunities and research trends are highlighted to provide some insights on the potential future work for researchers in this field. In addition, unlike the conventional way of designing a specific multiple access scheme separately and individually, we propose the concept of software defined multiple access (SoDeMA), in which several candidates among multiple access schemes can be adaptively configured to satisfy different requirements of diverse services and applications in future 5G networks. Finally, conclusions are drawn.

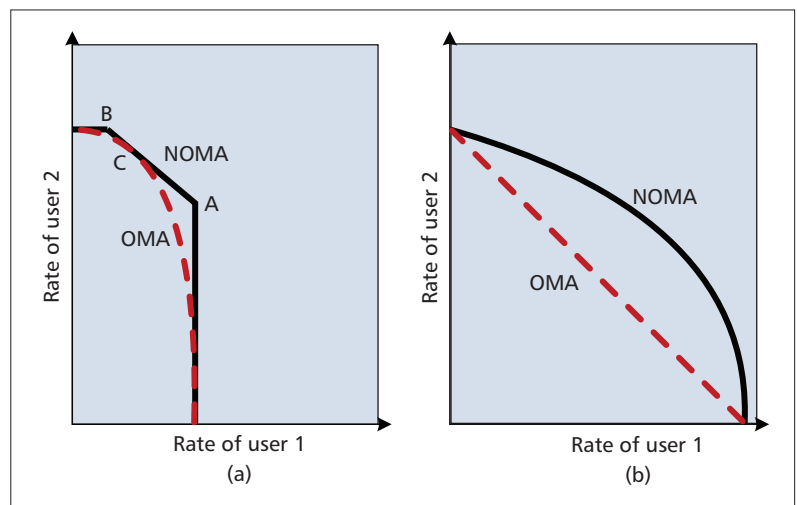
## FEATURES OF NOMA

In conventional OMA schemes, multiple users are allocated with radio resources which are orthogonal in time, frequency, or code domain. Ideally, no interference exists among multiple users due to the orthogonal resource allocation in OMA, so simple single-user detection can be used to separate different users' signals. Theoretically, it is known that OMA cannot always achieve the sum-rate capacity of multiuser wireless systems [10]. Apart from that, in conventional OMA schemes, the maximum number of supported users is limited by the total amount and the scheduling granularity of orthogonal resources.

Recently, NOMA has been investigated to deal with the problems of OMA as mentioned above. Basically, NOMA allows controllable interferences by non-orthogonal resource allocation with the tolerable increase in receiver complexity. Compared to OMA, the main advantages of NOMA include the following.

**Improved spectral efficiency:** According to the multi-user capacity analysis in the pioneering work [10], Fig. 1 shows the channel capacity comparison of OMA and NOMA, where two users in the additive white Gaussian noise (AWGN) channel are considered as an example without loss of generality. Figure 1a shows that the uplink NOMA is able to achieve the capacity bound, while OMA schemes are in general sub-optimal except at point C. However, at this optimal point, the user throughput fairness is quite poor when the difference of the received powers of the two users is significant, as the rate of the weak user is much lower than that of the strong user. In the downlink, Fig. 1b shows that the boundary of rate pairs of NOMA is outside of the OMA rate region in general. In multi-path fading channels with intersymbol interference (ISI), although OMA could achieve the sum capacity in the downlink, NOMA is optimal while OMA is strictly suboptimal if channel state information (CSI) is only known at the mobile receiver [10].

**Massive connectivity:** The non-orthogonal resource allocation in NOMA indicates that the number of supported users or devices is not



**Figure 1.** Channel capacity comparison of OMA and NOMA in an AWGN channel: a) uplink AWGN channel; b) downlink AWGN channel.

strictly limited by the amount of available resources and their scheduling granularity. Therefore, NOMA can accommodate significantly more users than OMA by using non-orthogonal resource allocation; for example, MUSA can still achieve reasonably good performance when the overloading is 300 percent [8].

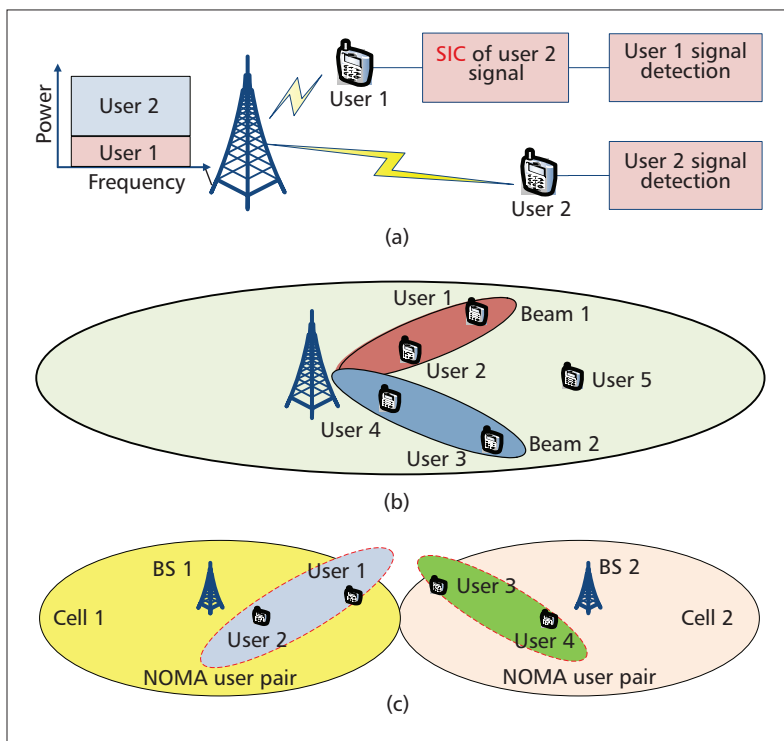
### Low transmission latency and signaling cost:

In conventional OMA with grant-based transmission, a user has to send a scheduling request to the base station (BS) at first. Then, based on the received request, the BS performs scheduling for the uplink transmission and sends a grant over the downlink channel. This procedure results in large latency and high signaling cost, which becomes worse or even unacceptable in the scenario of massive connectivity anticipated for 5G. In contrast, such dynamic scheduling is not required in some uplink schemes of NOMA, rendering a grant-free uplink transmission that can drastically reduce the transmission latency and signaling overhead.

Due to the potential advantages above, NOMA has been actively investigated as a promising technology for 5G. In the next section, existing dominant NOMA schemes are discussed and compared in detail.

## DOMINANT NOMA SOLUTIONS

In this section, we discuss dominant NOMA schemes by grouping them into two categories: power domain multiplexing and code domain multiplexing. Power domain multiplexing means that different users are allocated different power levels according to their channel conditions to obtain the maximum gain in system performance. Such power allocation is also beneficial to separate different users, where successive interference cancellation (SIC) is often used to cancel multi-user interference. In this article, power domain multiplexing is applied only to downlink NOMA. Code domain multiplexing is similar to CDMA or multicarrier CDMA (MC-CDMA), that is, different users are assigned different codes, and are then multiplexed over the same time-frequency resources. The difference



**Figure 2.** Illustration of NOMA via power domain multiplexing: a) basic NOMA with a SIC receiver; b) NOMA in MIMO systems; c) network NOMA.

between power domain multiplexing and code domain multiplexing is that code domain multiplexing can achieve certain spreading gain and shaping gain at the cost of increased signal bandwidth.

### NOMA VIA POWER DOMAIN MULTIPLEXING

**Basic NOMA with a SIC Receiver:** Figure 2a illustrates the basic NOMA scheme via power domain multiplexing with a SIC receiver in the downlink. Note that this NOMA scheme can also be applied in the uplink [2]. At the BS transmitter, signals for different users are linearly added up under certain power partitions to balance the sum rate of all multiplexed users and the throughput fairness among individual users.

At the receiver, SIC is commonly used to realize multi-user detection (MUD). Due to the near-far effect, the channel conditions may vary significantly among users. SIC is performed at users with relatively high signal-to-interference-plus-noise ratio (SINR), and should be carried out in descending order of SINR.

As we can see, the basic form of NOMA with SIC exploits SINR difference among users, either due to the natural near-far effect or by non-uniform power allocation at the transmitter. A similar scheme can be used for uplink to increase the uplink system capacity.

**NOMA in Massive MIMO Systems:** NOMA can be used in conjunction with multi-user multiple-input multiple-output (MU-MIMO) to further improve the system spectral efficiency [3]. As illustrated in Fig. 2b, multiple transmit antennas at a BS are used to form different beams in

the spatial domain, where each beam adopts the basic NOMA discussed above.

At the receiver, the inter-beam interference can be suppressed by spatial filtering [3], and then intra-beam SIC can be used to remove the inter-user interference. The extension of NOMA in massive MIMO systems can further improve spectral efficiency.

**Network NOMA:** When transmit power allocation is biased toward far away users in downlink NOMA, cell edge users experience increased interference from neighboring cells. As an example, a cellular system with two cells and four users is depicted in Fig. 2c, where a two-user NOMA scheme is assumed: user 1 and user 2 are served by BS 1, while user 3 and user 4 are served by BS 2. Strong interference is expected between users 1 and 3, which may degrade the performance of network NOMA, that is, multi-cell NOMA.

To mitigate the inter-cell interference, joint precoding of NOMA users' signals across neighboring cells can be considered. This requires that all users' data and CSI should be available at multiple BSs, but finding the optimal precoder is not trivial. Moreover, the multi-user precoding used for single-cell NOMA maybe not be feasible for the network NOMA scenario, since the precoder for geographically separate BS antennas does not actually form the physical beam that can readily be used for intra-beam NOMA. Based on the fact that large-scale fading would be quite different between the links to different cells, a complexity-reduced precoding scheme for network NOMA has been proposed in [4], where the multi-cell joint precoder is applied only to cell edge users (e.g., user 1 and user 3 as shown in Fig. 2c).

### NOMA VIA CODE DOMAIN MULTIPLEXING

**Low-Density Spreading CDMA:** The idea behind LDS-CDMA is to use sparse spreading sequences instead of dense spreading sequences in conventional CDMA [5] to reduce the interference at each chip. Therefore, LDS-CDMA can improve system performance by exploiting LDS sequences in CDMA [5], which is the key feature distinguishing conventional CDMA and LDS-CDMA. In this way, interference will be efficiently decreased among multiple users with appropriate spreading sequence design, and overloading can be achieved.

At the receiver, a message passing algorithm (MPA) can be used to realize MUD. MPA is very efficient for the factor graph [11], which is a bipartite graph including variable nodes and factor nodes as illustrated in Fig. 3. Messages are passed among variable nodes and factor nodes over edges, which can be interpreted as the soft-values that represent the reliability of the symbol associated with each edge. The marginal distribution of a variable node can be regarded as a function of the messages received by that node [11].

**Low-Density Spreading OFDM:** LDS orthogonal frequency-division multiplexing (LDS-OFDM) can be considered as a combined version of LDS-CDMA and OFDM, in which



the chips are subcarriers of OFDM in order to combat the multipath fading. In LDS-OFDM, the transmitted symbols are first mapped to certain LDS sequences, and then transmitted on different OFDM subcarriers. The number of symbols can be greater than the number of subcarriers, that is, overloading is allowed to improve spectral efficiency [6]. MPA in LDS-CDMA can also be used in an LDS-OFDM receiver. Essentially, LDS-OFDM can be viewed as an improved form of multi-carrier CDMA (MC-CDMA) by replacing the dense spreading sequences with LDS.

**Sparse Code Multiple Access:** The recently proposed SCMA [7] is an enhanced version of LDS-CDMA. Unlike LDS-CDMA, SCMA directly maps different bitstreams to different sparse codewords, as illustrated in Fig. 4, where each user has a predefined codebook (there are 6 users in Fig. 4). All codewords in the same codebook contain zeros in the same two dimensions, and the positions of zeros in different codebooks are distinct to facilitate the collision avoidance of any two users. For each user, two bits are mapped to a complex codeword. Codewords for all users are multiplexed over four shared orthogonal resources (e.g., OFDM subcarriers).

The key difference between LDS-CDMA and SCMA is that a multi-dimensional constellation for SCMA is designed to generate codebooks, which brings the “shaping gain” that is not possible for LDS [7]. Here, “shaping gain” is the gain in the average symbol energy when the shape of a constellation is changed. In general, the shaping gain is higher when the shape of a constellation is closer to a sphere, and the maximum achievable shaping gain by the optimization of a multi-dimensional constellation is 1.53 dB [7]. For the concatenated approach in high modulation order, the multi-dimensional constellation can be optimized to obtain shaping gain, and then codebooks are generated based on the multi-dimensional constellation [7]. The SCMA codebook design is a complicated problem, since different layers are multiplexed with different codebooks. As the appropriate design criterion and specific solution to the multi-dimensional problem are still unknown, a multi-stage approach has been proposed to realize a suboptimal solution [7]. Specifically, an  $N$ -dimensional complex constellation with  $M$  points (which is called the mother constellation) is first optimized to improve the shaping gain, and then some codebook-specific operations are performed to the mother constellation to generate the  $N$ -dimensional constellation for each codebook. Three typical operations are phase rotation, complex conjugate, and dimensional permutation of the constellation [7]. In the generated  $N$ -dimensional constellations after codebook-specific operations, each  $N$ -dimensional constellation point is multiplied with a projection matrix to generate a  $K$ -dimensional codeword ( $K \gg N$ ), which has  $N$  non-zero elements from the components of the  $N$ -dimensional constellation point. In this way, codebooks with  $M$  codewords can be obtained. Readers can find more details in [7].

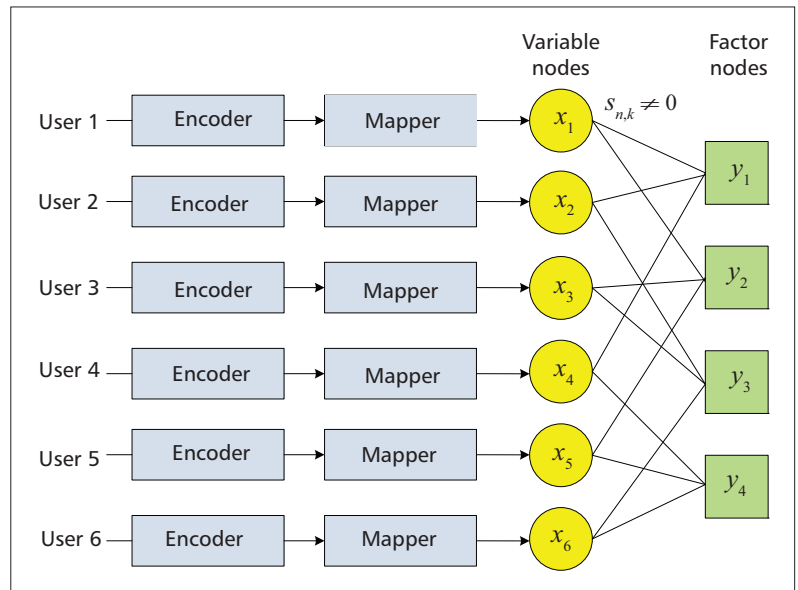


Figure 3. Illustration of LDS-CDMA: six users and four chips for transmission, which means 150 percent overloading.

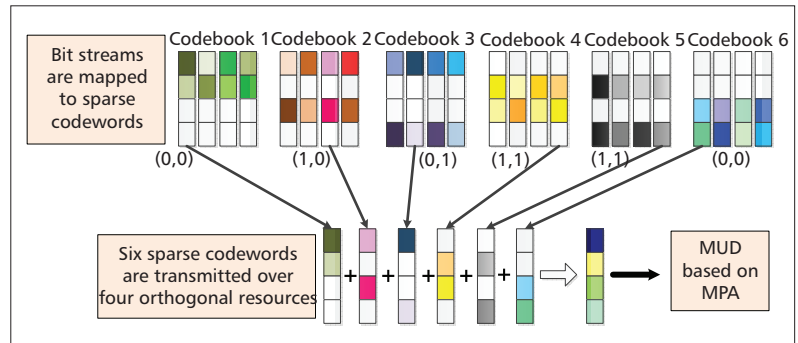


Figure 4. SCMA encoding and multiplexing.

**Multi-User Shared Access:** In the uplink MUSA system shown in Fig. 5 [8], symbols of each user are spread by a spreading sequence. Multiple spreading sequences constitute a pool from which each user can randomly pick one of the sequences. Note that for the same user, different spreading sequences may also be used for different symbols, which may further improve the performance via interference averaging. Then all spreading symbols are transmitted over the same time-frequency resources. The spreading sequences should have low cross-correlation and can be  $M$ -ary. At the receiver, codeword-level SIC is used to separate data from different users. The complexity of codeword-level SIC is less of an issue in the uplink as in any case the receiver needs to decode the data for all users. The only noticeable impact on the receiver implementation would be that the pipeline of processing may be changed in order to perform SIC operation. The difference between MUSA and MC-CDMA is that MUSA assumes that it is basically synchronous when users’ signals arrive at the BS, which is easier to realize SIC, while MC-CDMA does not require this synchronization in the uplink. In addition, MUSA uses non-binary spreading sequences, while binary

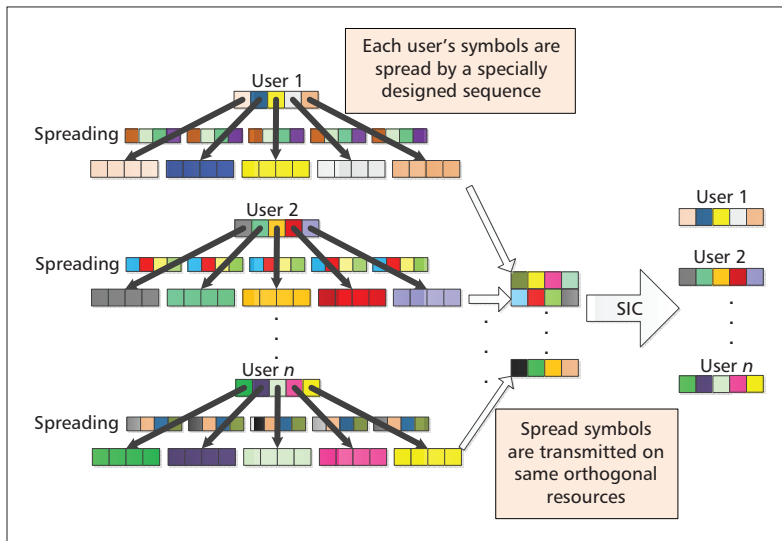


Figure 5. Uplink MUSA system.

spreading sequences are usually considered in classical MC-CDMA systems.

In downlink MUSA, users are separated into  $K$  groups. In each group, different users' symbols are mapped to different constellations in a way that can ensure Gray mapping in the combined constellation of superposed signals. The combined constellation is determined not only by the modulation order of each user, but also by the transmit power partition among multiplexed users. Orthogonal sequences can be used to spread the superposed symbols to get time or frequency diversity gain. Gray mapping of the combined constellation reduces the reliance on advanced receivers, so less processing-intensive receivers such as symbol-level SIC can be used.

#### OTHER NOMA SCHEMES

In addition to the power domain multiplexing and code domain multiplexing discussed above, a few other NOMA schemes are currently being investigated. Pattern-division multiple access (PDMA) is a NOMA scheme that can be realized in several domains. At the transmitter, PDMA uses non-orthogonal patterns, which are designed by maximizing the diversity and minimizing the overlaps among multiple users. Then the actual multiplexing can be carried out in the code domain, spatial domain, or a combination of the two. Multiplexing in the code domain corresponds to the case of successive interference cancellation amenable multiple access (SAMA) [12], which is similar to LDS-CDMA, with LDS sequences being replaced by non-orthogonal patterns. Hence, MPA can also be used for the sequence detection in PDMA. The multiplexing in the spatial domain, called spatial PDMA, requires multiple antennas at the BS. The diversity of PDMA can come from multiple transmit antennas, which is preferred for macrocell deployment. Different from multi-user MIMO, precoding is not needed in spatial PDMA since the aim is to increase the spatial diversity rather than spectral efficiency. PDMA can be used in both downlink and uplink transmissions.

Bit-division multiplexing (BDM) [9] is another form of NOMA particularly useful for downlink transmission. Its basic concept is based on hierarchical modulation, and the resources of multiplexed users are partitioned at the bit level. Although strictly speaking the resource allocation of BDM is orthogonal in the bit domain, multi-user signals share the same constellation (e.g., superposed in the modulation symbol domain).

Some other NOMA schemes were also proposed, such as interleave-division multiple access (IDMA), which performs interleaving of chips after symbols are multiplied by spreading sequences. As shown in [13], compared to CDMA, IDMA is able to achieve an  $E_b/N_0$  gain of about 1 dB when bit error rate (BER) performance of  $10^{-3}$  is considered in highly loaded systems with 200 percent overloading.

In many of the NOMA schemes mentioned above, especially when used for grant-free uplink transmission, there is an issue that the users' activity or instantaneous system loading is not readily known to the receiver. This would have a negative impact on the performance. Compressive sensing (CS) is a promising technique to estimate the resource occupancy. Some work on CS-based random access has been carried out recently, such as compressive random access [14].

#### COMPARISON OF NOMA SOLUTIONS

From the theoretical perspective, code-domain NOMA can obtain spreading gain due to the use of spreading sequences or codewords, which can be achieved only in the case that there is no CSI at the transmitter. Spreading gain is similar to that in CDMA, that is, the transmitted bandwidth can be spread by spreading sequences or codewords, and thus, according to Shannon's equation, signals can still be transmitted with the same capacity even when signal-to-noise ratio (SNR) is low. The spreading gain can be calculated by  $10 \log N$ , where  $N$  is the spreading factor. However, introducing redundancy through spreading will affect the system spectral efficiency [15]. In addition, SCMA can achieve extra "shaping gain" due to the optimization of multi-dimensional constellation [7].

We also compare these NOMA schemes in terms of the computational complexity of the multi-user signal detection algorithm. In power-domain NOMA, SIC is the key method for multi-user interference cancellation with complexity  $\mathcal{O}(K^3)$ , where  $K$  is the number of users. Therefore, the complexity of SIC is much less than that of the optimal maximum likelihood (ML) detection, whose complexity  $\mathcal{O}(|\mathbb{X}|^K)$  increases exponentially with the number of users  $K$ , where  $|\mathbb{X}|$  denotes the cardinality of the constellation set  $\mathbb{X}$ . On the other hand, in code-domain NOMA like LDS-CDMA, LDS-OFDM, and SCMA, spreading sequences or codebooks should be known at the receiver to realize MUD, and the complexity of the MPA-based receiver is  $\mathcal{O}(|\mathbb{X}|^w)$ , where  $w$  is the maximum number of nonzero signals superimposed on each chip or subcarrier. Thus, an MPA-based receiver usually has higher complexity than a SIC-based receiver as  $w$  is usually larger than 3 in typical 5G systems with massive connectivity.

## CHALLENGES, OPPORTUNITIES, AND FUTURE RESEARCH TRENDS

### THEORETICAL ANALYSIS OF ACHIEVABLE RATE AND OVERLOADING BOUNDS

In NOMA schemes, theoretical analysis is required to provide some insights for system design. Achievable rate of multiple access is a key metric of system performance. The achievable rate of code-domain NOMA with LDS needs to be studied, and can refer to the analytical approach of MC-CDMA. Particularly, due to the special structure of spreading sequences, some approximations can be used to simplify the calculation. It is expected to derive the closed-form expression to reveal the relationship between the achievable rate and LDS parameters such as sequence sparsity and overloading factor. Such theoretical results can shed light on how to design the system parameters according to the specific application requirements.

On the other hand, the interference cancellation capability and the affordable complexity at the receiver play an important role in the overall performance, for example, the maximum overloading factor that the system can support.

### DESIGN OF SPREADING SEQUENCES OR CODEBOOKS

In LDS systems, due to non-orthogonal resource allocation, interference exists among multiple users. A factor graph in MPA should be optimized to get good trade-off between overloading factor and receiver complexity.

In addition, it has been proved that MPA can obtain the exact marginal distribution with a cycle-free factor graph and the precise solution with a local tree-like factor graph. Graph theory can be used to design a cycle-free or local tree-like factor graph in NOMA without compromising spectral efficiency. In addition, the matrix design principle and methods in low-density-parity check (LDPC) can be considered when designing the factor graph for NOMA.

### RECEIVER DESIGN

For an MPA-based receiver, the complexity may still be high for massive connectivity in 5G. Therefore, simplified improvement of MPA can be used to reduce receiver complexity, such as Gaussian approximation of interference (GAI), which models the interference-plus-noise as Gaussian distributed, and such approximation tends to be more accurate as the amount of connectivity becomes larger in 5G. In addition, MPA can be used to jointly detect and decode the received symbols, in which the constructed graph consists of variable nodes, observation nodes, and check nodes corresponding to the check equations of the LDPC code. In this way, intrinsic information between the decoder and the demodulator can be used more efficiently to improve the detector's performance.

For a SIC-based receiver, error propagation may degrade the performance of some users. Therefore, at each stage of SIC, some nonlinear detection algorithms with higher detection accuracy can be considered to suppress the error propagation.

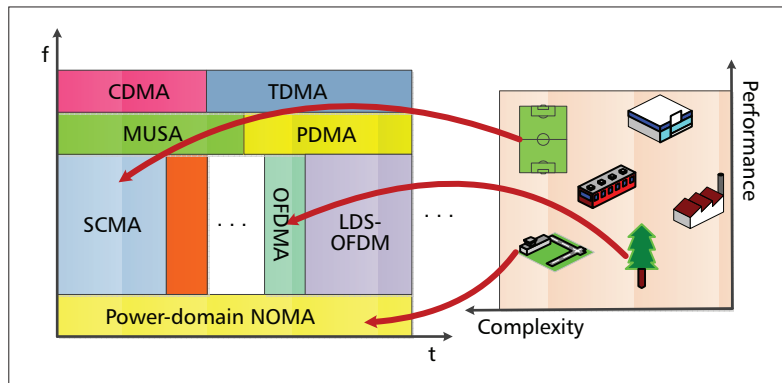


Figure 6. Illustration of the concept of software defined multiple access.

### OTHER CHALLENGES

There are also some other engineering aspects of NOMA, including reference signal design, channel estimation, and CSI feedback mechanism that can deliver robust performance when cross-user interference is severe, resource allocation signaling that can support different transmission modes for NOMA, extension to MIMO (especially massive MIMO) that can reap the performance benefits of both NOMA and multi-user MIMO, peak-to-average-power ratio (PAPR) reduction in multi-carrier NOMA, system scalability that can support different traffic loading and radio environment, and so on. These challenges need to be addressed before NOMA becomes part of 5G standards in the future.

### THE CONCEPT OF SODEMA

As discussed above, NOMA can be used for capacity improvement and massive connections in 5G. However, this does not mean that conventional OMA schemes will be completely replaced by NOMA in future 5G networks. For example, when the number of users is small and the near-far effect is not significant, such as in the case of small cells, OMA would be a better choice. In this sense, both OMA and NOMA will coexist in 5G to fulfill diverse requirements of different services and applications.

To this end, we borrow the idea of software defined radio (SDR) for multiple access design to propose the SoDeMA concept for 5G as shown in Fig. 6, where different NOMA schemes can coexist in a system assuming all of them will be specified in 5G standards. SoDeMA provides a very flexible configuration of multiple access schemes to support different services and applications in 5G. For example, for cell-center users or real-time services like ultra-high-definition video, conventional OMA schemes can be adopted to support high data rate transmission, which capitalizes on the orthogonality and synchronization. On the other hand, when high spectral efficiency, massive connectivity, and frequent access of small packets are required in some practical scenarios (e.g., dense population areas and mobile social applications), NOMA schemes can be selected. Moreover, different NOMA or OMA schemes have their own appropriate application situations, and can be adaptively configured to realize the trade-off between

Compared with conventional OMA, NOMA allows controllable interferences to realize overloading at the cost of tolerate increase of receiver complexity. Therefore, the demands of spectral efficiency and massive connectivity for 5G can be partially fulfilled by NOMA.

performance and implementation complexity. For instance, if a large difference among users' channel conditions exists due to the near-far effect or in moving networks, power-domain NOMA with a SIC receiver can be used with relatively low complexity. On the other hand, if high reliability should be guaranteed, especially when channel condition is bad or the location distribution of users is concentrated, SCMA is a feasible solution due to its shaping gain and near-optimal MPA detection. Of course, when the number of users is large enough, it may be difficult to design a codebook for each user, and in this case, LDS-OFDM or MUSA can also be used to reduce the design complexity at the transmitter or receiver, separately.

As elaborated in previous sections, certain signal processing modules are common to several NOMA schemes, for example, MPA at the receiver or spreading operation at the transmitter, which can be shared in hardware so that the hardware cost would be reduced at both user terminals and base stations. These general-purpose modules can be combined in different forms at the software level to implement different schemes. The switching between NOMA schemes is fast and flexible with software defined hardware architecture, and can quickly adapt to different deployment scenarios, that is, from capacity achieving to user loading improvement.

To enable SoDeMA, the frame structure should be flexible enough so that the time and frequency resources are partitioned into different blocks freely for different services and users. In each resource block, one specific multiple access scheme is configured with specific waveform, duplex mode, pilot signals, power level, and so on. Note that the inter-subcarrier interference between different resource blocks needs to be carefully mitigated. The proposed SoDeMA concept provides a flexible configuration of multiple access schemes to support different services and applications. It is highly expected that SoDeMA can be carefully designed to adapt to various application scenarios to support the system design goal of "anything as a service" in future 5G networks.

## CONCLUSIONS

In this article, we have discussed and compared several major NOMA schemes for 5G from the aspects of basic principles, key features, receiver complexity, engineering feasibility, and so on. Compared to conventional OMA, NOMA allows controllable interferences to realize overloading at the cost of a tolerable increase of receiver complexity. Therefore, the demands of spectral efficiency and massive connectivity for 5G can be partially fulfilled by NOMA. We have also highlighted key challenges, opportunities and future research trends for the design of NOMA, including theoretical work, optimal design of spreading sequences or codebooks, receiver design, a grant-free NOMA mechanism, and so on. The proposed concept of SoDeMA is able to flexibly support diverse services and applications with different requirements. It is expected that NOMA will play an important role in future 5G wireless communications.

## REFERENCES

- [1] F. Boccardi et al., "Five Disruptive Technology Directions for 5G," *IEEE Commun. Mag.*, vol. 52, no. 2, Feb. 2014, pp. 74–80.
- [2] Y. Saito et al., "Non-Orthogonal Multiple Access (NOMA) for Future Radio Access," *Proc. IEEE VTC-Spring '13*, June 2013, pp. 1–5.
- [3] K. Higuchi and Y. Kishiyama, "Non-Orthogonal Access with Random Beamforming and Intra-Beam SIC for Cellular MIMO Downlink," *Proc. IEEE VTC-Fall '13*, Sept. 2013, pp. 1–5.
- [4] S. Han et al., "Energy Efficiency and Spectrum Efficiency Co-Design: From NOMA to Network NOMA," *IEEE MMTTC E-Letter*, vol. 9, no. 5, Sept. 2014, pp. 21–24.
- [5] R. Hoshyar, F. P. Wathan, and R. Tafazolli, "Novel Low-Density Signature for Synchronous CDMA Systems over AWGN Channel," *IEEE Trans. Signal Proc.*, vol. 56, no. 4, Apr. 2008, pp. 1616–26.
- [6] M. Al-Imari et al., "Uplink Nonorthogonal Multiple Access for 5G Wireless Networks," *Proc. 11th Int'l. Symp. Wireless Commun. Sys.*, Aug. 2014, pp. 781–85.
- [7] H. Nikopour and H. Baligh, "Sparse Code Multiple Access," *Proc. IEEE PIMRC 2013*, Sept. 2013, pp. 332–36.
- [8] Z. Yuan, G. Yu, and W. Li, "Multi-User Shared Access for 5G," *Telecommun. Network Technology*, vol. 5, no. 5, May 2015, pp. 28–30.
- [9] J. Huang et al., "Scalable Video Broadcasting Using Bit Division Multiplexing," *IEEE Trans. Broadcast.*, vol. 60, no. 4, Dec. 2014, pp. 701–06.
- [10] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*, Cambridge Univ. Press, 2005.
- [11] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, "Factor Graphs and the Sum-Product Algorithm," *IEEE Trans. Info. Theory*, vol. 47, no. 2, Feb. 2001, pp. 498–519.
- [12] X. Dai et al., "Successive Interference Cancellation Amenable Multiple Access (SAMA) for Future Wireless Communications," *Proc. IEEE ICCS 2014*, Nov. 2014, pp. 1–5.
- [13] K. Kusume, G. Bauch, and W. Utschick, "IDMA vs. CDMA: Analysis and Comparison of Two Multiple Access Schemes," *IEEE Trans. Wireless Commun.*, vol. 11, no. 1, pp. 78–87, Jan. 2012.
- [14] G. Wunder, P. Jung, and C. Wang, "Compressive Random Access for Post-LTE Systems," *Proc. IEEE ICC '14*, June 2014, pp. 539–44.
- [15] V. V. Veeravalli and A. Mantravadi, "The Coding-Spreading Tradeoff in CDMA Systems," *IEEE JSAC*, vol. 20, no. 2, Feb. 2002, pp. 396–408.

## BIOGRAPHIES

LINGLONG DAI [M'11, SM'14] (dail@tsinghua.edu.cn) received his B.S. degree from Zhejiang University in 2003, his M.S. degree (with highest honor) from the China Academy of Telecommunications Technology (CATT) in 2006, and his Ph.D. degree (with the highest honor) from Tsinghua University, Beijing, China, in 2011. From 2011 to 2013, he was a postdoctoral fellow with the Department of Electronic Engineering, Tsinghua University, where he has been an assistant professor since July 2013. His research interests are in wireless communications, with a focus on multi-carrier techniques, multi-antenna techniques, and multi-user techniques. He has published over 60 journal and conference papers. He has received the Outstanding Ph.D. Graduate of Tsinghua University award in 2011, the Excellent Doctoral Dissertation of Beijing award in 2012, the IEEE ICC Best Paper Award in 2013, the National Excellent Doctoral Dissertation Nomination Award in 2013, the IEEE ICC Best Paper Award in 2014, the URSI Young Scientists Award in 2014, and the IEEE Scott Helt Memorial Award in 2015 (*IEEE Transactions on Broadcasting* Best Paper Award). He currently serves as Co-Chair of the IEEE Special Interest Group (SIG) on Signal Processing Techniques in 5G Communication Systems.

BICHAI WANG [S'15] (wang-bc11@mails.tsinghua.edu.cn) received her B.S. degree in electronic engineering from Tsinghua University in 2015. She is currently working toward her Ph.D. degree in the Department of Electronic Engineering, Tsinghua University. Her research interests are in wireless communications, with emphasis on new multiple access techniques. She received the Freshman Scholarship of Tsinghua University in 2011, Academic Merit Scholarships of Tsinghua University in 2012, 2013, and 2014, respectively, and the Excellent Thesis Award of Tsinghua University in 2015.

---

YIFEI YUAN (yifei.yuan@ztetx.com) received Bachelor's and Master's degrees from Tsinghua University, and a Ph.D. from Carnegie Mellon University, Pennsylvania. He was with Alcatel-Lucent from 2000 to 2008 working on 3G/4G key technologies. Since 2008, he has been with ZTE, responsible for standards research on LTE-Advanced physical layer and 5G technologies. His research interests include MIMO, iterative codes, resource scheduling, non-orthogonal access, and small cells. He was admitted to the Thousand Talent Plan Program of China in 2010. He has published extensively, including a book on LTE-A relay and a book on LTE-Advanced key technologies. He has over 30 granted patents.

SHUANGFENG HAN (hanshuangfeng@chinamobile.com) received his M.S. and Ph.D. degrees in electrical engineering from Tsinghua University in 2002 and 2006, respectively. He joined Samsung Electronics as a senior engineer in 2006 working on MIMO, multi-BS MIMO, and so on. Since 2012, he has been a senior project manager in the Green Communication Research Center at the China Mobile Research Institute. His research interests are green 5G, massive MIMO, full duplex, NOMA, and EE-SE co-design.

CHIH-LIN I (icl@chinamobile.com) received her Ph.D. degree in electrical engineering from Stanford University. She has been working at multiple world-class companies and research institutes leading R&D, including AT&T Bell Labs, AT&T HQ, ITRI of Taiwan, and ASTRI of Hong Kong. She received the *IEEE Transactions on Communications* Stephen Rice Best Paper Award and is a winner of the CCCP National 1000 Talent program. Currently, she is China Mobile's chief scientist of wireless technologies and has established the Green Communications Research Center, spearheading major initiatives including system architectures, technologies, and devices; green energy; and C-RAN and soft base

stations. She was an elected Board Member of IEEE ComSoc, Chair of the ComSoc Meetings and Conferences Board, and Founding Chair of the IEEE WCNC Steering Committee. She is currently an Executive Board Member of GreenTouch and a Network Operator Council Member of ETSI NFV. Her research interests are green communications, C-RAN, network convergence, and bandwidth active antenna arrays.

ZHAOCHENG WANG [M'09, SM'11] (zchwang@tsinghua.edu.cn) received his B.S., M.S., and Ph.D. degrees from Tsinghua University in 1991, 1993, and 1996, respectively. From 1996 to 1997, he was a postdoctoral fellow of Nanyang Technological University, Singapore. From 1997 to 1999, he was with OKI Techno Centre (Singapore) Pte. Ltd., where he was first a research engineer and later became a senior engineer. From 1999 to 2009, he was with Sony Deutschland GmbH, where he was first a senior engineer and later became a principal engineer. He is currently a professor with the Department of Electronic Engineering, Tsinghua University, and serves as director of the Broadband Communication Key Laboratory, Tsinghua National Laboratory for Information Science and Technology. He has authored or coauthored over 80 international journal papers (SCI indexed). He is the holder of 34 granted U.S./EU patents. He co-authored two books, one of which, *Millimeter Wave Communication Systems*, was selected for the IEEE Series on Digital & Mobile Communication and published by Wiley-IEEE Press. His research areas include wireless communications, visible light communications, millimeter-wave communications, and digital broadcasting. He is a Fellow of the Institution of Engineering and Technology. Currently he serves as an Associate Editor of *IEEE Transaction on Wireless Communications* and *IEEE Communications Letters*, and has also served as Technical Program Committee Co-Chair of various international conferences.

# Rethink Fronthaul for Soft RAN

Chih-Lin I, Yannan Yuan, Jinri Huang, Shijia Ma, Chunfeng Cui, and Ran Duan

## ABSTRACT

In this article we discuss the design of a new fronthaul interface for future 5G networks. The major shortcomings of current fronthaul solutions are first analyzed, and then a new fronthaul interface called next-generation fronthaul interface (NGFI) is proposed. The design principles for NGFI are presented, including decoupling the fronthaul bandwidth from the number of antennas, decoupling cell and user equipment processing, and focusing on high-performance-gain collaborative technologies. NGFI aims to better support key 5G technologies, in particular cloud RAN, network functions virtualization, and large-scale antenna systems. NGFI claims the advantages of reduced bandwidth as well as improved transmission efficiency by exploiting the tidal wave effect on mobile network traffic. The transmission of NGFI is based on Ethernet to enjoy the benefits of flexibility and reliability. The major impact, challenges, and potential solutions of Ethernet-based fronthaul networks are also analyzed. Jitter, latency, and time and frequency synchronization are the major issues to overcome.

## INTRODUCTION

With the maturity and wider deployment of fourth generation (4G) networks, future 5G technologies have become a research focus. As industrial progress accelerates, some achievements have been presented [1, 5]. Several white papers have been published by various organizations such as Next Generation Mobile Networks (NGMN), IMT-2020, Mobile and Wireless Communications Enablers for 2020 Information Society (METIS), 5G Infrastructure Public-Private Partnership (5G PPP), and others, while some proofs of concepts (PoCs) have been developed to allow people to have a quick grasp on 5G. Although the understanding of 5G may still differ among different people, there is a wide consensus that 5G should be a software defined network with the benefits of flexibility, quicker time to market, unified management, and flourishing applications [1]. Network functions virtualization (NFV) [2] is a strong technology candidate toward this end, and cloud radio access network (C-RAN) [3, 4] is an NFV

instance on the RAN side to achieve soft RAN. First proposed by China Mobile [3], C-RAN centralizes baseband processing units and virtualizes them into a resource pool. C-RAN has been viewed as a promising 5G RAN architecture. In addition to softness, C-RAN could also bring operators such benefits as quicker network deployment, system performance improvement, and energy savings.

On the road of C-RAN realization, the fronthaul (FH) issue has been one of the biggest challenges. An FH connection is a link between a baseband unit (BBU) and a remote radio head (RRH). Typical FH interfaces include the common public radio interface (CPRI), open base station architecture initiative (OBSAI), and open radio interface (ORI). The data rate of the FH connection for Long Term Evolution (LTE) is on the order of gigabits per second. The common FH solution in C-RAN is to use dark fiber. Due to the high FH data rate, centralization requires consumption of a number of fiber cores, which are scarce and not easy to afford. Although other transport technologies such as wavelength-division multiplexing (WDM) and optical transport network (OTN) could save fiber consumption, the cost of the introduction of additional transport equipment makes economic viability a concern of operators. Because of the concern regarding FH cost, some operators are still not very convinced of the merits of C-RAN deployment. Therefore, enabling large-scale C-RAN deployment in 5G requires reducing the FH bandwidth.

Current FH interfaces could also raise new issues for C-RAN in an NFV environment. C-RAN is supposed to run on general-purpose platforms (GPPs) consisting of standard IT servers, storage, and switches. However, the GPP platform does not provide an FH interface for telecom applications. To support FH, either a new interface should be created on the GPP platform or an adapter card is needed, both complicating the system and introducing additional cost. It would be desirable for the 5G FH interface to be based on existing GPP interfaces to maximize efficiency and save cost.

In addition, scalability issues exist for today's FH technology to support widely discussed 5G technologies, including large-scale antenna systems (LSASs), coordinated multipoint (CoMP)

The authors are with  
China Mobile Research  
Institute.

processing, and so on. Take LSAS as an example. It is possible that a 5G RRH could be equipped with 64 or even 128 antennas. With LTE, the FH bandwidth will rise to 100 Gb/s, at which point it is unaffordable. It is clear that the impact of the number of antennas on FH should be minimized to the greatest extent possible.

The industry is gradually realizing the deficiencies of current FH solutions as well as their importance to 5G, and making efforts on improvements. In NGMN, schemes of the BBU-RRH function split are analyzed, aiming to reduce the FH bandwidth to facilitate C-RAN deployment [5]. ORI is studying compression technology to reduce the CPRI data rate. The CPRI Forum has begun the discussion on radio over Ethernet, the idea of which is to use Ethernet to transport the CPRI stream, while in the IEEE, the IEEE 1904.3 Task Force was founded recently, targeting the design of CPRI encapsulation on Ethernet packets [6, 7]. In addition, IEEE 802.1 TSN decided very recently to initiate the development of a potential new work item on time-sensitive networking for FH. The IEEE 1588 Working Group is also considering adding optional specialized solutions to the next edition of IEEE 1588 to enable enhanced synchronization accuracy for FH.

Despite these continuous efforts, from China Mobile's perspective we think that much more effort is required. The current improvements have not touched the root of the FH itself. A new FH interface is required to better support C-RAN large-scale deployment, NFV realization in baseband virtualization, as well as serving other 5G key technologies.

In this article, we share our ideas on the future FH, the called next-generation FH interface (NGFI). From CMCC's perspective, a desirable NGFI should have dynamic bandwidth with traffic variation, be antenna-independent and packet-based, and support collaborative technologies. We describe the shortcomings of CPRI and the definition of NGFI. We elaborate on the design principles of NGFI, while the major impact and challenges of FH transport networks are presented, followed by conclusions.

## DEFINITION OF NGFI

In this section, we describe our view on what the NGFI should look like and the major advantages we expect from such an interface redesign. In this and the remaining sections, we use CPRI as an example of traditional FH protocols since it is the most widely used and the basis for other FH interfaces such as ORI.

Figure 1 shows a C-RAN network architecture combining the ideas of software defined network (SDN) and NFV, which consists of a radio cloud center (RCC), an NGFI-based FH network, and new types of RRHs. An RCC provides a cloud platform consisting of standard IT servers, storage, and switches. In an RCC all the radio access network functions appear as software applications running in virtual machines (VMs). Furthermore, in an RCC, the control plane and user plane can be separated based on the wireless performance requirements. Under this network architecture, an NGFI is proposed

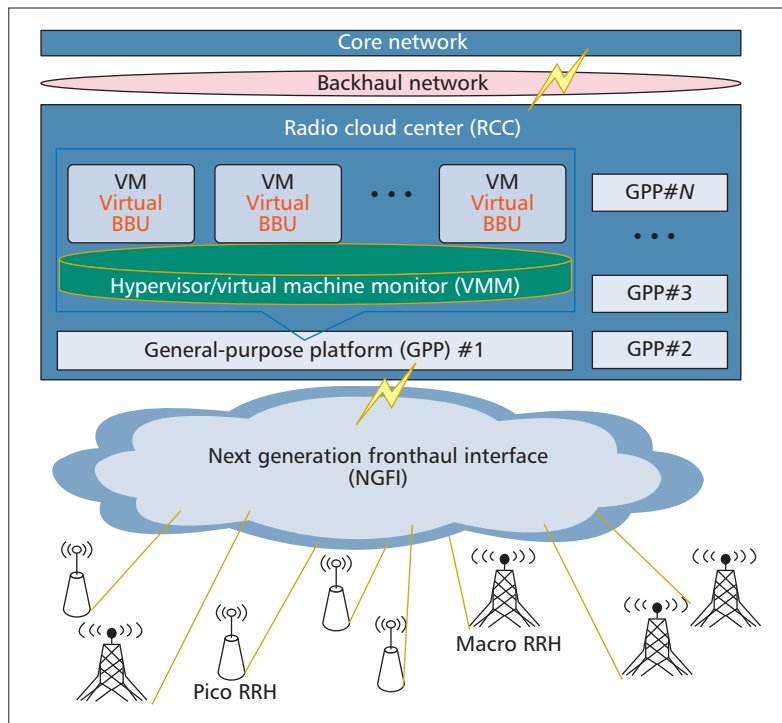


Figure 1. C-RAN network architecture with NGFI.

to connect an RCC and new RRHs. Compared to traditional RRHs, a new RRH contains not only radio processing but also partial baseband processing. What kind of baseband processing should be included in an RRH is discussed later.

## REVISITING CPRI

The CPRI helps separate the BBU and RRH to enable the deployment of distributed base stations. Although in traditional networks the CPRI is mainly deployed with short distances, usually on the order of several meters or several hundred meters, it can support up to 40 km between the BBU and RRH. CPRI has been working well for traditional mobile networks including 2G, 3G and 4G. With networks evolving to 5G, the following three factors make CPRI more and more unsuitable to accommodate such an evolution:

- Constant data rate due to synchronous digital hierarchy-based (SDH-based) transmission mode [6]
- Fixed one-to-one correspondence between RRH and BBU
- Sampling I/Q data rate dependent on the number of antennas

First, mobile traffic varies in the temporal dimension. For example, the data traffic in an office area is high in the daytime but plummets at midnight. For dense urban areas, the tidal wave effect is noticeable. However, the CPRI data stream is SDH-like [6], which means that it is constant regardless of the change of traffic even when there is no traffic at all. This then leads to low utilization efficiency.

Second, with CPRI an RRH has a one-to-one correspondence to a BBU. The relationship is configured offline. It may cause concern in the context of C-RAN. In C-RAN, the BBUs are

Site 1	Peak			Valley		
	RB (%)	BH (Mb/s)	Duration	RB (%)	BH (Mb/s)	Duration
DL	83.54	175	15: 30~ 16: 00	0.72	0.96	0: 30~ 6: 00
UL	20.19	12.5	15: 00~ 15: 15	3.02	0.23	1: 30~ 6: 30

**Table 1.** Statistical data of traffic load of site 1.

centralized and virtualized in a pool. Reliability becomes extremely important as each pool takes care of thousands of users. Therefore, for the sake of protection, it would be desirable if in C-RAN, one RRH could be automatically switched to another BBU pool. Current CPRI, however, does not support such flexible and automatic rerouting.

Finally, the CPRI bandwidth is dependent on the number of antennas. As the number of antennas increases, the CPRI data rate increases in proportion. This could become a major hindrance for CPRI's applicability in 5G as far as multiple-antenna technologies are concerned.

Based on the analysis above, we believe that a new FH interface is required to better support the 5G evolution.

#### NEXT-GENERATION FRONTHAUL INTERFACE

The purpose of designing a new FH interface is to facilitate C-RAN deployment, to make it compatible with GPP platforms and scalable enough to support the evolution to 5G.

We define the new FH as the NGFI between the BBU and RRH with the following four features:

- Its data rate should be traffic-dependent and therefore support statistical multiplexing.
- The mapping between BBU and RRH should be one-to-many and flexible.
- It should be independent of the number of antennas.
- It should be packet-based, that is, the FH data could be packetized and transported via packet-switched networks.

A traffic-aware FH interface could fully leverage the tidal wave effect of mobile networks so that the FH data rate is variable with traffic change. This way the FH transmission efficiency is improved. With a reduced average data rate, it will increase the likelihood of FH transport networks to deploy C-RAN. The antenna-independent feature ensures that the NGFI can support 5G antenna technologies. The packet-based feature makes it possible to use Ethernet to transport FH data. The benefits would be manifold. First, an Ethernet interface is the most common interface on standard IT servers, and the use of Ethernet makes C-RAN virtualization easier and cheaper. Then Ethernet can fully make use of the dynamic nature of NGFI to realize statistical multiplexing. The flexible routing capability could also be used to realize multiple paths between a BBU pool and an RRH.

With the novelty of NGFI, it is clear that it will have a great impact on the FH transport network. In the following sections, we elaborate on our design principles for NGFI, and analyze the impact and challenges confronting FH transport network design.

## RETHINKING THE BBU-RRH FUNCTION SPLIT

Traditionally, the baseband-related functions are processed by the BBU while the RRH processes radio frequency related functions. It is this simple partitioning that leads to the shortcomings of CPRI, as mentioned above. Therefore, the NGFI design should start with a paradigm shift by rethinking and redesigning the function split between BBU and RRH. Moreover, the function split between BBU and RRH may be different according to the bandwidth and latency of FH, which is adaptive to different scenarios.

#### DECOUPLING THE FH BANDWIDTH FROM THE NUMBER OF ANTENNAS

The air interface bandwidth per carrier on 2G, time-division synchronous code-division multiple access (TD-SCDMA) and time-division LTE (TD-LTE) are 0.2 MHz, 1.6 MHz, and 20 MHz, respectively. Correspondingly, FH transport bandwidth per carrier is 30 Mb/s, 400 Mb/s, and 10 Gb/s [8], respectively. At the same time, FH is facing a bandwidth explosion, considering the rapid traffic growth in 5G (potentially 1000× by 2020). Compared to the air interface bandwidth, the existing FH interface transportation efficiency is low. One of the most important reasons is that FH bandwidth is proportional to the number of antennas. In order to increase transport efficiency, a BBU/RRH function split scheme should enable NGFI to decouple FH bandwidth and the antenna number.

Taking TD-LTE and a large-scale antenna system (LSAS) as an example, an 8-antenna TD-LTE carrier FH bandwidth is 10 Gb/s based on the current BBU/RRH function split. If it is a 128-antenna LSAS system, a TD-LTE carrier FH bandwidth will increase to 160 Gb/s. Thus, the existing FH faces a big challenge with the increase of the number of antennas. In order to reduce the bandwidth, one potential idea is to redesign the BBU/RRH function split, that is, the antenna related functions should be moved to the RRH and the non-antenna related functions should remain on the BBU. In particular, in order to phase out the effect of the number of antennas, it is proposed that antenna related functions (e.g., downlink antenna mapping, fast Fourier transform [FFT], channel estimation, equalization) should be moved from the BBU to the RRH. It is shown that an LTE carrier FH bandwidth may decrease on the order of 100 Mb/s no matter how many antennas are used [9]. Therefore, FH bandwidth will decrease significantly if the BBU/RRH function split can decouple non-antenna-related processing and antenna-related processing.



## DECOUPLING CELL/UE PROCESSING

Dynamic variation is a major feature of wireless traffic. The tidal wave effect is obvious in many wireless deployment scenarios such as residential, office, and commercial districts. Moreover, the traffic load of most areas is usually in the valley between late night and early morning.

In order to quantify wireless traffic features, a TD-LTE traffic investigation on a commercial LTE network was done via network monitoring systems. There were six base station sites in the investigation, each having at least three carriers. Two of them are indoor distributed systems while the others are outdoor macro base station sites. The investigation period was 7 days during which the traffic load was sampled and collected by network monitoring systems every 15 minutes.

In Table 1, Table 2, and Table 3, the statistical traffic load of two TD-LTE base station sites is shown where DL means downlink and UL means uplink. Site 1 is an indoor site and site 2 is an outdoor site. In the table, RB means resource block utilization, which is expressed as a percentage. BH means backhaul transportation bandwidth, which is in the unit of megabits per second. The duration of peak/valley is the absolute time in 24 hours. The duration of average load is the time length during which traffic load distributes in [average load-1%, average load+1%]. Based on the statistical data, several traffic load features were observed:

- The tidal wave effect is obvious in the test districts, where the traffic load is almost zero for 12 hours.
- Even when the site is at peak status the RB usage is not high. The duration of a site at peak status is short, usually not exceeding 30 minutes.
- The probability that different sites are simultaneously in peak status is almost zero.
- When the site is at valley status the RB usage is low and the duration is long, usually exceeding two hours.
- Different sites are frequently in valley status at the same time. Moreover, the overlapping time is long among different sites.
- Most of the time the traffic load stays at the average level, which is low.

Based on the above observations, it is clear that constant-rate FH transport does not match the mobile traffic features, which results in a waste of resources. To address this issue, we first observe that the existing baseband processing can be divided into cell processing and user equipment (UE) processing. Cell processing is irrelevant to traffic load and is fixed no matter how many UEs are active. Some examples of such processing units in LTE include inverse FFT (iFFT)/FFT, cyclic prefix (CP) addition/removal, cell-specific reference signal/primary synchronization signal/secondary synchronization signal (CRS/PSS/SSS) generation, and physical broadcast channel (PBCH) processing. It is therefore proposed to move these cell processing functions from the BBU to the RRH, that is, decoupling the cell and UE processing.

If cell processing is moved from the BBU to the RRH, the FH bandwidth will be lower and

Site 2	Peak			Valley		
	RB (%)	BH (Mb/s)	Duration	RB (%)	BH (Mb/s)	Duration
DL	20.38	144	13:30~13:45	0.70	0	21:00~8:15
UL	47.06	34.2	14:00~14:15	3.00	0	21:00~8:15

**Table 2.** Statistical data of traffic load of site 2.

	Site1			Site2		
	RB (%)	BH (Mb/s)	Duration (hours)	RB (%)	BH (Mb/s)	Duration (hours)
DL	3.29	3.95	4.3	1.27	2.20	22
UL	4.22	0.31	15	4.00	0.40	22

**Table 3.** Average load distribution.

load-dependent. The load-dependent feature gives an opportunity to exploit the statistical multiplexing gain when it comes to FH transport network design for C-RAN deployment. Thanks to statistical multiplexing, the bandwidth needed for transport of a number of FH links in C-RAN could be greatly reduced, thereby diminishing the cost.

Cell/UE processing decoupling can further help reduce power consumption and enhance network reliability. This is because cell basic coverage signal processing is a kind of cell processing. Therefore, cell basic coverage will be provided by the RRH if cell processing functions are moved from the BBU to the RRH. On one hand, BBU software can be switched to a dormant state to save power when there is no active UE. On the other hand, RRH is able to provide continuous air interface coverage, even when a BBU breaks down. This way, it provides sufficient time for BBU fault processing.

### FOCUSING ON HIGH-PERFORMANCE-GAIN COLLABORATIVE TECHNOLOGIES

CoMP has been viewed as one of the important 5G technology candidates to improve system performance, which can be divided into two classes: medium access control (MAC) layer coordination and physical layer coordination. For example, collaborative scheduling is one of the MAC layer coordinated mechanisms. Joint reception (JR) and joint transmission (JT) are physical layer coordinated technologies. The design of NGFI should take into account support for CoMP. The above two principles lead to a low-bandwidth traffic-dependent FH. In the meantime, some physical-layer-coordinated technologies are difficult to implement since some collaborative information has been processed and terminated by an RRH. Fortunately, it is found that the performance gain of JR/JT decreases significantly as the number of antennas increases

Network load	JT/CoMP	CS/CoMP	CoMP/Total
20%	20.04%	79.96%	42.02%
50%	18.09%	81.91%	34.43%
70%	20.89%	79.11%	54.33%
100%	24.89%	75.11%	47.31%

**Table 4.** Utilization factor of JT and CS.

[10]. Moreover, it is also found that MAC-level collaborative technologies can bring comparable performance gains with lower complexity, easier implementation, and fewer constraints.

In order to verify this, a CoMP field trial was conducted in 2014 in which two CoMP schemes including JT and CS were examined and compared. The testing zone was a central business district including around 7000 active user equipments (UEs), which is a typical CoMP test scenario. It was covered by 35 base stations with different antenna heights in which inter-cell interference is serious because of the high-ratio overlapping area.

Test results show that cell edge UE throughput increased by 127.45~173.65 percent when the serving cell reference signal received power varied from -88 dBm to -106 dBm. Table 4 shows the utilization factor of JT, CS, and the CoMP application. In Table 4, CoMP/Total means the utilization factor of CoMP, which is defined as the ratio of the number of CoMP transmission time intervals (TTIs) to the total number of test TTIs. Similarly, JT/CoMP and CS/CoMP are defined as the ratio of the number of JT and CS TTIs to that of CoMP TTIs, respectively.

From Table 4, it is found that the network load growth results in a small reduction of CS usage. For example, when the network load increases to 100 percent, the usage ratio of CS is still around 75 percent, similar to other cases. The usage of JT is only around 25 percent, much lower than CS. It is therefore fair to say that most of the performance gain is contributed by CS. Compared to JT, CS does not need complex matrix computing. It is easier to implement CS with current base station equipment. On the contrary, JT performance is influenced by antenna calibration accuracy, channel estimation accuracy, and channel variation speed, all requiring high FH bandwidth.

The test results demonstrated that MAC-level collaborative technologies could solve most network interference. Therefore, NGFI design should focus on high-performance gain collaborative technologies rather than all the collaborative technologies. This principle provides guidance on how to make a trade-off between wireless and FH performance.

## RETHINKING FH TRANSPORT

Earlier, it was pointed out that a desirable FH interface should be packet-based, which makes it easy to transmit by packet-switched networks,

especially Ethernet. This could make full use of the advantages of Ethernet to achieve multi-point-to-multipoint connection, statistical multiplexing, flexible routing, and so on. However, the adoption of Ethernet also introduces new challenges. In this section, we analyze such challenges and propose potential solutions.

### RELIABLE SYNCHRONIZATION ON PACKETIZED NETWORKS

Time-division multiplexing (TDM) systems require strict synchronization that includes two aspects: frequency and time (or phase). For TD-SDMA and TD-LTE, the accuracy of frequency synchronization should be in the range of  $\pm 0.05$  ppm, while the accuracy of time synchronization should be in the range of  $\pm 1.5 \mu\text{s}$  [11].

In CPRI three types of data including wireless protocol data, synchronization data, and control and management data are packaged together and transmitted in TDM mode. Upon receiving the CPRI frames, the clock and data recovery (CDR) circuit of an RRH can extract the frequency information to achieve frequency synchronization. Meanwhile, the CPRI transport time is nearly constant and can be measured by the BBU. Based on the measurement, the timing between the BBU and RRH can be configured in advance. With the timing information extracted from CPRI frames, time synchronization at the RRH can be achieved.

For Ethernet-based NGFI, as opposed to CPRI, the transport time of FH data is no longer constant due to the packet-switched nature of Ethernet. As a result, frequency and time synchronization between the BBU and the RRH potentially becomes difficult. To address this issue, one potential solution is to use a synchronous Ethernet (SYNC-E)/1588v2 hybrid network. The working principle of SYNC-E is similar to CPRI as both use 8B/10B encoders in the physical layer. Therefore, high-accuracy frequency synchronization can still be achieved [12].

To achieve time synchronization, a potential solution is to use 1588v2, which is a high-accuracy time synchronization protocol based on packetized networks. In order to meet the high time synchronization requirement, a 1588v2 module should be added in both the BBU and RRH. At present, the accuracy of 1588v2 is on a magnitude order of 100 ns for one hop [13]. The major issue of 1588v2 when adopted in RRHs is the time hopping issue since an RRH only obtains time offset information between it and the BBU. Time hopping could result in discontinuous transmission on RRHs, which is intolerable for mobile communications. One potential solution is to use a frequency adjustable oscillator to calibrate time in the RRH, which can adjust the oscillator frequency gradually to ensure a continuous time variation.

When it comes to support for MIMO or TX diversity transmission technologies, the time synchronization requirement is stricter, and should be in the range of  $\pm 65$  ns [14]. This imposes a big challenge for Ethernet-based NGFI even when 1588v2 is leveraged. More efforts are needed to figure out how to meet the  $\pm 65$  ns time synchronization requirement such as

deploying a GPS antenna on the BBU, improving the accuracy of timestamps, increasing clock frequency, and optimizing the deviation adjusting algorithms.

### RELIABLE DATA TRANSPORT ON PACKETIZED NETWORKS

In Fig. 2, an example of FH topology of a C-RAN system is shown. All the RRHs are connected to a BBU pool through a ring Ethernet network. There are multiple routes between the BBU and RRH to help to enhance network reliability. When one of the routes fails, FH packets can be transported through another route.

FH downlink data is encapsulated with the Ethernet header in the BBU and de-encapsulated from the Ethernet header in the RRH and vice versa for the uplink. In Fig. 2, the structure of NGFI-supporting Ethernet packet is proposed. It includes the traditional Ethernet header, the NGFI header, and the payload. The source MAC address, destination MAC address, and packet type are filled in the traditional Ethernet header. The packet type here is the new NGFI type to distinguish NGFI packets from other packets. The NGFI header consists of the NGFI packet sub-type, the packet length, and the reserved field for protocol extension. There could be at least two NGFI packet sub-types, one for wireless data, and the other for control and management data. For the control and management types, it may include the link delay test, link status monitoring, RRH configuration, and RRH status report.

In CPRI jitter is negligible, while it is common and unavoidable in Ethernet-based NGFI since all the packets are processed in every network node based on the store-and-forward pattern. Therefore, transport latency fluctuates over a range. In order to meet RRH air interface timing, an appropriate circular buffering may be needed in the RRH. On one hand, data packets can be sorted in order. On the other hand, data is continuously sent to the air interface because the transport jitter can be isolated by the buffer.

The maximum transport latency is another big factor influencing FH transport performance. Take China Mobile's current packet transport network (PTN) as an example. For PTN equipment, the processing time of one hop is 50  $\mu$ s [15]. A typical PTN ring consists of 6 nodes (hops) and has 20 km length. The transmission time of fiber for 20 km is 100  $\mu$ s. In the case of 6 hops, the total delay is  $100 + 50 \times 6 = 400 \mu$ s. However, a RAN has strict timing requirements. For example, it is specified by the Third Generation Partnership Project (3GPP) that in LTE from the instant a RRH receives a frame from UE, within 3 ms it must respond by beginning to transmit the responding DL frame. The 3 ms time budget is consumed by BBU processing, RRH processing, and FH transportation, which includes the transmission latency on fiber and processing latency by FH nodes. The more time the BBU and RRH processing takes, the less budget can be allocated to FH transport. As a result, the maximum allowable transport latency for FH networks requires co-design from both the wireless and transport perspectives.

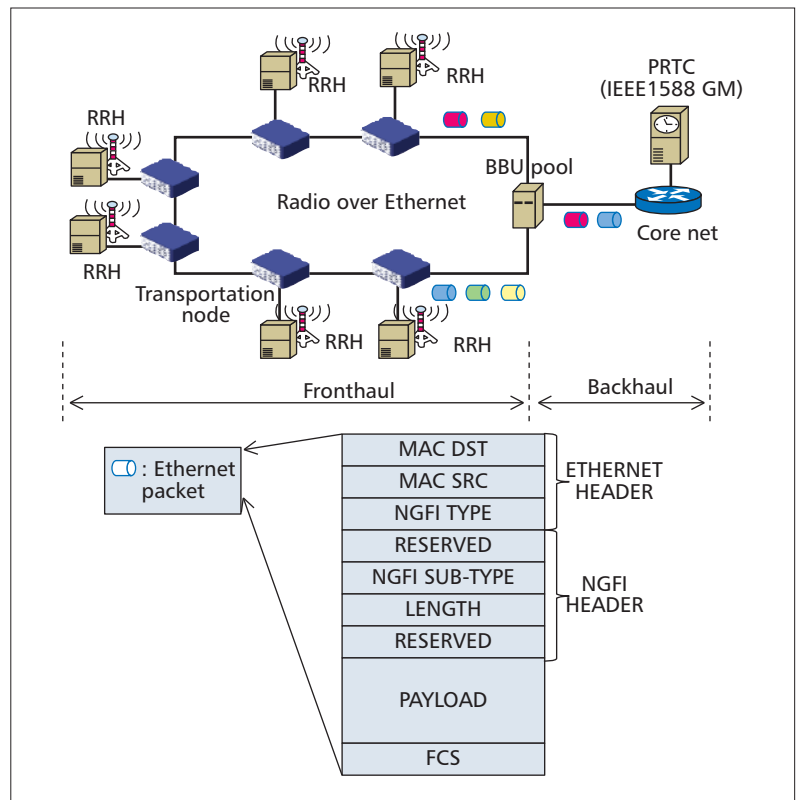


Figure 2. An FH topology example and illustrative Ethernet packet format to support NGFI.

### CONCLUSIONS

Traditional FH interfaces are not suitable for 5G technologies such as C-RAN, NFV, and LSAS due to the features of SDH-based transmission, point-to-point connection, and low transmission efficiency. In this article, a new FH interface called the next-generation fronthaul interface is proposed. The design principles are described, and the major impact, challenges, and potential solutions of and for FH transport networks are analyzed. NGFI requires redesign of the BBU-RRH function split and packetization of FH data. By decoupling the FH bandwidth from the antenna number, NGFI can better support large antenna technologies. With decoupling of cell and UE processing, the NGFI data rate varies with traffic change, which enables exploiting the statistical multiplexing gain to improve efficiency. The use of Ethernet for NGFI transmission brings the benefits of improved reliability and flexibility due to the packet switching nature of Ethernet. While SYNC-E and 1588v2 could be introduced to address the time and frequency synchronization issues, they still need careful design in order to support CoMP technologies. In the meantime, jitter and latency remain the other key difficulties to overcome to finally realize NGFI.

In addition to the challenges analyzed in this article, in the future there remains a lot of work to do to deeply understand NGFI. For example, the analysis of traffic performance with NGFI in 5G networks is necessary to evaluate the performance gain of NGFI. In addition, the control, data, and management channels that are trans-

While SYNC-E and 1588v2 could be introduced to address the time and frequency synchronization issues, they still need careful design in order to support CoMP technologies. In the meantime, jitter and latency remain the other key difficulties to overcome to finally realize NGFI.

ported via NGFI should be analyzed and carefully designed to make NGFI a better fit in different 5G architectures.

## REFERENCES

- [1] C.-L. I *et al.*, "Toward Green and Soft: a 5G Perspective," *IEEE Commun. Mag.*, vol. 52, no. 2, 2014, pp. 66–73.
- [2] ETSI NFV ISG. "Network Functions Virtualisation," Dec., 2012, <http://portal.etsi.org/portal/server.pt/community/NFV/367>
- [3] C. M. R. Institute, "C-ran: The Road Towards Green RAN," v. 1.0.0, Apr. 2010, [labs.chinamobile.com/cran](http://labs.chinamobile.com/cran)
- [4] C.-L. I *et al.*, "Recent Progress on C-RAN Centralization and Cloudification," *IEEE Access*, vol. 2, 2014, pp. 1030–39.
- [5] [www.ngmn.org](http://www.ngmn.org).
- [6] CPRI, "Common Public Radio Interface (CPRI) Specification (V6.0)," tech. rep. Aug. 2013, <http://www.cpri.info>.
- [7] [www.ieee1904.org](http://www.ieee1904.org).
- [8] H. Lee, Y. O. Park, and S. S. Song, "A Traffic-Efficient Fronthaul for the Cloud-RAN," *Info. and Commun. Technology Convergence*, Oct. 2014.
- [9] D. Wubben *et al.*, "Benefits and Impact of Cloud Computing on 5G Signal Processing: Flexible Centralization through Cloud-RAN," *IEEE Sig. Processing Mag.*, Nov. 2014, pp. 35–44.
- [10] A. Davydov *et al.*, "Evaluation of Joint Transmission CoMP in C-RAN based LTE-A HetNets with Large Coordination Areas," *IEEE GLOBECOM Wksp.*, Dec. 2013.
- [11] 3GPP TS 36.101, "User Equipment (UE) Radio Transmission and Reception (Release 9)," v. 8.9.0, Mar. 2009.
- [12] ITU-T Rec. G.8262/Y.1362, "Timing Characteristics of Synchronous Ethernet Equipment Slave Clock (EEC)," Aug. 2007.
- [13] ITU-T Rec. G.8273.2/Y.1368.2, "Timing Characteristics of Telecom Boundary Clocks and Telecom Time Slave Clocks," Apr. 2014.
- [14] 3GPP TS 36.104, "Base Station (BS) Radio Transmission and Reception (Release 11)," v. 11.4.0, Mar. 2013.
- [15] YD/T 1704-2007, "General Technical Requirement for Intelligentized PSTN Network," 2007.

## BIOGRAPHIES

CHIH-LIN I is the China Mobile chief scientist of wireless technologies, in charge of advanced wireless communication R&D effort of China Mobile Research Institute (CMRI). She established the Green Communications Research Center of China Mobile, spearheading major initiatives including 5G Key Technologies R&D; high energy efficiency system architecture, technologies, and devices; green energy; C-RAN; and soft base station. She received her Ph.D. degree in electrical engineering from Stanford University, and has almost 30 years experience in wireless communication area. She has worked in various world-class companies and research institutes, including the Wireless Communication Fundamental Research Department of AT&T Bell Labs; the headquarters of AT&T, as director of wireless communications infrastructure and access technology; ITRI of Taiwan, as director of wireless communication technology; and Hong Kong ASTRI, as VP and the Founding GD of the Communications Technology Domain. She received the *IEEE Transactions on Communications* Stephen Rice Best Paper Award, and is a winner of the CCCP National 1000 Talent program. She was an elected Board Member of IEEE ComSoc, Chair of the ComSoc Meeting and Conference Board, and Founding Chair of the IEEE WCNC Steering Committee.

She is currently Chair of the FuTURE Forum 5G SIG, a Scientific Advisor Board member of the Singapore NRF-Prime Minister's Office, an Executive Board Member of Green-Touch, and a Network Operator Council Member of ETSI NFV.

YANNAN YUAN received her B.S. degree in electronic engineering and M.S. degree in communication and information systems from Xiamen University in 2009 and 2012, respectively. She was a visiting researcher at Tsinghua University from 2010 to 2012. After graduation, she joined the Green Communication Research Center of China Mobile Research Institute. Her current interests include the next generation fronthaul evolution, open source software base station design and development, radio access network function virtualization, virtual base station accelerators, and so on.

JINRI HUANG received his B.S. degree from Xiamen University in electronic engineering in 2001, and his M.S. and Ph.D. degrees in wireless communication from Tsinghua University in 2004 and 2008 respectively. From 2008 to 2011, he was a project manager with the Network Institute of Technology at SK Telecom in South Korea. From 2011 until now, he has been a project manager with China Mobile Research Institute. He has been actively involved in various SDOs including NGMN, ITU-T, ETSI, and so on. His main interests include next generation RAN evolution, wireless resource management, next generation optical transportation technology, fiber-wireless convergence, network functions virtualization, green communication, and so on.

SHUIA MA received his B.S. degree from Beijing Jiaotong University in communication engineering in 2010, and his M.S. in communication and information systems from Beijing Jiaotong University in 2013. From 2013 to 2015, he was a project manager with China Mobile Research Institute. He has been actively involved in Cloud-RAN prototyping design and development, radio over Ethernet research, NGFI research, and so on.

CHUNFENG CUI is director of the Green Communications Research Center of China Mobile Research Institute. He obtained his Ph.D. degree from Beijing University of Posts and Communications in 2003 and has been with China Mobile R&D Center since then. His work focuses on standardization of wireless communications and research on advanced wireless communications. He has been active in 3GPP, IEEE, CCSA, IMT-Advanced Promotion Group, and Future Forum, and was formerly Chair of the requirement group of Future Forum, and Vice Chair of CCSA TC5 WG6 and requirement group of IMT-Advanced Promotion Group. Currently his research interests focus on 5G, cloud RAN, and wireless green communication technologies. He has published more than 20 papers and is a co-author of two books.

RAN DUAN received his B.S. degree from Beijing University of Posts and Telecommunications (BUPT) in electronic engineering in 2003 and his M.S. in wireless communication systems from BUPT in 2006 respectively. From 2006 to 2011, he was a development project manager with 3G BBU design and development in ZTE Telecom in Shanghai. From 2011 until now, he has been a project manager in China Mobile Research Institute. He has been actively involved in cloud RAN prototype design, CPRI over Ethernet analysis, centralization deployment, edge service, and so on.

**CALL FOR PAPERS**  
**IEEE COMMUNICATIONS MAGAZINE**  
**BIO-INSPIRED CYBER SECURITY FOR COMMUNICATIONS AND NETWORKING**

**BACKGROUND**

Nature is Earth's most amazing invention machine for solving problems and adapting to significant environmental changes. Its ability to address complex, large-scale problems with robust, adaptable, and efficient solutions results from many years of selection, genetic drift, and mutations. Thus, it is not surprising that inventors and researchers often look to natural systems for inspiration and methods for solving problems in human-created artificial environments. This has resulted in the development of evolutionary algorithms including genetic algorithms and swarm algorithms, and of classifier and pattern detection algorithms, such as neural networks, for solving hard computational problems.

A natural evolutionary driver is to survive long enough to create a next generation of descendants and ensure their survival. One factor in survival is an organism's ability to defend against attackers, both predators and parasites, and against rapid changes in environmental conditions. Analogously, networks and communications systems use cyber security to defend their assets against cyber criminals, hostile organizations, hackers, activists, and sudden changes in the network environment (e.g., DDoS attacks). Many of the defense methods used by natural organisms may be mapped to cyber space to implement effective cyber security. Some examples include immune systems, invader detection, friend vs. foe, camouflage, mimicry, evasion, and so on. Many cyber security technologies and systems in common use today have their roots in bio-inspired methods, including anti-virus, intrusion detection, threat behavior analysis, attribution, honeypots, counterattack, and the like. As the threats evolve to evade current cyber security technologies, similarly the bio-inspired security and defense technologies evolve to counter the threat.

The goal of this Feature Topic is twofold: (1) to survey the current academic and industry research in bio-inspired cyber security for communications and networking so that the ComSoc community can understand the current evolutionary state of cyber threats, defenses, and intelligence, and can plan for future transitions of the research into practical implementations; and (2) to survey current academic and industry system projects, prototypes, and deployed products and services (including threat intelligence services) that implement the next generation of bio-inspired methods. Please note that we recognize that in some cases, details may be limited or obscured for security reasons. Topics of interests include, but are not limited to:

- Bio-inspired anomaly and intrusion detection
- Adaptation algorithms for cyber security and networking
- Biometrics related to cyber security and networking
- Bio-inspired security and networking algorithms and technologies
- Biomimetics related to cyber security and networking
- Bio-inspired cyber threat intelligence methods and systems
- Moving-target techniques
- Network artificial immune systems
- Adaptive and evolvable systems
- Neural networks, evolutionary algorithms, and genetic algorithms for cyber security and networking
- Prediction techniques for cyber security and networking
- Information hiding solutions (steganography, watermarking) and detection for network traffic
- Cooperative defense systems
- Bio-inspired algorithms for dependable networks

**SUBMISSIONS**

Articles should be tutorial in nature and written in a style comprehensible and accessible to readers outside the specialty of the article. Authors must follow *IEEE Communications Magazine's* guidelines for preparation of the manuscript. Complete guidelines for prospective authors can be found at <http://www.comsoc.org/commag/paper-submission-guidelines>.

It is important to note that *IEEE Communications Magazine* strongly limits mathematical content, and the number of figures and tables. Paper length should not exceed 4500 words. All articles to be considered for publication must be submitted through the IEEE Manuscript Central site (<http://mc.manuscriptcentral.com/commag-ieee>) by the deadline. Submit articles to the "June 2016/Bio-inspired cyber security for communication and networking" category.

**SCHEDULE FOR SUBMISSIONS**

- Submission Deadline: November 1, 2015
- Notification Due Date: February 1, 2016
- Final Version Due Date: April 1, 2016
- Feature Topic Publication Date: June 2016

**GUEST EDITORS**

Wojciech Mazurczyk  
Warsaw University of Technology  
Poland  
[wmazurczyk@tele.pw.edu.pl](mailto:wmazurczyk@tele.pw.edu.pl)

Sean Moore  
Centripetal Networks  
USA  
[smoorephd@gmail.com](mailto:smoorephd@gmail.com)

Errin W. Fulp  
Wake Forest University  
USA  
[fulp@wfu.edu](mailto:fulp@wfu.edu)

Hiroshi Wada  
Unitrends  
Australia  
[hiroshi.wada@nicta.com.au](mailto:hiroshi.wada@nicta.com.au)

Kenji Leibnitz  
National Institute of Information and Communications Technology  
Japan  
[leibnitz@nict.go.jp](mailto:leibnitz@nict.go.jp)

# Baseband Unit Cloud Interconnection Enabled by Flexible Grid Optical Networks with Software Defined Elasticity

Jiawei Zhang, Yuefeng Ji, Jie Zhang, Rentao Gu, Yongli Zhao, Siming Liu, Kun Xu, Mei Song, Han Li, and Xinbo Wang

## ABSTRACT

The evolution toward 5G mobile networks is characterized by supporting higher data rate, excellent end-to-end performance and ubiquitous user-coverage with lower latency, power consumption, and cost. To support this, the RANs are evolving in two important aspects. One aspect is “cloudification,” which is to pool baseband units to be centralized for statistical multiplexing gain. The other aspect is to use advanced optical technologies for digital and analog signal transmission in a cloud-based RAN. In this article, we focus on BBU cloud interconnection with optical layer technologies. Flexible grid optical networks with the enabling technologies are introduced to provide elastic, transparent, and virtualized optical paths between the BBU pools. To improve the elasticity and intelligence of C-RAN, we propose a software defined centralized control plane to coordinate heterogeneous resources from three domains: the BBU domain, radio domain, and optical domain. We report an experimental demonstration of elastic lightpath provision for cloud radio-over-flexible grid optical networks in a software-defined-networking-based testbed.

## INTRODUCTION

The amount of traffic handled by wireless networks will have increased from 3 exabytes in 2010 to over 190 exabytes by 2018, on pace to exceed 500 exabytes by 2020 [1]. Mobile data traffic growth due to proliferation of smart mobile devices and high-definition video applications is accelerating the wireless networks from second/third generation (2G/3G) to 4G and beyond. Looking into the future, 5G wireless technology is on the horizon, which is characterized by supporting higher data rates, excellent end-to-end performance, and ubiquitous user

coverage with lower latency, power consumption, and cost.

To satisfy the 5G requirement, radio access networks (RANs) are evolving in two important aspects. One aspect is “cloudification,” in which the basic concept is to separate the digital baseband units (BBUs) of conventional cell sites from remote radio heads (RRHs), and move BBUs to a cloud-centric site for statistical multiplexing gain. With the benefit of the cloud RAN (C-RAN), conventional complicated and power-hungry cell sites can be simplified to RRH only, which reduces the capital and operating expenditure (CAPEX and OPEX) related to power consumption and site maintenance. From the performance perspective, virtualized BBU pools can be shared by different users and network operators, allowing them to rent the RAN as a service (RANaaS). The C-RAN architecture centralizes a large number of BBUs within a physical location or several geographically distributed sites to form a “BBU cloud.” To achieve high reliability in case of unit failure or to balance the dynamic traffic loads among BBUs, a distributed interconnection to combine multiple BBU pools into a scalable BBU cloud should be paid more attention. In addition, to terminate most traffic at a user’s site, a BBU cloud not only provides baseband processing capability but also supports content storage, which results in a large number of interactions in the BBU cloud. Figure 1 shows three possible scenarios of BBU cloud interconnection in the C-RAN. Therefore, inter-BBU cloud networks should support large scalability, high reliability, and flexible switch capacity.

The other aspect of RAN evolution is to use advanced optical technology for digital and analog signal transmission. First, as the baseband processing function moves to the “cloud,” a large amount of digitalized sampling data are generated by the RRHs and delivered to the BBU pools through the fronthaul for digital

Jiawei Zhang, Yuefeng Ji, Jie Zhang, Rentao Gu, Yongli Zhao, Siming Liu, Kun Xu, and Mei Song are with Beijing University of Posts and Telecommunications.

Han Li is with China Mobile Research Institute.

Xinbo Wang is with the University of California.

radio signal transmission. The bit rate will reach several gigabits per second, even tens of gigabits per second depending on the cell configuration and bandwidth compression [2]. Second, low-latency connections should be guaranteed in the C-RAN because the data transmission through the radio interface (e.g., common public radio interface, CPRI) is time-sensitive. Third, higher RF such as the millimeter-wave (mm-Wave) band is a new frontier for 5G wireless communication, and analog radio over fiber (RoF) has been well studied [3]. In cloud RoF access networks, a broad filter profile (e.g., 60 GHz) is required to switch mm-Wave signals. Since the digital and analog RF signals have different processing methods, they are selected to be sent to different functionalized BBUs in the cloud. Therefore, a flexible optical switch capability is required at the cloud access point to separate digital and analog RF signals. Based on the above, BBU cloud interconnection needs a high-bandwidth, low-latency, and elastic optical transport network to satisfy the requirement of C-RAN. Flexible grid optical networks have been introduced recently, and the term *flexibility* refers to the ability of a network to dynamically adjust its resources, such as bandwidth of light-path, transponder, and modulation format, according to the requirement of each connection [4, 5]. The flexible grid evolves the traditional International Telecommunication Union (ITU) grid toward high flexibility with fine-grained spectrum slots (e.g., 12.5 GHz vs. 50 GHz or 100 GHz) [6], which enables sub-wavelength, super-wavelength, and multi-rate accommodation in a highly spectrum-efficient way.

In this article, we focus on the BBU cloud interconnection with advanced optical layer technology. First, the benefits of flexible grid optical networks with the enabling technologies are introduced to provide elastic, transparent, and virtualized optical paths for inter-BBU cloud networks. Second, the cloudified and optical-enabled RAN is a multi-resource integrated environment, in which radio resources, optical resources, and BBU resources coexist. To improve the elasticity and intelligence of C-RAN, we propose a software defined centralized control plane to converge heterogeneous resources. Finally, we report an experimental setup of elastic lightpath adjustment for cloud radio-over-flexible grid optical networks (C-RoFlex) in a software defined networking (SDN)-enabled testbed.

## KEY BENEFITS BROUGHT BY FLEXIBLE GRID OPTICAL NETWORKS

Flexible grid optical networks have been demonstrated in the literature as a promising technology for inter-data-center networks [7]. Because of similar functions, such as computing and storage, a BBU pool can be seen as a data center. With the benefits of flexible grid optical networks for inter-data-center networks, a high-bandwidth, low-latency, and elastic optical transport network can be provided for BBU cloud interconnection. In this section, we mainly discuss three key benefits brought by flexible grid technology.

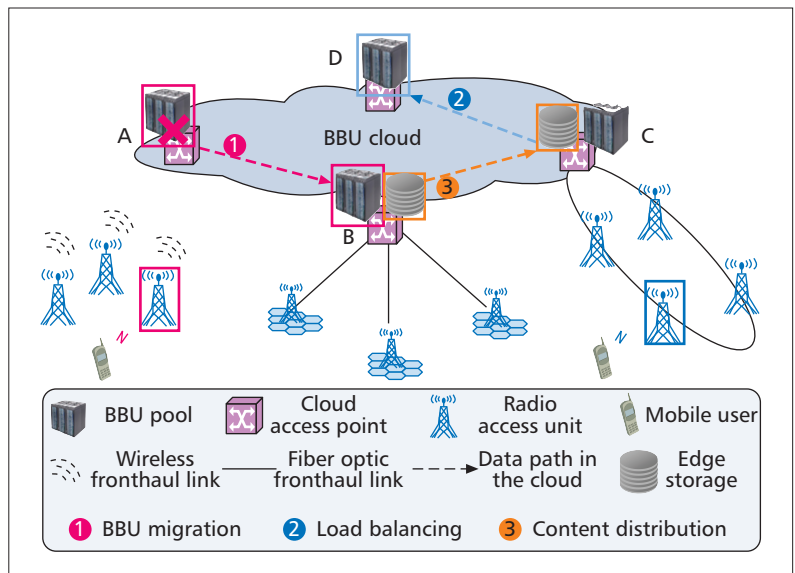
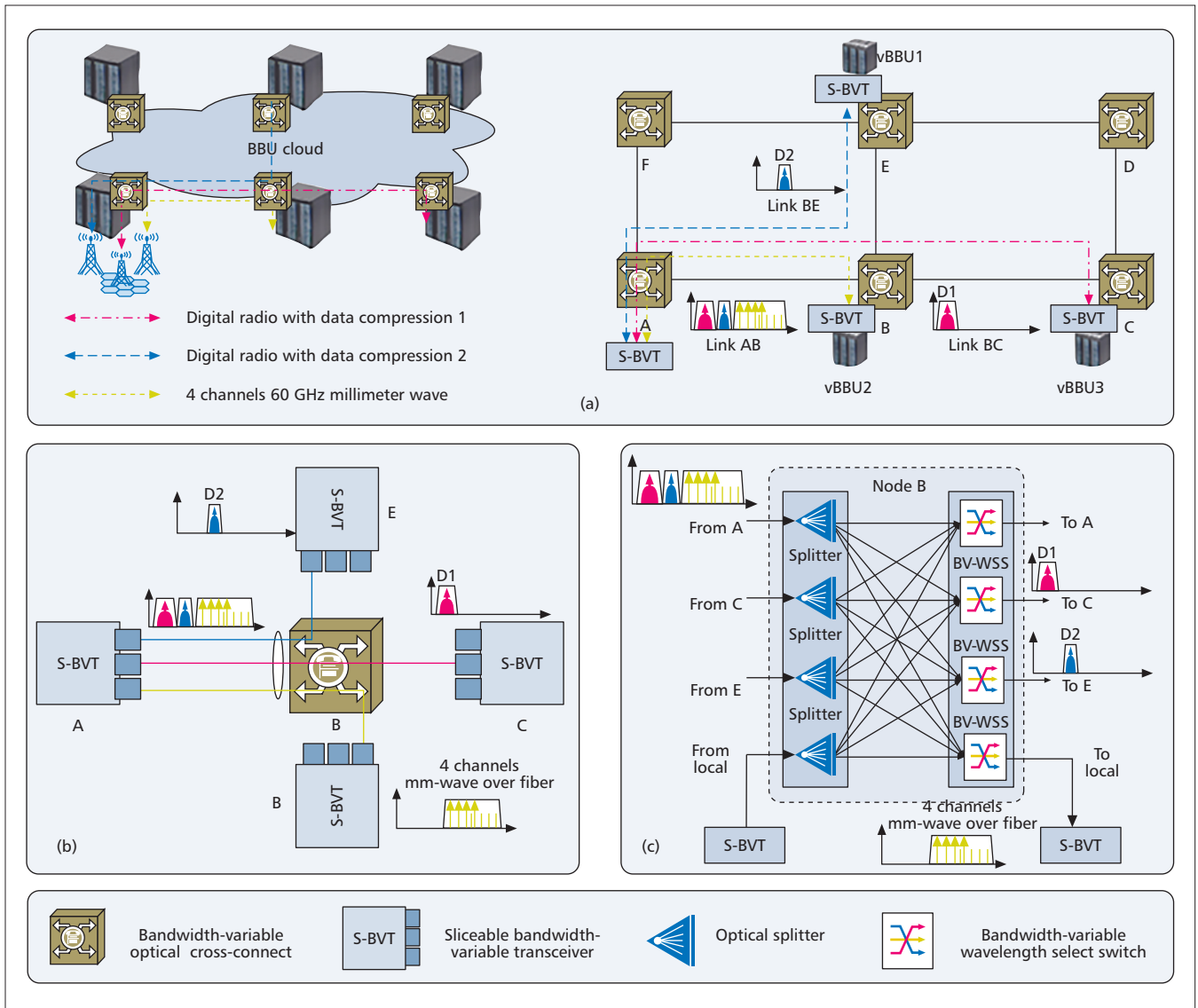


Figure 1. BBU cloud interconnection in the C-RAN.

## ELASTIC OPTICAL SPECTRUM ALLOCATION FOR MULTI-RATE TRAFFIC

C-RAN architecture brings a huge overhead to optical links between RRHs and BBUs. The bit rate has reached the gigabits per second level and presents a multi-rate property. Since different data compression technologies require different bandwidth [2], network operators expect to assign a “just-right-size” spectrum to each data path to improve the spectrum utilization. Besides, in order to fill in the blank (e.g., 60 GHz) between the optical carrier and frequency-band signal, multiple analog RoF signals at mm-Wave band are interleaved [3], and need a larger spectral range than a digital baseband to be filtered out. Therefore, the transport network that connects the BBU cloud requires flexible spectrum switch capability. Flexible grid optical networks can provide an elastic optical channel according to the requirement of each connection. 12.5 GHz is the minimal slot width for an individual optical channel [6], and the bit rate carried by this channel depends on the modulation format assigned to the traffic. Figure 2a shows a scenario of BBU cloud interconnection enabled by flexible grid optical networks. Two digital signals with different data compressions and an analog signal with four interleaved 60-GHz mm-Wave channels accessing the BBU cloud. Spectral resources are assigned properly according to the traffic demands and switched by a bandwidth-variable cross-connect (BV-OXC) to the corresponding virtualized BBU pools. The functionality of BV-OXC is shown in Fig. 2c, in which a bandwidth-variable wavelength selected switch (BV-WSS) is equipped at each direction to filter the variable spectral regions. Compared to conventional wavelength-division multiplexing (WDM) networks, flexible grid optical networks with elastic optical spectrum allocation can efficiently improve the spectral resource utilization to satisfy the multi-rate and multi-standard traffic demand in the C-RAN.



**Figure 2.** Enabling technologies of flexible grid optical networks for C-RAN: a) network scenario; b) sliceable bandwidth-variable optical transceiver; c) bandwidth-variable optical cross-connect.

### TRANSPARENT AND LOW-LATENCY CONNECTION PROVISION

Optical transport networks that connect the BBU cloud not only need to support high bandwidth and be cost efficient, but also need to support strict latency requirement, because the data transmission through the CPRI interface and C-RAN applications (e.g., coordinated multipoint and virtual BBU migration) are time-sensitive. Besides the latency of fiber length (transmission latency) and hardware processing time of BBU, which are essential, optical-electrical-optical (O-E-O) conversion is also an important latency contributor, which can cost about 13–15  $\mu\text{s}$  [8]. Therefore, the network operators expect to establish a transparent end-to-end direct lightpath without O-E-O conversions at intermediate nodes. However, establishing an end-to-end lightpath is inefficient in WDM networks because a small traffic volume cannot fill a WDM channel (e.g., 50 GHz), which results in the waste of spectral resources and transceiver resources.

Although electrical layer traffic grooming can improve network resource utilization in WDM networks, it also brings O-E-O conversions at intermediate nodes. In flexible grid optical networks, a direct lightpath can be established by using a sliceable bandwidth-variable transceiver (S-BVT) [9, 10]. The functionality of S-BVT is shown in Fig. 2b. A physical transceiver can be logically “sliced” into multiple sub-transceivers (or virtual transceivers), and each sub-transceiver can set up an independent lightpath for a connection with a “just-right-size” transceiver resource. With the benefits of S-BVT, the BBU cloud can be connected with each other through a transparent end-to-end lightpath, which can significantly reduce the network latency.

### OPTICAL LAYER VIRTUALIZATION FOR A MULTI-TENANCY SCENARIO

C-RAN typically comes with the virtualization concept, because the BBUs are centralized in a virtual resource pool. Since the cloud BBUs are



interconnected by an optical transport network, optical layer virtualization is also attracting much attention. Conventional WDM networks are tightly integrated with the underlying physical substrate (i.e., wavelength) and provide only a rigid granularity optical bandwidth of lightpath. This is inefficient when a virtual connection cannot fill the entire capacity of the wavelength, making it difficult to fully exploit the virtualization concept. Flexible grid optical networking is regarded as a promising technology for virtualization in the optical domain [11]. From the networking level viewpoint, a wide range of spectrum can be segmented to some finer granularities (i.e., spectrum slots) as sharable resources without the constraint of a fixed grid. From the equipment level viewpoint, a sliceable optical transceiver can be sliced into multiple virtual transceivers, and each virtual transceiver has the same ability as a physical transceiver (e.g., establishing a lightpath). The key benefit of the optical layer virtualization brought by flexible grid technology is the openness of the physical layer resources, which can be sliced and abstracted by upper-layer users. With this benefit, multiple wireless tenants can not only share an optical substrate to pay as you grow, but also obtain differentiated services by accessing different functionalized BBUs through virtual optical paths. Optical layer virtualization can coordinate with radio layer virtualization and BBU virtualization to achieve a complete virtualized C-RAN.

## SOFTWARE-DEFINED TECHNOLOGY FOR OPTICAL-ENABLED C-RAN

The cloudified and optical-enabled RAN is a multi-resource integrated environment, in which radio resources, optical resources, and BBU resources coexist. Therefore, the control and management of heterogeneous resources in the C-RAN becomes a big challenge for network operators. The key feature of the SDN-enabled control plane for 5G mobile networks has been discussed to support efficiently a heterogeneous set of resources integration with flexibility [12]. By exploiting SDN, two convergences in the general C-RAN fields are prominent:

1. Vertical convergence: convergence of multiple layers of the network stack (e.g., radio domain and optical domain)
2. Horizontal convergence: convergence of BBU and networking resources

### VERTICAL CONVERGENCE

From the network viewpoint, optical-enabled C-RAN can be seen as a multi-layer networking scenario, in which the optical layer is a carrier layer to the radio layer. Radio signals are sent to a BBU through a wavelength or a continuous spectrum path assigned by the network operators. The ultimate goal of multi-layer networking is to satisfy applications' quality of service (QoS) requirements while keeping the cost as low as possible. Multiple layers of networking should be jointly optimized and operated to achieve this ultimate goal. One important enabler of multi-layer convergence is SDN technology, which can

jointly consider the spectrum and radio frequency assignment (SRA) in a centralized way. SRA schemes can be implemented in the SDN controller to perform different optimization objects for the purposes of the network operator. For example, on the user's side, we expect to reuse the radio frequencies as much as possible without causing interference between adjacent RRHs. In the optical transport layer, we expect to use fewer spectral resources to carry more radio signals under the optical layer constraints (spectrum continuity and contiguity). After the SRA scheme is executed, the SDN controller can program both radio and optical devices through the southbound interfaces (e.g., extended OpenFlow protocol).

### HORIZONTAL CONVERGENCE

Another prominent trend is horizontal convergence. As the RAN shifts to a cloud-based architecture, the baseband processing functions are separated from a conventional base station to a centralized BBU cloud through an optical transport network. This paradigm shift is forcing the convergence of BBU and networking resources. In today's RAN architecture, BBU resources use the network infrastructure as a "dump pipe," which carries a significant amount of traffic but generates low BBU resource utilization. In addition, a network constructor has little knowledge of BBU resources distribution especially when traffic is dynamic, so it is hard in network planning to put the bandwidth where the BBU resources are. To achieve the global optimization goal, the network operators expect to jointly consider BBU and networking resources (radio and optical resources). SDN is a promising solution that aims to orchestrate multiple types of resources in a coordinated way. There are many use cases of horizontal convergence, such as dynamic network load balancing in the case of BBU overload and elastic lightpath provision for BBU cloud interconnection, which are experimented and discussed. Based on the concept of network virtualization and BBU virtualization, an SDN controller can abstract a physical C-RAN topology to a virtual resource graph. According to this graph, a controller can choose a proper network connection from an RRH to a BBU. Figure 3 shows a converged optical and cloud RAN enabled by software defined technology.

## EXPERIMENTAL SETUP OF SOFTWARE DEFINED ELASTIC LIGHTPATH PROVISION FOR BBU CLOUD INTERCONNECTION

Due to the mobility of end users, the bandwidth requirement of the optical channel in a BBU cloud is time-varying. One of the benefits brought by flexible grid optical network is the elastic lightpath provision, in which the spectral resources can be assigned flexibly according to the traffic demand. In this section, we experimentally demonstrate dynamic spectrum allocation and adjustment with increasing C-RAN traffic.

*As the RAN shifts to a cloud-based architecture, the baseband processing functions are separated from a conventional base station to a centralized BBU cloud through an optical transport network. This paradigm shift is forcing the convergence of BBU and networking resources.*

Lightpath adjustment application software, which runs in the network monitor, interacts with the SDN controller through a northbound interface. The SDN controller communicates with OpenFlow-enabled flexible grid optical nodes by using extended OpenFlow protocol

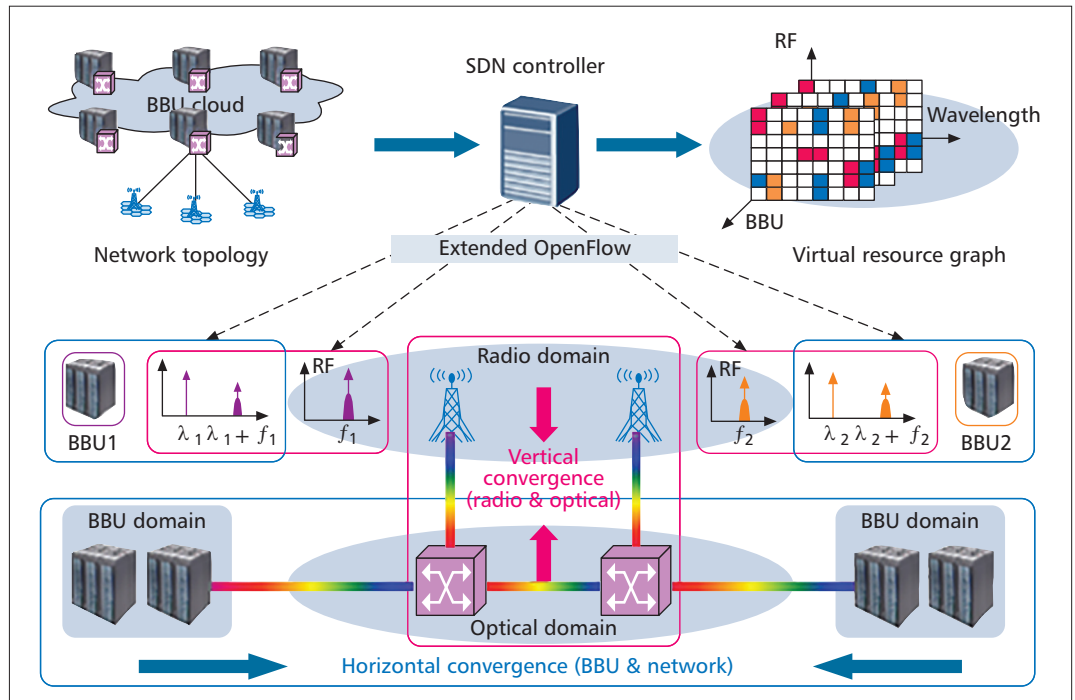


Figure 3. Converged optical and cloud radio access network with software-defined technology.

### EXPERIMENTAL SCENARIO

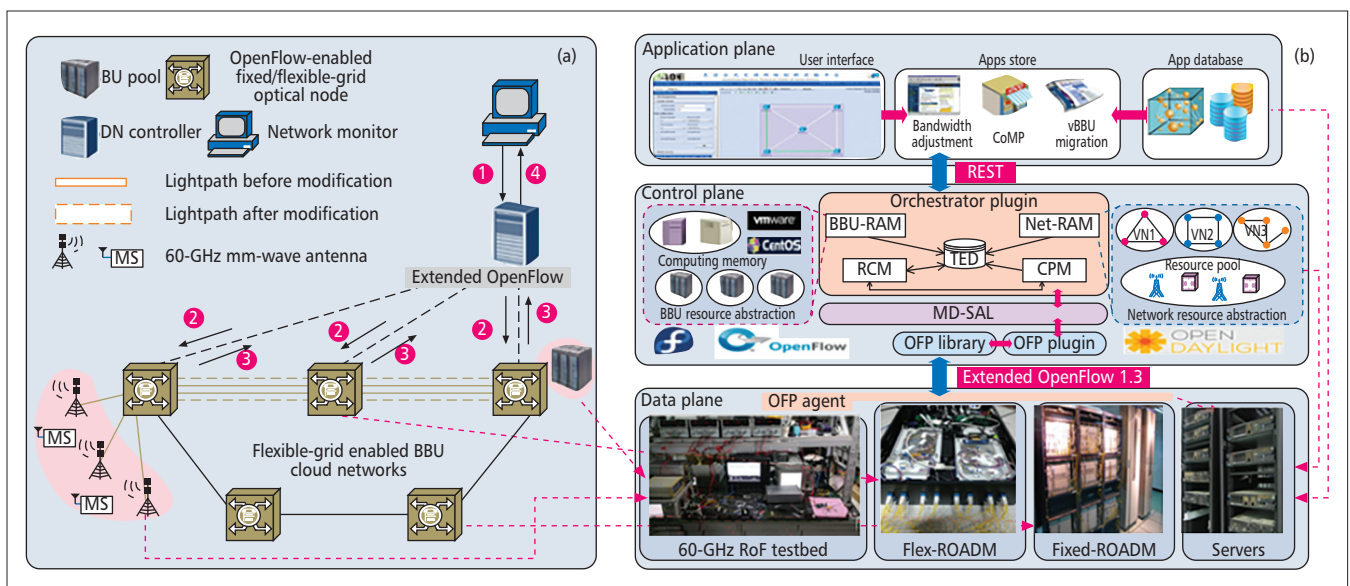
As shown in Fig. 4a, we set up a five-node network topology for BBU cloud interconnection; multiple interleaved 60 GHz mm-Wave signals are generated to access a BBU pool over an elastic lightpath. We simulate the time-varying property of traffic with an increasing step of 1 Gb/s/min. When the bandwidth requirement reaches the limit of an RF signal's capacity (1 Gb/s here), the controller adjusts the lightpath size and assigns a new RF onto the extended lightpath. Lightpath adjustment application software, which runs in the network monitor, interacts with the SDN controller through a northbound interface. The SDN controller communicates with OpenFlow-enabled flexible grid optical nodes by using extended OpenFlow protocol (OFP). The details of signaling procedures in Fig. 4a are described as follows:

- Step 1:** The network monitor sends a bandwidth adjustment request message to the controller.
- Step 2:** The SDN controller performs a spectrum allocation scheme, and sends lightpath modification request messages to OpenFlow-enabled optical switches.
- Step 3:** The OpenFlow-enabled optical switches return lightpath modification reply messages to the SDN controller after the lightpath is modified.
- Step 4:** The controller returns a bandwidth adjustment reply message to the network monitor.

### EXPERIMENTAL SETUP AND RESULTS

We build an SDN-based analog radio-over-flexible-grid optical networks testbed to experimentally verify the elastic lightpath provision for BBU cloud interconnection. As shown in Fig. 4b, the testbed includes three planes: the appli-

cation plane, control plane, and data plane. The original concept of SDN is to separate the control plane and data plane through an open southbound interface (SBI), and the application plane can manage/control the network through the northbound interface (NBI) of the control plane. The SBI and NBI are implemented by using an extended OFP 1.3 and representational state transfer (REST) application programming interface (API). In this experiment, the bandwidth adjustment application is developed as software running on an IBM server (IBM X3650) in the application plane. The network operators and users can access the application through a user interface. For the control plane, we develop an SDN controller based on OpenDaylight, which is an open source platform for SDN. Based on OpenDaylight, we can easily create a network function as a plug-in or a bundle of the controller by using a standardized model language (e.g., YANG). As shown in Fig. 4b, the OpenDaylight-based control plane consists of three parts. The model-driven service abstraction layer (MD-SAL) is the heart of the OpenDaylight controller; it provides the functions to abstract the plug-ins with their exposed features, and allows dynamic linking of the plug-ins that associate with it. Two developed plug-ins are introduced as follows. The OFP plug-in is to code/decode original or extended OFP messages, and interacts with the OFP library, which specifies the OFP message format. The other one is the orchestrator plug-in, which is the brain of the controller. The orchestrator plug-in includes five modules. The BBU resource abstraction module (BBU-RAM) abstracts BBU resources (e.g., analog/digital signal processing resources) and tells the traffic engineering database (TED) to update the BBU availability. Similarly, the network resource abstraction module (Net-RAM)



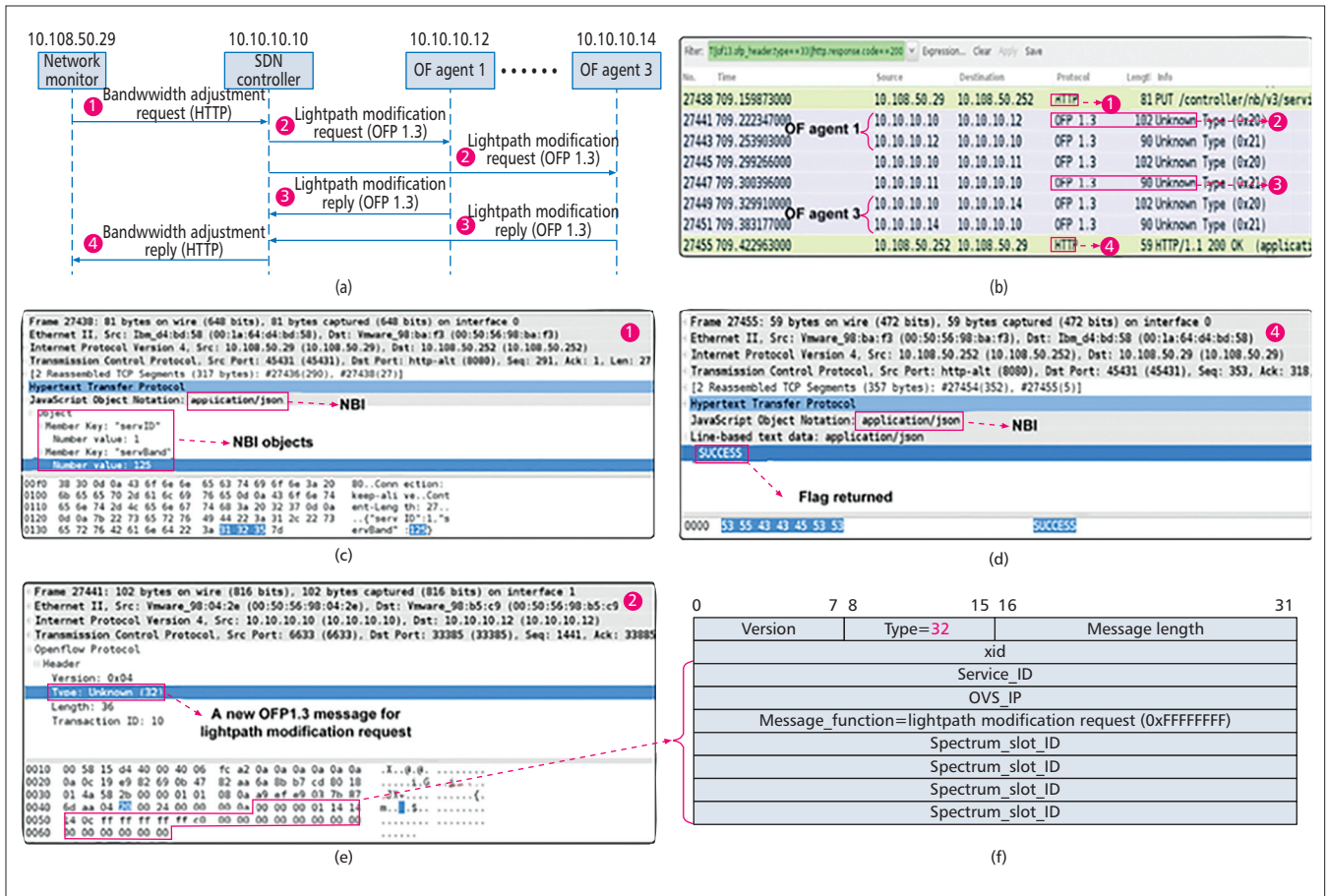
**Figure 4.** a) Elastic lightpath provision for C-RoFlex; b) experimental testbed and demonstrator setup.

virtualizes the radio and optical network resources (e.g., spectral and radio frequency resources), and informs the TED when the network environment changes. The connection provision module (CPM) interacts with the OPF plug-in to provide end-to-end connections related to an OpenFlow message. Upon receiving a connection request, the CPM instructs the resource computing module (RCM) to perform resources computation and allocation; it also tells the TED to update the network condition when the request is successfully provisioned. Different SRA schemes can be implemented in the RCM to determine how to assign optical spectrum and radio frequency according to different purposes (e.g., assigning the minimal number of spectrum slots to support the maximal radio coverage). In this experiment, we only configure the spectrum and radio frequency assignment by using OFP, while the BBU resources (60 GHz mm-Wave system) are preset before the experiment. For the data plane, we set up a five-node topology with three commercial fixed reconfigurable optical add-drop multiplexers (ROADMs) and two 4° OpenFlow-enabled flex-ROADMs which are developed with four Finisar BV-WSSs (Fig. 4b). The mm-Wave signals are generated by a 60 GHz RoF testbed, in which a phase modulator (PM) is driven by a microwave source working at 30 GHz frequency to generate double sideband. A 1 Gb/s electrical signal (custom signal) is added to one of the sidebands by an intensity modulator (IM), and an interleaver is used to separate the sidebands. At the receiving part, the optical signal is detected by a high-speed photodiode, and two rectangular horn antennas with a gain of 20 dBi, frequency range of 50–70 GHz are used to broadcast/receive the mm-Wave. Each physical node (optical switch and mm-Wave laser source) is attached with an OpenFlow agent, which communicates with the controller through extended OFP. We use Open vSwitch (OVS, v. 2.0.0) as an OpenFlow agent. Both the SDN controller and OpenFlow agents are developed in different virtual machines with

independent IP addresses. The outputs of the signaling procedure and protocol extensions are shown in Fig. 5.

Figure 5a shows the interworking procedures for lightpath adjustment. The network monitor, which has installed a bandwidth adjustment application, initiates a session by sending a *BANDWIDTH\_ADJUSTMENT\_REQUEST* message (#1) to an SDN controller. The SDN controller will execute a spectrum allocation scheme to obtain a range of free spectrum, and send a *LIGHTPATH\_MODIFICATION\_REQUEST* message (#2) to inform the OpenFlow-enabled optical switches to change the optical filters' status. After an optical switch finishes the action, it sends a *LIGHTPATH\_MODIFICATION\_REPLY* message (#3) back to the SDN controller. When the SDN controller receives the reply messages from all the optical switches along the lightpath, it returns a *BANDWIDTH\_ADJUSTMENT\_REPLY* message (#4) to the network monitor to terminate the session. The controller with IP address 10.10.10.10 sends *LIGHTPATH\_MODIFICATION\_REQUEST* messages to three OVSs with IP addresses 10.10.10.11, 10.10.10.12, and 10.10.10.14, respectively. In general, there are two methods of interaction between a controller and OVSs. One is the sequential method, in which a controller sends messages one by one to the OVSs. The other is the synchronous method, in which a controller sends messages simultaneously to all OVSs. For this experiment, we use the synchronous method, by which the controller will not send a *LIGHTPATH\_MODIFICATION\_REQUEST* until it receives a *LIGHTPATH\_MODIFICATION\_REPLY* from the last OVS. Figure 5b shows the signaling trace of the interworking procedures corresponding to Fig. 5a.

The JavaScript Object Notation (JSON) API is chosen as the NBI, since it is easy for humans to read and write. Two types of JSON-API are implemented: *BANDWIDTH\_ADJUSTMENT\_REQUEST* and *BANDWIDTH\_ADJUSTMENT\_REPLY*.



**Figure 5.** a) Interworking procedures for lightpath adjustment; b) control message trace of interworking procedures; (c–e) wire-shark captures of message #1, #2 and #4; f) TLV for lightpath adjustment message.

*REPLY*, which are shown in Figs. 5c and 5d, respectively. Two objects are added to the bandwidth adjustment request message, which are service ID and service bandwidth. One object is added for the bandwidth adjustment reply message, which is status flag (i.e., success or failure). For the SBI, OFP 1.3 is extended for the optical requirement. Instead of extending an existing OpenFlow message such as FLOW\_MOD, we create new messages attached to an OFP 1.3 header. Adding a new message is easily implemented in an OpenDaylight-based platform by using YANG tools. Figures 5e and 5f show a new OFP 1.3 message (type 32) for a lightpath modification request and its Type-Length-Value (TLV) extension. Besides the OFP 1.3 header, we add 32 bits for service ID to present a service request and 32 bits for OVS IP address to present an OVS host IP address. We add 32 bits for message function to indicate the different functions of the message (e.g., lightpath modification function). We also add 128 bits for spectrum slot ID to describe the spectrum usage for 128 spectrum slots (e.g., 1 is busy and 0 is free).

The filter outputs of lightpath adjustment for analog radio over a flexible grid optical network are shown in Figs. 6a and 6b. Three 60-GHz-band signals are multiplexed onto an elastic optical channel, and this optical channel can be adjusted to carry four mm-Wave signals by the SDN controller. Figure 6c shows the measured

bit error rate (BER) vs. received optical power of on/off keying (OOK) data carried on 60 GHz mm-Wave with 1.2 m wireless transmission. We also report the overall latency, which is shown in Fig. 6d. The overall latency includes three parts: data plane latency, control plane latency, and application plane latency. The data plane latency is the hardware response delay of a flexROADM, which was verified to be spectrum-dependent in our previous work [13]. The control plane latency is the time slot during the controller receiving message #1 and message #3 in Fig. 5a. The control plane latency includes the algorithm process delay and the SBI process delay between controller and OpenFlow agents. The application plane latency is the time slot during the network monitor sending out message #1 and receiving message #4. It is shown to be a major contributor to the overall latency, because the TCP stack delay and application software delay are large relative to the latencies described above. The presented overall latency is different from the latency required by a CPRI interface, which physical layer constraint dominates; it reflects how fast the customer can program an SDN-enabled network by using application software (user interface). It is an important value for network operators to evaluate the dynamic performance of SDN technology. Software/algorithm running time is a major part of the presented latency. Possible ways of

improving the system latency are to design optimal schemes with respective protocol layers from end to end, and grid computing and in-memory caching are also effective ways to accelerate the processing speed.

## CONCLUSION

The demands of 5G mobile networks require C-RAN to support a large, scalable, flexible switching, and highly reliable network architecture for BBU cloud interconnection. In this article, we introduce flexible grid optical networks and the enabling technologies for an inter-BBU cloud network. The benefits brought by flexible grid technology for C-RAN are discussed. To improve the network intelligence, we propose a software-defined centralized control plane to coordinate the multi-layer and heterogeneous network environment, which consists of BBU resources, radio resources, and optical resources. The elastic lightpath provision for C-RoFlex was experimented in a SDN-based testbed. The protocol extensions and testbed performance are reported. Through the experiment, it was verified that elastic lightpath provision as a C-RAN service can be achieved by exploiting SDN.

## ACKNOWLEDGMENT

This work has been supported by China's 973 program (2012CB315705), NSFC project (61372118), and Ministry of Education-China Mobile Research Foundation (MCM20130132). The authors would like to thank the Editor and reviewers for their detailed reviews and constructive comments, which have helped to improve the quality of this article.

## REFERENCES

- [1] J. G. Andrews et al., "What Will 5G Be?," *IEEE JSAC*, vol. 32, no. 6, June 2014, pp. 1065–82.
- [2] A. Checko et al., "Cloud RAN for Mobile Networks — A Technology Overview," *IEEE Commun. Surveys & Tut.*, vol. 17, no. 1, Mar. 2015, pp. 405–26.
- [3] J. Beas et al., "Millimeter-Wave Frequency Radio over Fiber Systems: A Survey," *IEEE Commun. Surveys & Tut.*, vol. 15, no. 4, 2013, pp. 1593–1619.
- [4] M. Jinno et al., "Spectrum-Efficient and Scalable Elastic Optical Path Network: Architecture, Benefits, and Enabling Technologies," *IEEE Commun. Mag.*, vol. 47, no. 11, Nov. 2009, pp. 66–73.
- [5] O. Gerstel et al., "Elastic Optical Networking: A New Dawn for the Optical Layer?," *IEEE Commun. Mag.*, vol. 50, no. 2, Feb. 2012, pp. S12–S20.
- [6] ITU-T Rec. G.694.1, "Spectrum Grids for WDM Applications: DWDM Frequency Grid," 2012.
- [7] J. Zhang et al., "First Demonstration of Enhanced Software Defined Networking (eSDN) over Elastic Grid (eGrid) Optical Networks for Data Center Service Migration," *Proc. OFC/NFOEC '13*, PDP5B.1.
- [8] S. Zhang et al., "Connecting the Clouds with Low-Latency, Low-Cost Virtual Private Lines Enabled by Sliceable Optical Networks," *Proc. IEEE GLOBECOM '13*, pp. 2370–75.
- [9] M. Jinno et al., "Multiflow Optical Transponder for Efficient Multilayer Optical Networking," *IEEE Commun. Mag.*, vol. 50, no. 5, May 2012, pp. 56–65.
- [10] J. Zhang et al., "Energy-Efficient Traffic Grooming in Sliceable-Transponder-Equipped IP-Over-Elastic Optical Networks [Invited]," *J. Opt. Commun. Networks*, vol. 7, no. 1, Jan. 2015, pp. A142–52.
- [11] M. Jinno et al., "Virtualization in Optical Networks from Network Level to Hardware Level [Invited]," *J. Opt. Commun. Net.*, vol. 5, no. 10, Oct. 2013, pp. A46–56.
- [12] R. Trivisonno et al., "SDN-Based 5G Mobile Networks: Architecture, Functions, Procedures and Backward Compatibility," *Trans. Emerging Telecommun. Tech.*, vol. 26, no. 1, Jan. 2015, pp. 82–92.

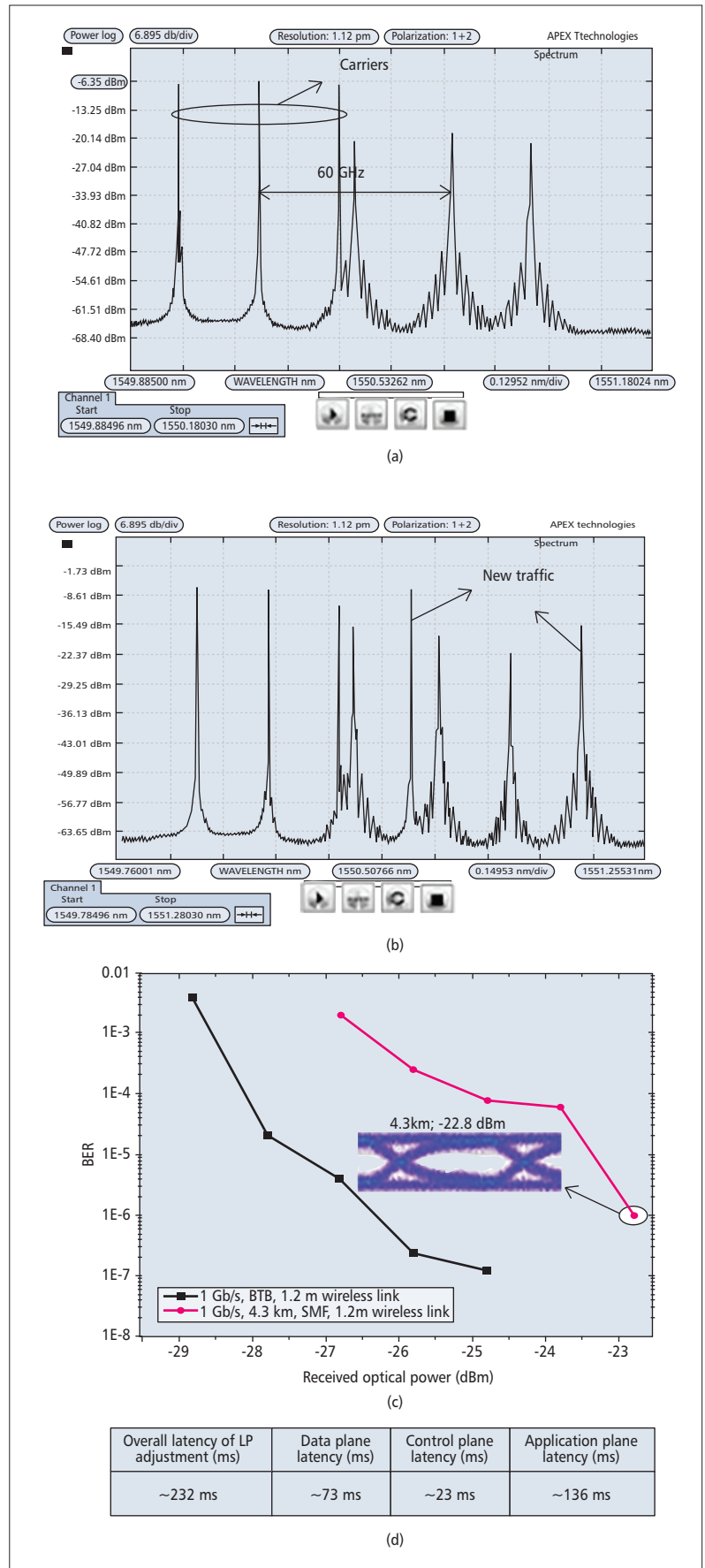


Figure 6. a) Lightpath before adjustment; b) lightpath after adjustment; c) BER performance of 1 Gb/s signal over fiber; d) overall latency of dynamic lightpath provision.

- [13] J. Zhang *et al.*, "Experimental Demonstration of Open-Flow-Based Control Plane for Elastic Lightpath Provisioning in Flexi-Grid Optical Networks," *Opt. Exp.*, vol. 21, no. 2, 2013, pp. 1364–73.

## BIOGRAPHIES

JIawei ZHANG received his Ph.D. degree from the State Key Lab of Information Photonics and Optical Communication at Beijing University of Posts and Telecommunications (BUPT) in 2014. He was a visiting student at the University of California, Davis for 15 months. He is currently a post-doctoral researcher at BUPT. He has served on the Technical Program Committees for workshops at IEEE GLOBECOM '14 and IEEE ICC '15. His research focuses on software-defined 5G mobile networks, with emphasis on front/backhaul networks based on various optical communication technologies.

YUEFENG Ji [SM] received his Ph.D. degree from BUPT. Now he is a professor and deputy director of the State Key Lab of Information Photonics and Optical Communications. His research interests are primarily in the area of broadband communication networks and optical communications, with emphasis on key theory, realization of technology, and applications. He is a Fellow of China Institute of Communications.

JIE ZHANG is a professor and vice dean of the Institute of Information Photonics and Optical Communications at BUPT. He is sponsored by over 10 projects of the Chinese government. He has published eight books and more than 200 articles. Seven patents have also been granted. He has served as Co-Chair of ACP '12/'13/'14, ICOCN '14, ChinaCom '12/'13, and so on. His research focuses on optical transport networks, packet transport networks, and so on.

RENTAO GU received his B.E. and Ph.D. degrees from BUPT in 2005 and 2010, respectively. From 2008 to 2009, he was a visiting scholar at Georgia Institute of Technology. He is an associate professor of the School of Information and Communication Engineering. His current research interests include software defined networking and broadband access networks. He is a senior member of the China Institute of Communications and Chinese Institute of Electronics.

YONGLI ZHAO is currently an associate professor of the Institute of Information Photonics and Optical Communications at BUPT. He received his B.S. degree in communication engineering and Ph.D. degree in electromagnetic field and microwave technology in 2005 and 2010, respectively. More than 150 journal and conference articles have been published. His research focuses on software defined optical networks, flexi-grid optical networks, network virtualization, and so on.

SIMING LIU received his B.E. degree in electronic and information engineering from Tianjin University of Technology, China, in 2013. He is currently working toward a Ph.D. degree at the School of Information and Communication Engineering, BUPT. His research interests are microwave photonics and digital processing.

KUN XU, professor is a TPC member of IEEE MTT-3, deputy director of the State Key Laboratory of Information Photonics and Optical Communication, vice dean of the Science School, and director of RF Photonics Laboratory at BUPT. His research interests include microwave photonic devices and integrated systems. He has served as TPC and Workshop/Session Co-Chair for several international conferences, including IEEE/OSA/SPIE ACP, IEEE MWP/APMP, IEEE GSMM, and IEEE ICC.

MEI SONG received her B.S. and M.S. degrees in electric engineering from Tianjin University in 1983 and 1986, respectively. She is currently a professor with BUPT. Her research interests are in future communication, wireless broadband networks, and heterogeneous networks.

HAN LI received his Ph.D. degree from BUPT. He is currently a vice director of China Mobile Research Institute. His research focuses on next generation mobile communication networks.

XINBO WANG is currently pursuing his Ph.D. degree in computer science from the University of California, Davis. He received his B.S. in Telecommunication Engineering from BUPT in 2013, and his M.S. in computer science from the University of California, Davis in 2015. His research interests include next generation mobile fronthaul, next generation Ethernet passive optical networks, and their integration with cloud radio access networks.

**CALL FOR PAPERS**  
**IEEE COMMUNICATIONS MAGAZINE**  
**WIRELESS COMMUNICATIONS, NETWORKING, AND POSITIONING WITH**  
**UNMANNED AERIAL VEHICLES**

**BACKGROUND**

Enabled by the advances in computing, communication, and sensing as well as the miniaturization of devices, unmanned aerial vehicles (UAVs) such as balloons, quadcopters, and gliders have been receiving significant attention in the research community. Indeed, UAVs have become an integral component in several critical applications such as border surveillance, disaster monitoring, traffic monitoring, remote sensing, and the transportation of goods, medicine, and first-aid. More recently, new possibilities for commercial applications and public service for UAVs have begun to emerge, with the potential to dramatically change the way in which we lead our daily lives. For instance, in 2013, Amazon announced a research and development initiative focused on its next-generation Prime Air delivery service. The goal of this service is to deliver packages into customers' hands in 30 minutes or less using small UAVs, each with a payload of several pounds. 2014 has been a pivotal year that has witnessed an unprecedented proliferation of personal drones, such as the Phantom and Inspire from DJI, AR Drone and Bebop Drone from Parrot, and IRIS Drone from 3D Robotics.

Among the many technical challenges accompanying the aforementioned applications, leveraging the use of UAVs for delivering broadband connectivity plays a central role in next generation communication systems. Facebook and Google announced in 2014 that they will use a network of drones that circle in the stratosphere over specific population centers to deliver broadband connectivity. Such solar-powered drones are capable of flying for several years without refueling. UAVs have also been proposed as an effective solution for delivering broadband data rates in emergency situations through low-altitude platforms. For example, the ABSOLUTE, ANCHORS, and AVIGLE projects in Europe have been investigating the use of aerial base stations to establish opportunistic links and ad hoc radio coverage during unexpected and temporary events. They can serve as a temporary, dynamic, and agile infrastructure for enabling broadband communications, and quickly localizing victims in case of disaster scenarios.

This proposed Feature Topic (FT) issue will gather articles from a wide range of perspectives in different industrial and research communities. The primary FT goals are to advance the understanding of the challenges faced in UAV communications, networking, and positioning over the next decade, and provide further awareness in the communications and networking communities of these challenges, thus fostering future research. Original research papers are to be solicited in topics including, but not limited to, the following themes on communications, networking, and positioning with UAVs:

- Existing and future communication architectures and technologies for small UAVs
- Delay-tolerant networking for cooperative UAV operations
- Design and evaluation of wireless UAV test beds, prototypes, and platforms
- Multihop and device-to-device communications with UAVs
- Interfaces and cross-platform communications for UAVs
- QoS mechanisms and performance evaluation for UAV networks
- Game-theoretic and control-theoretic mechanisms for UAV communications
- Use of civilian networks for small UAV communications
- Integrating 4G and 5G wireless technologies into UAV communications, such as millimeter wave communications, beamforming, moving networks, and machine type communications
- Use of UAVs for public safety and emergency communications, networking, and positioning
- Integration of software defined radio and cognitive radio techniques with UAVs
- Channel propagation measurements and modeling for UAV communication channels

**SUBMISSIONS**

Articles should be tutorial in nature, with the intended audience being all members of the communications technology community. They should be written in a style comprehensible to readers outside the specialty of the article. Mathematical equations should not be used (in justified cases up to three simple equations are allowed). Articles should not exceed 4500 words (from introduction through conclusions). Figures and tables should be limited to a combined total of six. The number of references is recommended not to exceed 15. In some rare cases, more mathematical equations, figures, and tables may be allowed if well justified. In general, however, mathematics should be avoided; instead, references to papers containing the relevant mathematics should be provided. Complete guidelines for preparation of the manuscripts are posted at <http://www.comsoc.org/commag/paper-submission-guidelines>. Please send a pdf (preferred) or MSWORD formatted paper via Manuscript Central (<http://mc.manuscriptcentral.com/commag-ieee>). Register or log in, and go to Author Center. Follow the instructions there. Select "May 2016 / Wireless Communications, Networking and Positioning with UAVs" as the Feature Topic category for your submission.

**SCHEDULE FOR SUBMISSIONS**

- Submission Deadline: November 1, 2015
- Notification Due Date: January 15, 2016
- Final Version Due Date: March 1, 2016
- Feature Topic Publication Date: May 2016

**GUEST EDITORS**

Ismail Guvenc  
Florida International Univ., USA  
iguvenc@fiu.edu

Walid Saad  
Virginia Tech, USA  
walids@vt.edu

Mehdi Bennis  
Univ. of Oulu, Finland  
bennis@ee.oulu.fi

Christian Wietfeld  
TU Dortmund Univ., Germany  
christian.wietfeld@tu-dortmund.de

Ming Ding  
NICTA, Australia  
ming.ding@nicta.com.au

Lee Pike  
Galois, Inc., USA  
leepike@galois.com

# Network Coded Software Defined Networking: Enabling 5G Transmission and Storage Networks

Jonas Hansen, Daniel E. Lucani, Jeppe Krigslund, Muriel Médard, and Frank H. P. Fitzek

## ABSTRACT

Software defined networking has garnered large attention due to its potential to virtualize services in the Internet, introducing flexibility in the buffering, scheduling, processing, and routing of data in network routers. SDN breaks the deadlock that has kept Internet network protocols stagnant for decades, while applications and physical links have evolved. This article advocates for the use of SDN to bring about 5G network services by incorporating network coding (NC) functionalities. The latter constitutes a major leap forward compared to the state-of-the-art store and forward Internet paradigm. The inherent flexibility of both SDN and NC provides fertile ground to envision more efficient, robust, and secure networking designs, which may also incorporate content caching and storage, all of which are key challenges of the upcoming 5G networks. This article not only proposes the fundamentals of this intersection, but also supports it with key use cases and a thorough performance evaluation on an implementation that integrated the Kodo library (NC) into OpenFlow (SDN). Our results on single-hop, multihop, and multi-path scenarios show that gains of  $3\times$  to  $11\times$  are attainable over standard TCP and multi-path TCP.

## INTRODUCTION

Communication and networking systems have been structured in a series of layers to ease design and, in principle, allow for novel technologies and services to be incorporated with minimal effort. Although this has been true for higher layers dealing with applications and services (e.g., incorporating multimedia content, social networking, and cloud applications) and lower layers dealing with access technologies (e.g., evolving from early mobile cellular technologies to 4G), the intermediate network layer control protocols have been stagnant for several decades. In fact, this deadlock is due in part to the fact that multiple technologies depend on the current state of the

network to operate [1] and in part to the challenges of configuring and managing networks based on network element operations instead of the network services to be provided. The downside to the current network design and protocols is that they were conceived with a fairly narrow set of goals, which is now limiting the effectiveness and feasibility of the more complex and resource demanding applications of the upcoming 5G networks. For example, the lack of a holistic view of the network can translate into reduced efficiency and higher congestion in specific paths or the inability to use multi-path techniques to provide ultra-reliable services. This structure also limits the adoption of verifiably optimal multicast techniques, such as network coding (NC), at the core of the network.

As a way to deal with some of the inefficiencies derived from these drawbacks, networking theory and the Internet are currently facing two large, game-changing trends: software defined networking (SDN) and NC. Each provides a disruptive concept to enable more efficient and flexible networking, and both have experienced a proliferation of academic and industrial applications. The goal of this article is to bridge the two concepts through fundamental analysis and understanding of their joint potential as well as testing their combination with a real implementation. The combination of both ideas goes beyond marrying two buzz words. SDN offers fertile ground to implement NC ideas and make them widely and readily available, as well as to help guide network coding research in terms of realistic restrictions, such as complexity, memory, and in-network routers, as well as routers capabilities (e.g., changing data operations on the fly). In other words, SDN is the key enabler for the timely implementation and deployment of NC, initially as a virtual function next to each communication node and in the future as an integrated part of the SDN system.

This article provides key examples to demonstrate the potential of NC's recoding capabilities and ability to mask losses and stabilize lossy links, even without altering the end-to-end use

---

Jonas Hansen and Daniel E. Lucani are with Aalborg University. Daniel Lucani is also with Chocolate Cloud ApS.

Jeppe Krigslund is with Steinwurf ApS.

Muriel Médard is with the Massachusetts Institute of Technology.

Frank H. P. Fitzek is with TU Dresden.

This work was financed in part by the Green Mobile Cloud project granted by the Danish Council for Independent Research (Grant No. 10-081621).



of TCP. This provides a simple but critical stepping stone to improve current systems and open the path to more complex protocols natively centered around network coding, such as coded TCP (CTCP) [2]. We show that our simple approach provides three-fold to 11-fold gains over standard TCP and multi-path TCP [3].

## SDN AND NC: THE KEY TO 5G SERVICES?

More than providing an evolution in network technologies, 5G networks are envisioned as a revolution. It is not just about increased data rates with a new radio access technology, but rather a large expansion of the network's goals to provide traditional services (e.g., voice, user data), as well as radically new services: machine-to-machine (M2M) communications with support for massive numbers of devices, millisecond latency for communications, cloud and caching services, high reliability, and energy efficiency. This vision requires a system that is able to judiciously allocate resources, treats storage and transport of data as a single process, and exploits the meshed nature of communication networks to guarantee these requirements. This article argues that a combination of SDN and NC is the key to addressing these challenges.

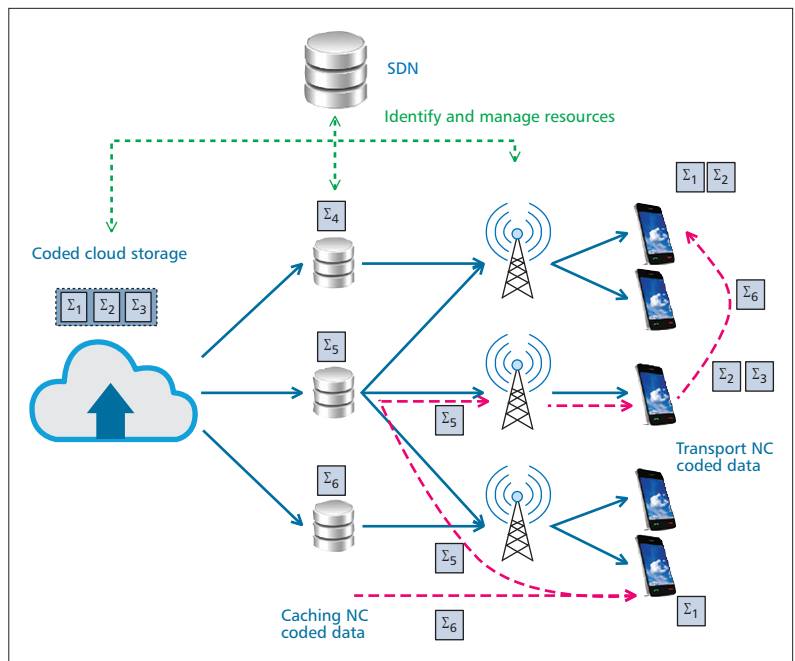
SDN allows for a network to perform flexible resource allocation (e.g., buffer management, dynamic routing, exploitation of multiple paths) beyond a single layer in the network stack (Fig. 1). To achieve this, SDN virtualizes services in the network by separating the data transmission and control of the network. In initial designs, the control plane was achieved with additional controller devices to allow it to evolve rapidly, while maintaining simple and cost-effective switching elements.<sup>1</sup>

Recent work has envisioned more refined and distributed control mechanisms and security aspects of SDN [4]. The de facto standard<sup>2</sup> for the time being for SDN is provided by OpenFlow, a commercially usable platform [1]. To date, OpenFlow supports a series of functionalities. First, it allows the control plane access to:

- The type of connection used (e.g., fiber optic, copper wire) along with negotiated connection speed
- Hardware description and available capabilities, and various statistics to each individual switch, specifically:
  - Port input/output in both packets and bytes
  - Data flow input/output in both packets and bytes
  - Packets dropped in both input and output queues of each port
  - Number of collisions and collision errors detected

Second, the control plane may instruct the data path to:

- Install data flow entries on switches and routers
- Modify existing flow entries on switches and routers
- Request and set features and configurations of switches and routers



**Figure 1.** 5G network with a network coded software defined network. Seamless identification of data source, cache management, device-to-device, or machine-to-machine communications, multi-path support, and multi-source exploitation. The SDN controller has an overall perspective of the network including data sources, edge caches, and base stations, and is able to identify and manage network resources, including computational ones for network coding. The SDN controller can identify coded fragments of files in caches as well as in peer devices for orchestrating the transmission of data to new devices interested in a specific content. Bold lines represent connections, while red dashed lines represent data flows, and green dashed lines represent control plane flows.

- Request update on input/output statistics for individual switches and routers
- Finally, the SDN data plane can inform the control plane about:
- Unidentified packet headers (and thus request a new flow entry)
  - Removed data flows
  - Modified port status
  - Errors in the data path

These functionalities allow a centralized control plane, shown in Fig. 1, to identify coding-capable devices, link characteristics (round-trip delay, packet loss rates), and topology to determine paths to be used and where/how much coding needs to be incorporated. The SDN controller could also identify coded fragments of popular files in caches as well as in peer devices for orchestrating the transmission of data to new devices.

On the other hand, NC breaks with the store-and-forward paradigm used in today's networks by encouraging intermediate nodes in the network to recode incoming data packets using algebraic operations over finite fields (e.g., a middle device in Fig. 1 recoding and sending a new coded packet to the upper devices). This contrasts with standard end-to-end erasure correcting codes (e.g., LT and Raptor codes), and allows the network to generate redundancy where it is needed instead of injecting it from the source. In a way, NC proposes a store-code-forward paradigm to networking. Random linear network coding (RLNC) provides a distributed,

<sup>1</sup> <https://irtf.org/sdnrg>

<sup>2</sup> There are other mechanisms to configure switches or new SDN software projects. For example, OpenDayLight (<http://www.opendaylight.org/>) is a recent and exciting project with a first software release in February 2014. However, OpenFlow was released in 2011 and is already supported in commercial switches.

*The comprehensive view of network conditions that is available through SDN can be pivotal to deploy and manage NC configurations and recoding potential within the network as well as identifying storage locations to bridge users and/or devices to their data.*

asymptotically optimal approach to employ NC. RLNC is based on choosing random coefficients to create linear combinations of incoming packets. The reason behind these gains comes from the fact that:

- The network itself does not need to transport packets without modification, but rather a linear combination of the original data packets, thus providing a richer set of options and actions available to the network.
- The receivers do not need to track individual packets, but instead focus on accumulating enough independent linear combinations in order to recover the original packets.

Although the gains have been shown in a variety of scenarios, and implementations have confirmed NC's potential in practice, NC's incorporation in standards and wide deployment has been limited with some exceptions; for example, CATWOMAN [5] is currently deployed as part of the Linux Kernel.<sup>3</sup> Part of the limitation lies in the difficulty of retrofitting routers and switches in the network with NC functionalities. However, enabling even a limited number of such devices with NC (e.g., new 5G equipment) can have a large impact on performance if we are able to identify and exploit them. SDN can ease this process and other functionalities. Finally, the presence of a high-performance network coding library (kodo [6]) would ease the deployment of a multiplicity of NC strategies. Although the base functionalities are simple — encode, recode, and decode data packets — NC supports a variety of code designs and configurations, from the classical block-by-block RLNC to online NC, which essentially allows the encoder and decoders to use a moving window for deciding which packets to include in the next linear combination. Kodo [6] also supports systematic NC codes, sparse RLNC, perpetual NC, and fulcrum NC [7], providing a wide range of configurations for deployment in networks.

#### CURRENT STATE OF AFFAIRS

Both technologies are currently discussed within the Internet Research Task Force (IRTF) on NC<sup>4</sup> and SDN.<sup>5</sup> However, the communities actively involved in SDN and NC have had little if any overlap in the past, which limits mutual understanding of the challenges and capabilities of each, and limits the opportunities to combine the two in a meaningful way.

#### WHERE CAN SDN AND NC HELP?

The combination of SDN and NC brings forth an interesting potential for the management and operation of 5G networks. In particular, the comprehensive view of network conditions that is available through SDN can be pivotal to deploy and manage NC configurations and recoding potential within the network as well as identifying storage locations to bridge users and/or devices to their data. The following benefits are possible by this combination:

**Exploitation of multiple communication paths:** NC is particularly well suited to exploiting multiple communication interfaces and routes [8], which can then exploit SDN's ability to recognize multiple communication paths between

source and destination. This is key in 5G networks to comply with reliability requirements as well as appropriate management of heterogeneous interfaces, such as millimeter-wave (mmWave) for increased speed, and another technology for continuous connectivity or for M2M purposes as in Fig. 1.

**Management of data storage and caching:** SDN's ability to virtualize and/or identify caching and storage nodes in the network are key to exploiting NC to enhance the impact and reduce the storage cost of caching/storage by relying on linearly coded packets instead of replication of the original data per storage location/device. NC also provides a single code for both storage and data transmission, which is key to treating data as a single holistic process, as seen in Fig. 1. This management can also include the preemptive caching of data of a user as it moves through the network using location information. The goal is to guarantee low latency for access to the user's data.

**Adaptation of redundancy based on link quality:** SDN provides simple mechanisms to identify the link quality, including packet losses, for the transmission routes used for a flow. This capability is particularly relevant for using NC's recoding capabilities to generate the right level of redundancy per link, instead of introducing end-to-end redundancy to compensate for packet losses, which is an inherently inefficient strategy.

**Assessment of system load and complexity allocation:** SDN is useful for identifying whether a device can commit resources to recoding and how many, since the control plane is able to access information about the hardware of each switch. This may allow us to choose the NC parameters to meet the current network demands. This is particularly relevant with novel NC schemes that provide fluid allocation of complexity, such as fulcrum network codes [7], by performing linear combinations using different finite fields end to end and at different nodes in the network. The choice of the finite field has a direct impact on the computational effort required by a given node. This added flexibility is key to dealing with energy efficiency in 5G networks, not only for the infrastructure but for connected machines (e.g., sensors, actuators) and end-user devices. Additionally, an SDN controller could provide a simple trade-off between computation (reduced coding load) and communication (increased communication load) in the network by deciding whether and where to code various flows depending on load statistics obtained by SDN, while utilizing the flexible SDN flow management to accommodate such changing conditions. In a sense, each switch can be configured to code (or not) the incoming flows, as illustrated in Fig. 2.

#### IMPLEMENTATION AND TESTBED

In order to advocate for the integration of NC into SDN, we include a set of simple network topologies that show the benefits of this merging of technologies. The simplicity of the topologies further emphasizes the usability of the merged technologies as large-scale and more complex network scenarios can be broken into a collection of these simple topologies. That is, the benefits found in the simple topologies are readily applica-

<sup>3</sup> <http://lwn.net/Articles/549477/>

<sup>4</sup> <https://irtf.org/nwcrng>

<sup>5</sup> <https://irtf.org/sdnrg>

ble in complex scenarios. The network topologies utilized for this each represent a scenario where NC provides a potential benefit, but where an implementation conflicts with the boundaries set up by the structure of the network protocol stack. Simple network scenarios include single-hop, multi-hop, and multi-path. Each of the scenarios are further elucidated below.

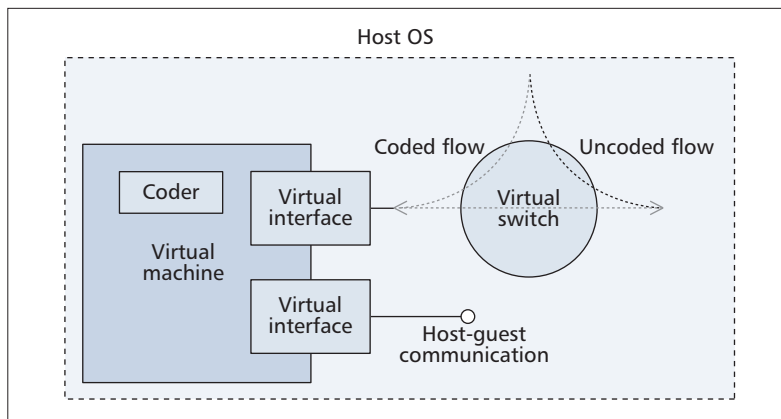
### CODING DEVICES

A functional software defined network is necessary in order to confirm the capabilities of the developed coding in such a network. However, network equipment capable of SDN (e.g., OpenFlow or similar protocols) only provide limited possibilities in terms of modification and configuration of capabilities. We propose a flexible approach to evaluate the potential of the combination of NC with SDN without initiating full integration of NC into SDN software. This is a crucial step not only to experiment with different strategies and schemes but to gain insight as to the key and most promising elements that could be included into OpenFlow or similar projects as well as motivating switch manufacturers to support new coding features.

For this reason, a virtual network environment has been set up using Open vSwitch,<sup>6</sup> an open source virtual switch supporting the OpenFlow protocol [9]. A coding device in the virtual network setup is based on the developed coding software. Instead of integrating the coding software directly into the Open vSwitch devices, the coding is deployed on virtual machines. This is a development decision based on both a limited timeframe for development, and the fact that a virtual machine can be substituted with a real device without changes to the code in case a real network scenario should be deployed. The possibilities of integrating the coding software on existing network equipment, such as a switch, is very limited. A (virtual) machine with the coding software deployed acts as a coding device on the virtual network, and can be used for either encoding, decoding, or recoding. The coding software distinguishes between packets that need to be coded, recoded, or decoded, and packets that should be ignored in terms of coding.

The virtual machine utilizes virtual network interfaces that ensures communication toward the host operating system (OS). That is, within the virtual machine, this device acts as a normal network interface connected to a network, while the host OS has the responsibility of handling traffic to and from this interface. The virtual interface can be included in virtual network scenarios specified by the host OS. Two of these network interfaces are included in each virtual machine. One is to be included in the virtual software defined network scenarios, and one is utilized for direct communication with the host OS. The latter eases the setup and configuration through host-side scripts without “polluting” the investigated network scenario. This relationship between host OS, virtual machine, and coding software is illustrated in Fig. 2.

Deploying the coding in virtual machines naturally curbs the obtainable performance compared to dedicated integration directly in the virtual switch, say, due to the introduction of



**Figure 2.** Integrating network coding into an Open vSwitch, showing the relation between coding software, virtual machine, and host OS. A data flow (the grey dashed line) is redirected through the coding device before it is forwarded onto the network. The flows that do not require coding follow the black dashed line straight through the switch.

additional delay for directing traffic to the virtual machine, processing, and sending back to the switch. However, we argue that a virtual implementation is an equal advocate for NC to be integrated into a software defined network. The virtual setup shows that only a limited amount of coding-related instructions is needed to create a beneficial coding approach that can be deployed on future SDN-capable equipment.

Intermediate network nodes in the constructed virtual network setup consist of an Open vSwitch along with an adjacent coding device. Data flows to be coded are then redirected from the switch through the coding device and back to the switch, which then forwards the coded packets appropriately. The coded data flows are specified in the coding devices using statistics data obtainable from the switch using the capabilities of the OpenFlow protocol, where key supported parameters for implementing efficient network coding are the packet loss statistics and topology information. Although packet loss statistics could be obtained by other means, such as using the Simple Network Management Protocol (SNMP), the SDN framework will allow for further system configuration and control in future developments, including route selection and complexity management for NC functionalities. From the perspective of the SDN controller, the controller would select the use of the coding switch for a specific flow by routing the data to the virtual machine and from the virtual machine back to the router for transmission to the next hop.

In general, this combination of switch and coding device imitates the behavior of a network-coding capable network node. This could be deployed either as an overlay, for example, computation on top of existing (but limited) SDN-capable switches, or within the switch itself. Figure 2 provides a simple illustration of how the two components, switch and coding device, cooperate to create coded flows within a network.

### CODING SOFTWARE

The software implementation<sup>7</sup> uses Kodo [6], which is a C++11 network coding library capable of random linear network coding (RLNC).

<sup>6</sup> <http://www.openvswitch.org>

<sup>7</sup> <https://github.com/14gr1010/software>

The coding scheme is used as forward error correction (FEC), that is, neither positive nor negative acknowledgement (ACK/NACK) is used to ensure delivery of every packet. The encoder uses a systematic code [10] where all the original packets are transmitted uncoded (but with a header added by Kodo and zero-padded to the RLNC symbol size) the first time. Throughout the performance measurements we use a generation size of 10 symbols, Galois field in use is  $GF(2^8)$ , and the symbol size is set to 1356 bytes. Some of the benefits of a systematic code is lower decoding complexity and delay [10]. Lower decoding delay is achievable since all uncoded packets can be forwarded directly, and having uncoded packets reduces the decoding complexity of Gaussian elimination, which is used to decode the RLNC generation. Kodo supports RLNC on the fly where packets can be added to the encoder as they arrive. Similarly, packets can be extracted from the decoder as soon as they are decoded, without the need for decoding the entire generation first. One of the benefits of RLNC on the fly is that coded packets can be transmitted before the entire generation is fed to the encoder (i.e., adding redundancy on the fly). The combination of systematic and on-the-fly coding allows the encoder to transmit uncoded and coded packets with minimal delay. However, using on-the-fly coding also reduces the decoding probability compared to a traditional block code. Since redundancy packets can only aid decoding the symbols that were added to the encoder at the time of their creation, they may become useless later. For example, a redundancy packet created at the

beginning of a generation cannot be used to recover a lost symbol in the last part of a generation. But without on-the-fly coding, slow or infrequent communication is at best problematic; for example, in a ping scenario it would be inconvenient to wait for an entire generation to be filled before sending any packets.

## MEASUREMENT RESULTS

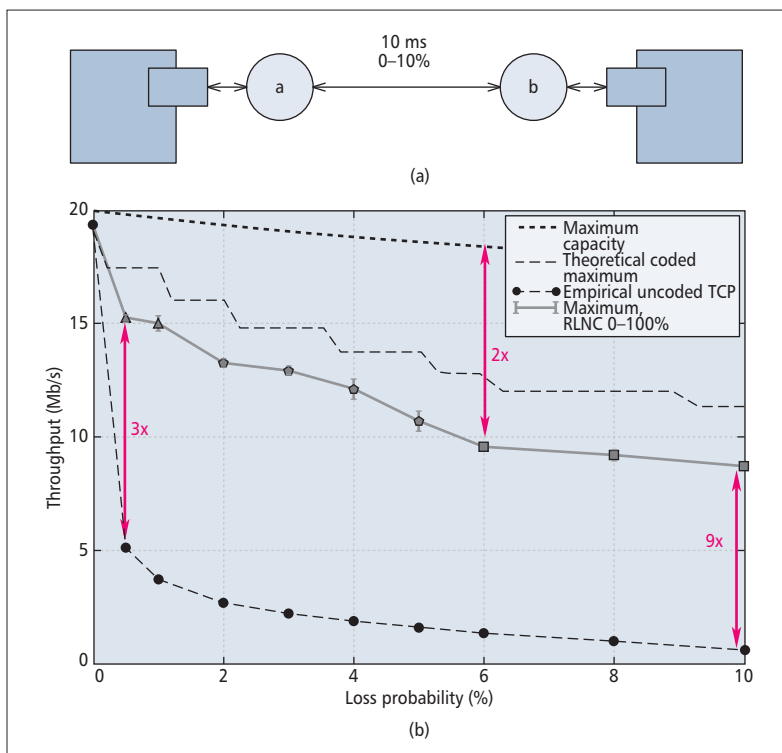
The developed software and the virtual network environment in which the software is deployed indicate that integration of NC as a functionality of a software defined network is indeed possible. However, this alone does not show that this combination of technologies is actually beneficial. In order to show that the proposed integration of network coding is plausible, a series of performance measurements has been carried out using the standard tool for network performance measurement, *iperf*. As a secondary result, this should also show that the coding approach can be applied without breaking functionality with the conventional TCP/IP network protocol stack.

We focus on links with a mean loss rate within 0 and 10 percent. Note that typical WiFi links' mean loss rates have been reported between 0.1 and 0.5 percent and as high as 4 percent depending on hardware and drivers [11] and higher in long-range links [12]. The mean packet loss rate for 3G, 3.5G, and 4G networks has been found to be in the range of 0.18–0.66 percent, 0.05–0.14 percent, and 0.03–0.30 percent, respectively [13]. The following results show that significant performance gains can be found even with loss rates below 0.5 percent using our techniques.

### SINGLE-HOP

The single-hop scenario consists of two virtual nodes, each connected to individual virtual switches. These switches are then connected with a virtual Ethernet connection on which delay and loss are introduced. This simple topology along with specified delay and packet loss is depicted in Fig. 3a. This scenario is representative of networks where only the destination and the source are capable of coding. Alternatively, the network itself could provide such functionality transparently for end devices in order to provide protection against lossy parts of the network (e.g., a satellite link).

By isolating the coding approach to a single link between coders, the consequence of the data recovery process within the developed coding approach is revealed. Despite the potential ability to recover every single erasure that may occur on the link, the performance of the transported TCP communication may not resemble error-free TCP communication. This is due to the inevitable delay of the error recovery phase, from when a lost packet should have been received to the point where it is successfully decoded. The amount of additional interference in terms of delay and jitter is reduced to a bare minimum in this single-hop scenario. The channel conditions on the investigated link are then adjusted to illustrate the tolerance of both the TCP communication and the deployed coding. Increasing packet loss on the link reveals the



**Figure 3.** Single-hop network comparing performance of TCP using network coding for erasure protection, TCP without coding, as well as the theoretical maximum of our simple coded approach and the theoretical capacity of the channel (without protocol effects). SDN is used to detect the link quality and adapt the redundancy introduced by the system.

robustness added by the coding. The bandwidth of an uncoded TCP connection is compared to TCP connections carried by the deployed coding approach using systematic on-the-fly RLNC. The achievable throughput for both uncoded and coded data flows is stated in Fig. 3b. This is compared to the channel capacity and the theoretical maximum throughput for a coded flow, found using a model for TCP throughput [14] modified to accommodate the utilized coding approach.

From the performance of the coded data flow a gain of 3x is obtained already at 0.5 percent packet loss. This performance boost increases up to 9x at 10 percent packet loss probability. Furthermore, the obtained performance of the coded flow follows a trend similar to the theoretical maximum coded throughput, showing coherence between theory and practice.

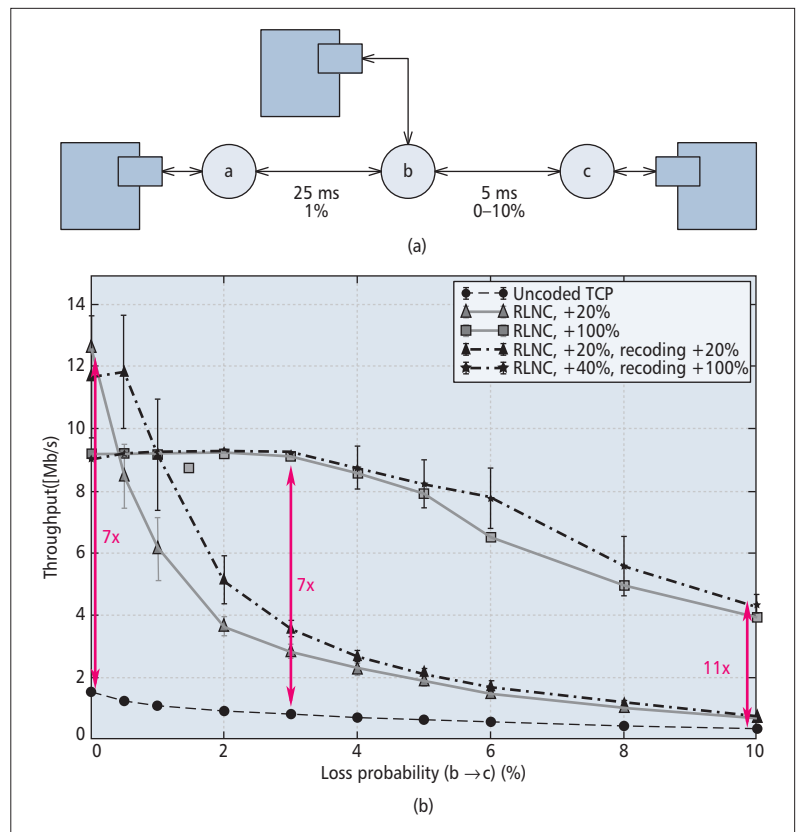
Finally, note that the latency (i.e., the number of packets in flight given the transmission rate and round trip delay) and the packet loss rate are the key factors to determine the optimal choice of redundancy to be added for coding. Thus, it is possible to develop a network coding control algorithm that can optimize for the current situation in the network.

It is important to note that the use of a more advanced and integrated coding technique (e.g., Coded TCP [14]) would require the optimal amount of redundancy for the specific network. Thus, coding would not introduce unnecessary redundancy and provide much better end-to-end service. In fact, the use of SDN would address the main limitation of Coded TCP, that is, the estimation of packet losses through the network not attributed to congestion, thus having the potential to improve its performance with respect to current practical demonstrations.

### MULTIHOP

In order to illustrate the benefits of recoding, we utilize a multihop setup where both links experience different loss probability and link delay. This network setup is presented in Fig. 4a. This particular setup is representative of Internet service provider (ISP) networks, where typically the last hop is a wireless and lossy link. Our assumption is that a switch at the last mile provides the necessary recoding for protecting against losses. Thus, studying a single recoder has a large implication in more complex topologies, where we abstract other hops within the network as they are error-free for all practical purposes.

This network scenario is also representative of mesh-like network structures, such as dedicated sensor networks, mobile ad hoc networks, and vehicular communication networks. While uncommon in consumer oriented networks, such networks are expected to gain popularity in the future with the growth of the Internet of Things (IoT). The setup illustrates the necessity of intermediate coding (recoding) when all links are prone to erasures. While prior research efforts have already drawn this conclusion, the setup tests the validity of this with a simple feedback-less coding approach using SDN to gain knowledge of the channel conditions and structure of the network. Additionally, recoding may introduce channel irregularities such as packet



**Figure 4.** Multihop network comparing a) uncoded TCP, RLNC without recoding; RLNC with recoding at the intermediate node. Recoding even after a moderate loss channel (1 percent loss rate) can provide some benefits in end throughput.

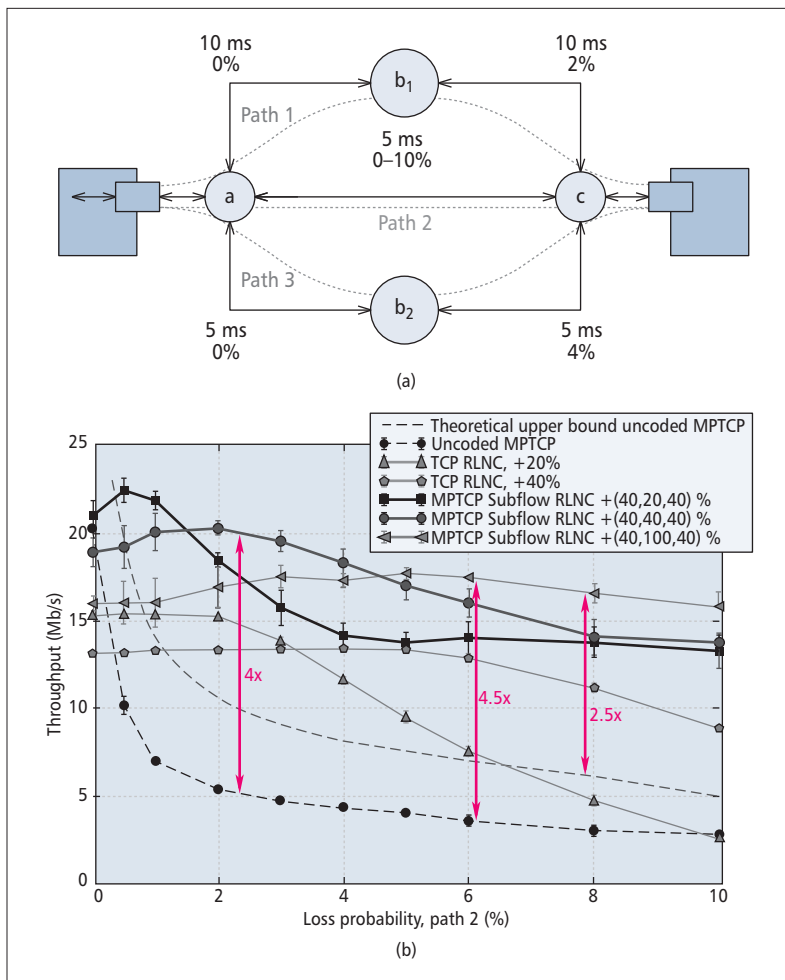
reordering and additional delay and jitter. By running TCP on top of a recoded data flow, these recoding issues are tested in practice.

Figure 4b reveals the strength of recoding. Apart from up to 11x gains over uncoded TCP, recoding also reduces congestion on the first link  $a \rightarrow b$  and introduces higher achievable throughput compared to the end-to-end RLNC coding approaches.

Similar to the single-hop case, network characteristics such as latency packet loss rate at each link determine the optimal network coding mechanism at the encoding and each of the recoding nodes. Thus, developing an NC control algorithm ran at the control plane will allow optimal operation of the system.

### MULTI-PATH

The final investigated network scenario is a multi-path setup, where multiple data paths span out between nodes. The conventional methods for communication in such scenarios is choosing the best of the available paths. This is naturally the simplest approach, and while the chosen data path provides adequate capacity, it is probably also the best approach. However, some communication scenarios may benefit from utilizing several of the available data paths. Multi-Path TCP (MPTCP) has been developed for such scenarios, but suffers from similar behavior toward packet loss and link delay as that of conventional TCP. Using a combination of network



**Figure 5.** Multi-path network comparing the use of standard MPTCP, individual TCP flows with RLNC loss protection, and MPTCP with protection of RLNC on individual paths. The plot shows that the theoretical maximum for MPTCP without coding can be outperformed more than two-fold by the use of coding, while the gain is over four-fold when comparing to a real MPTCP implementation: a) network structure; b) multi-path network performance.

coding and SDN to accommodate packet loss on each individual path may provide similar benefits for MPTCP as for conventional TCP in the single-hop and recoding scenarios.

SDN is used here not only for link quality statistics and discovery of coding nodes, but for topology discovery. In ISP networks, this can be representative for devices using separate technologies (e.g., 4G, WiFi) for accessing different networks but with a common destination (e.g., a server). In a sense, we consider that each of those networks can be represented by a simple model of delay and loss. Once the paths and their characteristics have been identified, the coding is set up to provide protection against losses in individual paths, allowing MPTCP to handle congestion control. This is not the most integrated solution, but it does allow the network to code and improve performance without changing the behavior of MPTCP.

Figure 5a illustrates the multi-path scenario investigated. This consists of one direct single-hop path (path 2) and two indirect paths with an additional hop. The two multi-hop paths,  $a \rightarrow b_1$

$\rightarrow c$  and  $a \rightarrow b_2 \rightarrow c$ , are denoted path 1 and path 3, respectively.

In a multi-path scenario, the achievable throughput of uncoded MPTCP is compared to that of a coded approach, where RLNC is utilized to protect each MPTCP subflow individually. The resulting performance is illustrated in Fig. 5b. This also includes a theoretical upper bound for the performance of MPTCP [15] along with a coded conventional TCP flow, carried over path 2. The various configurations of Coded MPTCP (top three curves) correspond to different levels of redundancy introduced in the direct link between  $a$  and  $c$ , that is, 20, 40, and 100 percent, while the redundancy of the two other paths is kept fixed at 40 percent.

Due to the high sensitivity toward packet loss and link delay, even the theoretical upper bound for MPTCP indicates poor performance in the multi-path network, and even the single-path coding approach outperforms MPTCP even though only a third of the total capacity is available to this approach. The multi-path coding provides a performance increase of up to 2.5x over the theoretical MPTCP upper bound and 4.5x over the obtained MPTCP performance.

## CONCLUSIONS

This article advocates the integration of network coding as part of software defined networking as a key to operate 5G networks more efficiently, with higher resiliency, providing higher throughput, and allowing control of data location to enable low-latency services. Furthermore, we show that the essential software packages from each concept, OpenFlow and Kodo, are already available and can be integrated to provide the required functionalities to current and future networks.

Using three basic topologies, we demonstrate not only this integration of concepts but also that simple coding strategies enable us to outperform standard TCP and multi-path TCP without modifying the underlying end-to-end transport protocols as a first step to understanding their potential. To achieve this, we exploit recoding at intermediate nodes in the network and show that gains of 3x to 11x are attainable.

## OUTLOOK

In order to validate the measurement results using real SDN-capable switches and high-end desktops, we recently built a testbed with 16 programmable high-performance network nodes and one real SDN-enabled 48-port switch. These nodes can be configured in a variety of scenarios and topologies for measurement and demonstration purposes. Each node consists of one NetFPGA with a 10 Gb/s PCI-Express programmable network interface (netfpga.org), which is highly configurable and has an active research community. The FPGA solution is needed to allow for fast switching and routing decisions, and high-end processors are needed to get the network coding speeds to satisfy the 10 Gb/s links. This equipment shall be at the core of the design and testing of 5G algorithms and network protocols.

## ACKNOWLEDGMENT

This work was partially financed by the Green Mobile Cloud project (Grant No. DFF — 0602-01372B) granted by the Danish Council for Independent Research and by the VELUX Visiting Professor Programme 2013-2014 granted by the VELUX Foundation.

## REFERENCES

- [1] N. McKeown *et al.*, "Openflow: Enabling Innovation in Campus Networks," *SIGCOMM Comp. Commun. Rev.*, vol. 38, no. 2, Mar. 2008, pp. 69–74, <http://doi.acm.org/10.1145/1355734.1355746>
- [2] M. Kim *et al.*, "Network Coded TCP (CTCP)," *Computing Research Repository*, vol. abs/1212.2291, 2012.
- [3] O. Bonaventure, M. Handley, and C. Raiciu, "An Overview of Multipath TCP," *USENIX login*; vol. 37, no. 5, Oct. 2012.
- [4] D. Kreutz, F. M. Ramos, and P. Verissimo, "Towards Secure and Dependable Software-Defined Networks," *Proc. 2nd ACM SIGCOMM Wksp. Hot Topics in Software Defined Networking*, 2013, pp. 55–60, <http://doi.acm.org/10.1145/2491185.2491199>
- [5] M. Hundebøll *et al.*, "Catwoman: Implementation and Performance Evaluation of IEEE 802.11 based Multi-Hop Networks Using Network Coding," *2012 IEEE VTC-Fall*, Sept. 2012, pp. 1–5.
- [6] M. V. Pedersen, J. Heide, and F. Fitzek, "Kodo: An Open and Research Oriented Network Coding Library," *LNCS*, vol. 6827, 2011, pp. 145–52.
- [7] D. E. Lucani *et al.*, "Fulcrum Network Codes: A Code for Fluid Allocation of Complexity," *CoRR*, vol. abs/1404.6620, 2014.
- [8] A. Moreira and D. Lucani, "On Coding for Asymmetric Wireless Interfaces," *2012 Int'l. Symp. Network Coding*, June 2012, pp. 149–54.
- [9] S. Vaughan-Nichols, "OpenFlow: The Next Generation of the Network?," *Computer*, vol. 44, no. 8, 2011, pp. 13–15.
- [10] J. Heide *et al.*, "Network Coding for Mobile Devices — Systematic Binary Random Rateless Codes," *IEEE ICC Wksp.*, 2009, June 2009, pp. 1–6.
- [11] D. C. Salyers, A. D. Striegel, and C. Poellabauer, "Wireless Reliability: Rethinking 802.11 Packet Loss," *2008 Int'l. Symp. A World of Wireless, Mobile and Multimedia Networks*, June 2008, pp. 1–4.
- [12] A. Sheth *et al.*, "Packet Loss Characterization in WiFi-based Long Distance Networks," *26th IEEE INFOCOM*, May 2007, pp. 312–20.
- [13] Y.-C. Chen *et al.*, "Characterizing 4G and 3G Networks: Supporting Mobility with Multi-Path TCP," *Dept. Comp. Sci.*, UMass Amherst, tech. rep. UM-CS-2012-022, 2012.
- [14] M. Kim, M. Médard, and J. a. Barros, "Modeling Network Coded TCP Throughput: A Simple Model and Its Validation," *Proc. 5th Int'l. ICST Conf. Performance Evaluation Methodologies and Tools*, Brussels, Belgium, 2011, pp. 131–40, <http://dl.acm.org/citation.cfm?id=2151688.2151704>
- [15] J. Cloud *et al.*, "Multi-Path TCP with Network Coding for Mobile Devices in heterogeneous networks," *2013 IEEE VTC-Fall*, Sept. 2013, pp. 1–5.

## BIOGRAPHIES

JEPPE KRIGSLUND (jepkri@es.aau.dk) is a software developer at Steinwurf ApS working on network coding protocols for wireless video multicast. He was a student in the Elite Masters Programme in Wireless Communication at Aalborg University (AAU), Denmark, where he also received his B.S. degree in electrical engineering in 2012. His research interests revolve around wireless communications and multimedia transmission with work including a mix of wireless mesh networks, network coding, video streaming, and cooperative protocol design.

JONAS HANSEN (jh@es.aau.dk) is an industrial Ph.D. student at Bang & Olufsen and AAU working on network coding code design and wireless protocols for audio signals. He was a student in the Elite Masters Programme in Wireless Communication at AAU, where he also received his B.S. degree in electrical engineering in 2012. His research interests are wireless communications and multimedia transmission with an emphasis on low-latency traffic and applications.

DANIEL E. LUCANI (del@es.aau.dk) is an associate professor in the Department of Electronic Systems, AAU. He was an assistant professor at the University of Porto from 2010 to 2012 before joining AAU. He received his B.S. and M.S. degrees in Electronics Engineering from Universidad Simón Bolívar, Venezuela, in 2005 and 2006, respectively, and his Ph.D. degree in electrical engineering from the Massachusetts Institute of Technology in 2010. His research focuses on communications, network theory, and network coding theory and applications.

FRANK H. P. FITZEK (frank.fitzek@tu-dresden.de) is the coordinator of the 5G Lab Germany and a professor at Technische Universität Dresden. He received his diploma (Dipl.-Ing.) degree in electrical engineering from RWTH-Aachen, Germany, in 1997, and his Ph.D. (Dr.-Ing.) in electrical engineering from the Technical University Berlin, Germany in 2002. He has received numerous awards, including the NOKIA Champion Award five times, the NOKIA Achievement Award (2008), the Danish SAPERE AUDE research grant (2010), and the Vodafone Innovation prize (2012). His research focuses on wireless and mobile networks, mobile phone programming, network coding, cross-layer and energy-efficient protocol design, and cooperative networking.

MURIEL MEDARD [F] (medard@mit.edu) is the Cecil H. Green Professor in the Electrical Engineering and Computer Science Department at MIT, and leads the Network Coding and Reliably Communications Group at the Research Laboratory for Electronics at MIT. She is the Editor-in-Chief of the *IEEE Journal on Selected Areas in Communications*. She was President of the IEEE Information Theory Society in 2012. She received the 2009 IEEE Communication Society and Information Theory Society Joint Paper Award, the 2009 William R. Bennett Prize in the Field of Communications Networking, and the 2002 IEEE Leon K. Kirchmayer Prize Paper Award. She is among the most highly cited researchers in her field. As a result, she was named one of the World's Most Influential Scientific Minds in 2014 by Thomson Reuters.

*The FPGA solution is needed to allow for fast switching and routing decisions, and high-end processors are needed to get the network coding speeds to satisfy the 10 Gb/s links. This equipment shall be at the core of the design and testing of 5G algorithms and network protocols.*

# Software Defined Service Migration through Legacy Service Integration into 4G Networks and Future Evolutions

Yeunwoong Kyung, Tri M. Nguyen, Kiwon Hong, Jongkwan Park, and Jinwoo Park

## ABSTRACT

With the ongoing worldwide deployment of 4G networks and consequent drop of the utilization ratio of legacy networks, network operators need a cost-effective way to flexibly operate and manage networks. The obligatory continuation of legacy services in their respective legacy networks consistently incurs OPEX. In this situation, we argue that network operators can benefit from the recent advances in SDN and NFV. This article introduces a blueprint for designing a novel network architecture that integrates legacy network services into 4G networks through SDN/NFV's software-based characteristics. In addition, our architecture reduces time to market for employing newly emerging services without mandating any changes in 4G infrastructures.

## INTRODUCTION

The explosive increase of mobile traffic along with the proliferation of smartphones, tablets, and laptops is fueling network evolution, from lower-generation network connectivity (second/third generation, 2G/3G) to higher-generation network connectivity (4G or Long-Term Evolution [LTE]/Evolved Packet Core [EPC]). With this transition, the number of 4G subscribers is continuously increasing, as 4G networks meet end-user demand with greater bandwidth. For example, Fig. 1 shows that the number of 4G subscribers already surpassed that of 2G and 3G subscribers in the months of September 2012 and July 2013, respectively, for SK Telecom subscribers in Korea [1].

Even in this situation, network operators continuously maintain and operate all the legacy networks (before 4G) because they are committed to providing services to satisfy subscribers' quality of service (QoS) requirements for legacy networks. This means that operational expenditure (OPEX) for legacy networks is constantly incurred for operation, maintenance, power supply, real estate, and so on, although their utilization decreases.

In addition, a current network has hardware/vendor-specific closed architecture with standardized interfaces. This feature leads to long time to market for deploying new services because operators have to wait for vendors to actually implement them after the standardization process, which is also time consuming. This means limited flexibility and openness.

On these grounds, software defined networking (SDN) emerged, which is a new networking paradigm where the forwarding plane is decoupled from the control plane with a certain programmable interface between them. In SDN, network intelligence is logically centralized in the control plane, and the forwarding plane is abstracted from the control plane to orchestrate service delivery. SDN introduces the benefits of programmability, which enables network configuration or new services to be launched at a much faster rate than what is currently expected using hardware-centric networks. In addition, SDN has triggered significant interest in network functions virtualization (NFV). NFV is a new concept that decouples network functions from proprietary hardware appliances. Therefore, network functions run flexibly in software to accelerate service innovation and provisioning. The European Telecommunications Standards Institute (ETSI) is currently working on the requirements and architecture as well as technical challenges to produce normative specification of NFV [2]. Both SDN and NFV aim to advance a software-based approach for more manageable, agile, and innovative networks.

With SDN/NFV's characteristics, many researchers have proposed new architectures for future network design to handle the inefficiency of the current networks, as mentioned above. They have shown that current and new network services can be provided according to operator demand using novel SDN/NFV-based platforms. However, most existing proposals have assumed a fully SDN/NFV-based new network architecture [3, 4]. This means that conventional network infrastructures should be replaced by SDN/NFV-enabled infrastructures.

Yeunwoong Kyung, Tri M. Nguyen, and Kiwon Hong are with Jinwoo Park Korea University.

Jongkwan Park is with SK Telecom.



To upgrade a network, a staged or transitional approach with the existing deployment should be considered rather than a green-field approach, as it is time consuming and expensive [5]. Therefore, we propose a new transitional approach called software defined transitional networking (SDTN). SDTN aims to provide legacy network services using the SDN/NFV paradigm on a current network.

In SDTN, a 4G network is assumed to be the underlying current network because it is a highly enhanced mobile network in terms of network functionalities such as security, QoS, mobility, and interactions with other networks. As a result, other network service functions can be migrated or mapped to the 4G network service functions. In addition, a 4G network has been invested in substantially to accommodate a larger amount of traffic demand than that of 2G/3G. This means that it is required that network operators take full advantage of the 4G network.

Therefore, in the SDTN architecture, legacy network service functions are implemented in the SDTN control plane and provided through the 4G network. Additionally, SDTN provides a suitable direction toward a fully SDN/NFV-based future network (e.g., 5G network).

The rest of this article is organized as follows. After describing early efforts at staged approaches for SDN adoption, we present the core elements of SDTN. We then discuss its operational benefits and future directions.

## RELATED WORKS

Along with the considerable attention on SDN/NFV, lots of research works have been placed on SDN/NFV-based future network design [3, 4, 6, 7]. Although the necessity for compatibility [6] and an evolutionary approach [7] was mentioned, there are no works that provide a specific method of realizing transitional architecture based on existing networks.

In addition, market leaders introduced their own SDN/NFV solutions [8, 9] to provide network services including 3G and 4G functions. Although it is conceptually clear that the network service functions can be virtualized in their solution and provided through the NFV infrastructure, a specific methodology is not explained for interacting with the forwarding plane to operate the network services and guarantee interoperability with existing network infrastructure.

Compared to these approaches, there were early proposals for SDN deployment considering existing networks. They insisted that a more realistic and cost-effective strategy than building a new network from scratch was needed.

There are three approaches to SDN adoption in an existing network infrastructure. The first is a dual stack approach [10], where all network nodes have both SDN and existing network interfaces (for normal layer 2 or 3 processing) to perform either SDN or normal processing. This approach necessitates a full deployment of network nodes capable of both SDN and normal operations. Therefore, it can become a complicated and costly solution. The second is a partial SDN approach called Panopticon [5]. Its authors

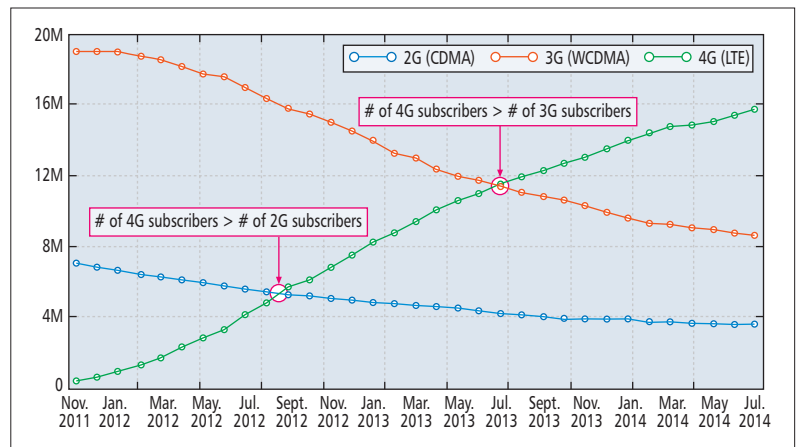


Figure 1. SK Telecom mobile subscribers by month: split per access technology.

showed that effective traffic steering could be performed with only some SDN nodes. However, network performance and operation can be affected by the number of SDN nodes and network size.

In addition, the network is actually responsible for not only packet delivery but also providing network services. However, both aforementioned approaches have not revealed a solid way to provide network services, focusing only on SDN deployment for traffic engineering.

On the other hand, in the Fabric approach [11], the authors proposed deploying SDN at the network access edge with a decoupled architecture, where the edge offers complex network services and the core only performs simple packet transport. In addition, separating the edge and core allows them to evolve independently. The authors focused on simplifying forwarding of core networks, but did not address network service provisioning and the problems of current networks mentioned in the previous section. The separation concept of the approach provides a hint as to the way to implement the SDN control plane for providing network services with an existing network infrastructure.

Our focus is on SDTN design, which applies the SDN/NFV principles to integrate the legacy network services on the current 4G network. We detail the development of our prototype to validate the key elements of the architecture.

## SOFTWARE-DEFINED TRANSITIONAL NETWORKS

To address the challenges introduced in the previous sections, we set out to define SDTN architecture. The fundamental elements of the architecture are depicted in Fig. 2. Our architecture consists of three main components: access nodes, which are already deployed to support the variety of wireless networks such as evolved NodeB (eNodeB) for 4G; NodeB with radio network controller (RNC) for 3G (defined by the Third Generation Partnership Project, 3GPP, standards) and access point for Wi-Fi (defined by the IEEE 802.11 standards); edge switches (ESs) with an edge control plane (ECP), which

Although basic network functions such as link discovery can be provided by the controller itself, there is no sharp distinction between network functions provided by VNFs and the controller itself because it depends on the controller model and can be an implementation issue. Therefore, we assume that all network functions are implemented as VNFs.

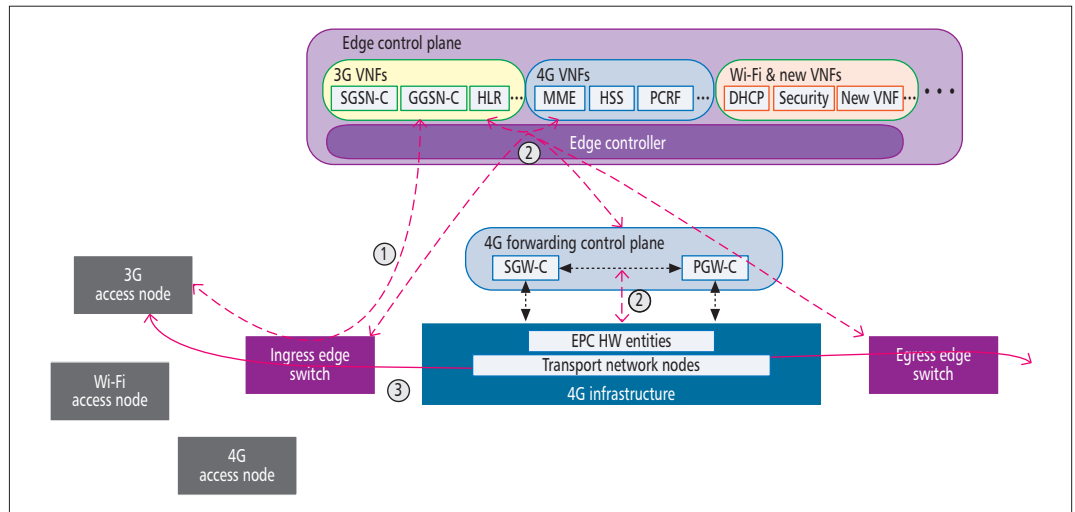


Figure 2. Software defined transitional network architecture.

serve as both ingress and egress elements for network service provisioning; and the 4G network, including the 4G forwarding control plane and infrastructure.

The key property of SDTN design is the separation of ECP and ESs from the 4G network, applying the Fabric approach [11]. Through this separation, we limit most of the intelligence to ECP, keeping the 4G network unchanged. This means that ECP is in charge of providing rich network services, while the 4G network only performs packet transmission through the 4G infrastructure. Additionally, this separation allows for the independent evolution of the ECP and 4G network, focusing on the specifics of each role.

Under this separation, incoming service requests are processed following three phases in Fig. 2. In the first phase, specific network functions are performed with the corresponding virtual network functions (VNFs) in the ECP. Then the ECP interacts with the 4G forwarding control plane and ESs to set the end-to-end data path. Finally, data is delivered. More specific service flow handling is described later in this article.

Like any system, network architecture design is guided by interfaces. As shown in Fig. 3, we use five interface types to design our architecture: application-operator, operator-operator, operator-network, host-network, and packet-switch. The application-operator interface determines how services can be provisioned by the network operator and the operator-operator interface decides how different domains inform each other of their requirements. Other interfaces are defined in [11]. In addition, the application-operator, operator-operator, and operator-network interfaces are respectively mapped to the northbound, horizontal, and southbound interfaces, as defined in the general SDN framework [12].

In this context, we assume that SDTN utilizes already deployed access nodes as they are, because splitting the control functionality from the forwarding part of the access nodes is a much more complex and challenging task than for other core nodes [3], and modifying all

deployed access nodes is not practical because of cost constraints. Extending SDTN to the wireless parts will be one of our future works.

We next investigate the details of each architectural component of SDTN.

### SDTN NETWORK FUNCTIONS

Although basic network functions such as link discovery can be provided by the controller itself, there is no sharp distinction between network functions provided by VNFs and the controller itself because it depends on the controller model and can be an implementation issue. Therefore, we assume that all network functions are implemented as VNFs.

Figure 2 shows that legacy network functions, such as serving general packet radio service (GPRS) support node (SGSN), gateway GPRS support node (GGSN), home location register (HLR), and dynamic host configuration protocol (DHCP), are implemented as VNFs through NFV at the ECP, and perform the control processes previously done at the physical boxes. During each process, a set of VNFs are configured, provisioned, and chained by utilizing physical or virtual resources (i.e., orchestrated) to operate the specific network service on demand. In addition to the legacy network services, it is possible to provide new network services from ECP by virtue of SDN/NFV without the need for specialized physical boxes that only perform specific functions.

In SDTN, pure 4G control functions such as the mobility management entity (MME), home subscriber server (HSS), and policy charging and rules function (PCRF) are also virtualized as VNFs at ECP. This means that only the control parts of serving gateway (SGW) physical box (SGW-C) and packet data network gateway (PGW) physical box (PGW-C) remain in the existing 4G control plane. There are two reasons for this approach. First, a network operator can manage all types of service information, such as user information, security, policy, and charging, simultaneously and centrally at the ECP. Furthermore, the network operator can flexibly control the MME and PCRF to communicate with

the eNodeB, SGW-C, and PGW-C to prepare user plane processes utilizing this information. Second, SDTN can use existing 4G networks with minimal modification because the SGW-C and PGW-C are usually integrated into the EPC hardware entities within the 4G infrastructure. We call this existing 4G control plane that includes only SGW-C and PGW-C the 4G forwarding control plane. This means that the 4G infrastructure in Fig. 2 includes 4G transport network nodes and EPC hardware entities without control parts where the 4G traffic is delivered based on 4G-defined interfaces [13] and transmission policy [14].

VNFs run in a cloud environment and interact with the edge controller (EC) using the application-operator interface. Although VNFs are grouped logically for each service (e.g., 3G and 4G in Fig. 2), there can be many issues with the implementation of VNFs because instantiating them in a cloud environment has a tremendous effect on service performance. For example, grouping VNFs can reduce control signaling traffic by 70 percent from that of non-grouped VNFs [15]. Therefore, the implementation strategy for VNFs should be considered carefully with sufficient analysis of the interactions among VNFs when diverse network events occur.

### EDGE CONTROLLER WITH EDGE SWITCHES

As shown in Fig. 3, the EC supports four interfaces to interact with other elements: application-operator interface, operator-network interface 1, and operator-network interface 2. The application-operator interface is used to provide network services. The operator-operator interface facilitates communication with other ECs in the operator's own network as well as different operators' networks for scalability and policy consistency. As the network grows, the centralized controller causes a scalability issue. Therefore, operators with a very large area must consider an interface that interconnects with other controllers to handle this problem [12]. In addition, when a mobile user moves to another controller's domain, a consistent service policy such as QoS should be guaranteed. To do this, policy-related information between controllers can be exchanged through the operator-operator interface. However, because information sharing between operators can be a very sensitive issue, the interface should be defined carefully with a contract between operators, and only essential information will be shared between them.

Operator-network interface 1 is used to control ESs. ESs handle both the host-network and packet-switch interfaces, and carry out interface translation between them. In SDTN architecture, the host-network interface can be packet header fields according to each host's network service protocol stack (e.g., 3G, 4G, or Internet) and the packet-switch interface can be expressed in the packet header that each 4G network node uses as an index for its forwarding table. For example, the GPRS Tunneling Protocol (GTP) header, User Datagram Protocol (UDP) header, and IP header fields can be used for the packet-switch interface because GTP tunneling is usually supported in 4G networks to transfer the data

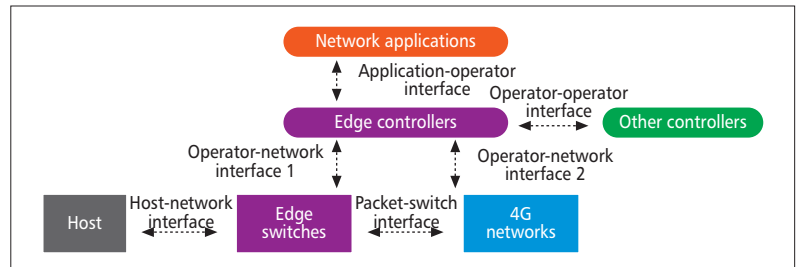


Figure 3. SDTN interfaces.

traffic through eNodeB, SGW, and PGW. This means that ESs act as tunneling points like eNodeBs in 4G networks when legacy network service traffic goes through them. ESs are controlled by EC to make each entry in the forwarding table, as MME manages eNodeBs to set up tunnels in 4G networks. In this way, other network services can be supported in our architecture.

Each entry is generated and updated using operator-network interface 1. For this interface, SDTN uses open or application programming interfaces (APIs) such as OpenFlow [10] that enable network software to be developed and upgraded much more easily than the current manual (re)configuration of hardware-centric machinery.

The complexity lies in mapping the host-network interface to the packet-switch interface at ESs controlled by EC. There are two mechanisms for this process: protocol stack translation and encapsulation. Protocol stack translation provides the mapping by swapping the header of the incoming packets with the internal 4G protocol header such as the GTP header through the ingress ESs. The reverse process is performed at the egress ESs. Using encapsulation, packets crossing the ingress ESs are encapsulated by the 4G protocol headers, and then decapsulated at the egress ESs. The latter approach is more popular and preferable due to its lower complexity. In either method, all the incoming packets must be mapped to at least one entry at ESs. The entry may also include sophisticated functions such as filtering, isolation, or policy routing, as well as 4G protocol header fields.

As SDTN aims for minimal modification of the existing 4G networks, EC supports 4G-defined interfaces [13] for operator-network interface 2 to interact with existing 4G network nodes. For instance, S1-MME, S11, and Gx are utilized to manage eNodeB, SGW-C, and PGW-C, respectively.

### EXISTING 4G NETWORKS

As explained above, the 4G forwarding control plane and infrastructure are included in 4G networks. Network nodes included in the infrastructure are controlled by the forwarding control plane to set the traffic path. After the mapping procedure at the ingress ESs, 4G network nodes can process the incoming packets with the packet-switch interface. As a result, 4G network nodes can deliver the packets while unaware of the original network service type of the packets.

In addition, all the traffic is transferred fol-

The transmission policy such as traffic classes and QoS parameters can differ among network services. Therefore, a mapping between 4G transmission policy and other services' transmission policies is required. Network operators should define this mapping carefully to ensure the QoS requirement for each service.

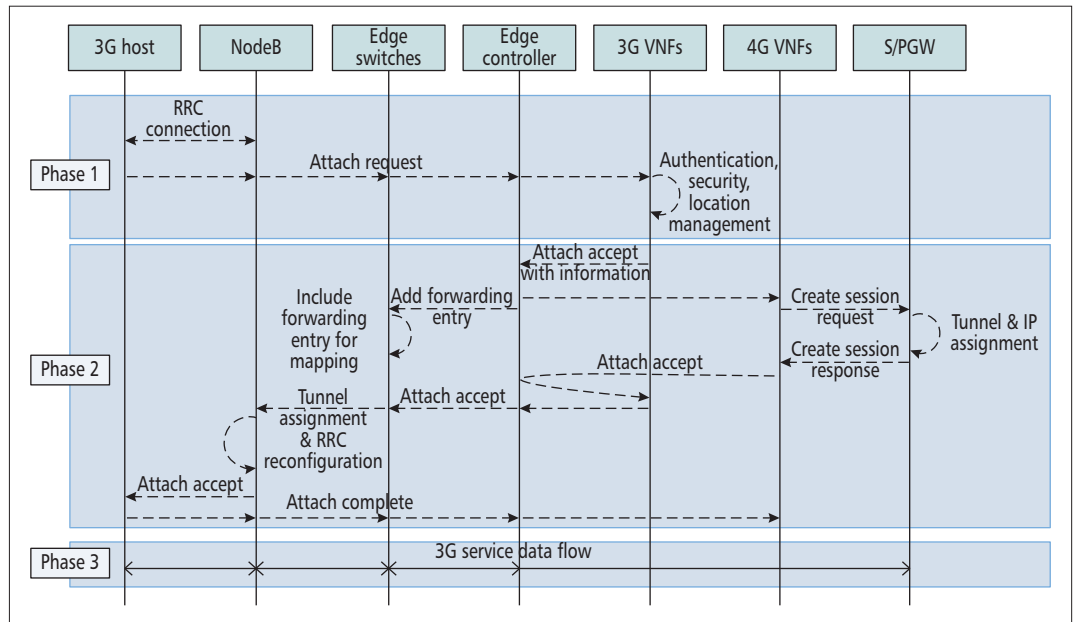


Figure 4. Information flow of 3G service in SDTN.

lowing a 4G-defined transmission policy such as the EPS bearer-based QoS policy. For example, incoming traffic into the ingress ES is assigned to a specific EPS bearer with bearer ID based on QoS requirements by using various QoS parameters defined in 3GPP standards [14]. Then bearer-based QoS control can be supported through the established bearer in 4G networks between ingress and egress ESs. In this way, traffic is transmitted applying the 4G-defined policy as if it is 4G service traffic. When the traffic goes out from 4G networks, QoS policy is converted to that of the original service at the egress ES. In this way, end-to-end QoS control can be provided. In fact, the transmission policy such as traffic classes and QoS parameters can differ among network services. Therefore, a mapping between 4G transmission policy and other services' transmission policies is required. Network operators should define this mapping carefully to ensure the QoS requirement for each service.

### USE CASE: FLOW HANDLING OF HETEROGENEOUS SERVICES

Support of heterogeneous services, based on 4G, other legacy, or new protocols, can be achieved by interacting between the EC and the corresponding VNFs.

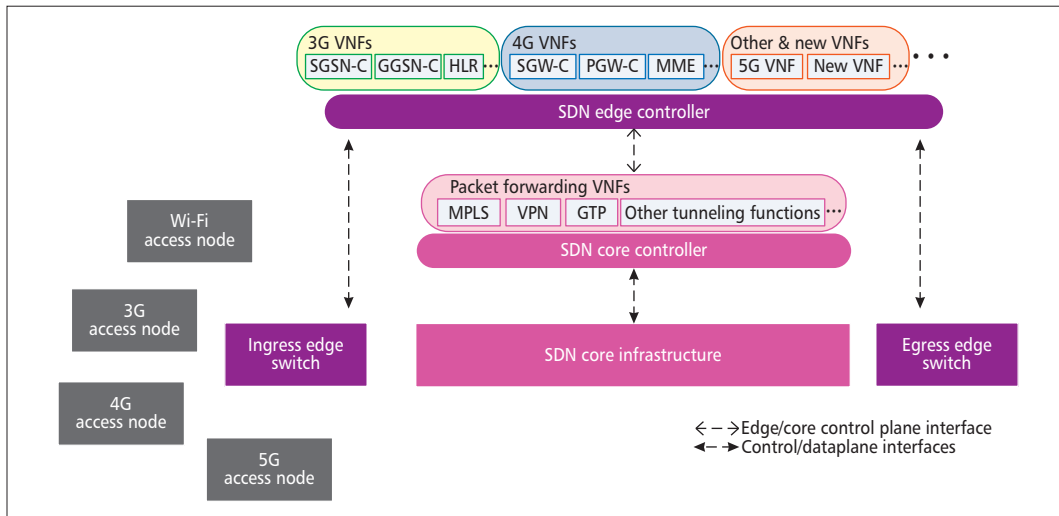
For instance, when 3G service flow is processed, 3G VNFs such as SGSN-C and GGSN-C act as the original 3G functions, that is, SGSN and GGSN. In this case, the control plane process through 3G VNFs is always performed in advance of the user plane process in the 4G infrastructure.

As introduced in Fig. 2, there are three phases for handling the 3G service flow, two control plane phases and a user plane phase. Figure 4 describes the abstracted information flow of these three phases for the initial attachment of a 3G host. In this figure, we assume that the NodeB and edge switches blocks

include NodeB and RNC, and ingress/egress ESs, respectively.

First, after the radio resource control (RRC) connection is performed, the Attach Request message of the 3G host is delivered to the EC via the ingress ES because the ES does not have an entry matched to process the message. EC then carries it to the associated 3G VNFs. Subsequent to the host authentication and security setup between the host and 3G VNFs (we omit the detailed signaling flow of this step), the host location is updated, and the virtual session is established. In the 3G service control plane process, both SGSN-C and GGSN-C perform session management, and both SGSN-C and HLR exercise mobility management.

In the second phase, all the information related to session and mobility management such as the host profile, policy, and location is stored centrally at the ECP and utilized to prepare the user plane process. To enable this, the EC handles S/PGW through MME and PCRF control using this information. In other words, via conventional 4G interfaces such as Gx and S11, the 4G VNFs control S/PGW using Session Request/Response messages. Thanks to the interaction between the 4G forwarding control plane and ECP, the packet forwarding behavior of the 3G service flow such as GTP tunneling for EPS bearer service can be defined at the S/PGW. To complete the traffic path from the ingress to egress ESs, the EC controls ESs to include the forwarding entry for the mapping between 3G and 4G service flows. After the session establishment in the 4G networks is finished, the 3G VNFs send an Attach Accept message to the host via NodeB. This way, the radio access bearer of 3G service is established between the host and ingress ES via NodeB, and the EPS bearer of 4G service is established between the ingress and egress ESs for end-to-end QoS provision in user plane. Finally, 3G service traffic is transferred through the established bearer in the last phase.



**Figure 5.** SDTN's future evolution.

In the case of Wi-Fi service, it follows the same procedures as 3G service in SDTN. It is noted that Wi-Fi traffic is transmitted through ESs via either 4G infrastructure or the Internet based on the network operator's policy.

## OPERATOR BENEFITS AND FUTURE DIRECTION OF SDTN

Along with the previously stated aspects, SDTN can provide new benefits that are not present in current networks. First, SDTN allows operators to reduce the OPEX of the core network for legacy network services because their core network equipment, such as physical boxes of GGSN and SGSN, the utilization of which gradually lessens with time, can be effectively discarded. OPEX reduction also occurs thanks to the lower maintenance, upgrade, real estate, and operational cost of a software-based approach. We expect that the OPEX could be further reduced when the majority of legacy services are virtualized on the SDTN architecture.

Second, by separating the edge and 4G networks, SDTN allows operators to use already deployed 4G networks to provide legacy network services. This means that additional infrastructure capital expenditure (CAPEX) for the transport network is not required. In addition, SDTN can provide not only legacy services, but also new network services on the 4G infrastructure thanks to the benefits of SDN/NFV.

Although additional CAPEX and OPEX of cloud infrastructure to provide VNFs are required, cost-efficient deployment and operation will be possible because cloud infrastructure can be less dependent on hardware appliances as well as physical locations, and optimally utilized in terms of resources and power even when the new services are introduced based on SDN/NFV's flexibility.

Third, SDTN can be a transitional approach to the future network. As introduced earlier, many researchers and vendors believe that a fully SDN/NFV-based network is the appropriate direction of the future network. If a transi-

tional approach is not considered, operators will have to discard all current network equipment and deploy a completely new SDN-enabled infrastructure. This is a time consuming process with large CAPEX. Instead, SDTN enables operators to prepare a fully SDN/NFV-based network with underlying current 4G networks as a staged approach. Additionally, SDTN architecture will be migrated to a fully SDN/NFV-based network.

We expect that the fully SDN/NFV-based network will also have decoupled architecture between the edge and core parts, as shown in Fig. 5. Similar to SDTN, the edge part is responsible for complex service provisioning, whereas the core part performs packet forwarding. The difference with SDTN is the evolution of the underlying infrastructure and a core controller. Infrastructure consists of SDN-enabled network equipment and is controlled by an SDN core controller where VNFs related to packet forwarding such as tunneling functions are implemented. Therefore, flexible configuration of infrastructure for end-to-end traffic management is possible in a fully SDN/NFV-based network. In addition, new interfaces between the infrastructure and core controller as well as between the edge and core controller should be defined. While traditional standardized interfaces are supported in SDTN, open interfaces or associated APIs will be utilized to foster network service innovation in the future network.

## CONCLUSION

This article introduces SDTN, an architecture and methodology for aiding network operators applying SDN/NFV in existing 4G networks. SDTN reaps the benefits of a software-based approach that not only integrates all the virtualized legacy network functions, but also increases the capability to roll out new network features. In addition, SDTN can evolve its network easily by updating or replacing network services on the edge control plane without mandating any changes in the existing 4G infrastructure. Therefore, the OPEX for the legacy networks and CAPEX for additional transport network infrastructure are not needed.

*New interfaces between the infrastructure and core controller as well as between the edge and core controller should be defined. While traditional standardized interfaces are supported in SDTN, open interfaces or associated APIs will be utilized to foster network service innovation in the future network.*

*SDTN takes significant steps toward fully SDN/NFV-based future networks. We expect that SDTN will be migrated easily into future networks where all complex network services are implemented from the edge and forwarding-related functions are provided in the core with SDN-enabled infrastructure.*

Furthermore, SDTN takes significant steps toward fully SDN/NFV-based future networks. We expect that SDTN will be migrated easily into future networks where all complex network services are implemented from the edge and forwarding-related functions are provided in the core with SDN-enabled infrastructure. In future work, SDTN will be validated to measure the practical benefits based on a virtualized network environment using OpenFlow-based components and real network events.

#### ACKNOWLEDGMENTS

This work was supported by ICT R&D program of MSIP/IITP. [B0101-15-1272, Development of Device Collaborative Giga-Level Smart Cloudlet Technology].

#### REFERENCES

- [1] Ministry of Science, ICT and Future Planning, "Statistics on Telecom Services (2014)," July 2014.
- [2] ETSI, "Network Functions Virtualization (NFV); Infrastructure Overview," ETSI Group Spec. NFV-INF 001 v. 1.1.1, Jan. 2015.
- [3] K. Pentikousis, Y. Wang, and W. Hu, "MobileFlow: Toward Software-Defined Mobile Networks," *IEEE Commun. Mag.*, vol. 51, no. 7, July 2013, pp. 44–53.
- [4] P. K. Agyapong et al., "Design Considerations for a 5G Network Architecture," *IEEE Commun. Mag.*, vol. 52, no. 11, Nov. 2014, pp. 65–75.
- [5] D. Levin et al., "Incremental SDN Deployment in Enterprise Networks," *Proc. ACM SIGCOMM*, Aug. 2013, pp. 473–74.
- [6] R. Trivisonno et al., "SDN-Based 5G Mobile Networks: Architecture, Functions, Procedures and Backward Compatibility," *Trans. Emerging Telecommun. Technologies*, vol. 26, no. 1, Jan. 2014, pp. 82–92.
- [7] C. J. Bernardos et al., "An Architecture for Software Defined Wireless Networking," *IEEE Wireless Commun.*, vol. 21, no. 3, June 2014, pp. 52–61.
- [8] Cisco, "Cisco Virtualized Packet Core," C45-730492-01, Jan. 2015.
- [9] Alcatel Lucent, "Alcatel Lucent Virtualized EPC: Delivering on the Promise of NFV and SDN," Alcatel-Lucent Application Note (NP2014014132EN), Feb. 2014.
- [10] N. McKeown et al., "OpenFlow: Enabling Innovation in Campus Networks," *Proc. ACM SIGCOMM Comp. Commun. Rev.*, Apr. 2008, pp. 69–74.
- [11] M. Casado et al., "Fabric: A Retrospective on Evolving SDN," *Proc. ACM HotSDN*, Aug. 2012, pp. 85–90.
- [12] M. Jarschel et al., "Interfaces, Attributes, and Use Cases: A Compass for SDN," *IEEE Commun. Mag.*, vol. 52, no. 6, June 2014, pp. 210–17.

- [13] 3GPP, "Digital Cellular Telecommunications System; Universal Mobile Telecommunications System (UMTS); LTE; Network Architecture," TS 23.002 v.12.6.0 Release 12 (2015-01), Jan. 2015.
- [14] 3GPP, "LTE; Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description; Stage 2," TS 136 300 v.12.3.0 (2014-09), Sept. 2014.
- [15] H. Hawilo et al., "NFV: State of the Art, Challenges and Implementation in Next Generation Mobile Networks (vEPC)," *IEEE Network*, vol. 28, no. 6, Nov. 2014, pp. 18–26.

#### BIOGRAPHIES

YEUNWOONG KYUNG (ywkyung@korea.ac.kr) received his B.S. degree in electrical engineering from Korea University, Seoul, in 2011. He is currently in an M.S. and Ph.D. integrated course in the School of Electrical Engineering at Korea University. His research interests include flow-based mobility management in wireless networks and software-defined networking (SDN), especially SDN scalability and network service integration.

TRI M. NGUYEN (nmtrivn@gmail.com) received his B.S. and M.S. degrees in computer science from Vietnam National University, Hanoi, in 2000 and 2006, respectively. He is currently working toward a Ph.D. degree in the School of Electrical Engineering at Korea University. His research interests include mobility management and resource management in wireless and broadband access networks.

KIWON HONG (shalaman@korea.ac.kr) received his B.S. degree in electrical engineering from Korea University in 2011. He is currently in an M.S. and Ph.D. integrated course in the School of Electrical Engineering at Korea University. His research interests include the QoS management of optical IP networks considering service level agreements.

JONGKWAN PARK (jongkwan.park@sk.com) is currently a leader of the Core Network Laboratory in the Network Technology R&D Center at SK Telecom. His research interests include software-defined networking and network functions virtualization for smart network deployment. In addition, he is focusing on service-centric networks to optimize network operation through flexible network traffic control based on the user profile, service types, and network context.

JINWOO PARK (jwpark@korea.ac.kr) received his B.S. in electronics engineering from Korea University in 1979 and his Ph.D. degree in electrical engineering from Virginia Tech in 1987. He is currently a professor with the School of Electrical Engineering at Korea University. His research interests are mobile service management in integrated wireless and wired networks, context-aware networking, content delivery networks, and software-defined networking.

**CALL FOR PAPERS**  
**IEEE COMMUNICATIONS MAGAZINE**  
**WIRELESS TECHNOLOGIES FOR DEVELOPMENT (W4D)**

**BACKGROUND**

We live in a world in which there is a great disparity between the lives of the rich and the poor. Using information and communication technologies for the purpose of development (ICT4D) offers great promise in bridging this gap through its focus on connecting human capacity with computing and informational content. It is well known that Internet access has the capability of fostering development and growth by enabling access to information, education, and opportunities. Wireless technology is a promising solution to this problem of digital exclusion and can be instrumental in democratizing access to the Internet by unfettering developing communities from the encumbering constraints of infrastructure (traditionally associated with broadband Internet provisioning). The focus of the proposed feature topic is on leveraging wireless technologies for development (W4D) to increase the quality of life for a larger segment of human societies by providing them opportunities to connect resources and capacity, especially by provisioning affordable universal Internet access. To reflect recent research advances in using W4D, this feature topic calls for original manuscripts with contributions in, but not limited to, the following topics:

- “Global access to the Internet for all” (GAIA) using wireless technologies
- Do-it-yourself (DIY) wireless networking (such as community wireless networks) for the developing world
- Cost-efficient wireless networked systems appropriate for use in underdeveloped areas
- Fault-tolerant resilient wireless networking technologies for the developing world
- Rural/remote area wireless solutions (that can work efficiently with resource constraints such as intermittent and unreliable access to power/ networking service)
- Simplified network management techniques (including support for heterogeneous service delivery through multiple solutions)
- Using cognitive radio technology and 5G standards (with possible native integration of satellites) for GAIA
- Techno-economic issues related to W4D (including development of flexible pricing and incentive structures as well as new spectrum access models for wireless)
- Techno-political and cultural issues related to using wireless communications for development
- Using emerging networking architectures and future Internet architectures [e.g., cloud computing, fog computing, network functions virtualization (NFV), information centric networking (ICN), software defined networking (SDN), and delay tolerant networking (DTN)] with wireless technologies for development.
- Using wireless access/ distribution technologies (such as the following) for development: TV white spaces (TVWS); satellite communications using advances in geostationary orbit (GEO) and low-earth orbit (LEO) satellites; low-cost community networks; cellular technologies (such as CDMA 450, the open-source OpenBTS, etc.); wireless mesh and sensor networks; Wi-Fi-Based Long-distance (WiLD) networks; and wireless based wireless regional access networks (WRANs).

Since our aim with this feature topic (FT) is to provide a balanced overview of the current state of the art of using wireless technologies for development, we solicit papers from both industry professionals and researchers, and we are interested in both reports of experience and in new technical insights/ideas.

**SUBMISSIONS**

Articles should be tutorial in nature and written in a style comprehensible to readers outside the specialty of the article. Authors must follow IEEE Communications Magazine's guidelines for preparation of the manuscript. Complete guidelines for prospective authors are found at: <http://www.comsoc.org/commag/paper-submission-guidelines>.

It is important to note that IEEE Communications Magazine strongly limits mathematical content, and the number of figures and tables. Paper length (introduction through conclusions) should not exceed 4,500 words. All articles to be considered for publication must be submitted through the IEEE Manuscript Central (<http://mc.manuscriptcentral.com/commag-ieee>) by the deadline.

**SCHEDULE FOR SUBMISSIONS**

- Submission Deadline: December 1, 2015
- Notification Due Date: March 1, 2016
- Final Version Due Date: May 1, 2016
- Feature Topic Publication Date: July 1, 2016

**GUEST EDITORS**

Junaid Qadir  
School of EE and CS (SEECs),  
National University of Sciences and  
Technology (NUST), Pakistan  
[junaid.qadir@seecs.edu.pk](mailto:junaid.qadir@seecs.edu.pk)

Marco Zennaro  
The Abdus Salam International Centre for  
Theoretical Physics (ICTP), Italy  
[mzennaro@ictp.it](mailto:mzennaro@ictp.it)

Saleem Bhatti  
University of St Andrews  
St Andrews, UK  
[saleem@st-andrews.ac.uk](mailto:saleem@st-andrews.ac.uk)

Arjuna Sathiaseelan  
Computer Laboratory,  
University of Cambridge,  
United Kingdom  
[arjuna.sathiaseelan@cl.cam.ac.uk](mailto:arjuna.sathiaseelan@cl.cam.ac.uk)

Adam Wolisz  
Technische Universität Berlin and  
University of California, Berkeley, USA  
[awo@ieee.org](mailto:awo@ieee.org)

Kannan Govindan  
Samsung Research, India  
[gkannan16@ieee.org](mailto:gkannan16@ieee.org)

## NETWORK TESTING



Ying-Dar Lin



Erica Johnson

The geographical scale of network testing depends on the size of the testbed. It could be as small as the desk area in a laboratory, medium on a campus or within a building, or as large as the wide area in the Internet. Traditionally, one would construct the testbed in a laboratory so that all the environmental parameters could be programmed and controlled, and the results could be reproduced. However, since some of the environmental parameters are too complicated to program and control, often we need to use the real environment as the testbed. This extends the testbed from a laboratory to a building, a campus, or even the Internet. As the parameters in the real environment are often not programmable and controllable, we thus lose the exact reproducibility of results. One common way to solve this problem is to repeat the experiments for a prolonged period and average over the huge number of repetitive runs. Then the testing period is usually weeks to months instead of hours. Another less common way is to “capture” the environmental parameters as traces and replay them onto a laboratory testbed. Precise reproducibility can be achieved, but the environment parameters in the captured traces can only be “replayed” instead of programmed and controlled. The results from such trace-driven testing are useful if the environments and captured traces are representative enough.

In this issue, we accepted three articles from 14 submissions. The major reasons to reject submissions are (1) not enough focus or percentage on network testing, (2) not significant enough in technical contributions, and (3) presentation quality (e.g., not enough comparison with the state of the art, not insightful enough in result interpretation). Starting from this issue, we have successfully shortened the time to publication by two months. The due dates are February 1 and August 1, with their corresponding publication dates in September and March, respectively.

Among the three accepted articles, there is one each on testing at the laboratory scale, building scale, and Internet scale. At the laboratory scale, the authors of the first one built a 10 Gb/s wire-speed packet sniffer from off-the-shelf

hardware and open source software. The authors of the second article designed a building-scale platform to benchmark various indoor localization solutions. At the Internet scale, the authors of the third article measured to find which Facebook servers or Akamai servers are accessed by global geographically distributed users.

In the article on the 10 Gb/s packet sniffer (“Testing the Capacity of Off-the-Shelf Systems to Store 10GbE Traffic”), off-the-shelf hardware and open source software are used to build a packet sniffer that needs to capture and store packets at 10 Gb/s. The authors identify the bottleneck as the storage subsystem, and thus develop low-level optimization techniques at the network interface card (NIC) driver, hard drives, and the interaction between them. The deployed techniques include system call minimization (processing packets in batches instead of individually), huge intermediate receive buffer (avoiding write throughput drops), memory alignment (direct memory access, DMA, with page alignment for packets transferred in blocks), sniffing and storage overlapping (isolated and parallelized execution), timestamping (accurate timestamping by low-level techniques), tuned redundant array of independent disks (RAID) volumes (parameter configurations for number of disks, strip size, write cache policy, file system), and sniffing optimization (pre-allocated and reuse of memory, memory mapping, parallel direct paths, batch processing, prefetching, affinity). Together they reduce packet capture loss significantly. As the industry moves forward to 40 Gb/s and 100 Gb/s interfaces, these techniques and yet-to-be-developed new techniques would be needed for wire-speed packet capture and sniffing.

The authors of the article on benchmarking indoor localization solutions (“Platform for Benchmarking of RF-based Indoor Localization Solutions”) tried to standardize the metrics and methodologies to evaluate the prevailing indoor localization solutions. Their approach is aligned with the upcoming ISO/IEC 18305 standard “Test and Evaluation of Localization and Tracking Systems.” The



metrics include not just accuracy, but also latency, energy efficiency, setup time, and required infrastructure. The developed test platform enables online live/real-time testing on wireless testbed facilities within buildings, and also offline replayed/non-real-time testing with the raw RF dataset including received signal strength indicator (RSSI), link quality indicator (LQI), and time of arrival (ToA). The raw RF dataset is generated from experiments defined by experiment specification, evaluation points, and interference specification, and are then fed into the system under test (SUT, i.e., an indoor localization solution), in either real-time mode or replayed mode, to calculate the performance metrics. The results illustrate the performance trade-off of three localization solutions in accuracy, latency, and energy efficiency. Another contribution is that the publically accessible RF datasets from multiple environments (office with brick walls, office with plywood walls, industrial open space) with evaluation points near a wall or in the middle of a room. It is possible for researchers to evaluate their localization solutions with these RF datasets without actually building the indoor testbeds.

The final article on Facebook server distribution (“How Far Is Facebook from Me? Facebook Network Infrastructure Analysis”) assesses the effectiveness of content distribution networks (CDNs) for services such as Facebook. The contents on Facebook are accessible to subscribers through either native Facebook servers or Akamai servers distributed world-wide, where Akamai is a CDN provider that helps content providers to push their contents near their subscribers. The authors conducted a comprehensive measurement-based analysis of the Facebook infrastructure, by:

1. Identifying URLs associated with Facebook
2. Measuring access delay and network path to the URLs
3. Establishing a distributed infrastructure, via Planet-Lab, to obtain results from different subscriber locations
4. Geo-locating subscriber nodes, Facebook servers, and Akamai servers

They found that users experience much lower delay with Akamai servers than with native Facebook servers. The delay reduction is over 10 times in countries with Akamai server coverage. Akamai covers 35 countries and reaches 41 countries (i.e., subscribers in 6 countries can access Akamai servers in nearby countries). The location of a subscriber matters in the experienced delay, because the location decides whether he or she would access a native Facebook server or an Akamai server and how far the accessed server is. This interesting study could be extended to services other than Facebook.

## BIOGRAPHIES

YING-DAR LIN [F] (ydlin@cs.nctu.edu.tw) is a Distinguished Professor at National Chiao Tung University, Taiwan. He received his Ph.D. in computer science from UCLA in 1993. He is the director of the Network Benchmarking Lab which reviews network products with real traffic and is an approved test lab of the Open Networking Foundation (ONF). He is an IEEE Distinguished Lecturer and ONF Research Associate. He co-authored *Computer Networks: An Open Source Approach* (McGraw-Hill, 2011).

ERICA JOHNSON (erica.johnson@iol.unh.edu) combines business acumen and an in-depth understanding of complex networking technology to direct the University of New Hampshire InterOperability Laboratory (UNH-IOL). In recognition of her ability to drive technical innovation, *Fierce Telecom* named her to the publication's 2011 Women in Wireline. She serves as an IPv6 Ready Logo Regional Officer, IPv6 Forum Fellow, and USGv6 Test Program lead. She received her Bachelor of Computer Science and M.B.A. from UNH in 2001 and 2011, respectively.

# Testing the Capacity of Off-the-Shelf Systems to Store 10GbE Traffic

Victor Moreno, Javier Ramos, José Luis García-Dorado, Ivan Gonzalez, Francisco J. Gomez-Arribas, and Javier Aracil

## ABSTRACT

The maturity of the telecommunications market and the fact that user demands increase every day leaves network operators no option but to deploy high-speed infrastructures and test them in an efficient and economical manner. A common approach to this problem has been the storage of network traffic samples for analysis and replay using different versions of what we have named NTSS. This type of task is particularly demanding in 10 Gb Ethernet links and has traditionally been addressed by closed solutions or NTSS built on top of high-end hardware. However, these approaches lack flexibility and extensibility, which typically translates into higher cost. This work studies how NTSS can be built using COTS: a combination of commodity hardware and open source software. To this end, we present the current limitations of COTS systems and focus on low-level optimization techniques at several levels: the NIC driver, hard drives, and the software interaction between them. The application of these techniques has proven crucial for reaching 10 Gb/s rates, as different state-of-the-art systems have shown after an extensive performance test.

## INTRODUCTION

Operators are currently deploying novel network architectures and equipment with bandwidth capabilities of multi-gigabit rates and beyond. Testing the performance and correct operation of such deployments is a challenging task that operators must face. This also applies to other players in the Internet arena, such as companies, third-party enterprises, and banks, which deploy new services for their customers and employees.

The most simple and efficient way to test such infrastructures and services is sniffing and storing all traversing test traffic for its subsequent analysis [1]. Such an analysis may focus not only on searching malformed or unexpected packets, for example, erroneous virtual LAN (VLAN) or multiprotocol label switching (MPLS) headers, or duplicated frames, but also on the network performance and quality of service (QoS) parameters' values — band-

width, packet loss, delay, or jitter. Additionally, stored traffic may be used not only passively but also actively when replaying the content of the stored traces for testing purposes [2]. We suggest referring to those systems that sniff and store traffic as network traffic storage systems (NTSS).

Even an intuitively simple task such as sniffing and storing traversing traffic is a challenge when dealing with 10 Gb/s rates or higher due to the great amount of resources and computational power needed. Traditionally, specialized hardware devices such as field programmable gate arrays (FPGAs), network processors, and high-end closed commercial solutions have been applied to tackle the traffic sniffing and storage problem. Such solutions address the performance part of the problem in a very effective way, and they also offer high degrees of both determinism and robustness, desirable for any industrial development. However, such positive features are obtained at the expense of flexibility and extensibility, turning the deployment, evolution, and maintenance processes into difficult tasks. Furthermore, the prices of those systems is elevated, ranging from \$85,000 to \$250,000 depending on their storage capacity.<sup>1, 2</sup>

As an alternative, the research community has recently focused on commercial off-the-shelf (COTS) solutions to accomplish high-performance tasks [3]. COTS systems have emerged as the combination of commodity hardware and open source software. Such systems provide flexibility, availability, and scalability while handling multi-gigabit rates and cutting expenditures in terms of both deployment and maintenance [4]. For example, the system used for the experiments presented in this article had a price of \$10,000, which is one order of magnitude lower than the price of its closed solution counterpart. With this in mind, this article explains the key aspects for COTS systems to sniff and store packets at multi-gigabit rates. Specifically, such keys comprise fine low-level tuning at the network interface card (NIC) driver, hard drives, and application levels. Subsequently, we provide an extensive performance evaluation of state-of-the-art NTSS systems that have successfully reached such a goal.

The authors are with Universidad Autónoma de Madrid.

<sup>1</sup> <http://www.netapp.com/products/storage-systems/ef-series/>

<sup>2</sup> <http://www.napatech.com/products/>

## PROBLEM STATEMENT AND CONTRIBUTIONS

Traffic storage has become a challenging task, as a fully saturated 10 Gigabit Ethernet (GbE) link in the worst case scenario (minimal size packets, i.e., 64 bytes on Ethernet with cyclic redundancy check, CRC, included) carries more than 14 million packets/s. In this demanding scenario, we first note how much traffic may be sniffed with standard software running on a commodity server. Specifically, our commodity server is a Supermicro X9DR3-F with two 6-core Xeon E5-2630 processors running at 2.30 GHz and hyper-threading disabled, with 96 GB of DDR3 RAM at 1333 MHz. The server is equipped with an Intel 82599 10 GbE NIC plugged into a Peripheral Component Interconnect Express (PCIe) 3.0 slot. The software is composed of an Ubuntu Server 14.04 configured with a 3.14 kernel and the default network stack, the vanilla Intel `ixgbe` NIC driver, and the de facto standard traffic sniffer `tcpdump`. On the sender side, we have an FPGA-based traffic transmission system capable of replaying at link-rate both fixed-size synthetic traffic and packet traces previously stored in the machine [5].

The results are shown in the left two-column group in Fig. 1, which shows the percentage of stored packets with this configuration (named `tcpdump vanilla`), for two traffic injection cases: namely, synthetic 64-byte packets (CRC included) and a real backbone trace from CAIDA,<sup>3</sup> both replayed at wire speed. We observe that this out-of-the-box scenario can only sniff and store less than 10 and 38 percent of the total sent packets for synthetic and real traffic, respectively. The out-of-the-box configuration has thus proven insufficient to capture full-rate 10 GbE traffic.

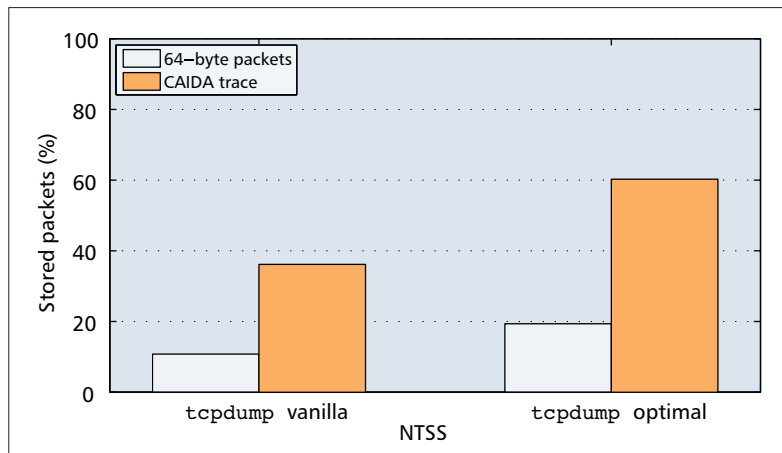
Nevertheless, the research community has first tried to improve the NIC vanilla driver to boost up its performance, as reviewed in the following section, but there is scarce knowledge on how to write such traffic in hard drives at multi-gigabit rates.

Consequently, this article studies how to tune such drives to increase their performance, a question that is dealt with in the following section. Then we explain how to optimally combine sniffing and storing, and provide a performance evaluation of the state-of-the-art NTSS in 10 GbE networks. Finally, the guidelines and take-away messages presented through this article are summarized together with proposed future work in the field.

### SNIFFING TRAFFIC FROM THE WIRE

The first task is to sniff the traffic from the wire. Alternatively, this task was traditionally termed *capture*; however, we avoid the term here as it is easily confused with the term *store*. The following performance optimizations techniques have been applied to the sniffing process in the literature.

**Pre-allocation and reuse of memory:** On vanilla drivers, for each received packet, a set of structures and buffers is allocated. Those



**Figure 1.** Percentage of stored packets (into a RAID-0 volume with nine disks) for a fully saturated 10 Gb Ethernet link for NTSS based on `tcpdump` with vanilla and optimized configurations for a 30-minute experiment.

resources are released once the packet is delivered to upper layers. It has proven more efficient to pre-allocate a pool of structures and reuse them for subsequent incoming packets.

**Memory mapping:** The use of these techniques (e.g., direct NIC access, DNA [6]) allows high-level applications to map the receiving buffers located at the driver level, thus reducing the number of copies.

**Use of parallel direct paths:** Modifying the network driver to bypass the operating system's network stack makes it possible to create parallel paths from the NIC (receive side scaling, RSS) to user-level applications. As collateral effects, more CPU cores are necessary, and packet reordering may occur [7].

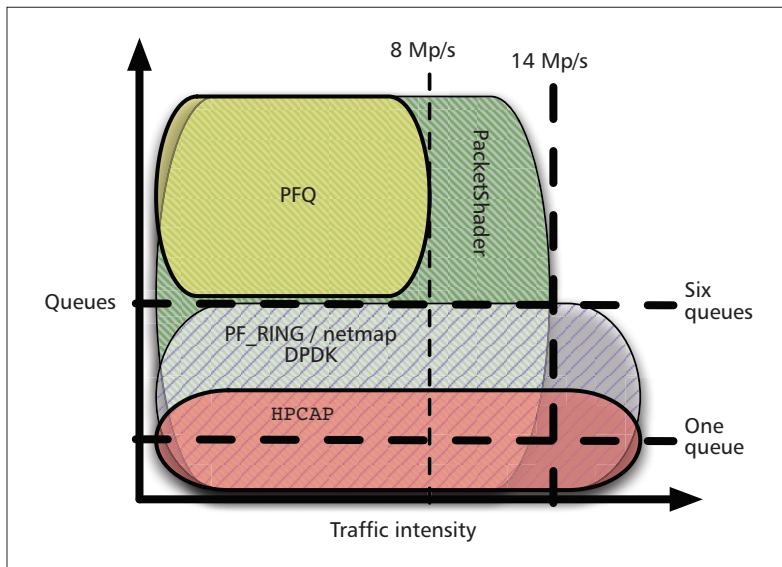
**Batch processing:** Typically, upper-layer applications make a system call to receive a packet. Processing a batch of packets per system call can reduce the resulting overhead. Nevertheless, such techniques may also entail a latency increment and inaccurate timestamping [8].

**Prefetching:** This technique consists of preloading memory locations in processors' caches in a predictive way so that they can be quickly accessed in the near future should they be needed, thus reducing cache misses.

**Affinity:** Non-uniform memory access (NUMA) architectures group processing cores along with independent memory banks to create a NUMA node. Performance is significantly increased if the sniffing process is placed on the same NUMA node as the driver-level receiving threads, thus reducing the data access overhead [4].

We note that the techniques to improve sniffing performance require changes in either the NIC driver or the operating system's network stack, or both. Nevertheless, some straightforward affinity-related adjustments may be performed in an out-of-the-box Linux machine with `tcpdump`, with significant performance improvements as shown in Fig. 1, labeled `tcpdump optimal`. The main difference between the vanilla and optimized execution of `tcpdump` lies in the CPU assignment of the sniffing process and the

<sup>3</sup> [http://www.caida.org/data/passive/passive\\_2009\\_dataset.xml](http://www.caida.org/data/passive/passive_2009_dataset.xml)



**Figure 2.** Qualitative comparison between the existing high-performance packet sniffing engines.

number of reception queues. In the vanilla case, the program runs in a random processor core selected by the operating system, and the number of reception queues is set to the total number of cores of the machine — 12 in our case. In the optimized case, the capture process runs in the NUMA node attached to the NIC with four reception queues with interrupts also attached to that node. This 4-queue configuration showed the maximum throughput in our setup. With this simple change, the number of stored packets increases up to 20 and 60 percent (from 10 and 38 percent) for synthetic and real traffic, respectively.

To fully sniff high-speed traffic, some engines were proposed [9–13]<sup>4</sup> that apply the aforementioned techniques over the driver and network stacks. Figure 2 shows a qualitative comparison between the most popular capture engines in terms of traffic intensity and number of receive queues needed to work. Specifically, PacketShader, PF\_RING, netmap, and DPDK achieve wire speed using one receive queue. HPCAP obtains similar results but it additionally provides accurate timestamping by using two CPU cores per receive queue, while the other engines use one core per queue. PFQ needs a larger number of queues and cores to give more flexibility and additional functionalities, such as customized packet aggregation, at the expense of performance. The reader is referred to [4] for further details.

## STORING DATA ON HARD DRIVES

According to manufacturers' specifications, a high-end SATA-3 mechanical disk allows theoretical rates up to 4.8 Gb/s for sequential reads and 1.2 Gb/s for sequential writes. A SATA-3 solid-state drive (SSD) may achieve speeds close to 3.2 Gb/s for both read and write operations, but the price per gigabyte is 10 times greater. Consequently, no matter if the hard drive is SSD

or conventional, a single disk is not enough to comply with the line rate of 10 GbE networks, and a RAID volume is in order.

In terms of packet sniffing the worst case scenario is when packets are of minimal size (64 bytes, CRC included, for 10 GbE), so the quantity of packets per second to be processed is maximized — 14.88 Mp/s. However, we note that the worst case for traffic storage is the opposite (maximum sized packets, i.e., 1514 for 10 GbE leading to 816 kp/s), as shown in Table 1. It turns out that the disc load increases with packet size and, to complicate matters, an additional header is aggregated to the packet with the timestamp and both packet and capture lengths. For instance, the de facto PCAP standard stores a 16-byte header per packet — 4 bytes for caplen, 4 bytes for len, and 8 bytes for timestamp. Other packet storage formats such as reduced PCAP (labeled RAW) reduces by 2 the amount of bytes required for each length field. Table 1 shows the storage overhead using PCAP and RAW formats.

Table 1 shows that the requested disk capacities are larger than 9 Gb/s with a worst case of 9.95 Gb/s assuming PCAP headers. As shown, even for the mean Internet packet size (ranging between 256 and 512 bytes according to CAIDA), the demand for store throughput ranges between 9.71 and 9.85 Gb/s, respectively. As a conclusion, once the sniffing engines proved capable of dealing fairly with all packet sizes, now the goal is for the RAID volume to attain rates of nearly 10 Gb/s.

Unfortunately, the amount of parameters involved in RAID configuration is large, and a wrong choice of values may lead to severe performance degradation. We have evaluated the actual write throughput of a RAID-0 volume composed of high-end mechanical hard disks (Hitachi HUA723030ALA640 with 3 TB of capacity) or SSD (Samsung 840 EVO with 250 GB of capacity), both with SATA-3 interfaces, using the out-of-the-box Linux server described previously. Specifically, we have conducted thorough testing with the following configuration parameters.

**Number of disks:** We have assessed how performance varies with number of disks ranging from 1 to 12 for mechanical drives and from 1 to 8 for SSD, merged into a RAID-0 volume; we used an Intel RS25DB080 RAID controller.

**Strip size:** The strip size is the amount of data per basic write operation. Thus, small strip sizes will be translated into a higher number of write operations into the RAID volume and may degrade the overall write throughput due to per-operation overheads. We have evaluated strip sizes of 64 kB, 256 kB, and 1 MB.

**RAID write cache policy:** This parameter refers to the use of the RAID controller's cache memory. The *Direct* policy disables the cache and performs poorly. The *Write Through Cache* (WTC) policy writes the cache content to disk, and then a new cache write operation proceeds. Thus, when using WTC the data has to be stored in both the disks and their caches before a new write operation is started. Finally, the *Write Back Cache* (WBC) policy is less conservative and does not require the cache to be flushed to the

<sup>4</sup> Intel Data Plane Development Kit (Intel DPDK) Release Notes are in <http://www.intel.com/content/dam/www/public/us/en/documents/release-notes/intel-dpdk-release-notes.pdf>

Max. throughput	Packet size (bytes, CRC included)								
	60	64	128	256	512	750	1024	1250	1514
Mp/s	14.88	14.21	8.22	4.46	2.33	1.62	1.19	0.98	0.82
Gb/s	7.14	7.27	8.38	9.13	9.55	9.69	9.77	9.81	9.84
Gb/s (PCAP header included)	9.05	9.09	9.46	9.71	9.85	9.90	9.92	9.94	9.95
Gb/s (RAW header included)	8.57	8.64	9.19	9.57	9.77	9.84	9.89	9.91	9.92

**Table 1.** Maximum rates generated by a fully saturated 10GbE link in terms of packets and bits for different packet sizes and header formats.

hard disks before a new write operation is performed.

**Disk cache:** Some hard drives feature a cache that performs bundling of write operations to a given sector, thus saving disk head movements. We have considered this option in our experiments.

**File system:** We have evaluated the `ext4` file system, which is the de facto standard for Linux systems. Additionally, we have tested the `xfs` and `jfs` file systems, as a previous analysis highlighted them as promising candidates. Specifically, `xfs` was designed with the goal of managing a large number of big files. We have additionally tested the RAID's write throughput when no file system is instantiated as a baseline.

The experiments were carried out by taking all possible combinations of the parameters, and results are shown in Table 2. For each combination, 100 2-GB-sized files were written using the Linux `dd` tool. Specifically, the table depicts which parameter combination offers the minimum and maximum write throughput for different numbers of disks. For each parameter combination we show the mean write throughput, the confidence interval with a 0.01 significance level, and the 5th/95th percentiles.

In practical terms, Table 2 shows that a single mechanical disk has an average write throughput of 1.26 Gb/s for its best configuration, with a narrow confidence interval and a percentile range of roughly tenths of megabits per second. Interestingly, average throughputs scale linearly with the number of disks when they are optimally configured, but this linearity is not observed for the worst-performing parameter combinations. Our findings show that eight disks suffice for all scenarios and packet sizes under study if a properly sized buffer is used to absorb peaks in the write throughput. In fact, the 5th percentile for the throughput obtained with eight disks, which is 9.77 Gb/s, is below the target for some of the most typical scenarios on the Internet assuming both RAW and PCAP headers.

Additionally, we note that a 9-disk RAID exceeds the target rate both in the mean rate and corresponding 5th percentile. This configuration presents a good trade-off to cope with the oscillations that commodity hard drives experience, and it can handle all packet storage sce-

narios for 10 GbE networks even with the largest packet size.

Interestingly, Table 2 also shows that a single SSD drive is capable of consuming nearly twice as much data as a mechanical one: an average of 2.21 Gb/s in contrast with 1.26 Gb/s obtained from the mechanical counterpart. However, our results show that the write throughput scaling ratio (with number of disks) decreases for SSD disks and remains almost constant for mechanical drives. By the time the set of mechanical disks has already achieved an average throughput of 10 Gb/s (i.e., 8 disks), the SSD alternative is far below with the same number of disks: as the table shows, it only reaches 7.52 Gb/s. Moreover, the throughput oscillations of the SSD RAIDs were far larger for all file systems. Specifically, for the best case, it was more than an order of magnitude larger than its mechanical alternative—see the width of the confidence interval for the mean. The above issues, together with cost, discourage the use of SSD disks for our packet storage purpose.

Once we have studied how to gauge a RAID volume, we turn our attention to explain how the configuration parameters and their interactions impact performance. To this end, we posed a balanced full-factorial analysis of the data for a RAID-0 array with nine mechanical hard drives. In such analysis, the response variable under study (in this case the write throughput) is explained as the outcome of a set of factors (i.e., *strip size*, *raid cache policy*, *disks cache*, and *file system*) and their respective levels — that is, possible values each factor takes: *strip size* = 64 kB, 256 kB, or 1 MB; *raid cache policy* = Direct, WB, or WT; *disks cache* = off or on; and *file system* = `ext4`, `jfs`, or `xfs`.

All factors and their interactions turned out to be statistically significant in our dataset. Thus, each sample is characterized by the addition of 16 terms: The overall sample mean (often referred as to intercept): one term per main-effect factor that accounts for the different response given by each level of such factors, and an additional term per each of the possible combinations of factors (six pairs, four trios, and a quartet), which account for the different impact that each combination of levels exerts on the sample.

Several conclusions arise from the analysis,

*Our findings show that eight disks suffice for all scenarios and packet sizes under study if a properly sized buffer is used to absorb peaks in the write throughput. In fact, the 5th percentile for the throughput obtained with eight disks, which is 9.77 Gb/s, is below the target for some of the most typical scenarios on the Internet assuming both RAW and PCAP headers.*

Number of disks	Technology	Scenario	Strip size	RAID cache policy	Disks' cache	FS	Throughput (Gb/s)			
							Average	Confidence interval ( $\alpha = 0.01$ )	Percentile	
									5th	95th
1	Mech	Min	1 MB	WTC	Off	jfs	0.69	(0.67, 0.70)	0.58	0.77
		Max	64 kB	WBC	On	xfS	1.27	(1.26, 1.27)	1.26	1.27
	SSD	Min	64 kB	WTC	Off	xfS	0.58	(0.58, 0.59)	0.55	0.62
		Max	1 MB	WBC	On	jfs	2.21	(2.18, 2.23)	2.19	2.21
8	Mech	Min	64 kB	Direct	Off	jfs	3.64	(3.51, 3.76)	2.90	4.19
		Max	1 MB	WBC	On	xfS	10.06	(10.01, 10.12)	9.77	10.35
	SSD	Min	64 kB	WTC	Off	jfs	2.15	(1.99, 2.31)	1.68	3.82
		Max	1 MB	WBC	On	xfS	7.52	(6.54, 8.51)	3.17	15.97
9	Mech	Min	1 MB	Direct	Off	jfs	4.14	(3.98, 4.30)	3.27	4.81
		Max	1 MB	WBC	On	xfS	11.31	(11.25, 11.37)	10.96	11.47

**Table 2.** Write throughput summary results.

especially the importance of hardware caches. Starting from an overall mean of roughly 4 Gb/s, the use of disk caches represented an average addition of 3.1 Gb/s, and similarly, the use of WBC policy gave an average gain of 3.3 Gb/s. The volume's strip size had a relatively marginal significance, which translated into a few hundred of Mb/s for the best configuration — a strip size of 1 MB. On the other hand, the choice of file system was also significant: *xfS* showed the best results with an increase of 0.7 Gb/s in mean compared to *jfs*, which showed the worst results. Finally, the contribution in absolute value of the terms accounting for all possible combinations of levels was limited, from a few tens to 100 Mb/s.

Furthermore, we found that the file system choice not only affects the average write rate, but also exerts a critical effect on its variance. Figure 3 shows the throughput obtained when writing the same files as in the previous experiment (100 2-GB-sized files) on a 9-disk RAID-0 volume with the optimal configuration for each file system. The figure shows that writing data on the raw volume with no file system present has a low-variance behavior. This is a non-practical scenario because data cannot be accessed afterwards, although it is of interest for baselining purposes. When a file system is incorporated, throughput oscillations happen, which may be severe. Interestingly, both Fig. 3 and Table 2 show that the *xfS* file system presents the smallest oscillation, which makes it the file system of choice. For *jfs* and *ext4*, the throughput oscillation may require adding more disks to the RAID volume to ensure that in case of oscillations the RAID meets the target rate.

## NETWORK TRAFFIC STORAGE SOLUTIONS

Once we have discussed how to optimize both network traffic sniffing and data storage processes separately, we proceed to optimally combine both and come up with a cost-effective high-performance NTSS. More specifically, we outline the fundamental techniques the NTSS may adopt.

**System call minimization:** The sniffing and storage processes imply data transfer and synchronization between user and kernel level contexts. We seek to minimize the quantity of system calls to reduce context switches and improve overall performance by using buffer mapping or accessing data in a bytestream or batch fashion rather than on a per-packet basis.

**Huge intermediate receive buffers:** As mentioned above, the target storage device may experience sudden write throughput drops, which can be avoided by means of large intermediate buffers.

**Memory alignment:** Maximum write performance is achieved when the transfers between system memory and the storage device are done via direct memory access (DMA) operations (`O_DIRECT` flag in the write options). This way, neither CPU cycles nor cache management nor memory bandwidth is spent in the data transfers. However, this efficient configuration requires the transfer operations to be page-aligned, making memory alignment a critical feature. Moreover, packets must be transferred in blocks with a multiple of the disk's sector size to reach maximum write performance and avoid quantization effects.

**Sniffing and storage overlapping:** Overall performance can be increased if the sniffing and storage processes are isolated, allowing their execution to be parallelized.

**Timestamping:** As explained in [1], accurate packet timestamping is a critical issue, which depends on the low-level NTSS sniffing techniques adopted [8].

Note that only two of the approaches explained previously, PF\_RING and HPCAP, gave rise to final NTSS, `n2disk` and `hpcapdd`, respectively. On one hand, `hpcapdd` was developed on top of the HPCAP driver. Both the driver and the application were designed with the goal of optimizing network traffic storage [13]. Regarding the aforementioned techniques, the HPCAP+`hpcapdd` system instantiates a 1 GB kernel-level buffer, limited by the kernel configuration. The driver is in charge of timestamping and copying the incoming packets into this buffer, so `hpcapdd` can access them on a byte-stream basis. This buffer is efficiently accessed as it is properly aligned and mapped at the user level. Furthermore, this buffer isolates the sniffing and storage processes, so the overall process is parallelized and pipelined.

On the other hand, `n2disk` has recently been developed by the authors of PF\_RING [14]. Specifically, `n2disk` instantiates one or more packet storage threads which are executed in parallel, leading to single-threaded (ST) and multi-threaded (MT) versions. Each thread has an independent memory buffer, and traffic is distributed among them using a hash function. Importantly, `n2disk` not only stores the incoming packets, but also creates additional index files for optimizing subsequent access to the stored data.

For both NTSSs, we measured the percentage of incoming traffic stored versus the number of cores with an optimized 9-disk RAID-0 volume — as discussed previously. Figure 4 shows the measured percentage of packets stored over the total link load, along with the amount of fully occupied CPU cores used by each NTSS. Those results are shown for the worst-case scenarios in terms of both packet sniffing and storage and in an average case.

Remarkably, `hpcapdd` is capable of storing 99.2 percent of the incoming traffic for synthetic 64-byte packets (CRC included) and all the real-traffic and synthetic maximum-sized packet traces, with two CPU cores.

The results show that `n2disk`'s single-thread version uses one thread for packet sniffing and one more thread for storage, whereas the multi-thread version uses one thread for sniffing and four threads for processing and storing. The ST version stores 92.4 percent of the packets for the 64-byte experiment, and 98.2 percent for both maximum-sized packets and real traffic. The MT counterpart stores 99.1 percent of the packets for the 64-byte experiment, and 98.0 percent for both maximum-sized packets and real traffic. These last results show that instantiating several threads helps when dealing with the worst case sniffing scenario (i.e., 64-byte packets) but does not solve demanding storage throughput scenarios.

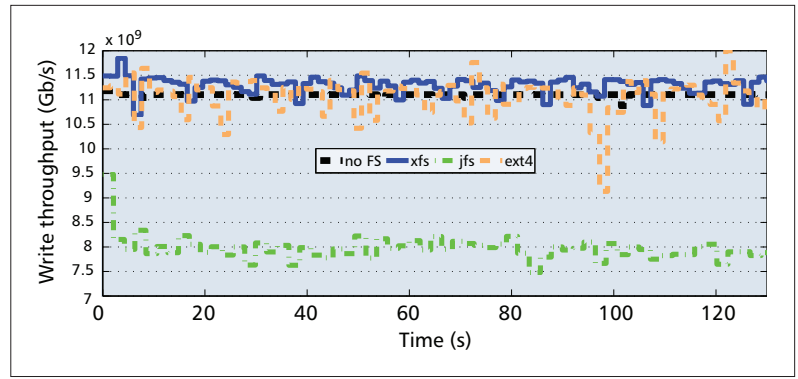


Figure 3. Influence of the file system on a 9-disk RAID-0 volume.

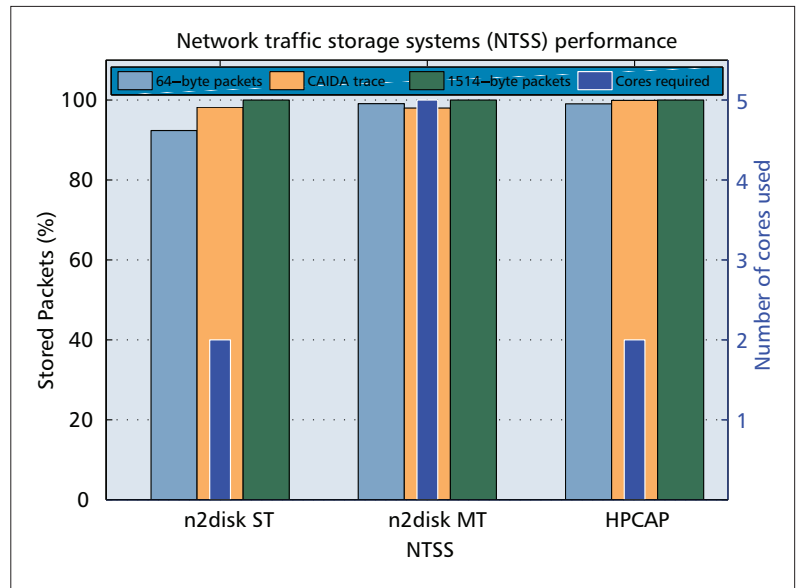


Figure 4. Percentage of stored packets (in a RAID-0 volume) for a fully saturated 10 Gb Ethernet link vs. the number of occupied CPU cores for a 30-minute experiment.

## FINAL REMARKS AND FUTURE WORK

Given the growing importance of COTS systems for network monitoring tasks, this work has addressed the problem of traffic sniffing and storage at 10 Gb/s rates following a bottom-up approach.

First, out-of-the-box tools (i.e., vanilla network drivers and `tcpdump`) have proven to be insufficient to sniff and store packets at 10 Gb/s using off-the-shelf systems. Although the obtained performance is far below 10 Gb/s, the application of some ideas discussed in this work improved the performance of such tools up to rates that may be useful in networks with low utilization.

Second, to improve the sniffing performance and make the most of commodity multi-core servers and modern NICs, several optimizations must be applied. For example, affinity planning is a key factor to improve the performance for both optimized and out-of-the-box applications. Additionally, some of these optimizations present collateral effects such as timestamp accuracy degradation or

To conclude, this work has provided both the research community and practitioners with a roadmap not only to understand and use state-of-the-art NTSS systems based on COTS, but also to implement and deploy their own systems.

packet reordering when applying batch processing and multi-queue reception, respectively. Thus, depending on the application requirements, such optimizations may not apply.

Third, although a single commodity hard drive is not able to achieve enough write throughput to cope with a 10 Gb Ethernet link, such storage performance may be improved by skillfully using tuned RAID volumes. After this tuning, a write throughput beyond 10 Gb/s is achieved using 9 high-end mechanical disks. The write throughput of commodity hard-drives presents significant oscillations over the mean across time. Such oscillations may be controlled using *xfs* file-system, whereas the use of other file-systems causes remarkable excursions. Such excursions are even more noticeable using SSDs. In fact, although one single SSD achieves more throughput than a rotational drive they do not scale proportionally. Thus, SSDs are left behind their mechanical counterpart.

Finally, one important fact is that the most demanding scenario in terms of packet sniffing (i.e., minimal-size packets) is the least demanding scenario in terms of packet storage throughput. Conversely, the best case for packet sniffing (i.e., maximal-size packets) becomes the most demanding scenario in terms of packet storage throughput. Thus, obtaining maximum performance in an NTSS implies not only properly tuning the sniffing and storage processes, but also their interaction.

As future work, the scaling of the mentioned packet storage solutions to higher link rates must be considered — that is, 40 Gb/s and 100 Gb/s. Current hardware limitations may prevent a linear scaling of the results shown along this work, and more complex techniques will be required. Along this line, we propose:

- The study of techniques that allow reducing the amount of information to be stored by smartly selecting the most interesting parts — for example, the first bytes of each flow or packet [1, 15].
- Reducing the magnitude of the stored data by aggregating through some criteria — that is, moving from packet traces to flow records reduces the storage by one order of magnitude.

Needless to say, these techniques come at the expense of potential information losses.

On the other hand, the use of a high-performance traffic distribution would open the possibility of replicating standalone systems to achieve the desired rates, although transferring the problem to the distribution system. Anyway, regardless of the approach chosen, most of the guidelines presented in this article still apply.

To conclude, this work has provided both the research community and practitioners with a roadmap not only to understand and use state-of-the-art NTSS systems based on COTS, but also to implement and deploy their own systems. We expect the lessons and ideas we share here may open new opportunities for the use of COTS systems in areas traditionally reserved for high-end and expensive hardware.

## REFERENCES

- [1] G. Maier et al., "Enriching Network Security Analysis with Time Travel," *Proc. ACM SIGCOMM*, 2008, pp. 183–94.
- [2] Y.-D. Lin et al., "Low-Storage Capture and Loss Recovery Selective Replay of Real Flows," *IEEE Commun. Mag.*, vol. 50, no. 4, 2012, pp. 114–21.
- [3] L. Braun et al., "Comparing and Improving Current Packet Capturing Solutions Based on Commodity Hardware," *Proc. ACM Internet Measurement Conf.*, 2010, pp. 206–07.
- [4] V. Moreno et al., "Commodity Packet Capture Engines: Tutorial, Cookbook and Applicability," *IEEE Commun. Surveys & Tutorials*, to appear.
- [5] J. Zazo et al., "TNT10G: A High-Accuracy 10 GbE Traffic Player and Recorder for Multi-Terabyte Traces," *Proc. Conf. Reconfigurable Computing and FPGAs*, 2014, pp. 1–6.
- [6] L. Deri, "Cap: Wire-Speed Packet Capture and Transmission," *Proc. IEEE/IFIP Workshop on End-to-End Monitoring Techniques and Services*, 2005, pp. 47–55.
- [7] W. Wu, P. DeMar, and M. Crawford, "Why Can Some Advanced Ethernet NICs Cause Packet Reordering?" *IEEE Commun. Lett.*, vol. 15, no. 2, 2011, pp. 253–55.
- [8] V. Moreno et al., "Batch to the Future: Analyzing Timestamp Accuracy of High-Performance Packet I/O Engines," *IEEE Commun. Lett.*, vol. 16, no. 11, 2012, pp. 1888–91.
- [9] S. Han et al., "PacketShader: A GPU-Accelerated Software Router," *Proc. ACM SIGCOMM*, 2010, pp. 195–206.
- [10] F. Fusco and L. Deri, "High Speed Network Traffic Analysis with Commodity Multi-Core Systems," *Proc. ACM Internet Measurement Conf.*, 2010, pp. 218–24.
- [11] L. Rizzo, "Revisiting Network I/O APIs: The Netmap Framework," *Commun. ACM*, vol. 55, no. 3, 2012, pp. 45–51.
- [12] N. Bonelli et al., "On Multi-Gigabit Packet Capturing with Multi-Core Commodity Hardware," *Proc. Passive and Active Network Measurement Conf.*, 2012, pp. 64–73.
- [13] V. Moreno et al., "Packet Storage at Multi-Gigabit Rates using Off-the-Shelf Systems," *Proc. IEEE Conf. High Performance and Commun.*, 2014, pp. 486–89.
- [14] L. Deri, A. Cardigliano, and F. Fusco, "10 Gb Line Rate Packet-to-Disk Using *n2disk*," *Proc. Traffic Monitoring and Analysis Workshop*, 2013, pp. 441–46.
- [15] V. Uceda et al., "Selective Capping of Packet Payloads for Network Analysis and Management," *Proc. Traffic Monitoring and Analysis Workshop*, 2015, pp. 3–16.

## BIOGRAPHIES

VICTOR MORENO (victor.moreno@uam.es) received his B.Sc. degree in mathematics and his M.Sc. degree in computer science from Universidad Autónoma de Madrid (UAM), both in 2010, and is currently a Ph.D. student in computer science. He joined the High Performance Computing and Networking group of the same university in 2008. Since then he has participated in both industrial and European Union projects. In 2011, he was awarded a four-year fellowship by the Ministry of Education of Spain (F.P.U scholarship). His research interest areas include high-performance heterogeneous computing (with GPUs, FPGAs, etc.), big data problems, computer architecture, performance tuning, and statistics.

JAVIER RAMOS (javier.ramos@uam.es) received his M.Sc. degree in computer science and Ph.D. degree in computer science and telecommunications from UAM in 2008 and 2013, respectively. Before that, he joined the Networking Research Group of the same university, where he participates in the European Union project OneLab2 and in the national research project PASITO. His research interests are in the analysis of network traffic, quality of service, bandwidth measurement, and next generation networks.

JOSÉ LUIS GARCÍA DORADO [M'12] (jl.garcia@uam.es) received his M.Sc. and Ph.D. degrees, both in computer and telecommunications engineering, from UAM, Spain, in 2006 and 2010, respectively. He is a member of the Networking Research Group at UAM since 2005 when he began collaborating in national and European R&D projects as an assistant researcher in the ePhoton/One+ Network of Excellence. Since then, he was awarded with a four-year fellowship (F.P.I. scholarship) by the Ministry of Education of Spain (2007), and was a visiting scholar in the Telecommunication Networks Group at Politecnico di Torino, Italy (2010), in the Internet Systems Lab at Purdue University, USA (2013), and in FICA at Universidad Técnica



---

del Norte, Ecuador (2014 and 2015). Currently, he is an assistant professor at UAM with research interests in the analysis of Internet traffic: its management, modeling, and evolution.

IVAN GONZALEZ (ivan.gonzalez@uam.es) received his M.Sc. degree in computer engineering in 2000 and his Ph.D. degree in computer engineering in 2006, both from UAM, Spain. From October 2002 to October 2006 he was a teaching assistant at the Computer Engineering Department of UAM. From November 2006 to January 2008 he was a postdoctoral research scientist at the High Performance Computing Laboratory (HPCL), Electrical & Computer Engineering Department, George Washington University (Washington, DC). He was a faculty member of the NSF Center of High Performance Reconfigurable Computing (CHREC) at George Washington University. His main research interests are heterogeneous computing (with GPUs, FPGAs, etc.), parallel algorithms, and performance tuning. Other interests include FPGA-based reconfigurable computing applications, with a special focus on dynamic partial reconfiguration, embedded systems, and robotics.

FRANCISCO J. GOMEZ-ARRIBAS (francisco.gomez@uam.es) received his Ph.D. from UAM, Spain, in 1996. From October 1996 until November 2000 he was an assistant professor at the Computer Engineering Department of UAM. He

is currently a professor of computer architecture and parallel computing courses at the same university. His research fields of interest concern reconfigurable computing applications based on FPGA circuits, with a special focus on the design of multiprocessor systems with reconfigurable architecture. Secondary fields of interest include network computing, cryptographic coprocessors, embedded system-on-a-chip, and experimental support of computer science and electrical engineering education on the Internet.

JAVIER ARACIL [SM] (javier.aracil@uam.es) received his M.Sc. and Ph.D. degrees (Honors) from the Technical University of Madrid in 1993 and 1995, both in telecommunications engineering. In 1995 he was awarded a Fulbright scholarship and was appointed a postdoctoral researcher of the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley. In 1998 he was a research scholar at the Center for Advanced Telecommunications, Systems and Services of the University of Texas at Dallas. He has been an associate professor for the University of Cantabria and Public University of Navarra and he is currently a full professor at UAM and a founding partner of the spin-off company Naudit HPCN. His research interests are in optical networks and performance evaluation of communication networks. He has authored more than 100 papers in international conferences and journals.

# Platform for Benchmarking of RF-Based Indoor Localization Solutions

Tom Van Haute, Eli De Poorter, Filip Lemic, Vlado Handziski, Niklas Wirström, Thiemo Voigt, Adam Wolisz, and Ingrid Moerman

## ABSTRACT

Over the last few years, the number of indoor localization solutions has grown exponentially, and a wide variety of different technologies and approaches are being explored. Unfortunately, there is currently no established standardized evaluation method for comparing their performance. As a result, each solution is evaluated in a different environment using proprietary evaluation metrics. Consequently, it is currently extremely hard to objectively compare the performance of multiple localization solutions with each other. To address the problem, we present the EVARILOS Benchmarking Platform, which enables automated evaluation and comparison of multiple solutions in different environments using multiple evaluation metrics. We propose a testbed-independent benchmarking platform, combined with multiple testbed-dependent plug-ins for executing experiments and storing performance results. The platform implements the standardized evaluation method described in the EVARILOS Benchmarking Handbook, which is aligned with the upcoming ISO/IEC 18305 standard “Test and Evaluation of Localization and Tracking Systems.” The platform and plug-ins can be used in real time on existing wireless testbed facilities, while also supporting a remote offline evaluation method using precollected data traces. Using these facilities, and analyzing and comparing the performance of three different localization solutions, we demonstrate the need for objective evaluation methods that consider multiple evaluation criteria in different environments.

Tom Van Haute, Eli De Poorter, and Ingrid Moerman are with Ghent University.

Filip Lemic and Vlado Handziski are with Technische Universität Berlin.

Niklas Wirström and Thiemo Voigt are with the Swedish Institute of Computer Science.

<sup>1</sup> <http://www.ict-fire.eu/home.html>.

## INTRODUCTION

This article addresses one of the major problems of indoor localization research: the lack of comparability between existing localization solutions, due to the fact that most of them have been evaluated under individual, thus not comparable and not repeatable, conditions. This situation is partially the result of the complexity required for the evaluation of an indoor localization solution, which requires technical expertise to efficiently set up large-scale experiments, control the exper-

imental environment, gather the necessary performance data, and calculate the output metrics using standardized methods. All these steps are time consuming, and more theoretically inclined researchers typically lack the necessary technical skills to perform these steps efficiently and accurately. We address these deficiencies by providing a platform that allows simple evaluation of indoor localization solutions. The main contributions of the presented article are as follows.

We describe a generic benchmarking platform that implements the standardized evaluation method described in the EVARILOS Benchmarking Handbook (EBH), and is aligned with the upcoming International Organization for Standardization/International Electrotechnical Commission (ISO/IEC) 18305 standard “Test and Evaluation of Localization and Tracking Systems.”

We further describe plug-ins that are available for instantiating the components of the EVARILOS Benchmarking Platform on multiple future Internet research and experimentation (FIRE) facilities.<sup>1</sup>

Finally, we provide open datasets that help in simplifying the process of benchmarking and evaluation of indoor localization solutions.

The rest of this article is structured as follows. The next section provides an overview of the related work. Then the EVARILOS Benchmarking Platform (EBP) is explained in detail. The integration of the EBP in a wireless test facility and the public datasets are then discussed. We then demonstrate the usage of the EBP in an experimental validation of multiple RF-based indoor localization solutions. Finally, we conclude the work.

## RELATED WORK

As the number of indoor localization solutions is growing, a more thorough procedure for evaluating and comparing them is necessary. As already observed in other fields [1], a well defined objective evaluation methodology needs to take into consideration a wide range of metrics. Some metrics are important from a theoretical point of view, and are well suited for analyzing and improving proposed algorithms, whereas others

focus on the performance of end solutions, and are more important for industry and end users. If only accuracy is taken into account, the results can give a distorted view. Such considerations have motivated M. Ficco *et al.* [2] to evaluate indoor localization solutions with respect to deployment metrics. They compare and calibrate the deployment and usage of access points (APs), and show that the quality of the radiomap has a direct influence on the accuracy. Furthermore, Hui Liu *et al.* state in [3] that precision, complexity, scalability, robustness, and cost should be included if a comprehensive performance analysis is required. Additionally, they also recognize the lack of an objective methodology for the evaluation of indoor localization solutions. Motivated by these circumstances, a number of organizations are trying to develop comprehensive standardized evaluation approaches for indoor localization solutions.

**EVARILOS Project:** In the scope of the FP7 EVARILOS project, focused on objective evaluation of RF-based indoor localization solutions, the EBH [4] has been published. The handbook describes a set of evaluation metrics that are important for the evaluation of indoor localization, including different notions of accuracy, functional metrics such as response delays, and deployment metrics such as setup time and required infrastructure. Furthermore, the handbook contains a set of scenarios that describe how to adequately evaluate an indoor localization solution. The project is also the first one to systematically address the effect of interference on indoor localization solutions, although interference is expected to be present at most sites where these solutions are deployed. The EBH includes a wide range of evaluation metrics, including functional metrics, such as response delays, and deployment metrics, such as setup time and required infrastructure.

**ISO:** Recently, the ISO and IEC established a joint technical committee, ISO/IEC JTC 1, focused on proposing a new ISO/IEC 18305 standard, “Test and Evaluation of Localization and Tracking Systems.”<sup>2</sup> Current drafts include evaluation methodologies for a single technology (e.g., Bluetooth), as well as methodologies for the evaluation of full localization solutions, which is in line with the methodology proposed in the EVARILOS project. While this effort is more general in that it also pertains to a wide range of non-RF-based technologies such as motion sensors, thus far it does not include non-accuracy-related metrics such as ease of use or energy consumption. At the time of writing, none of the drafts were publicly available.

**EvAAL:** Until now, the only attempts at direct comparison of different indoor localization solutions were indoor localization competitions. One popular series of indoor localization competitions has been organized by Microsoft as part of the Information Processing in Sensor Networks (IPSN) conference. During the 2014 edition of the competition [5], 22 different indoor localization solutions were evaluated (organized in two categories: infrastructure-free and infrastructure-

based). The evaluation process uses only a single metric: average localization error across 20 test points. The errors are measured manually using a handheld laser distance meter. In 2015, the evaluation process for the 23 competing solutions took more than one day. In 2014 we shadowed the official evaluation process using the EBP presented in this article, and demonstrated the viability and the benefits of a full automation of this process. The Evaluating AAL Systems through Competitive Benchmarking (EvAAL) project<sup>3</sup> uses a set of metrics as part of the evaluation process for its competition series. In addition to the accuracy of indoor localization, usability metrics are defined such as installation complexity, user acceptance, availability, and interoperability with AAL systems. The evaluation process is not automated, and involves deploying physical devices in the environment of interest.

Most scientific papers evaluate the solution they propose in an easily accessible environment in the development area of the authors. Typically, these are office environments with brick walls [6, 7]. Since evaluation is rather time consuming, most localization solutions are evaluated only in a single environment. Both the EVARILOS project and ISO/IEC JTC 1 refer to the fact that this evaluation is not representative for other environments. Therefore, our platform offers developers the possibility to evaluate their localization solutions using input datasets collected in multiple environments: an office environment with brick walls, an office environment with plywood walls, and finally, an industrial-like open-space environment. Since the accuracy strongly depends on the used evaluation points, for example, points near a wall vs. in the middle of a room or in an open space, our public datasets contain data measured at a wide range of measurement points.

## EVARILOS BENCHMARKING PLATFORM

This section describes the EBP.<sup>4</sup> The EBP has been created to address the fact that, although numerous experimental testbed facilities are available [8, 9], evaluating the performance of a localization solution under controlled conditions using standardized performance metrics has proven to be very complicated, in particular for researchers who have limited experience with experimental research. The EBP addresses this issue by providing an open software solution that implements user friendly methods to support the full performance evaluation cycle. The developed software components are independent of any experimental facilities and use open source principles, allowing researchers to download and modify any of the components.

An overview of the EBP architecture is shown in Fig. 1.

**The Rectangles:** Represent components that are available as web services. These components run on a cloud platform where they can be accessed remotely or downloaded to be modified and/or run locally.

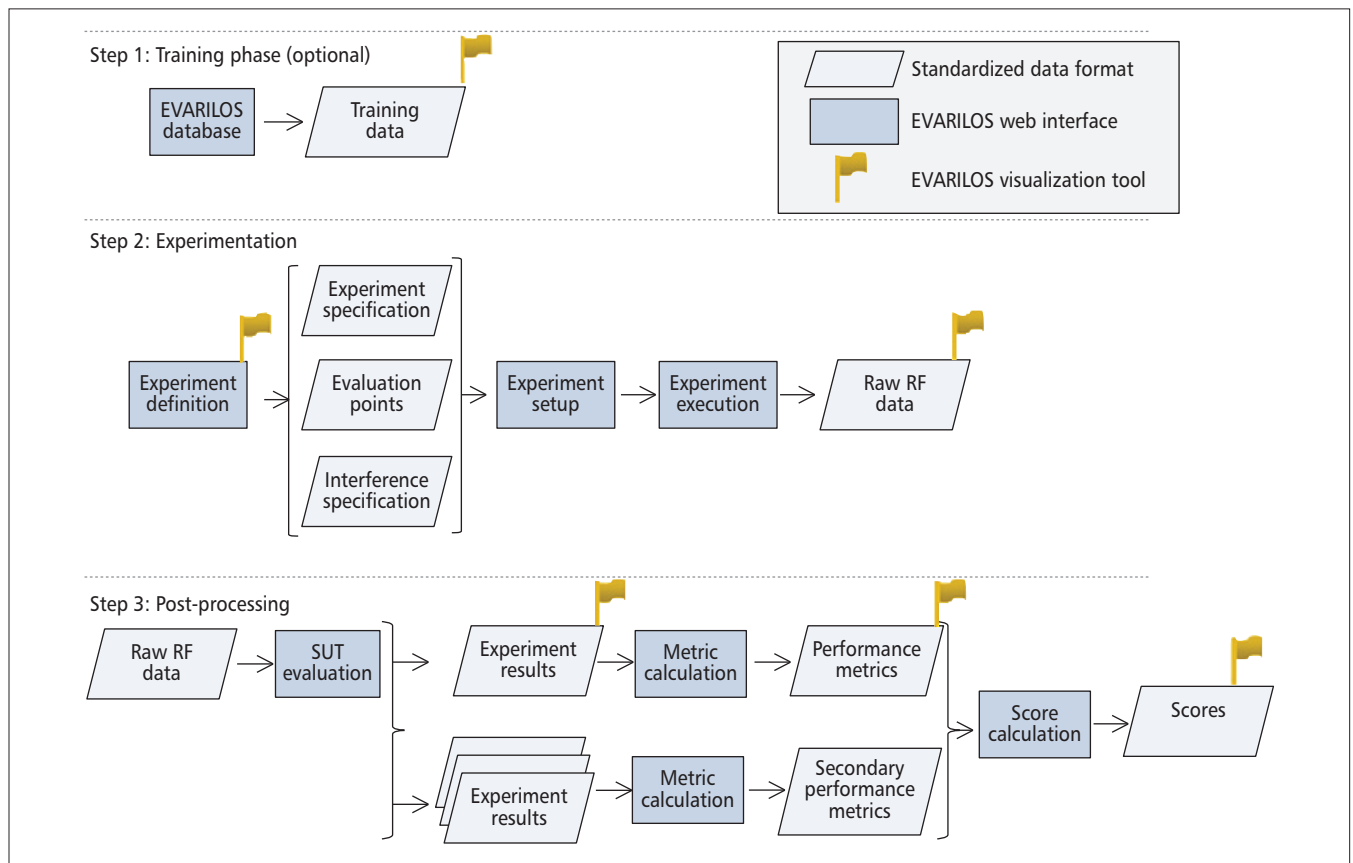
**The Parallelograms:** Represent data struc-

*Our platform offers developers the possibility to evaluate their localization solutions using input datasets collected in multiple environments: an office environment with brick walls, an office environment with plywood walls, and finally, an industrial-like open-space environment.*

<sup>2</sup> <http://www.iso.org>

<sup>3</sup> <http://evaal.aalooa.org>

<sup>4</sup> <http://ebp.evarilos.eu/>



**Figure 1.** Overview of the components of the EBP and the data structures used to exchange information between the components.

tures that are used to exchange data between the web services.

**The Flags:** Represent the tools that can be used to analyze and visualize the different steps of the process.

The architecture consists of a set of components that, when used sequentially, implement a workflow which represents three experimentation steps. A summary can be found below, while in the next subsections each step is discussed in detail.

**Pre-Experimentation Phase:** During a pre-experimentation phase, users can download environment-specific training datasets from public repositories. These datasets are typically used for training the localization solution.

**Experimentation Phase:** In the experimentation phase, all the components required for the experimentation are orchestrated, and the experiments are executed. The platform offers the possibility for automated generation of experiment configurations, including specifications of the used evaluation points, the interference patterns that will be generated, and so on. Based on these descriptions, experiment executables are created using testbed-specific tools with the Control and Management Framework (OMF),<sup>5</sup> which is used in many recent wireless testbeds,<sup>6</sup> and are automatically executed. Note that this step can be omitted if the next step utilizes pre-collected input (e.g., WiFi beacons) for a localization solution.

**SUT:** Finally, the environmental RF data is fed to the system under test (SUT), either in

real-time or using precollected measurements, depending on the experiment configuration. The estimated locations are stored together with additional performance metrics such as the response delay. It is also possible to combine results from multiple experiments to observe how certain evaluation metrics evolve.

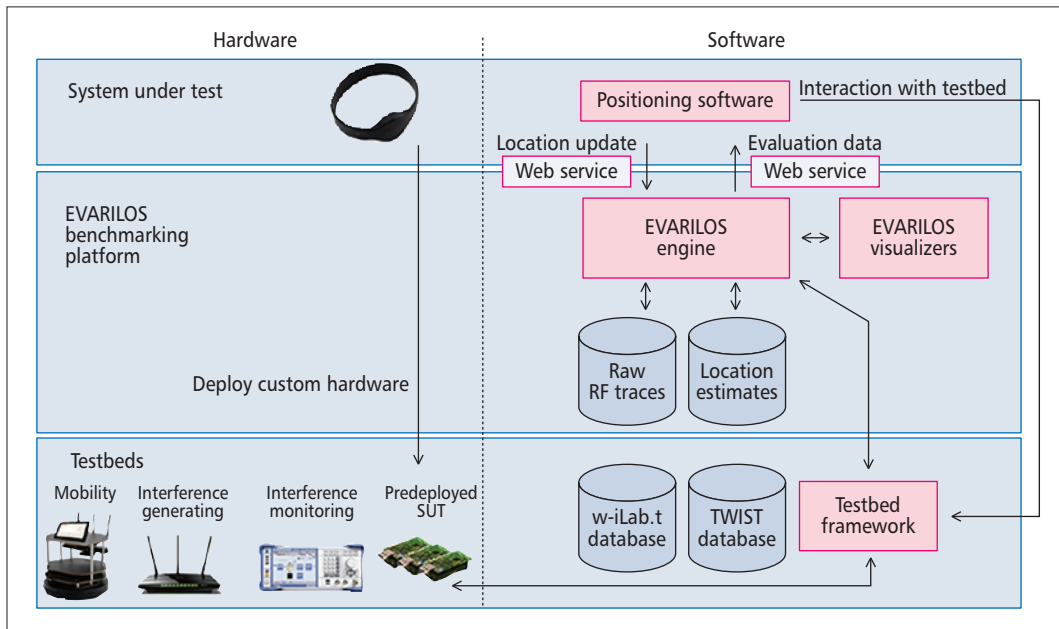
### TRAINING PHASE

The training phase offers experimenters the possibility to train their localization solutions based on measurements that are performed in advance on a representative location. The measurements currently offered represent raw data that can be used as input into an RF-based indoor localization solution, such as received signal strength indicator (RSSI), link quality indicator (LQI), or time of arrival (ToA). Measurements for training purposes are captured in an area that is representative for the experimentation phase. Typically, the data is captured in the same environment where the SUT will be evaluated. To prevent aliasing problems, the training data should not exactly correspond to the data that is used during the evaluation phase. Otherwise, the performance evaluation of step 2 of the evaluation process will be biased. To this end, users can use data that is:

- Captured at a different time
- Captured using devices from a different manufacturer
- Captured at evaluation points other than the one used during the performance evaluation

<sup>5</sup> <http://omf.mytestbed.net/projects/omf6/wiki/Wik>

<sup>6</sup> <http://mytestbed.net/projects/omf/wiki/DeploymentSite>



**Figure 2.** Deployment of the EBP.

The platform offers researchers a database to access previously measured environmental information relevant for their localization solution. Users can either download the data directly from the EVARILOS data repository or can access an EVARILOS application programming interface (API) that encapsulates the data and can serve the data at a finer granularity.

### EXPERIMENTATION PHASE

The experimentation phase offers experimenters the possibility to define setups for raw RF data collection or full localization experiments in FIRE facilities, as well as an interface for automatic execution. The user will start with an “experiment definition” (Fig. 1). The role of the experiment definition component is to configure all aspects of the experiment that will be used to evaluate a SUT. To this end, the experiment definition component requires the following input: the *experiment specification* (e.g., which nodes will be used as anchor points, when will the experiment be scheduled, which binary files to use), the *evaluation points* (at which locations is a SUT evaluated), and the *type of (artificial) interference* that should be generated. To assist with this process, a fully automated web service is available, where users can select among different preconfigured options. Of course, it is possible to modify any of the default settings to adjust the experiment behavior. This information is also stored in a standardized data format.

Next, the “experiment creation” component is executed, which is a fully automated step, whereby the testbed-independent information is translated into testbed-dependent executables using the appropriate plug-ins. The final step is the actual execution of an experiment. In this step, the executables are executed on the corresponding testbed, and the result of the execution is stored in an appropriate data structure together with additional metadata, describing a whole experiment in detail. The result of the execution

is raw data, such as WiFi or IEEE 802.15.4 beacon information, which is collected by a SUT at different locations in an environment.

### POST-PROCESSING PHASE

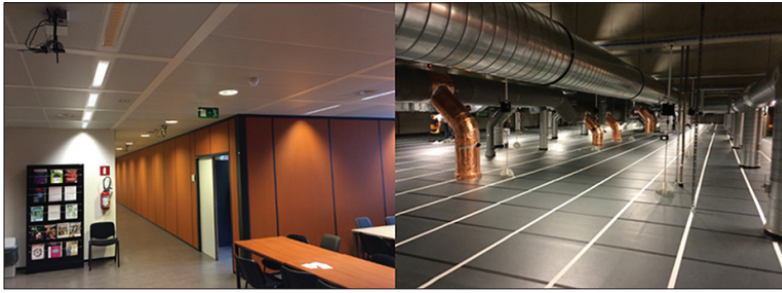
In this step the obtained raw data traces can be fed to the evaluated SUT, and location estimates can be produced. Furthermore, the metrics characterizing the performance of a SUT can be calculated. The experiment results are stored in an appropriate data structure, which consists of a set of ground truths and estimates for different measurement locations and a set of metrics characterizing the performance of a SUT for a given experiment.

Experiment results from multiple experiments can be combined to observe how certain evaluation metrics evolve, for example, for different scenarios or different parametrization of a SUT. These results are stored in a secondary metrics data structure. For comparability purposes, a final score can be assigned to the performance of each SUT. This score is an abstraction of the performance of a SUT in a specific environment and necessarily hides many intrinsic trade-offs. Finally, it is worth mentioning that the full post-processing phase can also be applied to location estimates from non-EBP-compliant solutions. As long as the experiment results are provided in the correct data format, the same tools can be used to analyze and rank the outcome of any localization solution.

## INTEGRATION OF EBP IN WIRELESS EXPERIMENTATION FACILITIES

The EBP is designed to simplify the evaluation of RF-based localization solutions. The components of the platform can be used “as is” by utilizing precollected data traces as input. However, as already mentioned, the platform components can also be used to *facilitate the evaluation of*

*It is worth mentioning that the full post-processing phase can also be applied to location estimates from non-EBP-compliant solutions. As long as the experiment results are provided in the correct data format, the same tools can be used to analyze and rank the outcome of any localization solution.*



**Figure 3.** Two examples of the testbeds (w-iLab.t I and II) where experiments can be executed.

*localization solutions in new environments.* The available deployment options for indoor localization benchmarking are presented in Fig. 2. Three main components can be identified:

- The bottom layer represents a wireless experimentation facility or testbed. The testbed-specific tools are installed on a server in a given test facility.
- The EBP includes services that facilitate testbed-independent definition of experimentation and the evaluation of localization solutions.
- Finally, the upper layer represents a SUT, which can include both hardware and/or software components.

As mentioned, the EBP is integrated in existing FIRE facilities. This integration is part of the experiment execution component illustrated in Fig. 1. Automatic conversion from experiment descriptions to testbed-dependent scripts is supported, thereby integrating and simplifying the complex steps that otherwise need to be taken for objective experimentation. Building on top of the CREW Cognitive Radio testbeds,<sup>7</sup> the infrastructure leverages a robotic mobility platform, which serves as a reference localization system and can transport the localized device in an autonomous and repeatable manner. In addition, the platform uses the capabilities of the CREW testbed infrastructure to generate typical interference scenarios in a reproducible manner. This further improves benchmarking of indoor localization solutions by testing the performance of a SUT under realistic and repeatable interference conditions.

The interaction between a SUT and the EBP is designed to be as simple as possible: at most two REST interfaces [10, 11] are required, depending on the requirements of an experiment. One interface provides location estimates and ground truth information to the EBP, and the other stores the raw data from a SUT or uses the precollected raw data as input to a SUT.

**During an Experiment:** The EBP can issue a request for the location estimate from a SUT through the first REST interface. As such, the minimum requirement for a SUT to comply with the EBP is to provide the location estimate over HTTP upon request.

**The EBP Can Also Request the Real-Time Environmental Data:** (RSSI values, ToA, etc.) from a SUT, which is then stored through a second REST interface. This data can be collected and at a later time be offered to future experimenters as an open data set.

This architecture allows experimenters to choose among different utilization options.

**Option 1:** The evaluation of a localization algorithm using precollected raw data traces that can be used as input to a SUT. In this scenario, the localization algorithms can be evaluated remotely using the EBP.

**Option 2:** The evaluation of a localization solution using software running on an existing wireless testbed. In this scenario, the localization algorithms can run on local hardware that is available at the experimentation facilities.

**Option 3:** The evaluation of localization hardware using a testbed. In this scenario, experimenters can install custom hardware at the experimentation facility while still using the EBP for the evaluation of their solution.

One of the major advantages of the EBP is that all three approaches make use of the same common components.

The feasibility of these options has been demonstrated through the EVARILOS Open Challenge [12], as well as during the Microsoft Indoor Localization Competition (IPSN 2014) [5].

## PUBLIC DATASETS

One of the features of the EBP is the capability to reuse previously collected RF data for offline evaluation of RF-based indoor localization solutions. This feature addresses one of the important challenges for the indoor localization research community: the complex and expensive process of obtaining relevant measurements of RF features from multiple environments. The EBP offers a wide range of available precollected RF data sources through its user interface. However, for those researchers who prefer to download full annotated datasets, the EBP also offers the possibility to download the datasets for research purposes. Two types of datasets are currently available: raw RF traces and performance information.

### RAW RF TRACES

Environmental RF data can be used as a basis for either training an algorithm (e.g., by creating propagation models) or offline evaluation of a SUT. The EBP makes available the measured raw RF traces from multiple environments, including a plywood office environment (w-iLab.t I [8]), a brick office environment (TWIST [13]), an industrial-like environment (w-iLab.t II [8]), a hospital environment, and an underground mine. A view of w-iLab.t I and II is available in Fig. 3. The details about the structure of the raw RF data, exact descriptions of the currently available datasets, and an overview of the services available for using the raw RF data for the evaluation of RF-based indoor localization algorithms can be found in [14].

To evaluate a solution for a wide range of conditions, the raw RF traces contain significantly more data than would be used in a typical operational environment. The datasets are rich in terms of number of collected samples per evaluation point (over 1000 samples per evaluation point), the captured data types (including WiFi beacons, sensor RSSI, and sensor time-of-flight information), the used configuration set-

<sup>7</sup> <http://www.crew-project.eu/>

tings (multiple frequencies, multiple transmission powers), and the used anchor points (data is collected from up to 60 anchor points per evaluation point). This richness of the dataset makes the data relevant for a wide range of interested researchers and allows investigation of how changing any of these parameters influences the performance of the solution. Transforming the over-dimensional dataset into a set that is more sparse (and more realistic from an operational point of view) can easily be done by removing any unnecessary information (sub-sampling). In addition, the available environment data is annotated with metadata describing the exact conditions in which the data was captured. This metadata describes characteristics including the used hardware, type of collected raw data, timestamps, measurement frequency, environment description, and so on.

### PERFORMANCE INFORMATION

The EBP gives a ranked overview of evaluated solutions on its web page. However, these performance indicators necessarily hide a number of low-level statistics. Researchers interested in also evaluating the temporal or spatial behavior of different solutions can analyze the performance datasets. EBP makes available the results from its own localization solutions, as well as of those solutions that participated in the EVARILOS Open Challenge [12]. Each of these datasets also has its associated experiment configuration settings, allowing detailed analysis not only of the performance but also of the conditions in which the solutions were evaluated.

## EXPERIMENTAL VALIDATION

In [14] we illustrate the benefits of leveraging the presented platform for the evaluation of RF-based indoor localization, in terms of time and complexity of usage, in comparison to using an infrastructure or performing a manual evaluation. In the following we demonstrate the need for a standardized evaluation method by showing that the performance of localization solutions depends strongly on its parametrization and can only be done objectively by considering multiple evaluation metrics.

### THREE INDOOR LOCALIZATION SOLUTIONS

In order to develop, test, and optimize our platform, three different types of indoor localization solutions were used as SUTs. The basic concept behind the first localization solution [15] is the following: measurements are performed by requesting a stationary node to transmit packets to the testbed nodes that then reply with a hardware acknowledgment (ACK). The initiating node measures the time between transmission of the packet and reception of the ACK, and stores the RSSI values associated with the ACK. These measurements are then processed using Spray, a particle-filter-based platform [15]. The basic idea of the ToF ranging is to estimate the distance between two nodes by measuring the propagation time, which is linearly correlated to the distance when the nodes are in the line of sight (LoS).

A second solution [16] is based on fingerprinting. Fingerprinting methods for indoor

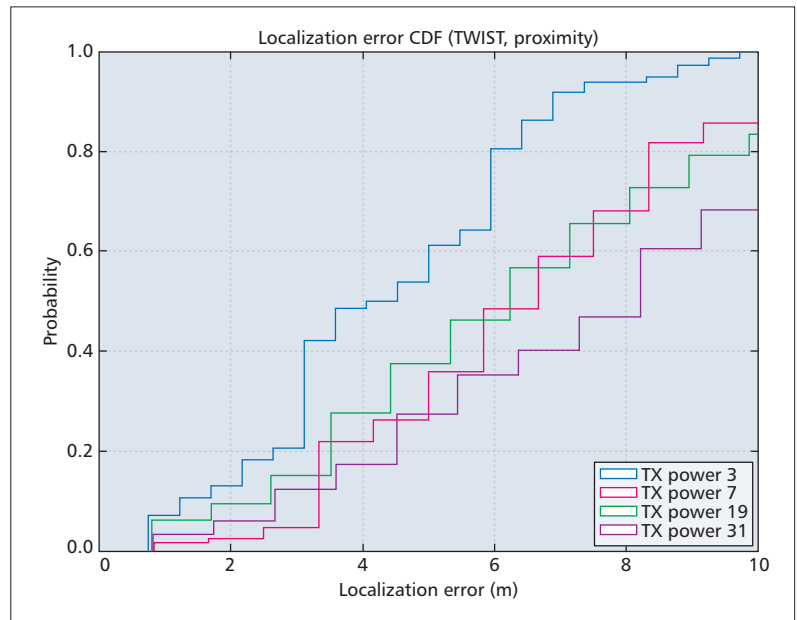


Figure 4. CDFs for the hybrid solution in the TWIST testbed.

localization are generally divided in two phases. The first phase is called training or offline phase. In this phase, the localization area is divided in a certain number of cells. Each cell is scanned a certain number of times for different signal properties, and using a methodology for processing the received data, a representative fingerprint of each cell is created. Using the obtained training fingerprints, the training database is created and stored on a localization server. In the second phase, known as the runtime or online phase, a number of scans of the environment are created using the user's device. From the scanned data, using the same predefined data processing methodology, a runtime fingerprint is created and sent to the localization server. At the server's side, the runtime fingerprint is compared to the training dataset using a matching method. The training fingerprint with the most similarities to the runtime fingerprint is reported as the estimated position.

A third localization solution [17] that has been implemented and evaluated is a hybrid combination range-based and range-free algorithm. It includes a range-based location estimator based on weighted RSSI values. Each RSSI value can be matched with a certain distance. The proposed algorithm in [17] not only uses the RSSI values to measure the distance between a fixed node and a mobile node, but also the distance between fixed nodes. These values function as weight factors for the distance calculation between the fixed and mobile nodes. Once the distances are known, triangulation can be applied in order to determine the final position of the person/object that needs to be localized. This approach is combined with a range-free algorithm, which does not take RSSI-values into account. If a mobile sensor node has a range of 10 m, a fixed node can only receive its messages if the mobile node is maximum 10 m away. This is the only information that is used to calculate the position of a mobile node. For this approach,

it is important that the transmission power is well configured. If the power is too low, the mobile node could be out of range between two fixed nodes. On the other hand, if the power is too high, too many fixed nodes will receive the beacon, and a wrong estimation could be made.

### ANALYSIS OF A SINGLE SOLUTION

An important feature of the EBP is its capability to streamline the process of obtaining better insight on the evaluated localization solution. Every solution contains a set of adjustable parameters, which can considerably influence the overall performance, implying that optimizing this set of parameters can be a hard task. Therefore, the EBP can easily compare the same solution using multiple values of a single parameter.

This can be demonstrated with an example. The hybrid solution [17] described in the section above states that the transmission power is an important value that needs to be configured well

in order to receive acceptable results. Therefore, the solution was evaluated using the EBP using multiple transmission powers, the outcome of which is shown using a cumulative distribution function (CDF) (Fig. 4) and a table with multiple metrics (Table 1). Based on these results, it is clear that this solution obtains the lowest average error when the transmission power equals three. But it also illustrates inherent trade-offs that are present in the solution: suppose the response time is the most important criteria; then a transmission power of 31 would be the best option. This example illustrates the advantages of the EBP for fast and efficient identification of an optimal operating point depending on adjustable parameters, and demonstrates the need for considering multiple metrics to identify trade-offs.

### COMPARISON BETWEEN MULTIPLE SOLUTIONS

Table 2 compares the performance of three different solutions evaluated using the EBP by considering multiple evaluation criteria. By utilizing the same evaluation points, objective comparisons are possible. Again, the results illustrate the presence of trade-offs that can only be observed by comparing multiple metrics. More specifically, it demonstrates that the approach taken in most current scientific papers, wherein point accuracy is considered as the only relevant metric, fails to take into account the associated costs in response time and energy consumption.

## CONCLUSION

The proliferation of RF-based indoor localization solutions raises the need for testing systems that enable objective evaluation of their functional and non-functional properties. Although a significant number of localization solutions are available, different approaches are used for the evaluation of these solutions in terms of used performance metrics and evaluation methodology. This article tries to address these shortcomings by providing tools for evaluating and comparing localization solutions using standardized evaluation methods, as described in the EBH.

We introduce a testbed-independent benchmarking platform for automatized benchmarking of RF-based indoor localization solutions. Using a well defined interface, the infrastructure obtains location estimates from the SUT, which are subsequently processed in a dedicated metrics computation engine. The components can be accessed through web services that are available for external users or can be downloaded for custom modifications. The benchmarking platform has proven to be useful for locations where no testbed facilities are available. Multiple components of the platform were extensively used during the Microsoft Indoor Localization Competition (IPSN 2014) as well as the EVARILOS Open Challenge. In these events, the components of the benchmarking platform improve the time efficiency and ease of use of the experiments, and also resulted in more objective comparability.

Finally, to accommodate the need for wider accessibility of experimental data, open datasets are provided. These datasets include both annotated localization data from multiple environments, as well as detailed descriptions of the

Metric	TX 3	TX 7	TX 19	TX 31
Average error (m)	4.63	7.08	6.93	8.31
Min. error (m)	0.75	0.83	0.80	0.82
Max. error (m)	10.20	17.52	18.93	19.31
Median error (m)	4.39	6.81	6.68	8.63
Room accuracy (%)	26.67	6.70	13.45	9.56
Response time (ms)	1503	1507	480	460

**Table 1.** Statistical information about the performance of the hybrid solution in the TWIST testbed.

Algorithm	Mean error (m)	Room acc. (%)	Latency (ms)	Energy eff. (mW)	
				Mobile	Fixed
<b>Particle filter solution</b>					
Using RSSI	4.35	45.00	14,285	~105	~105
Using ToA	5.56	30.00	14,282	~105	~105
<b>Fingerprinting solution</b>					
Using ED distance	2.2	80.0	~35,000	~7000	~500
Using PH distance	2.0	85.0	~35,000	~7000	~500
<b>Hybrid solution</b>					
TX Power = 3	4.6	26.7	1503	~30.9	~47.4
TX Power = 7	7.1	6.7	1507	~35.1	~47.4

**Table 2.** TWIST testbed: summarized results.



setup and outcome of the performed localization experiments from earlier experiments. These repositories can be used to quickly evaluate a SUT in different environments, analyze the effects of changing configuration settings, analyze the setup of different experiments, and compare the performance of a wide range of localization solutions.

### ACKNOWLEDGMENT

The research leading to these results has received funding from the European Union's Seventh Framework Program (FP7/2007-2013) under grant agreement no 317989 (STREP EVARILOS). The author Filip Lemic was partially supported by DAAD (German Academic Exchange Service).

### REFERENCES

- [1] M. Seltzer *et al.*, "The Case for Application-Specific Benchmarking," *Proc. 7th Wksp. Hot Topics in Op. Sys.*, 1999.
- [2] A. N. M. F. C. Esposito, "Calibrating Indoor Positioning Systems with Low Efforts," *IEEE Trans. Mobile Computing*, vol. 13, no. 4, 2014.
- [3] H. Liu *et al.*, "Survey of Wireless Indoor Positioning Techniques and Systems," *IEEE Trans. Systems, Man, and Cybernetics, Part C: Applications and Reviews*, vol. 37, no. 6, 2007.
- [4] T. V. Haute *et al.*, "The EVARILOS Benchmarking Handbook: Evaluation of RF-based Indoor Localization Solutions," *MERMAT 2013*, May 2013.
- [5] D. Lymberopoulos *et al.*, "A Realistic Evaluation and Comparison of Indoor Location Technologies: Experiences and Lessons Learned," *IPSN '15*, 2015.
- [6] K. Chintalapudi *et al.*, "Indoor Localization without the Pain," *Proc. 16th Int'l. Conf. Mobile Computing and Networking*, ACM, 2010.
- [7] E. Martin *et al.*, "Precise Indoor Localization Using Smartphones," *Proc. Int'l. Conf. Multimedia*, ACM, 2010.
- [8] S. Bouckaert *et al.*, "The w-ilab. t testbed," *Testbeds and Research Infrastructures, Development of Networks and Communities*, Springer, 2011.
- [9] F. Lemic *et al.*, "Infrastructure for Benchmarking RF-Based Indoor Localization under Controlled Interference," *Proc. UPINLBS'14*, 2014.
- [10] F. Lemic, "Service for Calculation of Performance Metrics of Indoor Localization Benchmarking Experiments," tech. rep. TKN-14-003, 2014.
- [11] F. Lemic and V. Handziski, "Data Management Services for Evaluation of RF-Based Indoor Localization," tech. rep. TKN-14-002, 2014.
- [12] F. Lemic *et al.*, "Experimental Evaluation of RF-Based Indoor Localization Algorithms under RF Interference," *Proc. ICL-GNSS'15*, 2015.
- [13] V. Handziski *et al.*, "TWIST: A Scalable and Reconfigurable Testbed for Wireless Indoor Experiments with Sensor Network," *Proc. RealMAN'06*, 2006.
- [14] F. Lemic *et al.*, "Web-Based Platform for Evaluation of RF-Based Indoor Localization Algorithms," *Proc. IEEE ICC Wksp.*, 2015.
- [15] N. Wirstrom, P. Misra, and T. Voigt, "Spray: A Multimodal Localization System for Stationary Sensor Network Deployment," *Proc. IEEE Wireless On-Demand Network Systems and Services*, 2014.
- [16] F. Lemic, "Benchmarking of Quantile based Indoor Fingerprinting Algorithm," Tech. Rep. TKN-14-001, 2014.
- [17] T. Van Haute *et al.*, "A Hybrid Indoor Localization Solution Using a Generic Architectural Framework for Sparse Distributed Wireless Sensor Networks," *Proc. IEEE Comp. Sci. and Info. Systems*, 2014.

### BIOGRAPHIES

TOM VAN HAUTE is a doctoral researcher at Ghent University. He received his M. Sc. degree (cum laude) in computer science engineering from Ghent University, Belgium, in 2012. In September 2012, he joined the Department of Information Technology (INTEC) at Ghent University. Within this department, he is working in the Internet Based Communication Networks and Services research group (IBCN). His research is focused on wireless sensor networks combined with indoor localization and indoor navigation in particular.

ELI DE POORTER is a postdoctoral researcher at Ghent University. He received his Master's degree in computer science engineering from Ghent University, Belgium, in 2006. He received his Ph.D. degree in 2011 from the Department of Information Technology at Ghent University through a Ph.D scholarship from the Institute for Promotion of Innovation through Science and Technology in Flanders (IWT-Vlaanderen). After obtaining his Ph.D., he received an FWO postdoctoral research grant and is now a postdoctoral fellow in the same research group.

FILIP LEMIC is a junior researcher and Ph.D. candidate in the Telecommunication Networks Group at the Technical University of Berlin. He finished his Bachelor's and Master's studies with the Faculty of Electrical Engineering and Computing at the University of Zagreb. His main scientific interests are in context awareness, with an emphasis on indoor localization.

VLADO HANDZISKI is a senior researcher in the Telecommunication Networks Group at Technische Universität Berlin, where he coordinates the activities in the areas of sensor networks, cyber-physical systems, and the Internet of Things. He is currently also serving as interim professor at the chair for Embedded Systems at Technische Universität Dresden. He received his doctoral degree in electrical engineering from TU Berlin (summa cum laude, 2011) and his M.Sc. degree from Ss. Cyril and Methodius University in Skopje (2002).

NIKLAS WIRSTRÖM is a researcher in the Networked Embedded Systems (NES) group at SICS and a Ph.D. student at Uppsala Universitet, Sweden. His research focus is on machine learning techniques for localization in WSNs and other resource constrained systems.

THIEMO VOIGT is a professor at Uppsala University. He also leads the NES group at SICS Swedish ICT. His main interests are networking and system issues in wireless sensor networks and the Internet of Things. He has published papers at flagship sensor networking conferences such as ACM SenSys and IEEE/ACM IPSN, and received awards for several of these publications. He has also been TPC Co-Chair for IEEE/ACM IPSN and EWSN.

ADAM WOLISZ received his degrees (Diploma 1972, Ph.D. 1976, Habil. 1983) from Silesian University of Technology, Gliwice, Poland. He joined TU-Berlin in 1993, where he is a chaired professor in telecommunication networks and executive director of the Institute for Telecommunication Systems. He is also an adjunct professor at the Department of Electrical Engineering and Computer Science, University of California, Berkeley. His research interests are in architectures and protocols of communication networks.

INGRID MOERMAN received her degree in electrical engineering (1987) and her Ph.D. degree (1992) from Ghent University, where she became a part-time professor in 2000. She is a staff member of the research group on Internet-Based Communication Networks and Services ([www.ibcn.intec.ugent.be](http://www.ibcn.intec.ugent.be)), where she leads the research on mobile and wireless communication networks. In 2006 she joined iMinds, where she coordinates several interdisciplinary research projects.

*The proliferation of RF-based indoor localization solutions raises the need for testing systems that enable objective evaluation of their functional and non-functional properties. Although a significant number of localization solutions are available, different approaches are used for the evaluation of these solutions in terms of used performance metrics and evaluation methodology.*

# How Far Is Facebook from Me? Facebook Network Infrastructure Analysis

Reza Farahbakhsh, Angel Cuevas, Antonio M. Ortiz, Xiao Han, and Noel Crespi

## ABSTRACT

Facebook is today the most popular social network with more than one billion subscribers worldwide. To provide good quality of service (e.g., low access delay) to their clients, FB relies on Akamai, which provides a worldwide content distribution network with a large number of edge servers that are much closer to FB subscribers. In this article we aim to depict a global picture of the current FB network infrastructure deployment taking into account both native FB servers and Akamai nodes. Toward this end, we have performed a measurement-based analysis during a period of two weeks using 463 Planet-Lab nodes distributed across 41 countries. Based on the obtained data we compare the average access delay that nodes in different countries experience accessing both native FB servers and Akamai nodes. In addition, we obtain a wide view of the deployment of Akamai nodes serving FB users worldwide. Finally, we analyze the geographical coverage of those nodes, and demonstrate that in most of the cases Akamai nodes located in a particular country service not only local FB subscribers, but also FB users located in nearby countries.

Reza Farahbakhsh, Xiao Han, and Noël Crespi are with Institut Mines-Telecom, Telecom Sud-Paris.

Ángel Cuevas is with Universidad Carlos III de Madrid and Institut Mines-Telecom, Telecom Sud-Paris.

Antonio M. Ortiz is with Montimage and Institut Mines-Telecom, Telecom Sud-Paris.

<sup>1</sup> <http://www.alex.com/topsites>

<sup>2</sup> <http://www.akamai.com>

<sup>3</sup> Akamai Facts & Figures, 2014, [www.akamai.com/html/about/facts\\_figures.html](http://www.akamai.com/html/about/facts_figures.html)

<sup>4</sup> We refer to any activity that a regular FB subscriber can perform when she/he is connected to FB as an FB service, including visualizing pictures, watching videos, gaming, chatting, and so on.

## INTRODUCTION

Facebook (FB) is the most popular online social network (OSN) with more than 1 billion subscribers all over the world. According to Alexa Ranking,<sup>1</sup> FB is the second most popular website in the world. A system of that dimension needs to be sustained by a robust and reliable architecture. Toward this end, FB owns and manages a number of centralized data centers located in the United States and Ireland [1]. However, those data centers are far from a large number of FB subscribers, who could incur very high delays to reach them. Access delay is a very sensitive parameter that impacts user experience and may have a very negative effect on online services if it is not bounded. Some illustrative examples of the actual relevance of delay reported in [2] are:

- 100 ms delay increment implies 1 percent sales loss for Amazon.
- An extra latency of 400 ms reduces Google search volume by 0.74 percent.
- 500 ms of delay decrements the revenue per user in Bing by 1.2 percent.

These numbers state that the lower the delay, the better the quality of experience of the users.

Therefore, to provide efficient service, a worldwide popular system like FB needs to rely on a distributed infrastructure that provides subscribers good quality of service (e.g., low access delay). To achieve this goal FB uses Akamai,<sup>2</sup> a content distribution network (CDN) with 170,000 servers deployed in 102 countries, which delivers between 15–30 percent of all web traffic.<sup>3</sup>

In this context, an intriguing question that motivates our research is how this complex infrastructure offers FB services<sup>4</sup> to FB subscribers, and whether all countries experience the same quality of service in terms of their delay in accessing those services. The goal of this article is to present a rigorous measurement study that allows us to construct the actual FB infrastructure (including Akamai servers) and see how it is being used to meet subscribers' demand.

To answer the previous question, it is essential to determine how the Akamai servers that offer FB services are distributed around the world, and to which Akamai locations FB subscribers are redirected when they access a particular service. Toward this end, we followed a systematic methodology that allows us to identify which Akamai servers are offering what FB services as well as geolocating them. This methodology is composed of four basic steps:

- Identify the URLs associated with FB services.
- Execute ping and traceroute commands from edge machines distributed worldwide to extract IP addresses associated with servers attending queries related to the discovered FB services.
- Geolocate those IPs and determine which ones are associated with native FB servers and which ones belong to Akamai servers.

- Determine which source nodes (in which locations) are assisted by which Akamai servers.

To apply this methodology we used 463 Plan- etLab (PL) [3] nodes distributed across 41 countries all over the world, which sent ping and traceroute probes to 47 different FB URLs (grouped into 16 different service categories) six times a day for two weeks, from May 7 to May 21, 2013. Overall we collected almost 2 million delay samples from PL nodes to FB native servers and Akamai nodes.

Based on the results obtained from our measurements, we present a discussion that mainly covers two aspects:

- The quality of service (in terms of delay) experienced by subscribers depending on their location
- The picture of where Akamai nodes offering access to FB services are located and which geographical areas they cover (i.e., whether an Akamai node located in country A only receives queries from nodes located in that country or if it also serves nodes in other countries, and in such cases whether these are neighboring countries or not)

The results of our research serve as a solid benchmark to understand the performance offered by CDNs to large demanding clients with hundred of millions of subscribers distributed all over the world. Therefore, researchers aiming to improve CDN services could use the results presented in this article to validate their solutions with respect to the performance offered by the largest commercial CDN. In addition, it opens a door to the networking community to analyze what are the main sources of delay in order to propose solutions that minimize end users' access delay to services like FB. Finally, the simple but efficient methodology employed in the article can be replicated with other online sites and CDNs to perform comparative analysis to our work.

## METHODOLOGY

The goal of this article is twofold: to analyze the user experience in accessing FB services from different countries in terms of latency, and to describe a geographical picture for the location of those servers (with a special focus on Akamai nodes) offering FB services, and, linked to that location, whether they only cover a local region or also serve users located in different countries. Toward this end, we have employed a simple yet meaningful methodology that could be replicated to evaluate the performance in terms of the access delay a CDN offers to a particular website. Next, we define in detail the steps followed in our methodology.

### **Step 1. Identify URLs Associated with the Service Offered by the Website (i.e., FB):**

We asked several Facebook subscribers to perform a number of activities in FB such as login to the site, access their profiles, access photos and videos, and access friends' content. In parallel, we used a network protocol analyzer tool [4] that collected all the traffic associated with each of the described actions. After a simple filtering

Service category	#URLs	Service provider
Access website	2	Facebook and Akamai
Authentication	4	Facebook and Akamai
Blog site	1	Facebook
Chat	2	Facebook
Developer site	1	Facebook
Error	1	Facebook
Friend finder	1	Akamai
Friend site	1	Facebook
Game applications	3	Facebook
Group site	1	Facebook
Multiple services	4	Facebook and Akamai
News feed	4	Facebook
Photo upload	1	Facebook
Photo view	19	Facebook and Akamai
Post site	1	Facebook
Video view	1	Akamai

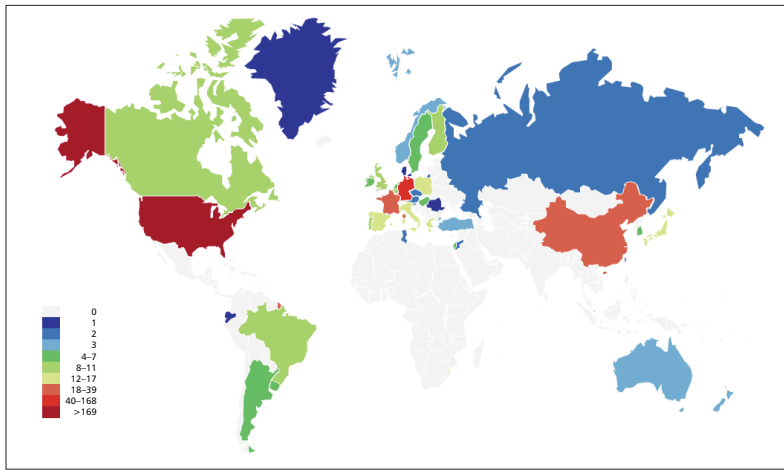
**Table 1.** Facebook service categories, number of URLs for each service, and the service provider (Facebook and/or Akamai).

of the network traces we could map each FB action to one (or more) URLs that could refer to either an FB native server (e.g., profile.facebook.com) or an Akamai server (e.g. photos-a.ak.fbcdn.net). We identified 47 URLs that correspond to 16 different FB services. To be sure that the URLs were not location-dependent we repeated this exercise on several machines at different geographical locations leading to the same results. Table 1 shows the 16 identified FB service categories included in this study as well as the information on which service provider, Akamai and/or FB, is in charge of replying to the queries for these services.

### **Step 2. Script to Measure Access Delay and Network Path to the URLs:**

We implemented a simple script, following a standard discovery method [5], that executes ping and traceroute operations from the machine where it is executed to all 47 identified URLs. The ping measures the latency from the source node to the queried server, which served us to evaluate the performance in terms of access delay. The traceroute reports the intermediate hops between the

*FB owns and manages a number of centralized data centers located in the United States and Ireland [1]. However, those data centers are far from a large number of FB subscribers, who could incur very high delays to reach them. Access delay is a very sensitive parameter that impacts user experience and may have a very negative effect.*



**Figure 1.** Presence and distribution of the 463 PlanetLab nodes (PL\_node) per country.

Country	Acr.	#PL_node	Country	Acr.	#PL_node
United States	US	169	Argentina	AR	4
Germany	DE	40	Hungary	HU	4
China	CN	19	Korea, Rep.	KR	4
France	FR	18	Netherlands	NL	4
Italy	IT	16	Australia	AT	3
Poland	PL	16	New Zealand	NZ	3
Spain	ES	16	Norway	NO	3
Greece	GR	12	Singapore	SG	3
Japan	JP	12	Slovenia	SI	3
Switzerland	SZ	12	Turkey	TR	3
Canada	CA	11	Austria	AT	2
United Kingdom	UK	11	Czech Rep.	CZ	2
Belgium	BE	9	Jordan	JO	2
Brazil	BR	8	Puerto Rico	PR	2
Finland	FI	8	Russia	RU	2
Portugal	PT	8	Taiwan	TW	2
Israel	IL	6	Tunisia	TN	2
Sweden	SE	6	Denmark	DK	1
Hong Kong	HK	5	Ecuador	EC	1
Ireland	IE	5	Romania	RO	1
Uruguay	UY	5			

**Table 2.** Distribution of the 463 PlanetLab nodes (PL\_node) per country.

source node and the server, and the delay to each hop (in case the intermediate router accepts ICMP traffic). The traceroute results may serve to dig into the particular reasons why a particular source node-server path is incurring unexpected delays and try to identify the elements in the paths leading to that situation. However, that individualized analysis goes beyond the scope of this article and would require an article itself.

**Step 3. Create a Distributed Infrastructure to Obtain Comprehensive Results from Different Geographic Locations:** The goal of this research required measuring access delay to the servers serving the 47 URLs from a large number of source machines distributed all over the world. For this purpose we relied on PL [3]. In particular, we distributed the script described in step 2 across 463 PL nodes located in 41 different countries (Fig. 1) as shown in Table 2. In addition, in order to have a large enough and robust dataset that avoids eventual network effects which could corrupt the average delay results, we ran the script six times a day (every four hours at the same time across all machines) in each PL node during a period of two weeks from May 7 to May 21, 2013. Our dataset contains more than 2 million ping and traceroute probes.

**Step 4. Source Nodes, FB Servers, and Akamai Servers Geolocation:** Until this step we have a large dataset in which each ping probe is associated with a source IP address (i.e., PL node), destination IP address (i.e., FB or Akamai server), and delay. However, in order to perform the study described in the introduction we have to geolocate each IP address so that for each ping entry in our dataset we also know location of source node and location of destination node. To geolocate each source node, FB server and Akamai server we used the Maxmind database<sup>5</sup> to bind each IP address to its respective location. The location included country and city (if available).

We would like to note that the final dataset employed in our research is publicly available for the research community.<sup>6</sup>

## END USERS' ACCESS DELAY TO FACEBOOK SERVICES

In this section we aim to understand the performance level experienced by end users in terms of the latency in accessing FB services located in either native FB or Akamai servers. Table 3a shows the detail of the average access delay (and its standard deviation) per country to access FB services in servers located in FB facilities, and Table 3b shows the same parameters in relation to Akamai servers. In addition, Fig. 2 shows the average access delay to access FB services in servers located at FB facilities (Fig. 2a)

<sup>5</sup> <http://www.maxmind.com/>

<sup>6</sup> [http://www.it.uc3m.es/acrumin/papers/FB\\_Arch\\_project.rar](http://www.it.uc3m.es/acrumin/papers/FB_Arch_project.rar)

(a) Facebook				(b) Akamai			
Country	Avg.Delay (ms) $\pm$ STD	Country	Avg.Delay (ms) $\pm$ STD	Country	Avg.Delay (ms) $\pm$ STD	Country	Avg.Delay (ms) $\pm$ STD
(1)				(1)			
Singapore	193.66 $\pm$ 59.41	Brazil	169.78 $\pm$ 60.93	China	174.59 $\pm$ 213.30	Argentina	124.98 $\pm$ 79.67
Romania	190.07 $\pm$ 50.55	Israel	167.14 $\pm$ 90.85	Uruguay	157.40 $\pm$ 78.98		
China	187.14 $\pm$ 227.29	Australia	164.11 $\pm$ 43.22	(2)			
Uruguay	179.96 $\pm$ 65.08	Argentina	155.38 $\pm$ 67.49	New Zealand	95.98 $\pm$ 83.41	Hong Kong	71.41 $\pm$ 80.92
Portugal	177.91 $\pm$ 69.02	New Zealand	152.02 $\pm$ 38.00	Korea, Rep.	90.14 $\pm$ 90.75	Jordan	68.72 $\pm$ 38.89
Slovenia	169.86 $\pm$ 48.50			Australia	87.03 $\pm$ 89.32	Tunisia	63.05 $\pm$ 27.39
(2)				Ecuador	79.62 $\pm$ 55.81	Israel	54.64 $\pm$ 78.04
Denmark	140.93 $\pm$ 38.52	Ecuador	106.69 $\pm$ 36.66	Brazil	78.22 $\pm$ 68.44		
Finland	137.12 $\pm$ 61.17	Tunisia	104.47 $\pm$ 50.99	(3)			
France	133.12 $\pm$ 61.04	Norway	104.01 $\pm$ 62.21	Portugal	49.43 $\pm$ 16.24	Canada	22.59 $\pm$ 38.57
Korea, Rep.	128.84 $\pm$ 76.56	Italy	102.57 $\pm$ 75.37	Singapore	45.32 $\pm$ 73.15	Finland	22.54 $\pm$ 17.42
Japan	126.96 $\pm$ 64.96	Taiwan	101.71 $\pm$ 85.34	Puerto Rico	41.86 $\pm$ 41.65	Slovenia	<b>18.70 <math>\pm</math> 17.11</b>
Sweden	114.28 $\pm$ 56.11	Spain	100.94 $\pm$ 73.13	Turkey	39.14 $\pm$ 45.65	U.S.	<b>15.90 <math>\pm</math> 25.02</b>
Jordan	109.95 $\pm$ 61.85	Hong Kong	100.58 $\pm$ 84.43	Taiwan	35.74 $\pm$ 57.75	Italy	<b>15.06 <math>\pm</math> 12.83</b>
Puerto Rico	108.42 $\pm$ 36.14	Hungary	100.05 $\pm$ 76.77	Greece	33.78 $\pm$ 24.88	Germany	<b>10.94 <math>\pm</math> 8.58</b>
(3)				Japan	30.42 $\pm$ 42.95	U.K.	<b>10.80 <math>\pm</math> 11.74</b>
Poland	99.69 $\pm$ 58.80	Russia	77.49 $\pm$ 52.54	Spain	27.25 $\pm$ 19.00	Belgium	10.68 $\pm$ 27.21
Greece	92.70 $\pm$ 69.36	Netherlands	59.52 $\pm$ 54.77	Russia	26.38 $\pm$ 20.41	Sweden	<b>10.20 <math>\pm</math> 10.98</b>
UK	90.46 $\pm$ 50.67	Austria	53.75 $\pm$ 50.77	Romania	26.36 $\pm$ 17.35	Hungary	<b>8.80 <math>\pm</math> 7.36</b>
Switzerland	88.40 $\pm$ 66.13	Turkey	51.37 $\pm$ 64.37	Ireland	24.35 $\pm$ 41.62	Switzerland	<b>8.56 <math>\pm</math> 12.02</b>
Germany	84.47 $\pm$ 61.80			France	24.34 $\pm$ 46.36	Netherlands	<b>7.77 <math>\pm</math> 13.00</b>
(4)				Norway	23.33 $\pm$ 15.62	Denmark	7.09 $\pm$ 6.02
Czech Rep.	48.36 $\pm$ 51.95	Canada	38.51 $\pm$ 46.15	Poland	23.15 $\pm$ 10.31	Austria	6.84 $\pm$ 5.76
Ireland	45.88 $\pm$ 50.55	US	36.81 $\pm$ 34.72				
Belgium	42.70 $\pm$ 56.02						

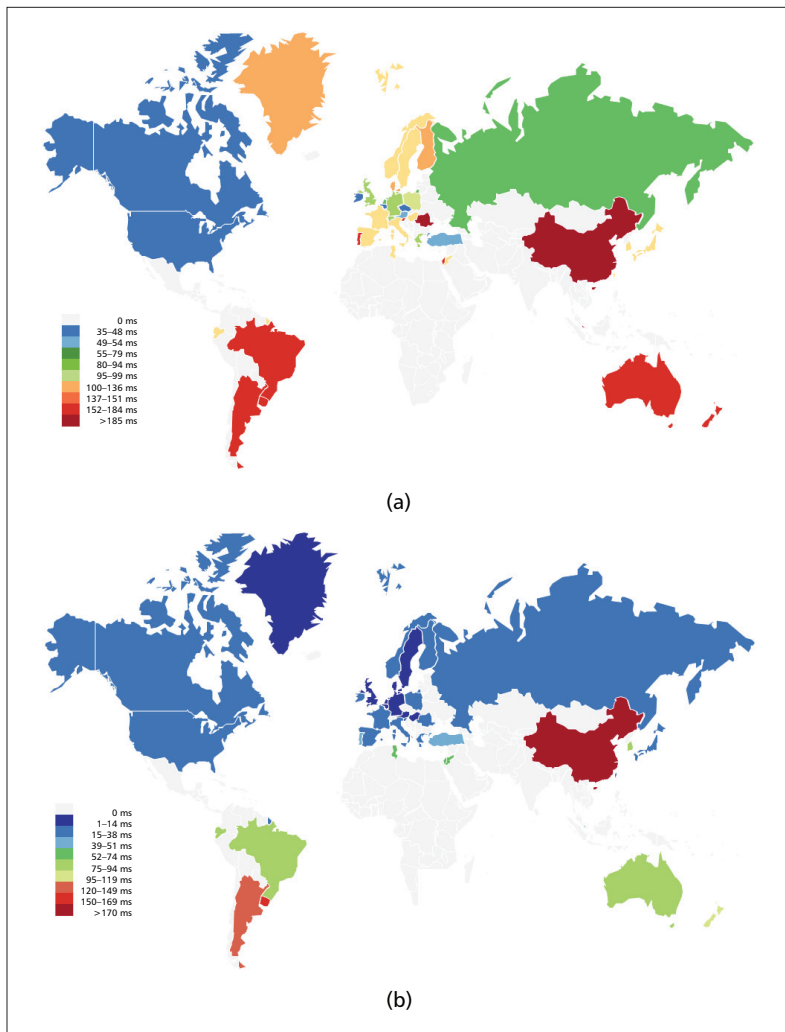
**Table 3.** Average delay (milliseconds)  $\pm$  standard deviation to access FB services from different countries for services located in a) FB servers; b) Akamai servers.

and at Akamai facilities (Fig. 2b). Overall, FB users need 113 ms in average to access native FB servers, but only 43 ms to reach Akamai nodes providing FB access. This means that accessing FB services in Akamai nodes reduces the delay 2.5 $\times$  the delay. Next, we provide a detailed analysis of the access delay performance per country.

#### ACCESS DELAY TO NATIVE FACEBOOK SERVERS

Based on the results of Table 3, we have defined four groups in terms of their access delay to FB servers, which are illustrated in different color range in Fig. 2 as well.

**The first group:** Refers to all those countries with an access delay longer than 150 ms



**Figure 2.** Average delay (milliseconds) to access FB services from different countries for services: a) located in FB servers; b) located in Akamai servers.

(red group in Fig. 2a). This group is formed by countries that are quite far from the United States (e.g., Australia, New Zealand), South American countries, and three countries we did not expect to find in this group (Portugal, Slovenia, and Israel) since their surrounding neighbors show a considerably lower delay.

**The second group:** Formed of those countries whose delay ranges between 100 and 150 ms (orange group). This group includes Northern European countries, Asian countries with deep penetration of high-speed access connections (e.g., Japan, South Korea, Hong Kong), countries from Central America, and Mediterranean countries including some important European ones such as France, Italy, and Spain.

**The third group:** Includes those countries with a delay greater than 50 ms but less than 100 ms (green group). This group is mainly formed by countries located in Central Europe plus Greece, Turkey, and the United Kingdom.

**The last group:** Contains those countries with access delay under 50 ms (blue group). This

includes the two countries hosting native FB servers, the United States and Ireland [1], and Canada due to its proximity and good connectivity with the United States. Surprisingly, this group also includes Belgium and the Czech Republic, which intuitively would have fit better in the third group.

#### ACCESS DELAY TO AKAMAI SERVERS

In the case of Akamai nodes we just define three groups for our discussion.

**The first group:** Formed by three countries that experience an average delay longer than 100 ms (red group in Fig. 2b). These countries are China, Argentina, and Uruguay. This happens because an important portion of the FB queries from these countries are redirected to remote Akamai nodes, which could be located, for instance, in the United States.

**The Second Group:** Consists of countries with an average access delay ranging between 50 and 100 ms (green group). This include far eastern countries like Australia, New Zealand, South Korea, and Hong Kong; two countries in South America, Brazil and Ecuador; and three countries from North Africa and the Middle East: Jordan, Tunisia, and Israel. As seen in the next section, the first six countries count on their own Akamai nodes, but a relevant portion of their demand is attended to by foreign Akamai servers. In addition, Jordan and Tunisia do not host any Akamai nodes, but are served by Akamai nodes located in Europe, which is relatively close. It is surprising that Australia (as a developed country) experiences quite bad performance in accessing FB services through Akamai nodes. To have better insight, we leveraged the FB ad planner<sup>7</sup> to retrieve the potential reach for ads in each country. We have found that Australia has a potential reach of 13 million FB users, while some of the countries in the third group, like Greece and Slovenia, which present 50 and 70 ms less average access delay, respectively. Another surprising case in this group is Brazil, a huge country with a population of more than 200 million and potential reach of an audience of 86 million for FB ads, and shows an average Akamai access delay around 78 ms.

**The third group:** Includes the countries with access delay below 50 ms (blue group). This group mainly includes developed countries from Europe, Asia (i.e., Japan and Singapore), and North America (United States and Canada). This is a good estimation of a short list of important countries for FB, where FB is interested in offering a better quality of service through Akamai nodes.

Furthermore, it is interesting to note that Akamai offers the best delay performance (i.e., below 10 ms) to small countries roughly located in Central Europe (Hungary, Switzerland, Netherlands, Denmark, Austria, and Czech Republic). This happens because these are very small countries (in size) that experience a very small delay due to the short distance to a large number of Akamai nodes located in Central Europe.

<sup>7</sup> <https://www.facebook.com/ads>

## AKAMAI NODES DISTRIBUTION TO PROVIDE ACCESS TO FACEBOOK SERVICES

This section provides a global picture of the deployment of Akamai nodes to serve FB services worldwide.

### LOCAL VS. EXTERNAL ACCESS

Figure 3 shows which portion of the queries for each country (i.e., ICMP echo) is managed by Akamai servers hosted in the same country than the source node(s) and which portion is served by Akamai nodes in foreign countries.

There are only two countries showing a higher portion of local access to Akamai servers compared to external access which are US and Singapore with 90 percent and 62 percent of the queries going to local Akamai servers. The case of Singapore might be unexpected, but as we will show later, Singapore has a high number of Akamai nodes (i.e., IPs). Close to Singapore performance, we find the case of Taiwan in which half of the queries are dealt with local servers and half by foreign servers.

We found that there are a limited number of countries that use local Akamai nodes to serve between 30 percent-40 percent of their queries. These are:

- The largest European countries by size (i.e., Germany, France and Spain) all of which have a large number of Akamai servers;
- Australia, another large country with high number of FB subscribers, that are located far from native FB servers and thus FB is motivated to use Akamai CDN to offer a good performance to Australian subscribers;
- Three European countries, Switzerland, Sweden and Romania, each particularly distributed geographically in the center, north and east of Europe respectively.

The Akamai infrastructures in Switzerland and Sweden bring them to have access delays to Akamai nodes in the order of 10ms. Finally, Romania has just six Akamai servers that service 35 percent of the queries generated in Romanian nodes.

Next, we found a large number of countries keeping between 7 percent and 30 percent of FB queries were responded locally, while most of them were serviced by foreign Akamai servers. Each of these countries have more or less Akamai nodes that allow keeping part of the queries locally, but their delay is mostly affected by how far those Akamai nodes are located from the major part of their queries.

Finally, there were 10 countries for which we could not identify any local Akamai server. Among these are five European countries (Belgium, Denmark, Hungary, Portugal and Slovenia), each with a population under 10M and close to countries with a significant deployment of Akamai nodes running FB services. The fact that these countries are experiencing a very good service by accessing Akamai nodes in nearby countries explains the low presence of Akamai servers.

This group of countries without Akamai

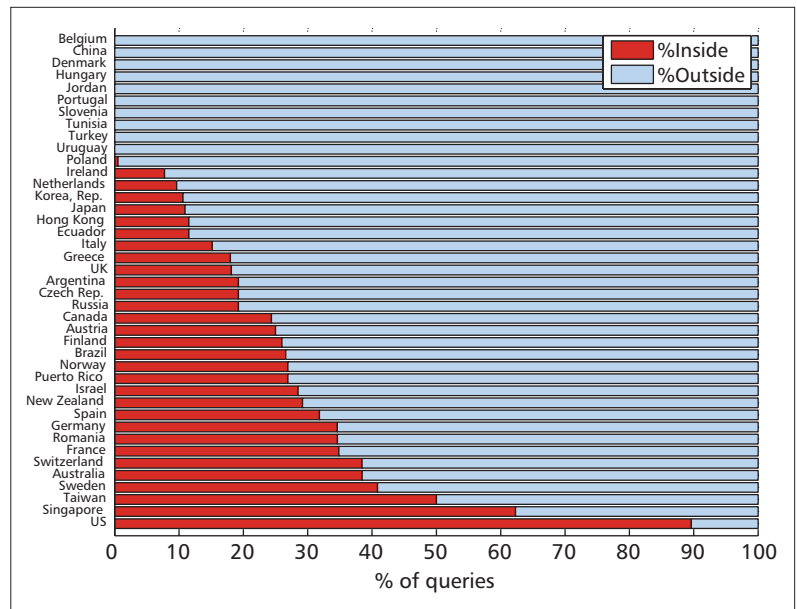


Figure 3. Portion of FB queries from each country served by local (% inside) and foreign (% outside) Akamai nodes.

servers also includes Turkey, which we found similar delay to some European countries like Greece or Portugal because the three PL nodes used for our experiments are located in the western part of Turkey (i.e., Istanbul and Izmir). Next, we discuss the case of Uruguay, a small South-American country surrounded by Argentina and Brazil that already contains some Akamai servers. Interestingly, the results in Table 4 show that the five PL nodes placed in Uruguay access Akamai servers located in Brazil as well as servers in Mexico and US that are far away, but never go to Argentina. Two small countries, Tunisia and Jordan, both of them are served by Akamai nodes located (mainly) in Europe. Finally, we find China which is currently blocking FB, and thus it does not make sense to deploy Akamai nodes to serve FB subscribers and they are served by Akamai nodes all over the world.

### COUNTRY COVERAGE BY AKAMAI SERVERS

Table 4 shows for each country which is hosting Akamai nodes, the overall number of IPs linked to Akamai nodes located in that country (column #IP), and the list of countries hosting nodes that access those IPs<sup>8</sup> (column Served To(#IP)). For each source-querying country we represent the overall number of IPs (between brackets) accessed in the destination country hosting Akamai nodes.

We found 35 countries that host Akamai nodes to provide FB access to the 41 countries represented by PL nodes. Among them, at the top of Table 4, we find 13 countries where Akamai nodes only serve local users. In the middle of the table we list four countries: Azerbaijan, Malaysia, Mexico and Panama, whose Akamai nodes only serve foreign countries. In fact, this behavior responds to the fact that we did not have any PL node located in those countries. Otherwise, we would very likely have observed that these Akamai nodes also serve local users.

<sup>8</sup> For simplicity during the discussion we use the number of IPs as the number of servers/nodes, even though we are aware that it is feasible that the same physical server could hold more than one IP (multiple network cards, virtualization, etc.)

Country	#IP	Served To (#IP)
Argentina	9	Argentina (9)
Canada	28	Canada (22)
Ecuador	3	Ecuador (3)
Greece	7	Greece (7)
HongKong	6	HongKong (6)
Israel	18	Israel (18)
Korea	7	Korea (7)
Poland	2	Poland (2)
Puerto Rico	8	Puerto Rico (8)
Romania	6	Romania (6)
Russia	7	Russia (7)
Spain	35	Spain (35)
Taiwan	9	Taiwan (9)
Azerbaijan	1	China (1)
Malaysia	21	HongKong (3), NewZealand (16), Singapore (2)
Mexico	4	Uruguay (4)
Panama	4	Canada (4)
Australia	15	Australia (10), Japan (4), Taiwan (1)
Austria	49	Austria (9), Greece (24), Hungary (26), Israel (2), Poland (37), Slovenia (27)
Brazil	26	Brazil (22), Uruguay (16)
Czech	11	Czech (6), Poland (7), Russia (4)
Finland	24	Finland (19), Norway (4), Russia (3), Sweden (12)
France	176	Belgium (20), Finland (4), France (60), Germany (4), Greece (1), Hungary (1), Ireland (7), Israel (3), Jordan (20), Poland (10), Singapore (3), Spain (41), Switzerland (5), Tunisia (4), Turkey (1), United Kingdom (48)
Germany	473	Australia (3), Austria (11), Belgium (30), China (6), Czech (13), Denmark (4), Finland (19), France (25), Germany (184), Greece (43), Hungary (11), Ireland (9), Israel (20), Italy (20), Jordan (5), Netherlands (12), Norway (5), Poland (31), Portugal (24), Romania (7), Russia (12), Slovenia (2), Spain (47), Sweden (11), Switzerland (86), Tunisia (9), Turkey (16), UnitedKingdom (14), United States (4)
Ireland	6	China (1), Ireland (5)
Italy	49	China (1), Greece (4), Hungary (2), Israel (1), Italy (20), Jordan (14), Switzerland (1), Tunisia (6), Turkey (1), United States (2)
Japan	36	China (17), Hong Kong (4), Japan (16), Korea (5)
Netherlands	39	Belgium (3), China (1), France (4), Ireland (14), Netherlands (6), Tunisia (6), UnitedKingdom (1), United States (3)
NewZealand	11	China (1), NewZealand (10)
Norway	8	Finland (2), Norway (5), Sweden (2)
Singapore	110	Argentina (3), Brazil (26), China (2), Ecuador (4), HongKong (13), Japan (1), Korea (3), New Zealand (1), Puerto Rico (15), Singapore (26), Taiwan (1), United States (30), Uruguay (12)
Sweden	77	Denmark (1), Finland (31), Ireland (7), Norway (6), Poland (7), Russia (10), Sweden (24), United Kingdom (19)
Switzerland	49	Australia (5), Poland (9), Sweden (8), Switzerland (33)
United Kingdom	246	Belgium (21), Denmark (4), France (19), Germany (22), Greece (34), Hungary (7), Ireland (17), Israel (34), Italy (2), Netherlands (23), Norway (13), Poland (29), Portugal (21), Romania (1), Spain (21), Sweden (5), Switzerland (11), Tunisia (9), Turkey (15), UnitedKingdom (45), UnitedStates (4)
UnitedStates	2505	Argentina (67), Australia (27), Austria (19), Belgium (39), Brazil (48), Canada (148), China (177), Czech (17), Denmark (13), Ecuador (17), Finland (16), France (69), Germany (117), Greece (32), HongKong (67), Hungary (26), Ireland (16), Israel (11), Japan (96), Jordan (1), Korea (66), Netherlands (18), NewZealand (21), Norway (17), Poland (47), Portugal (79), Puerto Rico (17), Romania (7), Singapore (21), Slovenia (14), Spain (24), Sweden (22), Switzerland (9), Taiwan (19), Tunisia (6), Turkey (13), UnitedKingdom (32), UnitedStates (1668), Uruguay (52)

**Table 4.** The first column shows the list of countries hosting Akamai nodes offering access to FB services. The second column shows the number of identified Akamai-related IPs in each country. The third column shows the list of countries including nodes querying Akamai IPs in the country referred to in the first column. The number between parentheses reflects the number of IPs accessed in the reference (first column) country.



Finally, at the bottom of the table, we find a major part of the countries (18 in total) with Akamai nodes that process queries from both local and foreign PL nodes. Next, we discuss the most interesting aspects for this group.

First, we observe that large countries with a relatively heavy weight in the geopolitical environment such as the United States, United Kingdom, France, Germany, and Italy have a high number of Akamai nodes (i.e., associated IPs) that serve a large number of countries. The four European countries mainly serve nodes from all over Europe, at a minor level nearby non-European countries like Israel, Jordan, Tunisia, and Turkey, and on a very small scale the United States and China. We also found a similar pattern in The Netherlands, although it has a lower Akamai presence. Furthermore, we discovered more Akamai nodes in the United States than in the rest of the countries together. These servers process queries from users located all over the world. This clearly has an impact on the delay for those countries that access Akamai nodes in the United States for a large portion of their queries, despite being far from the United States (e.g., Uruguay, Argentina, China, and Korea).

Next, we observe that Akamai nodes in Northern European countries (Norway, Finland, and Sweden) mainly respond to the demands of users located within those northern countries. A third observation is that Ireland and New Zealand should actually be located at the top of the table since they mostly attend to local FB demand, along with a few queries from China. Fourth, Akamai nodes located in small Central European countries such as Austria, Czech Republic, and Switzerland service FB demand mainly from local and nearby countries' users. We can find a similar pattern for Japan and Brazil, and additionally Australia, where the nodes mostly deal with internal demand for FB services but also receive some queries from nodes located in Japan and Taiwan. Finally, Singapore (the fourth country in terms of number of Akamai IPs) presents rarer results. On one hand, Akamai nodes in Singapore exhibit an expected behavior by serving users located in Asia. On the other hand, we discovered a very strange pattern in which Akamai nodes in Singapore attend quite a few nodes located all over America (including North and South America).

In summary, we can conclude that FB subscribers' queries are usually attended by Akamai nodes located either locally or in some nearby country. This provides a bounded access delay leading to the result presented above that indicates a delay  $2.5\times$  lower when an FB query is resolved by an Akamai node instead of a native FB server. However, we can still find some odd cases where source nodes access Akamai nodes located far away, which has a harmful impact on their access delay to FB services.

## RELATED WORK

We found a number of works related to our article that can be classified into two different categories: CDN infrastructure analysis and Facebook services analysis.

## CDN INFRASTRUCTURE ANALYSIS

There are some prior studies that analyzed different aspects of large CDNs like Akamai [6, 7] or the CDN used by Google to serve youtube videos [8]. In the latter study the authors aim at understanding from where videos are served, and how effective is this distribution. One of the main conclusions of this study is that round-trip time (RTT) is used to select the preferred data center to serve the video. The studies on Akamai CDN go from a general overview [6] to a more detailed analysis of Akamai's system components and architecture [7] in which the authors probe an Akamai network from 140 PlanetLab nodes during two months and characterize some aspects of Akamai architecture deployment such as server diversity, redirection dynamics, and latency. Finally, we found a study [9] in which the authors examined how CDNs are used to host and serve FB content from a network perspective. This work relies on a dataset including one month of HTTP traces collected in mid-2013 from the third generation (3G) mobile network of a large European ISP.

## FACEBOOK SERVICES ANALYSIS

There are also some research works that carried out different performance analyses on Facebook services. The authors in [1] look at the established connections when FB users login to the system. In particular, they identify different sections in the FB wall page of a user, and analyze how the information filling those sections is retrieved. An earlier work from 2010 [10] identified some performance degradation (delay, packet losses, etc.) for some users accessing FB from outside the United States. Finally, we have found another interesting study [11] which states that photo viewing is the most critical service for FB, and presents a detailed description on how FB photos are distributed to CDN Akamai servers. However, it does not perform a geographical analysis to understand how different regions of the world are being served as we do in our article.

## LESSONS LEARNED AND RECOMMENDATIONS

In this section we present the most important lessons extracted from our work and provide some recommendations that could improve the current delay experienced by users in some relevant countries.

**1:** Our study confirms the good work Akamai does for a large-scale web service such as Facebook. Our results show that FB is reducing delay  $2.5\times$  by using the Akamai nodes. This latency reduction is of great importance for Facebook or any other Internet service given the monetary implications associated the delay experience by end users [2].

**2:** At the time of our study, Akamai nodes were mostly responsible for serving heavy content mainly associated with photos and videos shared on Facebook. In contrast, Facebook native servers were in charge of processes like registration and authentication.

*We can conclude that FB subscriber queries are usually attended by Akamai nodes located either locally or in some nearby country. This provides a bounded access delay leading to the result that indicates a delay that is  $2.5\times$  lower when a FB query is resolved by an Akamai node instead of a native FB server.*

At the time of our study, Akamai nodes were mostly responsible of serving heavy content associated mainly to photos and videos shared in Facebook. In contrast, Facebook native servers were in charge of processes like the registration and authentication.

**3:** Akamai is very efficient ( $< 50$  ms delay) in serving Facebook content in Europe and North America, which is explained by two factors:

- Akamai is very well positioned there with a huge number of servers.
- A major part of the revenues obtained by FB from advertisement happens in Europe and North America; thus, it is very important to offer good quality of service to the subscribers in those locations.

**4:** There is some room for improving the current Facebook infrastructure in some countries like Australia and Brazil. These two countries have 13 and 86 million subscribers, respectively, according to the data reported by the FB ads planner, and experience much higher access delay (87 and 78 ms, respectively) than other countries with much lower numbers of subscribers like Slovenia. Therefore, we believe that Facebook should find a solution to improve the experience of Australian and Brazilian users by further exploiting Akamai nodes in those countries.

## CONCLUSIONS

This study presents a comprehensive measurement-based analysis of the FB network infrastructure with special emphasis on depicting how Akamai nodes replying to FB queries from subscribers are distributed throughout the world. In this context, we have analyzed the average access delay FB subscribers experience to access FB services delivered from native FB servers as well as Akamai servers depending on the country in which they are located. Moreover, we have thoroughly discussed the coverage offered by those Akamai nodes serving FB services.

## ACKNOWLEDGMENTS

This work is partially supported by the European Celtic-Plus project CONVINCe and eCOUSIN (EUF7-318398). This work is also funded by the Ministerio de Economía y Competitividad of SPAIN through the project Big-DatAAM (FIS2013-47532-C3-3-P).

## REFERENCES

- [1] W. Wongyai and L. Charoenwatana, "Examining the Network Traffic of Facebook Homepage Retrieval: An End User Perspective," *Proc. Int'l. Joint Conf. Computer Science and Software Engineering*, 2012.
- [2] A. Singla et al., "The Internet at the Speed of Light," *Proc. 13th ACM Wksp. Hot Topics in Networks*, 2014.
- [3] B. Chun et al., "Planetlab: An Overlay Testbed for Broad-Coverage Services," *ACM SIGCOMM Computer Commun. Review*, vol. 33, no. 3, 2003.
- [4] A. Orebaugh, G. Ramirez, and J. Beale, *Wireshark & Ethereal Network Protocol Analyzer Toolkit*, Syngress, 2006.
- [5] S. K. R. Siamwalla and R. Sharma, "Discovering Internet Topology," *Proc. IEEE INFOCOM*, 1999.
- [6] E. Nygren, R. K. Sitaraman, and J. Sun, "The Akamai Network: A Platform for High-Performance Internet Applications," *SIGOPS Op. Sys. Rev.*, 2010.
- [7] A.-J. Su et al., "Drafting Behind Akamai: Inferring Network Conditions Based on CDN Redirections," *IEEE/ACM Trans. Net.*, vol. 17, no. 6, 2009.
- [8] R. Torres et al., "Dissecting Video Server Selection Strategies in the Youtube CDN," *Proc. 31st Int'l. Conf. Distributed Computing Systems*, 2011.
- [9] P. Fiadino, A. D'Alconzo, and P. Casas, "Characterizing Web Services Provisioning via CDNS: The Case of Facebook," *Proc. Wireless Communications and Mobile Computing Conf.*, 2014.

- [10] M. P. Wittie et al., "Exploiting Locality of Interest in Online Social Networks," *Proc. Int'l Conf. Emerging Networking Experiments and Technologies*, 2010.
- [11] D. Beaver et al., "Finding a Needle in a Haystack: Facebook's Photo Storage," *Proc. 9th USENIX Conf. Operating Systems Design and Implementation*, 2010.

## BIOGRAPHIES

REZA FARAHABKSH [M] (reza.farahabksh@it-sudparis.eu) received his B.S. degree in computer engineering from Qazvin Azad University in 2006, his M.S. degree in 2008 from the University of Isfahan, Iran, in 2008, and his Ph.D from Institut-Mines Telecom, Telecom SudParis (CNRS Lab UMR5157) jointly with Paris VI (UPMC) in May 2015. He is now a postdoctoral researcher at Institut Mines-Telecom, Telecom SudParis since May 2015. His research interests are online social networks, P2P networks, Internet measurements, as well as Mobile IPv6 and IMS. He is a member of ACM.

ÁNGEL CUEVAS (acrumin@it.uc3m.es) received his B.Sc. in telecommunication engineering, his M.Sc. in telematics engineering, and his Ph.D. in telematics engineering from Universidad Carlos III de Madrid in 2006, 2007, and 2011, respectively. He is currently a tenure-track visiting professor at the Department of Telematic Engineering at Universidad Carlos III de Madrid. Prior to that he was a postdoctoral researcher at Institut Mines-Telecom, Telecom SudParis from March 2011 until January 2013. His research interests focus on online social networks, P2P networks, wireless sensor networks, and Internet measurements. He is a co-author of more than 30 papers in prestigious international journals and conferences such as IEEE/ACM Transactions on Networking, Elsevier Computer Networks, IEEE Network, IEEE Communications Magazine, ACM CONEXT, ACM MSWiM, IEEE ISCC, and IEEE ICC. He was co-recipient of the Best Paper Award at ACM MSWiM 2010.

ANTONIO M. ORTIZ received his Master's degree in advanced computer technologies (2008) and his Ph.D. in computer science (2011) from the University of Castilla-La Mancha, Spain. In 2012, he worked as a post-doctoral researcher at the Letterkenny Institute of Technology, Ireland, and at the Institut Mines-Telecom, Telecom SudParis, France during 2013 and 2014. He is currently a research engineer and project manager at Montimage, France. He actively participates in a number of European research projects and has contributed to the ITU-T and oneM2M standardization bodies. His topics of interest are focused on advanced monitoring and testing solutions applied to diverse aspects of network communications and cloud-based technologies, as well as on the Internet of Things, social networks, and artificial intelligence.

XIO HAN received her co-joint Ph.D degree from Telecom SudParis and Université Pierre et Marie Curie. She studied at the Automation School of Northwestern Polytechnical University, China, and received her B.Sc. and M.Sc. there in 2008 and 2011, respectively. Her research interests include social networks analysis and modeling, social-based applications, and social P2P networks.

NOËL CRESPI [M'07, SM'08] holds Master's degrees from the Universities of Orsay (Paris 11) and Kent (United Kingdom), a Diplôme d'Ingénieur from Telecom ParisTech, and a Ph.D and an Habilitation from Paris VI University (Paris-Sorbonne). From 1993 he worked at CLIP, Bouygues Telecom and then at Orange Labs in 1995. He took leading roles in the creation of new services with the successful conception and launch of Orange's prepaid service and in standardization. In 1999, he joined Nortel Networks as telephony program manager, architecting core network products for the EMEA region. He joined Institut Mines-Telecom in 2002 and is currently a professor and program director, leading the Service Architecture Lab. He coordinates the standardization activities for Institut Mines-Telecom at ITU-T, ETSI, and 3GPP. He is also an adjunct professor at KAIST, an affiliate professor at Concordia University, and on the four-person Scientific Advisory Board of FTW, Austria. He is the scientific director of the French-Korean laboratory ILLUMINE. His current research interests are in service architectures, services webification, social networks, and Internet of Things/Services. <http://noel-crespi.wp.tem-tsp.eu/>

**CALL FOR PAPERS**  
**IEEE COMMUNICATIONS MAGAZINE**  
**COMMUNICATIONS STANDARDS SUPPLEMENT**

**BACKGROUND**

Communications standards enable the global marketplace to offer interoperable products and services at affordable cost. Standards development organizations (SDOs) bring together stakeholders to develop consensus standards for use by a global industry. The importance of standards to the work and careers of communications practitioners has motivated the creation of a new publication on standards that meets the needs of a broad range of individuals, including industrial researchers, industry practitioners, business entrepreneurs, marketing managers, compliance/interoperability specialists, social scientists, regulators, intellectual property managers, and end users. This new publication will be incubated as a Communications Standards Supplement in *IEEE Communications Magazine*, which, if successful, will transition into a full-fledged new magazine. It is a platform for presenting and discussing standards-related topics in the areas of communications, networking, and related disciplines. Contributions are also encouraged from relevant disciplines of computer science, information systems, management, business studies, social sciences, economics, engineering, political science, public policy, sociology, and human factors/usability.

**SCOPE OF CONTRIBUTIONS**

Submissions are solicited on topics related to the areas of communications and networking standards and standardization research in at least the following topical areas:

Analysis of new topic areas for standardization, either enhancements to existing standards or in a new area. The standards activity may be just starting or nearing completion. For example, current topics of interest include:

- 5G radio access
- Wireless LAN
- SDN
- Ethernet
- Media codecs
- Cloud computing

Tutorials on, analysis of, and comparisons of IEEE and non-IEEE standards. For example, possible topics of interest include:

- Optical transport
- Radio access
- Power line carrier

The relationship between innovation and standardization, including, but not limited to:

- Patent policies, intellectual property rights, and antitrust law
- Examples and case studies of different kinds of innovation processes, analytical models of innovation, and new innovation methods

Technology governance aspects of standards focusing on both the socio-economic impact as well as the policies that guide them. These would include, but are not limited to:

- The national, regional, and global impacts of standards on industry, society, and economies
- The processes and organizations for creation and diffusion of standards, including the roles of organizations such as IEEE and IEEE-SA
- National and international policies and regulation for standards
- Standards and developing countries

The history of standardization, including, but not limited to:

- The cultures of different SDOs
- Standards education and its impact
- Corporate standards strategies
- The impact of open source on standards
- The impact of technology development and convergence on standards

Research-to-standards, including standards-oriented research, standards-related research, and research on standards

Compatibility and interoperability, including testing methodologies and certification to standards

Tools and services related to any or all aspects of the standardization life cycle

Proposals are also solicited for Feature Topic issues of the Communications Standards Supplement.

Articles should be submitted to the *IEEE Communications Magazine* submissions site at

**<http://mc.manuscriptcentral.com/commag-ieee>**

Select "Standards Supplement" from the drop-down menu of submission options.

## RADIO COMMUNICATIONS: COMPONENTS, SYSTEMS, AND NETWORKS



Amitabh Mishra



Tom Alexander

Computer networking in the 1970s and 1980s meant exactly that: interconnecting large timesharing systems to exchange information, principally email and files. As computing became more pervasive, however, it ultimately resulted in the personal data devices we carry around with us (also known as smartphones), and networking moved from a computer-centric focus to a people-centric focus. By some estimates, more than 3 billion people today, or about 40 percent of the world's population, are connected to the global Internet in one way or another. Clearly, the majority of these users must be networked via their personal devices; for most people the Internet represents a direct personal communications facility rather than anything to do with a “computer” in the traditional sense of the word.

In the last five years, however, the “Internet of Things” has started to gain momentum. This refers to the concept of endowing inanimate objects with Internet connectivity, hooking them up to a data center — typically in the cloud — and then harvesting data from them, directing commands to them, or both. The concept is not new; apparently a Carnegie Mellon University vending machine could display its status on the Internet as far back as 1982. However, the Internet of Things is now a commercial reality applied to several different areas, and it is now predicted that there could be 20 to 30 billion devices connected to the Internet by 2020. Current applications for the Internet of Things (customarily abbreviated as IoT) exist in transportation infrastructure management, energy management (the Smart Grid), home automation, and healthcare systems, and the technology is being applied to new areas every day.

It is relatively straightforward to expect that the IoT is actually a wireless IoT, as it would be patently absurd to run cables to these billions of devices: quite apart from being cost-prohibitive, many of these devices are not even stationary. Thus, radio communications will play a large role in facilitating the development of the IoT. However, the wireless links that serve well for IoT communications differ considerably from those that support human to Internet connectivity. The latter are optimized for very high bandwidth — up to 7 Gb/s for the recently ratified IEEE 802.11ac standard — and long-duration traffic, such as streaming video, sent to a relatively small number of sta-

tions. IoT links, on the other hand, must cope with thousands of stations, each of which may transfer short bursts of small packets with data rates in kilobits or even bits per second. On top of that, many IoT applications demand extremely low power consumption, which severely constrains the wireless radios and protocols that can be used.

This issue of the Radio Communications Series contains an article presenting a tutorial overview of key elements of the IEEE 802.11ah draft standard, which is specifically designed to support key use cases important to IoT applications. The article focuses on those aspects of IEEE 802.11ah that are pertinent to IoT systems, linking individual aspects of the medium access control (MAC) and physical (PHY) layers to specific IoT needs. For example, IoT devices such as sensors or meters can profitably trade off occupied channel bandwidth for range; 802.11ah takes advantage of this trade-off by shifting the PHY into the sub-1-GHz spectrum (e.g., 915 MHz in the United States), where the available bandwidth is low, but considerable range improvement accrues from the lower path loss.

The article particularly calls out the various power-saving aspects of the IEEE 802.11ah MAC. Power saving is critical for IoT devices that may not only need to run on battery power for years, but are also likely to have space constraints that limit the size of battery that can be integrated. Even more constrained are energy harvesting IoT devices that may not even have batteries at all, but have to subsist on tiny amounts of energy harvested from the surrounding environment. Such devices necessarily must spend virtually all their time sleeping, and only wake up at long intervals to exchange data in short but highly efficient bursts. This forms a particular challenge for IEEE 802.11 carrier sense multiple access (CSMA-CA) MAC protocols, and the article capably illustrates the new 802.11ah MAC functions that were introduced to deal with the challenge.

We would like to continue to place similar articles before you, in future issues, to cover topics of current interest. However, for this to happen we need contributions from our community of authors! We therefore encourage you to submit articles discussing emerging trends in wireless communications. We also thank our reviewers for their time and attention to helping us create a quality Series.

# IEEE 802.11ah: Sub-1-GHz License-Exempt Operation for the Internet of Things

Minyoung Park

## ABSTRACT

IEEE 802.11ah Task Group has been developing an amendment to the 802.11 standard to define sub-1-GHz license-exempt operation to support sensors and Internet of Things applications. This article presents an overview of major physical layer and MAC layer features of 802.11ah.

## INTRODUCTION

The IEEE 802.11 standard [1] is one of the most widely adopted wireless connectivity technologies for digital devices such as laptops, smartphones, tablets, digital TVs, DVD/Blu-ray players, gaming consoles, and portable music players. Since the introduction of IEEE 802.11b in 1999 with data rates up to 11 Mb/s, IEEE 802.11 has been evolving to support more and more use cases that require higher throughput. In 2003, 802.11g was introduced with the adoption of an orthogonal frequency-division multiplexing (OFDM) modulation technique, which increased the data rate to 54 Mb/s. In 2009, 802.11n was published, which supports data rates up to 600 Mb/s utilizing multiple-input multiple-output (MIMO) and doubling the channel width to 40 MHz. The latest amendment, 802.11ac, was published in 2013, which supports data rates up to 6.9 Gb/s in the 5 GHz band. The enhancement is achieved by increasing the channel widths to 80 and 160 MHz, and the maximum number of spatial streams to eight.

Most of the enhancements in the amendments have been optimized for a small number of high data rate devices such as laptops, tablets, and smartphones, but not for a large number of low data rate devices such as sensors and Internet of Things (IoT) devices. As an example, the use cases and requirements of 802.11ac are summarized in Table 1. The main usages of 802.11ac are wireless display, distribution of HDTV, and rapid upload/download of large files, which all require high data rates for streaming videos or large file transfers among small numbers of devices, mostly in indoor environments.

## 802.11AH USE CASES AND IOT APPLICATIONS

In 2010, the IEEE 802.11ah Task Group was formed in the IEEE 802.11 Working Group to define an amendment to the 802.11 standard that operates in the sub-1-GHz license-exempt frequency spectrum to support the following three use cases: sensors and meters, backhaul sensor and meter data, and extended range Wi-Fi [2]. Among these use cases, the sensors and meters use case includes the following sub use cases, which all fall into IoT use cases: smart grid, environmental/agricultural monitoring, industrial process sensors, healthcare, home/building automation, and home sensors. As summarized in Table 1, these IoT use cases typically support a much larger number of devices per access point (AP) for both indoor and outdoor environments with a much longer transmission range but at a much lower data rate than the use cases and requirements of 802.11ac.

## SUB-1-GHz LICENSE-EXEMPT OPERATION

Although sub-1-GHz bands have more limited frequency spectrum available (e.g., there is only a total of 26 MHz spectrum available in the 915 MHz industrial, scientific, and medical [ISM] band in the United States) than the 2.4 and 5 GHz ISM bands, as shown in Fig. 1, it is sufficient for low data rate applications such as IoT applications because such applications typically transmit small amounts of data infrequently. Moreover, since the 915 MHz ISM band (902–928 MHz) has 8.5 dB less free space propagation loss than the 2.4 GHz ISM band, this can be used to enhance the link budget between devices, and either enable long-range transmission for outdoor metering applications or reduce energy consumption of a device by lowering transmit power and supporting indoor sensor applications.

This article gives an overview of the major features adopted in the 802.11ah specification currently being developed in the IEEE 802.11 Working Group with details of how the features solve the challenges to support IoT use cases. In the following sections, the availability of sub-1-GHz spectrum, as well as channelization in key geographical locations, are described and then

Minyoung Park is with Intel Corporation.

	802.11ac use cases and requirements [3, 4]	802.11ah IoT use case and requirements [2]
Use cases	1) Wireless display 2) Distribution of HDTV 3) Rapid upload/download	Sensors and meters
Data rate requirement	20 Mb/s–3 Gb/s	100 kb/s
Single frame packet size	Large (e.g., 1500 bytes)	Small (e.g., few 100 bytes)
Traffic type	Video streaming/large file transfer	Periodic packet transmission every few to tens minutes
Distance between devices	5–60 m	Up to 1 km
Number of stations	3–20	Up to 6000
Location	Mostly indoor	Indoor and outdoor

**Table 1.** Comparison of 802.11ac and 802.11ah.

the major physical (PHY) layer and medium access control (MAC) layer features of 802.11ah are explained.

## 802.11AH PHY FEATURES

### SUB-1-GHZ FREQUENCY SPECTRUM AND CHANNELIZATION

The availabilities of frequency spectrum and channelization for key geographies are shown in Fig. 1. 802.11ah supports five frequency channel bandwidths ranging from 1 to 16 MHz. As shown in Fig. 1, there can be 26 1 MHz channels in the United States and fewer channels for the wider channel bandwidths.

Different countries have different rules for the allocated frequency bands in terms of maximum transmit power and channel width. For example, in the United States, an 802.11ah station can operate in the 915 MHz ISM band with the maximum transmit power of 1 W without any restrictions on channel bandwidth. On the other hand, Japan can use 915.9–929.7 MHz band for 802.11ah operation, but the signal frequency bandwidth is limited to 1 MHz and the maximum transmit power is limited to 250 mW (920.5–923.5 MHz). Please refer to [5, Annexes D and E] for more details.

### BASIC PHY DESIGNS OF 802.11AH

The basic PHY parameters of 802.11ah are summarized in Table 2. In order to support the sensor and meter use case as well as the extended range WiFi use case [2], 802.11ah supports data rates ranging from 150 kb/s to 347 Mb/s. The basic PHY design of 802.11ah is mostly inherited from 802.11ac. The 802.11ac signal waveforms with 20 to 160 MHz channel bandwidths are scaled down by 10 times to 802.11ah signal waveforms with 2 to 16 MHz maintaining the same number of subcarriers. Therefore, the tone spac-

ing between adjacent subcarriers is now 31.25 kHz for all bandwidths. This makes the inverse/discrete Fourier transform (IDFT/DFT) period of 802.11ah equal to 32  $\mu$ s, which is 10 times longer than that of 802.11ac. The OFDM symbol period is now 40  $\mu$ s with 8  $\mu$ s guard interval.

The IEEE 802.15.4 standard supports sub-1-GHz operation at the following frequency bands: 779–787 MHz, 868–868.6 MHz, 902–928 MHz, and 950–956 MHz with the following PHY data rates: 20, 40, 100, and 250 kb/s [6].

### TRANSMISSION RANGE ENHANCEMENT

In order to meet a transmission range of 1 km at a minimum data rate of at least 100 kb/s for outdoor IoT applications [2], the transmission range of 802.11ah operating in the 900 MHz band is significantly enhanced compared to that of 802.11n operating in the 2.4 GHz band. The transmission range is enhanced by increased link budget between two devices based on the following design choices.

**Lower operation frequency spectrum:** 802.11ah operates in sub-1-GHz bands, which improves link budget significantly compared to 2.4 GHz operation. For example, a client device operating in the 2.4 GHz band experiences 8.5 dB more free space path loss than one operating in the 900 MHz band.

**Narrower channel bandwidth:** 802.11ah transmits a signal in a 10 times narrower channel bandwidth than 802.11n. This reduces its noise bandwidth by 10 times and thus can increase the signal-to-noise ratio (SNR) at a receiver by 10 dB.

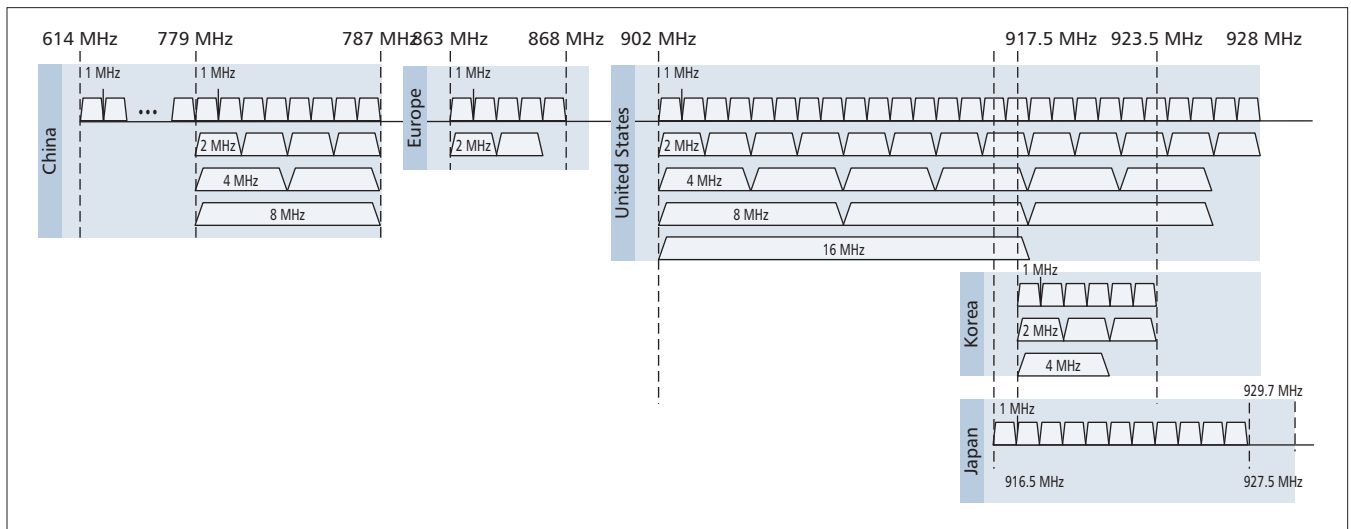
**1 MHz channel bandwidth support:** In addition to 10 times narrower channel bandwidths, 802.11ah supports 1 MHz channel bandwidth, which can increase the SNR by 3 dB compared to using the 2 MHz channel bandwidth.

**Robust coding scheme:** 802.11ah supports the repetition coding scheme for 1 MHz channel bandwidth when binary phase shift keying (BPSK) modulation is used with 1/2 coding rate, and this adds another 3 dB to the SNR.

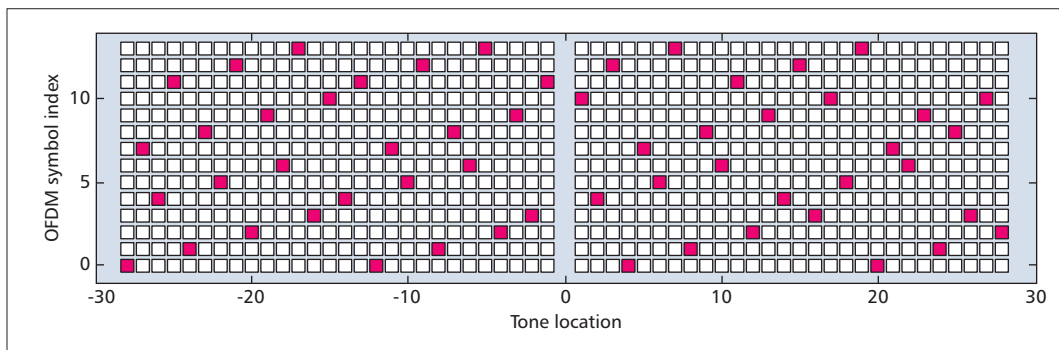
Table 3 summarizes the link budget enhancement of 900 MHz 802.11ah over 2.4 GHz 802.11n. 802.11ah operating in the 900 MHz band at 150 kb/s is expected to have approximately 24.5 dB better link budget than 802.11n operating in the 2.4 GHz band at 6.5 Mb/s.

### LOW-POWER AND LOW-COST SUPPORT FOR INDOOR SENSORS

As shown in Table 3, the link budget can be enhanced by more than 24 dB compared to an 802.11n device operating in the 2.4 GHz band. Instead of increasing the transmission range, this enhancement can be used to lower transmit power of a sensor device. The transmit power can easily be reduced to 0 dBm considering that the nominal transmit power of a typical 802.11 device is around 15–17 dBm. This can reduce the transmit energy consumption and also lower the cost of an 802.11ah radio of a small sensor device.



**Figure 1.** Frequency spectrum availabilities in sub-1-GHz and channelization in key geographies.



**Figure 2.** Traveling pilot positions for a signal transmitted in 2 MHz channel width and the number of space-time stream (STS) is 1 (red square = pilot tone, white square = data tone).

### SELECTIVE SUBCHANNEL TRANSMISSION

Although the 1 or 2 MHz channel bandwidth can enhance transmission range or reduce transmit power for IoT and sensor applications by reducing noise bandwidth, such a narrow bandwidth signal becomes more susceptible to flat fading than a 20 MHz 802.11ac signal. Considering that the 1 MHz channel width is approximately just 3 tones of the 20 MHz 802.11ac signal, the 1 MHz signal can easily be in a deep fade for an indoor environment with a short delay spread.

In 802.11ah, this problem is mitigated by selective subchannel transmission, which selects the best subchannel from a wider channel bandwidth. For example, the best 1 MHz subchannel can be selected from a 4 MHz channel bandwidth. This can improve signal power by approximately 7 dB for an indoor channel model with 50 ns root mean square (rms) delay spread as shown in [7].

### TRAVELING PILOTS FOR THE OUTDOOR ENVIRONMENT

Since 802.11ah is designed for outdoor IoT usage as well as indoor usage, it has to cope with a high Doppler case where a signal is reflected from a moving car. This means that channel estimation has to be updated throughout a packet.

The previous 802.11 amendments such as 802.11n and 802.11ac use the long training field in the preamble of a packet for channel estimation and rely on the initial channel estimation throughout the packet [8]. Although a 20 MHz OFDM symbol has 4 pilot tones out of 64 subcarriers, their locations are fixed and thus cannot cover the other 52 data tone locations for the channel estimation throughout a packet for a large delay spread in outdoor usage.

In order to update channel estimation throughout a packet for all 52 data tone locations, 802.11ah supports a traveling pilot scheme that shifts pilot tones every OFDM symbol such that the traveling pilot tones over multiple OFDM symbols can cover all data tone locations. Figure 2 illustrates the traveling pilot positions for a 2 MHz channel width signal when the number of space-time streams (NSTS) is 1. In this case, the traveling pilot pattern repeats every 14 OFDM symbols. The details of other traveling patterns are described in [5, 8].

### 802.11AH MAC FEATURES

802.11ah MAC features are mainly designed to enhance the efficiency of MAC protocols and MAC frame formats to reduce energy consumption of a client device, and to support a large

number of clients for sensor and IoT applications for both indoor and outdoor environments.

### MITIGATING CONTENTION AND REDUCING CHANNEL ACCESS DELAY

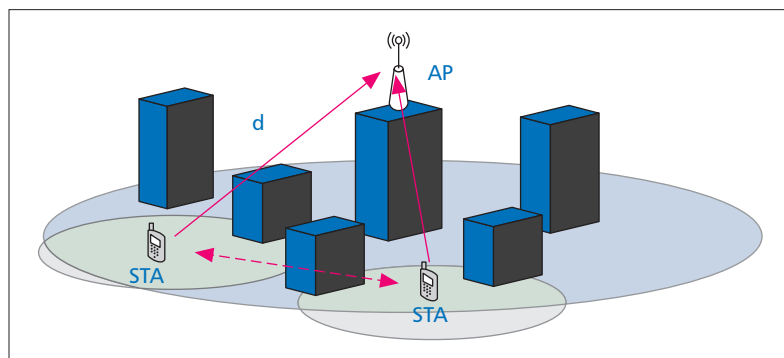
Supporting a transmission range of up to 1 km introduces a new challenge for an 802.11 station. As shown in Fig. 3, in a typical outdoor

PHY parameters	Supported values
Channel bandwidths	1 MHz, 2 MHz, 4 MHz, 8 MHz, 16 MHz
Modulation schemes	BPSK, QPSK, 16-QAM, 64-QAM, 256-QAM
Code rates	1/2 with 2 times repetition, 1/2, 2/3, 3/4, 5/6 in either convolutional or low-density parity check (LDPC)
Maximum number of spatial streams	Four spatial streams
Data rates	150 kb/s (1 MHz channel bandwidth, 1 spatial stream, BPSK, 1/2 coding rate, repetition coding) to 347 Mb/s (16 MHz channel bandwidth, 4 spatial streams, 256 QAM, 5/6 coding rate)

**Table 2.** 802.11ah PHY parameters.

Parameters	Link budget enhancements of 900 MHz 802.11ah over 2.4 GHz 802.11n
Free space path loss	+8.5 dB
Noise bandwidth	+10 dB
Sub-total link budget gain	+18.5 dB
1 MHz channel width	+3 dB
Repetition coding	+3 dB
Total link budget gain	+24.5 dB

**Table 3.** Link budget parameters between 900 MHz and 2.4 GHz.



**Figure 3.** Illustration of an outdoor environment where an AP is installed at a rooftop and client stations are located near the ground. The path loss between STAs is much higher than the path loss between the AP and a STA.

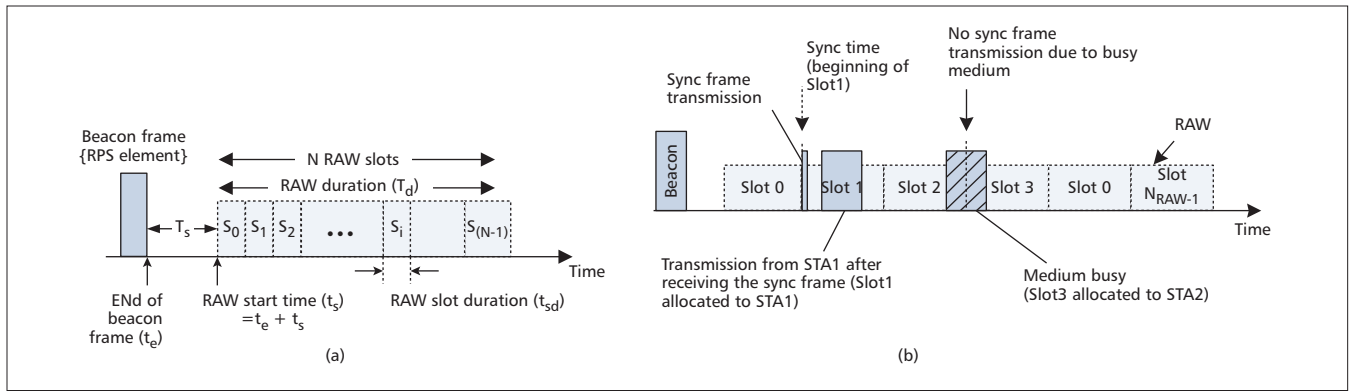
environment, an access point (AP) is installed at a much higher elevation, such as a rooftop, than client stations, which are mostly located near the ground. This creates much higher path loss between client stations than that between a client station and the AP because there are more building structures obstructing signal propagation between client stations than between a client station and the AP [9]. This makes the hidden node problem worse because it becomes harder for a station to hear transmissions from other stations to the AP, and this causes more collisions among transmissions from the stations to the AP. More collisions increase energy consumption of a station because the station needs to stay active longer for retransmissions to successfully deliver a data packet to the AP.

**Restricted Access Window:** 802.11ah mitigates the hidden node problem by restricting the time at which a station can start to contend for the medium so that packet transmissions from stations do not overlap with each other. 802.11ah defines a time window called a restricted access window (RAW) during which only a group of stations that are associated with an AP are allowed to access the medium. The structure of a RAW is illustrated in Fig. 4a. A RAW is divided into RAW slots, and each RAW slot is typically allocated to one station. The maximum number of slots in a RAW is 64. A RAW slot may also be allocated to more than one station to achieve statistical multiplexing among the stations in the RAW slot. A station in the group is allowed to contend for the medium at the beginning of the allocated RAW slot. Although the RAW slot is allocated to the station by the AP, the station still needs to perform the contention-based channel access because stations associated with other APs do not have the information about the RAW slot allocation and may still contend for the medium. The RAW parameter set (RPS) element contains information of one or more RAWs such as the RAW start time, RAW slot duration, and number of slots. The details of the RAW parameter set element are defined in [5]. The simulation results in [10] show that collisions among hidden nodes can be mitigated significantly using the RAW scheme.

**Synchronization Frame:** In 802.11, a station is allowed to start to contend for the medium only if it has valid information of the medium so that it does not interrupt packet transmissions and receptions of other stations in the network. The station obtains the information of the medium from the duration field of a packet that was received correctly; otherwise, it has to wait for a time duration called ProbeDelay to expire to avoid collision with a packet transmission of other station [1]. For example, if a station just woke up from sleep and has not received a packet, it has to either wait for a packet or wait until the ProbeDelay timer expires, which adds additional energy consumption to the station. The ProbeDelay problem is described in [10] in more detail.

In 802.11ah, the medium access delay due to





**Figure 4.** Illustration of a) a RAW; b) synchronization frame operation.

the ProbeDelay is reduced with help from an AP. As illustrated in Fig. 4b, the AP transmits a synchronization (sync) frame at the beginning of a RAW slot when the medium is idle so that a station can obtain the information of the medium from the sync frame and access the medium right after the end of the sync frame reception instead of waiting for the ProbeDelay time to expire. If the medium is busy at the AP, the sync frame is not transmitted, and the station waits for the ProbeDelay time to expire to prevent possible collision at the AP. The simulation results in [10] show that the sync frame technique can increase the battery life of a station by up to 31 percent when the ProbeDelay time is set to 5 ms.

#### EFFICIENT BIDIRECTIONAL PACKET EXCHANGES

802.11ah introduced the target wake time and bidirectional transmission opportunity (TXOP; BDT) schemes to enhance bidirectional packet exchanges between two stations by eliminating overheads between uplink and downlink transmissions, thereby minimizing energy consumption of sensor devices.

**Target Wake Time:** An AP buffers data destined for a station while the station is in a sleep state. The station periodically wakes up at beacon transmission times and receives a beacon to see if there is any buffered data at the AP based on the information in the traffic indication map (TIM) element contained in the received beacon. If the TIM element indicates that there is buffered data for the station, the station first sends a PS-Poll frame to the AP to indicate that the station is awake and is ready to receive the buffered data. The AP, however, needs processing time to find the buffered data for the station from its memory, and then has to contend for the medium and transmit the data to the station. This indefinite latency makes the station consume energy waiting for the buffered data.

Target wake time (TWT) addresses this problem by having an AP and a station schedule a future wake-up time (i.e., a TWT) of the station so that the AP knows when the station will be awake. The AP fetches buffered data from its memory before the TWT so that when a PS-Poll frame or a trigger frame is received from the station, the AP can transmit the buffered data with-

out the processing time and the medium access latency. The time information of the next TWT can be delivered explicitly to the station during the frame exchange in the current TWT or calculated implicitly by adding a fixed time interval to the current TWT.

**Bidirectional TXOP:** For a sensor device with a small battery, it is critical to minimize energy consumption of the device. A sensor device can minimize its energy consumption by coalescing transmit and receive activities in a single TXOP to increase its time in a sleep state.

BDT allows an AP and a station to exchange one or more uplink and downlink packets separated by a short inter frame space (SIFS) in a TXOP duration. In the BDT procedure, a station uses the More Data bit in the SIGNAL field of the PHY preamble to signal whether the station has more data to transmit following the current packet transmission.

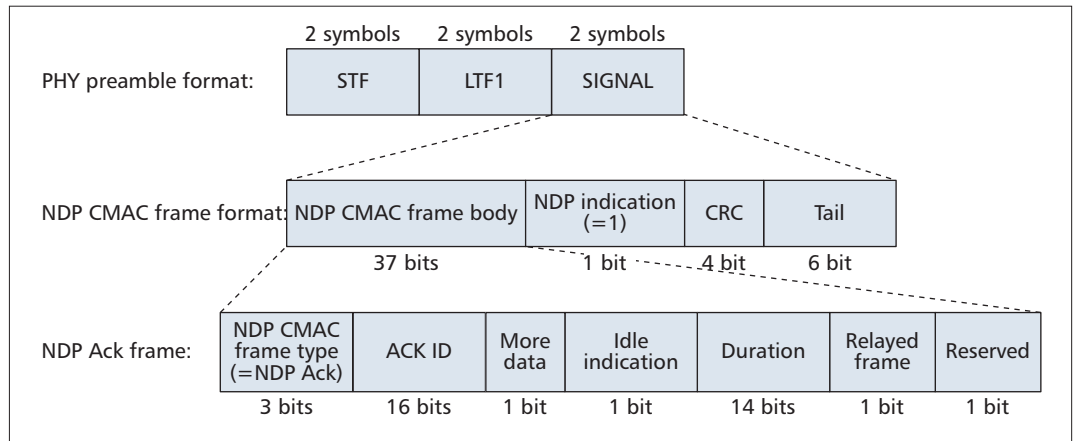
802.11ah also added the Response Indication field in the SIGNAL field of the PHY preamble of a packet to better protect a response packet. The Response Indication field in the transmitted packet indicates one of the four different lengths of a response packet that will follow the transmitted packet so that the third party stations near the transmitter can defer until the end of the response packet. Since the Response Indication field is in the SIGNAL field, it uses the most robust modulation and coding scheme, and the third party stations can defer correctly as long as they decode the SIGNAL field correctly.

#### NULL DATA PACKET CARRYING MAC FRAME

In 802.11, a Control frame such as the ACK frame, which is 14 bytes long, has not been considered to be a large overhead for email or file transfer that transmits large packets (e.g., 1500 bytes). For sensor and IoT applications, however, the ACK frame can be a large overhead considering that those applications typically transmit much shorter data packets (e.g., 100 bytes [2]).

Control frames are transmitted with the lowest modulation and coding scheme, and for the 802.11ah 2 MHz channel width, it takes 440  $\mu$ s to transmit an ACK packet (240  $\mu$ s for the PHY

The new Short MAC header of a Short QoS Data frame is reduced to 12 bytes. In order to differentiate the new short MAC header format from the legacy MAC frame format, the protocol version field of the frame control field in the short MAC header is increased from 0 to 1 for the first time in 802.11 history.



**Figure 5.** Illustration of the NDP CMAC frame format and the NDP Ack frame format in 2 MHz, 4 MHz, 8 MHz, and 16 MHz.

preamble and 200  $\mu$ s for the ACK frame) at the 650 kb/s data rate. To transmit a 100-byte data packet at the same data rate, it takes 1480  $\mu$ s (1240  $\mu$ s for the data frame). In this case, the ratio between the ACK packet transmission time and the ACK and data packets transmission time is approximately 23 percent, which shows a large overhead of ACK frame transmission.

802.11ah mitigates this problem by defining a null data packet (NDP) carrying MAC (CMAC) frame format that consists of only the PHY preamble and no data field, as shown in Fig. 5. The SIGNAL field of the NDP CMAC frame is redefined to contain all the necessary information previously contained in a Control MAC frame. The length of an NDP CMAC frame is now fixed to the length of the PHY preamble, which is 240  $\mu$ s (6 OFDM symbols) for a signal transmitted in a channel bandwidth larger than or equal to 2 MHz. For example, 802.11ah uses the NDP Ack frame shown in Fig. 5 instead of the ACK frame as a response to a data frame. The ratio between the NDP Ack frame transmission time and the NDP Ack frame and the data packet transmission time is now reduced to 14 percent. The details of the NDP CMAC frames are described in [5].

### SHORT MAC FRAME

In 802.11n, the MAC header of a quality of service (QoS) Data frame can be 30 bytes long. For an IoT application that transmits a small packet (e.g., a temperature sensor with 50 bytes of data [2]) infrequently, the MAC header can be a big overhead.

802.11ah mitigates this problem by defining a new Short MAC frame format. The new Short MAC header of a Short QoS Data frame is reduced to 12 bytes. In order to differentiate the new short MAC header format from the legacy MAC frame format, the protocol version field of the frame control field in the short MAC header is increased from 0 to 1 for the first time in 802.11 history.

### SUPPORT OF A LARGE NUMBER OF STATIONS

One of the use cases of 802.11ah is a smart grid use case where an AP has to support as many as 6000 stations within a 1 km<sup>2</sup> area [2, 11]. In

order to support this use case, in 802.11ah, the number of stations that an AP can support is increased from 2007 to 8191. During an association process, a station is assigned with an association identifier chosen from 1 to 8191 by an AP. The AP maintains a traffic indication virtual bitmap where the  $n$ th bit is mapped to a station that has the association identifier value of  $n$ . When there is buffered data for the station with the association identifier value of  $n$ , the AP sets the  $n$ th bit of the bitmap to 1 and signals the information in the Partial Virtual Bitmap field of the TIM element in a beacon. Signaling such a large traffic indication virtual bitmap becomes challenging because in some cases the length of the Partial Virtual Bitmap field can be a couple of hundred bytes long.

In order to minimize the size of the TIM element, the traffic indication virtual bitmap is structured in a hierarchical fashion, and 802.11ah defines four encoding modes to compress the traffic indication virtual bitmap [5]. The simulation results in [12] show that the bitmap size can be compressed by 30–98 percent with the encoding schemes compared to the scheme defined in the baseline 802.11 standard [1].

### INCREASED SLEEP TIME

IEEE 802.11ah is designed to support a sensor device that sleeps for a very long period of time (e.g., a couple of days), waking up infrequently to transmit or receive a short data packet and going back to sleep. In the baseline 802.11 standard [1], however, a station cannot sleep for more than approximately 18 hours because the Max Idle Period field that contains the parameter BSSMaxIdlePeriod is an unsigned 16-bit value, and this parameter determines how long a station can stay idle (in seconds) before the station is disassociated from the AP. If the station wants to stay associated with the AP, it has to transmit a packet at least every BSSMaxIdlePeriod time.

In 802.11ah, the Max Idle Period field is redefined such that the two most significant bits of the Max Idle Period field is used as a scaling factor. When a scaling factor of 10,000 is used, a station can sleep for approximately 5.2 years without being disassociated from the AP.

## 802.11AH DEVELOPMENT TIMELINE

In 2010, the IEEE 802.11ah Task Group was formed. The task group developed the first 802.11ah draft amendment, D1.0, in October 2013. After going through four iterations of the letter ballot and comment resolution process, the task group developed the 802.11ah draft amendment, D5.0, in April 2015. The final 802.11ah amendment is expected to be published in 2016 [13]. The WiFi Alliance has also formed a task group named Extended Range ah to study requirements for an interoperability certification program [14].

## CONCLUSIONS

In this article, the major PHY and MAC features of IEEE 802.11ah are presented. 802.11ah is designed to support a wide range of applications in sub-1-GHz frequency spectrum: from sensors and Internet of Things applications to extended WiFi applications by providing data rates from 150 kb/s to 347 Mb/s. 802.11ah is designed to support an outdoor application that needs a transmission range up to 1 km at 150 kb/s. These are achieved by using the PHY parameters listed in Table 2. For an indoor application, the long-range capability can be used for energy-efficient communications for sensors by using a very low transmit power. As summarized in Table 4, 802.11ah provides energy-efficient MAC protocols and MAC frame formats that are optimized for sensors and IoT applications, which need to support a large number of stations and infrequent small packet transmissions.

## REFERENCES

- [1] IEEE Std 802.11-2012, "IEEE Standard for Information Technology Telecommunications and Information Exchange between Systems — Local and Metropolitan Area Networks Specific Requirements Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications."
- [2] R. de Vegt, "Potential Compromise for 802.11ah Use Case Document," IEEE 802.11-11/0457r0, Mar. 2011.
- [3] R. de Vegt, "802.11ac Usage Models Document," IEEE 802.11-09/0161r2, Jan. 2009.
- [4] P. Loc and M. Cheong, "TGac Functional Requirements and Evaluation Methodology," IEEE 802.11-09/451r16, Jan. 2011.

Functions for sensors and IoT	MAC features
Support for a large outdoor IoT network	RAW, synchronization frame, hierarchical TIM
Support for energy-efficient communications for sensors	TWT, BDT, NDP CMAP frame, short MAC frame, increased sleep time

**Table 4.** Summary of 802.11ah MAC features.

- [5] IEEE 802.11ah/D5.0, "Draft for Information Technology — Telecommunications and Information Exchange between Systems — Local and Metropolitan Area Networks — Specific Requirements — Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications — Amendment 6: Sub 1 GHz License Exempt Operation."
- [6] IEEE Std 802.15.4-2011, "IEEE Standard for Local and Metropolitan Area Networks Part 15.4: Low-Rate Wireless Personal Area Networks (LR-WPANs)."
- [7] M. Fischer, *et al.*, "Frequency Selective Transmission," IEEE 11-12/1338r0, Nov. 2012.
- [8] R. Porat *et al.*, "Traveling Pilots," IEEE 802.11-12/1322r0, Nov. 2012.
- [9] R. Porat, S. Yong, and K. Doppler, "TGah Channel Model," IEEE 802.11-11/0968r3, Nov. 2011.
- [10] M. Park, "IEEE 802.11ah: Energy Efficient MAC Protocols for Long Range Wireless LAN," *Proc. IEEE ICC*, 10–14 June 2014, pp. 2388–93.
- [11] S. Kim *et al.*, "Association ID management for TGah," IEEE 802.11-11/88r1, Jan. 2011.
- [12] M. Park, "TGah Efficient TIM Encoding," IEEE 802.11-12/388r2, May 2012.
- [13] IEEE, "IEEE 802.11 Working Group Project Timelines," [http://www.ieee802.org/11/Reports/802.11\\_Timelines.htm](http://www.ieee802.org/11/Reports/802.11_Timelines.htm), accessed June 22, 2015.
- [14] WiFi Alliance, "Current Work Areas," <http://www.wifi.org/who-we-are/current-work-areas>, accessed: June 22, 2015.

## BIOGRAPHY:

MINYOUNG PARK received B.S. and M.S. degrees from Yonsei University, Seoul, Korea, in 1993 and 1995, and his Ph.D. degree from the University of Texas at Austin in 2005, all in electrical engineering. From 1995 to 2001, he was a senior research engineer at LG Electronics Inc. He is currently a research scientist at Intel Labs. He has participated in various standards developments in IEEE 802.11 and the WiFi Alliance, and served as the Technical Editor of the IEEE 802.11ah Task Group. His research interests are focused on the physical and MAC layers of wireless communications systems.

# ADVERTISERS' INDEX

COMPANY	PAGE
AR Modular RF .....	15
Fraunhofer.....	3
ICC 2016.....	47
IEEE Member Digital Library .....	55
IEEE Sales & Marketing.....	Cover 3
IEEE USA.....	10
Keysight.....	Cover 2, 1
MILCOM.....	13
National Instruments.....	5
PTC.....	14
Rohde & Schwarz .....	9
Siemens Industries .....	Cover 4
Tutorials.....	21
Wiley .....	11

## ADVERTISING SALES OFFICES

*Closing date for space reservation: 15th of the month prior to date of issue*

### NATIONAL SALES OFFICE

James A. Vick  
Sr. Director Advertising Business, IEEE Media  
EMAIL: [jv.ieeemediamedia@ieee.org](mailto:jv.ieeemediamedia@ieee.org)

Marion Delaney  
Sales Director, IEEE Media  
EMAIL: [md.ieeemediamedia@ieee.org](mailto:md.ieeemediamedia@ieee.org)

Mark David  
Sr. Manager Advertising & Business Development  
EMAIL: [m.david@ieee.org](mailto:m.david@ieee.org)

Mindy Belfer  
Advertising Sales Coordinator  
EMAIL: [m.belfer@ieee.org](mailto:m.belfer@ieee.org)

**NORTHERN CALIFORNIA**  
George Roman  
TEL: (702) 515-7247  
FAX: (702) 515-7248  
EMAIL: [George@George.RomanMedia.com](mailto:George@George.RomanMedia.com)

**SOUTHERN CALIFORNIA**  
Marshall Rubin  
TEL: (818) 888 2407

FAX: (818) 888-4907  
EMAIL: [mr.ieeemediamedia@ieee.org](mailto:mr.ieeemediamedia@ieee.org)

### MID-ATLANTIC

Dawn Becker  
TEL: (732) 772-0160  
FAX: (732) 772-0164

EMAIL: [db.ieeemediamedia@ieee.org](mailto:db.ieeemediamedia@ieee.org)

### NORTHEAST

Merrie Lynch  
TEL: (617) 357-8190  
FAX: (617) 357-8194

EMAIL: [Merrie.Lynch@celassociates2.com](mailto:Merrie.Lynch@celassociates2.com)

Jody Estabrook  
TEL: (77) 283-4528  
FAX: (774) 283-4527  
EMAIL: [je.ieeemediamedia@ieee.org](mailto:je.ieeemediamedia@ieee.org)

### SOUTHEAST

Scott Rickles  
TEL: (770) 664-4567  
FAX: (770) 740-1399  
EMAIL: [srickles@aol.com](mailto:srickles@aol.com)

### MIDWEST/CENTRAL CANADA

Dave Jones  
TEL: (708) 442-5633  
FAX: (708) 442-7620  
EMAIL: [dj.ieeemediamedia@ieee.org](mailto:dj.ieeemediamedia@ieee.org)

### MIDWEST/ONTARIO, CANADA

Will Hamilton  
TEL: (269) 381-2156  
FAX: (269) 381-2556  
EMAIL: [wh.ieeemediamedia@ieee.org](mailto:wh.ieeemediamedia@ieee.org)

### TEXAS

Ben Skidmore  
TEL: (972) 587-9064  
FAX: (972) 692-8138  
EMAIL: [ben@partnerspr.com](mailto:ben@partnerspr.com)

### EUROPE

Christian Hoelscher  
TEL: +49 (0) 89 95002778  
FAX: +49 (0) 89 95002779  
EMAIL: [Christian.Hoelscher@husonmedia.com](mailto:Christian.Hoelscher@husonmedia.com)

While the world benefits from what's new,  
IEEE can focus you on what's next.

IEEE *Xplore* can power your research  
and help develop new ideas faster with  
access to trusted content:

- Journals and Magazines
- Conference Proceedings
- Standards
- eBooks
- eLearning
- Plus content from select partners

**IEEE *Xplore*<sup>®</sup> Digital Library**

Information Driving Innovation

Learn More

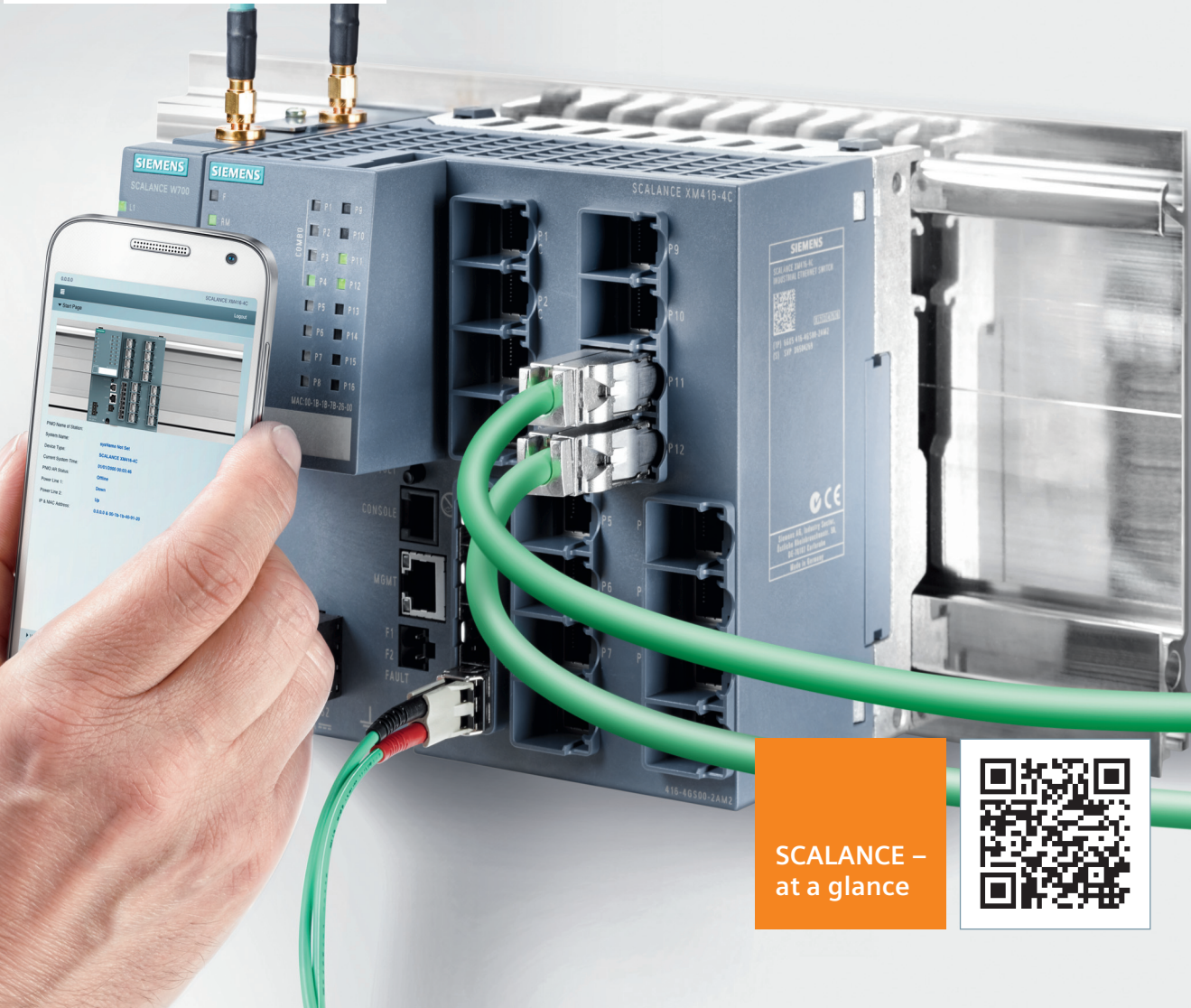
[innovate.ieee.org](http://innovate.ieee.org)

Follow IEEE *Xplore* on  



# SIEMENS

E20001-F730-P820-X-7600



SCALANCE –  
at a glance



## Industrial Ethernet Switches

# Switch to the future: SCALANCE XM-400

Secure and reliable data transmission is indispensable for industrial communication networks – not only today, but in the future as well. That makes it all the more important to be equipped with technology capable of keeping up with increasing requirements – as it is then much simpler to adapt or expand a system.

Industrial Ethernet Switches SCALANCE XM-400 are our innovative answer to the challenges of tomorrow – their modular design, fast mobile diagnostics for smartphones/tablets via NFC and existing WLAN, and upgradeable Layer 3 functionality using KEY-PLUG will ensure that your automation network is powerfully and flexibly equipped for the future.

[siemens.com/x-400](https://www.siemens.com/x-400)



# We're here to help you write your 5G future.

Unprecedented experience in wideband mmWave, 5G waveforms, and Massive MIMO design.



The world of wireless communications is about to change. Again. The fifth generation—5G—will mean “everything, everywhere, and always connected.” If you’re on the cutting edge of this emerging technology, we can help you. We have expertise in all areas of 5G research and development, including wideband mmWave, radio spectrum, ASIC, antenna technologies, and network architecture. So from design simulation and verification to wideband signal generation and analysis, from component characterization to optical solutions, we’ve got you covered.

## HARDWARE + SOFTWARE + PEOPLE = 5G INSIGHTS

Keysight engineers are active in the leading 5G forums and consortia

Keysight engineers are keynote speakers at 5G conferences and key contributors in top technical journals

Applications engineers are in more than 100 countries around the world



Unlocking Measurement Insights



Download our white paper *Implementing a Flexible Testbed for 5G Waveform Generation and Analysis* at [www.keysight.com/find/5G-Insight](http://www.keysight.com/find/5G-Insight)



USA: 800 829 4444 CAN: 877 894 4414

© Keysight Technologies, Inc. 2015



#### Director of Magazines

Steve Gorshe, PMC-Sierra, Inc (USA)

#### Editor-in-Chief

Osman S. Gebizlioglu, Huawei Tech. Co., Ltd. (USA)

#### Associate Editor-in-Chief

Zoran Zvonar, MediaTek (USA)

#### Senior Technical Editors

Nim Cheung, ASTRI (China)

Nelson Fonseca, State Univ. of Campinas (Brazil)

Steve Gorshe, PMC-Sierra, Inc (USA)

Sean Moore, Centripetal Networks (USA)

Peter T. S. Yum, The Chinese U. Hong Kong (China)

#### Technical Editors

Sonia Aissa, Univ. of Quebec (Canada)

Mohammed Atiquzzaman, Univ. of Oklahoma (USA)

Guillermo Atkin, Illinois Institute of Technology (USA)

Mischa Dohler, King's College London (UK)

Frank Effenberger, Huawei Technologies Co., Ltd. (USA)

Tarek El-Bawab, Jackson State University (USA)

Xiaoming Fu, Univ. of Goettingen (Germany)

Stefano Galli, ASSIA, Inc. (USA)

Admela Jukan, Tech. Univ. Carolo-Wilhelmina zu

Braunschweig (Germany)

Vimal Kumar Khanna, mCalibre Technologies (India)

Myung J. Lee, City Univ. of New York (USA)

Yoichi Maeda, TTC (Japan)

Nader F. Mir, San Jose State Univ. (USA)

Seshradi Mohan, University of Arkansas (USA)

Mohamed Moustafa, Egyptian Russian Univ. (Egypt)

Tom Oh, Rochester Institute of Tech. (USA)

Glenn Parsons, Ericsson Canada (Canada)

Joel Rodrigues, Univ. of Beira Interior (Portugal)

Jungwoo Ryoo, The Penn. State Univ.-Altoona (USA)

Antonio Sánchez Esguevillas, Telefonica (Spain)

Mostafa Hashem Sherif, AT&T (USA)

Tom Starr, AT&T (USA)

Ravi Subrahmanyan, InVisage (USA)

Danny Tsang, Hong Kong U. of Sci. & Tech. (China)

Hsiao-Chun Wu, Louisiana State University (USA)

Alexander M. Wyglinski, Worcester Poly. Institute (USA)

Jun Zheng, Nat'l. Mobile Commun. Research Lab (China)

#### Series Editors

##### *Ad Hoc and Sensor Networks*

Edoardo Biagioni, U. of Hawaii, Manoa (USA)

Silvia Giordano, Univ. of App. Sci. (Switzerland)

##### *Automotive Networking and Applications*

Wai Chen, Telcordia Technologies, Inc (USA)

Luca Delgrossi, Mercedes-Benz R&D N.A. (USA)

Timo Kosch, BMW Group (Germany)

Tadao Saito, Toyota Information Technology Center (Japan)

##### *Consumer Communications and Networking*

Ali Begen, Cisco (Canada)

Mario Kolberg, University of Sterling (UK)

Madjid Merabti, Liverpool John Moores U. (UK)

##### *Design & Implementation*

Vijay K. Gurbani, Bell Labs/Alcatel Lucent (USA)

Salvatore Loreto, Ericsson Research (Finland)

Ravi Subrahmanyan, Invisage (USA)

##### *Green Communications and Computing Networks*

Daniel C. Kilper, Univ. of Arizona (USA)

John Thompson, Univ. of Edinburgh (UK)

Jinsong Wu, Alcatel-Lucent (China)

Honggang Zhang, Zhejiang Univ. (China)

##### *Integrated Circuits for Communications*

Charles Chien, CreoNex Systems (USA)

Zhiwei Xu, HRL Laboratories (USA)

##### *Network and Service Management*

George Pavlou, U. College London (UK)

Juergen Schoenwaelder, Jacobs University (Germany)

##### *Networking Testing*

Ying-Dar Lin, National Chiao Tung University (Taiwan)

Erica Johnson, University of New Hampshire (USA)

##### *Optical Communications*

Osman Gebizlioglu, Huawei Technologies (USA)

Vijay Jain, Sterlite Network Limited (India)

##### *Radio Communications*

Thomas Alexander, Ixia Inc. (USA)

Amitabh Mishra, Johns Hopkins Univ. (USA)

#### Columns

##### *Book Reviews*

Piotr Cholda, AGH U. of Sci. & Tech. (Poland)

#### Publications Staff

Joseph Milizzo, Assistant Publisher

Susan Lange, Online Production Manager

Jennifer Porcello, Production Specialist

Catherine Kemelmacher, Associate Editor

# COMMUNICATIONS STANDARDS

A Supplement to IEEE Communications Magazine

## SEPTEMBER 2015

#### SUPPLEMENT EDITOR

GLENN PARSONS

#### MANAGING EDITOR

JACK HOWELL

#### STANDARDS NEWS CONTRIBUTORS

NASSER SALEH AL MARZOUQI • SANGKEUN YOO • ROBBY SIMPSON • PAT KINNEY • CLINTON POWELL  
SAMITA CHAKRABARTI • PERIKLIS CHATZIMISIOS

- 2 **EDITOR'S NOTE: RESEARCH AND STANDARDS**  
GLENN PARSONS, EDITOR-IN-CHIEF
- 3 **COMMENTARY**  
BRUCE KRAEMER, PRESIDENT, IEEE STANDARDS ASSOCIATION  
YATIN TRIVEDI, PAST-CHAIR, IEEE STANDARDS EDUCATION COMMITTEE
- 5 **COMMUNICATIONS STANDARDS NEWS**  
**IoT/M2M FROM RESEARCH TO STANDARDS: THE NEXT STEPS (PART I)**
- 8 **GUEST EDITORIAL**  
OMAR ELLOUMI, JAESEUNG SONG, YACINE GHAMRI-DOUDANE, AND VICTOR LEUNG
- 10 **MACHINE-TYPE COMMUNICATIONS: CURRENT STATUS AND FUTURE PERSPECTIVES TOWARD 5G SYSTEMS**  
HAMIDREZA SHARIATMADARI, RAPEEPAT RATASUK, SASSAN IRAJI, ANDRÉS LAYA, TARIK TALEB, RIKU JÄNTTI, AND AMITAVA GHOSH
- 18 **NEXT GENERATION M2M CELLULAR NETWORKS: CHALLENGES AND PRACTICAL CONSIDERATIONS**  
ABDELMOHSEN ALI, WALAA HAMOUDA, AND MURAT UYSAL
- 26 **CONNECTIONLESS ACCESS FOR MOBILE CELLULAR NETWORKS**  
COLIN KAHN AND HARISH VISWANATHAN
- 32 **UNDERSTANDING THE IoT CONNECTIVITY LANDSCAPE: A CONTEMPORARY M2M RADIO TECHNOLOGY ROADMAP**  
SERGEY ANDREEV, OLGA GALININA, ALEXANDER PYATTAEV, MIKHAIL GERASIMENKO, TUOMAS TIRRONEN, JOHAN TORSNER, JOACHIM SACHS, MISCHA DOHLER, AND YEVGENI KOUCHERYAVY
- 41 **WHAT CAN WIRELESS CELLULAR TECHNOLOGIES DO ABOUT THE UPCOMING SMART METERING TRAFFIC?**  
JIMMY J. NIELSEN, GERMÁN C. MADUEÑO, NUNO K. PRATAS, RENÉ B. SØRENSEN, ČEDOMIR STEFANOVIĆ, AND PETAR POPOVSKI
- 48 **PUBLISH/SUBSCRIBE-ENABLED SOFTWARE DEFINED NETWORKING FOR EFFICIENT AND SCALABLE IoT COMMUNICATIONS**  
AKRAM HAKIRI, PASCAL BERTHOU, ANIRUDDHA GOKHALE, AND SLIM ABDELLATIF
- 55 **SOFTWARE-DEFINED INTERNET OF THINGS FOR SMART URBAN SENSING**  
JIAQIANG LIU, YONG LI, MIN CHEN, WENXIA DONG, AND DEPENG JIN
- 64 **LOW POWER WIDE AREA MACHINE-TO-MACHINE NETWORKS: KEY TECHNIQUES AND PROTOTYPE**  
XIONG XIONG, KAN ZHENG, RONGTAO XU, WEI XIANG, AND PERIKLIS CHATZIMISIOS
- 72 **TOWARD BETTER HORIZONTAL INTEGRATION AMONG IoT SERVICES**  
ALA AL-FUQAHA, ABDALLAH KHREISHAH, MOHSEN GUIZANI, AMMAR RAYES, AND MEHDI MOHAMMADI



### INTERNET OF THINGS: FUNDAMENTAL TRANSFORMATION AGENT



*Glenn Parsons*

**B**y creating digital representations of the real world, the Internet of Things (IoT) allows us to interact, monitor, and control systems and environments as never before. Some predictions anticipate IoT will initiate a 100 fold increase in the number of connected devices, resulting in a demand for lower network latency and power consumption. IoT will become a fundamental transformation agent across a wide spectrum of industries and society. This rapid change and lack of standardization has resulted in a fragmentation of technology. To create a truly networked society requires market driven standardization and the deployment of standardized technology for IoT devices, networks, and software. In the end, users expect seamless connectivity, security, and manageability end-to-end.

The importance of standards to the work and careers of communications practitioners is the basis of this publication. It is a platform for presenting and discussing standards related topics in the areas of communications, networking, research, and related disciplines. This issue of the Communications Standards Supplement contains a very relevant feature topic on IoT. JaeSeung Song assembled an esteemed editorial team to develop the feature topic, “IoT/M2M from Research to Standards: The Next Steps (Part I).” This topic comprises all of the articles in this issue, and an editorial from the feature topic editors will introduce each paper. Given the overwhelming response to this feature topic, the selected papers have been split between this issue and the next one in December.

Readers will notice the ongoing Commentary section with a recurring view from the IEEE-SA President, and an introduction to the IEEE Standards Education activity. The Standards News section offers the current status of standards work in various SDOs relevant to IoT and pointers to SDO material. I trust that the reader will find these informative and illustrative of the fundamental role standards play in the communications networking ecosystem.

As the Communications Standards Supplement matures, each issue will be “anchored” around a topic of current market relevance to drive focus. Proposals for future standards-related feature topics are welcome. The next issue will contain the second part of the feature topics on the Internet of Things (IOT) and on Advanced Cloud & Virtualization Techniques for 5G Networks.

#### BIOGRAPHY

GLENN PARSONS [SM] ([glenn.parsons@ericsson.com](mailto:glenn.parsons@ericsson.com)) is an internationally known expert in mobile backhaul and Ethernet technology. He is a standards advisor with Ericsson Canada, where he coordinates standards strategy and policy for Ericsson, including network architecture for LTE mobile backhaul. Previously, he has held positions in development, product management and standards architecture in the ICT industry. Over the past number of years, he has held several management and editor positions in various standards activities including IETF, IEEE, and ITU-T. He has been an active participant in the IEEE-SA Board of Governors, Standards Board and its Committees since 2004. He is currently involved with mobile backhaul standardization in MEF, IEEE and ITU-T and is chair of IEEE 802.1. He is a Technical Editor for *IEEE Communications Magazine* and has been co-editor of several IEEE Communications Society Magazine feature topics. He graduated in 1992 with a B.Eng. degree in electrical engineering from Memorial University of Newfoundland.

## IEEE STANDARDS PHILOSOPHY

BY BRUCE KRAEMER, PRESIDENT, IEEE STANDARDS ASSOCIATION

IEEE is uniquely positioned as an independent, global, not-for-profit, technology-driven organization, with the flexibility to engage decision-making as a neutral body not representing any single constituency or government interest. In congruence with these principles, IEEE standards are intended to define technology concepts that can be used as blueprints for constructing technically feasible and economically attractive products that broadly benefit humanity.

A premise of IEEE standards activities is that the deliberations on how to construct and describe the standard will be open to anyone who wishes to contribute and, when the standard is complete, anyone who chooses to implement it will be allowed to do so. One portion of the ecosystem that applies to these stated goals is a policy for inclusion of intellectual property in the standards.

Contributions from IEEE Standards Association participants requires the expenditure of time, effort, and money. IEEE believes these participants should be reasonably rewarded for their contributions through the ability to market standards related goods and services. One form of compensation can be licensing of intellectual property (IP) included in a standard. To provide guidance to standards developers regarding the process for including IP, IEEE maintains a Patent Policy as part of the Standards Board Bylaws.

The role of a standards setting body like the IEEE Standards Association is not to establish or negotiate licensing arrangements but rather to set the ground rules that ensure a level playing field for all participants who engage in standards development, whether they are from small or big companies, incumbents or newcomers.

### PATENT POLICY HISTORY

The IEEE-SA has maintained a Patent Policy for many years, with the most recent prior update occurring in late 2006. Since then, the IEEE-SA and other standards developing organizations (SDOs) with similar IP policies noted that there remained problems in the business ecosystem. ITU observations resulted in a one-day patent roundtable in October 2012, where the expressed motivation for the event was "... recent patent disputes that have caused shipments of goods to be impounded in customs and recent worldwide increase in litigation involving standard essential patents (or SEPs)..." IEEE participated in that event and similarly noticed that too frequently negotiations between patent holders and potential licensees were unsuccessful, and the lack of agreement was followed by lengthy and expensive court litigations.

Discussions of the IEEE Standards Association (IEEE-SA) Patent Committee (PatCom) and observations of the general standards and intellectual property communities led to the creation of a PatCom ad hoc in February 2013 to study possible policy changes. Over a two year period the PatCom ad hoc considered court rulings, information from patent offices, and other standards development organizations, as well as statements from all interested stakeholders. The result of this process was that IEEE enacted an updated IEEE-SA Patent Policy on March 15, 2015.

### POLICY CONTENTS

This short article cannot completely describe the updates to the IEEE-SA Patent Policy, but we can briefly cover some highlights. (A significant body of explanatory and educational material is available for study at the links shown below.)

As part of the update, language was added clarifying what is meant by "reasonable" rate by explicitly defining the term to mean "appropriate compensation to the patent holder for the practice of an Essential Patent Claim excluding the value, if any, resulting from the inclusion of that Essential Patent Claim's technology in the IEEE Standard", as well as providing optional factors for the licensing party to consider. (Please note that the IEEE does not determine reasonable rates. Royalty rates are determined by negotiations between parties or litigation or arbitration, if necessary). The update also provides greater clarity on the topics of nondiscrimination, availability of prohibitive orders, and permissible demands for reciprocal license.

IEEE requests that holders of potentially essential patent claims voluntarily submit a letter of assurance (LOA) that states that, should the intellectual property be included in the standard, the licensor will provide a license to any implementer of the standard on fair, reasonable, and non-discriminatory (FRAND) terms. The submission of LOAs remains voluntary under the update.

### CONCLUSION

The updated IEEE-SA Patent Policy is a very carefully considered step toward clarity for owners and users of potentially essential IP in standards. By clarifying rules we expect all participants to fully understand, respect, and employ them.

The IEEE-SA and its patent policy has always intended to foster a healthy, inclusive standardization ecosystem that does not favor one party over another and gives everyone an equal seat at the table. We want innovation to remain open to all on a global basis, especially as we continue to explore the next generation of world-changing technologies.

### RESOURCES

Questions? Contact the IEEE-SA Standards Board Patent Committee Administrator at [patcom@ieee.org](mailto:patcom@ieee.org), or visit these resources:

- Patent Policy (IEEE-SA Standards Board Bylaws clause 6), <http://standards.ieee.org/develop/policies/bylaws/sect6-7.html#6Bylaws>
- Understanding Patent Issues During IEEE Standards Development, <http://standards.ieee.org/faqs/patents.pdf>
- PatCom Patent Materials, <https://standards.ieee.org/about/sasb/patcom/materials.html>
- Drafts, comments and responses publically available on pp-dialog, <http://grouper.ieee.org/groups/pp-dialog/index.html>
- Promoting Competition and Innovation: What You Need to Know about the IEEE Standards Association's Antitrust and Competition Policy, <http://standards.ieee.org/develop/policies/antitrust.pdf>
- IEEE-SA Standards Board Resolution (2 June 2015), <https://standards.ieee.org/about/sasb/resolutions.html>

## IEEE STANDARDS UNIVERSITY: NEW FRONTIER IN STANDARDS EDUCATION

By YATIN TRIVEDI, IEEE MEMBER AND PAST-CHAIR, IEEE STANDARDS EDUCATION COMMITTEE

Standards fuel the development and implementation of technologies that influence and transform the way we live, work, and communicate. They form the fundamental building blocks for product and new technology development around the world, by establishing consistent protocols that can be universally understood and adopted. Standards also enable compatibility and interoperability, simplify product development, hasten time-to-market for new products, and play a vital role in helping consumers understand and compare competing products. As standards are globally adopted and applied in many markets, they also drive international trade.

IEEE's commitment to educating students, faculty, and practicing professionals about technical standards takes many forms. As part of the effort to actively promote the integration of standards into academic programs, on 28 June 2009, the IEEE Board of Directors approved an IEEE Position Paper defining the desired role of technical standards in education within engineering, technology, and computing (ETC) academic curricula in the technical areas of interest to the IEEE.<sup>1</sup>

The IEEE Educational Activities Board and the IEEE Standards Association Board of Governors work together through a joint standing committee of volunteers, the IEEE Standards Education Committee, to oversee this commitment to standards education.

The IEEE Standards Education Committee in 2015 is overseeing an exciting new era for standards education by implementing an IEEE New Initiative Committee supported project called IEEE Standards University. The IEEE Standards University is a multi-track initiative intended to greatly expand standards education content. The program is designed to train, educate, and inspire the next generation of professionals about the critical importance standards will play in the workforce and society.

The new web site for IEEE Standards University will

launch before the end of 2015. It will bring IEEE standards-related content under one umbrella:

- Users will be able to find eLearning and other courses related to standards.
- Search for standards.
- Read case studies, articles, and student papers.
- Find ideas from academics who teach standards in their courses.
- Apply for student grants for design projects using standards.
- Request a guest lecturer or workshop on standards related topics.
- Watch videos.
- Read the e-Magazine, and much more.

In the spring of 2016, the first massive open online course (MOOC) focused exclusively on standards will be offered, titled "Innovation and Competition: Succeeding through Global Standards." Delivered through the edx.org platform on the IEEE's new site, IEEEEx.org, this six-week introductory course will provide a survey of fundamental standardization themes including history, basic concepts and classifications, impact on innovation, global markets, product design, implementation and strategy, building consensus, conformity assessment, regulation, intellectual property, and motivating factors in standards development. It is designed to reach a worldwide audience and to help students develop skills around technical standards enabling them to compete more effectively in the global economy. Engineering faculty interested in integrating technical standards in their classes may find this course ready for adoption. A little customization such as local government regulations may help them quickly communicate the importance of standards to the future professional engineers.

The Standards Education Committee is a collaborative effort between IEEE volunteers and professional staff. Your comments and critiques of our efforts are welcome. We invite you to join us in this effort. Please contact Jennifer McClain at [j.mcclain@ieee.org](mailto:j.mcclain@ieee.org).

<sup>1</sup> [http://www.ieee.org/education\\_careers/education/standards/standards\\_position\\_paper.html](http://www.ieee.org/education_careers/education/standards/standards_position_paper.html)

## ITU-T STUDY GROUP 20 ON IoT AND ITS APPLICATIONS, INCLUDING SMART CITIES AND COMMUNITIES

NASSER SALEH AL MARZOUQI (UNITED ARAB EMIRATES), CHAIRMAN OF ITU-T STUDY GROUP 20

The deployment of Internet of Things (IoT) technologies is expected to connect an estimated 50 billion devices to the network by the year 2020, impacting nearly every aspect of our daily lives. IoT is contributing to the convergence of industry sectors, with utilities, healthcare, and transportation among the many sectors with a stake in the future of IoT technologies.

The newly established ITU-T Study Group 20 on “IoT and its Applications, Including Smart Cities and Communities” provides the specialized IoT standardization platform necessary for this convergence to rest on a cohesive set of international standards. With participants representing the many stakeholders in the field of information and communication technologies, this Study Group will be influential in promoting the development of the highly efficient ‘systems of systems’ that will help bridge the digital divide and make possible a more connected world.

Building on the foundations set by the ITU-T Focus Group on Smart Sustainable Cities, the new ITU-T Study Group 20 will guide cities in upgrading their traditional infrastructures by integrating new digital technologies. The group will center its work on the coordinated development of IoT technologies, including machine-to-machine communications and ubiquitous sensor networks.

The focus of the research conducted by ITU-T Study Group 20 will be on identifying and analyzing emerging applications and global solutions for IoT and smart cities. The group’s research and standardization work will contribute to improving the interoperability of various IoT-based technologies, a key factor in ensuring end-user and market acceptance of IoT solutions.

ITU-T Study Group 20 is well-placed to analyze the potential paradigm shift in urban areas to occur as a consequence of the use of IoT in smart cities. The group will also assist in the standardization of end-to-end architectures for IoT and the mechanisms for the interoperability of IoT applications and datasets employed by various vertically oriented industry sectors.

The first meeting of ITU-T Study Group 20 will be held at ITU Head-

quarters in Geneva, 19-23 October 2015, preceded by a Forum on the topic ‘IoT: Empowering the New Urban Agenda’ that will take place on 19 October 2015. For more information, please contact [tsbsg20@itu.int](mailto:tsbsg20@itu.int).

## INTERNET OF THINGS IN ISO/IEC JTC 1/WG 10

SANGKEUN YOO, ETRI, WG 10 CONVENOR

ISO/IEC JTC 1/WG 10 was established at the 2014 JTC 1 Plenary with a standard project, ISO/IEC 30141 (IoT Reference Architecture). Among other IoT standards, WG 10 mainly covers foundational standards including terms and definitions and reference architecture. WG 10 formed five internal groups to conduct ToR given by the JTC 1 Plenary, and to make efficient progress in ongoing work.

- Group 1: Standardization gaps.
- Group 2: Network level technologies for IoT.
- Group 3: IoT Identification.
- Group 4: Systems integration guidelines.
- Group 5: Conceptual reference model for IoT reference architecture.

Group 1 is collecting IoT standards data from various SDOs and finding gaps among standardization works. Network level technologies are one of the topics in which WG 10 has interests. Group 2 explores interoperable IoT technologies at the network level. In IoT systems, identification is one of the foundational technologies that WG 10 should take. Object identifier-based identification for IoT has been discussed in Group 3. Group 4 deals with the systems integration approach that is under development by JTC 1/SWG-M (Special Working Group on Management) to cover standard topics that extend over wide areas such as IoT and Big data. Group 5 is developing a conceptual reference model for the IoT reference architecture embedded in ISO/IEC 30141.

WG 10 approaches the IoT reference architecture in terms of following:

- Characteristics of IoT systems, characteristics.
- Characteristics, principles, and requirements of the IoT RA.
- IoT reference architecture framework.
- IoT conceptual reference model.
- Reference architecture (RA) of IoT Systems.

Besides ISO/IEC 30141, WG 10 has submitted a NWIP on the definition and vocabulary of IoT to JTC 1, and this NWIP is under JTC 1 ballot for

approval. Also, WG 10 initiated a technical report on IoT use cases to identify IoT scenarios and use cases based on real-world applications and requirements.

WG 10 has three meetings in a year. This year the scheduled meetings are: January in Berlin, May in Brussels, and August in Ottawa.

## IEEE 2030.5: AN IoT STANDARD FOR SMART GRID CONSUMER COMMUNICATIONS

ROBBY SIMPSON, MEMBER, IEEE-SA CAG, SYSTEM ARCHITECT, GE DIGITAL ENERGY

IEEE 2030.5, Standard for Smart Energy Profile Application Protocol, defines an interoperable profile for consumer interaction with the smart grid. This standard is gaining popularity as both a smart grid and Internet of Things (IoT) standard with its focus on interacting with consumer devices.

Recently, IEEE 2030.5 has been incorporated as part of California Rule 21 for integration of smart inverters for distributed energy resources (e.g. solar). The P2030.5 working group is working closely with implementers to incorporate any feedback from field experience and new requirements into future revisions of the standard.

In addition to DER integration, IEEE 2030.5 is used to communicate with a variety of consumer devices such as smart phones, smart meters, electric vehicles, thermostats, water heaters, and appliances. In general, IEEE 2030.5 can be used to inform these devices (e.g. energy usage, pricing) as well as request actions from those devices (e.g. thermostat changes, electric vehicle charging scheduling).

Although recently much emphasis has been on energy, IEEE 2030.5 also supports other commodity types such as water, natural gas, and steam.

As an IoT standard, IEEE 2030.5 runs over the Internet Protocol (IP) to enable a variety of communication paths and link layers. Further, the standard is designed to be lightweight, reducing the cost to integrate the standard into consumer products as well as to enable applications that may be battery-powered. And of course, IEEE 2030.5 is designed with security in mind, incorporating the latest in cybersecurity technology and best practices.

The P2030.5 welcomes new contributors and encourages other activities in the IoT space to consider their approach to M2M communication for their IoT activities.

## IEEE 802.15 – WIRELESS PERSONAL AREA NETWORKS

PAT KINNEY, IEEE 802.15 WG VICE CHAIR,  
CLINTON POWELL, IEEE 802.15 – TG10  
(LAYER 2 ROUTING) CHAIR

For more than a decade, the IEEE 802.15 Working Group has been in the forefront of the Internet of Things (IoT) with standards focused on machine to machine (M2M) applications with little if any human interaction. The increasing recognition of the benefits of IoT have dramatically increased the usage of 802.15 standards.

IEEE 802.15 includes a family of complementary standards dealing with diverse applications in the consumer, industrial, and commercial spaces. Industry-wide uptake and ongoing developments are continuing for several of the 802.15 work streams (primary Task Groups) including 802.15.3, 802.15.4, 802.15.6, 802.15.7, and 802.15.8.

- 802.15.3 provides high data rate over short distances for the efficient transfer of large data payloads with very little overhead.

- 802.15.4 offers low data rate, low cost, and low energy, with simple implementation for small payloads and periodic communication needs over a long period of time unattended.

- 802.15.6 focuses on body area networks, featuring low energy and high reliability for implant devices and other medical sensing, all of which enable remote monitoring of patients by doctors helping to curb the increases health care costs.

- 802.15.7 defines optical link networks using devices such as the flash, display, and image sensor as the transmitting and receiving devices for point of service areas, driver information using street lights, and commercial communication using light fixtures.

- 802.15.8 standardizes peer aware communications (PAC) for peer to peer and infrastructure-less communications with fully distributed coordination.

Among the 802.15 standards, 802.15.4 is by far the most popular and prolific due to its focus on remote sensors and monitoring, with hundreds of millions of devices deployed. To accommodate the anticipated large quantities of these types of devices, 802.15.4 was the first 802 standard to use the 64-bit IEEE organizationally unique identifier (OUI), making it ideal for use with IPv6 networks.

The 802.15.4 standard is focused on low cost and complexity devices, with low energy consumption for long battery life or use of an energy harvesting

power source. This standard was designed from the beginning to work within mesh networks to extend network coverage and link reliability while maintaining low energy consumption.

This standard consists of a media access control (MAC) sublayer and a variety of physical layers. The MAC is very configurable and flexible, allowing applications to elide unnecessary fields and disable optional behaviors as appropriate. To address diverse applications, the physical layers cover frequencies ranging from 433 MHz to 6500 MHz, along with modulations such as FSK, BPSK, O-QPSK, OFDM, and UWB. Finally, 802.15.4 is also very diverse in that it supports communication distances from a few feet up to several kilometers.

As a result of its capabilities and flexibility, the 802.15.4 standard has received widespread acceptance from consortiums and alliances such as ZigBee, Thread, HART, and Wi-SUN, as well as other standards development organizations such as ISA (62734), TIA (TR-51), IETF (6LoWPAN and 6tisch), ISO/IEC (24730-62), and ETSI (TS 102 887-1).

In its liaison letter to 802.15, the ETSI TC ERM stated that “802.15.4 standards are also important components of the M2M infrastructure leading to the Internet of Things.”

Applications and uses for 802.15.4 devices and networks include:

- Utility applications for both the HAN and FAN spaces, including AMI and smart grid for electric meters, gas meters, and water meters.

- Industrial wireless networks for industrial automation and process control.

- Lighting control networks for commercial buildings that are easily configured, energy efficient, and user friendly.

- Home automation applications such as remote controls, lighting, and HVAC controls.

- Medical applications including patient care connecting sensors to medical devices, and remote patient monitoring (reducing hospital stays).

- Critical infrastructure such as monitoring pipelines, electrical vaults, and perimeters.

- Railways, in applications such as positive train control.

- RFID, for asset tracking including precision location awareness.

- International space station, for asset tracking and control.

- Consumer electronics such as cable/set top boxes, integrated A/V, and home theatre control.

Current efforts within IEEE 802.15.4

include providing increased support of 802.15.4 to IoT applications developments such as defining key management protocols for security, an advanced Layer 2 routing protocol for mesh networks, and a new LLC focused on simplifying the use of 802.15.4 devices and networks.

## IETF IPv6-OVER-LOW POWER AND OTHER IOT STANDARDS ACTIVITIES

SAMITA CHAKRABARTI, ERICSSON, SAN JOSE, USA

IEEE released the 802.15.4 low power wireless personal area network standard in 2003 as the stepping stone for the global low power radio standard for small embedded devices. Soon after this release, that IETF formed a new workgroup to standardize ‘IPv6-over-IEEE 802.15.4’ or the 6LoWPAN WG in order to integrate IP on sensor devices with IEEE 802.15.4 radio. Given the special requirements for low power devices with limited processing, bandwidth, radio power, etc., the 6LoWPAN working group had a set of unique requirements that are quite different from regular IPv6 standardization on the standard PC or IP-enabled devices. One of them was the need for a simple and stateless compression mechanism for the IPv6 header (40 bytes) which was perhaps carrying only 10-20 bytes of IoT data over the low power and lossy networks.

The choice of IPv6 addressing over IPv4 on the IoT devices is clear as IPv6 naturally offers a large range of IP addresses over a subnet considering the billions of such interconnected devices. 6LoWPAN produced the basic framework of IPv6-over-IEEE 802.15.4 devices and produced three main documents: RFC4944, RFC6282, and RFC6775. RFC4944 describes the frame format for IPv6 packets, methods of forming the IPv6 addresses on the IEEE 802.15.4 networks, and the 6lowpan adaptation layer frames. RFC6282 followed RFC4944 describing the compression technique for 6LoWPAN packets, while RFC6775 provides a set of optimizations for saving neighbor discovery control messages and making the booting process reliable in the lossy and low-power radio network. The 6LowPAN stack is widely accepted in the industry. ZigBeeIP, Threads, NIST Smart Energy profile version 2, and ITU-T are the other standards forums that have adopted 6LoWPAN technology.

The 6LoWPAN WG has successfully ended their work on IEEE 802.15.4,

and IETF formed a successor WG in 2013 called ‘6lo’ (IPv6 over Networks of Resource Constrained Nodes) as the popularity of the 6LoWPAN stack principle continues to grow and the intent of running the 6LoWPAN stack over other low power radio technologies strengthened. 6lo is chartered to facilitate IPv6 connectivity over different constrained node networks with characteristics of limited power, memory, processing power, code space and processing cycles. 6lo is built upon 6LoWPAN technologies with modifications for the specified wireless low power radio network. The 6lo group has already produced IPv6-over-Zwave (RFC 7428), Generic Header Compression over 6LoWPAN (RFC7400), and Definitions of Managed Objects for 6LoWPAN (RFC 7388). IPv6-over-Bluetooth-Low Energy RFC will be published soon. IPv6 over DECT-ULE, IPv6-over Low Power BACNET, and IPv6-over NFC are a few of the notable works in this area. Other work includes further optimization for running constrained node networks in terms of security, bootstrapping, etc., and improvement of auto-configuration are encouraged. The 6lo IETF WG is currently working closely with ITU-T and IEEE and it looks forward to working with other standards organizations.

6lo nodes can connect to each other by L2 and L3 routing protocols. IETF has developed RPL (RFC6550) for connecting nodes in directed graph mesh networks using L3 routing. It can form the network dynamically as the nodes leave and join. This is the product of the IETF ROLL working group (WG).

As interest in the Internet of things is growing in the industry, there are multiple new working groups have been formed at IETF, and many other working groups are considering ‘constrained

nodes’ behavior in their protocol considerations. 6Tisch is working on defining protocols and a shim-layer between 6LoWPAN and IEEE 802.15.4e (TSCH mode). 6Tisch provides a reliable transmission over 6LoWPAN for industrial level reliable networks. The CORE (constrained restful environments) WG defines API and COAP (RFC7252) protocols for constrained nodes communications. ACE (authentication and authorization for constrained environments) and DICE (DTLS in constrained environments) are the security related working groups that define application and transport level security for the constrained networks. Recently IETF has formed a team of IOT directorates to co-ordinate with internal constrained WG communications and information sharing.

Today the IoT deployment space is mostly vertical with many devices with proprietary protocols and their applications in the cloud. Ideally from the architectural point of view, the integration of IoT networks with common IP(v6) infrastructure is necessary irrespective of their L2 radio technologies for efficient transport networks and effective delivery of IoT data along with the advantage of existing service and network management infrastructure.

## UPDATES ON THE IEEE COMSoc STANDARDIZATION RESEARCH GROUP ON SOFTWARE DEFINED AND VIRTUALIZED WIRELESS ACCESS

PERIKLIS CHATZIMISIOS, ALEXANDER TEI OF THESSALONIKI, GREECE

The Standardization Research Group (RG) on Software Defined and Virtualized Wireless Access was formed during the Rapid Standardization Activity Working Meetings that were recently

organized by the IEEE Communications Society Standards Activities Council. This RG is chaired by Prof. Fabrizio Granelli (University of Trento) and works on identifying and addressing the research issues that need to be solved, and assess the feasibility of launching an IEEE standardization effort on software defined and virtualized wireless access, with specific focus on different future scenarios.

The Standardization RG recently completed the first draft of a white paper, entitled “Software Defined and Virtualized Wireless Access in Future Wireless Networks: Scenarios, Standards and Standardization Opportunities”. The purpose of the white paper is to provide a review of possible scenarios for standards in the area of SDN over wireless and to identify possible standardization opportunities. The white paper will probably be completed by the end of 2015.

Some technical content about the progress of the research group is available in the following publication: F. Granelli, A. A. Gebremariam, M. Usman, F. Cugini, V. Stamati, M. Alitska, and P. Chatzimisios, “Software defined and virtualized wireless access in future wireless networks: scenarios and standards,” *IEEE Commun. Mag.*, vol. 53, no. 6, pp. 26–34, June 2015 (DOI: 10.1109/MCOM.2015.7120042).

For more information or to contribute to the activities of the research group, please visit [1].

### ACKNOWLEDGMENTS

The author would like to thank Prof. Granelli for providing information and details about the Standardization Research Group on Software Defined and Virtualized Wireless Access.

### REFERENCES

[1] <http://community.comsoc.org/groups/rg-software-defined-and-virtualized-wireless-access>

## IoT/M2M FROM RESEARCH TO STANDARDS: THE NEXT STEPS (PART I)



Omar Elloumi



JaeSeung Song



Yacine Ghamri-Doudane



Victor C.M. Leung

As the pace of IoT deployments accelerate, IoT standards are undergoing major evolutions, sometimes revolutions. For instance, cellular networks standards are now adding techniques to improve network performance to address traffic patterns generated by an increasing number of IoT devices. Ongoing discussions around 5G requirements may become game changing for M2M communications because the standard will be designed, from the ground-up, for massive scale IoT deployments. This is a radical shift compared to the “quick-fixes” 3GPP and 3GPP2 have been adding to 2G/3G and 4G standards so far. Another example of this radical shift is related to IoT service platforms (such as the platform standardized by oneM2M) and IoT applications. Semantic interoperability is now emerging as a major trend that allows data exchange between applications, an increased level of interoperability, analytics, and reasoning. With ontologies engineering, researchers will soon overcome the limitations of static data models and bridge the gap between the currently deployed vertical silos. Other areas that will see intense standardization activity are IoT security and low power wide area connectivity.

The articles selected for this feature topic can broadly be grouped into four categories: networking including 5G, semantic interoperability, security, and low power wide area connectivity. This Feature Topic, which attracted no less than 40 submissions, out of which we have selected 13 articles, is structured into two parts. Part I, published in this issue, addresses network and connectivity topics, while Part II, to be published in the December issue, will mostly cover semantic interoperability and security topics.

The first article, “Machine-Type Communications: Current Status and Future Perspectives Toward 5G Systems” by Shariatmadari *et al.*, provides an overview of machine type communication (MTC) challenges and potential solutions toward fully MTC-capable 5G mobile systems.

The second article, “Next Generation M2M Cellular Networks: Challenges and Practical Considerations” by Ali *et al.*, provides a good overview of the challenges for cellular networks to address M2M communications require-

ments such as low power and low cost. The authors argue that cognitive radio provides a cost optimized solution addressing the requirements.

The third article, “Connectionless Access for Mobile Cellular Networks” by Khan *et al.*, shows the limitation of the LTE connection-oriented approach and proposes an efficient alternative based on a connection-less access method that is efficient for small burst transmissions in cellular networks. This subject is attracting a lot of interest from academia and standards organizations alike, as they study requirements and architectures for 5G.

Another interesting analysis pertaining to IoT connectivity is provided in the fourth article, “Understanding the IoT Connectivity Landscape: A Contemporary M2M Radio Technology Roadmap” by Andreev *et al.* The authors provide useful recommendations for the ongoing 3GPP efforts for LTE and 5G optimizations pertaining to MTC.

Concluding the cellular network articles, the fifth article, “What Can Wireless Cellular Technologies Do about the Upcoming Smart Metering Traffic?” by Nielsen *et al.*, provides a study of the optimization of 3GPP cellular networks to address communication needs of smart utility meters.

Software defined networking, a popular research topic, is also addressed in the sixth article, “Publish/Subscribe-enabled Software Defined Networking for Efficient and Scalable IoT Communications” from Hakiri *et al.* The authors describe an IoT architecture that combines SDN and a data distributed service middleware aiming at improving the network flexibility.

In the seventh article, “Software-Defined Internet of Things for Smart Urban Sensing” by Li *et al.*, the authors present an IoT network architecture, also using SDN, for smart urban sensing. The architecture aims at flexible control and management of physical infrastructures and provides a framework to expedite application development.

Low power wide area networking is emerging as a hot topic providing an alternative to both short range and cellular network communications. In the eighth article, “Low Power Wide Area Machine-to-Machine Networks: Key



Techniques and Prototype” by Xiong *et al.*, the authors succeed in describing the key techniques and describe a prototype of low power wide area networks.

Finally in the ninth article, “Towards Better Horizontal Integration Among IoT Services” by Al-Fuqaha *et al.*, the authors suggest IoT protocols enhancements to support QoS by exploiting techniques such as multicast, intelligent broker queue management, and traffic analytics.

We hope the readers will enjoy this issue and find the articles useful. We would like to thank all the authors for submitting their proposals and the reviewers for their tremendous help in reviewing the articles. Special thanks to Chonggang Wang for his help in structuring this Feature Topic. We also like to express our thanks for the great leadership and help of Sean Moore (previous Editor-in-Chief), Osman S. Gebizlioglu (current Editor-in-Chief) and Glenn Parsons (Editor-in-Chief of the Standards Supplement) who have guided us in this endeavor, and Joseph Milizzo from the Communications Society staff for his support.

## BIOGRAPHIES

OMAR ELLOUMI [M] (omar.elloumi@alcatel-lucent.com) is head of M2M and smart grid standards within Alcatel-Lucent CTO. He is the chair of the oneM2M Technical Plenary. He joined Alcatel-Lucent in 1999 and has held several positions in areas including research, strategy, system architecture, and more recently standards. He holds a Ph.D. degree in computer science, and he has served on several consortia Board of Directors and chaired several standards committees.

JAESEUNG SONG [M] (jssong@sejong.ac.kr) is an assistant professor in the Computer and Information Security Department at Sejong University. He holds the position of oneM2M Test Working Group Chair. Prior to his current position, he worked for NEC Europe Ltd. and LG Electronics in various positions. He received a Ph.D. at Imperial College London in the Department of Computing, United Kingdom. He holds B.S. and M.S. degrees in computer science from Sogang University.

YACINE GHAMRI-DOUDANE [M] (yacine.ghamri@univ-lr.fr) is currently a full professor at the University of La Rochelle (ULR), France. Previously, he received an engineering degree in computer science from the National Institute of Computer Science (INI), Algiers, Algeria, in 1998, an M.S. degree from the National Institute of Applied Sciences (INSA), Lyon, in 1999, and a Ph.D. degree in computer science from the University Pierre & Marie Curie in 2003. He has been an officer for two IEEE ComSoc committees.

VICTOR C. M. LEUNG [S75, M'89, SM'97, F'03] (vleung@ece.ubc.ca) is a professor of electrical and computer engineering and holder of the TELUS Mobility Research Chair at the University of British Columbia. He has co-authored more than 800 technical papers in the area of wireless networks and mobile systems. He is a Fellow of the Royal Society of Canada, the Canadian Academy of Engineering, and the Engineering Institute of Canada.

## CALL FOR PAPERS IEEE COMMUNICATIONS MAGAZINE SEMANTICS FOR ANYTHING-AS-A-SERVICE

### BACKGROUND

Services (including anything-as-a-service) are the buzz in the industry. Networks are morphing to utilize new technologies like Network Functions Virtualization and Software-Defined Networking that are changing the way Services are ordered, configured and monitored. To support the evolving infrastructure, new network and service management platforms need to support standard mechanisms for communication within and across administrative domains. In order to support on-demand, dynamic, configuration and monitoring, both common application programming interfaces (APIs) and a common language that has agreed semantics are required. Standards bodies are using information and data modeling to describe the abstract representations and the detailed structured data needed by the orchestrators and controllers in the ecosystem.

This Feature Topic addresses the standards industries usage and advancements in the area of Information and Data Modeling that support the semantics needed for End-to-End Service Management. Comparing and contrasting the top-down vs. bottom-up approach to API (Application Programming Interface) development is also invited. Solicited topics include (but are not limited to):

- Information modeling
- Data modeling
- Transforming information models to data models
- Service development lifecycle aspects
- End-to-End service management frameworks
- Model driven development
- Modeling tools
- Landscape of YANG models
- Survey of modeling work from industry groups
- Advances needed in network management protocols
- Interaction of Open Source and Traditional Industry Fora and Standards Development Organizations

### SUBMISSIONS

Articles should be tutorial in nature and written in a style comprehensible to readers outside the specialty of the article. Authors must follow the IEEE Communications Magazine's guidelines for preparation of the manuscript. Complete guidelines for prospective authors can be found at <http://www.comsoc.org/commag/paper-submission-guidelines>. It is very important to note that the IEEE Communications Magazine strongly limits mathematical content, and the number of figures and tables. Paper length should not exceed 4,500 words. All articles to be considered for publication must be submitted through the IEEE Manuscript Central site (<http://mc.manuscriptcentral.com/commag-ieee>) by the deadline. Submit articles to the “March 2016/Semantics for Anything as a Service” category.

### SCHEDULE FOR SUBMISSIONS

- Manuscript Submission Due: September 15, 2015
- Decision Notification: November 15, 2015
- Final Manuscript Due: January 1, 2016
- Publication Date: March 2016

### GUEST EDITORS

Scott Mansfield  
Ericsson Inc.  
scott.mansfield@ericsson.com

Hing-Kam Lam  
Alcatel-Lucent  
kam.lam@alcatel-lucent.com

Nigel Davis  
Ciena  
ndavis@ciena.com

Yuji Tochio  
Fujitsu  
tochio@jp.fujitsu.com

---

# MACHINE-TYPE COMMUNICATIONS: CURRENT STATUS AND FUTURE PERSPECTIVES TOWARD 5G SYSTEMS

The authors provide a clear mapping between the main MTC service requirements and their associated challenges. Their goal is to develop a comprehensive understanding of these challenges and the potential solutions. This study presents, in part, a roadmap from the current cellular technologies toward fully MTC-capable 5G mobile systems.

*Hamidreza Shariatmadari, Rapeepat Ratasuk, Sassan Iraj, Andrés Laya, Tarik Taleb, Riku Jäntti, and Amitava Ghosh*

---

## ABSTRACT

Machine-type communications (MTC) enables a broad range of applications from mission-critical services to massive deployment of autonomous devices. To spread these applications widely, cellular systems are considered as a potential candidate to provide connectivity for MTC devices. The ubiquitous deployment of these systems reduces network installation cost and provides mobility support. However, based on the service functions, there are key challenges that currently hinder the broad use of cellular systems for MTC. This article provides a clear mapping between the main MTC service requirements and their associated challenges. The goal is to develop a comprehensive understanding of these challenges and the potential solutions. This study presents, in part, a roadmap from the current cellular technologies toward fully MTC-capable 5G mobile systems.

## INTRODUCTION

Machine-type communications (MTC) or machine-to-machine communications (M2M) refer to automated data communications among devices and the underlying data transport infrastructure. The data communications may occur between an MTC device and a server, or directly between two MTC devices [1]. MTC has great potential in a wide range of applications and services. The potential applications are widespread across different industries, including healthcare, logistics, manufacturing, process automation, energy, and utilities.

Communications among MTC devices can be handled through different network technologies. Point-to-point and multi-hop wireless networks, such as ad hoc networks, sensor, and mesh networks have been considered as a means to provide Internet access for devices, forming the so-called Internet of Things (IoT) [2]. For instance, IEEE 802.15.x with its different amendments has been developed to serve a variety of applications in personal area networks [3]. IEEE

802.11ah is another technology that supports low-power transmissions with extended coverage range in Wi-Fi networks [4]. However, these technologies suffer from some fundamental limitations that limit their wide implementation for MTC. The main drawback is the lack of efficient backhaul, which limits network scalability and coverage. Another issue is their operation over unlicensed frequency bands, making the communication links unreliable and susceptible to interference. Therefore, it is challenging to support applications requiring a high degree of reliability. Cellular systems, such as Long Term Evolution (LTE), are considered as alternative solutions for the wide provision of MTC applications. Their ubiquitous presence reduces network installation cost and provides widespread coverage and mobility support. In addition, since cellular systems are regulated and interference controlled, their communication links are more reliable.

Until recently, cellular systems have been mainly designed and optimized to serve traffic from human-to-human (H2H) communications, which are generally characterized by bursts of data during active periods with a higher demand on downlink. However, major MTC applications have different traffic characteristics: usually small and infrequent data generated from a mass of MTC devices imposing a higher traffic volume on uplink. Examples of these infrequent data transmissions, which are uplink-centric, include advanced metering infrastructure and vending machines. Furthermore, MTC devices are also different from

H2H equipment: most MTC devices are inexpensive with limited computational or power resources. These distinct features have raised new technical challenges to enable the widespread deployment of cellular-based MTC [2, 5]. Therefore, these challenges must be effectively addressed for the future broadband wireless communications toward 5G systems in order to fully support MTC. In this vein, the Third Generation Partnership Project (3GPP) has also launched numerous activities to support MTC for future releases of LTE networks, referred to as LTE-Advanced (LTE-A). 3GPP has already specified the general requirements for MTC applications and identified issues and challenges related to them. Network and device modifications have been considered in upcoming releases of LTE standardization to facilitate and better support the integration of MTC [6].

Various methods, use cases, and requirements for 5G systems have been studied lately to support a diverse set of communications. Some of the visions and early results have been summarized in [7]. In this article our goal is to present a thorough overview of the current status of MTC in 4G systems, or more specifically in LTE and LTE-A systems (in the rest of this article, LTE and LTE-A are used interchangeably), and in particular to present perspectives toward 5G mobile systems.

The remainder of this article is organized as follows. We review different MTC service functions and address the main requirements that they impose on cellular systems. Then we

COMMUNICATIONS  
STANDARDS

---

*Hamidreza Shariatmadari, Sassan Iraj, Tarik Taleb, and Riku Jäntti are with Aalto University.*

*Rapeepat Ratasuk and Amitava Ghosh are with Nokia Networks.*

*Andrés Laya is with KTH Royal Institute of Technology.*

describe the current and envisioned cellular-based MTC network architectures. We provide some potential solutions and enhancements for meeting the requirements. Finally, conclusions are drawn.

## MTC SERVICE FUNCTIONS AND THEIR REQUIREMENTS

With MTC, a diverse range of new services and applications can be offered. The potential MTC applications have very different features and requirements that imply constraints on the network technology as well as on MTC devices. Table 1 provides some notable MTC service functions and application examples, including their imposed requirements on cellular systems to be served as radio access technologies.

Metering applications facilitate automatic collection of utility measurements. The gathered information can be utilized for online system optimization and billing purposes. In general, the metering devices generate infrequent and small amounts of data. In addition, the density of metering devices is usually high. These features require that the network technology be capable of handling small bursts of data from a large number of devices. The metering devices may be deployed in indoor environments; hence, they require enhanced coverage for their connectivity. Further examples of MTC services are control and monitoring systems, which enable remote system control or optimization. Reliable communications are required for most of these systems, while low-latency data transmissions are essential for many real-time control systems. Among others, tracking applications assist in managing fleets, locating assets, and preventing theft of equipment. Large-scale connectivity is necessary to support these applications. Furthermore, MTC devices used for these applications are generally equipped with batteries and are expected to operate for a long period of time without the need to replace the batteries; hence, very low power consumption is vital for their operations. For payment applications, security is the most important concern. Security and public safety services need reliable communications with low-latency in addition to a high level of data security to ensure flawless operations.

## NETWORK ARCHITECTURES FOR MTC

3GPP has defined two kinds of communications for MTC applications: communications between an MTC device and a server, and communications between two MTC devices. The required connectivity for the MTC devices and the servers can be provided by cellular systems. Figure 1 shows a typical MTC architecture underlying the current design of LTE networks, along with proposed enhancements toward 5G mobile systems. This architecture consists of three parts: an MTC device domain, a network domain, and an MTC application domain. The MTC devices can be connected to base stations (e.g. eNodeBs) directly or through MTC gateways (MTCGs). The direct connection requires that the devices directly interact with the eNodeBs and are nec-

Service functions	Application examples	Main requirements
Metering	Electric power, gas, and water metering	Support of a massive number of MTC devices with small data bursts and high coverage
Control systems and monitoring	Industrial and home automation, and real-time control	High mobility and low-latency data transmissions
Tracking	Fleet management and asset tracking	High mobility and low power consumption
Payment	Point of sale and vending machines	High level of security
Security and public safety	Surveillance systems, home security, and access control	High reliability, high security, and low latency

Table 1. Examples of MTC applications and their requirements.

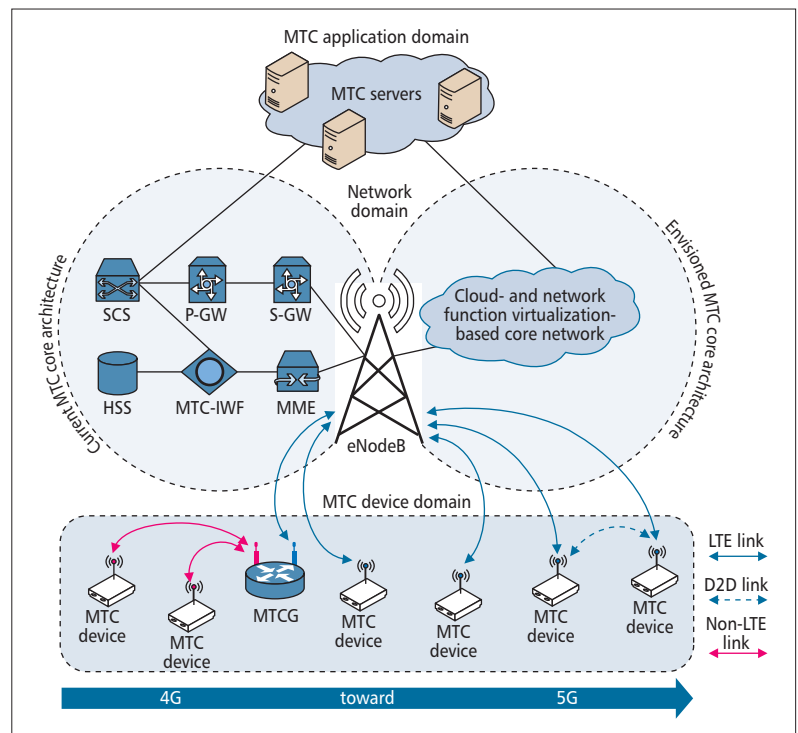


Figure 1. Current and envisioned cellular-based network architectures for supporting MTC services.

essarily compatible with the cellular air infrastructure. The use of MTCGs enables the MTC devices to form capillary networks benefiting from other wired or wireless communication technologies. The MTCGs exchange data between the cellular and capillary networks [8].

The cellular network provides connectivity for MTC devices and MTC servers to exchange data. It comprises two subsystems: a radio access network (RAN) and a core network (CN). The RAN connects the MTC devices to the CN through air interfaces, while the CN manages the overall control of the MTC and the actual data delivery. In LTE systems, data transmissions are performed utilizing the radio resources in a dedicated or shared manner. The available

LTE networks support Internet Protocol version 6 (IPv6), which provides a large address space for devices connected to the network. However, handling a large number of MTC devices simultaneously, required for metering and monitoring applications, causes problems in connection establishment and radio resource allocation.

radio resources are divided into various uplink and downlink physical channels to segregate the different data types. The CN, known as evolved packet core (EPC) in the context of LTE, consists of a number of entities that are used for all types of communications, H2H as well as MTC. The packet data network gateway (PGW) and the serving gateway (SGW) are responsible for forwarding the user data traffic to and from the network via the so called bearers, i.e. channels created with the end users. The mobility management entity (MME) is responsible for all mobility related functions, paging, authentication, and bearer management in the network. The home subscriber server (HSS) functions as a main database containing subscription-related information. Two newly defined entities, service capability servers (SCSs) and the MTC-interworking function (MTC-IWF), are specified for MTC. The SCS entity is mainly designed to offer services for MTC applications hosted in external networks. The MTC-IWF hides the internal public land mobile network (PLMN) topology and relays or translates signaling protocols to invoke specific functionalities in PLMN. This entity is also responsible for relaying trigger requests from the SCS after checking authorization and reporting the acceptance or denial of these requests [5].

The current design of cellular networks can accommodate many MTC applications by meeting their requirements. However, some requirements, stated in Table 1, would need additional considerations before the actual deployment of their relevant applications. Some of these enhancements are illustrated in Figure 1 as part of the ongoing evolution toward 5G mobile systems. An important enhancement for MTC is the support of direct communications between nearby devices, which is known as device-to-device (D2D) communications [7]. Direct communication refers to a radio link establishment between devices without transiting through the network. This feature is already specified in LTE Release 12. Another MTC enabler is network function virtualization (NFV) which allows running network functions on virtual machines instead of utilizing sophisticated and expensive infrastructures. This approach provides virtual instances of the required hardware according to the demand [9]. In addition, the equipment and operation cost of eNodeBs can be further reduced by moving higher layer functionalities to a network cloud. These enhancements have been partially considered for network enhancements in future LTE releases, and for 5G systems. The next section discusses the main challenges associated with the deployment of MTC applications over an LTE network infrastructure and presents some possible remedies.

## CURRENT MTC CHALLENGES AND PROPOSED SOLUTIONS

As mentioned earlier, the adoption of MTC devices into LTE systems has introduced new challenges. The main consideration in addressing the challenges is to simplify the complexity of the network to support MTC devices without

jeopardizing the security or quality of services for both MTC devices and H2H users that share the same network infrastructure. This section describes the challenges listed in Table 1 and the corresponding solutions under discussion to:

- Support the deployment of a massive number of MTC devices.
- Accommodate small data bursts.
- Ensure a high level of security.
- Provide ultra-reliable communications with low-latency.
- Achieve low power consumption.
- Support low cost devices.
- Enhance coverage.

### SUPPORT A MASSIVE NUMBER OF MTC DEVICES

LTE networks support Internet Protocol version 6 (IPv6), which provides a large address space for devices connected to the network. However, handling a large number of MTC devices simultaneously, required for metering and monitoring applications, causes problems in connection establishment and radio resource allocation.

In LTE, each eNodeB hosts radio resource control (RRC) to maintain the RRC state and perform radio resource allocation for all active users. A device can be either in RRC\_IDLE or RRC\_CONNECTED states. In the former state, the device is not connected to any eNodeB and is not granted radio resources for data transmissions. The power consumption of the device is low in this state as the radio transceiver is mostly off. However, the device must transit to the RRC\_CONNECTED state to be able to communicate with an eNodeB. Transiting to this state is initiated by a random access procedure, sending a random preamble over a shared physical random access channel (PRACH) [10]. Performing the random access procedure simultaneously by a large number of MTC devices may congest and degrade the performance of the channel. The system overload is not only on the RAN and CN, but also on the interconnecting nodes such as gateways. This results in undesirable delays and waste of radio resources in the network, and lower performance for both MTC and H2H communications.

*Solutions:* In order to avoid congestion in PRACH, the offered traffic to the channel should be proportional to the allocated resources. Various schemes have been considered to improve the performance of the PRACH and resource allocation to handle a large number of devices [10]. These schemes adopt some of the following approaches.

- Defining new PRACH resources, particularly for MTC, to avoid congestion for H2H.
- Dynamic PRACH resource allocation to adjust available resources based on traffic conditions.
- Performing priority based channel access, granting access to a user with the highest priority among all users requesting access. The desired quality for H2H can be guaranteed by associating higher priorities to them compared to MTC devices [3].
- Access class barring, enabling eNodeB to control traffic on PRACH by setting a barring factor for non-time-critical MTC devices.

Depending on their class, MTC devices execute a different back-off timer before attempting to access PRACH.

- A pull-based scheme to page MTC devices that have permission to send their requests on PRACH.

- A group-based communication protocol for MTC application. For instance, instead of having individual bearers established for each MTC device, a single bearer can be created for a group of MTC devices that have certain characteristics in common or carry on the same MTC service. This group of MTC devices is identified by a unique identifier that can be known a priori to the network or is dynamically created by the network, similar in spirit to the work presented in [11]. A bearer, dedicated to a group of MTC devices, is created following the standard procedure when the first MTC device of the group connects to the network. For other MTC devices of the group, they simply are notified of the availability of the bearer immediately after it is created or when they attempt to establish one. The group bearer remains established as long as the actual data is being transferred between the MTC devices and the MTC servers, or following a particular policy that depends on the behavior of the MTC devices and the underlying MTC service. After a predetermined timeout, or if the actual delivery of data from the MTC group stops for a predetermined period of time, the group bearer may be released.

- To address the system overload, there are solutions based on separating the MTC and H2H traffic. One approach is to use NFV for MTC traffic, which enables allocating the required resources on virtual machines to run network functions [9].

### ACCOMMODATE SMALL BURSTS OF DATA

The current design of LTE systems requires that a user perform the connection establishment procedure before sending the information data. This approach adds signaling overhead to the information data. The efficiency of this scheme is low for handling small amounts of data, which is the case for some MTC applications, because that amount of signaling overhead is high compared to the amount of information data. As an example, Figure 2 illustrates the network access procedure and signaling messages for transmitting 100 bytes of data from a device to a serving eNodeB. The access procedure includes a random access procedure followed by RRC connection establishment and security procedures. When the RRC connection is established, the device transmits the information data and then releases the connection [12]. The signaling transmissions for this process take approximately 59 bytes of overhead on the uplink and 136 bytes on the downlink.

*Solutions:* Increasing the efficiency of data transmissions helps in reducing power consumption and latency. The following solutions are considered to achieve this goal.

- The network access procedure for transmission of small amount of data can be redesigned. For example, a portion of radio resources in the uplink can be dedicated to a group of MTC

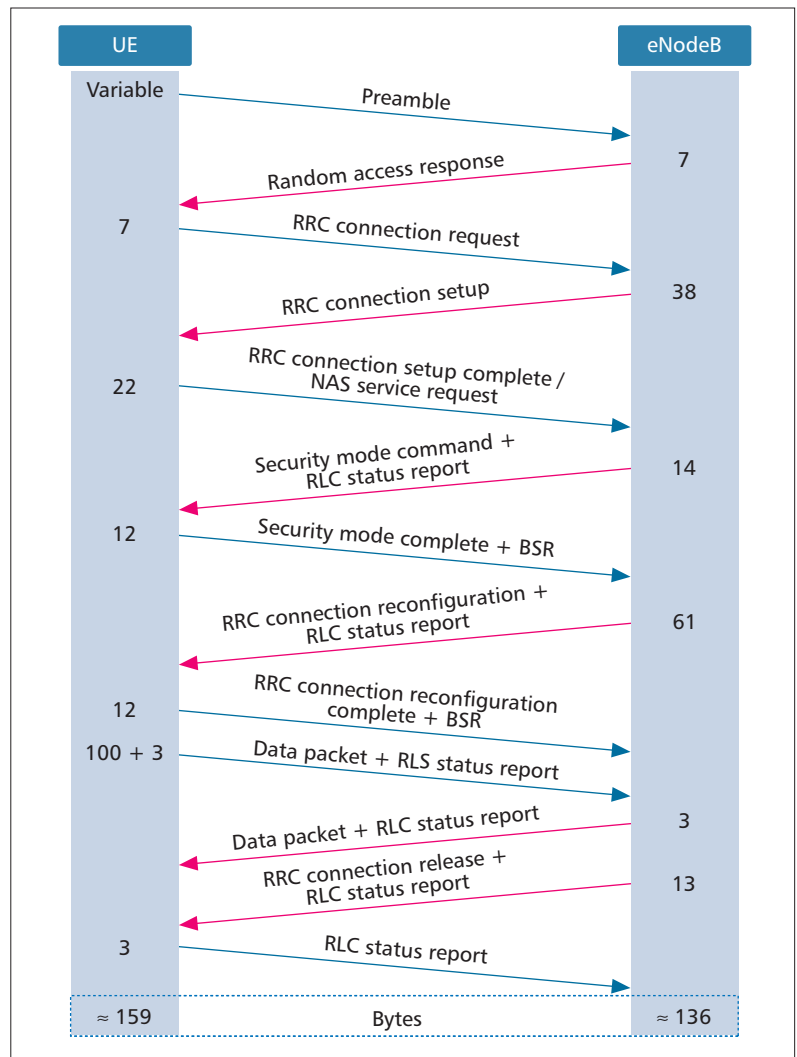


Figure 2. Network access procedure for data transmissions in LTE system.

devices to transmit their data in a contention manner without establishing a link in advance [13].

- There are solutions based on modifications of the random access procedure to handle small data bursts from detached MTC devices. The small data bursts can be carried either by implementing predetermined preambles dedicated to this purpose, or by sending the data load in the initial uplink resource allocated for RRC connection requests. However, sending the data along with RRC connection requests has security implications.

- Data aggregation can improve the efficiency of data transmissions. Data or signaling message aggregation may occur at different locations in the network (e.g. MTC device, MTC gateway, eNodeB, or MME) [11, 14]. Intuitively, this incurs some additional delays and is only applicable to non delay-sensitive MTC applications.

### ENSURE HIGH LEVEL OF SECURITY

LTE systems provide security for communications by integrating various security algorithms, such as authentication, integrity, and encryption. In addition to security for communications, most of the envisioned MTC applications require security for data. This concern is mainly due to

Using asymmetric security schemes, the burden of required computations can be moved to the network domain or gateways, which usually have high computation power. This helps in performing the security computations faster, and reduces energy consumption in MTC devices.

the deployment of MTC devices in unprotected environments, which increases the risk of tampering or fraudulent modification. For instance, moving sensor nodes can disrupt the proper functionality of the MTC applications. In addition, low-cost MTC devices may not be able to perform existing security schemes since they have limited computation power.

*Solutions:* Some solutions for the security relevant challenges are as follows.

- A higher level of security for MTC devices is achievable by utilizing new security mechanisms. For example, the operation of MTC applications can be restricted only to authenticated devices. Also, embedding subscriber identity module (SIM) cards reduces the risk of fraud and SIM theft.

- Employing physical-layer security adopting radio-frequency (RF) fingerprinting. An MTC server should monitor the signal emission characteristics of devices and exploit this information to detect abnormal activities [15]. If the signal from an MTC device changes significantly, it can be deduced that the device has been tampered with or moved to another place. The MTC server then might deactivate the account for the MTC device. However, an RF fingerprinting scheme requires that signal measurements are delivered to the MTC server from the network. Furthermore, variations in the environment and device aging can affect the performance of the scheme.

- Using asymmetric security schemes, the burden of required computations can be moved to the network domain or gateways, which usually have high computation power. This helps in performing the security computations faster, and reduces energy consumption in MTC devices.

### PROVIDE ULTRA-RELIABLE COMMUNICATIONS WITH LOW-LATENCY

Ultra-reliable communications are vital for safe operation of some MTC applications, such as control and monitoring systems, cloud-based systems, and vehicle-to-vehicle wireless coordination. Additionally, some of the mentioned applications also require data transmissions with low-latency. LTE systems currently do not support these features.

*Solutions:* The following solutions are considered to provide ultra-reliable communications with low-latency.

- In order to bypass the link establishment procedure and reduce its associated delay, semi-persistent scheduling can be utilized. This scheme can provide low-delay connection for real-time control applications by periodically providing dedicated radio resources for MTC devices. The semi-persistent scheduling can also eliminate scheduling delays in feedback control applications, as sensors and controllers usually produce data periodically.

- In some applications, reliable communications are required mainly for exchanging data between MTC devices that are located in close proximity. A network enhancement for such applications is to support D2D communications

between nearby devices, which can significantly reduce transmission latency and improve link quality. D2D communications have been studied in LTE Release 12, referred to as proximity services (ProSe), which is one of the areas of 3GPP LTE enhancements to address public safety applications. ProSe allows the identification of mobile devices in physical proximity and enables optimized communications between them. It generally includes two main elements:

- Network assisted discovery of devices that wish to communicate and are located in close physical proximity.

- The establishment of direct communication between such devices with or without supervision from the network.

- Low-latency communications cannot use the classical coding methods that rely on long code words. Instead, they need techniques for handling short packets using new coding methods designed for finite block length. Thus, the coded symbols need to be transmitted within the given time budget using more bandwidth or spatial dimension of freedom.

- Massive multi-input multi-output (MIMO) antennas can provide reliable links by benefiting from spatial diversity and the mitigating effects of fast fading, beamforming, and zero forcing caused by multi-user interference.

- A new medium access control (MAC) scheme is needed for handling event-based ultra-reliable communications. Here, joint coding of the metadata (header) and information payload can help. The challenge is to dimension the random access resources in such a way that delay bounds can be met with high probability.

- The most challenging cases for ultra-reliable communications are safety critical systems, which need low latency and also require high availability of the infrastructure. In some cases, D2D communications and ad hoc networking could be utilized to provide additional links and maintain connectivity when the primary link fails.

### ACHIEVE LOW POWER CONSUMPTION

Energy-efficiency is a decisive metric for choosing a radio technology, particularly for MTC applications whereby MTC devices have a limited energy budget. This is a common constraint for many sensing and monitoring applications whereby MTC devices are expected to operate in scattered areas without the possibility of regular battery replacement. Deploying such MTC devices in LTE networks requires that the network support low-power consumption mode.

*Solutions:* In order to accommodate low power consumption devices in LTE systems, different modifications have been proposed, including the following.

- Modifying signaling and MAC protocols, which can boost energy efficiency by reducing the time that a radio should be turned on. However, the modifications should not greatly affect network performance. For instance, disabling tracking area updates for static MTC devices is a way to reduce signaling transmissions.

- Aggregating the control information and sending them only when data transmission is scheduled.

- Allowing longer discontinuous reception (DRX) sleeping periods.

- Deploying MTC gateways and supporting multi-hop communications can reduce power consumption by allowing MTC machines to transmit with lower power.

- Defining a new communication state for MTC devices. Indeed, LTE has been designed with the assumption that user equipments (UEs) have to be always connected to the network. This “always on” concept gave rise to two main states for UEs: active and idle. During idle state, UEs still have to carry on some control procedures to keep their respective bearers and to update the network of their points of presence. In the case of MTC services whereby MTC devices are triggered only at specific points in time (e.g. the end of the month for utility meters), these MTC devices do not have to be always in idle mode. Rather, they should be “off” with the ability to switch to idle or active state very quickly when required. Such an “energy saving” state will exempt MTC devices from carrying on many standard procedures that would otherwise waste their energy budget.

- Supporting D2D communications between nearby MTC devices can reduce power consumption by allowing transmissions with lower power. However, distributed network management is the main concern for supporting D2D links.

### SUPPORT LOW COST DEVICES

Typical LTE devices have been designed to provide broadband services and are therefore overdesigned for low-rate and delay-tolerant MTC services. For MTC uses, it is desired that LTE devices with bill of materials cost be comparable to that of Global System for Mobile Communications (GSM).

*Solutions:* In Release 12, 3GPP has introduced a new low-cost MTC UE category (called Category-0 UE) with the following reduced capabilities:

- Support of a single antenna instead of at least two receive antennas for other UEs. It is estimated that this could provide a cost reduction of approximately 24 to 29 percent of the bill of materials cost for the modem when compared to the baseline Category-1 UE [16]. However, a single-receive antenna will reduce downlink coverage and the capacity of the system due to the loss of receiver combining gain and lack of channel diversity.

- Reduced peak data rates of one Mbps in both downlink and uplink. The overall achievable cost saving using this technique is 10.5 to 21 percent of the modem electronic bill of materials, and has only a small impact on system efficiency due to small MTC data bursts.

- Use half-duplex operation for frequency division duplex (FDD). In this mode, a UE can only either transmit or receive data at one time (i.e. time-division multiplexing operation). This allows the duplexers to be removed, leading to an estimated cost saving of 7 to 10 percent in the modem. Half-duplex is already an optional feature in LTE, and it is well suited for MTC applications with low data rates.

As a result, Category-0 UE can achieve a cost

saving of approximately 50 percent over the baseline LTE UE. In Release-13, further device cost reduction will be achieved by considering additional modifications, including:

- Reducing the RF bandwidth of the UE from 20 MHz to 1.4 MHz in both downlink and uplink. This reduces both cost and power consumption. However, the UE will still be able to operate in wideband systems by operating only in a portion of the system bandwidth.

- Reducing the maximum transmit power to allow for an integrated power amplifier implementation. However, this power reduction will reduce uplink coverage.

When considered together, Release 13 low-complexity UEs will be able to provide a cost saving of approximately 75 percent over the baseline LTE UE. In terms of deployment, LTE Release 12 and 13 networks will be able to support both legacy and low-cost devices simultaneously. The network will be informed about the capability of the devices earlier during the access procedure, so that they can be treated appropriately. However, it is important to note that low-cost MTC devices will not be able to access legacy (pre-Release 12) networks.

Furthermore, as noted earlier, low-cost MTC devices will have decreased performance due to reduced capabilities. For instance, in the downlink, it is estimated that the performance of low-cost MTC devices will be worse by up to 5 dB due to the single receive antenna deployment and lack of frequency diversity. In the uplink, coverage will be smaller for UEs with reduced maximum transmit power. To compensate for these losses, coverage enhancement features will be standardized.

### ENHANCE COVERAGE

Coverage is an important issue for wide area MTC deployment where machines are installed in challenging locations (e.g. indoor environments, basements, or meter closet). In addition, as discussed earlier, low-cost MTC UE will decrease coverage due to reduced capabilities, such as single receive antenna and lower power. Thus, network access must be extended to ensure ubiquitous coverage throughout the service area. Coverage enhancement features for MTC will be specified in Release 13. The goal is to provide coverage improvement compared to Release 8. The amount of coverage enhancement in each cell is configurable and scalable. Thus, it would be up to the network to configure the level of coverage enhancement.

Coverage analysis via link budget calculation is performed to determine the normal LTE footprint as shown in Table 2. This footprint is determined by the channel with the worst coverage as determined by the maximum coupling loss (MCL). In Release-13, the target is to extend the coverage to 155.7 dB MCL. From Table 2 it can be observed that different amounts of coverage enhancement will be required for different channels.

*Solutions:* The potential coverage improvement techniques applicable to different physical channels are listed in Table 3 and correspond to the following.

Coverage enhancement features for MTC will be specified in Release 13. The goal is to provide coverage improvement compared to Release 8. The amount of coverage enhancement in each cell is configurable and scalable. Thus, it would be up to the network to configure the level of coverage enhancement.

LTE modes	Physical uplink control channel (PUCCH)	Physical random access channel (PRACH)	Physical uplink shared channel (PUSCH)	Physical downlink shared channel (PDSCH)	Physical broadcast channel (PBCH)	Synchronization channel (SCH)	Physical downlink control channel (PDCCH)
FDD 2Tx-2Rx	147.2 dB	141.7 dB	140.7 dB	146.1 dB	149.0 dB	149.3 dB	145.4 dB
TDD 8Tx-8Rx	149.4 dB	146.7 dB	147.4 dB	146.9 dB	149.0 dB	149.3 dB	148.1 dB

Table 2. Maximum coupling loss for LTE channels.

Coverage enhancement techniques	PUCCH	PRACH	PUSCH	PDSCH	PBCH	SCH	PDCCH
Repetition/subframe bundling	✓	✓	✓	✓	✓		✓
Frequency hopping	✓	✓	✓	✓			✓
PSD boosting				✓	✓	✓	✓
Receiver based techniques			✓	✓	✓		
Retransmission			✓	✓			
Relaxed requirement		✓				✓	
Increased RS density	✓		✓	✓	✓		✓

Table 3. Potential coverage enhancement techniques.

- Repetition or subframe bundling. Longer transmission time allows additional energy to be accumulated at the receiver, thus improving the strength of the desired signals.

- Frequency hopping across wideband systems to provide frequency diversity gain to narrowband UEs.

- Power spectral density (PSD) boosting and receive-based techniques such as multi-subframe channel estimation, multiple decoding attempts, and advanced receiver.

- Retransmission using hybrid automatic repeat request (HARQ).

- Relaxing the current specification requirements, such as reducing the missed probability of the random access transmissions.

- Increasing reference signal density to improve channel estimation performance.

To improve performance under coverage enhancement, several techniques may be jointly utilized. As an example, a smart meter installed in a basement may experience penetration loss of 10 dB to 15 dB. To ensure the coverage for the meter, at a maximum coupling loss of 155.7 dB, the downlink data channel can be transmitted employing 40 repetitions, frequency hopping, and 3 dB PSD boosting. This assures that the smart meter receives the data reliably.

## CONCLUSION

In this article we have presented some of the most important requirements imposed by specific MTC applications. Such requirements cause tech-

nical challenges when incorporating MTC in cellular systems. Throughout the description of the challenges, we elaborated on the feasible solutions and their status in terms of maturity and integration to the technical specifications. Even though there has been extensive progress in terms of standardization regarding MTC, currently only some of these challenges have been partially addressed. Hence, further research and exploration are essential to move from current cellular systems toward fully MTC-enabled 5G networks.

## REFERENCES

- [1] R. Ratasuk *et al.*, "Performance of Low-Cost LTE Devices for Advanced Metering Infrastructure," *Proc. IEEE VTC*, June 2013, pp. 1-5.
- [2] K. Chen and S. Lien, "Machine-to-Machine Communications: Technologies and Challenges," *Ad Hoc Networks*, vol. 18, July 2014, pp. 3-23.
- [3] K. Zheng *et al.*, "Challenges of Massive Access in Highly Dense LTE-Advanced Networks with Machine-to-Machine Communications," *IEEE Wireless Commun.*, June 2014, pp. 12-18.
- [4] W. Sun *et al.*, "Wi-Fi Could Be Much More," *IEEE Commun. Mag.*, vol. 52, no. 11, Nov. 2014, pp. 22-29.
- [5] T. Taleb and A. Kunz, "Machine Type Communications in 3GPP Networks: Potential, Challenges, and Solutions," *IEEE Commun. Mag.*, vol. 50, no. 3, March 2012, pp. 178-84.
- [6] 3GPP TS 22.368, "Service Requirements for Machine-Type Communications (MTC)," V13.1.0, Dec. 2014.
- [7] A. Osseiran *et al.*, "Scenarios for 5G Mobile and Wireless Communications: The Vision of the METIS Project," *IEEE Commun. Mag.*, vol. 52, no. 5, May 2014, pp. 26-35.
- [8] K. Zheng *et al.*, "Radio Resource Allocation in LTE-Advanced Cellular Networks with M2M Communications," *IEEE Commun. Mag.*, vol. 50, no. 7, July 2012, pp. 184-92.
- [9] T. Taleb *et al.*, "Lightweight Mobile Core Networks for Machine Type Communications," *IEEE Access*, vol. 2, Sept. 2014, pp. 1128-37.
- [10] A. Laya *et al.*, "Is the Random Access Channel of LTE and LTE-A Suitable for M2M Communications? A Survey of Alternatives," *IEEE Commun. Surveys Tutorials*, vol. 16, no. 1, Dec. 2014, pp. 4-16.



- [11] T. Taleb and A. Ksentini, "On Alleviating MTC Overload in EPS," *Ad Hoc Networks*, vol. 18, July 2014, pp. 24–39.
- [12] M. Hasan *et al.*, "Random Access for Machine-to-Machine Communication in LTE-Advanced Networks: Issues and Approaches," *IEEE Commun. Mag.*, vol. 51, no. 6, June 2013, pp. 86–3.
- [13] S. Andreev *et al.*, "Efficient Small Data Access for Machine-Type Communications in LTE," *Proc. IEEE ICC*, June 2013, pp. 3569–74.
- [14] K. Zhou and N. Nikaein, "Packet Aggregation for Machine Type Communications in LTE with Random Access Channel," *Proc. IEEE WCNC*, Apr. 2013, pp. 262–67.
- [15] D. Reising *et al.*, "Improving Intra-Cellular Security Using Air Monitoring with RF Fingerprints," *Proc. IEEE WCNC*, Apr. 2010, pp. 1–6.
- [16] 3GPP TR 36.888, "Study on Provision of Low-Cost MTC UEs based on LTE," V 12.0.0, June 2013.

## BIOGRAPHIES

HAMIDREZA SHARIATMADARI (hamidreza.shariatmadari@aalto.fi) received the B.Sc. degree in electrical engineering from the University of Tabriz in 2009, and the M.Sc. degree (with distinction) in communications engineering from Aalto University in 2013. He is currently pursuing a doctorate degree at Aalto University in the Department of Electrical Engineering. His current research interests are tied to the development of machine-to-machine communications over wireless systems.

RAPEEPAT RATASUK (rapeepat.ratasuk@nokia.com) received the Ph.D. degree in electrical engineering from Northwestern University in 2000. He is currently a principal research specialist with Nokia Networks. He has extensive experience in cellular system design and analysis, including algorithm development, performance analysis and validation, physical-layer modeling, and simulations. He is a co-author of the book *Essentials of LTE and LTE-A*. His current research interests are in the areas of 5G wireless networks and machine-to-machine communications.

SASSAN IRAJI [M'00, SM'05] (sassan.iraji@aalto.fi) received the Ph.D. degree from Tampere University of Technology, Finland, the M.Sc. (with distinction) from Helsinki University of Technology, Finland, and the B.Sc. from Tehran University, Iran, all in electrical engineering, in 2005, 1999, and 1992, respectively. He worked for Nokia from 1997 to 2012, holding various positions from senior research engineer to senior technology manager, principal researcher, and research leader. He was the founder of the Internet-of-Things team at the Nokia Research Center in 2007. He is currently with Aalto University as a research scientist, where he teaches, manages research projects, and guides graduate students. He has numerous patents and papers published in the field of wireless communications. His current research interests fall within the field of wireless communication systems toward 5G, and in particular machine-to-machine communications and ultra-reliable communications.

ANDRES LAYA (laya@kth.se) is a Ph.D. student in the Communication Systems Department at KTH Royal Institute of Technology. He received his M.Sc. degree in information and communication technologies in September 2009 from UPC-BarcelonaTECH. His research objective is to understand the technical and business implications of services based on connected devices. His research interests are in the areas of radio technologies and micro-economic aspects related to machine-to-machine communications as a key enabler for the Internet-of-Things.

TARIK TALEB [S'04, M'05, SM'10] (tarik.taleb@aalto.fi) received the B.E. degree in information engineering and the M.Sc. and Ph.D. degrees in information sciences from Tohoku University, Sendai, Japan, in 2001, 2003, and 2005, respectively. He is currently a professor at the School of Electrical Engineering, Aalto University, Finland. He has been a senior researcher and 3GPP standardization expert with NEC Europe Ltd., Heidelberg, Germany. He previously worked as assistant professor at the Graduate School of Information Sciences, Tohoku University, Sendai, Japan. His current research interests include architectural enhancements to mobile core networks, mobile cloud networking, mobile multimedia streaming, and social media networking. He has also been directly engaged in the development and standardization of the Evolved Packet System as a member of 3GPP's System Architecture working group. He has received many awards, including the IEEE ComSoc Asia Pacific Best Young Researcher Award in June 2009, and some of his research work has also received best paper awards at prestigious conferences.

RIKU JÄNTTI [M'02, SM'07] (riku.jantti@aalto.fi) is an associate professor (tenured) in Communications Engineering and the head of the Department of Communications and Networking at Aalto University's School of Electrical Engineering, Finland. He received his M.Sc. (with distinction) in electrical engineering in 1997, and the D.Sc. (with distinction) in automation and systems technology in 2001, both from Helsinki University of Technology (TKK). Prior to joining Aalto (formerly known as TKK) in August 2006, he was professor pro tem in the Department of Computer Science, University of Vaasa. He is an associate editor of *IEEE Transactions on Vehicular Technology*. The research interests of Prof. Jäntti include radio resource control and optimization for machine type communications, cloud based radio access networks, spectrum and co-existence management, and RF inference.

AMITABHA (AMITAVA) GHOSH [F] (amitava.ghosh@nokia.com) joined Motorola in 1990 after receiving his Ph.D. in electrical engineering from Southern Methodist University, Dallas. Since joining Motorola he has worked on multiple wireless technologies starting with IS-95, cdma-2000, 1xEV-DV/1XTREME, 1xEV-DO, UMTS, HSPA, 802.16e/WiMAX/802.16m, Enhanced EDGE, and 3GPP LTE. He has 60 issued patents and numerous external and internal technical papers. Currently he leads North America Radio Systems Research within the Technology and Innovation office of Nokia Networks. He is currently working on 3GPP LTE-Advanced and 5G technologies. His research interests are in the area of digital communications, signal processing and wireless communications. He is a co-author of the book titled *Essentials of LTE and LTE-A*.

---

# NEXT GENERATION M2M CELLULAR NETWORKS: CHALLENGES AND PRACTICAL CONSIDERATIONS

The authors present the major challenges of future machine-to-machine cellular networks such as spectrum scarcity, and support for a large number of low-power, low-cost devices. As an integral part of the future Internet-of-Things (IoT), the true vision of M2M communications cannot be reached with conventional solutions that are typically cost inefficient.

---

*Abdelmohsen Ali, Walaa Hamouda, and Murat Uysal*

---

## ABSTRACT

In this article we present the major challenges of future machine-to-machine (M2M) cellular networks such as spectrum scarcity, and support for a large number of low-power, low-cost devices. As an integral part of the future Internet-of-Things (IoT), the true vision of M2M communications cannot be reached with conventional solutions that are typically cost inefficient. The cognitive radio concept has emerged to address spectrum under-utilization and scarcity. The heterogeneous network model is another alternative to relax the number of covered users. To this extent, we present a complete fundamental understanding and the engineering details of cognitive radios, the heterogeneous network model, and power and cost challenges in the context of future M2M cellular networks.

## INTRODUCTION

Internet technology has undergone enormous changes since its early stages and it has become an important communication infrastructure targeting anywhere, anytime connectivity. Historically, human-to-human (H2H) communication, mainly voice communication, has been the primary focus. Therefore, current network protocols and infrastructure are optimized for human-oriented traffic characteristics. Lately, an entirely different paradigm of communication has emerged with the inclusion of “machines” in the communications landscape. In that sense, machines/devices that are typically wireless, such as sensors, actuators, and smart meters, are able to communicate with each other, exchanging information and data without human intervention. Since the number of connected devices/machines is expected to surpass the number of human-centric communication devices by tenfold, machine-to-machine (M2M) communication is expected to be a key element in future networks [1]. With the introduction of M2M, the next generation Internet, or the Internet-of-Things (IoT), must be able to connect different objects together whether they belong to humans or not.

The ultimate objective of M2M communications is to construct comprehensive connections among all machines distributed over an extensive coverage area. Recent reports show that the projected number of connected machines/devices in the IoT will reach approximately 50 billion by the year 2020 (Fig. 1). This massive introduction of communicating machines requires planning and applications that have a wide range of requirements and characteristics such as mobility support, reliability, coverage, required data rate, power consumption, hardware complexity, and device cost. Other planning and design issues for M2M communications include the future network architecture, the massive growth in the number of users, and the various device requirements that enable the concept of IoT. In terms of M2M, the future network has to provide machine requirements such as power and cost that are critical aspects of M2M devices. For instance, a set-and-forget type of application in M2M devices such as smart meters requires very long battery life where the device has to operate in an ultra low-power mode. Moreover, the future network should allow for low complexity and low data rate communication technologies that provide low cost devices that promote the large scale of the IoT. The network architecture, therefore, needs to be flexible enough to provide these requirements and more. In this regard, a considerable amount of research

## COMMUNICATIONS STANDARDS

has been directed toward available network technologies such as Zigbee (IEEE 802.15.4) or WiFi (IEEE 802.11b) by interconnecting devices through large heterogeneous networks [2]. Furthermore, solutions for the heterogeneous network architecture (connections, routing, congestion control, energy-efficient transmission, etc.) have been presented to suit the new requirements of M2M communications. However, it is still not clear whether these sophisticated solutions can be applied to M2M communications due to constraints on hardware complexity.

With the large coverage and flexible data rates offered by cellular systems, research efforts from industry have recently been focused on optimizing the existing cellular networks considering M2M specifications. Among other solutions, scenarios defined by the 3rd Generation Partnership Project (3GPP) standardization body have emerged as the most promising solutions to enable wireless infrastructure of M2M communications [3]. In this area, a special category that supports M2M features has been incorporated by the 3GPP to Long-Term-Evolution (LTE) specifications. Due to the M2M communication challenges and the wide range of supported device specifications, developing the features for M2M communication also refers to machine-type-communication (MTC) in the context of LTE, started as early as release 10 (R10) for the advanced LTE standard. This evolved to future releases including release 13 (R13) that is currently developed and expected to be released in 2016. For these reasons, in this article we will focus on the cellular MTC architecture based on the LTE technology as a key enabler with a wide range of MTC support.

Due to the radical change in the number of users, the network has to carefully utilize the

---

*Abdelmohsen Ali and  
Walaa Hamouda are with  
Concordia University.*

*Murat Uysal is with  
Ozyegin University.*

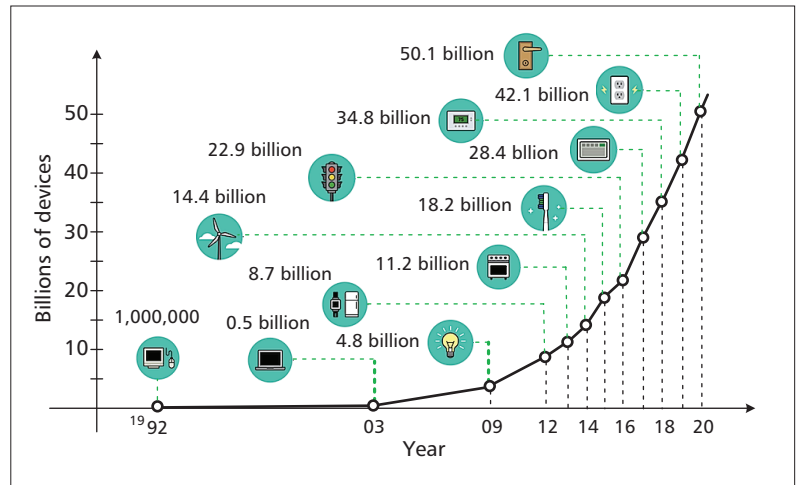
available resources in order to maintain reasonable quality-of-service (QoS). Generally, one of the most important resources in wireless communications is the frequency spectrum. To support a larger number of connected devices in the future IoT, it is likely to add more degrees of freedom represented in more operating frequency bands. However, the frequency spectrum is currently scarce, and requiring additional frequency resources makes the problem of supporting this massive number of devices even harder to solve. In fact, this issue is extremely important, especially for the cellular architecture since the spectrum scarcity problem directly influences the reliability and the QoS offered by the network. To overcome this problem, *small cell design*, *interconnecting the cellular network to other wireless networks*, and *cognitive radio (CR) support* are three promising solutions.

In this article, we address the issues that facilitate the existence of cellular MTC including the network architecture, the spectrum scarcity problem, and the device requirements. We review different approaches, including *small cell design*, *interconnecting the cellular network to other wireless networks*, and *cognitive radio (CR) support*, based on research efforts and industrial technologies to tackle these issues. Furthermore, we provide a comparison of the potential solutions and the challenges and open issues that require future work to allow for practical development of each solution.

The article is organized as follows. We provide an introduction to cellular MTC as well as the technological scenario of M2M communications based on the available standards. In the context of MTC, a description of the spectrum scarcity problem is discussed. This is followed by a description of the cognitive radio solution to solve this problem. We also present the cellular heterogeneous network concept. Then important open issues and future directions are discussed. Finally, we draw our conclusions.

## MACHINE-TYPE-COMMUNICATION IN LTE TECHNOLOGY

Current M2M markets are highly fragmented and most vertical M2M solutions have been designed independently and separately for each application, which inevitably impacts large-scale M2M deployment [4]. However, when it comes to standardizations, global coverage, cellular network stability and maturity, together with the speed offered by recent cellular networks (LTE rates up to 150 Mbps for mobile objects), render wireless cellular technologies as the best candidate for the implementation of secure and reliable business critical M2M services. Several working groups in radio-access-networks (RAN) contribute very actively to the work on MTC-related optimization for 3GPP LTE networks. From day one, the support for MTC was one of the major concerns for the 3GPP, and the development of a robust MTC design was divided across different releases [5]. Figure 2 shows the development steps and features for MTC in different releases. Since LTE has the ability to support high performance, high throughput devices,



**Figure 1.** Expected number of connected devices to the Internet. This chart is obtained from recent reports developed by both Cisco and Ericsson. The reports discuss the expected growth in the number of connected devices by 2020 due to the introduction of the M2M market.

the objective was to develop high volume, low cost, low complexity, and low throughput user-equipment (UE) LTE-based MTC devices.

From the history of MTC/LTE development, the first generation of a complete feature MTC device has emerged in R12. In this release, the 3GPP committee has defined a new profile, referred to as category 0 or CAT-0, for low-cost MTC operation. Also, a full coverage improvement is guaranteed for all LTE duplex modes. On the other hand, R13 is a future release for LTE-A in which MTC has the main weight of contribution. Its main goal is to further enhance the MTC LTE-based UE beyond R12. The main objectives for the MTC improvements are:

- Supporting ultra low-power, low-cost, and narrow-band UE.
- Enhancing the monitoring of service quality.
- Cooperation with other service delivery platforms represented in only oneM2M [6].

Recall that the main objective of oneM2M is to minimize M2M service layer standards market fragmentation by consolidating currently isolated M2M service layer standards activities and jointly developing global specifications. In fact, seven of the world's leading standards bodies, including the European Telecommunications Standards Institute (ETSI) and the Association of Radio Industries and Businesses (ARIB), have come together to create oneM2M. Although this solution considers some test cases for predefined devices such as smart metering, smart grids, eHealth, and automotive applications, not much attention has been given to the underlying connectivity layer since oneM2M leverages current and future technologies such as LTE networks.

## SMALL CELL VS HETEROGENEOUS NETWORK MODEL

The next generation cellular MTC network has to efficiently interconnect several billion wireless machines to support IoT. The traditional method to support these devices is to employ a well-

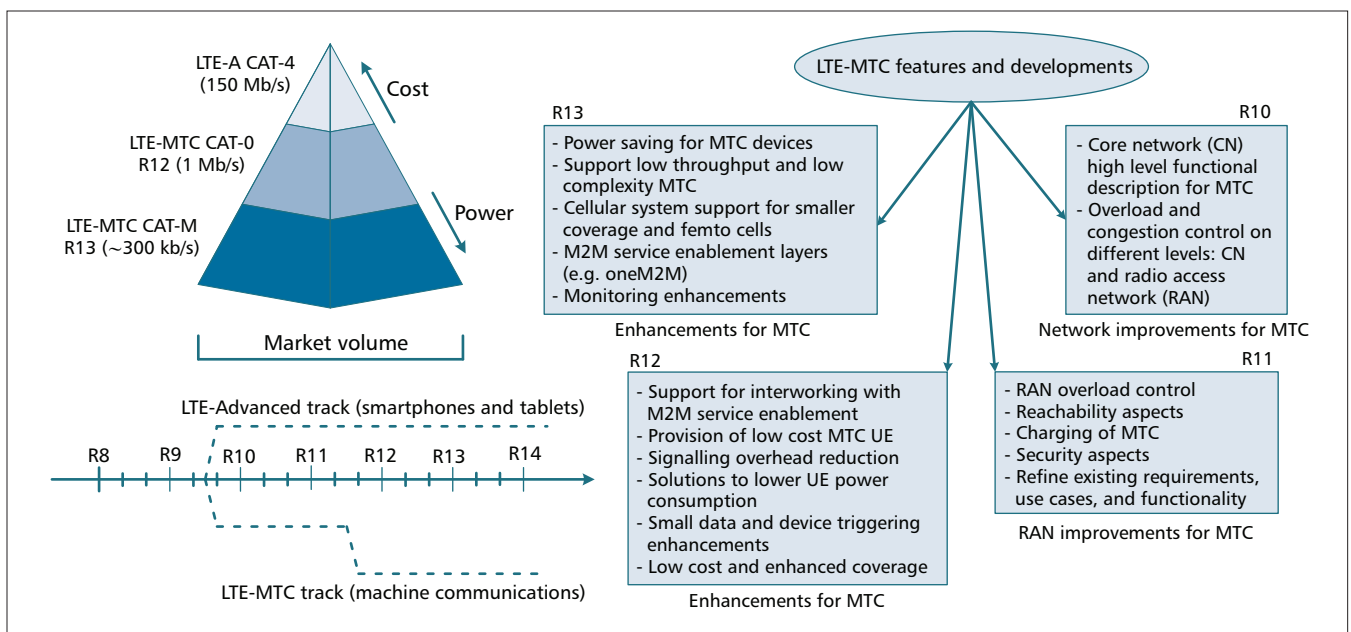


Figure 2. MTC in 3GPP LTE networks: releases and features.

designed M2M technology over a small cell structured system. In this case, cellular network providers need to deploy several thousand base stations (eNodeB in the LTE context) each with a smaller cell radius rather than full-power transmitters with large cells. Of course, this solution is cost-inefficient. Moreover, with such a large number of small cells, co-channel interference is a limiting factor, and complex designs are needed to maintain the required QoS. Another major drawback of this approach is the significant traffic increase due to signalling congestion and network management.

Although the “heterogeneous network” model is not currently recommended for MTC due to the limited capabilities of the machines, research efforts [7] have been invested to support the idea of utilizing the cellular network itself as a small type of a heterogeneous network. The concept is that, in many applications, machines can be clustered geographically where the members of each cluster can be interconnected together through certain technology. To reduce the number of machines connected to the cellular network, each cluster would select a representative, a cluster head, to connect with the cellular network. Inside the cluster, the cellular network is transparent to all machines and only the cluster head will be responsible for relaying the aggregate traffic of the entire cluster. For example, if all machines have WiFi interfaces, then WiFi technology can be utilized to interconnect cluster members. In that sense, the cluster head will be communicating over its WiFi interface inside the cluster while using the LTE interface, for example, to connect to the cellular network (Fig. 3). In this model, the cellular network has offloaded part of its traffic to the individual clusters, and therefore reduces the effective number of covered users. The main benefit of this approach is the relaxation of congestion that would result if no clusters are formed.

## COGNITIVE CELLULAR M2M NETWORKS

The idea of cognitive radios was originally proposed to offer more efficient utilization for the RF spectrum [8]. In this context, there are two approaches to apply the CR concept in cellular M2M networks. The first approach [9] assumes that there can be two types of eNodeB stations, one for typical UEs and other for MTC UEs coexisting with each other (Fig. 3), to relax signalling congestion and management burden. In this case, M2M devices seek to opportunistically use the spectrum when the H2H devices are idle. Therefore, M2M and H2H devices are not allowed to simultaneously operate over H2H links. This can be done through coordination between the corresponding eNodeB stations. Once a radio resource is occupied by M2M communications, this radio resource is regarded as suffering from server interference and will not be utilized by H2H communication. Even though this approach is simple to apply, it can degrade the QoS of H2H applications, especially when the number of MTC devices is very large.

To overcome the aforementioned problems, we propose a second approach that supports unlicensed bands in addition to existing licensed bands. Here it is assumed that the network will sense unlicensed bands to find extra vacant bands. If complexity permits, more than one unlicensed band per cell can be utilized by a smart-eNodeB (S-eNodeB), a coined term to differentiate between the traditional eNodeB and the proposed eNodeB, to further increase the number of devices (Fig. 3). Indeed, this solution leverages the huge amount of free spectrum available around the 5 GHz and TV white space bands. However, current radio access standards such as IEEE 802.22 and IEEE 802.11af already allow the use of this free unlicensed spectrum. Therefore, spectrum sensing and monitoring is a must. This can be implemented by introducing a new layer for spectrum management to support

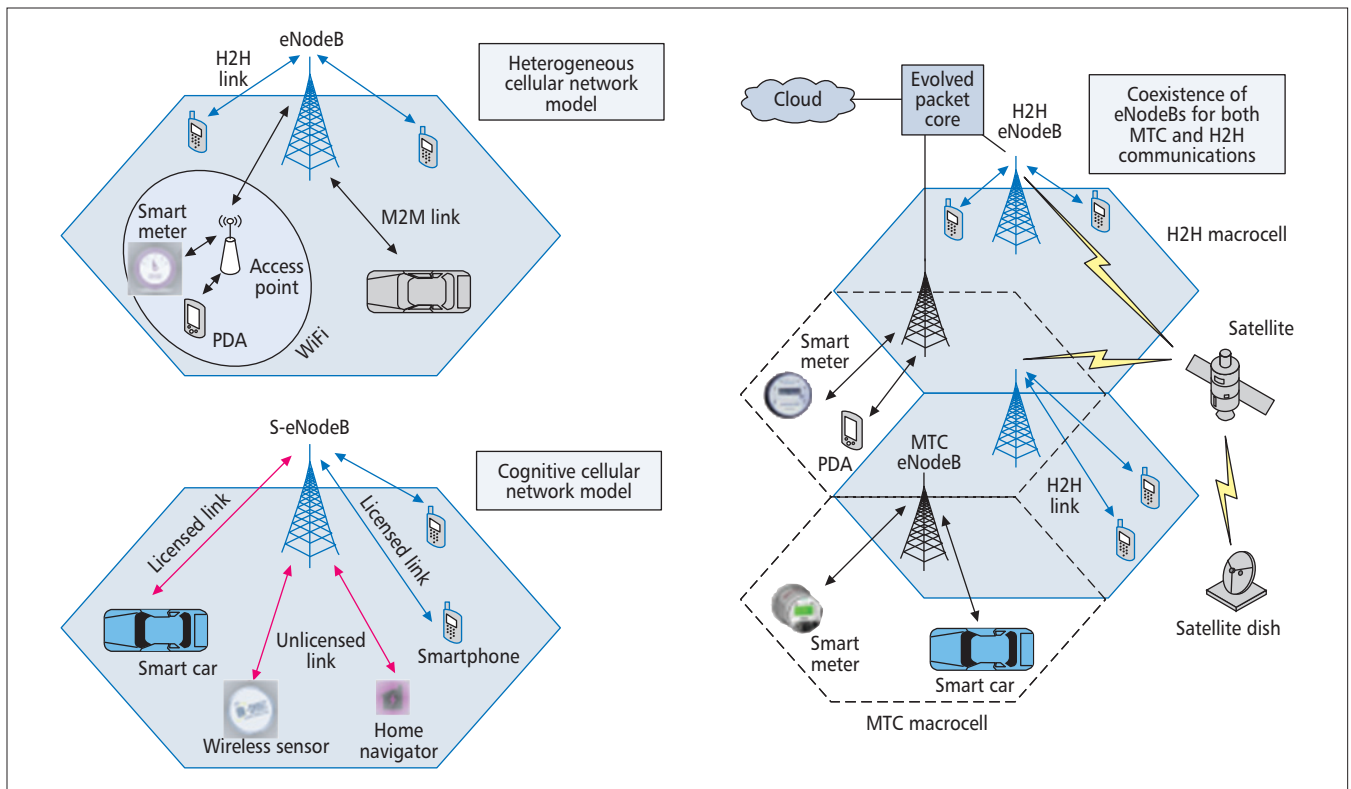


Figure 3. Various network models to interconnect numerous number of machines to IoT.

cognition over the unlicensed bands. That is, the S-eNodeB should be capable of:

- Sensing the spectrum.
- Gathering information about the available suitable bands.
- Making decisions on the conditions of these bands.
- Informing the neighboring S-eNodeBs about the allocated unlicensed band.
- Monitoring the allocated unlicensed band.
- Always providing an alternative band.

If the S-eNodeB handles multiple unlicensed bands, then it should classify the machines based on their performance tolerance so that a machine is switched to the proper unlicensed band that meets its requirement. Of course, this assumes that the machine would have a group ID to declare its needs, which in turn has to be shared with the S-eNodeB during call setup. To clarify how machines and S-eNodeBs can work in this scenario, a detailed call procedure is demonstrated to show how a machine can access the unlicensed band. Once the machine is switched on, it goes to the calibration process in which the RF front-end adjusts or even estimates the IQ mismatch parameters. The following procedure is shown in Fig. 4 and is discussed below.

- The machine would start the usual frequency scanning over the licensed LTE carriers. Once it locates a strong serving cell, a synchronization procedure is followed so that the machine is locked to the base station. It further decodes the master information block to recognize the cell specification.

- The machine sends a random access request to connect to the cell. The S-eNodeB then requests the group ID, which will be sent over the uplink control channel.

- The S-eNodeB will request the machine to switch to another carrier in the unlicensed band. Full information about the carrier, such as modulation, coding, and relative timing to the licensed carrier, are also sent to the machine. Afterward, the S-eNodeB assumes that the session is complete and the machine has been configured.

- The machine will then switch its RF to the desired carrier and enter the synchronization mode to lock itself to the S-eNodeB at the unlicensed carrier.

- The machine defines itself one more time by sending a random access request over this carrier. If it is permitted, the machine can exchange data with the S-eNodeB over the physical uplink and downlink shared channels.

- The S-eNodeB can interrupt the machine by scheduling a measurement gap in which the machine measures and reports the power of a certain carrier in the unlicensed or licensed bands.

- The unlicensed carrier can be dynamically changed based on the collected measurements at the S-eNodeB. In this case, machines have to be informed about the new carrier and its settings.

## ULTRA LOW-POWER AND LOW-COST NETWORKS

To save battery life, low-power design is always desired for wireless communication systems. However, power reduction is not an easy task as it is related to system reliability, the rate of data exchange, and the radio chip design and implementation constraints. When the communication link is unreliable, higher layers translate this into retransmissions, which results in longer active times and hence, high power consumption. Similarly, if the system continuously exchanges data

When a number of machines is able to form clusters, the cellular network becomes lightly loaded. This conclusion has been investigated by many researchers and even practically demonstrated on WiFi as the internal technology inside the cluster. However, it is hard to judge if the machines can really form clustering or not.

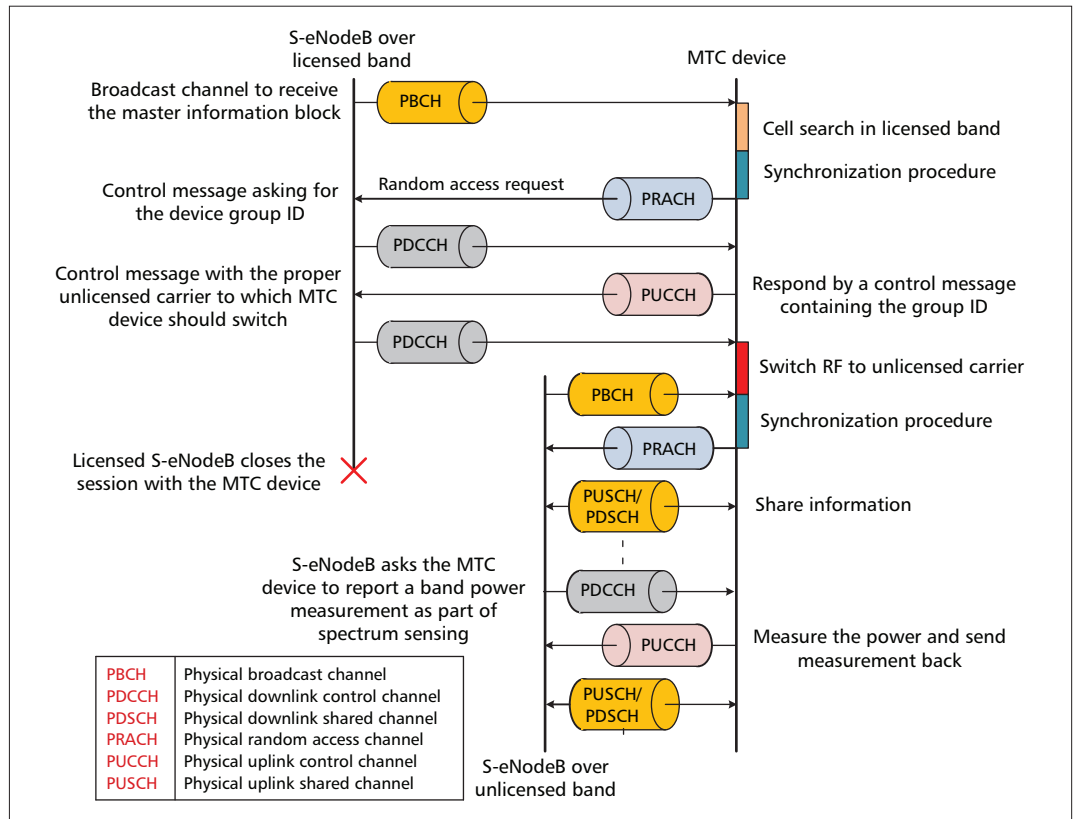


Figure 4. Handshake messaging for MTC device over cognitive cellular network.

then it will consume more power. Based on a case study for ZigBee [4], it is shown that if the radio is switched on all the time, it will deplete a typical AA battery within a week. However, turning the radio duty cycle to 25 percent extends the lifetime to about a month. Turning it further down to 1 percent yields years of lifetime. Therefore, low power can be achieved through a reliable communication link with small duty cycle. In LTE-advanced, the concept of discontinuous reception (DRX) cycles is applied where the eNodeB schedules a silent period (DRX cycle) to encourage the device to switch off the radio chip so that low duty cycle is achieved. To support ultra low-power design in recent releases, a long DRX cycle mode has been employed (with a maximum period of 2.56 sec. in R12) to further reduce the duty cycle.

Another important aspect of future MTC devices is their low-cost design, which is typically provisioned by reducing the complexity of the system while providing the same coverage. The communication system architecture usually involves a general processor to run the software, memory to hold both instructions and data, and a physical-layer modem to handle the communication protocol. As expected, most of the complexity reduction comes from the physical-layer modem features along with a small portion of data memory reduction. Therefore, a low-cost design is typically related to a feature reduction while the coverage is carefully kept unchanged. For specific applications, low data rates and/or low latency are acceptable. In this case, the modem features can be relaxed to target low-cost design. In recent LTE releases, a special category has been defined to support MTC for low data rates which leads to

complexity reduction. In LTE-R12, this category supports only one receive antenna and a maximum data rate of 1 Mbps. However, those features will be further reduced in LTE-R13 with the expected maximum data rate being 300 Kbps and only one operating bandwidth of 1.4 MHz.

## CHALLENGES, OPEN ISSUES, AND FUTURE DIRECTIONS

### HETEROGENEOUS NETWORKS

When a number of machines is able to form clusters, the cellular network becomes lightly loaded. This conclusion has been investigated by many researchers and even practically demonstrated on WiFi as the internal technology inside the cluster. However, it is hard to judge if the machines can really form clustering or not. In fact, clusters are formed only if the WiFi connectivity between cluster members is acceptable (data rates are higher than the LTE load generated in the cluster). Also, clustering allows machines to enjoy seamless connectivity to the cellular system while spending more time on a secondary, WiFi-based interface, which consumes less power than LTE. On the other hand, shifting the responsibility of the aggregate traffic from all cluster members to the cluster head can be challenging, especially if the link from the cluster head to the eNodeB is poor. Since the architecture assumes a centralized control at the head node, it is expected that the full cluster will fail. Therefore, more research effort is required to investigate the possibility of dynamically selecting the head node based on the channel quality with the cellular system. One challenge with this solu-

tion is to select the optimum period after which a rescheduling has to be done.

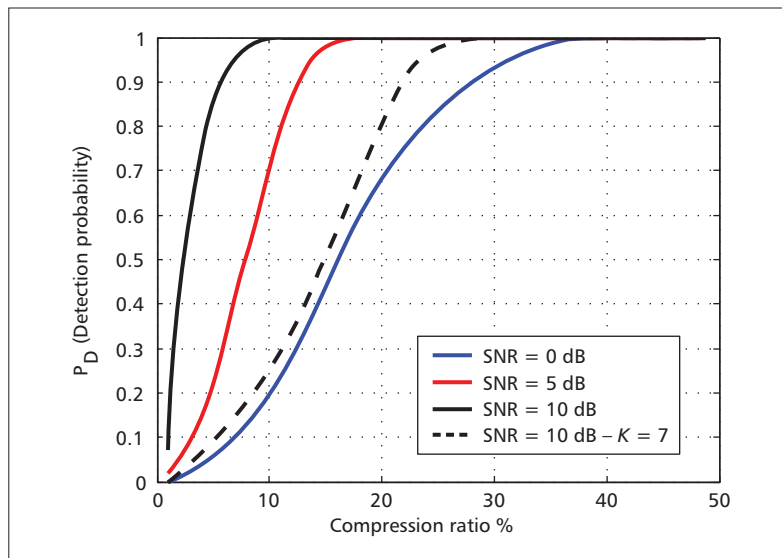
### COGNITIVE RADIO NETWORK

As discussed earlier, spectrum sensing and monitoring are essential to utilize the cognitive radio concept in which some of the machines operate over an unlicensed band. However, there are many challenges to address this problem.

**Spectrum Sensing Techniques:** Sensing can be either centralized at the S-eNodeB or done in cooperation with the machines. Better performance is expected from the latter case since more spatial diversity is utilized. Generally, cooperation is achieved by sending either local decisions [10], which can be either hard or soft decisions, or by sending the useful portion of the received data set. The processing power of the machines limits the first approach, while high traffic over the control channel is the main challenge for the second approach. Moreover, the link between the machine and the S-eNodeB is not ideal and the sensing decisions/data can be received incorrectly, which may alter the sensing accuracy at the S-eNodeB.

**Wideband Sensing Methodology:** During the initial sensing stage, a very wide band (approximately 1 GHz) has to be assessed to locate a suitable vacant band. This can be implemented by scanning different bands one after another and measuring the in-band power. This technique is simple but it requires time and power to find a suitable band. Another alternative is to examine the power spectral density of the entire wide band at once. Since this method requires high speed analog-to-digital conversion, compressive sensing (CS) [10] is a promising technique to obtain the power spectral density of the wideband spectrum while sampling at rates lower than the Nyquist rate. The concept is to capture a few measurements of the *sparse* spectrum. The wideband spectrum is related to those raw measurements by a linear under-determined system of equations. Optimization techniques can be employed to solve this set of equations in order to find the best solution that satisfies the original assumption for the spectrum that is being *sparse*. Fig. 5 shows the detection performance as a function of the ratio of non-uniform sampling frequency to the typical Nyquist rate. It is clear that CS is able to detect spectrum occupancy by a ratio of 1/10 of the Nyquist rate at high signal-to-noise ratio (SNR). Although CS is very promising in this context, many challenges exist due to the current algorithmic complexity as well as the basic assumptions. For example, the spectrum is dynamically loaded and the *sparse* assumption may not be valid, which results in performance degradation ( $K = 4, 7$  cases in Fig. 5). Cooperation may be utilized to enhance the accuracy; however, finding a high-performance low-complex/low-data rate cooperative sensing technique is not a trivial task. More research efforts are needed in that direction to develop efficient algorithms to render CS possible with reasonable complexity, especially for MTC, where complexity is a real challenge.

**Narrowband Sensing Techniques:** A signal processing algorithm is needed to decide on the activities within each of the wideband slices (vacant or not). Conventional algorithms/detectors [11] include the energy detector, the cyclo-



**Figure 5.** Effect of compression ratio (i.e., ratio of the non-uniform sampling rate to the conventional Nyquist rate) on the detection performance when the false alarm rate is 1 percent. 16 contiguous non-overlapped bands are investigated where each has a bandwidth of 1MHz. Only four active bands are considered, therefore the sparsity level  $K = 4$  out of the available 16.

stationary detector, and the matched-filter detector. In all cases, a decision statistic is computed and compared to a threshold to decide whether a specific band is occupied or not. Complexity, performance, and prior information about the signal to be detected are the main metrics to judge the quality of the detector. Among those detectors, the energy detector is known to be the only simple non-coherent detector. From the performance perspective, the matched-filter is known to be the optimal detector. However, it requires full knowledge of the detected signal.

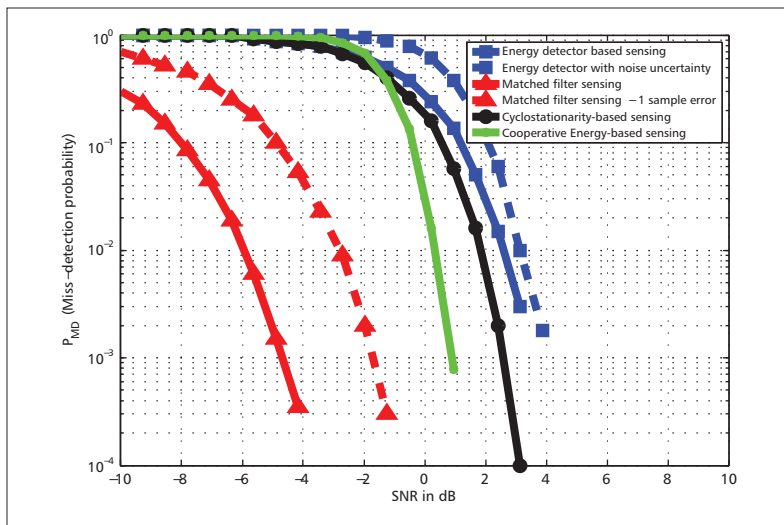
The cyclostationary detector can be used only if the signal possesses the cyclostationarity property where its statistics, mean, and autocorrelation are periodic with some known period. Therefore, it requires partial information about the detected signal, which is typically the period of cyclostationarity. Figure 6 shows the probability of miss-detection for various narrowband sensing techniques against SNR. The effects of timing errors, noise uncertainty, and hard decision cooperative sensing have been included. The performance results show that:

- Any uncertainty of the noise level will significantly alter the performance of the energy detector.
- Matched-filter detection is very sensitive to timing errors.
- Cooperation involves high diversity gain.

However, these results assume an ideal channel (no noise and no fading) between the machines and the S-eNodeB. The conclusion is that, improvements and/or new sensing techniques are needed to provide less-complex, non-coherent, and robust practical algorithms.

### LOW-POWER LOW-COST NETWORKS

Although a longer DRX cycle significantly reduces power consumption, it also introduces some challenges to the system design. Since the radio chip



**Figure 6.** These curves are plotted for a false alarm rate of 1 percent. The window size for the energy detector is the same as the matched filter length. Both agree with the cyclostationary detector period which is 32 samples. For the cooperative sensing, hard decision is used with K-out-of-N rule where  $K = 5$  users and  $N = 10$  users. The noise uncertainty error is 0.5 dBs for the energy detection case.

will be off during the DRX cycle, the device/UE has no way to synchronize itself to the eNodeB. Therefore, the typical behavior for the device/UE would be to wake-up as early as required to quickly resynchronize itself to the eNodeB before receiving further packets. One of the issues is to determine the best wake-up time so that the synchronization performance is met and no additional power is lost. Another issue is related to the cooperative sensing architecture, if applicable, where the device/UE will not be able to sense or monitor any band while it is in a deep sleep mode. The band can suffer from high interference levels caused by other networks that attempt to access the same band. Finally, power consumption can be minimized by properly designing power domains in the hardware to decide which module is not needed to be switched off.

A low-cost design always comes at the cost of system performance with less features provided. For instance, reducing the number of receive antennas from two to one would reduce the spatial diversity of the modem. Therefore, advanced signal processing algorithms for synchronization, cell detection, and decoding will need to be revised to guarantee the same performance with less diversity gain. Indeed, reducing the cost is not only related to the required features from the network, but it also depends on the hardware design process and underlying technology. For example, optimizing the internal word sizes of the various hardware modules inside the modem will result in a low gate count and low power consumption. However, the optimization algorithms that can achieve this are not unique as signal statistics across various modules are system dependent.

## CONCLUSION

We presented the challenges that are expected from the next generation MTC network as an integral part of the future IoT. It is argued that

the cognitive radio concept is a possible solution from the cost and performance perspectives. However, there are more practical challenges that need efforts from researchers. The application of the heterogeneous network concept was investigated where cellular MTC networks can utilize other networks such as WiFi to reduce the number of directly connected machines/ users. Future standards are encouraged to provide both options (i.e. the cognition concept and the heterogeneous network model). Finally, a design of a low-power low-cost machine is discussed. However, there are important design challenges that must be solved to make it possible. For example, an extended DRX cycle is a valid option to significantly reduce power consumption. No matter how, a feasibility study is scheduled in Release 13 to provide a solid solution in which extended DRX cycle implementation challenges can be overcome, if possible. The trade-offs between cost, feasibility, and performance have also been discussed.

## REFERENCES

- [1] J. Wan *et al.*, "M2M Communications for Smart City: An Event-Based Architecture," *IEEE Int'l. Conf. Computer and Information Technology (CIT)*, Oct 2012, pp. 895-900.
- [2] A. Damnjanovic *et al.*, "A Survey on 3GPP Heterogeneous Networks," *IEEE Wireless Commun.*, vol. 18, no. 3, June 2011, pp. 10-21.
- [3] S. Y. Lien, K. C. Chen, and Y. Lin, "Toward Ubiquitous Massive Accesses in 3GPP Machine-to-Machine Communications," *IEEE Commun. Mag.*, vol. 49, no. 4, Apr. 2011, pp. 66-74.
- [4] C. A. Haro and M. Dohler, *Machine-to-Machine (M2M) Communications: Architecture, Performance and Applications (Google eBook)*, Elsevier, Dec 2014.
- [5] 3GPP, "Standardization of Machine-type Communications," 3rd Generation Partnership Project, Tech. Rep. V0.2.4, June 2014.
- [6] J. Swetina *et al.*, "Toward a Standardized Common M2M Service Layer Platform: Introduction to oneM2M," *IEEE Wireless Commun.*, vol. 21, no. 3, June 2014, pp. 20-26.
- [7] A. Pyattaev *et al.*, "3GPP LTE Traffic Offloading onto WiFi Direct," *IEEE Wireless Commun. and Networking Conf. Wksp.*, Apr. 2013, pp. 135-40.
- [8] B. Wang and K. Liu, "Advances in Cognitive Radio Networks: A Survey," *IEEE J. Sel. Topics Signal Process.*, vol. 5, no. 1, Feb 2011, pp. 5-23.
- [9] S. Y. Lien *et al.*, "Radio Resource Management for QoS Guarantees in Cyber-Physical Systems," *IEEE Trans. Parallel and Distributed Systems*, vol. 23, no. 9, Sept 2012, pp. 1752-61.
- [10] I. F. Akyildiz, B. F. Lo, and R. Balakrishnan, "Cooperative Spectrum Sensing in Cognitive Radio Networks: A Survey," *Phys. Commun.*, vol. 4, no. 1, Mar. 2011, pp. 40-62.
- [11] T. Yucek and H. Arslan, "A Survey of Spectrum Sensing Algorithms for Cognitive Radio Applications," *IEEE Commun. Surveys Tutorials*, vol. 11, no. 1, 1st 2009, pp. 116-30.

## BIOGRAPHIES

ABDELMOHSEN ALI (ali\_abde@ece.concordia.ca) received the B.Sc and M.A.Sc. degrees in electronics and communication engineering from Cairo University, Egypt, in 2004 and 2008, respectively. He is currently pursuing the Ph.D. degree at Concordia University, Montreal, Canada. From 2008 to 2012 he was with Wasiela working in the area of digital communications modem design, including DVB-T2 and LTE. His research interests include communication signal processing and VLSI architectures for communication systems.

WALAA HAMOUDA (hamouda@ece.concordia.ca) received the M.A.Sc. and Ph.D. degrees in electrical and computer engineering from Queen's University, Kingston, Canada, in 1998 and 2002, respectively. Since July 2002 he has been with the Department of Electrical and Computer Engineering, Concordia University, Montreal, Canada, where he is currently a full professor. Since June 2006 he has been Concordia University's Research Chair in Communications and Networking. His current research interests include space-time processing, cooperative communications, wireless networks, and cross-layer design.

MURAT UYSAL (murat.uysal@ozyegin.edu.tr) is a full professor at Ozyegin University, Istanbul, Turkey, where he leads the Communication Theory and Technologies (CT&T) Research Group. Prior to joining Ozyegin University, he was a tenured associate professor at the University of Waterloo, Canada, where he still holds an adjunct faculty position. His research interests are in the broad areas of communication theory and signal processing, with a particular emphasis on the physical layer aspects of wireless communication systems in radio, acoustic, and optical frequency bands.



**CALL FOR PAPERS**  
**IEEE COMMUNICATIONS MAGAZINE**  
**COMMUNICATIONS STANDARDS SUPPLEMENT**

**BACKGROUND**

Communications standards enable the global marketplace to offer interoperable products and services at affordable cost. Standards development organizations (SDOs) bring together stakeholders to develop consensus standards for use by a global industry. The importance of standards to the work and careers of communications practitioners has motivated the creation of a new publication on standards that meets the needs of a broad range of individuals, including industrial researchers, industry practitioners, business entrepreneurs, marketing managers, compliance/interoperability specialists, social scientists, regulators, intellectual property managers, and end users. This new publication will be incubated as a Communications Standards Supplement in *IEEE Communications Magazine*, which, if successful, will transition into a full-fledged new magazine. It is a platform for presenting and discussing standards-related topics in the areas of communications, networking, and related disciplines. Contributions are also encouraged from relevant disciplines of computer science, information systems, management, business studies, social sciences, economics, engineering, political science, public policy, sociology, and human factors/usability.

**SCOPE OF CONTRIBUTIONS**

Submissions are solicited on topics related to the areas of communications and networking standards and standardization research in at least the following topical areas:

Analysis of new topic areas for standardization, either enhancements to existing standards or in a new area. The standards activity may be just starting or nearing completion. For example, current topics of interest include:

- 5G radio access
- Wireless LAN
- SDN
- Ethernet
- Media codecs
- Cloud computing

Tutorials on, analysis of, and comparisons of IEEE and non-IEEE standards. For example, possible topics of interest include:

- Optical transport
- Radio access
- Power line carrier

The relationship between innovation and standardization, including, but not limited to:

- Patent policies, intellectual property rights, and antitrust law
- Examples and case studies of different kinds of innovation processes, analytical models of innovation, and new innovation methods

Technology governance aspects of standards focusing on both the socio-economic impact as well as the policies that guide them. These would include, but are not limited to:

- The national, regional, and global impacts of standards on industry, society, and economies
- The processes and organizations for creation and diffusion of standards, including the roles of organizations such as IEEE and IEEE-SA
- National and international policies and regulation for standards
- Standards and developing countries

The history of standardization, including, but not limited to:

- The cultures of different SDOs
- Standards education and its impact
- Corporate standards strategies
- The impact of open source on standards
- The impact of technology development and convergence on standards

Research-to-standards, including standards-oriented research, standards-related research, and research on standards

Compatibility and interoperability, including testing methodologies and certification to standards

Tools and services related to any or all aspects of the standardization life cycle

Proposals are also solicited for Feature Topic issues of the Communications Standards Supplement.

Articles should be submitted to the *IEEE Communications Magazine* submissions site at

**<http://mc.manuscriptcentral.com/commag-ieee>**

Select "Standards Supplement" from the drop-down menu of submission options.

---

# CONNECTIONLESS ACCESS FOR MOBILE CELLULAR NETWORKS

The authors propose a connectionless access method for efficient small burst transmission for future cellular networks. With the proposed protocol for connectionless access, over-the-air signaling can potentially be reduced substantially while meeting the security requirements of mobile networks.

---

Colin Kahn and Harish Viswanathan

---

## ABSTRACT

The traffic characteristics of Internet of Things and machine type devices differ significantly from that of smart phone applications. The connection oriented approach in Long Term Evolution (LTE) that involves establishing a bearer prior to data transmission is inefficient for small burst traffic because of the signaling involved. We propose a connectionless access method for efficient small burst transmission for future cellular networks. With the proposed protocol for connectionless access, over-the-air signaling can potentially be reduced substantially while meeting the security requirements of mobile networks.

## COMMUNICATIONS STANDARDS

## INTRODUCTION

A wireless industry goal for evolved 4G and 5G is to more efficiently support wide area machine type communications (MTC) and thereby enable a highly diverse ecosystem of devices and applications that can be connected anywhere. This is a significant departure from the roots of 3G and 4G with its focus on human communication and more recent support for high-end smart devices. The Internet of things (IoT) and MTC has the potential to eclipse human communication in mobile networks in terms of number of devices and bursty traffic. This traffic and device growth stems from the extremely wide range of vertical applications and the large number of devices each can involve, unconstrained by human interaction.

Examples of verticals include health monitoring, shipping/tracking, connected car, inventory management, and smart cities, to name but a few where the benefits of deploying IoT are well established. The diversity of device types is also vast. Communication capability may be integrated in complex monitoring equipment for high value applications (for example, healthcare and industrial infrastructure), or be embedded in simple sensors for temperature, smoke, or location tracking.

Many of these applications and device types share a set of key attributes that are important to consider in the design of evolved 4G and 5G networks. These include:

- Sporadic transmission of small data bursts [1, 2].

- Large numbers of devices.
- Low power consumption to extend battery life.

Smart phones and tablets also generate small data bursts when applications are operating in the background for reasons such as TCP keep-alive's and notifications, and when only applications such as text messaging are active [3].

The 3GPP forum has studied and proposed incremental modifications to the LTE standard to enhance the network to efficiently support these devices and traffic types. Optimizations have also been proposed for IoT [4, 5], and specifically to improve radio resource management using group based mechanisms [6] and methods targeted for heterogeneous networks[7]. However a dramatic improvement in efficiency requires major changes to the air interface and core network. In LTE, to transmit a short packet the device first has to establish a radio connection to the base station, which involves establishing a signaling bearer, activating a data bearer, and obtaining dedicated resources from a scheduler to transmit the packet. In contrast, we propose a connectionless approach where the device can send the small packet directly on common channel resources without bearer establishment, thereby reducing the amount of signaling that is required to

transmit a small packet. We describe the architecture and associated signaling required for such a connectionless access approach. We also describe how security can be handled for connectionless access, a critical aspect without which any solution will not be considered viable for mobile operator networks.

This article is organized as follows. We first describe the traffic types under consideration along with some example applications. Then we give an overview of the process and associated overhead involved in sending small bursts over 4G networks. We describe our 'connectionless access' architecture and signaling and provide estimates of overhead. We devote a section to discussing security issues. We conclude with a summary.

## SMALL DATA BURSTS TRAFFIC

The majority of smart phone and tablet use involves applications such as audio and video streaming, web page downloads, or picture and video uploads. For such applications, traffic typically consists of large data bursts composed of multiple packets with short inter-packet times. The bursts themselves may be separated by larger inter-burst durations. This type of traffic model is depicted in Fig. 1a. Since the data bursts are large, air-interface efficiency is critical and thus additional signaling is used in 3GPP networks to establish one or more radio bearers where scheduled transactions can occur. This is warranted as the signaling overhead to establish the radio bearers is small in proportion to the user data. Note that in 3GPP networks, for base station scalability and device energy saving, radio bearers are released at the end of the burst if there is no additional traffic within a predetermined period, which is typically on the order of a few seconds.

---

The authors are with Alcatel-Lucent, Bell Labs.

Figure 1b shows a second type of traffic where there are many small packets per burst. For example, these might correspond to TCP acknowledgment packets that are sent on the uplink during a large download. Although each individual packet is small, there are a large number of them with small inter-packet times, resulting in a large overall burst size. Hence establishing a radio bearer for efficiency is again warranted.

Figure 1c shows a third type of traffic model in which the overall burst size is small, with only one or two packets per burst. This type of traffic is typical of many sensor and wearable devices reporting sensed value, events, or status [1, 2]. Light smart phone background traffic also fits this type of model [3]. For such a traffic model, setting up a radio bearer per burst will incur a significant overhead. The overall efficiency, including the signaling, will be lower than if a less spectrally efficient, common channel access with reduced signaling is employed.

In summary, the traffic is characterized by packet size, burst size or number of packets per burst, inter-packet time, and inter-burst time. Since it is not efficient to retain the radio bearer during long periods of inactivity, for traffic with small burst size and long inter-burst time, significant signaling overhead relative to the burst size is incurred for each burst to establish and tear down the radio bearer. Hence this approach is very inefficient. The connectionless access mechanism proposed in this article is targeted at this type of traffic.

## CONNECTION ORIENTED SMALL PACKET TRANSMISSION

When a 4G device is powered on, it has to “attach” to the LTE network by going through an attach signaling procedure in which the network authenticates the device, assigns temporary identifiers, and allocates one or more IP addresses to the device [8]. While attached to the network, the device transitions to the active state to transmit or receive data. When there are no packets exchanged for a period of about five seconds (configurable by the operator), the device enters the idle state. This conserves computing and memory resources in both the eNodeB and the device, and also some air-interface resources even if no data packets are exchanged. Hence, eNodeB scalability and device battery life hinges on transitioning the device to the idle state after a period of inactivity.

Transmission of new data in LTE is based on a connection oriented approach in which the device first transitions to the active state by establishing bearers using radio resource control (RRC) signaling to the eNodeB and non-access stratum (NAS) signaling to the mobility management entity (MME). A typical sequence of information exchange before the first transmission is highlighted in Fig. 2. Numbers in parentheses indicate the size of the messages in bytes, not including lower layer overheads. Radio link control (RLC) status and RLC acknowledgement messages, physical uplink and downlink common control channel messages used for resource

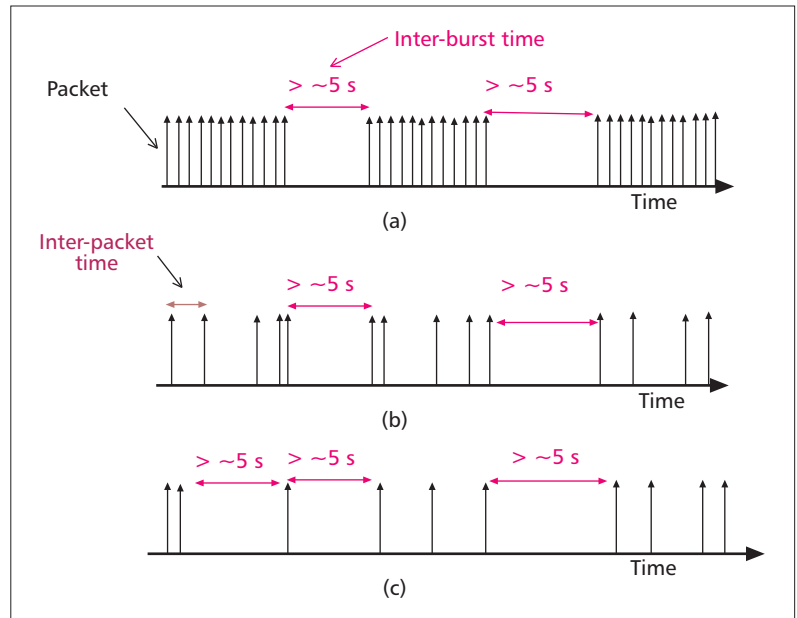


Figure 1. Representation of different traffic types.

requests by the UE, and resource grants by the eNodeB, respectively, are not shown. The first message is an uplink message on the random access channel in which a UE sends a preamble that allows the base station to synchronize to the UE’s transmission. This message carries no data, but the preamble occupies six physical resource blocks. Resources are allocated in the response message to transmit the first RRC signaling message with UE-specific information. In response, the eNodeB sends the RRC connection setup complete message with all the information required to configure the RRC state in the UE. The UE acknowledges and adds the NAS service request that is sent to the MME to retrieve the security context and GSM Packet Radio Service (GPRS) tunneling protocol (GTP) tunnel end point identifiers that will be used by the eNodeB to tunnel the data packet to the upstream serving gateway router node. Based on the security context received from the MME, the eNodeB chooses a security algorithm and then informs the UE. The UE responds by accepting the security mode command. After this, RRC reconfiguration is performed to finalize the RRC parameters, set up data radio bearers, and indicate the measurement reports to be sent by the UE. Finally, the UE responds with the reconfiguration complete message and uplink flows may begin. The eNodeB then informs the MME that bearers have been established and provides tunnel identifiers for sending downlink data. The MME sends the tunnel identifiers to the serving gateway (SGW), which can then forward downlink data toward the UE.

The total signaling overhead is thus in excess of 103 bytes in the downlink and 64 bytes in the uplink to establish a radio bearer with appropriate security. While this is not an issue for large bursts with multiple packets, for small bursts with small packet sizes, the overhead is comparable to or even more than the size of the data packet. Hence, for such traffic an alternate con-

nectionless access method described in the next section is more favorable.

## CONNECTIONLESS TRANSMISSION

Connectionless transmission allows the UE to transmit small packets of less than approximately 300 bytes without establishment of radio bearers. Dedicated service request signaling between the network and UE that is used in 4G to set up radio bearers and establish device and security contexts is eliminated, saving device battery power and reducing the number of network control plane transactions. Connectionless access may be coupled with contention based access or may have uplink data transmissions scheduled depending on the physical layer capabilities of the air-interface. In the case of contention based access, an attached device may simply awake, locate, and synchronize to an acceptable local base station based on broadcast information, and send a user plane packet. In the case of scheduled access, the device sends a resource request first and is assigned uplink transmission resources in the response from the base station. In both cases, the device may subsequently set a timer and await a response from the network, or immediately return to its sleep state.

We foresee connectionless access applying to an evolved 4G network or a 5G network where the air interface can simultaneously support both

connection oriented access and connectionless access. For 4G the supporting access and core networks would contain an evolved eNB, MME, and S/PGW. For 5G the corresponding elements may be a base station, a controller, and an IP/mobility anchor.

Prior to initiating a user plane transaction using connectionless access, a device must first attach/register with the network. For an IoT device such as a sensor, this can be a procedure that occurs occasionally, when the device is deregistered, the connectionless access procedure is unsuccessful, or when there is intervention from an operator. During the attach process, a controller is selected, which subsequently authenticates the device, and verifies the subscriber is authorized for wireless service. The controller also establishes and stores context information, including security keys that will be used during subsequent connectionless accesses for integrity protection and encryption of user data. The network may also assign mobility and IP anchors at nodes in the network and establish enforcement points for charging and QoS. Finally, an IP address/prefix is assigned to the device to be used for subsequent connectionless transactions.

When the attach process completes, the device may transition to an idle state or immediately engage in a connectionless transaction.

As described in the previous section, with 3GPP access, radio bearers are established and access stratum security activated via an exchange of RRC signaling between the mobile and the eNB for each small burst communication. The benefits of connectionless access stem from eliminating that signaling, and optionally lower layer scheduling requests and grants for uplink transmissions as shown in Fig. 3. An architecture that supports this call flow is shown in Fig. 4a, and the corresponding protocol processing is shown in Figure 4b. The architecture contains a BTS (eNodeB in 4G) which receives and processes the packet from the mobile, a controller (MME in 4G) that maintains UE state information and with which the BTS interacts, and optionally an IP and mobility anchor (S/PGW in 4G).

In the first step in the call flow, an access request is transmitted by the UE, and a response is received from the network. The exchange provides timing information so that the UE can align its physical layer frame boundaries, compensating for the round-trip transmission delay to the eNB, and optionally scheduling information so the UE can send uplink transmissions without colliding with other devices. For devices for which timing is known a priori from prior transactions because the device is static, a lower overhead random access transmission can be used, avoiding the access request and response messages.

For contention based access, a physical layer random access channel (PRACH) specifically designed for connectionless access could be more efficient in 5G networks than the PRACH in LTE. This is achieved through the use of new waveforms, through different parameter choices such as preamble length targeting a different operating regime, and through superior detection techniques [9, 10].

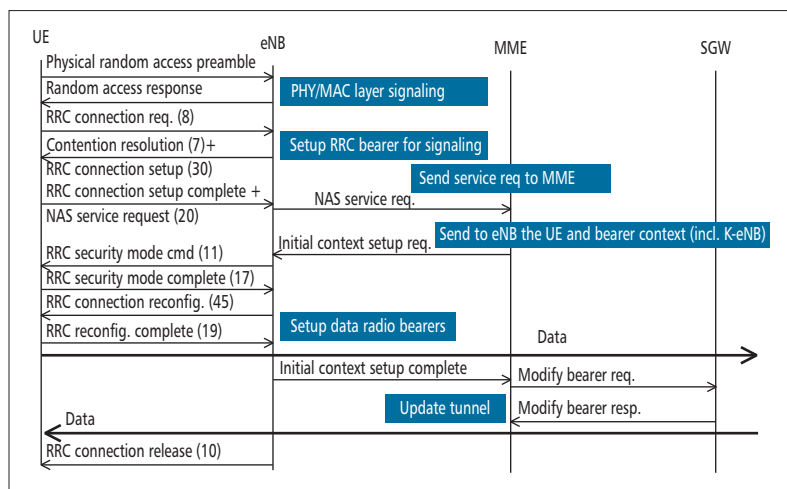


Figure 2. Signaling in LTE for connection oriented data transfer.

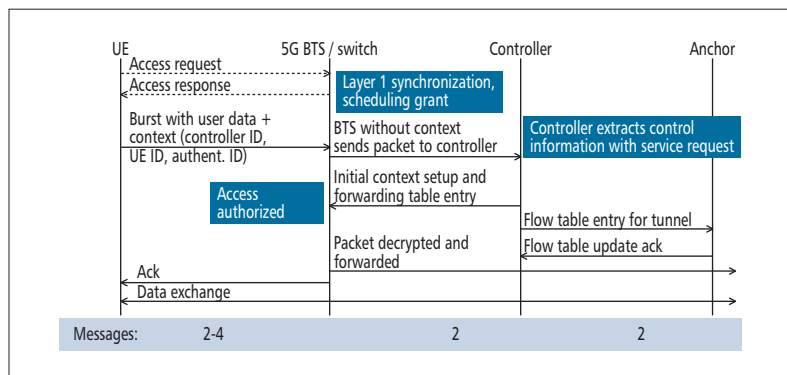


Figure 3. Connectionless access call flow.

Whereas in 4G, subsequent RRC messages from the eNB to the UE configure signaling and data radio bearers, including logical channel, Packet Data Convergence Protocol (PDCP), RLC, and medium access control (MAC) layer parameters, in connectionless access either pre-selected default values are used, or sub-functions within each layer are disabled. For example, RRC configures PDCP parameters for “discard timer” to determine when to discard transmitted PDUs, maximum allowed number of “context identifiers”, and maximum “sequence number” length among other parameters. With the restricted size of connectionless bursts, these normally configurable parameters may be assigned fixed values, obviating the need for RRC configuration.

Without the need for further parameter configuration, the UE transmits the small data burst containing user data. However, it is not sufficient to simply transmit a user plane IP packet without providing additional context information in the transmission as the receiving base station may have had no prior interaction with the device, or previously stored information may have expired due to inactivity. The information provided with the burst must include a globally unique identifier of the controller with the stored UE context information, and a UE identifier unique within the domain of the controller. For sizing purposes the comparable 4G controller identifier is the 6 byte globally unique MME ID (GUMMEI) and the 4 byte M-TMSI (temporary mobile station ID).

The context information may be conveyed by the UE as a MAC Layer 2 (L2) header appended to the user plane IP packet, as indicated in Fig. 4b. When the data burst is received by the BTS, the context header is removed, encapsulated in an IP packet, and sent to the controller indicated by the context information.

The controller uses the UE ID as a key to retrieve the stored state information and determine the forwarding path for the data. The controller sends the base station at which the packet was received both the forwarding instructions and the relevant state information, and may update anchor points to indicate the new point of attachment for the UE. After service is authorized the uplink packet is forwarded to the anchor, and onward to the corresponding node.

The BTS sends an acknowledgement to the UE indicating whether the transaction was successful. If the controller has indicated that the device is likely to send additional bursts, the BTS may assign the UE a new identifier unique within the scope of the BTS, similar to the 2 byte cell radio network identifier (C-RNTI) used in 4G. Until a dormancy timer expires, this identifier may be used in the context header of subsequent bursts sent to the same BTS in place of the globally unique M-TMSI+GUMMEI like identifier, saving approximately 8 bytes of overhead in each transaction. Using the new identifier, the BTS can retrieve the locally stored state information and the packet may be processed without further involvement of the controller.

When a downlink response is received in the network for an active device, the packet(s) are forwarded by the network to the BTS where the

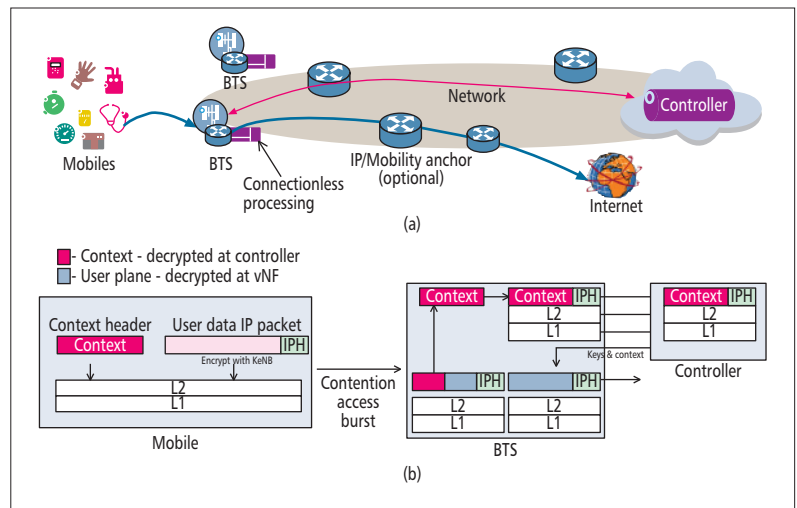


Figure 4. a) Connectionless architecture; b) protocol for connectionless access.

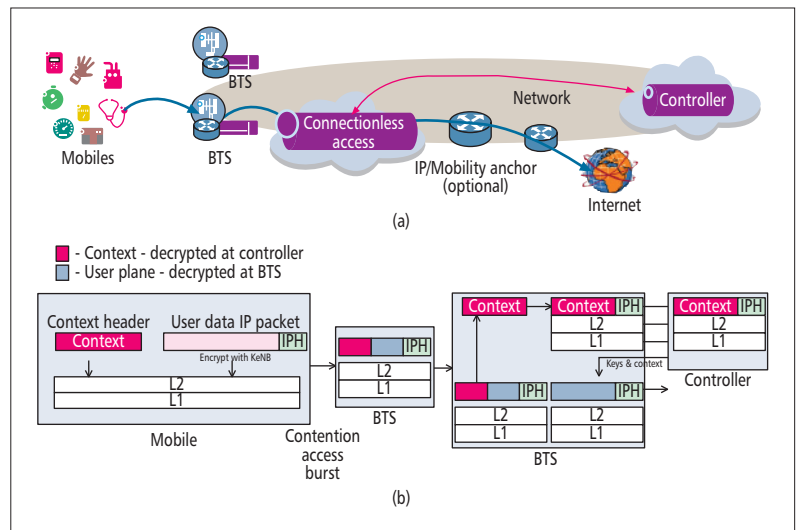


Figure 5. a) Connectionless architecture with vNF processing; b) protocol for connectionless access with vNF processing.

UE initiated access. The packet(s) are encrypted using the previously obtained security context and sent to the UE via a shared physical channel dedicated to connectionless access. Downlink transmission may be scheduled among the UEs using a simple first-in first-out (FIFO) queue, with no additional RRC signaling or channel quality reporting from the UEs. This simplification of downlink packet flow in the BTS allows connectionless access to scale to a large number of devices.

Figures 5a and 5b show an alternative architecture and protocol where the processing for security and forwarding, previously instantiated in the base station, is instead performed by a more centrally located connectionless access virtualized network function (CAvNF), instantiated in a data center and common to many BTSs. This simplifies the base station, and allows state information established in the vNF for a burst sent at one BTS to be used when the UE moves to a new BTS.

The BTS simply encapsulates the connection-

One of the critical aspects of any new protocol for data communication is security. 3GPP has standardized procedures for authentication, integrity protection, and data encryption for both signaling and data exchange at the access stratum between the UE and the eNodeB, and also in the non-access stratum between the UE and the MME.

less packet and sends it to a previously assigned CAVNF through a common IP tunnel across all devices, substituting the globally unique UE ID for any local, BTS assigned UE ID in the context header. The CAVNF extracts the context header and if necessary queries the controller to retrieve state information for the UE. Upon receiving a response from the controller the CAVNF decrypts the packet and forwards it according to instructions provided by the controller.

Downlink packets received by the network in response are forwarded to the CAVNF, where they are encrypted using the stored security context. They are then sent in a tunnel to the appropriate BTS. The BTS maps the IP address to the UE ID inferred from the uplink packet and forwards the packet to the UE.

Although the connectionless access described above still requires the BTS to maintain state for the device during the transaction as in the current connection oriented approach, we emphasize that the BTS state required is primarily restricted to the security context, which is above the physical (PHY) and MAC layers. The lower layers of the BTS are thus truly connectionless with UE state only to map a locally assigned UE ID to a globally unique ID for devices likely to transmit multiple short data bursts. Thus, the BTS can more easily scale to a large number of UEs. Furthermore, since lower-layer states are minimized and no air-interface resources are consumed when no packets are transferred, the dormancy timer can be set to a large value and security context retained for a longer time than the duration for which a typical dormant radio bearer is retained in 4G LTE, which in turn reduces the amount of network signaling (BTS to controller messages) required for connectionless access. A similar advantage applies with the CAVNF. The BTS need only maintain a mapping of local UE ID to global ID. All higher-layer (above RLC) processing is instantiated in a data center, allowing better scalability and longer dormancy timers compared to 4G.

The gains for connectionless access relative to the connection oriented approach in LTE can be assessed by estimating the overhead associated with connectionless access. The first data burst sent on a base station requires in the context header 10 bytes to identify the UE and controller, plus 3 bytes for security related information as detailed in the next section (13 bytes in total). The acknowledgement sent on the downlink adds another 4 bytes. Subsequent bursts on the same BTS require only 2 bytes for the UE ID. This compares favorably to the 103 downlink and 64 uplink bytes needed to setup radio bearers in 4G access for connection oriented access.

## SECURITY IN CONNECTIONLESS ACCESS

One of the critical aspects of any new protocol for data communication is security. 3GPP has standardized procedures for authentication, integrity protection, and data encryption for both signaling and data exchange at the access stratum between the UE and the eNodeB, and also in the non-access stratum between the UE and the MME. These procedures are executed

in LTE systems when the device attaches/registers with the network and during radio bearer establishment prior to transmission of data. Any new solution such as the proposed connectionless access should incorporate mechanisms to ensure the same level of security. In this section, we begin by describing the LTE security architecture and then discuss how the proposed connectionless access achieves the same level of security.

In LTE, a security context is established between the UE and the MME through the Extensible Authentication Protocol/Authentication and Key Agreement (EAP/AKA) protocol, resulting in access stratum security management entity keys (KASME) in the MME and UE. Security between the UE and MME is then activated with an exchange that selects the security algorithm or mode, and additional keys are derived from KASME for signaling integrity protection and encryption. This process usually occurs during the attach procedure.

When the UE has data to send, it transitions from idle to active and sends a service request message containing an authentication code for integrity protection among other fields. When the MME receives the message, it verifies the integrity of the message and derives a fresh base station key (KeNB) for user plane and RRC signaling between the UE and the eNB. KeNB is based on KASME and a NAS message count (NAS COUNT), a portion of which is sent by the UE and MME as a sequence number. The initial NAS message is sent without ciphering, but with integrity protection. Integrity protection uses a KASME derived key, a NAS sequence number count, and the received message among other inputs to calculate an expected authentication code. Integrity protection is successful if the code calculated at the receiver matches with one appended to the received message. Similarly for encryption, a count and key are used by the sender to derive a keystream block that is XORed with a data block to produce a ciphertext block. The process is repeated in reverse at the receiver to decrypt the data. The UE independently derives the KeNB after receiving the AS security mode command indicating the chosen ciphering and integrity protection algorithms.

A similar structure, with separate security associations between the UE and the controller, and between the UE and the base station or CAVNF, can be used to support security for connectionless access. However, we propose the modification that encryption and integrity protection are always provided in the base station or CAVNF with a pre-selected algorithm, eliminating signaling to communicate information such as the key set identifier and the selected security algorithm. We do not see any major issues with standards prescribing a security algorithm that both the UE and base station must implement in order to simplify the signaling.

As in 4G, for connectionless access a security context is initially established between the UE and the controller when the UE attaches to the network and authenticates. A persistent, shared key similar to KASME is established in the UE and controller.

After exiting the idle state, prior to sending a connectionless burst, the UE derives a KeNB from its stored persistent key and a count similar to the NAS counter used in 4G. In the first connectionless burst, the UE integrity protects the context header based on the controller security association, and encrypts the user data IP packet based on a KeNB derived key, as described below.

For integrity protection, a two byte authentication code similar to that used for NAS integrity protection, and a one byte sequence number similar to that derived from NAS COUNT, is appended to the context header containing the UE and controller identification information described earlier. The code is calculated in the UE from a key derived from the persistent (KASME-like) key, the context header, and other inputs such as key length and direction (uplink). The authentication code and sequence number, along with the context header, are forwarded by the base station to the controller. The controller retrieves the stored context based on the provided UE ID, calculates the expected authentication code, and verifies that it matches that provided by the UE. The integrity protection achieved is thus on par with the NAS signaling in the radio bearer establishment procedure. If successful, the controller calculates a new KeNB from its copy of the persistent key and the sequence number, and sends the resultant to the base station in the "initial context setup" message indicated in Figure 3, or to the CAvNF in Figure 5. The provided KeNB is used to calculate the key for decrypting the user data IP packet before forwarding it onward.

Once the security context has been established in the base station, subsequent bursts need not contain the authentication code or sequence number. The user data is encrypted/decrypted using the locally available keys associated with the security context available in the base station and the burst count at that base station. Similarly the security context established in the CAvNF may be used for subsequent bursts sent on the same BTS. The context may also be used if the mobile moves to a neighbor BTS and sends a subsequent burst. In this case the same key may be maintained, or a new key calculated from the old key using a mechanism similar to 3GPP horizontal key derivation [11].

When the network idle timer expires, the base station or CAvNF removes all context. A similar timer in the UE triggers a transition to the idle state. The process of retrieving the security context from the controller and establishing it at the new base station is repeated when the UE once again has a short burst to send.

## SUMMARY

The explosion in IoT devices on cellular networks is expected to significantly increase the proportion of small burst traffic. While incremental modifications are being made to LTE to accommodate such traffic, a more radical modification is needed to substantially increase the efficiency. In this article we proposed a connectionless access approach for handling small burst traffic. The connectionless access call flow and two implementation alternatives, where security and forwarding functions are at the BTS or in a CAvNF, were discussed. Initial estimates on the signaling overhead show potential for substantial reduction relative to LTE. We believe that future cellular networks could incorporate such a protocol in combination with efficient mechanisms for contention based or scheduled access at the physical layer.

tionless access approach for handling small burst traffic. The connectionless access call flow and two implementation alternatives, where security and forwarding functions are at the BTS or in a CAvNF, were discussed. Initial estimates on the signaling overhead show potential for substantial reduction relative to LTE. We believe that future cellular networks could incorporate such a protocol in combination with efficient mechanisms for contention based or scheduled access at the physical layer.

## ACKNOWLEDGMENTS

We thank Alistair Urie, Alessio Cassati, Sudeep Palat, Dietrich Zeller, Mark Doll and Bozo Cesar at Alcatel-Lucent for many fruitful discussions.

## REFERENCES

- [1] M. Zubair Shafiq *et al.*, "A First Look at Cellular Machine-to-Machine Traffic: Large Scale Measurement and Characterization," SIGMETRICS'12, June 11-15, 2012, London, UK.
- [2] H. Viswanathan, K. Mun, and F. Scholler, "Modeling and Analysis of Cellular Wireless Machine-to-Machine Traffic," *Bell Labs Technical Memorandum*, 2012.
- [3] Third Generation Partnership Project (3GPP) TSG RAN, "LTE RAN Enhancements for Diverse Data Applications (Release 11)," 3GPP TR 36.822 v1.0.2, 2012.
- [4] F. Xia *et al.*, Internet of Things, *International Journal of Communication Systems*, vol. 25, no. 9, Sept. 2012.
- [5] K. Zheng *et al.*, "Radio Resource Allocation in LTE-Advanced Cellular Networks with M2M Communications," *IEEE Commun. Mag.*, vol. 50, no. 7, July 2012, pp. 184-92.
- [6] S.-Y. Lien, K.-C. Chen, and Y. Lin, "Toward Ubiquitous Massive Accesses in 3GPP Machine-to-Machine Communications," *IEEE Commun. Mag.*, vol. 49, no. 4, Apr. 2011, pp. 66-74.
- [7] Y. L. Lee *et al.*, "Recent Advances in Radio Resource Management for Heterogeneous LTE/LTE Networks," *IEEE Commun. Surveys and Tutorials*, vol. 16, no. 4, 2014.
- [8] 3GPP TS23.401, 3rd Generation Partnership Project; Technical Specification Group Services and System Aspects; General Packet Radio Service (GPRS) enhancements for Evolved Universal Terrestrial Radio Access Network (E-UTRAN) access (Release 12).
- [9] V. Vakilian *et al.*, "Universal Filtered Multi-Carrier technique for wireless systems beyond LTE," *Proc. 9th Int'l. Wksp. Broadband Wireless Access*, IEEE Globecom'13, Atlanta, GA, USA, Dec. 2013.
- [10] H. Dhillon, H. Huang, H. Viswanathan, and R. Valenzuela, "Throughput Optimal Communication Strategy for the Random Access Channel," *IEEE Globecom'13*, Atlanta, GA, USA, Dec. 2013.
- [11] 3GPP TS33.402, 3rd Generation Partnership Project; Technical Specification Group Services and System Aspects; 3GPP System Architecture Evolution (SAE); Security architecture (Release 12).

## BIOGRAPHIES

COLIN KAHN [M] (colin.kahn@alcatel-lucent.com) is a member of the Corporate CTO Organization at Alcatel-Lucent. He currently supports 5G and LTE access and core network architecture initiatives, focusing on the development of new solutions that leverage Alcatel-Lucent's traditional strengths in network systems. Over the past 20 years he has worked at ALU, Lucent Technologies, and AT&T, providing systems engineering, standards, and customer support for IS-136 TDMA, CDMA (IS-95, 3G1X and EV-DO), GSM, UMTS, and LTE. Prior to joining the AT&T wireless business unit he spent six years at AT&T Federal Systems conducting acoustics related research. Prior to joining AT&T he conducted fusion energy research at General Atomic Corp. and the Princeton University Plasma Physics Laboratory. He holds electrical engineering degrees from M.I.T. and Cornell University.

HARISH VISWANATHAN [F] (harish.viswanathan@alcatel-lucent.com) is a CTO Partner in the Corporate CTO Organization at Alcatel-Lucent. He received the B. Tech. degree from the Department of Electrical Engineering, Indian Institute of Technology, Chennai, India, and the M.S. and Ph.D. degrees from the School of Electrical Engineering, Cornell University, Ithaca, NY. He was a recipient of the Sage Fellowship at Cornell. Since joining Bell Labs in October 1997 he has worked on multiple antenna technology for cellular wireless networks, network optimization, network architecture, mesh networking, and M2M communications. His research interests are in the areas of communications theory, networking, and information theory.

Initial estimates on the signaling overhead show potential for substantial reduction relative to LTE. We believe that future cellular networks could incorporate such a protocol in combination with efficient mechanisms for contention based or scheduled access at the physical layer.

---

# UNDERSTANDING THE IoT CONNECTIVITY LANDSCAPE: A CONTEMPORARY M2M RADIO TECHNOLOGY ROADMAP

The authors address the market-changing phenomenon of the Internet of Things (IoT), which relies on the underlying paradigm of machine-to-machine (M2M) communications to integrate a plethora of various sensors, actuators, and smart meters across a wide spectrum of businesses.

*Sergey Andreev, Olga Galinina, Alexander Pyattaev, Mikhail Gerasimenko, Tuomas Tirronen, Johan Torsner, Joachim Sachs, Mischa Dohler, and Yevgeni Koucheryav*

## ABSTRACT

This article addresses the market-changing phenomenon of the Internet of Things (IoT), which relies on the underlying paradigm of machine-to-machine (M2M) communications to integrate a plethora of various sensors, actuators, and smart meters across a wide spectrum of businesses. Today the M2M landscape features an extreme diversity of available connectivity solutions which, due to the enormous economic promise of the IoT, need to be harmonized across multiple industries. To this end, we comprehensively review the most prominent existing and novel M2M radio technologies, as well as share our first-hand real-world deployment experiences, with the goal to provide a unified insight into enabling M2M architectures, unique technology features, expected performance, and related standardization developments. We pay particular attention to the cellular M2M sector employing 3GPP LTE technology. This work is a systematic recollection of our many recent research, industrial, entrepreneurial, and standardization efforts within the contemporary M2M ecosystem.

## COMMUNICATIONS STANDARDS

## INTRODUCTION AND OPPORTUNITIES

In the 1990s the word “Internet” had the connotation of a physical system of computers networked by means of an Ethernet cable; today, this is forgotten and the Internet is synonymous with the likes of Facebook, eBay, and LinkedIn. The Internet has thus undergone an enormous transformation from being technology-driven to becoming market-driven. The decoupling of underlying technologies from the services able to run on top of them has been a painful but instrumental shift in unlocking what is now often referred to as the 3rd Industrial Revolution.

Going beyond the 3rd Industrial Revolution, we are rapidly moving toward a world of ubiquitously connected objects, things, and processes. It is the world of the emerging Internet of Things (IoT), which has the potential to produce a new wave of technological innovation. Indeed, the range of IoT applications is extremely broad, from wearable fitness trackers to connected cars,

spanning such industries as utilities, transportation, healthcare, consumer electronics, and many others (Fig. 1).

However, we are only beginning to witness the true explosive growth of the IoT, with 10 billion M2M devices connected presently and 24 billion to 50 billion total connections expected within the next five years. Thus, over the following decade we may see our everyday furniture, food containers, and even paper documents accessing the Internet. Futurists have also coined a number of new keywords to emphasize the IoT’s ongoing transformation, including the Internet of Everything (by Cisco), the Industrial Internet (by General Electric *et al.*), as well as the Networked Society (by Ericsson).

Today, machine-to-machine (M2M) technologies are an integral part of the IoT connectivity ecosystem [1] and serve as the underlying facilitator for the IoT phenomenon. But they are just a small part; they are the beginning; they are, in a sense, the new (mostly wireless and feature-rich) “Ethernet cable” able to connect objects with other objects, with people, and the enormous computing nervous system spanning the globe. Surprisingly, the design efforts related to M2M span back a few decades.

Indeed, driven by industrial needs, early forms of M2M connectivity trace back to supervisory control and data acquisition (SCADA) systems of the 1980s, all being highly isolated and proprietary connectivity islands [2]. Along the way of its rapid development, the connectivity landscape has embraced legacy radio frequency identification (RFID) technologies (starting in the late 80s), as well as wireless sensor network (WSN) technology (starting in the 90s). Marked by the very attractive application scenarios in both business and consumer markets, the first decade of the 21st century was thus dedicated to the development of standardized low-power M2M solutions, through either industry alliances or standards developing organizations (SDOs).

Notable examples tailored to a range of industry verticals are ISA100.11a, WirelessHART, Z-Wave, and KNX. More generic (horizontal) connectivity technologies were developed within the leading SDOs, i.e. the IEEE, ETSI, 3GPP, and IETF (even though strictly not an SDO). Low-power short-range solutions available today include Bluetooth (promoted by the Bluetooth SIG) and IEEE 802.15.4 (promoted by the Zigbee alliance) [3]. In subsequent years, the IEEE 802.15.4 physical (PHY) and medium access control (MAC) layers have been complemented by IP-enabled (networking), as well as the web-enabled IETF stacks. In parallel, capitalizing on the ability to provide global coverage, 3GPP developed cellular-enabled machine-type connectivity modules [4] tailored to markets with inherent mobility (e.g. car telemetry).

Despite decade-long developments by some of the best engineering teams in the world, none of the above technologies has emerged as a clear market leader. The reasons, a mix of technology shortcomings and business model uncertainties, are rather important and thus discussed subse-

---

*Sergey Andreev, Olga Galinina, Alexander Pyattaev, Mikhail Gerasimenko, and Yevgeni Koucheryav are with Tampere University of Technology.*

*Tuomas Tirronen, Johan Torsner, and Joachim Sachs are with Ericsson.*

*Mischa Dohler is with King’s College London.*

This work is supported by GETA, TISE, and the Internet of Things program of DIGILE, funded by Tekes. The work of the first author is supported with a Postdoctoral Researcher grant by the Academy of Finland as well as with a Jorma Ollila grant by Nokia Foundation.



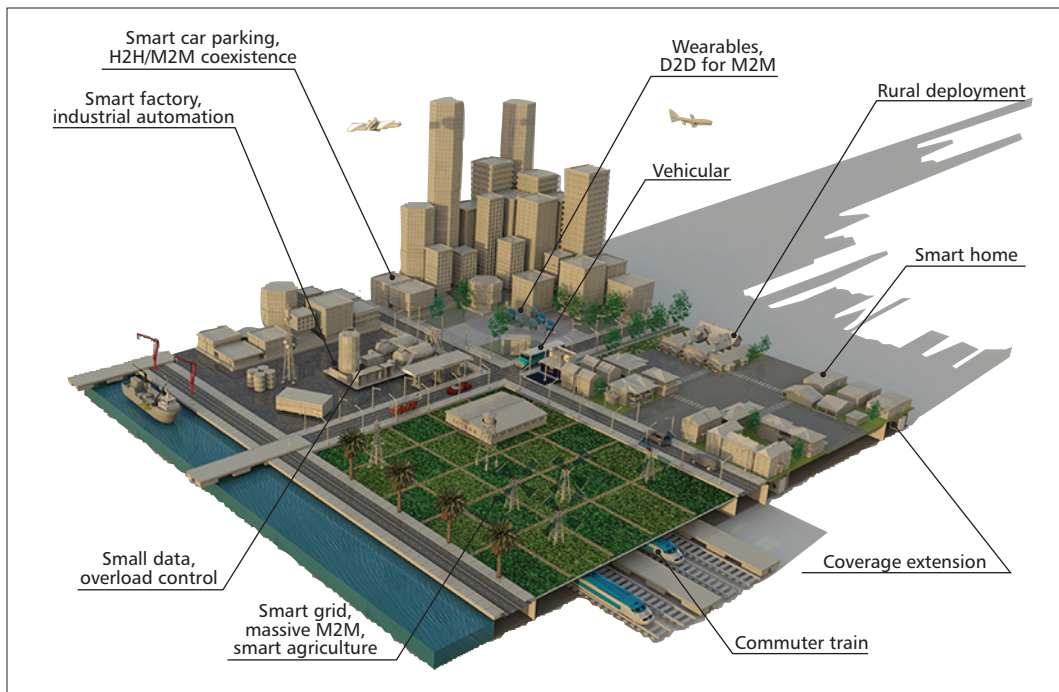


Figure 1. A vision of diverse M2M use cases across a variety of industries.

The first challenge was to identify suitable locations to place repeaters and gateways. Given the sheer density of repeaters needed and the high uncertainty of the propagation conditions, this involved very complex planning. A further uncertainty involved who pays the electricity for repeaters and gateways.

quently. A key consequence, however, is that the field of the IoT connectivity is now at a turning point, with many promising radio technologies emerging as true M2M connectivity contenders: Low-Power WiFi, Low-Power Wide Area (LPWA) networks, and various improvements for cellular M2M systems.

These afterthought solutions may be significantly more attractive for the prospective IoT deployments from both the availability and reliability points of view, and, given their emerging nature, we focus on characterizing these in the remainder of this article. With our hands-on experience in design, standardization, as well as roll-out of these technologies, we share our most essential findings in this work. We believe that these solutions may allow for a decisive transformation of the global M2M industry and thus enable a truly dynamic and sustainable IoT ecosystem at par with the Internet of today.

## SMART CITY IoT: THE AWAKENING REALITY CHECK

From 2010 to 2012 we have been using a set of aforementioned technologies in various Smart City deployments around the world. After a decade of theoretical, design, standards, and engineering work, it was thus an opportunity to prove the viability of IoT connectivity solutions. We share our experience in what follows, which forms the rationale for the subsequent sections.

### REAL-WORLD IoT ROLL-OUTS

Worldsensing pioneered the concept of Smart Parking which, as shown in Fig. 2, involves the placement of sensors in every parking space to detect the presence or absence of cars in real-time [5]. This information is relayed to the driv-

er who thus avoids circulating in the city in a quest for a parking space.

Smart parking systems are of interest to cities as they reduce pollution, traffic volumes, and thus traffic jams. Given that congestion currently costs Europe about 1 percent of its GDP every year, smart parking is also seen as boosting the economy. The above projections, however, unlock only if some technical key performance indicators (KPIs) are met. Notably, the system should not be in outage for more than 0.1 percent per annum, which translates to about 9h per year. Furthermore, the parking information needs to be relayed within a few seconds.

Over the years of 2010–2012, smart parking systems of various sizes have been trialed in various cities, such as Moscow, Barcelona, and Sant Cugat (satellite city to Barcelona), among others. The topology of these early roll-outs included the Zigbee-powered smart parking nodes connected in a multihop mesh network until a repeater. The Zigbee-powered repeaters also networked in a multihop fashion until the gateway. The gateway was connected to the Internet via an Ethernet connection or a cellular 3G modem. There was a repeater every 5–10 parking nodes, and a gateway every 100–150 parking nodes. A trial typically involved at least 100 live parking spaces in a city, thereby giving a realistic picture of the technology's capabilities at a reasonable network scale.

### OBSERVED IoT DEPLOYMENT CHALLENGES

The first challenge was to identify suitable locations to place repeaters and gateways. Given the sheer density of repeaters needed and the high uncertainty of the propagation conditions, this involved very complex planning. A further uncertainty involved who pays the electricity for repeaters and gateways. To put this into perspec-

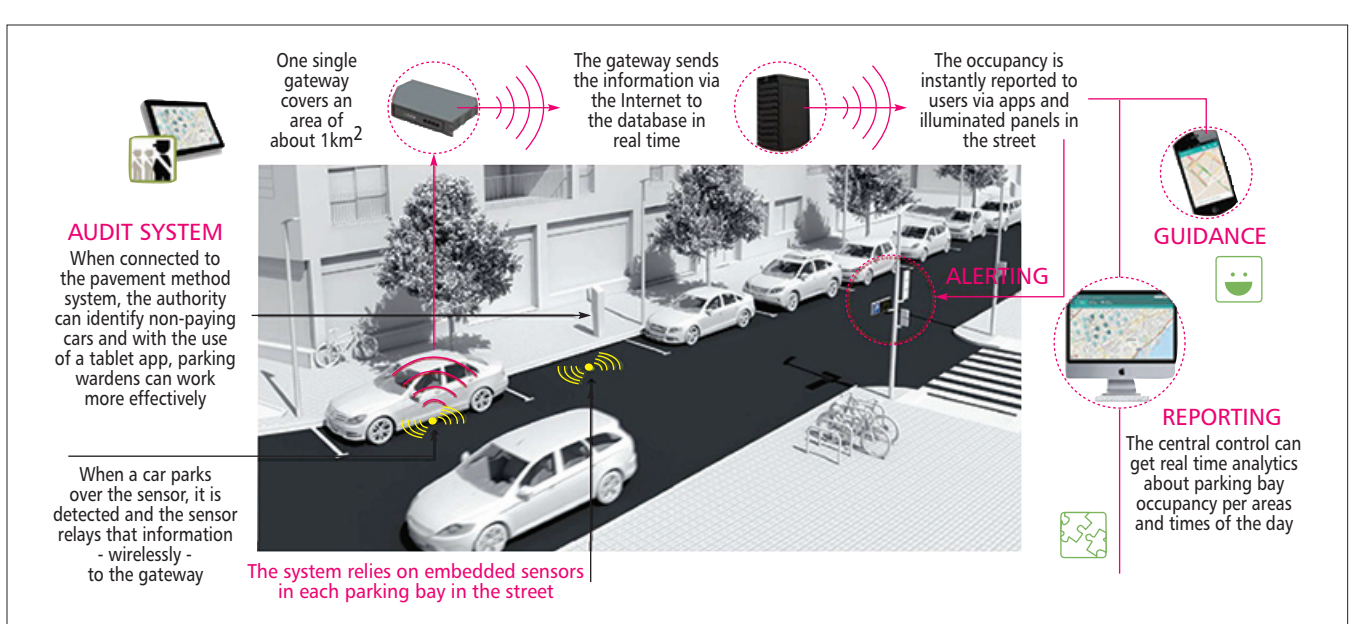


Figure 2. Functionality of Worldsensing's real-world IoT deployment of smart parking technologies.

tive, a city like Moscow with 60K parking spaces would require 6K lamp posts to be equipped with repeaters.

The second challenge was to ensure reliable and robust connectivity no matter what the car parking and interference situation is, so as to ensure the above KPIs could be met. Given the high degree of freedom of the network due to mesh deployment and the high dynamics of the channel due to cars' movement, the system was very often in outage. In some unfavorable testing conditions, the outage reached 10 percent or more and would thus have been in straight violation of a service level agreement (SLA).

The third challenge was to ensure that the KPIs in terms of delay were met. Given the multi-hop nature and the high channel dynamics with frequent outages, delays of minutes often occurred. Again, this would have violated any SLA.

### LESSONS LEARNED

Arguably, the biggest mistake of the M2M community in the early days was to believe in the need for low-power technologies, when we actually needed a high-transmission-power low-energy solution. This seems contradictory, but remember that power gives you range and energy drains your battery; and energy = power x time. IEEE 802.15.4 systems, however, only offer low power, which leads to short transmission range and thus the need for multihop. This, in turn, yields poor reliability due to the high degree of freedom as we have experienced in the above-described roll-outs. WiFi/3GPP technologies, transmitting at high power (with the advantage of a high communication range), are able to be energy efficient as long as the transmission is done within the shortest time.

Another lesson we learned is that system reliability matters, and not only link reliability. To be able to rely on a functioning M2M deployment, the underlying end-to-end system must be reliable and available, and not only singular links. In light of these requirements, it is apparent that cover-

age, support of mobility, and roaming are very poor with Zigbee, at least in large-scale Smart City deployments. Given that Zigbee did and still (mid-2015) enjoys roll-outs, it is far from reaching a critical mass. From a reliability point of view, it is extremely susceptible to interference (particularly in urban environments), has no throughput guarantees, and often produces lengthy system outages. Major companies have finally realized this and stopped producing Zigbee chips, while ramping up on low-power WiFi chip ranges.

## EMERGING M2M TECHNOLOGIES

Whereas Zigbee-like solutions may still find their market niche with simple and static applications, the large IoT market is, according to our experience, well beyond their reach. We are thus witnessing a shift in M2M connectivity technologies that will be discussed in the remainder of this article.

### LOW-POWER WiFi TECHNOLOGY

In recent years WiFi (IEEE 802.11) technology has experienced tremendous growth and has become a de-facto solution for home and corporate connectivity. However, WiFi has mostly been out of reach for M2M communications due to its fairly large energy consumption. This has changed as of late, i.e. when the IEEE 802.11 community started to apply duty cycling and hardware optimization, resulting in an extremely energy efficient system.

On the downside, support of mobility and roaming in WiFi is currently rather poor. In terms of reliability, there is neither guaranteed QoS support, nor adequate tools to combat severe interference typical for unlicensed bands. To this end, it was soon recognized that the favorable propagation properties of low-frequency spectrum at sub-1 GHz may provide improved communication properties when compared to, e.g. conventional WiFi protocols operating at the

2.4 GHz and 5 GHz bands. However, the available spectrum at sub-1 GHz license-exempt ISM bands is extremely scarce, and hence required careful system design considerations.

With this in mind, after outlining the purpose and the technical scope of the novel IEEE 802.11ah project, the standardization work of the corresponding TGah task group commenced in November 2010. The prospective technology is generally based on a variation of the IEEE 802.11ac standard, but down-clocked by a factor of 10. It is currently being developed to enable low-cost long-range (up to 1 km) connectivity across massive M2M deployments with high spectral and energy efficiencies. Today, thousands of M2M devices may already be found in dense urban areas, which required providing support for up to 6,000 machines connecting to a single access point.

Fortunately, IEEE 802.11ah technology does not need to maintain backward compatibility with the other representatives of the IEEE 802.11 family. Operating over different frequencies, 802.11ah could thus afford defining novel compact frame formats, as well as offering more efficient mechanisms to support a large population of devices, advanced channel access schemes, as well as important power saving and throughput enhancements [6]. As a result, 802.11ah is believed to significantly enrich the family of 802.11 protocols, which already receive increasing attention from mobile network operators willing to introduce low-cost connectivity in unlicensed bands.

### UNLICENSED LOW-POWER WIDE AREA NETWORKS

Given that Zigbee-like solutions have not lived up to their expectations, whereas Low-Power WiFi and Cellular M2M systems have commenced to take shape only recently, a novel class of M2M technologies has emerged lately, termed Low-Power Wide Area (LPWA), which operate in unlicensed spectrum. However, only low data rates and small daily traffic volumes are foreseen [7], which limits the application to a subset of M2M services with infrequent small data transmissions.

LPWA technology today is proprietary with multiple non-compatible alternatives. There are also initiatives to propose LPWA technology concepts into the cellular M2M direction within a recent new study item that has been initiated in the 3GPP GERAN (GSM/EDGE Radio Access Network) specification group. Similar to the standardization targets of LTE evolution in 3GPP Radio Access Networks (RANs), the motivation is to enable extended coverage beyond GSM coverage today, low device complexity, and long battery lifetimes.

Our experience with LPWA shows that it works successfully for large projects, such as the Moscow Smart City deployment [5], where almost 20,000 sensors have been connected to a modest number of access points. In the trial deployments, we have seen suburban and rural ranges of over 20 km, the typical urban ranges of around 5 km, and the “difficult” urban ranges of 1-2 km. Mobile network operators may hence become the early adopters of this emerging technology, building on their well-developed network infrastructure and strong customer trust. For

instance, a possible deployment model for an operator may be to install LPWA systems complementary to existing cellular technology and the cell sites that they already have [8].

Despite the time-to-market benefits of LPWA, there are also clear downsides of using unlicensed spectrum for long-range communication. Typical regulation imposes several restrictions on radio transmitters in unlicensed spectrum [9] in terms of effective radiated power (ERP), allowed duty cycles, and listen-before-talk requirements. For long-range transmissions, the limited ERP causes asymmetric link budgets between the uplink and downlink directions, because the ERP is limiting the radiated power after the antenna gain has been applied. However, antennas have significantly different performance between simple devices using a single antenna with around 0 dB antenna gain and a base station with an antenna gain of around 19 dB. This means that the uplink signal has an additional antenna gain at the receiver in contrast to the downlink signal. For European regulation, this can be partly compensated by selecting a subband for downlink with 13 dB higher output power. But even then a link asymmetry of at least 6 dB remains.

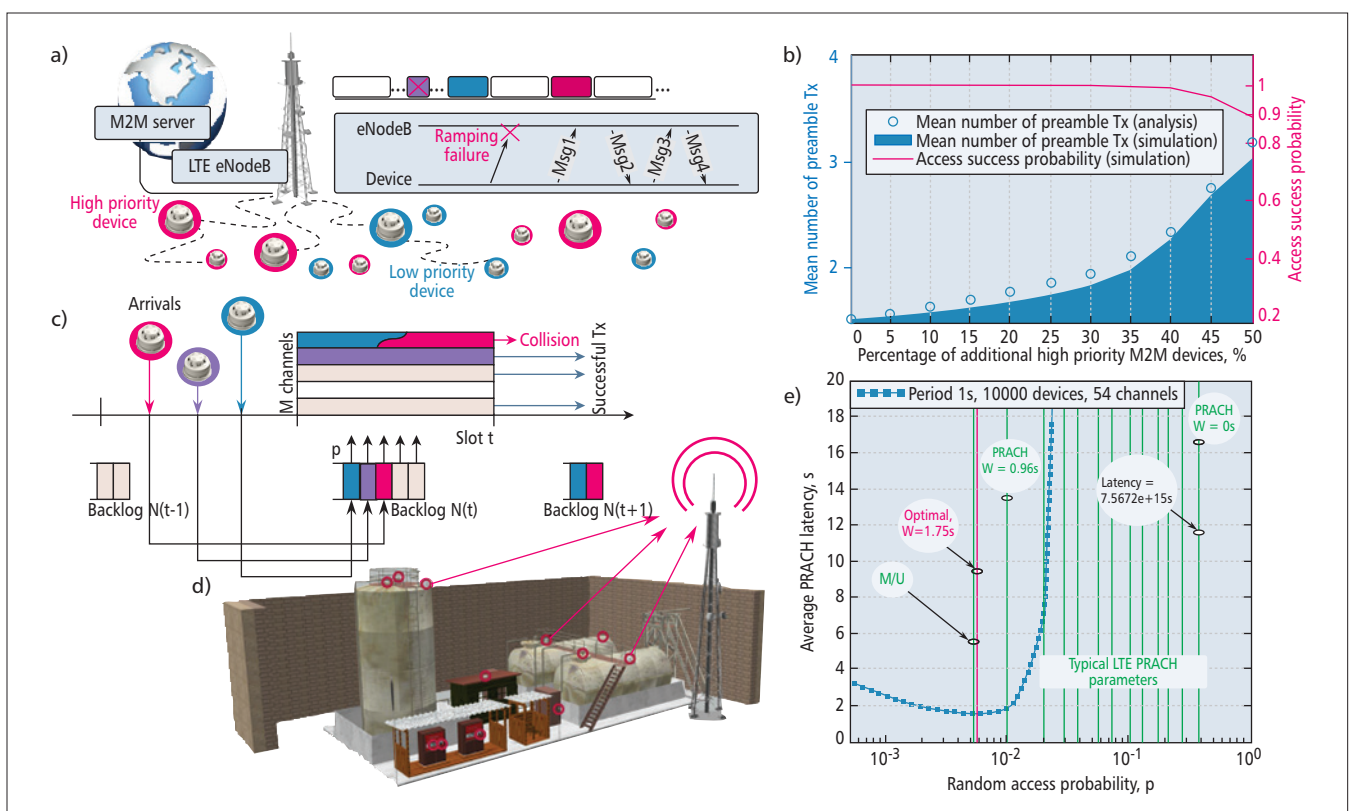
As a consequence, at least 50 percent of devices experience only uplink connectivity under non-line-of-sight propagation conditions. This is unreliable in the sense that no acknowledgements for successful uplink data delivery are possible. Further, scalability limitations come from the range covered by a single LPWA base station [7]. Projecting that the total number of connected M2M devices is to become approximately 10 times larger than the number of people, easily millions of devices may appear within the coverage area of a single LPWA base station. Many of those will use other radio technologies that share the spectrum with LPWA, such as low-power WiFi (IEEE 802.11ah), Z-Wave, Zigbee, IEEE 802.15.4g, etc. With its low receiver sensitivity for long-range communication, the LPWA device will perceive all of these other transmissions as interference.

We therefore foresee that LPWA will only remain viable at the early stage of IoT development when the number of devices is still moderate. However, LPWA can play an important role to support the early IoT market up-take until standardized cellular M2M solutions enter the market, which can handle the anticipated IoT scale in terms of numbers of devices, but also the variety of M2M services.

### CELLULAR M2M

Cellular technologies, and especially 3GPP LTE, are becoming increasingly attractive for supporting large-scale M2M installations due to their wide coverage, relatively low deployment costs, high level of security, access to dedicated spectrum, and simplicity of management. However, LTE networks have been neither historically designed with link budget requirements of M2M devices, nor optimized for M2M traffic patterns. Therefore, several improvements targeting M2M solutions have been initiated in 3GPP aiming at augmenting LTE to become more suitable for M2M applications.

We foresee that LPWA will only remain viable at the early stage of IoT development when the number of devices is still moderate. However, LPWA can play an important role to support the early IoT market up-take until standardized cellular M2M solutions enter the market.



**Figure 3.** Performance results on handling very large numbers of devices: a) motivating M2M scenario; b) connected-mode performance of different device priorities; c) proposed multi-channel M2M contention model; d) characteristic industrial automation application; e) analytical random-access latencies.

Given that the numbers of connected machines are expected to grow dramatically, LTE technology requires respective mechanisms to handle a very large number of devices [10]. Correspondingly, an overload control scheme named Enhanced Access Barring has been introduced as part of LTE Rel-11 to avoid overload in RAN, whenever there is a surge in near-simultaneous network entry attempts. Further, accounting for the fact that typical M2M data transmissions are infrequent and small, simplified signaling procedures for radio-bearer establishment are necessary to offer energy consumption savings for such M2M devices. In connection to lightweight signaling for small data, M2M device energy consumption can be reduced significantly for infrequent traffic by allowing for longer cycles of discontinuous reception (DRX).

In the following sections we review some of these important improvements in more detail. We intentionally focus our description on LTE, which we believe will become the major technology for M2M connectivity even though M2M-centric improvements are being discussed for other 3GPP technologies as well.

## M2M PERFORMANCE IMPROVEMENTS BY 3GPP HANDLING VERY LARGE NUMBERS OF DEVICES

Our research indicates that smart grid is one of the key M2M use cases incorporating a large number of metering devices that autonomously report their information to grid infrastructure.

The motivating smart metering use case therefore serves as a valuable reference “massive M2M” scenario covering many characterizing M2M features (see Fig. 1). Correspondingly, the involved M2M devices may be divided into several classes according to the priority of their information, e.g. high-priority (alarm messages) and low-priority (measurement data). Potentially, alarm messages constitute a bigger challenge for the network to handle, as they are typically highly synchronized and in addition may require certain latency guarantees [11].

Currently, the 3GPP LTE system defines a number of communication channels to deliver uplink transmissions from M2M devices to the network. In particular, the physical random access channel (PRACH) is employed by a device for its initial network entry, as well as to demand system resources if it does not already have a dedicated resource allocation. In case of many M2M devices connecting to the network near-simultaneously, we expect that the use of PRACH would be preferred, but may result in congestion due to its insufficient capacity.

More specifically, the PRACH procedure features two distinct stages (Fig. 3a). The former is the uplink timing synchronization stage (known as Msg1/Msg2), where the power ramping technique may be used to adjust the transmit power of a random-access preamble to particular channel conditions. Further, Msg3 is used to transmit a meaningful uplink message to the base station (termed eNodeB or eNB) and Msg4 is utilized for subsequent contention resolution.

To understand the impact of a large population of M2M devices on their network access, we construct an event-driven protocol-level PRACH simulator and thoroughly calibrate it against the reference 3GPP methodology documents [12]. Our simulation yields important conclusions on overloaded PRACH performance, when numerous *connected-mode* M2M devices of different priorities send their information into the network (Fig. 3b). In particular, we learn that around 40 percent of high-priority M2M devices, added to the original (typical) population of 30,000 (30k) low-priority devices, produce a sharp degradation in network access success probability.

Interestingly, PRACH preambles selected by the M2M devices randomly may be regarded as non-interfering code-based “channels” (Fig. 3c), where the case when two or more devices select an identical preamble (channel) would correspond to a conventional “packet” collision. This opens the door to assessing contention-based M2M behavior by relying on past knowledge of multichannel random-access protocols.

First, careful custom-made approximations can be forged for particular given ranges of PRACH parameters (Fig. 3b), such as the number of available preambles ( $M$ ) and contending devices ( $U$ ), backoff window size, etc. However, these may not be counted as adequate universal solutions, and another alternative is straightforward numerical analysis of contention behavior, which would only remain feasible for moderate numbers of users/channels due to high computational complexity.

More recently, we have demonstrated the feasibility of applying powerful fluid approximation techniques to rigorously characterize M2M performance, as well as the stability regions of a multichannel random-access system. As our target scenario, we have chosen an industrial automation application (Fig. 3d), which may require certain data access latency and reliability guarantees (e.g. for supporting priority or critical alarm messages). Along these lines, Fig. 3e indicates *analytical* PRACH latencies as evaluated with our method, which allows optimizing channel access by properly selecting the  $M$  over  $U$  retransmission probability value for arbitrary numbers of devices and channels.

More specifically, in the figure we compare our optimized latency against the values produced with the use of existing PRACH backoff indicator parameters. Our solution thus helps the base station regulate PRACH access by having system-wide knowledge across all connected M2M devices. However, if such knowledge is not available, simpler heuristic access control procedures (such as when the retransmission probability is chosen as  $M$  over  $U$ ) may be employed by the M2M devices locally, which sometimes results in close to optimal performance. These results allow for tighter control of important performance indicators, such as data access latency, which may benefit LTE in supporting constrained automation scenarios on the way to the Industrial Internet.

### ENERGY EFFICIENCY AND SMALL DATA TRANSMISSION

In tight connection with access latency and success rate of M2M transmissions goes their energy efficiency, which is accentuated by the fact

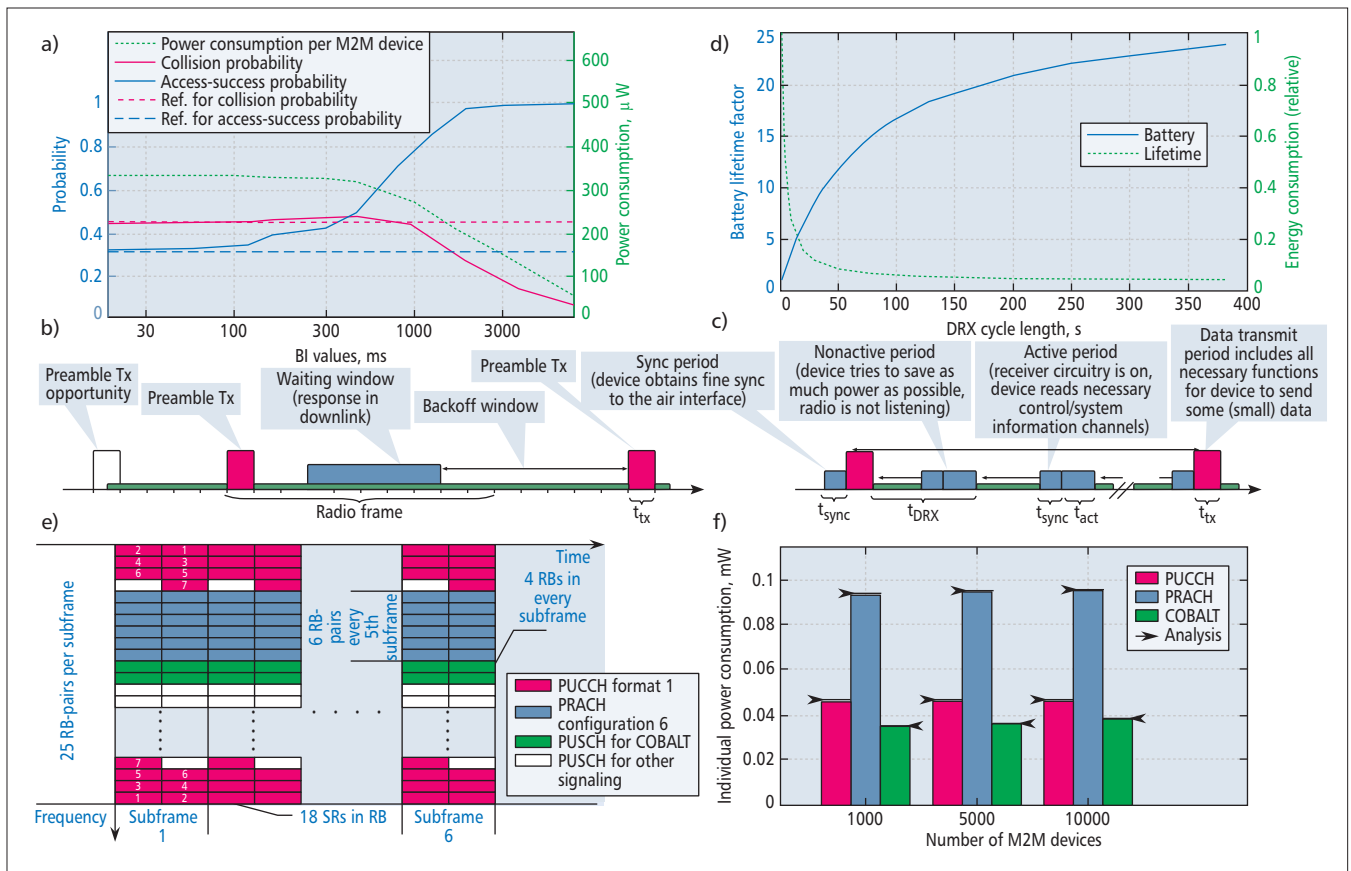
that M2M devices are typically small-scale and battery-powered. We continue to study the scenario when an IoT application requires a large number of M2M devices to perform a particular action near-simultaneously (e.g. smart meter data readings), or when an unexpected surge, outage, or failure occurs (massive power outage or restoration of power, network failure, etc.) causing multiple devices to (re)connect to the network within a short period of time. In this case, the transmitting devices would still be using the PRACH contention-based random access procedure to obtain uplink synchronization for initial network access or respective data transmission.

Along these lines, Fig. 4a illustrates the simulated initial network entry performance of 30K M2M devices with respect to their power consumption, collision probability, and access success probability across different PRACH backoff indicator (BI) values. These results are produced for characteristic beta-distributed M2M device activation patterns (traffic type 2: beta distribution over 10 seconds), as suggested by 3GPP evaluation methodology (see Table 6.1.1 in [12]), since uniformly-distributed activations (traffic type 1: uniform distribution over 60 seconds) do not cause actual network overloads. Our evaluation framework accounts for the main M2M device power consumption levels (inactive, idle, Rx, and Tx) at all states of PRACH signaling procedure (Fig. 4b) and sheds light on the feasibility of candidate network overload control solutions. In particular, Fig. 4a suggests that a combination of M2M-specific backoff (larger non-standard BI values) and initial backoff (pre-backoff) may successfully alleviate congestion caused by highly-correlated beta-distributed M2M device activation patterns.

Further, the focus of our investigation shifts to the dedicated power consumption aspects of M2M devices. Currently, short paging cycles in 3GPP LTE may be highly sub-optimal for M2M devices, especially given lengthy M2M traffic inter-arrival times and the delay tolerant nature of many M2M applications. Hence, extending paging cycle durations in the idle state may help delay-tolerant devices sleep for longer periods of time, thus extending their battery lifetimes. The corresponding studies require an appropriate power consumption model (see Fig. 4c), which is capable of capturing typical M2M traffic patterns. Correspondingly, our results in [13] indicate that increasing the current maximum DRX (discontinuous reception) and paging cycle lengths would indeed lead to significant gains in the energy consumption (over 20x) of M2M devices (Fig. 4d).

If additional delay can be tolerated by M2M devices, D2D (device-to-device) communication techniques may further reduce the consumption of power. For instance, one M2M device may act as an aggregation point and relay data from other proximate M2M devices with poor communication link (to avoid excessive retransmissions and associated energy costs). D2D-based “client relay” mechanisms may dramatically reduce energy expenditures of cell-edge M2M devices, especially when those only send small data packets, and additionally help relieve a

Smart grid is one of the key M2M use cases incorporating a large number of metering devices that autonomously report their information to grid infrastructure. The motivating smart metering use case therefore serves as a valuable reference “massive M2M” scenario covering many characterizing M2M features.



**Figure 4.** Performance results on energy efficiency and small data: a) initial network entry performance; b) M2M device power consumption levels; c) proposed power consumption model; d) assessment of M2M improvements for energy efficiency; e) basic LTE frame structure; f) benefits of COBALT mechanism.

surge in uncontrolled near-simultaneous M2M transmissions. However, if adding more delay is not acceptable, such as in critical control applications, further data access improvements would need to be made, reducing PRACH latencies to a few milliseconds by potentially shortening the signaling sequence in Fig. 4b. Other areas for data access enhancements concern the size of typical M2M payloads (on the order of several bytes), as existing coding mechanisms in LTE are not optimal for short data blocks. Small data also creates inefficiency in control and channel estimation procedures, causing excessive overheads, as well as in existing frame structures and resource allocation schemes.

In more detail, Fig. 4e illustrates the current frame structure of LTE (for 5 MHz bandwidth) as a rectangular grid of resource blocks (RBs). The groups of RBs compose different data access channels, including periodic PRACH allocations and continuous physical uplink control channel (PUCCH) resources to carry the uplink control information. As both PUCCH and PRACH capacities may be very limited to serve small M2M data from numerous sources, we propose to allocate a part of physical uplink shared channel (PUSCH), which is otherwise employed for actual human-to-human (H2H) data transmissions, for a dedicated M2M use. Our scheme, named contention-based LTE transmission (COBALT), takes advantage of fewer LTE signaling messages and a simple colli-

sion resolution procedure (see a detailed description in [14]). It thus yields better utilization of network resources, lower latencies, and most importantly, significantly reduced power consumption for M2M devices (Fig. 4f).

### 3GPP STANDARDS UPDATE AND FUTURE OUTLOOK

We conclude this work by offering the current standardization perspective on cellular M2M technologies. Even though presently there are a number of industrial alliances and research projects pursuing their own standards and technologies in the M2M space, in this work we concentrate primarily on the cellular sector due to the major promise that it holds for the entire IoT industry.

M2M or machine-type communication was identified as one of the major topics for further enhancement in 3GPP due to an extreme diversity of prospective applications and corresponding consumer demands. In particular, the RAN technical specification group of 3GPP has been very active on M2M-related features across several releases of LTE as well as legacy cellular technologies. To this end, LTE Rel-11 has introduced improvements in handling a large number of M2M devices with delay-tolerant traffic, including RAN overload control mechanisms (such as enhanced access barring) and device power preference indication. In LTE Rel-12,

3GPP had a number of M2M-related study items (e.g. UEPCOP on power saving optimizations and SDDTE on signaling enhancements). The most significant outcomes of this work are a dedicated power saving mode (providing possible battery lifetime savings for M2M, potentially over 10 years of operation) and a new low-complexity device category, named Cat-0.

As the result of past 3GPP work, M2M devices can be flagged as low-priority and barred from accessing the network in case it is congested. In addition, scheduling prioritization and service differentiation mechanisms have been ratified to efficiently handle different traffic types by minimizing the impact of M2M data on H2H traffic. Further, low-complexity device modules support LTE operation with a single receiver chain and antenna, reduced peak data rates of 1 Mbps, and optional half-duplex operation. This work is planned to continue in Rel-13 along the lines of further coverage and power saving enhancements, as well as a new device category based on Cat-0, but with even lower complexity (reduced RF bandwidth of 1.4 MHz and maximum transmit power of 20 dBm compared to 23 dBm of today). Small data transmission and coverage enhancements are considered by some study/work items, together with architectural latency-related modifications.

Consequently, cellular networks are becoming increasingly equipped to support a diversity of M2M use cases and associated technical requirements. They already offer sufficient bandwidths and nearly-ubiquitous coverage (further improved by 15-20 dB in Rel-12/13), support full mobility, and provide precise location information. With the ongoing LTE evolution and the corresponding thorough standardization,

- M2M traffic can coexist efficiently with H2H mobile broadband applications [15].
- M2M modem complexity drops by 50 percent (Rel-12) and by up to 75 percent (Rel-13) compared to today's cheapest Cat-1 UE, thus resulting in lower modem costs.
- Battery lifetimes extend over 10 years for downlink delay-tolerant traffic (Rel-12) and other use cases (planned in Rel-13).

In turn, lower complexity UE categories (Rel-12 and the upcoming Rel-13 work) provide attractive cost saving opportunities for chipset manufacturers.

As our world is moving toward a fully-integrated networked society, where everything that may benefit from interacting and sharing information will become connected, 5G-grade M2M systems are expected sometime around the year 2020. They should generally deliver ubiquitous M2M connectivity with an edge-free experience, either on a stand-alone cellular M2M carrier or multiplexed with other services (e.g. mobile broadband). Given that wireless connectivity is becoming a new commodity, the same as water or electricity, M2M-based applications are likely to become a centerpiece of the emerging 5G ecosystem by enabling ubiquitous interworking between various communicating objects, as well as collection and sharing of the massive amounts of data. However, for the research community, these emerging future systems come with their associated unique challenges, such as extreme

heterogeneity of services, large-scale unattended wireless connectivity, and unprecedentedly large volumes of information to handle.

## REFERENCES

- [1] Finnish Strategic Centre for Science, Technology, and Innovation, Internet of Things Strategic Research Agenda (IoT-SRA), Sept. 2011.
- [2] I. Stojmenovic, "Machine-to-Machine Communications with In-Network Data Aggregation, Processing, and Actuation for Large-Scale Cyber-Physical Systems," *IEEE Internet of Things Journal*, vol. 1, no. 2, pp. 122-128, 2014.
- [3] J. Sachs, N. Bejar, P. Eimdahl, J. Melen, F. Militano, and P. Salmela, "Capillary Networks - A Smart Way to Get Things Connected," *Ericsson Review*, vol. 8, pp. 1-8, 2014.
- [4] D. Astely, E. Dahlman, G. Fodor, S. Parkvall, and J. Sachs, "LTE Release 12 and Beyond," *IEEE Commun. Mag.*, vol. 51, no. 7, pp. 154-160, 2013.
- [5] I. Vilajosana, J. Llosa, B. Martinez, M. Domingo-Prieto, A. Angles, and X. Vilajosana, "Bootstrapping Smart Cities through a Self-Sustainable Model Based on Big Data Flows," *IEEE Commun. Mag.*, vol. 51, no. 6, pp. 128-134, 2013.
- [6] IEEE, "TGah Functional Requirements and Evaluation Methodology".
- [7] ETSI GS LTN 001, "Low Throughput Networks (LTN): Use Cases, Functional Architecture and Protocols".
- [8] Machina Research, The Need for Low Cost, High Reach, Wide Area Connectivity for the Internet of Things, 2014.
- [9] ETSI EN 300 200-1, "Electromagnetic Compatibility and Radio Spectrum Matters (ERM); Short Range Devices (SRD); Radio Equipment to be Used in the 25 MHz to 1 000 MHz Frequency Range with Power Levels Ranging up to 500 mW; Part 1: Technical Characteristics and Test Methods".
- [10] K. Zheng, S. Ou, J. Alonso-Zarate, M. Dohler, F. Liu, and H. Zhu, "Challenges of Massive Access in Highly Dense LTE-Advanced Networks with Machine-to-Machine Communications," *IEEE Wireless Commun.*, vol. 21, no. 3, pp. 12-18, 2014.
- [11] M. Condoluci, M. Dohler, G. Araniti, A. Molinaro, and J. Sachs, "Enhanced Radio Access and Data Transmission Procedures Facilitating Industry-Compliant Machine-Type Communications over LTE-Based 5G Networks," in *IEEE Wireless Commun.*, 2015.
- [12] Study on RAN Improvements for Machine-Type Communications. 3GPP TR 37.868, 2011.
- [13] T. Tirronen, A. Larmo, J. Sachs, B. Lindoff, and N. Wiberg, "Machine-to-Machine Communication with Long-Term Evolution with Reduced Device Energy Consumption," *Trans. Emerging Telecommun. Technologies*, vol. 24, no. 4, pp. 413-426, 2013.
- [14] S. Andreev, A. Larmo, M. Gerasimenko, V. Petrov, O. Galinina, T. Tirronen, J. Torsner, and Y. Koucheryavy, "Efficient Small Data Access for Machine-Type Communications in LTE," in *Proc. IEEE International Conference on Communications (ICC)*, pp. 3569-3574, 2013.
- [15] M. Condoluci, M. Dohler, G. Araniti, A. Molinaro, and K. Zheng, "Toward 5G DenseNets: Architectural Advances for Effective Machine-Type Communications over Femtocells," *IEEE Commun. Mag.*, vol. 53, no. 1, pp. 134-141, 2015.

## BIOGRAPHIES

SERGEY ANDREEV (sergey.andreev@tut.fi) is a senior research scientist in the Department of Electronics and Communications Engineering at Tampere University of Technology, Finland. He received the specialist degree (2006) and the Cand.Sc. degree (2009), both from St. Petersburg State University of Aerospace Instrumentation, St. Petersburg, Russia, as well as the Ph.D. degree (2012) from Tampere University of Technology. He has (co-)authored more than 90 published research works on wireless communications, energy efficiency, heterogeneous networking, cooperative communications, and machine-to-machine applications.

OLGA GALININA (olga.galinina@tut.fi) is a Ph.D. candidate in the Department of Electronics and Communications Engineering at Tampere University of Technology, Finland. She received her B.Sc. and M.Sc. degrees in applied mathematics from the Department of Applied Mathematics, Faculty of Mechanics and Physics, St. Petersburg State Polytechnical University, Russia. Her research interests include applied mathematics and statistics, queuing theory and its applications, wireless networking and energy efficient systems, machine-to-machine and device-to-device communication.

ALEXANDER PYATTAEV (alexander.pyattaev@tut.fi) is a Ph.D. candidate in the Department of Electronics and Communications Engineering at Tampere University of Technology, Finland. He received his B.Sc. degree from St. Petersburg State University of Telecommunications, Russia, and his M.Sc. degree from Tampere University of Technology. He has publications on a variety of networking-related topics in internationally recognized venues, as well as several technology patents. His primary research interest lies in the area of future wireless networks: shared spectrum access, smart RAT selection and flexible, adaptive topologies.

MIKHAIL GERASIMENKO (mikhail.gerasimenko@tut.fi) is a researcher at Tampere University of Technology in the Department of Electronics and Communications Engineering. He received the specialist degree from Saint-Petersburg University of Telecommunications in 2011. In 2013 he obtained a master of science

Given that wireless connectivity is becoming a new commodity, the same as water or electricity, M2M-based applications are likely to become a centerpiece of the emerging 5G ecosystem by enabling ubiquitous interworking between various communicating objects, as well as collection and sharing of the massive amounts of data.

---

degree from Tampere University of Technology. He started his academic career in 2011, and has since been a (co-)author on multiple scientific journal and conference publications, as well as several patents. His main subjects of interest are wireless communications, machine-type communications, and heterogeneous networks.

TUOMAS TIRRONEN (tuomas.tirronen@ericsson.com) is a senior researcher at Ericsson Research, which he joined in 2012. He received his D.Sc. in communications engineering in 2010 from Aalto University. His current research interests include 4G and 5G wireless access technologies, Internet of Things, performance evaluation, radio protocols and resources. He is also active in 3GPP standardization work and innovation and patenting.

JOHAN TORSNER (johan.torsner@ericsson.com) is a research manager at Ericsson Research and is currently leading Ericsson's research activities in Finland. He joined Ericsson in 1998 and has held several positions within research and R&D. He has been deeply involved in the development and standardization of 3G and 4G systems and has filed over 100 patent applications. His current research interests include 4G evolution, 5G and machine-type communication.

JOACHIM SACHS (joachim.sachs@ericsson.com) is a principal researcher at Ericsson Research. He joined Ericsson in 1997 and has worked on a variety of topics in the area of wireless communication systems. He holds a diploma in electrical engineering from Aachen University (RWTH), and a doctorate in electrical engineering from the Technical University of Berlin, Germany. Since 1995 he has been active in the IEEE and the German VDE Information Technology Society (ITG), where he is currently co-chair of the technical committee on communication.

MISCHA DOHLER (mischadohler@kcl.ac.uk) is chair professor in wireless communications at King's College London, Director of the Centre for Telecommunications Research, co-founder and member of the Board of Directors of the smart city pioneer WorldSensing, Fellow (2014) and Distinguished Lecturer of the IEEE, and Editor-in-Chief of the *Transactions on Emerging Telecommunications Technologies*. He is a frequent keynote, panel, and tutorial speaker. He has pioneered several research fields, contributed to numerous wireless broadband, IoT/M2M, and cyber security standards, holds a dozen patents, organized and chaired numerous conferences, has more than 200 publications, and has authored several books. He has a citation h-index of 39 (top 1 percent). He acts as a policy, technology, and entrepreneurship adviser, examples being Richard Branson's Carbon War Room, House of Lords UK, UK Ministry BIS, EPSRC ICT Strategy Advisory Team, European Commission, ISO Smart City working group, and various start-ups. He is also an entrepreneur, angel investor, and passionate pianist, and he is fluent in six languages. He has spoken at TEDx. He has received coverage from national and international TV and radio, and his contributions have been featured by BBC News and the *Wall Street Journal*.

YEVGENI KOUCHERYAVY (yk@cs.tut.fi) is a full professor and lab director in the Department of Electronics and Communications Engineering of Tampere University of Technology (TUT), Finland. He received his Ph.D. degree (2004) from TUT. He is the author of numerous publications in the field of advanced wired and wireless networking and communications. His current research interests include various areas in heterogeneous wireless communication networks and systems, the Internet of Things and its standardization, as well as nanocommunications. He is an associate technical editor of *IEEE Communications Magazine* and an editor of *IEEE Communications Surveys and Tutorials*.



# WHAT CAN WIRELESS CELLULAR TECHNOLOGIES DO ABOUT THE UPCOMING SMART METERING TRAFFIC?

The authors characterize the current traffic generated by smart electricity meters, and they discuss the potential traffic requirements resulting from the introduction of enhanced smart meters, i.e. meters with PMU-like capabilities. Their study shows how GSM/GPRS and LTE cellular system performance behaves with current generation and next generation smart meter traffic, where it is clearly seen that the PMU data will seriously challenge these wireless systems.

*Jimmy J. Nielsen, Germán C. Madueño, Nuno K. Pratas, René B. Sørensen, Čedomir Stefanović, and Petar Popovski*

## ABSTRACT

The introduction of smart electricity meters with cellular radio interfaces has placed an additional load on wireless cellular networks. Currently, these meters are designed for low duty cycle billing and occasional system check, which generates low-rate sporadic traffic. As the number of distributed energy resources increases, household power will become more variable and thus unpredictable from the viewpoint of the distribution system operator (DSO). Therefore, it is expected that in the near future there will be an increase in the number of wide area measurement system (WAMS) devices with phasor measurement unit (PMU)-like capabilities in the distribution grid, thus allowing utilities to monitor the low voltage grid quality while providing information required for tighter grid control. From a communication standpoint, the traffic profile will change drastically toward higher data volumes and higher rates per device. In this paper we characterize the current traffic generated by smart electricity meters, and we discuss the potential traffic requirements resulting from the introduction of enhanced smart meters, i.e. meters with PMU-like capabilities. Our study shows how GSM/GPRS and LTE cellular system performance behaves with current generation and next generation smart meter traffic, where it is clearly seen that the PMU data will seriously challenge these wireless systems. We conclude by highlighting the possible solutions for upgrading the cellular standards, in order to cope with the upcoming smart metering traffic.

## INTRODUCTION

Smart power grids represent an important group of devices and applications within the umbrella of Internet of Things (IoT). In particular, the large number of network-connected smart electricity meters that already are or will be located

in all households and commercial/industrial locations are representative examples of IoT devices. At present, smart electricity meters are primarily used by electricity providers only for availability monitoring and billing. However, with the increasing number of distributed energy resources (DERs) such as wind turbines, solar panels, and electric vehicles, strong and sometimes unpredictable variations in power quality are introduced, leading to an increased need for monitoring and control. Specifically, distribution system operators (DSOs) need to be able to observe the circumstances in the low voltage (LV) power grid by introducing more frequently-sampled measurement points. Such wide area measurement systems (WAMS) exist already in the transmission grid, whereas in the distribution grid the DSOs rely mainly on open loop control beyond the substation level, i.e. without real-time feedback from consumers. As the number of DERs increases, this control loop must be closed by providing the feedback from measurements in the LV grid, enabling state estimation and prediction of the grid behavior, and ultimately ensuring stable operation [1]. It is expected that in the future LV grid, in addition to the traditional smart meter (SM), which so far has primarily been used for billing purposes with hourly or daily reports, another more advanced monitoring node will be needed, here referred to as an enhanced smart meter (eSM). The eSM is largely similar to a wide area measurement system (WAMS) node, as it integrates phasor measurement unit (PMU)-like capabilities; in other words, the eSM measures power quality parameters (such as power phasors) more frequently and in more detail compared to SMs [2]. While it is generally

expected that not all smart meter locations need to be equipped with eSM devices, the fraction of eSMs needed in the distribution grid to achieve satisfactory state information is still an open research question [3].

Today, SM devices are typically connected to the DSO backend system using either:

- A concentrator that gathers the data from the SMs in its neighborhood, via local Wi-Fi or PLC connections, and then relays it via a cellular or a wired connection to the DSO backend.
- Direct connections from each SM through the cellular network to the DSO backend.

While the concentrator based approaches reduce the load on the access networks by aggregating data locally, they are not suited for real-time monitoring from eSMs. The reason for this in PLC is the limited bandwidth, reliability, and the delays related to the daisy-chain topology. Wi-Fi is challenged due to the issues of shared spectrum and uncontrollable interference. Therefore, we assume that SMs and eSMs are equipped with cellular interfaces, so as to eliminate the potential delays, ease deployment, and reduce maintenance costs associated with network connectivity.

The traffic profile generated by smart meters falls into the category of machine-to-machine (M2M) traffic. A main characteristic of M2M traffic is that it consists of transmissions of small

## COMMUNICATIONS STANDARDS

*The authors are with Aalborg University.*

Since there are differences in which use cases and applications are offered by the DSO or electricity retail company and which of those the individual customers are using, a one size fits all traffic model does not exist.

Event	Frequency (events per meter)
On-demand meter read requests	25/1000 per day
Meter capped energy mode request	5 per year
DR load management request to HAN devices	15/1000 per day
HAN device join/unjoin	5 per year
Real-time price (RTP) update	96 per day
Metrology firmware update	4 per year
Metrology program update	4 per year
NIC firmware update	4 per year
NIC program update	4 per year

**Table 1.** Assumptions for deriving traffic model.

amounts of data from a very high number of devices, differing significantly from the bursty and high data rate traffic patterns in human-oriented services, and instead requiring network reliability and availability. Further, M2M traffic is more demanding in the uplink and less focused on downlink performance, as typical use cases encompass monitoring and control functions.

With LTE gaining an increasing market share, it is expected that within a number of years one of the existing 2G or 3G systems will be taken out of service in order to re-harvest the spectrum to use for newer technologies. Current reports on active M2M cellular devices indicate that 64 percent of them are GSM/GPRS-only, 25 percent both 3G and GSM/GPRS compatible, 10 percent 3G-only, and only 1 percent is LTE capable [4]. It is clear that GSM/GPRS (hereafter denoted GPRS) dominates the M2M industry, therefore in this article we analyze how well this technology can support the connectivity demands of SM and eSM devices. Given the promise of LTE, we also investigate its potential use as the access network for eSM devices.

Specifically, in this article we offer the following four contributions:

- Extraction and classification of smart meter traffic models from relevant specifications, as well as predicted future traffic growth.
- Comprehensive simulation model of radio access systems that includes all phases in the access, in contrast to [5] and the NIST PAP2 guidelines for assessing wireless standards for smart grid application v1.0 that use simplified models.
- Quantitative assessment of how many smart meter devices can be supported in cellular systems, comparing the simplified and comprehensive simulation model results.
- Recommendations for standardization and a future roadmap of radio access technologies.

The rest of this article is organized as follows.

We characterize the traffic models of the SM and eSM devices. We describe the access bottlenecks in the cellular access reservation protocol and provide numerical results that show how the proposed traffic models affect the performance of GPRS and LTE networks. Then we provide guidelines that future cellular network standards should take into account when designing the network system. Finally, we conclude the article with the main take-home points.

## SMART METER TRAFFIC MODEL

In the literature there are different examples of traffic models for traditional smart meters. Of these, the OpenSG *Smart Grid Networks System Requirements Specification* (described in [6]) from the Utilities Communications Architecture (UCA) user group is the most coherent and detailed network requirement specification, and it has therefore been used in this article as input for the SM traffic model. The UCA OpenSG is a relevant consortium of 190 companies, and the considered smart grid use cases are in line with those studied by other standardization organizations such as ETSI and USEF. Since there are differences in which use cases and applications are offered by the DSO or electricity retail company and which of those the individual customers are using, a one size fits all traffic model does not exist. In the following we consider a comprehensive configuration where all use cases that involve communication from smart meters to the core network will be in operation and note that actual deployments with different configurations may lead to different results. To calculate the message frequency in the uplink SM traffic model, the event occurrence frequencies listed in Table 1 have been used. In addition to the values listed in Table 1, we assume that a commercial/industrial SM sends a 2400 bytes meter reading packet every hour, whereas a residential SM sends a 1200 bytes report every four hours.

The SM uplink traffic model, resulting from the above assumptions, is presented in Fig. 1. The gray boxes represent the different use cases and the boxes span the latency and payload size requirements of the corresponding messages. The white box represents eSM traffic (defined later). Nearly all use cases have reliability requirements of 98 percent, with the exceptions being two alarm messages in the IDCS use case requiring 99 percent, and the periodic meter reading, which has time-dependent reliability requirements ranging from 98 percent to 99.5 percent. In relation to the figure, Table 2 shows the average estimated uplink/downlink bandwidth for each use case.

The  $\lambda$ -values in Fig. 1 show the number of generated messages per day per SM. The use cases grouped in the dash-dotted box transmit very infrequently with a combined rate of only approximately 0.5 messages per day. Further, they are relatively similar in terms of latency and payload size. In addition to this group, two other OpenSG use cases from the figure stand out, namely the real-time pricing (RTP) that causes 96 messages per day, and the periodic meter reading at the top right. For periodic meter

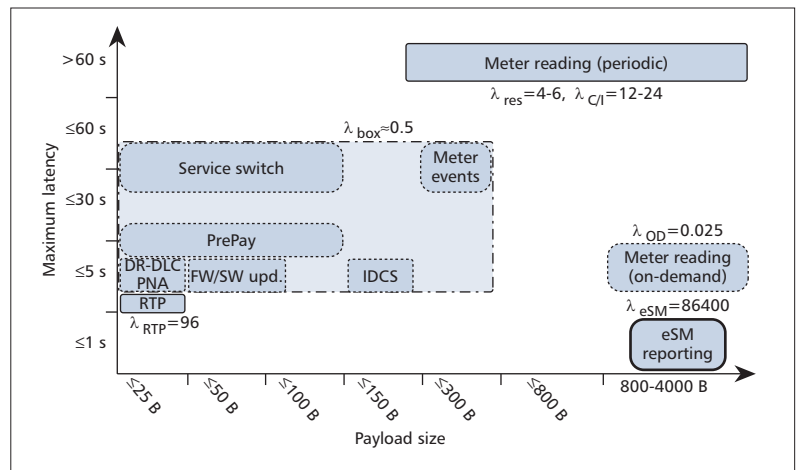
readings, a commercial/industrial (C/I) SM sends reports more often than a residential SM. Notice for the eSM reporting that, in addition to the stricter latency requirement of  $\leq 1$  sec, the number of generated messages per day is many orders of magnitude higher than any of the SM use cases.

Table 2 shows that the raw data rate requirements of SMs with a default reporting interval (RI) are quite modest, with an average uplink data rate of approximately 13.4 KB per day per SM and an average downlink data rate of approximately 40.4 KB per day per SM. While the total downlink data rate is actually higher than the uplink, it is constituted primarily of software updates, which are large low-priority data transfers that occur infrequently during the night, when they do not interfere with the day-to-day operation of the smart grid. Given the modest traffic requirements, it is expected that GPRS networks, which are deployed ubiquitously and offer reliable coverage but are gradually becoming less suitable for human-oriented traffic, can easily satisfy the default SM traffic requirements.

Further, an option to increase observability in the power grid is to reduce the meter reading reporting interval. We investigate how capable the current cellular systems are to support this earlier, when the report packet sizes are respectively 300 bytes and 600 bytes for residential and commercial/industrial, and reporting intervals range from 5 min, 1 min, 30 sec, to 15 sec. As shown in Table 2, in the case of these reduced RIs, the uplink data rate requirements become much larger than in the downlink.

## ENHANCED SMART METER TRAFFIC MODEL

The eSM is a PMU-like device for the distribution grid, which is able to measure voltage and current phasors. However, it has less strict real-time requirements than transmission grid PMUs, since it is used to increase observability rather than for protection purposes. Being deployed not only in DSO substations, but also in and close to customer homes, the eSM reports measurements through cellular networks, since this allows a mobile network operator to prioritize and dedicate resources to eSM traffic, thus achieving QoS, which may not be possible with third party consumer-grade wired Internet connections. Phasor measurements can be used on different time scales, ranging from a few milliseconds (e.g. for protective relays) up to several seconds (e.g. real-time monitoring and state estimation) [1]. The eSMs are intended to improve observability and enable state estimation and real-time control [1], with the suggested lower bound of 1 sec for the reporting interval [7]. Since eSM features and requirements are not yet standardized, the eSM traffic model considered in this study is based on the requirements of the transmission grid PMU and WAMS related standards, IEEE 1588, IEEE C37.118, and IEC 68150. Specifically, we assume that every second an eSM sends a measurement report that consist of concatenated PMU measurements (50 Hz sample rate) from the preceding 1 sec measurement interval. The samples are, as specified in



**Figure 1.** Classification of OpenSmartGrid traffic originating from an SM.  $\lambda$  -values show the number of generated messages per day per device. Use case short names: demand response–direct load control (DR-DLC); premise network administration (PNA); firmware and software updates (FW/SW upd.); real-time price (RTP); islanded distributed customer storage (IDCS).

PMU standards IEEE 1588 and C37.118, timestamped using GPS time precision. Assuming that the floating point PMU frame format from IEEE 1588 is used, and that each sample covers six phasors, one analog value and one digital value, each PMU frame accounts for 76 bytes. Adding the UDP header (8 bytes) and IPv6 header (40 bytes) to each report of 50 PMU samples, an eSM packet is 3848 bytes, with a bit rate of 30.8 kbit/s. Since it may be an exaggeration to send all 50 PMU samples per measurement interval, we also consider in our performance analysis the case of eSM reduced report sizes.

## CELLULAR SYSTEMS PERFORMANCE

From the communications perspective, it is important to investigate which cellular technologies can support the current billing-only smart meter use cases, but also the use cases/services that go beyond the current ones. In [5] and NIST PAP2 guidelines for assessing wireless standards for smart grid application v1.0, performance analyses were carried out to determine the number of smart grid devices supported by different wireless technologies. However, they only evaluated the data capacity of the systems and neglected to account for the bottlenecks in the access reservation protocol used in cellular systems. As shown in [8], the access reservation bottlenecks are particularly prone to exposure with M2M traffic such as smart grid traffic, meaning that a pure data capacity based analysis may lead to overly optimistic results. Therefore, our analysis will include all aspects of the access reservation procedure and compare those results to a data capacity only analysis. For the analysis we will consider the traffic patterns for SM and eSM devices described earlier. From those traffic models it is clear that the communication requirements of these two device types are orders of magnitudes apart in terms of message frequency and bandwidth, meaning that for eSM deployment a more capable technology than

Use case	RI	Downlink		Uplink			
		Default	Default	5 min	1 min	30 sec	15 sec
Meter reading		1.25	11K	95 K	475 K	950 K	1.9M
Service switch		3	6	6	6	6	6
PrePay		3.5	8	8	8	8	8
Meter events		0	50	50	50	50	50
Islanded distr. cust. storage		2	5	5	5	5	5
DR-DLC		400	0.5	0.5	0.5	0.5	0.5
Premise network admin		1	1	1	1	1	1
Price		10 K	2.4K	2.4K	2.4K	2.4K	2.4K
Firmware/program update		30 K	5	5	5	5	5
Total		40.4K	13.4K	97K	477K	952K	1.9M

**Table 2.** Average downlink/uplink raw data rate as (bytes/meter/day) for the considered use cases. Default value of RI is four hours for residential and one hour for commercial/industrial SMs.

GPRS is needed. With its integrated PMU unit, the eSM is already a more complex and expensive device than the SM, and since fewer eSMs than SMs will be needed, a higher unit price can be better tolerated, and thus we will assume that the eSM uses LTE.

### ACCESS RESERVATION PROTOCOL OPERATION AND LIMITATIONS

In cellular networks, a device with no active connection to the network first has to establish one in order to perform data transmission. This is accomplished via an access reservation protocol, which in general consists of three main stages: random access, granting access, and data access. In the first stage the cellular devices perform a random access request in one of the random access opportunities (RAOs). In the second stage, the base station grants access to the network if:

- The random access request is received without error by the base station.
- No other device has transmitted in the same RAO (i.e. collision free).
- There are data resources available to the device.

Otherwise, the access reservation procedure must be restarted and the device will transmit a new random access request until it is granted by the base station or until the maximum number of retransmissions is reached and the request fails. In GPRS there are 217 RAOs/s per carrier while in LTE there are 10.8k RAOs/s.<sup>1</sup> On the other hand, only 32 grants/s and 3k grants/s are offered in GPRS and LTE, respectively [8, 9]. Therefore, when the random access stage is heavily loaded, the grant stage becomes a decisive limitation in cellular networks. Furthermore,

<sup>1</sup> Assuming the contention resources occur every 5 ms, each with 54 contention preambles available.

in GPRS and LTE the data stage is not only limited by the amount of the actual data resources, but also by the amount of the uplink identifiers used to coordinate transmissions from active devices, which limits the amount of simultaneously active M2M communication links.

### OUTAGE PERFORMANCE EVALUATION

To evaluate the performance of the cellular access, we used the outage rate, i.e. the probability of a device failing to deliver a report before the report deadline expires, while accounting for the access reservation protocol. This can be regarded as a measure of cellular access reliability, which is the paramount performance indicator for the wide-area distribution supervision and control applications [10]. To evaluate the outage, we used event-driven simulators developed in MATLAB, which cover the complete access reservation protocol, as defined in 3GPP Release 12. In particular, the GPRS simulator considers the amount of available access granted messages in the access granted channel (AGCH), with a typical configuration of 28 AGCH/s [8], the limited number of identifiers used to coordinate the uplink transmissions, i.e. the uplink stage flag (USF), and the amount of data resources available. The LTE simulator considers the restricted amount of access grant messages (RAR messages) due to the physical downlink control channel (PDCCH) limitations [11, 12]. Finally, both in the GPRS and the LTE simulator, data resources are shared with the signaling required in the access reservation procedure and the actual data transmissions.

The evaluation scenario is set in a single cell with 1000 m radius, which includes 4500 smart meters [8], of which 90 percent correspond to residential customers and the remaining 10 percent to commercial/industrial customers. In the case of GPRS, we consider a single carrier corresponding to a 200 KHz system. The considered LTE bandwidth is 1.4 MHz (six PRBs), in line with the reduced capabilities for LTE devices [13]. In addition, the control channel and data channel probability of error are, respectively,  $10^{-2}$  and  $10^{-1}$  [9]. In both systems, we assume the devices always transmit with the highest modulation scheme available, in order to focus the evaluation on the performance of the access reservation protocol. In these conditions, we observed that the SM traffic, provided earlier, is supported by both GPRS and LTE with near 0 percent outage, as the total number of messages per hour from each SM only amounts to approximately 125.

We start by considering for GPRS the scenario of reducing SM report intervals (RI), as described earlier. Figure 2 depicts the outage probability for increasing the number of SMs and different RIs. Taking as reference a cell population of 4500 SMs, we can see that for  $RI > 5$  min, GPRS can provide a significant increase in the distribution network observability from hourly intervals to every five minutes. For smaller report intervals to be supported in GPRS, then, the options are either to reduce the cell size and/or increase the number of carriers.

We proceed by considering in Fig. 3 the cellular network outage as a function of the eSM

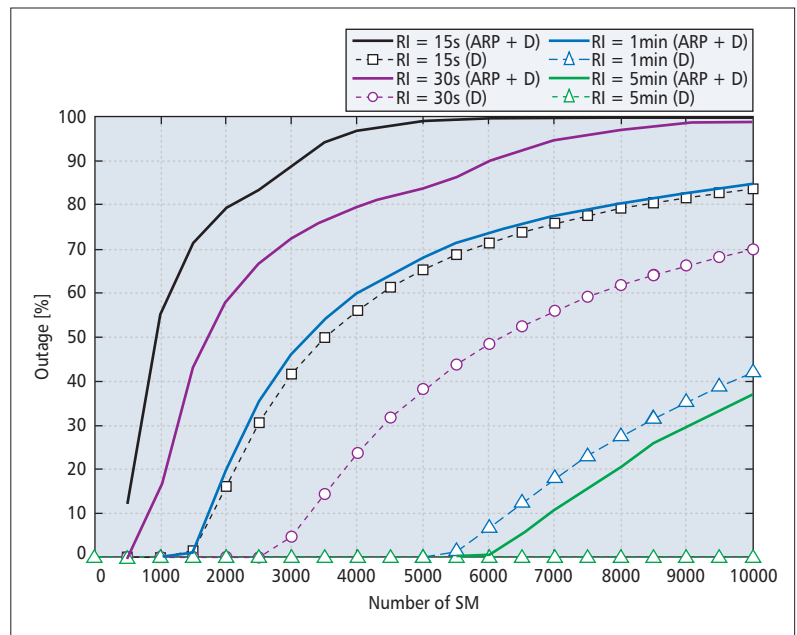
penetration, i.e. of how many eSMs are deployed per every 100 smart meter locations. As described earlier, each eSM report contains 50 samples of the power phasors measured since the last report with an expected payload of 3848 bytes. Since this large payload has severe implications on cellular network performance, we also consider the impact of smaller payloads on system performance, which can be motivated by the introduction of pre-processing to extract statistics, data compression, and/or a reduced number of samples. Specifically, we consider reduced report sizes (RS) of 3848, 400, and 115 bytes, where the last two values correspond respectively to a payload reduction of approximately 10 percent and 3 percent of the original payload size.

The outage results for LTE and GPRS are shown in Fig. 3. We note that GPRS is not able to support eSM traffic irrespective of the chosen RS, while LTE for RS of 3848 bytes only supports up to two percent eSM penetration. When a 10 MHz bandwidth is completely dedicated in LTE to serve the eSM traffic, then it is possible to reach 30 percent of penetration with less than 10 percent outage, which means a large amount of resources dedicated to a potentially low profit application. On the other hand, if we assume lower RS, already at 400 bytes LTE supports up to 20 percent of eSMs. Further, when comparing the results that correspond to the case when only the data phase is taken into account with the results obtained by considering the access reservation phase as well, it can be observed that the access reservation protocol impacts the number of supported eSMs. Particularly, the limitations of the access reservation protocol become substantial as the report size decreases, and it could be shown that this is mainly due to the lack of access grant messages required to complete the access reservation procedure. Note that this effect has been overlooked in previous works [5, 14].

The presented results allow us to conclude that the RS of the eSM nodes must be small to support a high percentage of nodes. In addition, we emphasize that small data traffic cannot be analyzed only in terms of the system data capacity, but that the bottlenecks of the access reservation protocol itself must be considered, as observed in the gap between the two types of analysis depicted in Fig. 3. We conclude by noting that in practice, when deploying eSMs, due to the required communication reliability, good coverage should be ensured, e.g. by careful selection of the placement location and/or by adding an external antenna if needed. In the above presented study, it is assumed that all SMs and eSMs are under a cellular coverage.

## STANDARDIZATION OUTLOOK

Although the traffic resulting from smart meters can be easily accommodated into current cellular systems, the same is not observed for the traffic generated by the eSM. In the following, we discuss the challenges and possible solutions that need to be tackled by standardization bodies to ensure that the observability of the distribution network can be improved efficiently.



**Figure 2.** GPRS outage evaluation for increasing number of SM with different report interval values and RS = 300 bytes for residential and RS = 600 bytes for commercial/industrial, where ARP+D denotes the access reservation protocol plus data phase, while D denotes only data phase.

## SMART METER

The inclusion of additional phasor measurement units into the distribution grid, so as to increase its observability, is being discussed specifically at the last mile to the customer premises [3]. Currently, it is not yet clear if that will imply the same level of detail (in number of samples and report frequency) as in the transmission grid PMUs, where the reporting is done by SCADA over dedicated wired links.

As discussed earlier, if the eSMs generate the same amount of traffic as transmission grid PMUs, then the cellular networks will require an extensive overhaul to be able to support both eSM and human centric traffic, leading to substantial investment in the cellular infrastructure. On the other hand, eSMs will most likely be lighter versions of PMUs, both sampling and reporting less frequently. Therefore, if local processing and compression of the monitoring data is allowed and/or the required level of detail lowered, then the amount of generated traffic will be much lower. Another viable option, as discussed earlier, is to increase the report frequency of current smart meters without introducing local PMU functionality. The generated small packets could then be handled by the network, as long as the bottlenecks at the access protocol level are addressed.

It seems likely that the standardization for the eSM's PMU functionality falls within the scope of the IEEE C37.118 and IEC 68150 standards, since these specify the measurement and communications requirements for traditional PMU units. Therefore, it is of paramount importance that standardization bodies reach a consensus on the eSM communication requirements allowing the affected stakeholders to make informed actions.

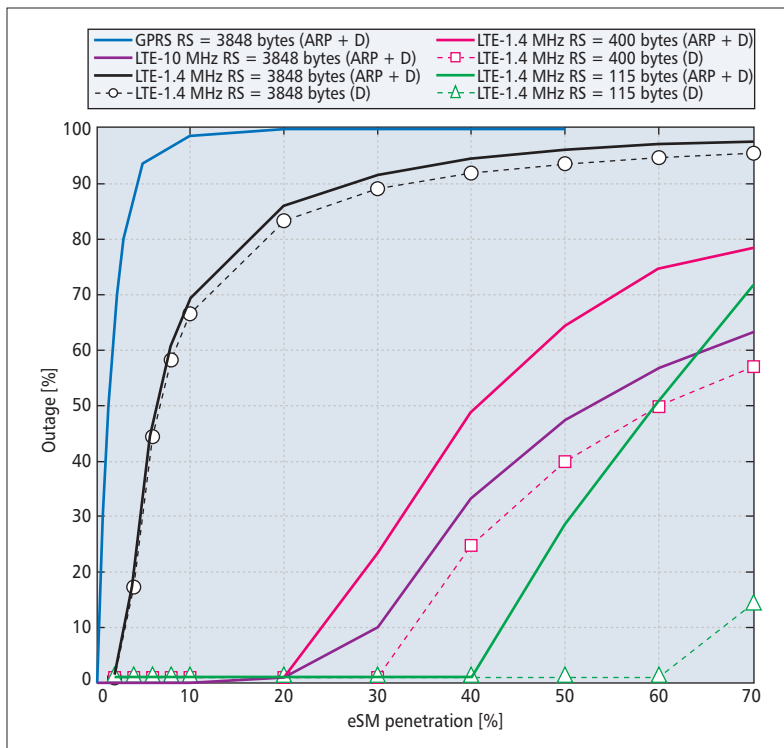


Figure 3. LTE and GPRS outage evaluation for increasing penetration of eSMs, where ARP+D denotes the access reservation protocol plus data phase, while D denotes only data phase.

### CELLULAR NETWORK

In 3GPP, the standardization body responsible for the cellular air interface and core network functionality, there are two activities that will affect how the traffic from SMs and eSMs will be handled [13].

We start by noting that, although GRPS is seen as an outdated communication technology [13], there is an ongoing effort to continue to reengineer GPRS to serve M2M applications, in which the SM traffic can be classified. One of the goals of this initiative is to achieve<sup>2</sup> 160 bit/s. Concurrently, there is a push from the industry (both utilities and vendors) to keep GPRS networks and their associated resources active, while facing the pressure to re-harvest the GPRS spectrum to be used in the next cellular network generation. A viable solution to keep the GPRS connectivity is to *virtualize* its air interface into the next generation cellular systems.

The second effort is to define a low complexity LTE user equipment category with respect to the cellular interface, which supports reduced bandwidth and transmit power while extending coverage operation [13]. Specifically, the goal of reduced bandwidth is to specify 1.4 MHz operation within any LTE system bandwidth, allowing operators to multiplex reduced bandwidth MTC devices and regular devices within their existing LTE deployments. In terms of extended coverage, the goal is to improve the coverage of delay-tolerant MTC devices by 15 dB, thereby allowing operators to reach MTC devices in poor coverage conditions, such as smart meters located in basements [13].

<sup>2</sup> Considering the minimum SDU size, i.e. 80 bytes, with 4 seconds latency.

To further improve the support of the traffic generated by SMs and eSMs with very low duty cycle and latency requirements in the order of seconds, the inclusion of *periodic reporting and discontinuous transmission* functionality into cellular standards should be considered. In here, the network provides periodic communication resources so that devices can perform their short data transmission. This allows devices to go to sleep and save energy, since they have prior knowledge of when the next transmission time slot can occur. A solution based on this concept has been proposed through the reengineering of the LTE access protocol [9].

To cope with eSM traffic demands and increase the network capacity, *localized aggregation of traffic* should be considered. In this solution the traffic generated by multiple SMs and eSMs in a geographical area could be aggregated, at eSMs or cellular relays, and then trunked to the cellular network [15]. With aggregation and relaying, contention pressure could be decreased at the base station and single link connection could be improved, providing connectivity and coverage enhancements to SMs and eSMs with poor propagation conditions.

Finally, to support *massive asynchronous access of small packet transmissions*, access reservation protocols in cellular systems are just the first step of asynchronous access to the network. After it has been completed, the device starts exchanging signaling information via the higher layers with the entities in the core network, which leads to high signaling overhead and possible air interface and core network congestion. Although there are already efforts to reduce the signaling exchanges with the core network [11], when the payloads are small enough, the facility to perform the data transmission already in the third step of the access reservation protocol should be in place.

### CONCLUSION

In this article we have evaluated two approaches to increase the observability of the network:

- Decreasing the report interval of the meter reading.
- Introducing enhanced smart meters with phasor measurement units (PMUs).

We provided details on the characteristics of the traffic generated by smart meters and enhanced smart meters and highlighted the associated challenges in supporting it from a cellular network point of view. The obtained results show that GPRS can support traditional smart meter traffic, as well as more frequent measurements down to 5 min report intervals. Further, it is shown that LTE can support distribution grid PMUs if the report payloads are appropriately dimensioned. These results can be used as input for both smart meter and cellular system standardization bodies to enable the introduction of current and future smart grid devices into cellular networks. The current main open issue is the uncertainty associated with the eSM communication requirements, which will lead to different cellular system optimizations.

## ACKNOWLEDGMENT

This work is partially funded by the EU under grant agreement no. 619437 “SUNSEED.” The SUNSEED project is a joint undertaking of nine partner institutions, and their contributions are fully acknowledged. The work of Č. Stefanović was supported by the Danish Council for Independent Research, grant no. DFF-4005-00281.

## REFERENCES

- [1] J. Sexauer, P. Javanbakht, and S. Mohagheghi, “Phasor Measurement Units for the Distribution Grid: Necessity and Benefits,” *2013 IEEE PES Innovative Smart Grid Technologies (ISGT)*, IEEE, 2013, pp. 1–6.
- [2] “TR 102 935: M2M Applicability of M2M Architecture to Smart Grid Networks Impact of Smart Grids on M2M Platform V2.1.1,” ETSI, Tech. Rep., 2012.
- [3] Y.-F. Huang *et al.*, “State Estimation in Electric Power Grids: Meeting New Challenges Presented by the Requirements of the Future Grid,” *IEEE Signal Proc. Mag.*, vol. 29, no. 5, 2012, pp. 33–43.
- [4] Ericsson AB, “Interim Ericsson Mobility Report,” Feb. 2014.
- [5] C. Hagerling, C. Ide, and C. Wietfeld, “Coverage and Capacity Analysis of Wireless M2M Technologies for Smart Distribution Grid Services,” *2014 IEEE Int'l. Conf. Smart Grid Communications (SmartGridComm)*, IEEE, 2014, pp. 36–73.
- [6] E. Hossain, Z. Han, and H. V. Poor, *Smart Grid Communications and Networking*, Cambridge University Press, 2012.
- [7] M. Adamiak, B. Kasztenny, and W. Premerlani, “Synchrophasors: Definition, Measurement, and Application,” *GE Digital Energy Protection and Control J.*, 2006.
- [8] G. C. Madueno, C. Stefanovic, and P. Popovski, “Reengineering GSM/GPRS Towards a Dedicated Network for Massive Smart Metering,” *2014 IEEE Int'l. Conf. Smart Grid Communications (SmartGridComm)*, Nov 2014, pp. 338–43.
- [9] G. Madueno, C. Stefanovic, and P. Popovski, “Reliable Reporting for Massive M2M Communications with Periodic Resource Pooling,” *IEEE Wireless Commun. Lett.*, vol. 3, no. 4, Aug. 2014, pp. 429–32.
- [10] S. Goel *et al.*, *IEEE Vision for Smart Grid Communications: 2030 and Beyond*, IEEE, 2013.
- [11] “TR 23.887 — Study on Machine-Type Communications (MTC) and Other Mobile Data Applications Communications Enhancements,” 3GPP, Tech. Rep., 2013.
- [12] P. Osti *et al.*, “Analysis of PDCCH Performance for M2M Traffic in LTE,” *IEEE Trans. Vehic. Tech.*, vol. 63, no. 9, Nov. 2014, pp. 4357–71.
- [13] “Overview of 3GPP Release 13,” 3GPP, Tech. Rep., 2015.
- [14] NIST, “NIST PAP2 Guidelines for Assessing Wireless Standards for Smart Grid Application,” 2012.
- [15] G. Rigazzi *et al.*, “Aggregation and Trunking of M2M Traffic via D2D Connections,” *2015 IEEE Int'l. Conf. Commun.*, 2015.

## BIOGRAPHIES

JIMMY JESSEN NIELSEN (jijn@es.aau.dk) obtained his M.Sc. in computer engineering in 2007 and the Ph.D. in wireless communications from Aalborg University in 2011. He is currently a postdoctoral researcher in the Department of Electronic Systems, Aalborg University. His research interests are in the areas of traffic models for M2M systems, reliable communications, and performance analysis of communication systems. He has served as reviewer for various IEEE journals and conferences.

GERMÁN C. MADUENO (gco@es.aau.dk) obtained his M.Sc. in mobile communications and his Ph.D. in reliable radio access for massive machine-to-machine (M2M) communications from Aalborg University. He is currently working as a research assistant in the MassM2M research group at Aalborg University, where he is also responsible for the M2M laboratory. He has served as a technical program committee member for JSAC. He also served as reviewer for *IEEE Wireless Communication Letters*, *IEEE Communications Magazine*, and various IEEE conferences.

NUNO K. PRATAS (nup@es.aau.dk) is a postdoctoral researcher on wireless communications in the Department of Electronic Systems, Aalborg University. He has been awarded twice with the best student conference paper award. He has authored and co-authored more than 30 publications in conferences, journals, books, and patent applications. His research interests are in the areas of wireless communications, networks, and the development of analysis tools for communication systems currently focused on machine-to-machine and device-to-device applications.

RENÉ B. SØRENSEN obtained his B.Sc. in computer engineering from Aalborg University in 2014. He expects to receive his M.Sc. in wireless communications from Aalborg University in 2016. He is currently working as a student assistant in the Department of Electronic Systems, Aalborg University. His interests include communication systems, protocol design, and coding theory.

ČEDOMIR STEFANOVIĆ (cs@es.aau.dk) received Dipl.-Ing., Mr.-Ing. and Ph.D. degrees in electrical engineering from the University of Novi Sad, Serbia. Since 2012 he has been affiliated with the Department of Electronic Systems, Aalborg University. His research interests are in the areas of communication theory, including design and analysis of enhanced random access mechanisms and distributed algorithms for wireless ad-hoc networks.

PETAR POPOVSKI (petarp@es.aau.dk) received his Dipl.-Ing. (1997) and Magister Ing. (2000) in communication engineering from Sts. Cyril and Methodius University, Skopje, Republic of Macedonia, and his Ph.D. from Aalborg University (2004). He is currently a professor at Aalborg University. He is an editor for *IEEE Transactions on Communications* and has served in the past as an editor for *IEEE Communications Letters*, the *IEEE JSAC Cognitive Radio Series*, and *IEEE Transactions on Wireless Communications*. He is a Steering Committee member for the *IEEE Internet of Things Journal* and *IEEE SmartGridComm*. In 2015 he received a Consolidator Grant from the European Research Council. His research interests are in the areas of communication theory, wireless communications, and networking.

The obtained results show that GPRS can support traditional smart meter traffic, as well as more frequent measurements down to 5 min report intervals. Further, it is shown that LTE can support distribution grid PMUs if the report payloads are appropriately dimensioned.

---

# PUBLISH/SUBSCRIBE-ENABLED SOFTWARE DEFINED NETWORKING FOR EFFICIENT AND SCALABLE IoT COMMUNICATIONS

The authors outline the most important issues related to standardization efforts, mobility of objects, networking and gateway access, and QoS support. They describe a novel IoT network architecture that integrates SDN and DDS middleware. The proposed architecture will improve service delivery of IoT systems and will bring flexibility to the network.

*Akram Hakiri, Pascal Berthou, Aniruddha Gokhale, and Slim Abdellatif*

---

## ABSTRACT

The Internet of Things (IoT) is the result of many different enabling technologies such as embedded systems, wireless sensor networks, cloud computing, big-data, etc., which are used to gather, process, infer, and transmit data. Combining all these technologies requires a research effort to address all the challenges of these technologies, especially for sensing and delivering information from the physical world to cloud-hosted services. In this article we outline the most important issues related to standardization efforts, mobility of objects, networking and gateway access, and QoS support. In particular, we describe a novel IoT network architecture that integrates software defined networking (SDN) and data distribution service (DDS) middleware. The proposed architecture will improve service delivery of IoT systems and will bring flexibility to the network.

## COMMUNICATIONS STANDARDS

## INTRODUCTION

The widespread evolution of the Internet of Things (IoT) concept imposes complex requirements on both the underlying networks and communication mechanisms between heterogeneous smart objects that communicate over the Internet. The current estimate for the number of deployed things are on the order of 50 billion connected devices, and by the year 2020 there will be 1000 times more connected mobile devices, all with different requirements, which will have a big impact on how people will interact with the surrounding things. However, this is a very difficult goal to achieve because today's network is limited in its ability to address the requirements of even current IoT deployments. In the future, leveraging the IoT potential will be the key enabler to the creation of different applications areas, including smart cities, e-health, industrial, transportation, retail, safety, and environmental services.

In order to achieve this goal, it is necessary to

provide an efficient network architecture that aids in developing many scalable, interoperable, and predictable IoT applications. This network architecture must be flexible enough to be reprogrammed in accordance with any change in IoT application needs. IoT usually adopts different communication schemes compared to the traditional Internet, stemming primarily from the variety of resource-constrained hardware platforms. In particular, an IoT network is a highly unstructured cloud of wireless and potentially mobile devices, whose states change dynamically (e.g. sleeping and waking up, connected and/or disconnected), and whose locations and speeds change as well. These smart applications should run in resource-constrained devices with limited computation resources, small data-storage capabilities, and low-power consumption. Since these devices sense the environment and communicate, they should be autonomic to respond to different scenarios without human interaction.

Since this large number of smart objects are being provided as services connected to the Internet, the current IPv4 addressing no longer suffices since it has reached its limit, and hence efforts are underway to deploy IPv6 addresses. Additionally, the heterogeneous and dynamic aspects of IoT systems pose major challenges for the underlying network by requiring support for handling heterogeneity, dynamic changes, device discovery, and context-awareness. Contemporary network protocols are designed in isolation to solve a specific problem and are often retrofitted to

address a new requirement. Unfortunately, this approach has limitations for IoT. Moreover, existing protocols often lack the right abstractions that address the requirements of IoT communication. Thus, if network resource utilization is a concern, the network must be flexible enough to be reprogrammed in accordance with any change in IoT application needs. Current network provisioning approaches neither address the dynamicity of IoT applications nor care about resource utilization. Recently, software defined networking (SDN) [1] was introduced to deliver dramatic improvements in network agility and flexibility. The SDN paradigm is thus a promising solution to solve the resource management needs of the IoT environment, but it cannot address the heterogeneous and dynamic needs of IoT applications.

At the distributed systems-level, the Web of things (WoT) [2] concept was introduced to alleviate the heterogeneity issues by allowing smart devices to speak the same language based on open web technologies such as HTTP and REST principles for information sharing. However, to fully explore the potential of WoT, many challenges remain unresolved, such as security and scalability. Moreover, WoT cannot address the need for distributed, peer-to-peer, publish/subscribe semantics, which are a key requirement for IoT applications.

Many solutions at the middleware level were introduced for IoT applications. LinkSmart [3] supports resource discovery, description, and access based on XML and web protocols. Open-IoT [4] realizes on-demand access to cloud-

---

*Akram Hakiri is with  
ISSAT Mateur, Tunisia.*

*Pascal Berthou and Slim  
Abdellatif are with  
LAAS-CNRS, France.*

*Aniruddha Gokhale is  
with Vanderbilt University,  
USA.*



based IoT services through Internet-connected smart objects. These contributions focus on the upper layer problems such as enabling IP-based radio communication over middleware, but without paying attention to the underlying network. The Object Management Group's data distribution service (DDS) provides real-time, scalable, data-centric publish/subscribe capabilities. However, although DDS has been used to develop many scalable, efficient, and predictable applications at a local area network-scale, its QoS mechanisms are rarely propagated to the network layer, which impedes its use across a wide area scale. Moreover, DDS is not a network-level solution.

To address both the network and distributed systems-level concerns in IoT, we propose combining ideas from SDN with the message-oriented publish/subscribe DDS middleware to define a powerful and simple abstract layer that is independent of the specific networking protocols and technology.

The remainder of this article is organized as follows. We briefly introduce the concepts of SDN and DDS. We introduce the architecture of IoT and provide a glimpse into the most important open issues related to the deployment of IoT networks. We describe the architecture of our SDN solution for efficient use of IoT and discuss the role of the proposed solution in solving the above mentioned issues. Finally, we provide concluding remarks describing potential future directions and open research problems in this realm.

## BACKGROUND

### SOFTWARE DEFINED NETWORKING (SDN) AND NETWORK FUNCTION VIRTUALIZATION (NFV)

SDN has emerged as a new intelligent architecture for network programmability. It moves the control plane outside the switches to enable external centralized control of data through a logical software entity called the controller. The controller offers northbound interfaces to network applications that provide higher-level abstractions to program various network-level services and applications. It also uses southbound interfaces to communicate with network devices. OpenFlow is an example of a southbound protocol. OpenFlow's behavior is simple but it can allow complex configurations. The hardware processing pipeline from legacy switches is replaced by a software pipeline based on flow tables. These flow tables are composed of simple rules to process packets, forward them to another table, and finally send them to an output queue or port.

One complementary technology to SDN called network function virtualization (NFV) has the potential to dramatically impact future networking by providing techniques to refactor the architecture of legacy networks by virtualizing as many network functions as possible. NFV advocates the virtualization of network functions as software modules running on standardized IT infrastructure (like commercial off-the-shelf servers), which can be assembled and/or chained to create services.

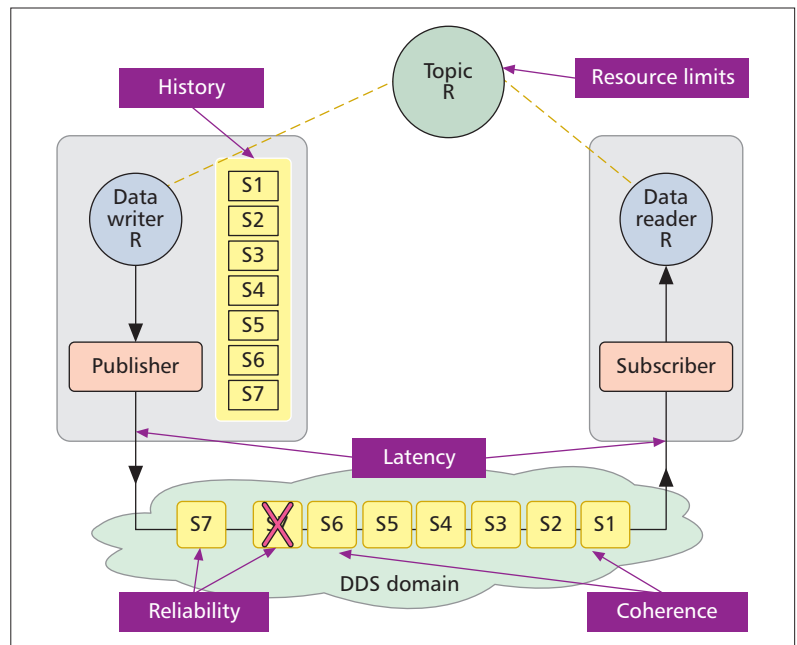


Figure 1. A view of the data distribution service.

### OMG DATA DISTRIBUTION SERVICE (DDS)

The Object Management Group has standardized the DDS [5] middleware as a protocol for the IoT to enable network interoperability between connected machines, enterprise systems, and mobile devices. DDS can be deployed in platforms ranging from low-footprint devices to the cloud, and supports efficient bandwidth usage as well as agile orchestration of system components. DDS provides a flexible and modular structure by decoupling:

- Location, via anonymous publish/subscribe.
- Redundancy, by allowing any number of readers and writers.
- Time, by providing asynchronous, time-independent data distribution.
- Message flow, by providing message-based data-centric connection management.
- Platform, by supporting a platform-independent model that can be mapped to different platform-specific models, such as C++ running on VxWorks or Java running on Real-Time Linux.

Figure 1 shows the relation between DDS entities: domains, topics, publishers, data writers, subscribers, and data readers. A DDS domain represents a virtual global data-space; information provided in the domain are accessible by the applications registered to that domain. Publishers manage one or more data writers, and subscribers manage one or more data readers. Also, DDS recognizes the importance of discovery and of meta-data, two areas addressed by IoT systems. The OMG DDS standard has evolved over time. Initially, it provided only the platform-independent and language-independent mechanisms to build distributed publish/subscribe systems with QoS capabilities. Subsequently, to promote interoperability, the standard provided a discovery mechanism via the Real-Time Publish-Subscribe (RTPS) protocol [6].

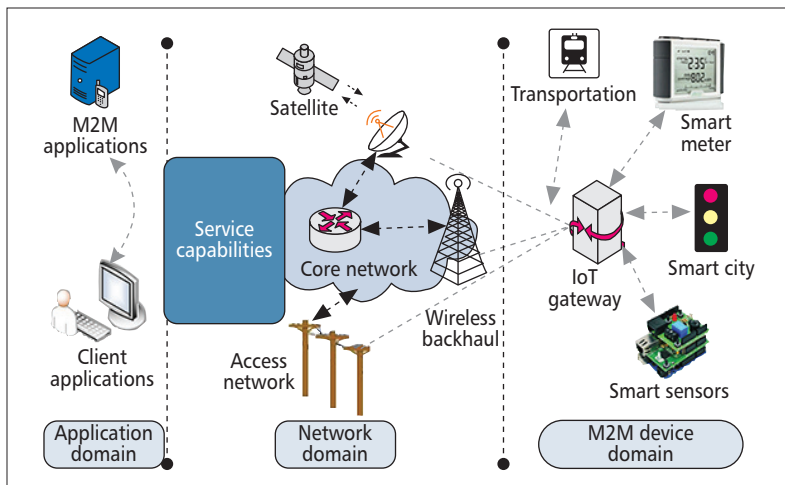


Figure 2. High-level IoT architecture.

During the discovery process each domain participant maintains a local database about all the active data writers and data readers that are in the same DDS domain. Several additional enhancements to the standard are being discussed, including the extended and dynamic types for topics, remote method invocation, integration with Web protocols, security extensions, and integration with SDN at the northbound APIs.

## NETWORK COMMUNICATION CHALLENGES IN IOT SYSTEMS

In this section we first describe the high-level architecture of IoT systems. Then we describe their most challenging networking open issues.

### IoT ARCHITECTURE

Figure 2 shows a high-level architecture of IoT systems, which is composed of three domains: the device domain, the network domain, and the application domain. In the device domain, the device provides direct connectivity to the network domain via access networks, which may include limited range PAN technologies such as Bluetooth, ZigBee, etc. or via a gateway that acts as a network proxy for the network domain.

Such a gateway must be flexible enough to efficiently manage available resources, QoS, security, and multimedia data exchange. These gateway concepts are prevalent in home ADSL models and WiFi access points found in cyber cafes and wireless hotspots. Because IoT systems integrate heterogeneous smart objects, the design of the gateway is quite different because it must not require each IoT subnetwork to have its own gateway. Thus, a convergent architecture toward a unique solution that integrates traffic incoming from heterogeneous smart devices should be designed. Furthermore, as smart objects are resource-constrained and energy-constrained, the gateway should be aware of the context of each process being managed. It should also employ intelligent routing protocols and caching techniques to route the traffic across the less constrained paths.

The network domain includes different access networks, which provide connectivity through

diverse technologies, such as xSDL, satellite, etc., to devices and/or gateways. It also provides connectivity to the core network, which includes heterogeneous and multi-technology connectivity, such as 3GPP, TISPAN, and LTE-A. Finally, the application domain includes the IoT applications and server/cloud infrastructures. The latter have to share their content, possibly back them up to other devices, analytic programs, and/or people who need to monitor real-time responses to events. They also include service capabilities, which provide functions shared between different applications through open, high-level abstractions and interfaces that hide the specificities of the underlying networks.

### NETWORK-LEVEL CHALLENGES FOR IoT

We now describe five key network-level challenges, active research areas, and standardization efforts for IoT.

*Current Standardization Efforts:* Several initiatives around the world taken by academic organizations, industries, standardization bodies, and governments have emerged during the past few years to enable IoT deployment in everyday life. Among them, the 6LoWPAN protocol was introduced to reside between the IPv6 and MAC protocol layers to make the IPv6 protocol compatible with low-capacity devices. The routing over low-power and lossy networks (ROLL) effort focuses on the routing issues for lower-energy consumption smart object networks. The Constrained Application Protocol (CoAP) is a specialized web transfer protocol for constrained smart objects and networks. Similarly, machine-to-machine (M2M) promotes the development and the maintenance of an end-to-end architecture for M2M, including sensor network integration, naming, addressing, location, QoS, security, etc. These efforts remain isolated, however.

*Mobility Management:* A large number of IoT devices are not fixed but are found in diverse mobile scenarios. Therefore, supporting and managing efficient mobility for those smart objects is of primary interest to the IoT community. Mobility management in IPv6 networks may be based on either a home agent such as in mobility IP or home location registrar/visitor location registers. The former does not require the use of a central server, which is critical in terms of scalability and flexibility. The latter is widely used in cellular networks, but is not suited for IoT communication since it requires additional and dedicated infrastructure to manage mobility. Both approaches generate a huge amount of signaling traffic to discover devices and maintain their up to date positions/locations, which degrades the network performance. Given the frequent mobility of IoT devices, it is challenging for a SDN controller to have a network view about the mobility of IoT objects to manage their spatial-temporal requests, collaborate with other controllers for adaptive handover, and dynamic flow scheduling in IoT networks

*Recurring Distributed Systems Issues:* Middleware is required in the IoT environment to provide reusable solutions to frequently encountered problems like heterogeneity, interoperability, security, and

dependability. The Message Queuing Telemetry Transport (MQTT) protocol [7] was introduced as a reliable and lightweight messaging protocol for low-bandwidth, high-latency smart devices. MQTT enables device-to-device communication through a centralized broker. However, the centralized broker presents a single point of failure. Moreover, maintaining a TCP connection between the client and the broker at all times becomes problematic for environments where packet loss is high and computation resources are scarce.

Of particular interest is the pub/sub data-centric middleware because it offers common services to help deliver events from source nodes to interested destinations in an asynchronous way. This kind of middleware integrates diverse types of communication patterns expected for IoT systems. The key challenge is how to cope with congestion control that arises in many M2M scenarios. A novel solution should be investigated to take into account the requirements of IoT devices such as short packet length and short-lived.

*Communication Protocols:* The major goals of widely used transport protocols such as TCP are to guarantee end-to-end reliability as well as congestion control. However, TCP is unsuitable for IoT scenarios since it incurs substantial overhead during the connection setup and severe performance degradation over wireless channels during congestion control. It also requires data buffering at both sources and sinks for packet retransmission and/or in-order delivery. Conversely, UDP has been used for CoAP due to its low overhead and connectionless properties. With gateways acting as repeaters and/or protocol translators, devices are expected to contribute to QoS management by optimizing the resource utilization of IoT networks. Thus, significant research efforts must be devoted to address the particular issues related to reliable transport, routing, and QoS provisioning that arise in many IoT scenarios.

*Security and Privacy:* The current landscape of IoT uses diverse wireless and/or cellular communication technologies, which makes eavesdropping and vulnerability of channels extremely simple. Also, IoT devices have limited power and computation resources, making complex security schemes infeasible. In particular, various properties, such as confidentiality, integrity, and privacy, must be ensured. Security policies and mechanisms should be specified to guarantee privacy. They should be configurable to suit individual ways to control which of their personal data is being collected, who is collecting those data, and which operations will be performed on such a data. However, in IoT there is no such infrastructure or servers to manage authentication that achieve the appropriate security policies. Thus, none of the existing solutions that have been proposed for wireless sensor networks can be applied to IoT.

## SDN-DDS ARCHITECTURE FOR IOT

This work considers the typical architecture of an IoT system described in Fig. 2, where sensors and actuators are connected to local processing and to the Internet. Internet access is provided through a core router connected to a service

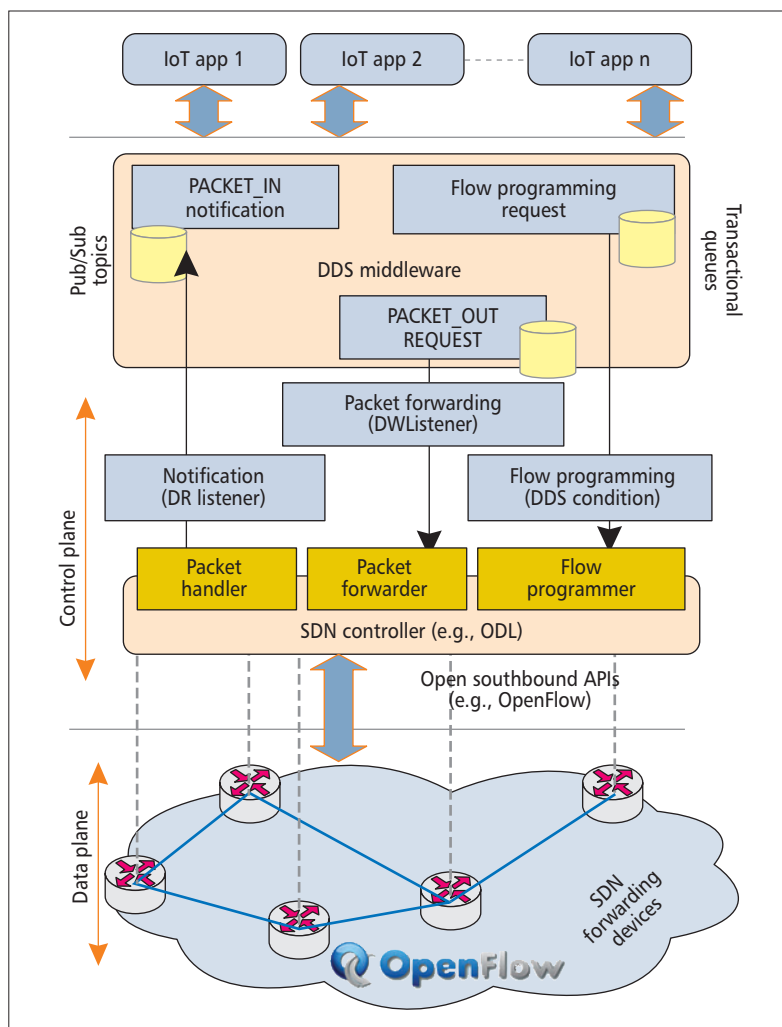


Figure 3. A view of the SDN controller with the DDS layer.

provider through terrestrial or mobile access (i.e. 3G/4G, ADSL, etc.). The router communicates with IoT gateways to enable connectivity over a wide range of communication technologies to support Wi-Fi, Ethernet, VPN, USB, Serial, Zig-Bee, etc. Both smart objects and gateways use DDS middleware to publish/subscribe data. To allow IoT systems to support SDN, we added a DDS northbound interface to the SDN controller that exposes all the necessary functionalities of IoT network applications to the controller to build network-agnostic support for IoT systems.

### THE CONTROLLER ARCHITECTURE

The controller integrates a DDS messaging layer that in turn acts as a mediator between the IoT system and the network. It supports both the proactive and reactive flow programming without the need for complex RESTful interfaces nor OSGI services. In contrast to RESTful interfaces that use synchronous client-server point-to-point communication, which may hinder scalability, DDS pub/sub messaging allows anonymous, asynchronous, and many-to-many communication semantics.

As depicted in Fig. 3, three services are implemented by the mediation layer: the packet Handler uses the DDS DataReader (DR) listen-

er to receive the PACKET\_IN events captured by the controller after being forwarded by the SDN data plane. The packet forwarder service enables forwarding packets received through PACKET\_IN events or new packets created by IoT applications. The latter are encapsulated in DDS topics, then sent to the packet forwarder service through the publisher and its data writer (DW) listener. Finally, the flow programming service is used to define the flow programming rules on the OpenFlow switches.

The mediation layer transports messages between the network control application and SDN services in the controller. That is, the DDS middleware natively fulfills the requirements of the future proactive IoT in SDN communication, which makes it an attractive technology to satisfy the features they need, such as scalability, reliability, flexibility, security and real-time data delivery.

### ADDRESSING OPEN RESEARCH CHALLENGES

The proposed architecture addresses the research questions discussed earlier as follows.

*Standardization and Open Innovation:* The proposed architecture is based on SDN and DDS, which are open and standardized, which will minimize or

even prevent the lock-in effect. As the middleware choice in our architecture, DDS can interoperate with existing standardized protocols such as CoAP, which can be easily plugged in as an extensions to the middleware without adding any complexity to the underlying network. Moreover, network virtualization enabled SDN as exemplified by OpenFlow will allow experimenting with new ideas. Using a modular design with some inexpensive general hardware platform and customized switch control software, it will be possible to improve the performance of the network.

*Addressing and Mobility:* Both DDS and SDN help enable a programmable wireless data plane to allow mobility management, dynamic channel configuration, and rapid client association. First, DDS enables many-to-many, broker-less IoT communication so that smart objects can benefit from self-addressable and self-routable data. It also enables the bounded use of resources over potentially intermittent links, and makes it possible to replay messages upon reconnection. Second, SDN techniques are being applied to wireless sensor networks to improve resource and mobility management [8].

We propose to virtualize IoT gateways to maintain active session continuity for IoT devices. Gateway virtualization helps in improving IoT multi-homing by enabling dynamic attachment of sensors to multiple networks. Figure 4 illustrates an example of handover migration of mobile smart devices from an SDN-enabled gateway (GW1) to another gateway (GW2) in a wide area network. If at any time a smart devices moves between different networks, then the centralized SDN controller modifies flow-tables in OpenFlow switches to take into account the mobility of the device.

*Efficient Middleware Integration:* Using a single middleware that supports multiple communication patterns is a very cost-effective way to develop and maintain large, distributed IoT systems. DDS provides message-oriented publish/subscribe semantics that implement various communication patterns, including transactional queues for request/response interaction (proactive flow programming), guaranteed delivery, and publish/subscribe topics for event-based interaction (reactive flow programming). Thus, we can program flows through a request-response interface implemented by message queues and react to packet-in events by subscribing to events through message topics. Moreover, DDS supports traffic differentiation and prioritization along with selective dawn-sampling.

*Service Discovery and Interoperability:* DDS discovery helps different smart devices find each other. The IoT gateway maintains a string formatted local objects list to build a complete picture about all its surrounding devices. The list includes the location and the description of each IoT object in the network, e.g. IP addresses, CoAP URI, and multicast address. However, in some cases multicast is not available in IoT networks. So if discovery frames sent with multicast addresses cannot be handled by the gateway, they will automatically be redirected to the controller.

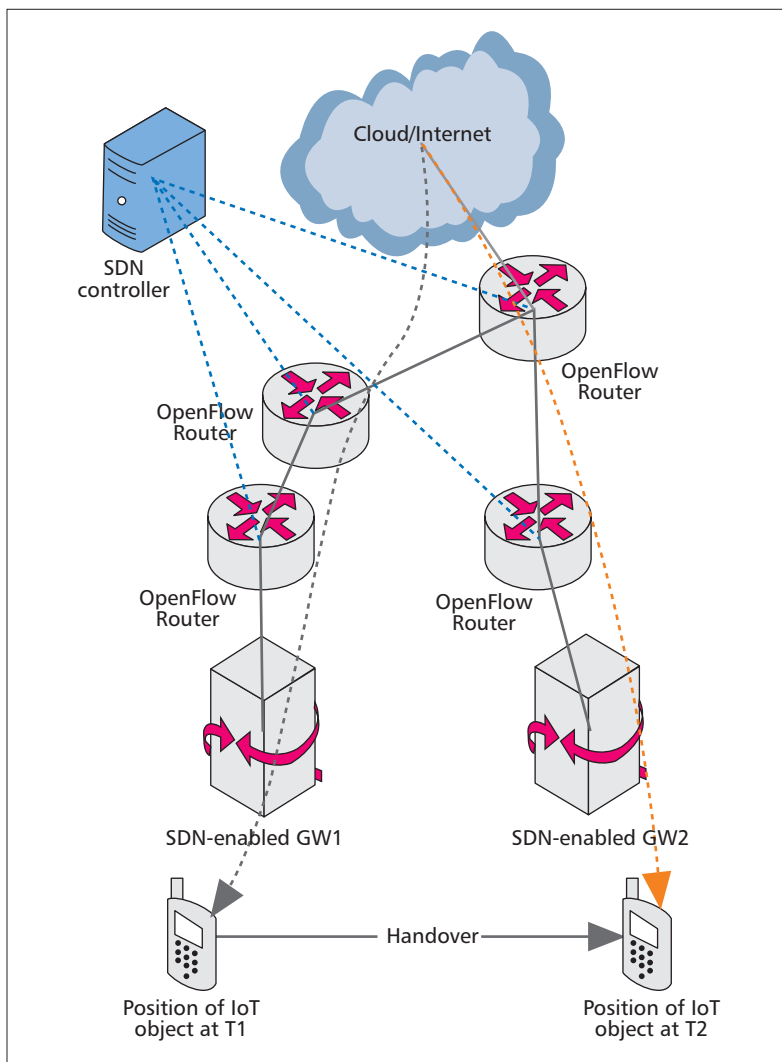


Figure 4. Use of SDN controller to improve mobility management.

The later will install a correct rule in the gateway to simulate a broadcast network so the next frame will be handled without being sent to the controller. This concept is extended to improve the interoperability of incompatible protocols. For instance, when a non DDS compliant IoT gateway is connected to the wireless AP, the gateway frames are redirected to the controller, which also redirects it to an SDN application that encapsulates data into DDS frames according to a predefined scheme.

*Scalability Handling:* A usual way to connect IoT devices consists of bridging devices with a gateway that collects data. This simple architecture, although frequently used, is not scalable when multiple gateways are used. The main problem stems from the bottleneck when several nodes have to send data at the same time. Usually, data filtering and data fusion are used to cope with this problem. Data filtering aims at dropping unnecessary packets according to a given policy. Data fusion consists of combining data to reduce the number of forwarded packets while keeping almost the same information quality.

DDS offers fusion and filtering mechanisms to reduce useless traffic. The former, called batching mode, enables collecting multiple small data samples and sending them in a single network packet to increase the bandwidth effectiveness. The latter is called content filtered topic, which makes it possible for a smart device to subscribe to a given topic and at the same time specify that it is interested in a subset of this topic data. For example, suppose a topic that contains temperature sensor samples is published with values from 0° to 100°, but the subscribing device needs only values that exceed some threshold. The content filtered topic mode can be used to limit the number of data samples the subscriber has to process and reduce the amount of data sent by the network.

As shown in Fig. 5, OpenFlow completes these mechanisms with an efficient monitoring system. The controller can collect traffic statistics at any granularity and is configured with flow entries downloaded to the switches. Byte or packet counters are associated with every OpenFlow entry. It is then easy to be informed of network flooding and react accordingly.

*Context Aware Network Management:* IoT, and especially embedded smart sensors, tend to generate huge volumes of data that must be transmitted over the network. It is not always easy to configure the transmission scheme to select exactly the relevant sensor data. DDS offers a useful mechanism called multi-channel.

With multichannel data writers, network traffic can be reduced by subdividing a data flow into a set of streams to be sent over multiple channels. Each channel maintains a set of multicast addresses, defined by an application-specified predicate called the filter expression. It is then possible to associate different multicast addresses with filter expressions, so that only the topics data that match the expression are delivered to the subscriber.

However, multichannel should be configured inside each sending device and cannot be easily

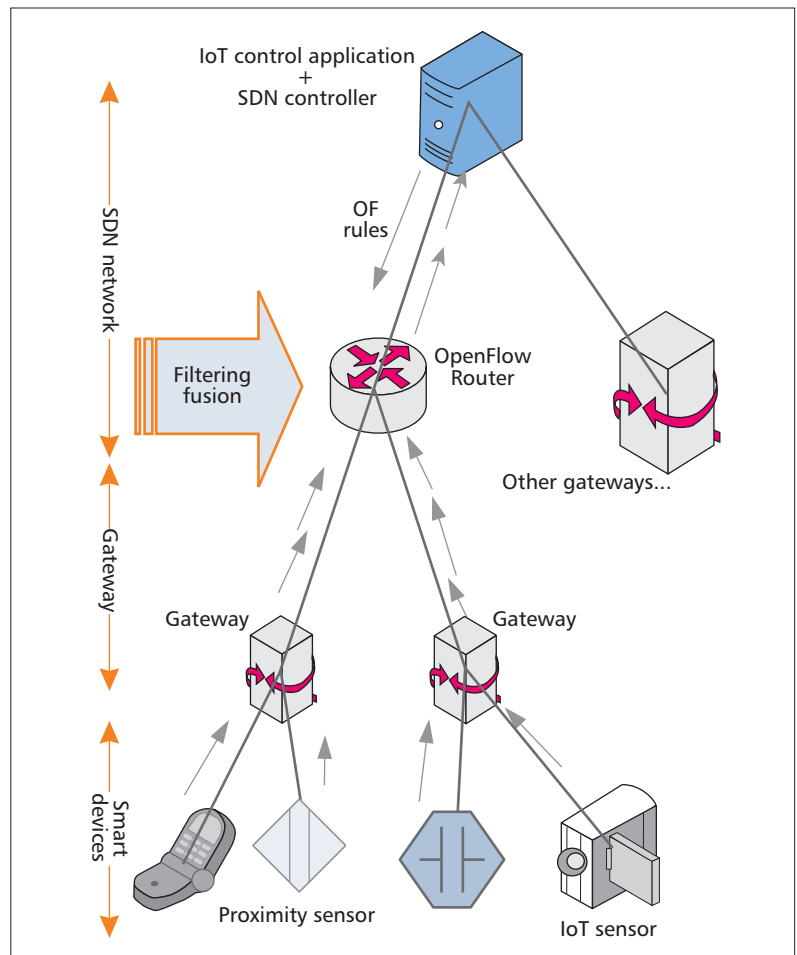


Figure 5. Scalability handling.

modified on-the-fly. OpenFlow could be used to offer a similar service. In Fig. 6, the OpenFlow match-action (Openflow rule) treatment chain could be set to filter data according a keyword contained inside the sent data as close as possible to the source. The decision can be taken at the highest application level, i.e. the application inside the controller. Owing to the OpenFlow monitoring system, the relevant information can be set.

*Solving Security Issues:* Security issues can be handled by both DDS and SDN. At the middleware level, smart devices can leverage the DDS security model (SM) [9], which offers simple and interoperable security policies without compromising the flexibility, scalability, performance, and QoS-awareness offered by DDS. Thus, we can build a fine-grained secure system that grants permissions to DDS domains, topics, or even data object instances within the topics. Also, DDS “partition” QoS provides another way to create isolated subdomains to defend IoT networks against attacks. Similarly, with SDN the IoT infrastructure can be sliced into multiple virtual partitions, so that if any attack occurs in a given partition, other IoT partitions remain isolated and secure. Along with network slicing, the controller can program the IoT gateways to conduct fine-grained packet inspection on traffic passing through IoT devices [10]. These statistics collected periodically offer a

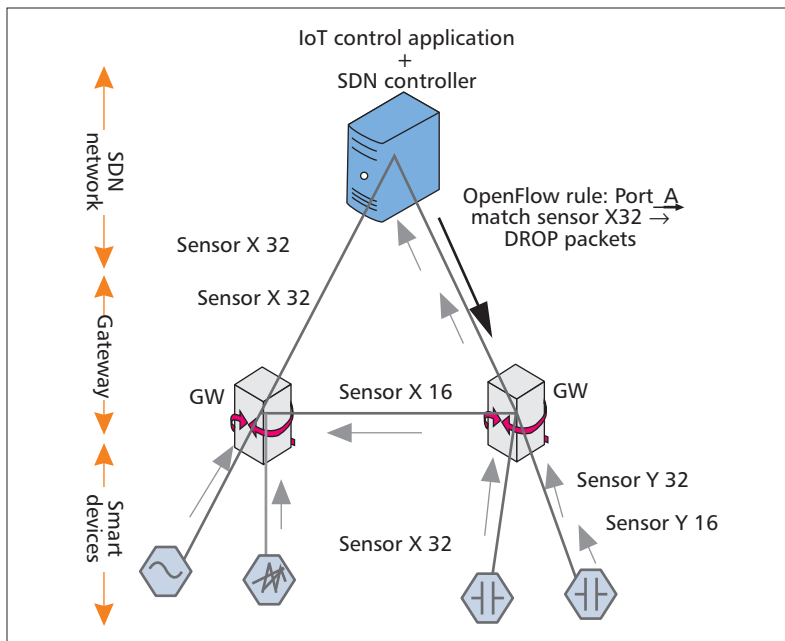


Figure 6. Multi-channel and data filtering.

real-time view of the network state. If any vulnerability occurs in the IoT system, a security SDN application can be deployed on-the-fly at the controller for detecting and driving mitigation of malware and DDoS attacks.

## CONCLUSION

The Internet of Things (IoT) promises to have a big impact by adding a new dimension in the way people will interact with surrounding things, and to form a virtual continuum of interconnected smart objects in a worldwide dynamic network. In this article we have surveyed the most important challenges that need to be tackled for efficient support of IoT systems. Then we introduced a data-centric architecture based on a symbiotic relationship between DDS and SDN that enables agile and flexible network orchestration.

As IoT covers a huge range of industries and scales of applications, IoT solutions and architectures are undergoing an evolution with a convergence of technologies, and new innovative solutions such as SDN will be considered, particularly with 5G networks on the horizon (2020 and beyond). As part of our ongoing work on DDS/SDN architectures, we believe that this study can shed light on how we can integrate DDS middleware and SDN to improve IoT communication and promote the adoption of SDN for future IoT networks.

## ACKNOWLEDGMENT

This work was partially funded by the Fulbright Visiting Scholars Program and the French National Research Agency (ANR), the French Defense Agency (DGA) under the project ANR DGA ADN (ANR-13-ASTR-0024), and the French Space Agency (CNES). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the

author(s) and do not necessarily reflect the views of the funding parties.

## REFERENCES

- [1] A. Hakiri *et al.*, "Software-Defined Networking: Challenges and Research Opportunities for Future Internet," *Computer Networks*, vol. 75, part A, Dec. 2014, pp. 453-71.
- [2] C. Benoit *et al.*, "The Web of Things Vision: Things as a Service and Interaction Patterns," *Bell Labs Technical J.*, vol. 16, no. 1, June 2011, pp. 55-61.
- [3] P. Kostelnik, M. Sarnovsky, and K. Furdik, "The Semantic Middleware for Networked Embedded Systems Applied in the Internet of Things and Services Domain," *Scalable Computing: Practice and Experience*, vol. 12, no. 3, Aug. 31, 2011.
- [4] J. Kim and J. W. Lee, "OpenIoT: An Open Service Framework for the Internet of Things," *2014 IEEE World Forum on Internet of Things (WF-IoT)*, Mar 2014, pp. 89-93.
- [5] Object Management Group, "Data-Distribution Service for Real-Time Systems," OMG, version 1.4. Sept. 2014.
- [6] Object Management Group, "The Real-Time Publish-Subscribe Wire Protocol DDS Interoperability Wire Protocol Specification," OMG, Version 2.2. Sept. 2014.
- [7] A. Stanford-Clark and H. L. Truong, "MQTT for Sensor Networks (MQTT-SN) Protocol Specification," Version 1.2. IBM Zurich Res. Lab., Zurich, Nov. 2013.
- [8] Y. Li *et al.*, "Software Defined Networking for Distributed Mobility Management," *2013 IEEE Globecom Wkspcs. (GC Wkshps.)*, 9-13 Dec. 2013, pp. 885-89.
- [9] Object Management Group. "DDS Security Specification," OMG, version 1.0 - Beta 1, June 2014.
- [10] C. Yoon *et al.*, "Enabling Security Functions with SDN: A Feasibility Study," *Computer Networks*, vol. 85, July 2015, pp. 19-35.

## BIOGRAPHIES

AKRAM HAKIRI (akram.hakiri@laas.fr) is an associate professor of electrical engineering and computer science at ISSAT Mateur, and a research scientist at SYSCOM Labs' ENIT, both in Tunisia. He is also a senior scientist and researcher at LAAS-CNRS, Toulouse, France, and a visiting research scientist at the Institute for Software Integrated Systems (ISIS) at Vanderbilt University, Nashville, TN, USA. His current research focuses on developing novel solutions to emerging challenges in network virtualization, software defined networking, QoS and resource management in wireless networks, and DRE middleware running over data networks and embedded system interconnects. He obtained his engineer's degree (computer engineering) from the National Institute of Applied Science and Technology (INSAT) in Tunisia in 2007, the M.S. (computer science) from Paul Sabatier University in Toulouse in 2008, and the D.Sc. (computer science) from Paul Sabatier University in Toulouse in 2012.

PASCAL BERTHOU (pascal.berthou@laas.fr) received a Ph.D. in computer science and telecommunication from the National Polytechnic Institute of Toulouse in 2001. He is an associate professor at the University Paul Sabatier. He joined the Laboratory for Analysis and Architecture of Systems of the French National Centre for Scientific Research (LAAS-CNRS) in 1998 as a research staff member, where he works in the area of high-speed networks and protocols and multimedia communications. Since then he has covered two major areas of activity. The first area deals with satellite communication systems. The second research area is sensor networks, particularly communication systems and their application, which in recent years has been directed toward the design of WSN for instrumentation networks. Within this domain he has focused on hardware/software network interface issues and cross layering interactions to reduce the energy consumption in WSN. He recently started research contributions in the area of SDN.

ANIRUDDHA S. GOKHALE (a.gokhale@vanderbilt.edu) is an associate professor in the Department of Electrical Engineering and Computer Science, and a senior research scientist at the Institute for Software Integrated Systems (ISIS), both at Vanderbilt University, Nashville, TN, USA. He has more than 175 technical articles to his credit focusing on topics pertaining to model-driven engineering (MDE), middleware solutions involving design patterns for quality of service (QoS) assurance, and correct-by-construction design and development of distributed real-time and embedded systems. His current research focuses on developing novel solutions to emerging challenges in mobile cloud computing, real-time stream processing, publish/subscribe systems, and cyber physical systems. He is also working on using cloud computing technologies for STEM education. Dr. Gokhale obtained his B.E. (computer engineering) from the University of Pune, India, in 1989, the M.S. (computer science) from Arizona State University in 1992, and the D.Sc (computer science) from Washington University in St. Louis in 1998. Prior to joining Vanderbilt, he was a member of technical staff at Lucent Bell Laboratories in NJ. He is a senior member of both IEEE and ACM, and a member of ASEE. His research has been funded in the past by DARPA, DoD, and NSF, including a NSF CAREER award.

SLIM ABDELLATIF (slim@laas.fr) received the M.S. (1998) and the Ph.D. (2002) in computer science, both from the University of Toulouse, France. He is currently an assistant professor at the National Institute of Applied Science (INSA) of Toulouse, and a research scientist at the French National Centre for Scientific Research (LAAS-CNRS). His current research interests include network virtualization, software defined networking, QoS, and resource management in wireless networks.

# SOFTWARE-DEFINED INTERNET OF THINGS FOR SMART URBAN SENSING

With more people living in cities, urban sensing is urgently required to create a comfortable and convenient living environment. IoT is the fundamental infrastructure to realize urban sensing, it should be flexible to support various application requirements and convenient management of infrastructure. Inspired by software-defined networking, which aims to make networks more flexible, the authors propose a software-defined IoT architecture for smart urban sensing.

Jiaqiang Liu, Yong Li, Min Chen, Wenxia Dong, and Depeng Jin

## ABSTRACT

With more people living in cities, urban sensing is urgently required to create a comfortable and convenient living environment. As Internet of Things (IoT) is the fundamental infrastructure to realize urban sensing, it should be flexible to support various application requirements and convenient management of infrastructure. Inspired by software-defined networking, which aims to make networks more flexible, the authors propose a software-defined IoT architecture for smart urban sensing. This architecture decouples urban sensing applications from the physical infrastructure. Centralized controllers are designed to manage physical devices and provide APIs of data acquisition, transmission, and processing services to develop urban sensing applications. With these properties, various applications can coexist on the shared infrastructure, and each application can request controllers to customize its data acquisition, transmission, and processing on-demand by generating specific configurations of physical devices. This article discusses the background, benefits, and design details of the proposed architecture as well as open problems and potential solutions to realize it, which opens a new research direction for IoT and urban sensing.

## COMMUNICATIONS STANDARDS

Urban sensing is one of the most promising solutions to address the above problems. As a simple example, if the real-time traffic load is known, the efficiency of existing transportation systems can be enhanced significantly. The Internet of Things (IoT) has the potential to make urban sensing a reality. With an increasing number of various sensor devices connected to the Internet, it is possible to obtain the infrastructural and environmental data in real time that would enable an efficient approach to perceive and manage urban facilities. Many cities have deployed sensor platforms to support urban sensing. For example, London has deployed various sensor nodes to obtain traffic, environmental, and utilities data [2], and various experimental platforms of IoT have been developed for research [3]. In addition to dedicated sensor platforms, as human carried smart phones are equipped with a rich set of sensors like cameras, digital compasses, GPS, etc., they can also be exploited to realize urban sensing. This is referred to as mobile crowd sensing [4].

Currently, IoT is still in the initial stages of development and deployment. However, there is no doubt that the IoT will have an important impact on people's lives, just like Internet does today. The Internet has had great success and changed our lives, but it still faces some problems. On one hand, as the control intelligence, which is implemented by various routing and management protocols, is embedded in every router/switch and is hard to change, Internet infrastructure becomes ossified and therefore evolves slowly. Also, vendor-dependent interfaces make the infrastructure management complex and error-prone. On the other hand, it only provides best-effort service and thus prevents the development of highly personalized applications with specific requirements on service quality and user experience. The design of the future IoT architecture should avoid these problems to support sustainable evolution, convenient management, and various application requirements.

Software-defined networks (SDNs) [5] offer the ability to address the above mentioned problems. In SDN, the control intelligence is moved from data plane devices (switches, routers) and implemented in a logically centralized controller, which interacts with data plane devices through standard interfaces. The network operator runs software programs on the controller to automatically manage data plane devices and optimize network resource usage. They can further develop up-to-date control schemes to provide different network services for applications, e.g. providing QoS guaranteed forwarding services.

Inspired by SDN, this article proposes a software-defined IoT (SD-IoT) architecture for smart urban sensing. In accordance with SDN, SD-IoT also decouples the control logic from functions of the physical devices through a logically centralized controller that manages the devices via standard interface. In particular, SD-IoT extends the spirit of the software-defined approach from network devices to sensor platforms and the cloud, and combines them to support urban sensing applications together by

## INTRODUCTION

The number of people living in cities has increased dramatically in recent years, and the trend is expected to continue. The United Nations Population Fund estimates that by the year 2030 nearly 60 percent of the world's population will live in an urban environment [1]. More convenient and comfortable living condition, as well as more opportunities for work and career development, are the main motivations for urbanization. However, the explosion of city populations resulting from urbanization is straining existing daily and public facilities such as transportation, healthcare, and security, and creating new problems like environmental pollution. These issues need to be solved to provide sustainable urbanization.

Jiaqiang Liu, Yong Li, and Depeng Jin are with Tsinghua University.

Min Chen is with Huazhong University of Science and Technology.

Wenxia Dong is with Huawei Technologies Co. Ltd.

This work is supported by the National Basic Research Program of China (973 Program) (No. 2013CB329001), and the National Nature Science Foundation of China (No. 61301080, No. 91338203, No. 91338102, and No. 61321061).

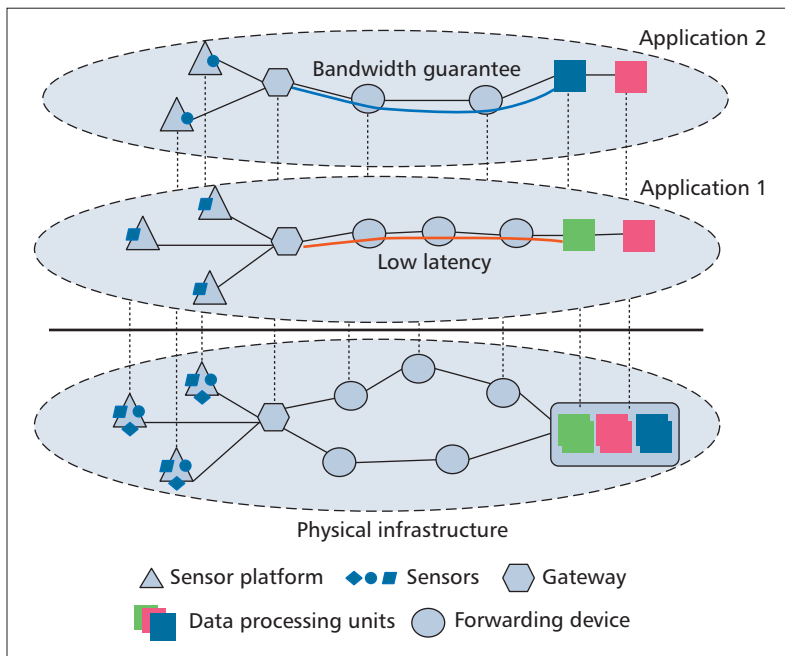


Figure 1. Illustration of software-defined IoT solution.

providing well-defined service APIs in terms of data acquisition, transmission, and processing. Figure 1 conceptually illustrates the usage of this architecture. The physical infrastructure consists of sensor platforms, forwarding devices, and servers. On top of this infrastructure, multiple urban sensing applications are deployed, and each application customizes its data acquisition, transmission, and processing through the service APIs. The standard service API reduces the complexity and developing cycle for deploying a new application, while the sharing of physical infrastructure greatly reduces the capital and OAM (operation, administration and maintenance) costs. These characteristics empower SD-IoT to efficiently support various application requirements and thus enable smart urban sensing.

In the rest of this paper, we first analyze current urban sensing applications along with the problems and trends. We then introduce the proposed software-defined IoT architecture, followed by open problems and potential solutions. After that, we present a quantitative analysis to show the benefits of SD-IoT. We then conclude the paper.

## IoT: STATE OF THE ART AND TRENDS AND SOLUTIONS

### URBAN SENSING: STATE-OF-THE-ART

Figure 2 presents the three most common urban sensing applications: temperature, noise, and PM 2.5 monitoring. These applications can be logically divided into three subsystems: data acquisition, transmission, and processing. Specifically, different sensors are deployed in the city to obtain temperature, noise, and PM 2.5 data. The obtained data are then transmitted to remote servers to be stored/processed. Usually, the sensor node first transmits the data to a gateway

through a wireless sensor network (WSN) [6]. The gateway then transmits data to the remote server through wireless or wired networks. Data processing may happen during the whole course, e.g. filtering undesired data at the sensor nodes [7], compressing and encrypting data at the gateway, further analyzing acquired data at the server to obtain the statistical information, etc.

Currently, an application-oriented approach is utilized to develop these three subsystems [11]. That is, application developers customize the sensor platform, gateway, network, and remote server from slate state, according to the application requirements. Specifically, the developers need to buy or develop a sensor platform according to the application requirements, which usually includes the sensors to obtain the required data, the radio modules to transmit the data, the power supply modules, and the microcontroller to coordinate the peripheral modules and execute data processing functions. The firmware also needs to be customized for this specific application. As an example, Downes *et al.* [8] introduced the design of a platform for wireless image sensor networks. Then the developers need to take a similar path to customize the network and computing infrastructure, e.g. to determine how to access the Internet, whether to cache the data or not, where to store and process the obtain data, etc.

### PROBLEMS ANALYSIS

While the above application-oriented approach seems quite direct, it has many drawbacks. We summarize them as follows:

**High Capital and Maintenance Cost:** As each application needs to deploy and manage its own sensor platforms, it requires a huge investment in hardware deployment and maintenance, while in fact it is possible for many applications to share sensor platforms if they require the same type of data. Even when the required data are different, many modules in a sensor platform, like the radio, power supply, and microcontroller, can be shared to reduce the overall cost.

**Inflexible for Potential Application Changes:** Under this approach, the infrastructure and the application are closely coupled, i.e. the intelligence of the application is hard-wired in the sensor platform, the gateway, and the server. Any change related to an application requires re-developing or re-customizing the physical infrastructure, which is complex, error-prone, and sometimes even impractical.

**Inefficient Resource Usage:** As the control logic of applications is embedded in hardware devices, it is difficult to improve resource utilization by dynamically optimizing data acquisition, transmission, and processing. For example, as there is no approach to dynamically control data collection and transmission in sensor platforms, they would continuously transmit the data to remote servers, even though such data is undesired during some time periods, and thus the energy of sensor platforms and the bandwidth of the network are wasted.

**Long Development and Deployment Cycle:** As each application needs to develop and deploy its own sensor platform, gateway, and remote server from scratch, the overall time to introduce a new application is long. The long development and



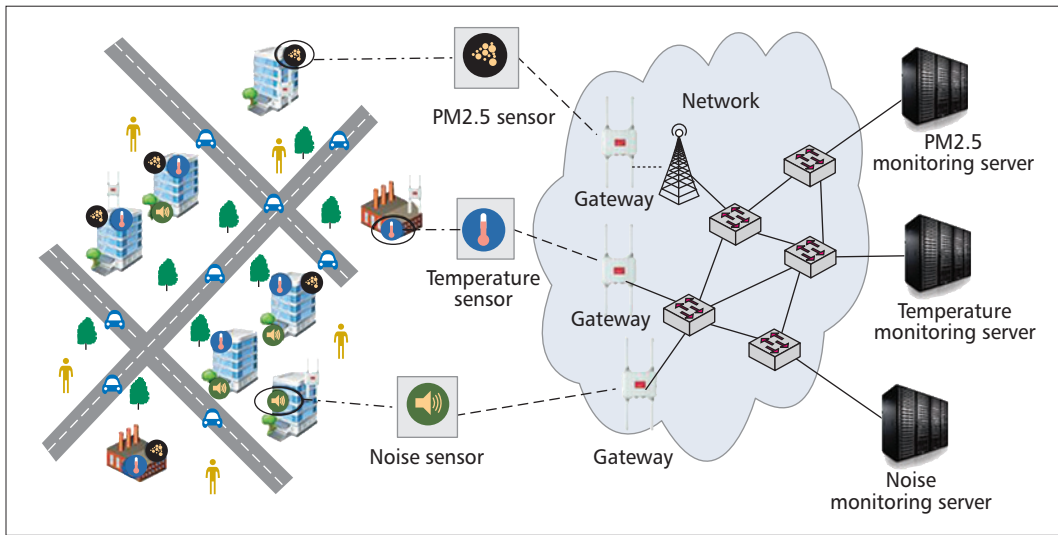


Figure 2. Illustration of urban sensing applications, which include three subsystems of data acquisition, transmission and processing.

With the increase in users and applications, Internet architecture and related infrastructure are also evolving in order to deal with encountered challenges. There are some apparent trends in this evolution, which should be considered in the design of IoT architecture.

deployment cycle, as well as the high investment required, definitely increase the barrier for deploying new applications and thus prevent potential innovations in applications.

### THE TRENDS

With the increase in users and applications, Internet architecture and related infrastructure are also evolving in order to deal with encountered challenges. There are some apparent trends in this evolution, which should be considered in the design of IoT architecture.

**The Sharing of Physical Infrastructure:** Sharing means that the underlying physical infrastructure simultaneously supports multiple applications belonging to multiple parties. The popularity of cloud computing best explains the trend of sharing physical infrastructure. Through cloud computing, the application developers deploy their applications in cloud data centers rather than build their own physical infrastructures. In addition to cloud computing, there is also a trend of sharing network infrastructures, such as base stations, access points, etc. As the sharing of physical infrastructure has one general benefit of reducing capital and maintenance costs, we envision that the IoT infrastructure should be shared to obtain this benefit.

**The Rising of Software-Defined Architecture:** SDN enables flexible network control by separating the control plane and the data plane. Inspired by these benefits, SDN has been extended to mobile access networks [9] and wireless sensor networks [10, 11]. As the physical infrastructure will become increasingly complex in the era of IoT, it is necessary to borrow the insight of SDN to realize flexible control and management of IoT infrastructure.

**The Prevalence of Application Programming Interfaces:** Providing application programming interfaces (API) is a growing trend to share physical infrastructure. Cloud providers like Google APP Engine have offered such APIs, and network controllers like OpenDaylight [12] also provide northbound APIs to develop control applications.

In addition to enabling the sharing of physical infrastructure, APIs also hide the complexity and heterogeneity of the physical infrastructure, which significantly reduces the difficulty of application development and shortens the time to market of new applications. This trend suggests that IoT, especially sensor platforms, should provide APIs for applications to exploit their abilities in a flexible and efficient manner.

## SD-IoT: ARCHITECTURE OVERVIEW AND SYSTEM DESIGN

### ARCHITECTURE OVERVIEW

In this paper we propose a software defined IoT (SD-IoT) architecture. As illustrated in Fig. 3, SD-IoT consists of three layers: a physical infrastructure layer, a control layer, and an application layer.

**Physical Infrastructure Layer:** This layer is composed of various kinds of physical devices, including sensor platforms, gateways, base stations, switches/routers, and servers. These devices possess the essential functions and resources to sense an urban environment, transmit data from one node to another, and process them to extract required information. However, they do not determine what to do by themselves. Instead, they leave the decision-making to the control layer by interacting with it through standard interfaces, i.e. a southbound interface named in SDN.

**Control Layer:** The control layer acts as the intermediary between the infrastructure layer and the application layer. On one hand, the control layer manages the physical devices with various characteristics and functions through different southbound interfaces. On the other hand, the control layer provides services to the application layer through APIs known as northbound interfaces. For urban sensing applications, the control layer will provide data acquisition, transmission, and processing service. We will explain these services in detail in the following subsections.

	State of the art	Problems	The trends
<b>Data acquisition</b>	Application-oriented wireless sensor platforms. The control functions are preset in the firmware.	Hard to customize in run time. Hard to implement dynamic optimization. High capital and maintenance cost.	Over the air programming to update sensor firmware.
<b>Data transmission</b>	Distributed protocols, such as WiFi, ZigBee, TCP/IP. The control protocols embed in each forwarding device.	Hard to control and evolve. No QoS guarantee.	Software-defined network. Network as a service with QoS guarantee.
<b>Data processing</b>	Each application developing data processing pipelines from the scratch.	The time cycle to develop a new application is long. Hard to share data processing resources.	Cloud based data processing to provide various data processing software, platform, and tool.

Table 1. Summary of state of the art, problems and the trends.

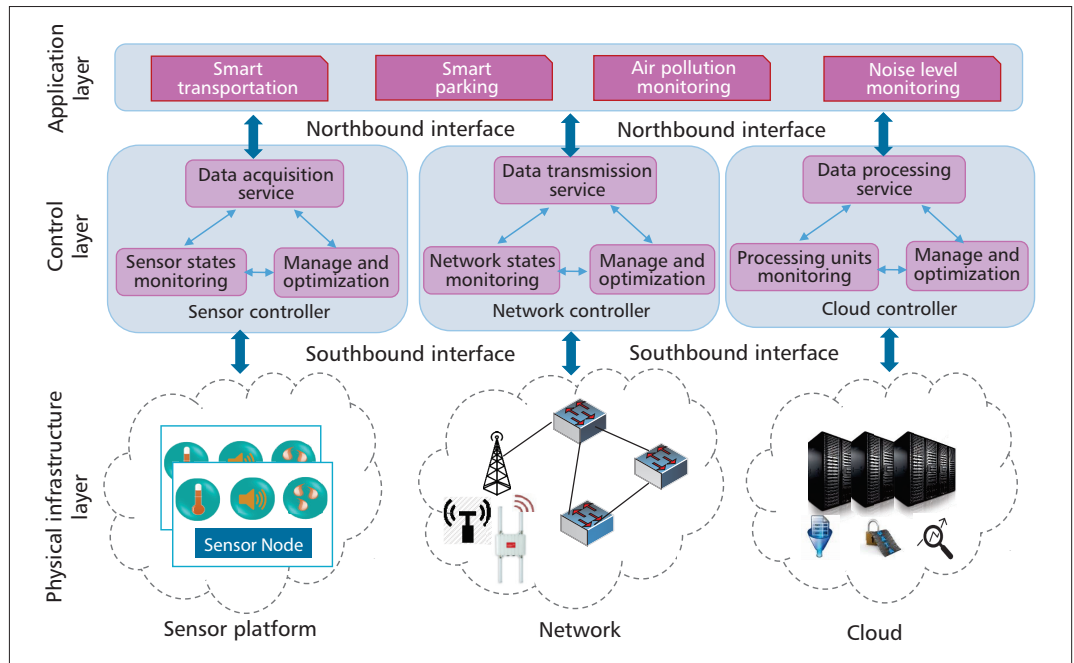


Figure 3. Architecture of software-defined IoT.

**Application Layer:** In this layer, developers build urban sensing applications using the provided APIs. In particular, they can customize data acquisition, transmission, and processing without worrying about the required change of configurations in physical devices, which greatly simplifies the procedure of developing new applications. Also, as the physical infrastructure is shared by multiple applications, the overall capital and maintenance costs are reduced.

#### SENSOR PLATFORM AND DATA ACQUISITION SERVICE

The data acquisition service provides APIs for applications to specify their data requirements. The controller automatically configures sensor platforms to obtain the required data. Data specification includes general attributes, such as data type, targeted geographical areas, and time duration. For example, as shown in Fig. 4a, an application would request PM 2.5 data at Tsinghua University. Applications can also specify type-dependent attributes, e.g. sampling rate can be set for PM 2.5 and noise data. The data acquisition service also provides APIs for applications to query the properties of available data such as data types and geographical

areas as well as optional attributes for each type of data.

Under SD-IoT, each sensor platform is equipped with more than one sensor with the same or different types and shared by many applications. For example, a sensor platform may include a PM 2.5 sensor and a noise sensor simultaneously, significantly reducing the total number of sensor platforms that need to be deployed. As a result, the overall investment for hardware, deployment, and maintenance is also reduced. The sensor controller has a global view of the underlying sensor platforms. Specifically, it knows the location and embedded sensors of every sensor platform. Based on the global view, the sensor controller can dynamically activate/deactivate sensors and customize their configurations to satisfy application requirements and simultaneously reduce energy consumption.

#### NETWORK AND DATA TRANSMISSION SERVICE

The network is used to transmit data from the sensor platforms to servers in the cloud. As applications may prefer different cloud data centers, they should have the ability to specify the destination of data transmission. Also, applica-

tions may have specific performance requirements for data transmission. For example, a smart transportation application that provides path planning suggestions must be aware of current traffic load, and thus requires low latency data transmission. In contrast, a video application that provides real time street views must guarantee that the video is fluently transmitted, and thus requires reservation of bandwidth.

The data transmission service provides APIs for applications to specify their requirements, which mainly include two dimensions: destination and QoS parameters. An IP address can be used to specify destination, while several options can be provided for QoS specification: basic transmission, latency sensitive transmission, and bandwidth guaranteed transmission.

Generally, basic transmission is carried in a best-effort manner, latency sensitive transmission has high priority during traffic scheduling, and the controller reserves bandwidth for bandwidth guaranteed transmission. In addition, with the advance of network function virtualization (NFV), the network will also provide on-path data processing services, e.g. data compressing and encryption. Specifically, the data transmission service API will also allow applications to specify the service chain [13], i.e. the pipelines of virtual network functions that a specific flow needs to go through. Figure 4b illustrates two examples of requests for data transmission service.

To realize the data transmission service, the network follows a software-defined network architecture. The forwarding devices are programmable, e.g. OpenFlow-enabled, and the controller is responsible for implementing traffic steering and scheduling. Specifically, based on the collected global network view, the controller steers packets to different destinations, and dynamically schedules traffic to satisfy application requirements for network quality and optimize the usage of network resource.

### CLOUD DATA CENTER AND DATA PROCESSING SERVICE

Urban sensing data is further stored and processed using resources provided by three main cloud computing models: IaaS, SaaS, and PaaS. Currently, a cloud usually uses one of these for service provisioning. However, as the application of urban sensing would require them simultaneously, we argue that they should be integrated to support data processing services, i.e. a cloud should simultaneously provide software service, platform service, and infrastructure service, and offer APIs for users to flexibly utilize them together. Figure 4c illustrates two examples. One application requires mining the received data in real time and exploiting visualization software to illustrate the mining results. Thus, it requires a data mining platform and visualization software, which can be provided by PaaS and SaaS, respectively. Another application aims to store the received data at first and then exploit their own programs for data processing. Thus, it needs storage and VMs, which can be provided by IaaS.

Data processing service APIs allow applications to specify the required resources, which includes running submitted programs on specific platforms, deploying existing software entities, and providing VMs. The cloud controller knows

the state of the underlying server resource pools, such as which servers are used to support a specific platform and the residual resource in them, and maps the application's resource request to underlying server pools based on it.

## OPEN PROBLEMS AND POTENTIAL SOLUTIONS

### SOUTHBOUND INTERFACE DESIGN

To implement SD-IoT, the southbound interface should be designed for controllers to interact with the physical infrastructure. Some interfaces have been designed to address the challenge. For forwarding devices, OpenFlow is the most widely used interface that abstracts the forwarding behavior in heterogeneous switches and routers. For servers, the interface usually depends on the cloud control system. However, in general, middleware software applications are utilized to deal with device heterogeneity. Compared to forwarding devices and servers, designing a southbound interface for sensor platforms is much more difficult due to higher device heterogeneity. Besides, sensor platforms are energy limited and thus the energy consumption for interacting with the controller should be reduced as much as possible.

As an initial idea for designing a southbound interface for sensor platforms, we propose to combine the strategy of abstraction and middleware software. First, by providing an abstraction on the data collection, processing, and transmission procedure in the sensor platform, the implementation of the sensor controller is decoupled from sensor platforms. Galluccio *et al.* [11] has proposed a finite state machine based abstraction for data processing and transmission. In Fig. 4a we also show an example to abstract the data collection ability based on the included sensors and their types and IDs. Second, before the standardization of the abstraction, the actual control interface of different sensor platforms varies across each other. Middleware software then can be exploited to carry out the transformation. Particularly, to save energy at the sensor platform, middleware software can be placed in the controller, and the controller should decrease the frequency of interactions with sensor platforms when they are inactive. Despite having these benefits, the design and implementation of the abstraction and middleware software needs more discussion and study.

### CONTROL LAYER DESIGN

The design of a logically centralized control layer should achieve three objectives: high scalability, high performance, and high robustness. First, as more and more physical devices and applications will be added over time, the control layer should scale at the same time to support them. Besides, in SD-IoT, application performance and control flexibility depends on the performance of the interaction between the control layer and the physical layer, e.g. communication delay. Further, the control layer must be robust enough to work normally under various possible failures.

Deploying multiple controllers is a general approach to achieve these objectives. On one hand, the controller can be replicated to increase

The cloud controller knows the state of the underlying server resource pools, such as which servers are used to support a specific platform and the residual resource in them, and maps the application's resource request to underlying server pools based on it.

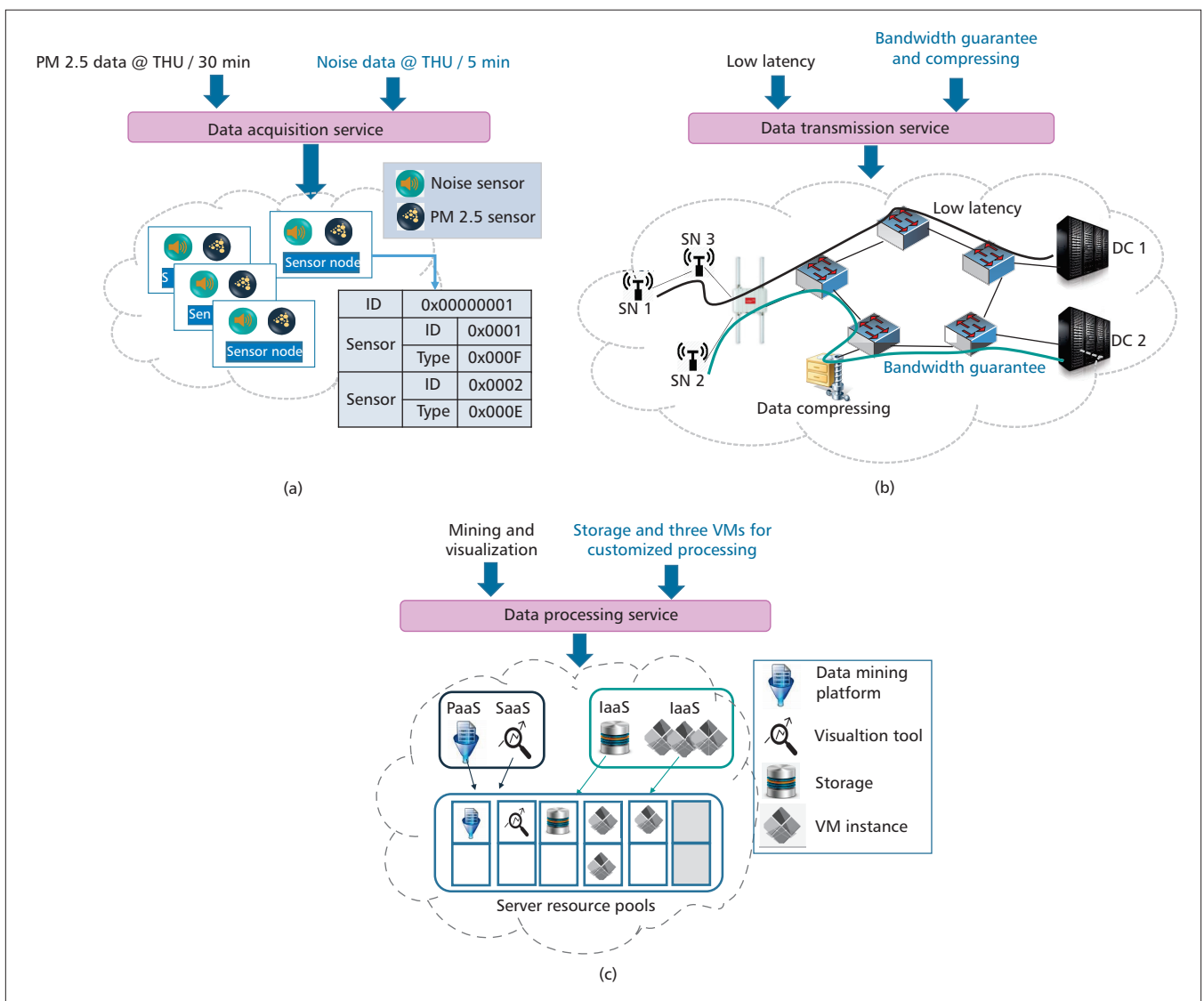


Figure 4. Illustration of data acquisition, transmission and processing services: a) data acquisition service; b) data transmission service; and c) data processing service.

its robustness. On the other hand, each controller can manage only part of the devices and thus the control layer can be scaled by increasing the number of controllers. Further, the controllers can be placed at different locations to reduce the average communication delay to the physical devices. Ahmed and Boutaba [14] proposed to exploit a vertical approach to organize multiple controllers for software defined wide area networks, which can also be extended to the control of sensor platforms and clouds. For example, we can use one controller to control one WSN and use an orchestration controller to coordinate them.

#### MOBILITY MANAGEMENT

In SD-IoT, the control plane needs to implement mobility management to support the mobile sensor platform handover from one gateway to another. Particularly, as the gateways are controlled by multiple physical controllers for scalability, the control plane needs to coordinate these controllers to implement the mobility management function.

Wu *et al.* [15] have introduced a solution to mobility management by maintaining a controller network based on structured overlay when the physical controllers are distributed. In addition, we can also employ an orchestration controller to coordinate the mobility management by recording the controller currently managing each mobile sensor platform. When a mobile sensor platform attaches to the gateway managed by a new controller, the controller reports the event to the orchestration controller, which then coordinates the original controller and the new controller to carry out the handover.

#### CONFLICT RESOLUTION AND OPTIMIZATION FOR THE SENSOR PLATFORM

Under SD-IoT, the sensor platforms are shared by different applications, which may lead to potential configuration conflicts. For example, one application may request noise data in Tsinghua with a sampling rate of once every five minutes, while another application may also

request the noise data in Tsinghua but with a different sampling rate. When conflicts happen, the sensor controller needs to decide whether to accept the application's request, how to resolve the conflicts, and simultaneously minimize the sensor platform's energy consumption.

To avoid the conflicts, the controller can allocate each sensor to at most one application. While this strategy is simple, it is inefficient because sometimes the sensor can still be shared by applications with different settings. For example, if Application A wants to set the sampling rate to once every five minutes, while Application B wants to set it to once every seven minutes, the controller can make Application A and Application B share the sensor by generating a series of sampling time points, e.g. five minutes, seven minutes, 10 minutes, 14 minutes, 15 minutes, etc. To extend this approach to a general case, the controller sets up a resolver for each sensor. The resolver records the sensor's current configuration. After receiving a service request and transforming it to the configuration of the sensor, the controller sends the configuration to the corresponding resolver. The resolver generates appropriate configurations according to the new configuration and current configuration, or it will reject the configuration request if there are unresolvable conflicts.

### QoS ENABLED TRAFFIC SCHEDULING

In SD-IoT, the network controller can provide end-to-end QoS guaranteed data transmission. However, several challenges need to be addressed to achieve this. First, the forwarding devices have a limited number of queues for QoS enforcement, and thus it is difficult for them to support a huge number of QoS requirements. Second, it is a challenge to design efficient traffic scheduling algorithms to satisfy QoS requirements in a large-scale network.

We now explain two strategies to deal with the above challenges. The first is reducing the demand for queues in forwarding devices by quantization of QoS requirements, which can be conducted based on statistics of application requirements. The second is considering the number of available queues in each forwarding device when scheduling traffic. As this strategy further increases the complexity of the traffic scheduling problem, some approximate algorithms should be developed to solve it efficiently.

### RESOURCE MAPPING IN CLOUD DATA CENTERS

In SD-IoT, the cloud controller needs to decide how to map application service requests to the physical devices. For example, considering that the cloud allows applications to store data at first and then rent VMs to process them, the cloud controller needs to decide where to store the data and which servers should be used to host VMs for post-data processing. Generally, several objectives are expected, such as increasing the cloud provider's revenue by accepting more service requests, saving energy by using less servers, and balancing the server's load by equally mapping service requests to different servers. The constraints include server capacity, storage capacity, and the type of softwares/platforms/VMs a server can host.

One challenge to implement the above optimizations is that sometimes different objectives are in conflict with each other, and hence it is impossible to achieve them simultaneously. As an example, energy saving and load balancing are conflicting objectives. Efficient heuristic or approximate algorithms should be developed to achieve a trade-off between conflicting objectives. As an example, a threshold-based strategy to activate and shut down servers can be an effective trade-off between load balancing and energy saving, i.e. activate more servers when the average workload exceeds a pre-set maximum, and shut down some servers when the workload goes below a pre-set minimum.

## CASE STUDY AND QUANTITATIVE ANALYSIS

### SELECTED SCENARIO

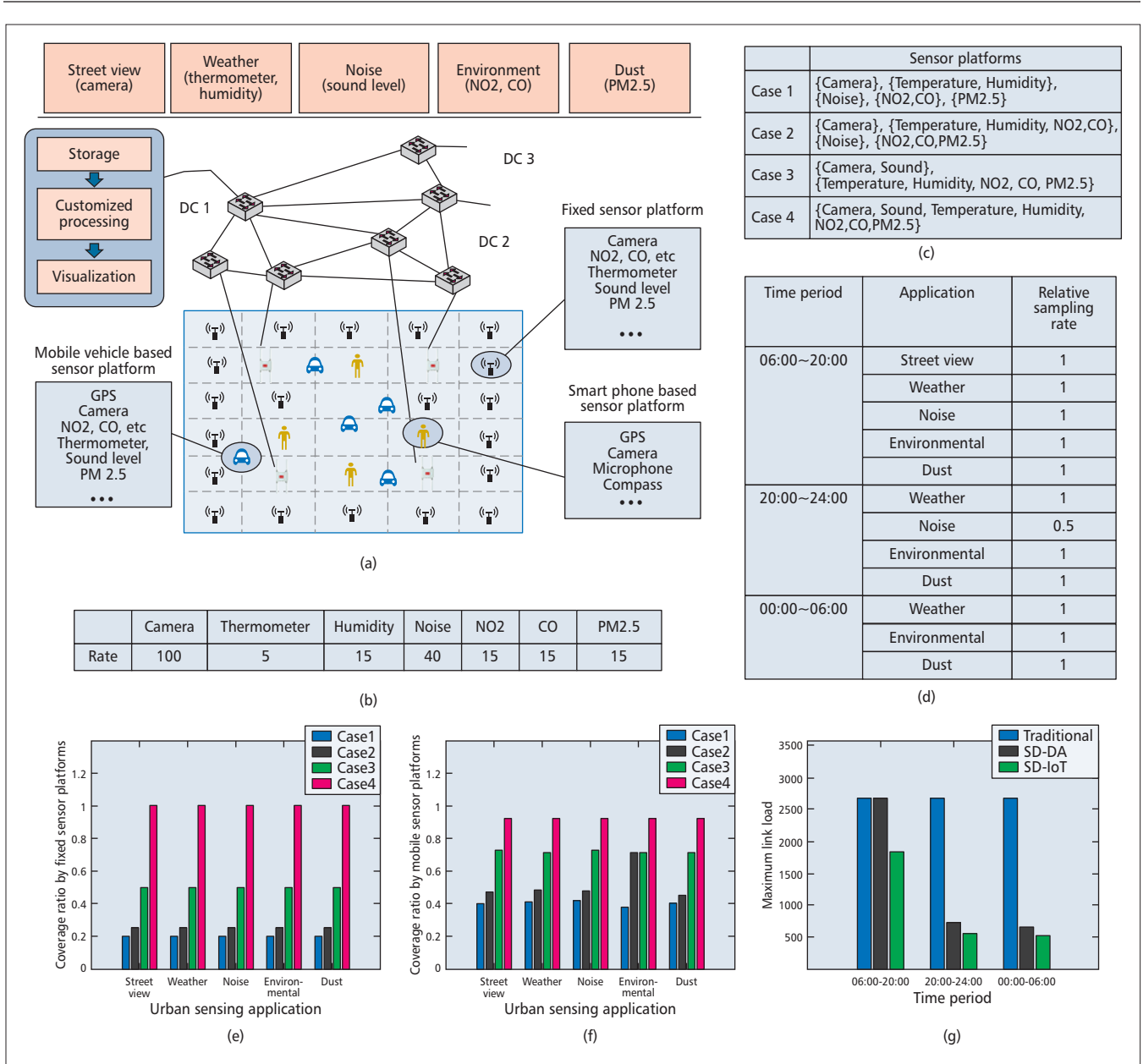
In order to further illustrate the benefits of SD-IoT, we conducted a case study and quantitative analysis, described in this section. Figure 5a illustrates the selected scenario. The target region consists of 5x6 rectangular urban areas. There are three types of sensor platforms: fixed sensor platform, user smart phone based sensor platform, and mobile vehicle based sensor platform. We considered five urban sensing applications: street view, weather monitoring, noise monitoring, environmental monitoring, and dust monitoring. The essential sensors to support each application are shown in the figure. We assume that each rectangle and vehicle can deploy at most one sensor platform. During a time period, each vehicle has a constant probability to appear in one specific rectangular area. In our simulation there are 2500 vehicles, and the constant probability is set to 1/1000. Further, there are three data centers connected by a network of seven forwarding devices. The data rate of each sensor under a standard sampling rate is shown in Figure 5b. The five applications employ a similar data processing procedure. Data is stored at first; then, VMs are used to execute customized processing; finally, the visualization software applications are exploited to illustrate the processing results.

### QUANTITATIVE ANALYSIS

*Data Acquisition:* To illustrate the benefits of SD-IoT, we consider four deployment cases, which are illustrated in Figure 5c. The first case reflects current IoT architecture, where there are five sensor platforms, each embedding sensors required by corresponding applications. The other three cases reflect the proposed SD-IoT architecture, where some sensor platforms embed more sensors and are shared by different applications. For example, in case 3 the sensor platform equipped with camera and sound level sensors can be shared by street view and noise monitoring applications. In each deployment case, every rectangular area and participating vehicle randomly selects one sensor platform to deploy.

We then observe the average coverage ratio of each application, which is defined as the ratio of areas covered by sensors of that application over the total target area. We show the results in

When a mobile sensor platform attaches to the gateway managed by a new controller, the controller reports the event to the orchestration controller, which then coordinates the original controller and the new controller to carry out the handover.



**Figure 5.** Quantitative analysis scenario and results: a) selected scenario; b) data rate of each sensor under standard sampling rate; c) sensor platforms exploited in four deployment cases; d) the relative sampling rate (compared to standard sampling rate) of each application during different time periods; e) the coverage ratio by fixed sensor platforms with four deployment cases; f) the coverage ratio by vehicle based sensor platforms with four deployment cases; and g) the maximum link load during different time periods under three different scenarios.

Fig. 5e and Fig. 5f. From the figure, we can find that the coverage ratio of each application increases from case 1 to case 4. For example, in Fig. 5f, when sensor platforms in each vehicle can be used by one application, as case 1 indicates, the coverage ratio is only about 40 percent. However, if each sensor platform can be shared by five applications, as case 4 indicates, the coverage ratio increases to 90 percent. Such enhancement is due to the increase in the number of sensor platforms that can be used by each application. As the cost to deploy a sensor platform with seven sensors is not much higher than the cost with fewer sensors, these results suggest that it is efficient for different applications to share the underlying sensor platforms.

*Data Transmission:* We assume there are four gateways deployed in the target area. The sensor platform first transmits the data to the nearest gateway, then the gateway transmits the data to the corresponding data center. In Fig. 5d we show the sampling rate of each application at different times of day. We consider three scenarios. The first scenario corresponds to the current IoT architecture, and we use “traditional” to refer to it, where the configuration of sensor platforms is unchangeable and the network uses the shortest path to transmit data from the gateway to the data center. The second scenario is denoted by “SD-DA,” where the configuration of sensor platforms can be dynamically changed, but the network still uses the shortest path for

data transmission. The third scenario, referred as “SD-IoT,” corresponds to the proposed software-defined IoT architecture, where both the sensor platform and network are software defined and thus can be exploited to dynamically optimize data transmission.

Figure 5g shows the maximum link load under the above three scenarios. For the traditional architecture, the maximum link load remains the same over time because the configuration of each sensor platform is fixed and thus the total data rate does not change. In contrast, the SD-DA scenario reduces the maximum link load at night, since it can switch off sensors when not required or lower their sampling rate to reduce the total data rate.

Moreover, as software defined networks enable dynamic and global optimization of traffic forwarding, SD-IoT further decreases the maximum link load by distributing the traffic equally over multi paths. Specifically, the maximum link load is reduced by 32 percent, 25 percent, and 22.6 percent during the time periods 6:00-20:00, 20:00-24:00, and 00:00-06:00, respectively.

*Data Processing:* In the selected scenario, each application requires data storage, the platforms, or VMs to execute customized data processing, and visualization software to illustrate the data processing results. Currently, these resources can only be separately provided by different cloud data centers. Therefore, the data processing procedure needs to traverse multiple data centers, which results in significant overhead on the network, and non-ignorable increase in the delay of data processing. In contrast, the proposed SD-IoT aims to provide these resources in the same data center to reduce the delay and mitigate the overhead on the network.

## CONCLUSION

This article focused on the design of a flexible IoT architecture for smart urban sensing. Specifically, we proposed a software-defined IoT architecture that decouples the applications from underlying physical infrastructures. With this architecture, urban sensing applications can customize their own data acquisition, transmission, and processing through well-defined APIs, and multiple applications coexist on the shared infrastructure to further reduce the overall capital and maintenance cost. As a result, this architecture enables flexible control and management of physical infrastructure, and accelerates application innovation.

## REFERENCES

- [1] M. Naphade *et al.*, “Smarter Cities and Their Innovation Challenges,” *Computer*, vol. 44, no. 6, 2011, pp. 32–39.
- [2] D. Boyle *et al.*, “Urban Sensor Data Streams: London 2013,” *IEEE Internet Comp.*, vol. 17, no. 6, 2013, pp. 12–20.

- [3] A. Gluhak *et al.*, “A Survey on Facilities for Experimental Internet of Things Research,” *IEEE Commun. Mag.*, vol. 49, no. 11, 2011, pp. 58–67.
- [4] H. Ma, D. Zhao, and P. Yuan, “Opportunities in Mobile Crowd Sensing,” *IEEE Commun. Mag.*, vol. 52, no. 8, 2014, pp. 29–35.
- [5] N. McKeown *et al.*, “OpenFlow: Enabling Innovation in Campus Networks,” *ACM SIGCOMM CCR*, vol. 38, no. 2, 2008, pp. 69–74.
- [6] J. Yick, B. Mukherjee, and D. Ghosal, “Wireless Sensor Network Survey,” *Computer Networks*, vol. 52, no. 12, 2008, pp. 2292–330.
- [7] A. Papageorgiou *et al.*, “Smart M2M Data Filtering Using Domain-Specific Thresholds in Domain-Agnostic Platforms,” *Proc. IEEE BigData Congress*, 2013, pp. 286–93.
- [8] I. Downes, L. B. Rad, and H. Aghajan, “Development of a Mote for Wireless Image Sensor Networks,” *Proc. COGIS'06*, 2006.
- [9] A. Gudipati *et al.*, “SoftRAN: Software Defined Radio Access Network,” *Proc. 2nd ACM HotSDN*, 2013, pp. 25–30.
- [10] T. Luo, H.-P. Tan, and T. Q. Quek, “Sensor OpenFlow: Enabling Software-Defined Wireless Sensor Networks,” *IEEE Commun. Lett.*, vol. 16, no. 11, 2012, pp. 1896–99.
- [11] L. Galluccio *et al.*, “SDN-WISE: Design, Prototyping and Experimentation of a Stateful SDN Solution for Wireless Sensor Networks,” *Proc. Infocom*, 2015, pp. 513–21.
- [12] J. Medved *et al.*, “Opendaylight: Towards a Model-Driven SDN Controller Architecture,” *Proc. 15th IEEE WoWMoW*, 2014, pp. 1–6.
- [13] Z. A. Qazi *et al.*, “Simple-Fying Middlebox Policy Enforcement Using SDN,” *ACM SIGCOMM CCR*, vol. 43, no. 4, 2013, pp. 27–38.
- [14] R. Ahmed and R. Boutaba, “Design Considerations for Managing Wide Area Software Defined Networks,” *IEEE Commun. Mag.*, vol. 52, no. 7, 2014, pp. 116–23.
- [15] D. Wu *et al.*, “UbiFlow: Mobility Management in Urban-Scale Software Defined IoT,” *Proc. Infocom*, 2015, pp. 208–16.

## BIOGRAPHIES

JIAQIANG LIU received his B.S. degree in electronic engineering from Tsinghua University, Beijing, China, in 2012, and is now a Ph.D. student at Tsinghua University. His research interests include software defined network, data center network, network function virtualization, and Internet of Things, etc.

YONG LI [M'09] received the B.S. and Ph.D degree from Huazhong University of Science and Technology and Tsinghua University in 2007 and 2012, respectively. During 2012 and 2013 he was a visiting research associate with Telekom Innovation Laboratories and Hong Kong University of Science and Technology, respectively. From 2013 to 2014 he was a visiting scientist with the University of Miami. He is currently a faculty member in the Department of Electronic Engineering, Tsinghua University. His research interests are in the areas of mobile computing and social networks, urban computing and vehicular networks, and network science and future Internet. Dr. Li has served as general chair, Technical Program Committee (TPC) chair, and TPC member for several international workshops and conferences. He is currently an associate editor of the *Journal of Communications and Networking* and *EURASIP Journal of Wireless Communications and Networking*.

MIN CHEN [SM'09] (minchen@ieee.org) is a professor at the School of Computer Science and Technology at HUST. He was an assistant professor in the School of Computer Science and Engineering at SNU from September 2009 to February 2012. He worked as a post-doctoral fellow in the Department of Electrical and Computer Engineering at UBC for three years. Before joining UBC he was a postdoctoral fellow at SNU for one and a half years. His research focuses on Internet of Things, machine-to-machine communications, body area networks, body sensor networks, e-healthcare, mobile cloud computing, cloud-assisted mobile computing, ubiquitous networks and services, mobile agents, multimedia transmission over wireless networks, and so on.

WENXIA DONG received the master's degree from Southwest Jiaotong University in 2008. She then joined Huawei, where she is now a network research engineer. Her research interests are in the areas of vehicular networks, SDN north-bound interface, social networks, and application of big data in networks. She has 13 patents as the first author.

DEPENG JIN received the B.S. and Ph.D. degrees from Tsinghua University, Beijing, China, in 1995 and 1999, respectively, both in electronics engineering. He is a professor at Tsinghua University and the chair of the Department of Electronic Engineering. Dr. Jin was awarded the National Scientific and Technological Innovation Prize (second class) in 2002. His research fields include telecommunications, high-speed networks, ASIC design, and future Internet architecture.

With the proposed architecture, urban sensing applications can customize their own data acquisition, transmission, and processing through well-defined APIs, and multiple applications coexist on the shared infrastructure to further reduce the overall capital and maintenance cost.

---

# LOW POWER WIDE AREA MACHINE-TO-MACHINE NETWORKS: KEY TECHNIQUES AND PROTOTYPE

While M2M communications have been developed for many years, major challenges still remain with their efficient implementation from the perspective of low energy consumption and wide coverage. To address these challenges, low power wide area (LPWA) technology is investigated as one of the potential candidate solutions.

*Xiong Xiong, Kan Zheng, Rongtao Xu, Wei Xiang, and Periklis Chatzimisios*

---

## ABSTRACT

As one of the fastest growing technologies, machine-to-machine (M2M) communications are expected to provide ubiquitous connectivity. M2M devices can be used for a wide range of emerging applications that have various communications requirements. While M2M communications have been developed for many years, major challenges still remain with their efficient implementation from the perspective of low energy consumption and wide coverage. To address these challenges, low power wide area (LPWA) technology is investigated as one of the potential candidate solutions. In this article, we first introduce some typical LPWA M2M application scenarios. Given their requirements, we highlight key techniques and standards that are explicitly designed for LPWA M2M communications. Finally, we present an LPWA prototype system to evaluate its performance and demonstrate its potential in bridging a technological gap for future Internet-of-Things (IoT) applications.

## INTRODUCTION

Machine-to-machine (M2M) communications enable direct connectivity among devices, which can be organized as a network in order to exchange information and perform actions without human intervention. It is an integral part of the Internet-of-Things (IoT), which can benefit end users from its countless range of applications [1].

There are some notable requirements and properties of M2M networks, such as low costs, low energy consumption, wide coverage, tolerable low latency, relatively low data throughput, etc. Among them, low energy consumption and wide coverage are the most desired features since M2M devices can be deployed at locations without main power and operate only on battery power. Moreover, numerous M2M devices are widely distributed in a wide variety of locations, some of which are difficult to reach either because they are remote or because they are underground or located deep inside buildings.

However, these requirements and properties cannot be well supported by existing wireless network technologies that are designed for human users, such as cellular mobile networks [2]. As a result, various systems have been specifically designed for M2M communications, such as Zigbee/IEEE 802.15.4 in wireless sensor networks (WSNs), Bluetooth, radio frequency identification (RFID), etc [3]. However, almost all these techniques are short-range connectivity solutions, which have difficulty meeting some of the requirements of M2M applications. In addition, another promising system based on IEEE 802.11ah, also known as Low-Power Wi-Fi, is now being developed to meet specific requirements of M2M networks, e.g., transmission range up to 1 kilometer in outdoor scenarios, low data rates, and low energy consumption [4]. It is still not suited for remote and underground areas.

Due to the limitation of the above short-range communications systems, low power wide area (LPWA) technology has been specifically designed with the objectives of low energy consumption and wide coverage. In particular, ultra-narrowband (UNB) and direct sequence spread spectrum (DSSS) modulation schemes have been proposed for the physical layer of LPWA M2M systems thanks to their excellent coverage performance. Meanwhile, in order to enable low energy consumption, the star topology and the random access method can be

## COMMUNICATIONS STANDARDS

employed in the MAC layer. A key feature of these LPWA techniques is that they provide a trade-off between the data rate, battery life, and deployment costs.

Such a trade-off is acceptable for most M2M applications, which do not require high data rates and low latency [5]. This is a significant difference between LPWA M2M systems and traditional wireless communication systems. Thus, LPWA M2M networks are expected to play a crucial role in the IoT. Meanwhile, standardization efforts for LPWA M2M networks are underway. In this paper, we first present an overview of typical LPWA M2M application scenarios. Then some key techniques that are currently under consideration for LPWA M2M systems are studied in detail. Several preliminary standards are compared with an emphasis on highlighting their key differences.

To the best of our knowledge, there still remain many open problems with the implementation of LPWA M2M networks, which are the main focus of this paper. Open source software-defined radio (SDR) platforms are appealing and efficient for fast prototyping LPWA M2M designs, where most modules are programmed by software and easily upgradable [7]. We implement an LPWA M2M prototype based upon the specifications of IEEE 802.15.4k by using an open source SDR based GNU radio real-time signal processing framework. To demonstrate the advantages of LPWA M2M networks, field experiments are conducted to evaluate the performance of our prototype in several deployment scenarios.

---

*Xiong Xiong and Kan Zheng are with Beijing University of Posts & Telecommunications.*

*Rongtao Xu is with Beijing Jiaotong University.*

*Wei Xiang is with the University of Southern Queensland.*

*Periklis Chatzimisios is with Alexander Technological Educational Institute of Thessaloniki (ATEITHE).*



## APPLICATION SCENARIOS

LPWA M2M networks are suitable for deployment in a broad range of application scenarios, which can be classified into the five categories (see Fig. 1): infrastructure monitoring, transportation, asset tracking, security, and healthcare.

**Infrastructure Monitoring:** Under this scenario, an LPWA M2M device equipped with sensors (i.e. smart meters) plays a significant role in measuring and reporting usage, functionality, and environmental data to an LPWA M2M network. The resources under monitoring (such as water, electricity, oil, or gas) are crucial for our modern society, so ensuring an adequate supply of such resources is currently a major focus for countries around the world.

**Transportation:** Transportation related applications enable monitoring of critical transportation conditions, e.g. road conditions, traffic congestion, public transport, and so on. LPWA M2M networks can offer the opportunity to establish a direct interaction between vehicles and information systems/centers. The vehicle's status information can help service providers develop new applications, and help relevant government agencies improve the efficiency of managing traffic and public transportation in city premises.

**Asset Tracking:** LPWA M2M networks offer a low-cost and reliable solution to asset tracking by owners, who do not have to actively manage their assets. Asset tracking is widely used in the manufacturing industry to achieve real-time visibility throughout the supply and distribution chains, monitor inventory levels and condition, and manage the overall processes and workflows, while improving product quality and reducing the costs of waste and disruption.

**Security:** Security services are imperative for end users to minimize their personal losses and to predict risks in advance. When security systems are breached, LPWA M2M networks can facilitate secure communications for emergency services by allowing them to transmit critical emergency information to building owners or national security agencies.

**Healthcare:** With the growing demand for disease treatment, the healthcare industry now faces many challenges, including a vast cost burden. LPWA M2M networks can provide a solution that decentralizes medical costs and hospital care. This trend toward user-centric healthcare and individually tailored medicine may drive the market for M2M healthcare applications in the foreseeable future.

As mentioned above, all these applications require wide coverage and extremely low energy consumption, which can be provided by LPWA M2M networks. The performance requirements of different applications are summarized in Table 1.

## KEY TECHNIQUES AND EARLY STANDARDS IN LPWA M2M NETWORKS

One of the most important requirements of an LPWA system is to provide ubiquitous communications. This implies that the architecture similar to a cellular network may be a possible solution. However, network components must be

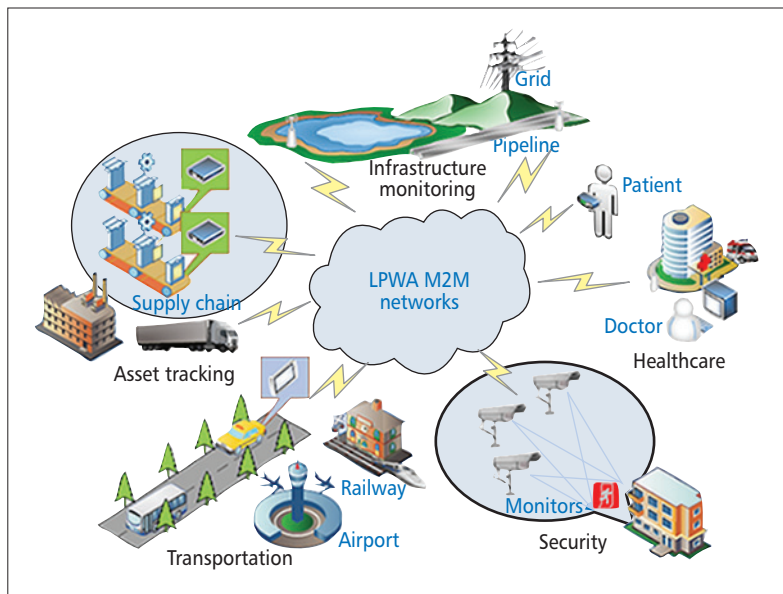


Figure 1. Illustration of five typical LPWA M2M application scenarios.

substantially simplified in comparison to cellular networks for easy deployment and maintenance [8]. Moreover, outdoor, indoor, and underground areas need to be covered by LPWA systems. In order to meet these performance demands of LPWA M2M networks, new advanced techniques have been proposed in both the physical (PHY) and medium access control (MAC) layers of LPWA M2M networks.

### PHY TECHNIQUES

The PHY layer design of LPWA M2M networks is instrumental in enabling wide coverage. A key requirement is to design a transceiver with an ultra-low receiver sensitivity threshold in dBm. As such, two candidate techniques, i.e. UNB modulation or DSSS, have become attractive solutions lately.

**UNB Modulation:** UNB modulation was initially proposed as a promising technique to improve spectral efficiency. However, the UNB technique is used for long distance communications in LPWA M2M networks. The typical procedure of UNB modulation consists of two significant stages, i.e. the abrupt phase shift modulation and UNB filtering stages. An input signal is first modulated and then passed through a UNB filter. Two key issues arise in designing a UNB system: the signal waveform of the modulator and the UNB filter.

For abrupt phase shift modulation, the phase of a carrier is abruptly switched to represent a digital one or zero. There are various abrupt phase shift modulation schemes that can be chosen for UNB, e.g. very maximum sideband keying (VMSK), pulse position phase shift keying (3PSK), and pulse position phase reversal keying (3PRK), also known as missing cycle modulation (MCM) [11]. Thanks to the abrupt phase shift property, the power spectral density (PSD) of the modulated signal waveform consists of a higher discrete component and a lower continuous component. All the information required for detection exists only in a single sideband of the discrete component without

Application category	Typical user case	Coverage	Power consumption	Data traffic	Periodicity	Mobility	Real-time requirement	Security/reliability requirement
Infrastructure monitoring	Water/Electric/Gas meter	Urban areas	Low	Medium	Tens of minutes	No	Medium	Low
	Agriculture/soil & oil/gas pipeline monitoring	Open fields	Low	Low	Event driven	No	Low	Medium
Transportation	Traffic congestion monitoring	Urban areas	Low	High	Tens of minutes	High	Medium	Low
	Public transport management	Urban areas	Low	Medium	Event driven	High	Medium	Medium
Asset tracking	Supply chain monitoring	Urban areas/ in-building	Low	Low	Event driven	Medium	Low	Low
	Vehicle tracking	Urban areas/ open fields	Low	High	Several minutes	High	Medium	Low
Security	Access control & building security systems	In-building	Low	Low	Event driven	No	High	High
	Natural disasters preparedness	Urban areas/ open fields	Low	Low	Event driven	No	High	Medium
Healthcare	Health status monitoring	Urban areas/ in-building	Low	Medium	Tens of minutes	Medium	Low	Low
	Medical alert	Urban areas/ in-building	Low	Low	Event driven	Medium	High	Low

**Table 1.** Requirements of LPWA M2M applications.

measurable bandwidth, which is a single frequency containing phase reversal modulation. Then a UNB filter with a zero or negative group delay is able to filter all the other sidebands and harmonics so as to keep the only single sideband. Conventional filters with a group delay disrupt the abrupt phase shift information in the modulated signal. Thanks to the UNB filter, the output UNB signal contains a single frequency requiring a transmission bandwidth of only 1 Hz (in theory) or several Hz in practice. As a result, the noise power is greatly reduced. Thus, the receiver sensitivity threshold in dBm of the UNB receiver can be extremely low.

To date there is no open acceptable method to implement the UNB band filter and baseband UNB modulation. At the RF level, the filters are very complex and must be hand tuned. It is impossible to configure a zero group delay narrow band filter using finite impulse response (FIR) or infinite impulse response (IIR) filters [11]. This limitation poses a big challenge for commercial UNB products.

*DSSS:* The DSSS technique has been widely used in commercial wireless communications systems, e.g. the third generation (3G) mobile communications networks. In the DSSS system, the large process gain allows a DSSS receiver to successfully detect signals with a very low carrier-to-interference ratio, which is essential for implementing the LPWA M2M network. Such a feature can be exploited thoroughly by employing a very long spreading sequence, e.g. the receiver sensitivity threshold may decrease by 3 dB when the length of the spreading sequence doubles.

Unlike conventional DSSS systems, the spreading sequences employed in LPWA M2M networks always utilize a much larger SF to make the receiver sensitivity threshold low enough to extend the communications range. However, employing long spreading sequences brings about its own challenges. For example, with an increase in the length of the spreading sequence, the computational complexity of the DSSS transceiver increases exponentially. Therefore, more efficient digital signal processing algorithms become necessary in the hardware design. Nevertheless, DSSS is still one of the most promising techniques for the LPWA system.

### MAC TECHNIQUES

An efficient MAC layer design plays a vital role in improving the energy efficiency of LPWA M2M networks [12, 13]. It is also responsible for overcoming the challenge of potential massive accesses in LPWA M2M networks. To achieve these goals, the topology and channel access techniques have to be carefully considered.

*Star Topology:* In wireless networks, topology directly impacts the performance of the network including scalability, energy efficiency, reliability, data latency, overhead, etc. There are several topologies widely adopted in traditional WSNs, such as the mesh, cluster, tree, and chain topologies [14]. These topologies share a common feature, where an end device (ED) doubles as a router to form multi-hop links in order to extend the communications coverage.

However, the features of LPWA M2M networks are greatly different from those of WSNs.

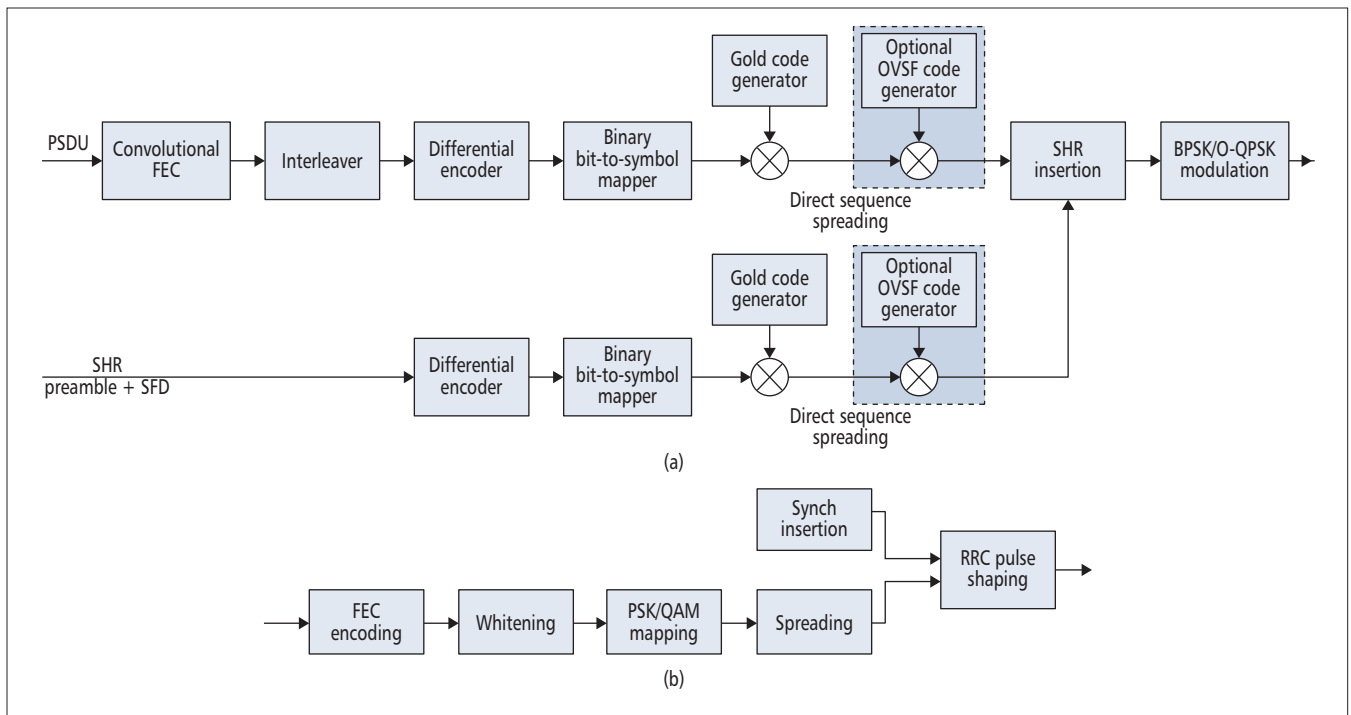


Figure 2. LPWA transmitter block diagram: a) IEEE 802.15.4k [6]; b) Weightless [10].

Since the physical layer techniques of the LPWA system are capable of providing coverage wide enough, the main concerns for the LPWA system when selecting a proper topology are the low costs and energy consumption. Compared to the aforementioned topologies, the star topology with only a single hop is considered to be the best choice for LPWA M2M networks. In this topology, each ED distributed in the coverage area communicates directly with a centralized access point (AP). There are various advantages of using such a simple topology. In particular, direct communications between an endpoint and the AP help minimize transmission latency. There are also no unnecessary packets for routing and multi-hop communications in the star topology, which help the EDs to reduce and balance their energy consumption. Besides, the star topology is easy for deployment. In terms of reliability and robustness, the star topology is not as good as the mesh and cluster topologies, but better than its tree and chain counterparts.

**Channel Access:** Channel access is probably the most crucial issue for LPWA M2M networks, since a large number of devices need to be simultaneously served in a large coverage area. There are two main categories of channel access methods for sharing access to the wireless medium: reservation-based access and contention-based access.

Specifically, reservation-based methods require the channelization of radio resources in different dimensions, which offer various advantages (e.g. reduced collisions and guaranteed delivery delays, etc). However, reservation-based access methods are unable to accommodate a vast number of devices in LPWA M2M networks due to limited channel resources.

Compared to reservation-based access methods, contention-based access methods, also

dubbed random access methods, are more appealing for LPWA systems. In a contention-based approach, devices compete for access to the shared medium as required, which is a relatively simple and flexible process. With this on-demand characteristic, reservation-based access methods may have to handle massive and unpredictable access activities. Besides, synchronization is no longer needed, which helps reduce energy consumption due to synchronization signaling, and enables the devices to sleep so as to save energy. However, contention-based access methods have to deal with other sources of energy waste, e.g. idle listening, overhearing, collisions, overhead, etc.

Several well-known solutions have been proposed for the above issues. A well-known protocol is Carrier Sense Multiple Access with Collision Avoidance (CSMA/CA), which employs the carrier sensing mechanism to reduce channel collisions. However, CSMA-based protocols cannot work well all the time in LPWA M2M networks. In some extreme cases, the signals transmitted from other EDs may not be detected by a target ED because of the large path loss caused by a long distance between the EDs. As a result, the carrier sensing mechanism is neither practical nor valid under this scenario. Another challenge is the massive connectivity requests that the LPWA M2M network has to handle, which renders CSMA-based protocols inefficient due to simultaneous massive requests for channel access.

Compared with CSMA-based protocols, Aloha-based protocols offer better performance when coupled with a properly designed PHY technique in LPWA M2M networks. For instance, adopting the DSSS technique in the PHY layer provides the LPWA M2M network with the capability of detecting and identifying multiple EDs arriving simultaneously. When

Our prototype network is deployed with the star topology. The AP can be configured as a multiple-baseband receiver with various physical parameters, e.g. different spreading factors, different seeds of the Gold code, and so on. The AP can collect all packets from EDs with different PHY configurations within its serving area.

	Attribute	IEEE 802.15.4k (DSSS)	Weightless
PHY	Operation frequency band	470 ~ 510 MHz; 779 ~ 787 MHz; 863 ~ 870 MHz; 902 ~ 928 MHz; 915 ~ 928 MHz; 917.1 ~ 923.5 MHz; 920 ~ 928 MHz; 921 ~ 928 MHz and 2.4 ~ 2.4835 GHz in different counties	470 ~ 790 MHz in Europe; 470 ~ 698 MHz in U.S.
	Channel bandwidth	100 kHz; 200 kHz; 400 kHz; 600 kHz; 800 kHz and 1 MHz	8 MHz in Europe; 6 MHz in U.S.
	Effective isotropic radiated power (EIRP)	Minimum: -3 dBm; Maximum: limited by local regulatory bodies	4 ~ 32 dBm
	FEC	Convolutional encoding; rate 1/2, constraint length 7	Convolutional encoding; rate 3/4 or 1/2, constraint length 7
	Interleaving	Pruned bit reversal interleaving algorithm	Matrix interleaving with 8 columns
	Spreading sequence	Gold code; SF 16 ~ 32768	Gold code and Kasami code: SF 15 ~ 1023
	Modulation	BPSK; OQPSK	16-QAM; pi/4-QPSK; pi/2-BPSK; pi/4-DQPSK; pi/2-DBPSK
	Frequency hopping	No	Yes
	Minimum receiver sensitivity threshold	-148 dBm	Downlink: -128 dBm; Uplink: -140 dBm
	Typical coverage	Up to 20 km in LoS and 5 km in NLoS	Up to 10 km
	Data rate	0.00153 ~ 125 kb/s	Downlink: 0.0025 ~ 16.0 Mb/s; Uplink: 0.00025 ~ 0.5 Mb/s
Sync sequence length	Preamble: 0/2/4 octets; SFD: 0/1 octets	8 ~ 2048 Symbol (No need to multiply by spreading sequence any more)	
MAC	Packet length	16/24/32 octets	0 ~ 255 octets
	Topology structure	Star	Star, with multi-hop relay capability
	Channel access method	CSMA/CA; CSMA/CA with PCA; Aloha with PCA	TDMA/FDMA
	Traffic priority	Yes	Yes

Table 2. Comparison between IEEE 802.15.4k and Weightless systems.

multiple EDs transmit packets with the Aloha mode, the receiver can detect and decode each packet successfully with the specified spreading sequences, even if the transmission periods of the EDs overlap completely. Therefore, the potential massive access issue in the LPWA M2M network may be properly dealt with. Moreover, the unnecessary overhead of the carrier sensing mechanism, such as the request to send (RTS) and clear to send (CTS) control messages, is eliminated in Aloha-based protocols, helping meet the energy efficiency requirement of LPWA M2M networks.

Also, interference cancellation techniques can be used by the receiver to improve the efficiency of channel access. This brings another new perspective of MAC protocol design for massive access in LPWA M2M networks.

## EARLY STANDARDS

Several early standards for M2M communications have been developed by different standardization organizations, such as the ETSI technical committee (TC), 3GPP, IETF, and so on. Among them, IEEE 802.15.4k and Weightless have been proposed for LPWA applications.

IEEE 802.15.4k aims at low energy critical infrastructure monitoring (LECIM) networks, to facilitate point-to-multi-point communications for monitoring and managing critical infrastructure applications [6]. In the standard, two PHY modes are specified to support LECIM applications, i.e. DSSS and frequency shift keying (FSK). The transmitter block diagram of DSSS PHY is illustrated in Fig. 2a.

The Weightless draft standard is developed

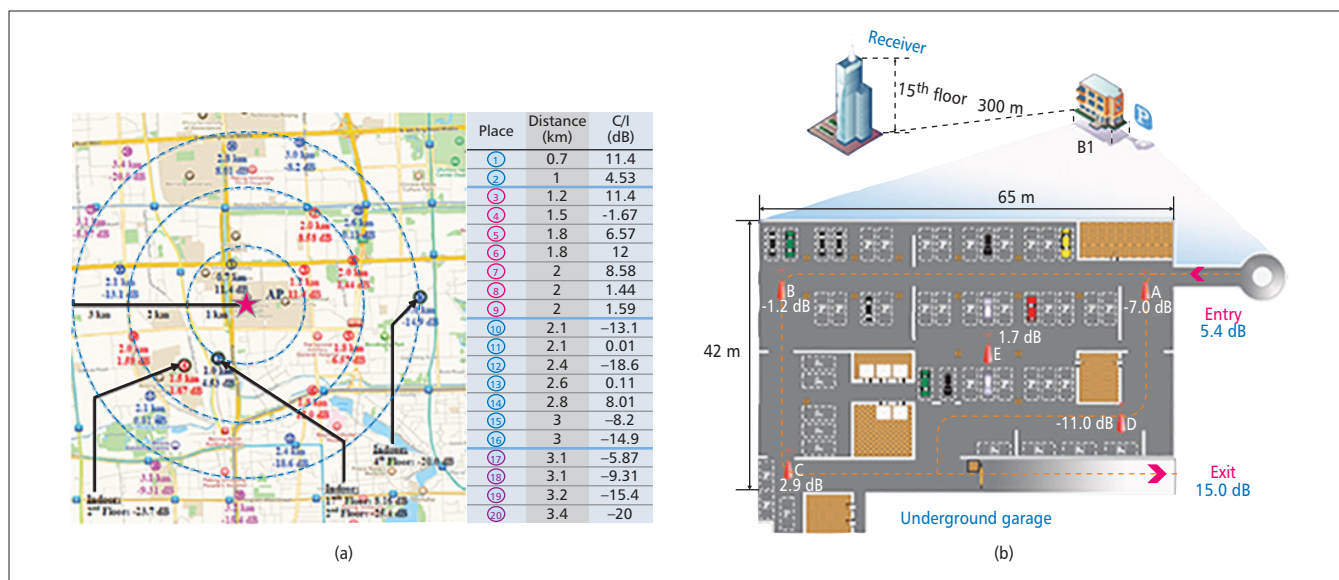


Figure 3. Coverage performance of the developed LPWA prototype: a) outdoor/indoor; b) underground.

by the Weightless special interest group (SIG) [9]. The Weightless specifications define not only the PHY and MAC layers, but also an upper layer, dubbed the server layer. A basic transmitter block diagram of PHY is depicted in Fig. 2b.

The IEEE 802.15.4k and Weightless standards have different attributes and also share some common features. To elaborate, a detailed comparison between IEEE 802.15.4k and Weightless is summarized in Table 2.

## PROTOTYPE SYSTEM AND RESULTS

### PROTOTYPE SYSTEM BASED ON SDR

In order to evaluate the performance of the LPWA M2M network, we have developed an IEEE 802.15.4k prototype based on SDR, which consists of GNU radio and universal software radio peripheral (USRP) [15].

The AP consists of an industrial personal computer (IPC) for baseband signal processing, and an Ettus USRP B210 as the RF transmitter and receiver. Powered by IPC via USB3.0, the USRP B210 can operate on different carrier frequency ranging from 70 MHz to 6 GHz covering all the operating frequency bands of LPWA devices. The IPC is assembled with a high-performance multi-core CPU, and runs the programs for both transmitter and receiver signal processing of 802.15.4k DSSS PHY.

At the receiver, non-coherent detection is used for symbol demodulation to eliminate the effect of the frequency and phase offsets. Then a parallel preamble detection scheme is designed to find the beginning of a packet. We apply a fast Fourier transform (FFT)-based algorithm to implement the correlation process for preamble detection. To decode traffic data, a dynamic timing adjustment algorithm is proposed to obtain the optimal sampling time. The FFT transform and the Viterbi algorithm for convolutional decoding are derived from the IT++ library.<sup>1</sup>

In our developed prototype, Gaussian frequency shift keying (GFSK) is chosen as the

modulation scheme, which is different from the design in IEEE 802.15.4k, because GFSK can overcome the effect of the frequency or phase offset. Although performance degradation is unavoidable due to non-coherent demodulation, the DSSS PHY is able to tolerate an acceptable frequency offset at very low SNRs.

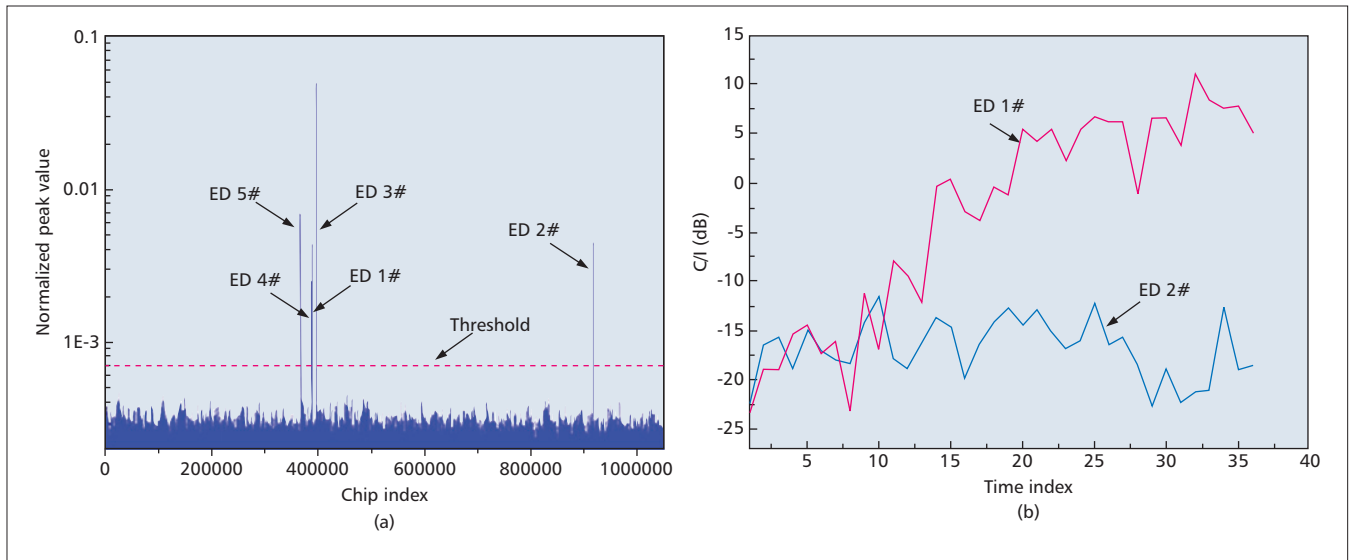
Our prototype network is deployed with the star topology. The AP can be configured as a multiple-baseband receiver with various physical parameters, e.g. different spreading factors, different seeds of the Gold code, and so on. The AP can collect all packets from EDs with different PHY configurations within its serving area.

### EXPERIMENT RESULTS AND DISCUSSIONS

Our experiments are carried out in an urban environment. The wide antenna of the AP is installed on the roof of a 15-floor building in the campus. Then the received RF signal is fed into the USRP B210 located in our laboratory with a 10-meter cable connected to the antenna. Preamble detection is performed at the chip level in parallel in order to capture all possible frame heads. Also, the cyclic redundancy check (CRC) bits of the decoded payload data are checked to ensure data integrity so that false preamble detection can be disregarded. In our field tests, the carrier-to-interference (C/I) value is approximated as the ratio of the normalized peak output value of the preamble detector to the received signal power, which can be used as a metric measuring communications reliability in the LPWA M2M network. Moreover, it is worth mentioning that received packets can be decoded without errors in most cases where C/I is greater than -30 dB.

Some common parameters of our field tests are stated as follows. Both the AP and EDs operate at 433 MHz with a symbol rate 200 k symbols per second. In all packets sent by the EDs, the physical protocol data unit (PPDU) that consists of a four-octet preamble, and a 16-octet physical service data unit (PSDU) is spread by the Gold code with a SF of 32,768.

<sup>1</sup> <http://itpp.sourceforge.net/devel/tutorial.html>



**Figure 4.** Multi-user performance of the developed LPWA prototype: a) MU-case 1 for preamble detection; b) MU-case 2 for the near-far effect.

*Coverage Performances:* For the first time, coverage in various application environments is evaluated in our experiments, in which the transmit power of the ED is fixed to 15 dBm. Apart from the outdoor environment, tests are also performed in both the indoor and underground environments.

**Outdoor/Indoor:** Figure 3a shows the experiment results for the outdoor/indoor scenarios, where 20 spots are chosen to conduct the field tests. In each spot, the upper number stands for the distance from the AP, while the lower number indicates the average C/I. In all the cases, there is no LoS between the transmitter and receiver.

Since the radio propagation characteristics in urban areas are too complicated to predict, the received C/I in the testing spots is not always larger with a smaller communications distance. For instance, spot 15 is further than spot 10, but the former has a higher C/I. However, in all the spots in the outdoor scenarios, the packets sent by the EDs can be successfully decoded at the AP. Moreover, the measured C/I value at the farthest spot with 3.4 km distance is still larger than  $-30$  dB, i.e. the C/I threshold to ensure correct decoding of the received packets. It is believed that wide coverage can be provided by our LPWA prototype, which will be validated in our next steps. Furthermore, due to the low receiver sensitivity threshold in dBm of the LECIM device, the EDs close to the AP can decrease their transmit power to reduce energy consumption.

When the EDs are installed inside buildings for monitoring or industry metering, building penetration is more likely to be a major performance constraint. Therefore, the indoor environments in buildings located in 2, 4, and 16 in Fig. 3a are selected to perform the indoor tests. In the building with 23 floors located in spot 2, we place the testing ED on different floors and collect the measured C/I, i.e. 8.16 dB on the 17th floor and  $-25.4$  dB on the 2nd floor, respectively. It can be found that the height of the ED location significantly affects performance. That is, the C/I increases with the height of the position. Although

the distances of buildings 2, 4, and 16 are different, i.e. 1.0 km, 1.5 km, and 3.0 km, respectively, the received packets are all decoded successfully at the AP, which demonstrates the effectiveness of our prototype in indoor environments.

**Underground:** To further evaluate the coverage performance of LPWA M2M networks, we carry out several field tests in an underground environment. Figure 3b gives the measurement results in an underground car park, which is located at a distance of 300 meters from the AP. EDs are installed in five typical places at the basement level for measurement purposes. The wireless signal can traverse through the exit or entry of the car park to reach the AP. The C/I value is shown for each place in Fig. 3b. The results indicate that LPWA M2M networks are also able to ensure reasonable underground coverage.

*Multi-User Performance:* Due to the large SFs in DSSS, signals from other devices are likely to be undetectable because they may fall below the effective noise floor. As a result, sensing the wireless channel is neither practical nor useful, and thus the MAC protocol works under the Aloha mode. Due to thousands of EDs in the LPWA M2M network, the packets from the EDs randomly arrive at the AP asynchronously. Collisions seldom occur even if packets arrive at the AP synchronously within one chip difference in time.

The multi-user (MU) performance of LPWA M2M networks is evaluated through two field tests, i.e. *MU-Case 1* for preamble detection and *MU-Case 2* for the near-far effect test. In *MU-Case 1*, five EDs are deployed on the second floor of the same building as shown in Fig. 3b, located 300 meters from the AP. Due to the limited number of EDs, the transmission timing of each ED is adjusted manually to ensure that the packets from all the EDs are able to arrive at the AP nearly simultaneously. Thus, this is similar to the highly concurrent traffic in an LPWA M2M network with a large number of EDs. Fig. 4a shows that the preamble of each packet

sent from each ED is properly detected in the first half of the FFT window. Although the packets from ED 1 and ED 4 arrive at the receiver almost at the same time, their preambles are completely detectable and distinguishable.

The near-far effect of DSSS may degrade detection performance. However, this impairment has little impact on the LPWA M2M network because of the high spreading factors. The C/I values in *MU Case 2* are plotted in Fig. 4b. The transmit power of ED 1 is fixed, while the transmit power of ED 2 increases step by step. The dynamical fluctuation of the two curves is due to channel fading and variation. It can be observed that the transmit power of ED 2 has little effect on the C/I of ED 1.

## CONCLUSIONS

LPWA M2M networks can meet the requirements of a broad range of IoT applications. Due to the ubiquitous coverage and low power consumption of the LPWA M2M network, endpoint devices are able to remain connected for an extended period of time. We first discussed and analyzed several application scenarios of LPWA M2M networks as well as the key techniques and early standards. Then field experiments were carried out with our developed SDR prototype to evaluate the field trial performances in urban environments. The large coverage performances are shown, e.g. more than 3 km in the outdoor environment and approximately 1 km in the indoor environment. Furthermore, the multi-user performance results are also presented to validate the MAC design. Our results clearly have demonstrated the feasibility and effectiveness of our developed prototype, as well as the advantages of LPWA M2M networks.

## ACKNOWLEDGMENT

This work was supported by the National High-Tech R&D Program under grant 2015AA01A705, the China Natural Science Funding under grant 61271183, the National Key Technology R&D Program of China under grant 2015ZX03002009-004, and the Fundamental Research Funds for the Central Universities under grant 2014ZD03-02.

## REFERENCES

- [1] K. Zheng *et al.*, "Radio Resource Allocation in LTE-Advanced Cellular Networks with M2M Communications," *IEEE Commun. Mag.*, vol. 50, no. 7, July 2012, pp. 184–92.
- [2] L. Lei *et al.*, "Performance Analysis of Device-to-Device Communications with Dynamic Interference Using Stochastic Petri Nets," *IEEE Trans. Wireless Commun.*, vol. 12, no. 12, Dec. 2013, pp. 6121–41.
- [3] Z. M. Fadlullah *et al.*, "Toward Intelligent Machine-to-Machine Communications in Smart Grid," *IEEE Commun. Mag.*, vol. 49, no. 4, Apr. 2011, pp. 60–65.
- [4] T. Adame *et al.*, "IEEE 802.11AH: The WiFi Approach for M2M Communications," *IEEE Wireless Commun.*, vol. 21, no. 6, Dec. 2014, pp. 144–52.
- [5] J. Morrish, "Low Power Wide Area Wireless Network Technologies for M2M: A Technical Comparison of On-Ramp, Sigfox and Weightless," *Machina Research*, June 2013.
- [6] "IEEE Standard for Local and Metropolitan Area Networks Part 15.4: Low-Rate Wireless Personal Area Networks (LR-WPANs) Amendment 5: Physical Layer Specifications for Low Energy, Critical Infrastructure Monitoring Networks," *IEEE Std 802.15.4k-2013*, Aug. 2013.
- [7] T. Ullersoy, "Software Defined Radio: Challenges and Opportunities," *IEEE Commun. Surveys & Tutorials*, vol. 12, no. 4, 4th Quarter 2010, pp. 531–50.
- [8] S. Bayat *et al.*, "Distributed Data Aggregation in Machine-to-Machine Communication Networks Based on Coalitional Game," *IEEE Wireless Communications and Networking Conf. (WCNC)*, Istanbul, Turkey, Apr. 2014, pp. 2026–31.

- [9] W. Webb, *Dynamic White Space Spectrum Access*, Webb Search Limited, Oct. 2013.
- [10] *Weightless System Specification (version 0.8)*, Weightless SIG, July 2012.
- [11] H. R. Walker, *Ultra Narrow Band Modulation Textbook*, Dec. 2014, <http://www/vmsk.org/>.
- [12] K. Zheng *et al.*, "Challenges of Massive Access in Highly Dense LTE-Advanced Networks with Machine-to-Machine Communications," *IEEE Wireless Commun.*, vol. 21, no. 3, Jun. 2014, pp. 12–18.
- [13] A. Bachir *et al.*, "MAC Essentials for Wireless Sensor Networks," *IEEE Commun. Surveys & Tutorials*, vol. 12, no. 2, 2nd Quarter 2010, pp. 222–48.
- [14] Q. Mamun, "A Qualitative Comparison of Different Logical Topologies for Wireless Sensor Networks," *Sensors*, vol. 12, no. 11, Nov. 2012, pp. 14998–13.
- [15] X. Xiong *et al.*, "Implementation and Performance Evaluation of LECIM for 5G M2M Applications with SDR," *IEEE Globecom Wks.ps (GC Wksps.)*, Austin, USA, Dec. 2014, pp. 612–17.

## BIOGRAPHIES

XIONG XIONG received his B.S. degree from Beijing University of Posts&Telecommunications (BUPT), China, in 2013. Since then he has been working toward a Ph.D. degree at BUPT. His research interests include M2M networks and software defined radio.

KAN ZHENG [SM'09] is currently a full professor at Beijing University of Posts & Telecommunications (BUPT), China. He received the B.S., M.S., and Ph.D. degrees from BUPT, China, in 1996, 2000, and 2005, respectively. He has extensive experience in the research and standardization of the new emerging technologies. He is the author of more than 200 journal articles and conference papers in the field of wireless networks, M2M networks, VANET, and so on. He holds editorial board positions for several journals. He has organized several special issues in famous journals including *IEEE Communications Surveys & Tutorials*, *IEEE Communication Magazine*, and *IEEE System Journal*.

RONGTAO XU received his bachelor degree in radio technology from Xian Jiaotong University in 1997, the master degree in communication and information system from Beijing University of Posts and Telecommunications in 2000, and the Ph.D. in electronic and information engineering from Hong Kong Polytechnic University in 2007. From 2000 to 2003 he worked at Siemens Ltd. China as a system engineer. In 2007 he joined Beijing Jiaotong University, where he is currently an associate professor with the State Key Laboratory of Rail Traffic Control and Safety. His research areas include wideband mobile communications, railway communications, and wireless sensor network.

WEI XIANG [SM'10] received the B.Eng. and M.Eng. degrees, both in electronic engineering, from the University of Electronic Science and Technology of China, Chengdu, China, in 1997 and 2000, respectively, and the Ph.D. degree in telecommunications engineering from the University of South Australia, Adelaide, Australia, in 2004. Since January 2004 he has been with the School of Mechanical and Electrical Engineering, University of Southern Queensland, Toowoomba, Australia, where he is currently an associate professor. He has been awarded a number of prestigious fellowship titles, including the Queensland International Fellow (2010–2011) by the Queensland Government of Australia, the Endeavour Research Fellow (2012–2013) by the Commonwealth Government of Australia, the Smart Futures Fellow (2012–2015) by the Queensland Government of Australia, and the JSPS Invitational Fellow (2014–2015) by the Japan Society for the Promotion of Science (JSPS). He received the Best Paper Award at the 2011 IEEE WCNC, and the USQ Excellence in Research Award in 2013. He is an editor of *IEEE Communications Letters*. His research interests are in the broad area of communications and information theory, particularly coding and signal processing for multimedia communications systems.

PERIKLIS CHATZIMISIOS serves as an associate professor at the Computing Systems, Security and Networks (CSSN) Research Lab of the Department of Informatics at the Alexander TEI of Thessaloniki (ATEITHE), Greece. Recently he was a visiting academic/researcher at the University of Toronto (Canada) and Massachusetts Institute of Technology (USA). Dr. Chatzimisios is involved in several standardization activities serving as a member of the Standards Development Board for the IEEE Communication Society (ComSoc) (2010 to today), and recently as an active member of the IEEE Research Groups on IoT Communications & Networking Infrastructure and on Software Defined & Virtualized Wireless Access. He is also very active in IEEE activities such as serving as the vice chair of the Emerging Technical Subcommittee on Big Data (TSCBD) and the secretary of the IEEE Technical Committee on Cognitive Networks (TCCN) (2012–2014). He has served as organizing/TPC committee member for more than 150 conferences and as a founder/organizer/co-chair for many workshops that are co-allocated with major IEEE conferences. He also holds editorial board positions for several IEEE and non-IEEE journals, and he is the director (co-director from 2012 to 2014) for the E-letter of the IEEE Technical Committee on Multimedia Communications (MMTC). He is the author/editor of eight books and more than 100 peer-reviewed papers and book chapters on the topics of performance evaluation and standardization activities of mobile/wireless communications, quality of service/quality of experience and vehicular networking. His published research work has received more than 1500 citations by other researchers. He received his Ph.D. from Bournemouth University (UK) (2005) and his B.Sc. from Alexander TEI of Thessaloniki (Greece) (2000).

LPWA M2M networks can meet the requirements of a broad range of IoT applications. Due to the ubiquitous coverage and low power consumption of the LPWA M2M network, endpoint devices are able to remain connected for an extended period of time.

# TOWARD BETTER HORIZONTAL INTEGRATION AMONG IOT SERVICES

Several divergent application protocols have been proposed for IoT solutions. Each protocol focuses on a specific aspect of IoT communications. The lack of a protocol that can handle the vertical market requirements of IoT applications has resulted in a fragmented market between many protocols. In turn, this fragmentation is a main hindrance in the development of new services that require the integration of multiple IoT services to unlock new capabilities and provide horizontal integration among services.

*Ala Al-Fuqaha, Abdallah Khreishah, Mohsen Guizani, Ammar Rayes, and Mehdi Mohammadi*

## ABSTRACT

Several divergent application protocols have been proposed for Internet of Things (IoT) solutions including CoAP, REST, XMPP, AMQP, MQTT, DDS, and others. Each protocol focuses on a specific aspect of IoT communications. The lack of a protocol that can handle the vertical market requirements of IoT applications including machine-to-machine, machine-to-server, and server-to-server communications has resulted in a fragmented market between many protocols. In turn, this fragmentation is a main hindrance in the development of new services that require the integration of multiple IoT services to unlock new capabilities and provide horizontal integration among services.

In this work, after articulating the major shortcomings of the current IoT protocols, we outline a rule-based intelligent gateway that bridges the gap between existing IoT protocols to enable the efficient integration of horizontal IoT services. While this intelligent gateway enhances the gloomy picture of protocol fragmentation in the context of IoT, it does not address the root cause of this fragmentation, which lies in the inability of the current protocols to offer a wide range of QoS guarantees. To offer a solution that stems the root cause of this protocol fragmentation issue, we propose a generic IoT protocol that is flexible enough to address the IoT vertical market requirements. In this regard, we enhance the baseline MQTT protocol by allowing it to support rich QoS features by exploiting a mix of IP multicasting, intelligent broker queuing management, and traffic analytics techniques. Our initial evaluation of the lightweight enhanced MQTT protocol reveals significant improvement over the baseline protocol in terms of the delay performance.

## INTRODUCTION

IoT devices can be classified into two major categories: resource-constrained devices and resource-rich devices. We define resource-rich devices as those that have the hardware and

software capabilities to support the TCP/IP protocol suite. For devices that support the TCP/IP protocol suite, IoT applications are implemented on top of a variety of application-level protocols and frameworks, including the Constrained Application Protocol (CoAP), Representational State Transfer (REST), Extensible Messaging and Presence Protocol (XMPP), Advanced Message Queuing Protocol (AMQP), Message Queue Telemetry Transport (MQTT), MQTT for Sensor Networks (MQTT-SN), Data Distribution Service (DDS)<sup>1</sup>, and others. On the other hand, devices that do not have the required resources to support TCP/IP cannot interoperate easily with resource-rich devices that support the TCP/IP suite. For example, microcontroller-based appliances and gadgets should have the capability to interoperate with other IoT elements that are TCP/IP enabled. Beyond the interoperability issues between devices that support TCP/IP and those that do not, IoT devices utilize a variety of protocols leading to a myriad of interoperability issues that limit the potential applications of the IoT. This fragmentation between the protocols utilized for communication within and across resource-constrained devices and resource-rich devices is not foreseen to change in the near future.

Furthermore, the communication patterns between the different entities that comprise the IoT can be classified into three main categories based on the nature of the communicating devices: machine-to-machine (M2M), machine-to-server (M2S), and server-to-server (S2S). M2M communications between IoT devices that can be resource-constrained or resource-rich typically require timely data delivery. Depending on the application, the response time might be required to be as low as 10  $\mu$ s. DDS is an object management group (OMG) data-centric middleware standard for real-time M2M communication. DDS utilizes a broker-less architecture that relies on multicasting to offer various quality of service (QoS) guarantee and reliability control capabilities even in scenarios that involve massive fan-outs (i.e. a large number of destinations for a given message). On the other hand, M2S and S2S communication does not typically require the stringent delay performance needed for M2M communication. Therefore, IoT typically employs broker-based protocols like MQTT and AMQP to offload the IoT devices from having to handle a large number of upstream server requests. Moreover, M2S and S2S protocols, like MQTT, offer simplistic QoS services in the form of at-least-once, at-most-once, and exactly-once deliveries.

There are many middleware proposals in the literature that address the interoperation of the IoT protocols. The common design patterns used in these protocols are: publish/subscribe oriented, service oriented [1], virtual machine based [2], and software defined network/radio (SDN/SDR) based designs [3, 4].

MQTT-SN<sup>2</sup>, which is the extension to MQTT for wireless sensor networks, uses the publish/subscribe model in the same way as in the baseline MQTT protocol. A main component of this

## COMMUNICATIONS STANDARDS

*Ala Al-Fuqaha and Mehdi Mohammadi are with Western Michigan University.*

*Abdallah Khreishah is with the New Jersey Institute of Technology.*

*Mohsen Guizani is with the University of Idaho.*

*Ammar Rayes is with Cisco Systems.*

<sup>1</sup> <http://www.omg.org/spec/DDS/1.2/>

<sup>2</sup> <http://mqtt.org/documentation>



architecture is the MQTT-SN gateway. Its main task is to translate MQTT and MQTT-SN messages between MQTT-SN clients and MQTT brokers. One main drawback of this model is that even the simple levels of QoS presented by MQTT are not supported. Things Broker [1] is a service oriented gateway that provides a REST-Ful interface to smart objects by a set of Twitter-based abstractions and communication models. In [2], instead of gateways, the authors propose a solution to integrate smart resource-constrained objects into the Internet using virtual networks. This work can provide end-to-end communication between devices, but scalability and binding to specific protocols are the main challenges. The other method is the combination of the publish/subscribe model and the virtual machine-based design that has been used in the ICSI<sup>3</sup> project to implement a flexible IoT middleware for smart city environments.

Software-defined networking (SDN) and software-defined radio (SDR) have drawn the attention of researchers toward using these brand new technologies to cover different IoT communications. For example, an approach based on SDN was proposed for IoT tasks in [3]. In their approach, the authors developed a middleware with a layered IoT SDN controller to manage dynamic and heterogeneous multi-network environments. As another example, the work in [4] used SDR technology to build a communications infrastructure for IoT applications.

There are some attempts to extend CoAP targeting QoS features and conditional observation. Among them, in [5] the authors proposed an extension to the CoAP protocol to support QoS for timeliness based on delivery priorities. In this work an observer can request a level of priority that establishes the order of notification transmission. In conditional observation, a notification condition can be determined by the clients. An extension to XMPP was presented in [6] in which XMPP was improved by a topic based filter for enhanced performance. Figure 1 illustrates the recent research work toward the IoT gateways as well as enhancements to the current application protocols. The state of our proposed approach is highlighted so that it uses the publish/subscribe pattern in conjunction with DDS (multicasting) and MQTT (intelligent queue management) to provide flexible QoS features and to bridge the gap between the different protocols.

The dissimilar QoS constraints for M2M and M2S/S2S communication and the failure of the current protocols to cover a wide range of QoS requirements is the primary culprit behind protocol fragmentation in the IoT. This gloomy picture of protocol fragmentation and interoperation between IoT devices calls for an intelligent protocol that is capable of meeting the mentioned challenges and an intelligent protocol gateway that is capable of bridging the gap between different technologies. There are other protocols to address the QoS requirements in the Internet (e.g. SNMP), but these protocols are too heavy to be used in resource-constrained devices and are rarely used within the context of IoT. Moreover, the QoS features offered by other application protocols such as XMPP,

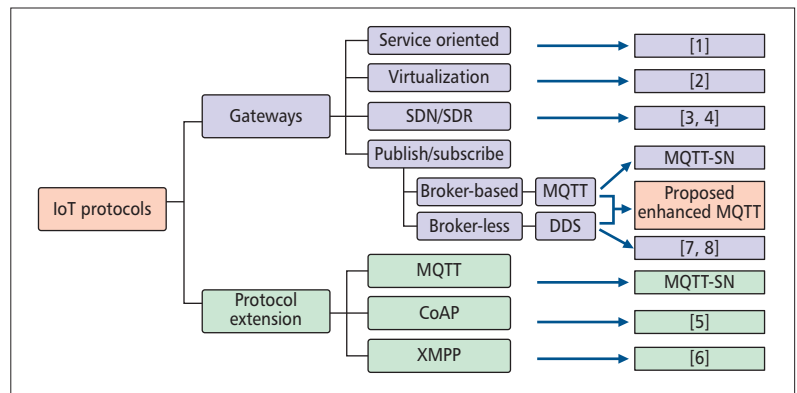


Figure 1. IoT gateway taxonomy, recent related literature and research position.

CoAP, and AMQP are not broader than MQTT. Therefore, MQTT as a light protocol is a good candidate to be extended to address the aforementioned challenges. In addition, the Enhanced MQTT protocol can be used in the fog layer of IoT environments and in collaboration with other gateways to support a full range of QoS requirements.

Although MQTT is considered to be a light-weight protocol, our extension is not going to make it heavy since our extension relies on multicasting to deliver M2M traffic and intelligent queue management for M2S/S2S communications. Therefore, our extension modifies MQTT while being backward compatible with the baseline protocol and without negatively impacting its clean and simple architecture.

In the next section we elaborate on the proposed intelligent IoT protocol gateway. Then in the next section we discuss the motivations and the details of our proposed approach for Enhanced MQTT relative to the current state-of-the-art research in the field.

## INTELLIGENT IOT PROTOCOL GATEWAY

### STATE OF THE ART

IoT applications bring about some horizontal integration situations in which a conversion between protocols is required. These situations usually happen when smart devices are talking to each other using different communication technologies (like Bluetooth and ZigBee) or when they are using different application protocols (like MQTT and CoAP). Several attempts have been made in the recent literature to address this issue. Paramount among these attempts is Ponte, which was initially developed as QEST and currently is under the Eclipse IoT project.<sup>4</sup> Ponte offers uniform open APIs for the programmer to enable the automatic conversion between the various IoT application protocols such as CoAP and MQTT.

In addition to Ponte, the Eclipse IoT Project includes other sub-projects that facilitate IoT application development such as Kura, Eclipse SCADA, Eclipse SmartHome, and Krikkit. Kura provides a Java/OSGi-based container for M2M applications running in service gateways and facilitates I/O access, data services, watchdog, network configuration, and remote management for M2M applications. Eclipse SCADA's focus is

<sup>3</sup> <http://ict-icsi.eu/>

<sup>4</sup> <https://projects.eclipse.org/projects/iot>

Gateways		Application protocols					Management		Connectivity			
		RESTful HTTP	COAP	MQTT	XMPP	DDS	OMA-DM	BBFTR-069	Cellular	Zigbee	Bluetooth	WiFi
Practice	Ponte	✓	✓	✓								
	oneM2M	✓	✓	✓			✓	✓	✓	✓	✓	✓
	SmartM2M	✓	✓				✓		✓	✓	✓	✓
	Intel IoT Platform			✓			✓	✓	✓	✓	✓	✓
	LWM2M								✓	✓		✓
Research	IoT communication gateway [9]								✓	✓		
	IoT gateway centric model [10]	✓								✓		✓
	HTTP-CoAP crossprotocol proxy [11]	✓	✓									
	Semantic gateway [12]		✓	✓	✓							
	Proposed enhanced MQTT	✓	✓	✓	✓	✓	✓		✓	✓	✓	✓

Table 1. IoT gateways and their supported protocols.

to provide a common communication mechanism as well as post-processing and visualization of the data. Eclipse SmartHome aims to build a heterogeneous environment to integrate different protocols and standards and to bring a uniform access while supporting different kinds of interactions. The Krikkit architecture uses the publish/subscribe model by which data acquiring is possible by registering rules or policies on edge devices like sensor gateways.

An ongoing standard specification, oneM2M<sup>5</sup>, aims to provide a common M2M/IoT service platform that supports secure, reliable, and efficient operations of M2M/IoT services. oneM2M tries to leverage existing specifications and bridge the gap between the current M2M service layer protocols. oneM2M utilizes REST to represent and manage smart things while using the common application protocols including HTTP, CoAP, OMA-DM, and MQTT to facilitate interworking with non-oneM2M systems. For device management, oneM2M relies on the OMA-DM and BBF TR-069 specifications. In order to have a heterogeneous environment for different access technologies like ZigBee, BACnet, or Bluetooth, oneM2M has started to provide device abstraction and semantic interworking.

The SmartM2M standard is developed by ETSI to provide a horizontal M2M service platform for developing network-independent services and deploying vertical applications.

Intel also announced the launch of Intel IoT<sup>6</sup> middleware, connectivity, and security components to simplify deploying IoT applications. The Intel IoT gateway enables connectivity up to the cloud and down to the smart devices. In terms of connectivity and communications, Intel IoT Platform supports ZigBee, Cellular 2G/3G/4G, Bluetooth, Serial, USB, VPN, Wi-Fi, and MQTT. Its device management module consists of OMA-

DM, BBF TR-069, and web-based configuration interfaces.

In addition to the aforementioned projects, academic research projects and proposals have presented partial solutions or addressed a specific application to the problem. They may also be limited to specific hardware devices or rely on resource-rich TCP/IP based devices. For example, the gateway that is proposed in [9] supports the conversion between ZigBee and GPRS protocols in wireless sensor networks and mobile communication networks. The Light-weight M2M<sup>7</sup> (LWM2M) device management protocol has been developed to ease unified remote device management. This protocol is also limited to those devices that support IP.

In [10] the authors proposed a wireless gateway by which mobile clients are able to interact with M2M devices. Their proposed gateway just supports REST queries in the application layer and provides dynamic device discovery, connection management of non-smart things, and binding metadata to sensor and actuator measurements. Other conversions like HTTP to CoAP [11] are also presented in the literature. In [12] the authors proposed a gateway and Semantic Web enabled IoT architecture to bring the interoperability of messaging protocols such as XMPP, CoAP, and MQTT using a multi-protocol proxy. Table 1 compares recent implementations of IoT gateways and their coverage of different IoT protocols.

The anecdotal data that we have obtained so far about the diverse needs of IoT applications and the capabilities and restrictions of the underlying hardware reveals that the existing IoT gateway solutions like Ponte or Intel IoT Platform are not sufficiently extensive to bridge the gap between the different IoT protocols, and therefore a more intelligent and comprehensive solu-

<sup>5</sup> <http://www.onem2m.org>

<sup>6</sup> <http://www.intel.com/content/www/us/en/internet.of-things>

<sup>7</sup> <http://technical.openmobilealliance.org/Technical/technical-information/omna/lightweight-m2m-lwm2mobject-registry>

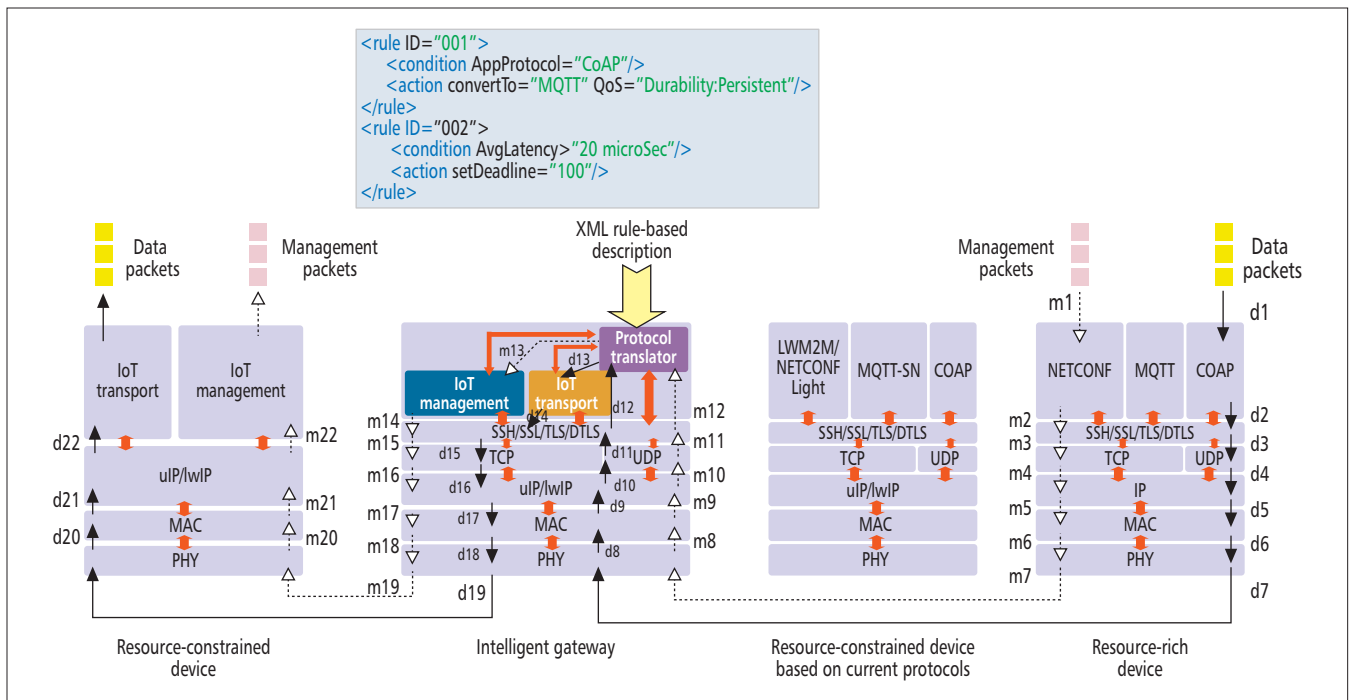


Figure 2. The architecture of a rule-based Intelligent IoT gateway.

tion is needed. For example, Ponte’s feature that automatically converts between a number of protocols comes at a price due to longer packets and verbose communication. Moreover, as with many other protocol gateways, Ponte assumes that the underlying devices rely on TCP/IP. Intel IoT Platform, although supporting a wide range of transmission protocols, currently does not provide full conversion capabilities between application protocols.

### FUNCTIONALITIES OF THE INTELLIGENT GATEWAY

Approaches like that followed by Ponte assist programmers to write a single version of the same application while supporting different protocols, but the programmers have no control over the underlying wire-protocol, and consequently the degradation of performance and efficiency will be uncontrolled. A more important drawback is that resource-constrained devices have no chance to be integrated with other objects in this solution, bringing to mind that they are considered to be “second-class citizens.” Therefore, we are motivated by the following three main observations to illustrate the need for a new intelligent IoT gateway.

- IoT gateways should include mechanisms by which programmers are allowed to control the wire protocol. Application-agnostic message passing for resource-constrained devices causes unnecessary packet exchanges. In such environments, an intelligent gateway provides facilities for programmers to optimize the performance by controlling the underlying wire protocol based on the specific application.

- Resource-constrained devices should be treated the same as resource-rich devices. An intelligent gateway should support true interoperability between these two types of IoT devices.

- The introduction of an intelligent protocol

gateway into the IoT has the potential to open up new opportunities to reduce the market fragmentation between IoT protocols.

The aforementioned observations provide strong evidence that an intelligent IoT gateway is essential to IoT development so that it offers “smart” services and it is deeply re-programmable through a high-level rule-based language written by the programmer. It is remarkable that our proposal for deeper re-programmability of the IoT gateway through a rule-based language not only does not conflict with the current interoperability and management standardization efforts, but also complements the existing protocols like IEEE 1905.1, LWM2M, and NETCONF Light. While the intelligent gateway can utilize default templates of rules out-of-the-box for convenience, the deep rule-based programmability of the gateway allows application developers to exploit application-specific communication patterns to achieve more efficient protocol translations. Therefore, application developers have to write the application protocol translation rules once in return for more efficient protocol translations.

The overall concept of the proposed gateway is demonstrated in Fig. 2. Here, the protocol stack that is installed on resource-constrained devices relies on the current technologies pertaining to resource-constrained devices that utilize the intelligent gateway. The flow sequence of data packets and the management packets are depicted in the figure by labels  $d_i$  and  $m_i$ , respectively. In the gateway, the logic of the rules is applied to the stream of data and management packets (i.e.  $d_{12}$  and  $m_{12}$  in Figure 2) leading to generate appropriate data and management packets (i.e.  $d_{13}$  and  $m_{13}$  in Figure 2). Autonomic management and data aggregation services are also possible by feeding rules to the

Research	Pub/sub	QoS-enabled	Application domain
QoS-aware platform [7]	✓	✓	WSN (healthcare)
Sensor network pub/sub protocol [8]	✓	✗	WSN
Decentralized pub/sub model [13]	✓	✓	Peer-to-peer
Real-time pub/sub service [14]	✓	✓	Real-time applications
DDS/SDN based communication [15]	✓	✗	Real-time applications

**Table 2.** Summary of recent research utilizing the DDS protocol.

gateway while the gateway can generate and transmit new data and management packets from the received ones. Inside the gateway, encapsulation and decapsulation of data packets are carried out by a lightweight protocol named “IoT transport,” while “IoT Management,” another lightweight protocol, does the same for management packets.

The light protocol stack presented by the intelligent gateway relies only on uIP/lwIP protocols without the need for other transportation and security protocols on the resource-constrained device. Security services can also be delegated to the gateway so that confidentiality, authenticity, and integrity traffic can be exchanged with the gateway and not directly with resource-constrained devices.

By supporting rule-based language, the programmer will have to describe the transformations in a high-level language by which the gateway should perform the required conversion. The advantage of this approach comes from the flexibility of the gateway by which a programmer can examine and fulfill wire protocol optimizations at the cost of describing required transformations in a high-level domain-specific language. The programmer can also use standard transformation templates instead of specifying application-specific rules. In this case, the wire protocol messaging will be less efficient but yet will be similar to that of Ponte.

The proposed intelligent gateway can also be used for localized autonomic management of the IoT devices automatically. Considering a real deployment scenario consisting of thousands or millions of IoT nodes, it is vital to have self-management fault, configuration, accounting, performance, and security (FCAPS) capabilities. Moreover, a deeply re-programmable gateway can also provide the option to accomplish data and flow aggregation and consequently reduce the number of flows handled by the network elements and to achieve high performance. As another benefit, the existence of several gateways within the IoT can also be utilized for load balancing of the potentially massive IoT traffic among the available gateways.

## ENHANCED MQTT

As discussed above, DDS is a quintessential IoT protocol that employs a broker-less publish/subscribe architecture to meet the stringent real-time constraints for M2M communications.

Nevertheless, it is only utilized in academic labs and, to the best of our knowledge, has not been used in IoT released products. In the following subsections we review recent implementations of DDS, followed by our proposal for the Enhanced MQTT and its motivations and opportunities.

## DDS IN THE LITERATURE

In recent years DDS has drawn the attention of many IoT researchers to design and develop more reliable and real-time IoT systems. Most of the studies focus on the QoS features of the DDS protocol [7, 13, 14]. In [13] a decentralized content-based publish/subscribe model based on DDS is introduced in which a peer-to-peer communication is guaranteed to have individual QoS requirements based on the bandwidth constraints. Exploiting the DDS protocol for publish/subscribe communication to support various QoS levels for real-time applications has been proposed in [14] as well. The Sensor Network Publish-Subscribe protocol (SNPS) [8] has been proposed to use the DDS concepts in WSNs, but it does not have the complete QoS features offered by DDS. Achieving QoS and reliability requirements in WSNs for healthcare environments has been investigated in [7] in which DDS has been used as the main building block. The feasibility of managing DDS communications dynamically using SDN is described in [15].

The common drawback of these research activities is that they do not support all types of IoT communications (i.e. M2M, M2S, and S2S). Specifically, they assume a real-time environment for their system, while real-time response is not necessary everywhere. Moreover, the current proposals do not provide sufficient QoS support. The characteristics of these attempts that are leveraging DDS in their core have been summarized in Table 2.

Broker-based architectures like MQTT and AMQP can offload devices from having to handle a large number of upstream server requests. However, this relief comes with the cost of losing quality of service and reliability, which are vital in most IoT applications.

## ENHANCED MQTT: THE BIG PICTURE

In this research we propose to revisit the MQTT protocol with the aim of providing it with a hybrid architecture that would allow it to operate in various M2M as well as M2S and S2S scenarios. Therefore, we provide a solution that stems the root cause of the protocol fragmentation issue encountered within the context of the IoT. A horizontal integrated IoT framework needs to handle a wide range of different communication and application protocols and to facilitate their interworking. To realize such a framework, we target the MQTT protocol by improving it to support rich QoS features, intelligent broker queue management, and traffic analytics techniques. The hybrid architecture will allow the protocol to seamlessly utilize direct and broker-less multicast communication for M2M communication while utilizing the broker for M2S and S2S communications. This hybrid architecture would allow MQTT to extend its role in the IoT to handle M2M communications and would allow for multiple MQTT brokers to

cooperate to deliver better QoS and reliability capabilities. Beyond this architectural change to MQTT, we believe that the QoS features currently offered by MQTT are very limited and there is a need to extend these features significantly to support various development and deployment scenarios. Also, MQTT does not offer a management interface that allows for the prioritization, preemption, and collection of analytics on the IoT traffic going through the broker. A full set of traffic analysis services is to be supported by the broker. The proposed architecture considers a *subscribers table* and a *management table* inside each broker. Brokers are registered with the desired topics in the subscribers table. The management table tracks the statistics of each topic in terms of QoS and reliability metrics. Through these extensions, our proposed Enhanced MQTT remains backward compatible with MQTT.

Figure 3 depicts the proposed new MQTT capabilities. These capabilities can be summarized as follows.

- Allow for broker-less multicast communication in support of M2M communication (cf. interaction d on Fig. 3). The broker-less multicast communication will allow MQTT to be used for M2M communication as well and to enhance its delay and fan-out performance.

- Allow multiple brokers to receive multicast communication from IoT devices in support of reliability (cf. interaction d on Fig. 3). The multicast communication will allow multiple MQTT brokers to listen to the device generated traffic allowing for failure recovery in the event of broker failures.

- Allow brokers to move subscribers to other brokers in support of QoS (cf. interaction c on Fig. 3). This will allow for QoS features that are beyond the naïve QoS features currently offered by MQTT.

- Allow brokers to re-prioritize (i.e. preempt) the distribution of MQTT packets in support of QoS (cf. interaction a on Fig. 3). This will enable the inclusion of rich QoS features even in deployments that do not involve multiple brokers.

- Allow subscribers to register for and obtain analytics on the IoT traffic going through the broker (cf. interaction d on Fig. 3). By far this is the most important feature that is enabled by our proposed approach, as it would allow for a new class of IoT applications that are otherwise infeasible using current protocols such as MQTT. The new class of applications enabled by our proposed protocol allows the IoT infrastructure itself to evolve based on collected analytics. For example, the collected analytics can be used as the input to a mathematical formulation or a reinforcement learning strategy that would reconfigure the number of brokers and their locations to obtain better QoS support for the given application. Therefore, the application designers do not have to worry about the optimal deployment for the given application offline; instead they can have their application monitor its own performance and consequently evolve its configuration.

We have performed a simulation experiment to compare the performance of our proposed model to the baseline MQTT model in terms of queuing delay. We compared the queue time

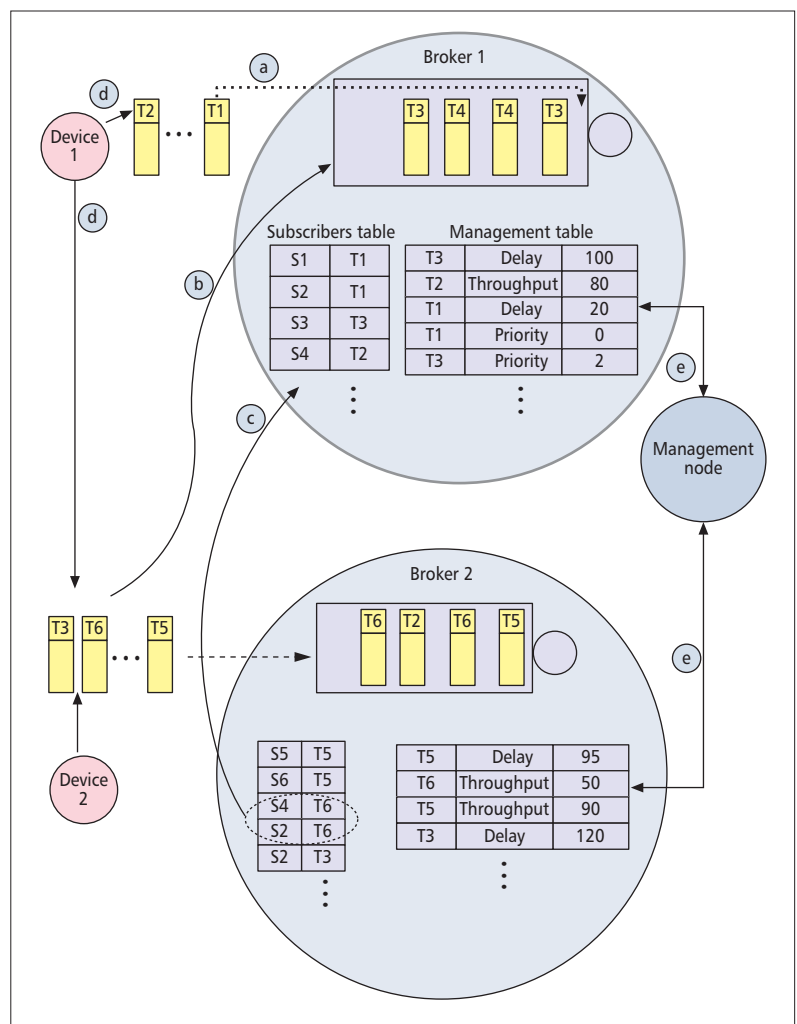


Figure 3. Enhancements of broker functions.

only since it is a bottleneck in broker-based approaches, and our proposal is an attempt to resolve this issue while leaving all the other parameters intact (i.e. processing, transmission, and propagation delays). Figure 4 illustrates the results of this evaluation. In our model, M2M and M2S traffic streams are separated to be handled differently. M2M traffic is directly delivered to the destination using multicasting, while M2S traffic goes through a queuing system with four different levels of priorities. Figure 4 reveals that the baseline MQTT has resulted in more queuing delays compared to our proposed model, even though the service rate was increased by 5 percent and 10 percent, respectively, compared to the service rate used in our proposed model. On the other hand, the proposed model shows a decrease in the queuing time when the fraction of M2M traffic constitutes a higher percentage of the overall IoT traffic.

## MOTIVATION AND SCENARIOS

The proposed QoS and reliability capabilities of MQTT can be utilized in a variety of IoT applications. As a matter of fact, our proposal was initially motivated by an actual need that appeared in an interdisciplinary study that we are involved with to alleviate some of the feed-

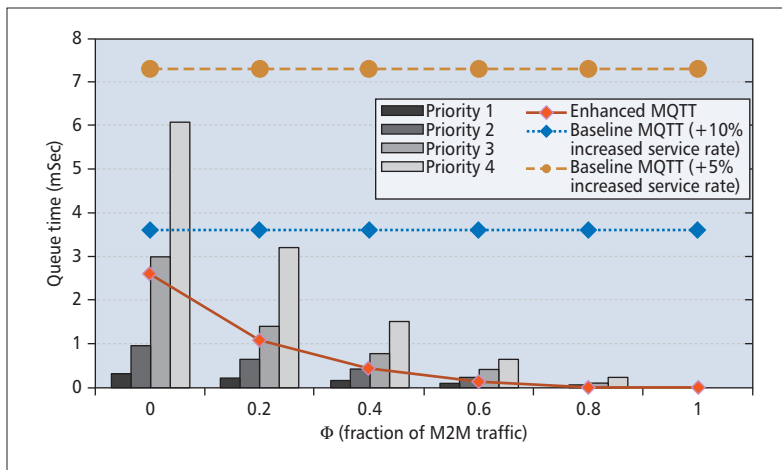


Figure 4. The Enhanced MQTT protocol compared to the baseline MQTT protocol.

ing difficulties experienced by adults with essential tremors, Parkinson's disease, and Alzheimer's disease. As an example of M2M communication with stringent delay requirements, in order to allow adults with essential tremors or Parkinson's disease to eat without spilling food, a glove can be equipped with tiny vibrating motors to counteract the hand movement instability measured by the accelerometers. In this application, the accelerometer sensors and vibrating motors have to communicate with the minimum delay possible to deliver the required functionality.

The monitoring of the patients' vital signs in a nursing home environment provides a quintessential scenario that requires brokers to deliver collected M2M data efficiently and reliably to different servers (i.e. M2S communication with stringent delay requirements). In this scenario, the vital sign measurements must be delivered to multiple nursing stations with minimum delay for accurate visualization and actuation of smart infusion pumps and other bed-size medical equipment even in the presence of broker failures. Another interesting scenario within a nursing home environment includes the use of accelerometers to monitor the eating behavior of Alzheimer's patients and activating audio-visual cues recorded by a close relative through a wearable device. These cues aim to stimulate the patient to eat better, exploiting the fact that most Alzheimer's patients continue to remember close relatives during the different stages of the disease.

Beyond a nursing home environment, a scenario that can benefit from the proposed enhancements includes exchanging time difference of arrival (TDoA) measurements between embedded devices deployed within the infrastructure, and wearable devices to provide precise localization details to blind and visually impaired people to navigate an indoor environment. In this scenario, the wearable device can also communicate with minimum delay to obtain the best estimated location and provide tactile cues to the patient to navigate to their destination through querying the meta-data associated with the location fingerprints.

In this article the restrictions and shortages of the current application protocols involved in IoT systems have been elaborated. We then identified the major driving forces for progression toward an intelligent IoT gateway to overcome the fragmentation and enable the efficient integration of horizontal IoT services. The opportunities that will be made possible by this evolution have been discussed as well.

To offer a solution that stems the root cause of this protocol fragmentation issue, an enhanced version of the MQTT protocol is proposed in this article that alleviates the deficiency of the existing MQTT protocol, especially in support of QoS and reliability. The enhanced protocol is strengthened by an intelligent queuing mechanism and real-time multicasting so that it is able to handle the different forms of communications required by IoT applications including M2M, M2S, and S2S. The performance of the current work is promising to be better than existing published results. Furthermore, the proposed protocol will contribute to minimizing the myriad of protocol options that application developers have to evaluate before developing IoT services and applications.

There are many potential future evaluation studies of the proposed Enhanced MQTT. In our future work we plan to evaluate our proposed extension in a wide range of deployment, emulation, and simulation scenarios. The application of the Enhanced MQTT model to implement real-time IoT systems is another potential direction for this research. Investigating other IoT challenges like security, scalability, availability, and management are also in our future research plan.

## REFERENCES

- [1] R. A. Perez de Almeida *et al.*, "Thing Broker: A Twitter for Things," *Proc. 2013 ACM Conf. Pervasive and Ubiquitous Computing Adjunct Publication*, Zurich, Switzerland, 2013, pp. 1545-54.
- [2] I. Ishaq *et al.*, "Internet of Things Virtual Networks: Bringing Network Virtualization to Resource-Constrained Devices," *2012 IEEE Int'l. Conf. Green Computing and Commun. (GreenCom)*, 2012, pp. 293-300.
- [3] Z. Qin *et al.*, "A Software Defined Networking Architecture for the Internet-of-Things," *2014 IEEE Network Operations and Management Symp. (NOMS)*, 2014, pp. 1-9.
- [4] Y. H. Lin *et al.*, "Wireless IoT Platform Based on SDR Technology," *2013 IEEE and Internet of Things (IThings/CPSCoM) Green Computing and Commun. (GreenCom), IEEE Int'l. Conf. and IEEE Cyber, Physical and Social Computing*, 2013, pp. 2245-46.
- [5] A. Ludovici *et al.*, "Adding QoS Support for Timeliness to the Observe Extension of CoAP," *2012 IEEE 8th Int'l. Conf. Wireless and Mobile Computing, Networking and Communications (WiMob)*, 2012, pp. 195-202.
- [6] R. Klauk and M. Kirsche, "Chatty Things Making the Internet of Things Readily Usable for the Masses with XMPP," *2012 8th Int'l. Conf. Collaborative Computing: Networking, Applications and Worksharing (CollaborateCom)*, 2012, pp. 60-69.
- [7] A. Agirre *et al.*, "QoS Aware Platform for Dependable Sensory Environments," *2014 IEEE Int'l. Conf. Multimedia and Expo Wksp. (ICMEW)*, 2014, pp. 1-5.
- [8] K. Beckmann and M. Thoss, "A Wireless Sensor Network Protocol for the OMG Data Distribution Service," *2012 Proc. 10th Wksp. Intelligent Solutions in Embedded Systems (WISES)*, 2012, pp. 45-50.
- [9] Q. Zhu *et al.*, "IoT Gateway: Bridging Wireless Sensor Networks into Internet of Things," *2010 IEEE/IFIP 8th Int'l. Conf. on Embedded and Ubiquitous Computing (EUC)*, 2010, pp. 347-52.
- [10] S. K. Datta, C. Bonnet, and N. Nikaen, "An IoT Gateway Centric Architecture to Provide Novel M2M Services," *2014 IEEE World Forum on Internet of Things (WF-IoT)*, 2014, pp. 514-19.
- [11] A. Castellani, T. Fossati, and S. Loreto, "HTTP-CoAP Cross Protocol Proxy: An Implementation Viewpoint," *2012 IEEE 9th Int'l. Conf. Mobile Adhoc and Sensor Systems (MASS)*, 2012, pp. 1-6.

- [12] P. Desai, A. P. Sheth, and P. Anantharam, "Semantic Gateway as a Service Architecture for IoT Interoperability," *CoRR*, vol. abs/1410.4977, 2014.
- [13] M. A. Tariq *et al.*, "Meeting Subscriber-Defined QoS Constraints in Publish/Subscribe Systems," *Concurrency and Computation: Practice and Experience*, vol. 23, 2011, pp. 2140–53.
- [14] X. Lu *et al.*, "A Novel QoS-Enabled Real-Time Publish-Subscribe Service," *Int'l. Symp. Parallel and Distributed Processing with Applications*, 2008, *ISPA '08*, pp. 19–26.
- [15] L. Bertaux *et al.*, "A DDS/SDN Based Communication System for Efficient Support of Dynamic Distributed Real-Time Applications," *2014 IEEE/ACM 18th Int'l. Symp. Distributed Simulation and Real Time Applications (DS-RT)*, 2014, pp. 77–84.

## BIOGRAPHIES

ALA AL-FUQAHA [S'00, M'04, SM'09] (Ala.Al-Fuqaha@wmich.edu) is an associate professor and director of the NEST Research Lab in the Computer Science Department at Western Michigan University. His research interests fall in the areas of vehicular networks, dynamic spectrum access etiquettes in cognitive radio networks, smart services in support of the Internet of Things, and planning of software defined networks (SDN). He is currently serving on the editorial board of multiple journals.

ABDALLAH KHREISHAH [M] (abdallah@njit.edu) is an assistant professor in the Department of ECE at the New Jersey Institute of Technology. His research interests are in the areas of visible-light communication, green networking, wireless networks, and network security. He received his Ph.D. degree in electrical and computer engineering from Purdue University in 2010. While pursuing

his Ph.D. studies, he worked with NEESCOM. He is the chair of the North Jersey IEEE EMBS chapter.

MOHSEN GUIZANI [S'85, M'89, SM'99, F'09] (mguizani@ieee.org) is a professor and chair of the ECE Department at the University of Idaho. His research interests include wireless communications and mobile computing, computer networks, cloud computing, cyber security, and smart grid. He currently serves on the editorial boards of several international technical journals. He is the author of nine books and more than 400 publications. He is a senior member of ACM.

AMMAR RAYES (rayes@cisco.com), Ph.D., is a distinguished engineer at the CTO Office of Cisco Services Technology Group focusing on IoT, analytics and smart services. He is the Founding President of the International Society of Service Innovation Professionals. Ammar has authored more than 100 papers and patents on advances in telecommunications-related technologies. He is the editor-in-chief of the journal *Advances of Internet of Things*, and he has served as an associate editor of *ACM Transactions on Internet Technology* and the *Journal of Wireless Communications and Mobile Computing*.

MEHDI MOHAMMADI [S'14] (Mehdi.Mohammadi@wmich.edu) is a Ph.D. student in the Department of Computer Science, Western Michigan University (WMU), Kalamazoo, MI, USA. His research interests include Internet of Things, Future Internet, software-defined networking, data mining, and natural language processing. He served as a reviewer for the Wiley journal *Security and Wireless Communication Networks*. He is a recipient of the Graduate Doctoral Assistantship from the WMU Libraries, as well as four travel grants from the National Science Foundation.

**CALL FOR PAPERS**  
**IEEE COMMUNICATIONS MAGAZINE**  
**WIRELESS COMMUNICATIONS, NETWORKING, AND POSITIONING**  
**WITH UNMANNED AERIAL VEHICLES**

**BACKGROUND**

Enabled by the advances in computing, communication, and sensing as well as the miniaturization of devices, unmanned aerial vehicles (UAVs) such as balloons, quadcopters, and gliders have been receiving significant attention in the research community. Indeed, UAVs have become an integral component in several critical applications such as border surveillance, disaster monitoring, traffic monitoring, remote sensing, and the transportation of goods, medicine, and first-aid. More recently, new possibilities for commercial applications and public service for UAVs have begun to emerge, with the potential to dramatically change the way in which we lead our daily lives. For instance, in 2013, Amazon announced a research and development initiative focused on its next-generation Prime Air delivery service. The goal of this service is to deliver packages into customers' hands in 30 minutes or less using small UAVs, each with a payload of several pounds. 2014 has been a pivotal year that has witnessed an unprecedented proliferation of personal drones, such as the Phantom and Inspire from DJI, AR Drone and Bebop Drone from Parrot, and IRIS Drone from 3D Robotics.

Among the many technical challenges accompanying the aforementioned applications, leveraging the use of UAVs for delivering broadband connectivity plays a central role in next generation communication systems. Facebook and Google announced in 2014 that they will use a network of drones which circle in the stratosphere over specific population centers to deliver broadband connectivity. Such solar-powered drones are capable of flying several years without refueling. UAVs have also been proposed as an effective solution for delivering broadband data rates in emergency situations through low-altitude platforms. For example, the ABSOLUTE, ANCHORS, and AVIGLE projects in Europe have been investigating the use of aerial base stations to establish opportunistic links and ad-hoc radio coverage during unexpected and temporary events. They can serve as a temporary, dynamic, and agile infrastructure for enabling broadband communications, and quickly localizing victims in case of disaster scenarios.

This proposed Feature Topic (FT) issue will gather articles from a wide range of perspectives in different industrial and research communities. The primary FT goals are to advance the understanding of the challenges faced in UAV communications, networking, and positioning over the next decade, and provide further awareness in the communications and networking communities on these challenges, thus fostering future research. Original research papers are to be solicited in topics including, but not limited to, the following themes on communications, networking, and positioning with UAVs.

- Existing and future communication architectures and technologies for small UAVs
- Delay-tolerant networking for cooperative UAV operations
- Design and evaluation of wireless UAV test beds, prototypes, and platforms
- Multi-hop and device-to-device communications with UAVs
- Interfaces and cross-platform communication for UAVs
- QoS mechanisms and performance evaluation for UAV networks
- Game-theoretic and control-theoretic mechanisms for UAV communications
- Use of civilian networks for small UAV communications
- Integrating 4G and 5G wireless technologies into UAV communications, such as millimeter wave communications, beamforming, moving networks, and machine type communications
- Use of UAVs for public safety and emergency communications, networking, and positioning
- Integration of software defined radio and cognitive radio techniques with UAVs
- Channel propagation measurements and modeling for UAV communication channels

**SUBMISSIONS**

Articles should be tutorial in nature, with the intended audience being all members of the communications technology community. They should be written in a style comprehensible to readers outside the specialty of the article. Mathematical equations should not be used (in justified cases up to three simple equations are allowed). Articles should not exceed 4500 words (from introduction through conclusions). Figures and tables should be limited to a combined total of six. The number of references is recommended not to exceed 15. In some rare cases, more mathematical equations, figures, and tables may be allowed if well-justified. In general, however, mathematics should be avoided; instead, references to papers containing the relevant mathematics should be provided. Complete guidelines for preparation of the manuscripts are posted at <http://www.comsoc.org/commag/paper-submission-guidelines>. Please send a pdf (preferred) or MSWORD formatted paper via Manuscript Central (<http://mc.manuscriptcentral.com/commag-ieee>). Register or log in, and go to Author Center. Follow the instructions there. Select "May 2016 / Wireless Communications, Networking and Positioning with UAVs" as the Feature Topic category for your submission.

**SCHEDULE FOR SUBMISSIONS**

- Submission Deadline: November 1, 2015
- Notification Due Date: January 15, 2016
- Final Version Due Date: March 1, 2016
- Feature Topic Publication Date: May 2016

**GUEST EDITORS**

Ismail Guvenc  
Florida International Univ., USA  
iguvenc@fiu.edu

Walid Saad  
Virginia Tech, USA  
walids@vt.edu

Mehdi Bennis  
Univ. of Oulu, Finland  
bennis@ee.oulu.fi

Christian Wietfeld  
TU Dortmund Univ., Germany  
christian.wietfeld@tu-dortmund.de

Ming Ding  
NICTA, Australia  
ming.ding@nicta.com.au

Lee Pike  
Galois, Inc., USA  
leepike@galois.com





## IEEE ICC 2016 CALL FOR PAPERS AND PROPOSALS

The 2016 IEEE International Conference on Communications (ICC) will be held from 23-27 May 2016 at Kuala Lumpur Convention Center, Malaysia, conveniently located in the middle of Southeast Asia, the region home to many of the world's largest ICT industries and research labs. Themed "Communications for All Things," this flagship conference of IEEE Communications Society will feature a comprehensive Technical Program including 13 Symposia and a number of Tutorials and Workshops. IEEE ICC 2016 will also include an attractive Industry Forum & Exhibition Program featuring keynote speakers, business and industry panels, and vendor exhibits.

### TECHNICAL SYMPOSIA

We invite you to submit original technical papers in the following areas:

#### Symposium on Selected Areas in Communications

##### - Access Systems and Networks

Ahmed E. Kamal, Iowa State University, USA

##### - Cloud Communications and Networking

Dzmitry Kliazovich, University of Luxembourg, Luxembourg

##### - Communications for the Smart Grid

Lutz Lampe, University of British Columbia, Canada

##### - Data Storage

Edward Au, Huawei Technologies, Canada

##### - E-Health

Joel Rodrigues, University of Beira Interior, Portugal

##### - Internet of Things

Antonio Skarmeta, University of Murcia, Spain

##### - Satellite and Space Communications

Song Guo, University of Aizu, Japan

##### - Social Networking

Pan Hui, HKUST, Hong Kong

#### Ad-Hoc and Sensor Networks

Abdelhakim Hafid, University of Montreal, Canada  
Cheng Li, Memorial University of Newfoundland, Canada  
Pascal Lorenz, University of Haute-Alsace, France

#### Communication and Information System Security

Kejie Lu, University of Puerto Rico, Mayaguez, Puerto Rico  
Yu Cheng, Illinois Institute of Technology, USA

#### Communications QoS, Reliability and Modelling

Kohei Shiimoto, NTT, Japan  
Christos Verikoukis, CTTC, Spain  
Charalabos Skianis, Aegean University, Greece

#### Cognitive Radio and Networks

Norman C. Beaulieu, BUPT, China  
Linyang Song, Peking University, China

#### Communications Software, Services and Multimedia Applications

Shingo Ata, Osaka City University, Japan  
Fen Hou, University of Macau, China

#### Communication Theory

Marios Kountouris, Supelec, France  
Marco Chiani, University of Bologna, Italy  
Xu (Judy) Zhu, University of Liverpool, UK

#### Green Communications Systems and Networks

Sumei Sun, Institute for Infocomm Research, Singapore  
Anura Jayasumana, Colorado State University, USA

#### Mobile and Wireless Networks

Adlen Ksentini, University of Rennes, France  
Mohammed Atiqzaman, University of Oklahoma, USA  
Jalel Ben-Othman, University of Paris 13, France

#### Next Generation Networking and Internet

Rami Langar, University of Paris 6, France  
Shiwen Mao, Auburn University, USA  
Abdelhamid Mellouk, University of Paris-Est, France

#### Optical Networks and Systems

Walter Cerroni, University of Bologna, Italy  
Krishna Sivalingam, IIT Madras, India

#### Signal Processing for Communications

Hsiao-Chun Wu, Louisiana State University, USA  
Shaodan Ma, University of Macau, China  
Tomohiko Taniguchi, Fujitsu Labs, Japan

#### Wireless Communications

Xiaohu Ge, Huazong University of Science and Technology, China  
Dimitrie Popescu, Old Dominion University, USA  
Hossam Hassanein, Queen's University, Canada  
Rui Zhang, National University of Singapore

### INDUSTRIAL FORUM AND EXHIBITION PROGRAM

IEEE ICC 2016 will feature several prominent keynote speakers, major business and technology forums, and a large number of vendor exhibits. Submit your proposals to the IF&E Chair.  
Khaled B. Letaief (eekhaled@ee.ust.hk)

### TUTORIALS

Proposals are invited for half- or full-day tutorials in all communication and networking topics. For enquiries, please contact Tutorial Program Co-Chairs.  
Mike Devetsikiotis (mdevets@ncsu.edu)  
Koichi Asatani (asatani@ieee.org)

### WORKSHOPS

Proposals are invited for half- or full-day workshops in all communication and networking topics. For enquiries, please contact Workshop Program Co-Chairs.  
Tarek El-Bawab (telbawab@ieee.org)  
Fabrizio Granelli (granelli@disi.unitn.it)

### ORGANIZING COMMITTEE

#### General Chair

**Dato' Sri Jamaludin Ibrahim**  
CEO, Axiata Group, Malaysia

#### Executive Co-Chairs

**Hikmet Sari**  
Supelec, France  
**Borhanuddin Mohd Ali**  
Universiti Putra, Malaysia

#### Technical Program Co-Chairs

**Stefano Bregni**  
Politecnico di Milano, Italy  
**Nelson Fonseca**  
State University of Campinas, Brazil

#### Technical Program Vice-Chair

**Jiang Linda Xie**  
University of North Carolina,  
Charlotte, USA

#### Industry Forums & Exhibition Chair

**Khaled B. Letaief**  
Hong Kong University of Science  
and Technology, Hong Kong

#### Tutorial Program Co-Chairs

**Mike Devetsikiotis**  
North Carolina State University, USA  
**Koichi Asatani**  
Kogakuin University, Japan

#### Workshop Program Co-Chairs

**Tarek El-Bawab**  
Jackson State University, USA  
**Fabrizio Granelli**  
University of Trento, Italy

#### Conference Operations Chair

**Hafizal Mohamad**  
MIMOS Berhad, Malaysia

#### Advisory Executive Vice-Chair

**Datuk Hod Parman**  
Past Communication Commission  
General Director, Malaysia

#### Exhibition Chair

**Nordin Ramli**  
MIMOS Berhad, Malaysia

### IMPORTANT DATES

Paper Submissions:  
**16 October 2015**

Tutorial Proposals:  
**13 November 2015**

IF&E Proposals:  
**13 November 2015**

Workshop Proposals:  
**17 July 2015**

Paper Acceptance Notification:  
**29 January 2016**

Camera-Ready Papers:  
**29 February 2016**

# MILCOM2015

LEVERAGING TECHNOLOGY – THE JOINT IMPERATIVE

OCTOBER 26–28, 2015 • TAMPA, FLORIDA

Register now and join your colleagues in government, military, academia and industry at MILCOM 2015! Experience an in-depth technical program with paper presentations, panel discussions, tutorials, and technology exhibits at the state-of-the-art Tampa Convention Center. *New this year*, one tutorial included with each paid conference registration! For more information – including registration details, technical program outline, and schedule of events – visit [www.milcom.org](http://www.milcom.org).

COHOSTED BY  
AFCEA AND IEEE COMMUNICATIONS SOCIETY

