

IEEE Communications

www.comsoc.org

MAGAZINE

protect information

secu

- *Security and Privacy in Emerging Networks*
- *Energy Harvesting Communications*
- *Energy-Efficient Optical Networks*



protect information



IEEE



IEEE COMMUNICATIONS SOCIETY

A Publication of the IEEE Communications Society

THANKS OUR CORPORATE SUPPORTERS



IEEE Communications

www.comsoc.org

MAGAZINE

protect information

secu

- *Security and Privacy in Emerging Networks*
- *Energy Harvesting Communications*
- *Energy-Efficient Optical Networks*



protect information



IEEE



IEEE COMMUNICATIONS SOCIETY

A Publication of the IEEE Communications Society

PAM-4 insights don't schedule meetings.

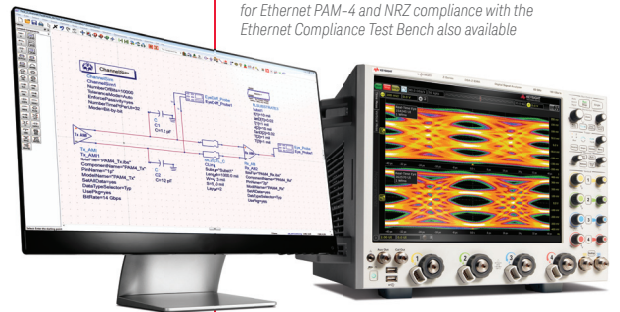
They come when they're good and ready.

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



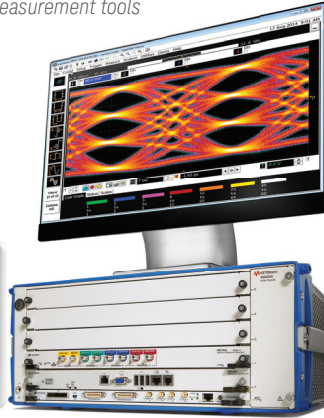
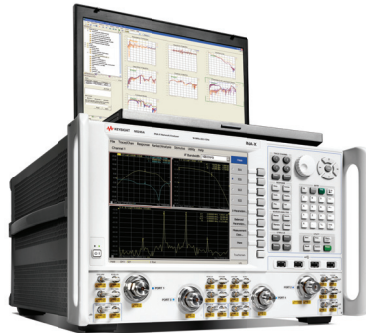
USA: 800 829 4444 CAN: 877 894 4414

© 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 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 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

AUGUST 2015, Vol. 53, No. 8

www.comsoc.org/commag

- 6 THE PRESIDENT'S PAGE
- 13 SOCIETY NEWS/NEWLY APPROVED AMENDMENTS TO THE IEEE COMSOC CONSTITUTION
- 16 SOCIETY NEWS/NEWLY APPROVED AMENDMENTS TO THE IEEE COMSOC BYLAWS
- 28 SOCIETY NEWS/THE FIRST IEEE COMSOC SUMMER SCHOOL
- 30 CONFERENCE REPORT/IEEE ICC 2015 SETS ATTENDANCE RECORD IN LONDON, UK
- 34 CONFERENCE CALENDAR
- 35 GLOBAL COMMUNICATIONS NEWSLETTER
- 208 ADVERTISERS' INDEX

SECURITY AND PRIVACY IN EMERGING NETWORKS: PART II

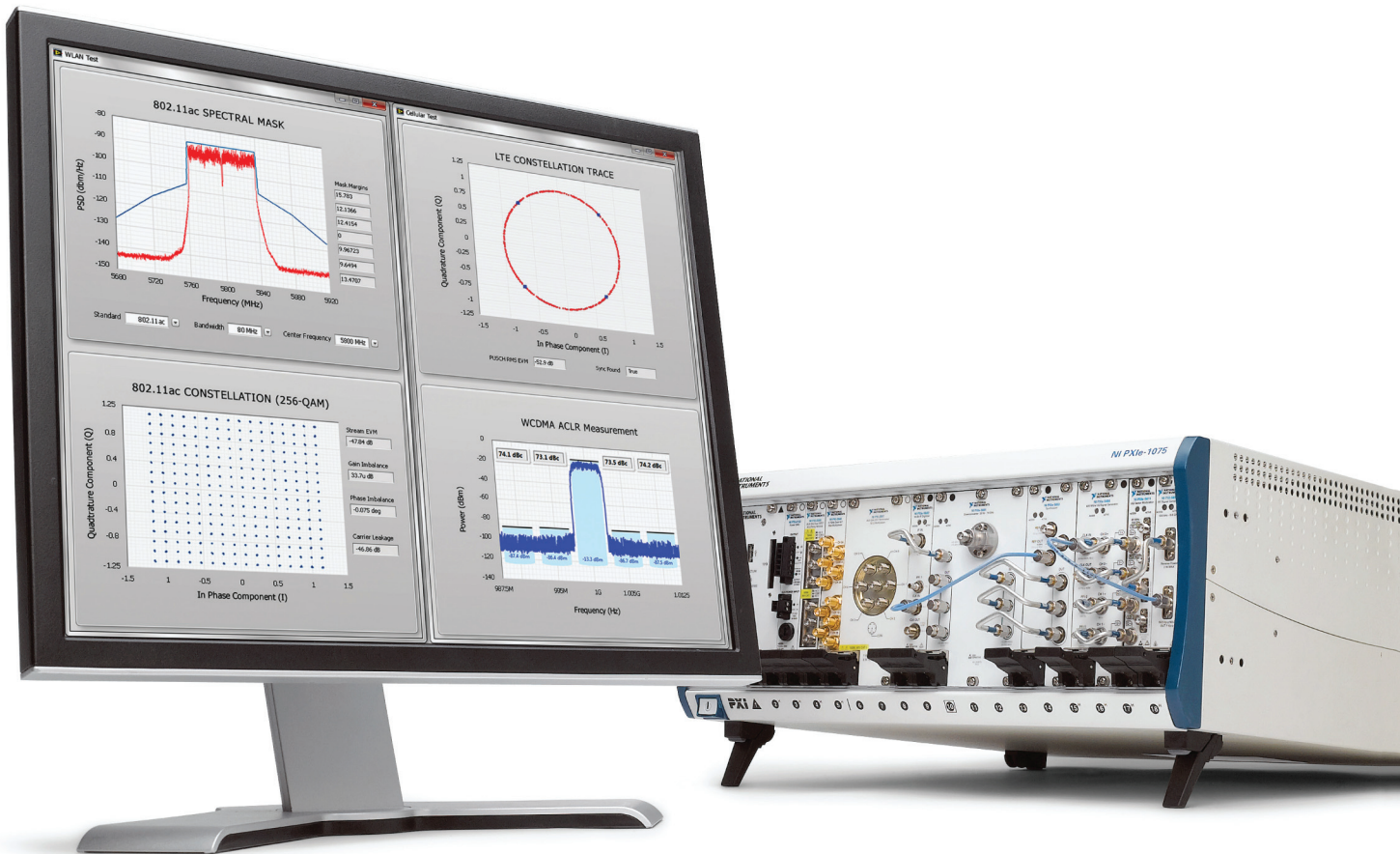
GUEST EDITORS: MOHSEN GUIZANI, DAOJING HE, KUI REN, JOEL J. P. RODRIGUES, SAMMY CHAN, AND YAN ZHANG

- 40 GUEST EDITORIAL
 - 42 TOWARD TRUSTWORTHY VEHICULAR SOCIAL NETWORKS
QING YANG AND HONGGANG WANG
 - 48 VERIFIABLE, PRIVACY-ASSURED, AND ACCURATE SIGNAL COLLECTION FOR CLOUD-ASSISTED WIRELESS SENSOR NETWORKS
CHIA-MU YU, CHI-YUAN CHEN, AND HAN-CHIEH CHAO
 - 54 SECURE COMMUNICATION FOR UNDERWATER ACOUSTIC SENSOR NETWORKS
GUANGJIE HAN, JINFANG JIANG, NING SUN, AND LEI SHU
 - 61 FROM BOTNETS TO MOBIBOTS: A NOVEL MALICIOUS COMMUNICATION PARADIGM FOR MOBILE BOTNETS
ABDERRAHMEN MTIBAA, KHALED A. HARRAS, AND HUSSEIN ALNUWEIRI
 - 68 PRIVACY-PRESERVING PARTICIPATORY SENSING
QINGHUA LI AND GUOHONG CAO
 - 75 SECURITY AND PRIVACY IN MOBILE CROWDSOURCING NETWORKS: CHALLENGES AND OPPORTUNITIES
KAN YANG, KUAN ZHANG, JU REN, AND XUEMIN (SHERMAN) SHEN
 - 82 SECURITY IN SPACE INFORMATION NETWORKS
CHUNXIAO JIANG, XUEXIA WANG, JIAN WANG, HSIAO-HWA CHEN, AND YONG REN
- ## ENERGY HARVESTING COMMUNICATIONS: PART III
- GUEST EDITORS: CHAU YUEN, MAGED ELKASHLAN, YI QIAN, TRUNG Q. DUONG, LEI SHU, AND FRANK SCHMIDT
- 90 GUEST EDITORIAL
 - 92 A HIERARCHICAL PACKET FORWARDING MECHANISM FOR ENERGY HARVESTING WIRELESS SENSOR NETWORKS
DAPENG WU, JING HE, HONGGANG WANG, CHONGGANG WANG, AND RUYAN WANG
 - 99 WIRELESS INFORMATION AND POWER TRANSFER: FROM SCIENTIFIC HYPOTHESIS TO ENGINEERING PRACTICE
RONG ZHANG, ROBERT G. MAUNDER, AND LAJOS HANZO
 - 106 DELAY-SENSITIVE DYNAMIC RESOURCE CONTROL FOR ENERGY HARVESTING WIRELESS SYSTEMS WITH FINITE ENERGY STORAGE
FAN ZHANG AND VINCENT K. N. LAU
 - 114 TOWARD SECURE ENERGY HARVESTING COOPERATIVE NETWORKS
JIAWEN KANG, RONG YU, SABITA MAHARJAN, YAN ZHANG, XUMIN HUANG, SHENGLI XIE, HANNA BOGUCKA, AND STEIN GJESSING



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

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



**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.

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 Erie, 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.

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). All submissions will be peer reviewed.



ENERGY-EFFICIENT OPTICAL NETWORKS

SERIES EDITORS: OSMAN GEBIZLIOGLU AND VIJAY JAIN

122 SERIES EDITORIAL

124 HIGH DATA RATE COHERENT OPTICAL SLOT SWITCHED NETWORKS: A PRACTICAL AND TECHNOLOGICAL PERSPECTIVE

YVAN POINTURIER, GUILHEM DE VALICOURT, JESSE E. SIMSARIAN, JÜRGEN GRIPP, AND FRANCESCO VACONDIO

130 CSO: CROSS STRATUM OPTIMIZATION FOR OPTICAL AS A SERVICE

HUI YANG, JIE ZHANG, YONGLI ZHAO, YUEFENG JI, JIANRUI HAN, YI LIN, AND YOUNG LEE

140 OPTICAL INTERCONNECTS AT THE TOP OF THE RACK FOR ENERGY-EFFICIENT DATA CENTERS

JIAJIA CHEN, YU GONG, MATTEO FIORANI, AND SLAVISIA ALEKSIC

150 THE WATCHFUL SLEEP MODE: A NEW STANDARD FOR ENERGY EFFICIENCY IN FUTURE ACCESS NETWORKS

RAISA O. C. HIRAFUJI, KELVIN B. DA CUNHA, DIVANILSON R. CAMPELO, AHMAD R. DHAINI, AND DENIS A. KHOTIMSKY

SISTER SOCIETY REPRINTED ARTICLE

158 CROWDSENDING BASED PUBLIC TRANSPORT INFORMATION SERVICE IN SMART CITIES
KÁROLY FARKAS, GÁBOR FEHÉR, ANDRÁS BENCZÚR, AND CSABA SIDLÓ

ACCEPTED FROM OPEN CALL

166 NETWORK DISTANCE PREDICTION FOR ENABLING SERVICE-ORIENTED APPLICATIONS OVER LARGE-SCALE NETWORKS

HAOJUN HUANG, HAO YIN, GEYONG MIN, DAPENG OLIVER WU, YULEI WU, TAO ZUO, AND KE LI

176 CAINE: A CONTEXT-AWARE INFORMATION-CENTRIC NETWORK ECOSYSTEM

GEORGE KAMEL, NING WANG, VASSILIOS VASSILAKIS, ZHILI SUN, PIRABAKARAN NAVARATNAM, CHONGGANG WANG, LIJUN DONG, AND RAHIM TAFAZOLLI

184 ON THE USE OF RADIO ENVIRONMENT MAPS FOR INTERFERENCE MANAGEMENT IN HETEROGENEOUS NETWORKS

JORDI PEREZ-ROMERO, ANDREAS ZALONIS, LILA BOUKHATEM, ADRIAN KLIKS, KATERINA KOUTLIA, NIKOS DIMITRIOU, AND REBEN KURDA

192 VEHICLE-TO-VEHICLE COMMUNICATION IN C-ACC/PLATOONING SCENARIOS

ALEXEY VINEL, LIN LAN, AND NIKITA LYAMIN

198 GREATER RELIABILITY IN DISRUPTED METROPOLITAN AREA NETWORKS: USE CASES, STANDARDS, AND PRACTICES

MING-TUO ZHOU, MASAYUKI OODO, VINH DIEN HOANG, LIRU LU, XIN ZHANG, AND HIROSHI HARADA

CURRENTLY SCHEDULED TOPICS

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

www.comsoc.org/commag/call-for-papers

10 MHz Rubidium Frequency Standard

- **5 MHz and 10 MHz outputs**
- **Ultra-low phase noise**
(< -130 dBc/Hz at 10 Hz)
- **0.005 ppm aging over 20 years**
- **Built-in distribution amplifier**
(up to 22 outputs)
- **1 pps input and output**

The FS725 Benchtop Rubidium Frequency Standard is ideal for metrology labs, R&D facilities, or anywhere a precision frequency standard is required.

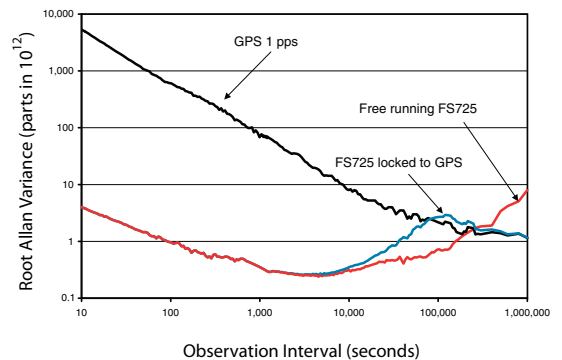
It has excellent aging characteristics, extremely low phase noise, and outstanding reliability. A 1 pps input is provided for phase-locking to GPS, providing Stratum 1 performance.

With a built-in 5 MHz and 10 MHz distribution amplifier, the FS725 is the ultimate laboratory frequency standard.

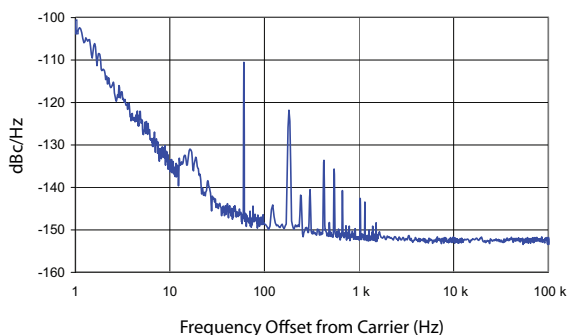
FS725 ... \$2695 (U.S. list)



Allan Variance vs. Time



FS725 Single Sideband Phase Noise



FS725 rear panel

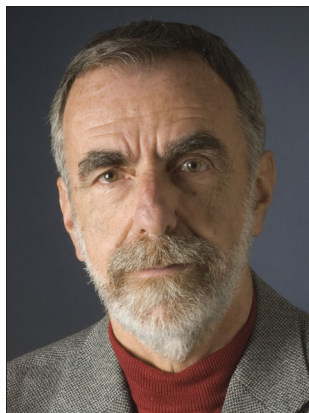
OUR JOURNEY TOWARD COMSoc 2020

The origin of the IEEE Communications Society (ComSoc) was the IRE Professional Group on Communication Systems (PGCS), which was founded in 1952 with a membership of about 600. In 1972 the PGCS was converted into a Society, ComSoc, with a membership of about 10,000 soon after the IEEE was founded by merging the IRE and the AIEE. Since then ComSoc continued to grow in membership and revenue as well as in all activities, including publications and conferences, until 2000, when ComSoc's membership peaked at about 58,000. However, membership began to decline right after that in accordance with the decline of the telecommunications industry, which continued until the membership declined to about 43,000. To make it worse, ComSoc's finances were seriously impacted during the global economic downturn of 2008.

ComSoc's decline in both membership and revenue was a warning signal for the Society to better align itself with the rapidly changing communications industry. The divestiture of AT&T in 1984 and the subsequent privatization of the telecommunication industry in many countries were the start of a tsunami-like global transformation of the communications industry. The telecom industry began to face competition in the 1980s and, in the 1990s, began to confront the challenge of the Internet, which was energized by the introduction of personal computers, the world-wide web, Navigator/Explorer, and the commercial Internet. In fact, the challenge of the Internet was the start of a new transformational dimension that led to the convergence of communications and computers. Despite such a drastic transformation, ComSoc did not seem to sense its impact, nor attempt to change its operations.

The first steps toward rebuilding our Society were taken by the 2008-09 President, Doug Zuckerman, with the "ComSoc 2.0" vision. This vision, which was realized during his two-year term, was that of a member-driven Society that closely matched member needs and maximized the value offered to members in the newly transformed worldwide community. (See "ComSoc 2.0 – Member Driven," President's Page, January 2008, *IEEE Communications Magazine*).

In 2010, on this member-driven base, ComSoc began to keenly respond to the serious impact of the industry transformation by creating future visions and initiating the necessary changes to ComSoc. It was the starting point of our journey toward ComSoc 2020. As the first step, Byeong Gi Lee, the 2010-2011 President, declared a "ComSoc's Golden Triangle" Vision in January 2010 with the goal of making internal innovation a fundamental driver of ComSoc's operation, with emphasis on three areas: globalization, young leaders, and industry. By building consensus among ComSoc's volunteer leaders on one side and by taking various actions to implement the vision on the other, Byeong was able to end 2010 with a big success in membership recovery and financial stabilization. ComSoc's membership recovered to the level of



SERGIO BENEDETTO



BYEONG GI LEE

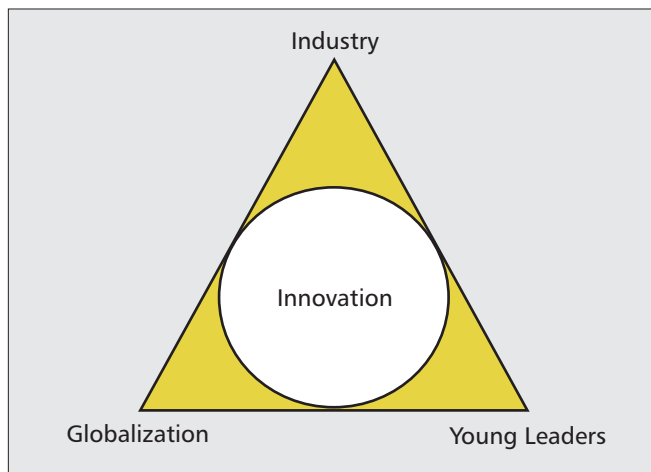
50,000 and ComSoc's finances recovered a stable state with a surplus of approximately \$500,000 on a \$17M budget. In December 2010, Byeong declared a "ComSoc's Growth Engine" Vision with the goal of creating new sources for ComSoc's growth and, at the same time, organized a "ComSoc 2020" Vision team to study ComSoc's vision toward the year 2020 and invited Roberto Saracco to chair the team. ComSoc's healthy state of membership and revenue of 2010 continued into 2011 with membership increasing to 51,000 and the finances stabilized with a surplus of about \$350,000. (Unfortunately, however, the membership increase and the financial stability did not continue after that.) In December 2011, after a full year's in-depth study, Roberto's vision team submitted to the ComSoc Board of Governors (BoG) a masterpiece "ComSoc's 2020" Vision report. With that report, the goal and direction of our journey toward ComSoc 2020 took a concrete form.

Once the vision toward ComSoc 2020 was set, the next step was implementation of the vision. Without implementation, a vision is no more than a dream. The 2012-2013 President, Vijay Bhargava, started to set a strategic plan for implementing the vision by inviting Roberto Saracco to chair the Society's Strategic Planning Committee (SPC). After two years, the 2014-2015 President, Sergio Benedetto, restarted the implementation process by inviting Byeong to chair the Strategic Planning Committee. Byeong and his SPC

team exerted dedicated efforts to successfully produce a strategic plan for implementing the ComSoc 2020 Vision. The SPC conducted the job in four steps. First, it identified four strategic goals out of the ComSoc 2020 Vision report. Second, it then set up four business plans to implement them. Third, after reporting the business plans to the BoG meeting of December 2014, it studied and proposed how to change ComSoc's leadership structure to effectively support the implementation of the strategic business plans. Fourth, it studied and recommended how to incorporate the structural changes into ComSoc's Bylaws. The SPC's proposal and recommendations were approved at the June 2015 BoG meeting. It was a memorable moment, setting the stage for making the ComSoc 2020 Vision real, and was the first step of our journey toward ComSoc 2020.

"COMSoc's GOLDEN TRIANGLE" VISION

"ComSoc's Golden Triangle" vision was presented in January 2010 by Byeong at the start of his term as the 2010-2011 President. (See "ComSoc's Golden Triangle: Globalization, Young Leaders, Industry," President's Page, January 2010, *IEEE Communications Magazine*.) It was a vision set to trigger a change in ComSoc's operations by stressing internal innovation that matched the changing communications and global environment so that ComSoc can better perform its honorable mission of serving humanity. The vision identified three areas



“ComSoc’s Golden Triangle.”

for concentrated efforts, the so-called “three vertices”: Globalization, Young Leaders, and Industry.

The essence of the “Globalization” vertex in the Golden Triangle was the so-called “level-2 globalization”. Globalization had already progressed in the demographic sense at that time, with non-U.S. membership having increased to a 60 percent level, but the number of Members-at-Large (MaL) on the Board of Governors was far less than that percentage, with almost no members representing the Asia-Pacific and Latin American Regions. This meant that ComSoc’s decision-making was not fully reflecting the Society’s cultural diversity and, therefore, ComSoc’s operations were not matching the global trend as well as it could if it did have more diverse leadership. So, for ComSoc to keep abreast with ecosystem shifts, it needed to incorporate more global cultures to innovate its operations and decision-making. This is what the term “level-2” globalization meant. Though simple and clear in concept, it took a long time to actually realize this “globalization” vertex. After about three years’ effort by the Nominations and Elections Process Committee, chaired by Larry Greenstein, and the Nominations and Elections Committees chaired by Doug Zuckerman and then by Byeong Gi Lee, we were able to attain the current rule of balanced geographical representation, with one, one, and two MaL members elected respectively from the Asia-Pacific, Europe-Middle East-Africa (EMEA), and the (North and Latin) American Regions every year.

The “Young Leaders” vertex was intended to attract capable young people to ComSoc and then open up opportunities to develop their careers by participating in ComSoc activities and grow themselves into becoming ComSoc’s future leaders. By embracing young people for roles in ComSoc’s leadership structure, we are assuring a vibrant and healthy future for the Society that is well suited to meet the needs of the greater communications community. A key theme of this vertex was the “open calls” that aimed at attracting more young people to ComSoc activities and leadership roles. The “open call” approach was successfully adopted for journal and magazine editorial appointments (thanks to VP-Publications Len Cimini and his Directors and Editors-in-Chief) and for conference steering committee membership (thanks to VP-Conferences Khaled Ben Letaief and his Directors). For example, we adopted the so-called “Open-Call Based 3-2-3 System” in appointing editors of all ComSoc journals and magazines from January 2011, and regulated a charter including a chair election rule to each steering committee of ComSoc conferences in 2011.

The “Industry” vertex was intended to exert dedicated effort to move ComSoc closer to industry and vice versa, thereby reestablishing industry as a vibrant part of ComSoc’s membership and activities. The essential idea was that engineering without industry is no longer engineering, and that academia without industry does not mean much. We developed several programs such as the Corporate Patron Program (CPP) and the Industry Now Program (INP), which both contributed much to ComSoc’s recovery from the global economic downturn. The most impactful venture, however, was the creation of VP-Standards Activities, which was included as a part of “ComSoc’s Growth Engines” Vision. Standards are a most important, and sometimes crucial element to industry, but ComSoc had not been playing any active or important role in the broader IEEE standards activities. Though the IEEE 802 standards were well established as the dominating standards for computer communications, ComSoc had never been a part of the related standards activities. As a symbolic expression of ComSoc’s resolution to come close to industry and also as a public announcement of ComSoc’s resolution to actively get involved in standards activities, in 2011 we created the position of VP-Standards Activities. The first VP-SA, Alex Gelman, dedicated his efforts to launch a new standards organizational structure to enable working with the IEEE Standards Association. His and industry colleagues’ efforts resulted in ComSoc playing a key role in creating new IEEE standards on power-line communications.

“COMSOC 2020” VISION

With “ComSoc’s Golden Triangle” Vision vigorously being implemented and stabilizing, Byeong announced “ComSoc’s Growth Engines” Vision at the end of 2010. (See “ComSoc 2020 in the Converged Communications Era.” President’s Page, January 2010, *IEEE Communications Magazine*.) The goal of this new vision was to develop new sources of revenue generation for ComSoc’s future growth. The vision had three key elements: “Content, including education and training content”; “Industry in converged communications”; and “Standards activities in emerging technologies”. The first element was intended to develop the content area, including positioning education and training as a third revenue-generating “pillar” in addition to the existing two pillars, publications and conferences. The other two vision elements, both related to industry, were intended to engage industry as a fourth revenue-generating pillar. However, one year was not long enough to develop plans and implement them, but nevertheless we were able to create a VP position for standards activities. The two revenue-generating pillars in the Growth Engine Vision were not developed beyond 2011, but fortunately they were revived in 2014-2015 in the process of strategic planning for the ComSoc 2020 Vision.

In December 2010, Byeong announced the creation of the “ComSoc 2020” Vision team with the hope that ComSoc would survive and surpass the environmental changes of the decades to come. Byeong invited Roberto Saracco to chair the team and called for participation among BoG members on the “ComSoc 2020” Vision team. Those who volunteered and then participated in the vision team were Vijay Bhargava, Vincent Chan, Celia Desmond, Andrzej Jajszczyk, Gabe Jakobson, Nelson Fonseca, Rob Fish, Alex Gelman, Shri Goyal, Russel Hsing, Mark Karol, Stan Moyer, Zhisheng Niu, John Pape, Vince Poor, Parag Pruthi, Sara Kate Wilson, and Doug Zuckerman. After a full year of study, Roberto and his vision team completed a “ComSoc 2020” Vision Report and reported it to the BoG in December 2011. The report (pub-

STRATEGIC BUSINESS PLANS FOR COMSOC 2020

Sergio invited Byeong to chair the 2014-2015 Strategic Planning Committee and organized the SPC with the following members: Merrily Hartmann (Secretary), Doug Zuckerman, Rob Fish, Stefano Galli, James Hong, Mark Karol, Khaled Letaief, Zhisheng Niu, Roberto Saracco, Steve Weinstein, Stefano Bregni*, Hikmet Sari*, Sara Kate Wilson*, and Harvey Freeman* (*non-voting). At the start, all SPC members reviewed the ComSoc 2020 Report and other related documents to identify key strategic items to focus on to implement the ComSoc 2020 Vision. As a result, four strategic goals were identified:

- Education: Education program, content, services.
- Industry: Standards, industry programs and services.
- Technical: New technical values creation.
- Member: Serving members globally at individual level.

The SPC organized four groups to study the four strategic items individually. Group 1, chaired by Khaled with the members Stefano Bregni, Stefano Galli, Mark, Roberto, Steve, Katie, and Michele Zorzi, studied the education issue. Group 2, chaired by Rob with the members Harvey, Stefano Galli, Hikmet, Steve, and Katie, studied the industry issue. Group 3, chaired by Mark with the members Harvey, Stefano Galli, Khaled, Zhisheng, Hikmet, Katie, and Doug, studied the technical issue. Group 4, chaired by Zhisheng with the members Stefano Bregni, Stefano Galli, James, Roberto, Hikmet, Koichi Asatani, and Elena Neira, studied the member issue. The four groups studied their allocated study issues individually and the results were discussed at the monthly SPC teleconferences. In June 2014 the SPC held a one-day workshop to finalize the four strategic plans of the four groups. The resulting key strategies are as described below. (See “ComSoc Strategic Planning for the “Smart Revolution,” President’s Page, October 2014, *IEEE Communications Magazine*.)

Education: In light of the changing communications landscape and the potential of online education, ComSoc needs to participate in this new value ecosystem by providing valuable services to its members and the broader community. In the long run, it may be desirable for ComSoc to run a comprehensive online interactive education system that embraces continuing education, training, certification, and an educational content portfolio in the communications field. As a near-term strategy, however, it will be desirable to start with a new “Continuing Professional Education and Training (CPET)” program targeted at specific industries or industry associations. Continuing education is expected to be very important for practicing engineers in industry to keep in tune with the “Smart Revolution”. ComSoc should leverage this opportunity to enhance services to our members by developing courses and training programs for them. It will then help us to build a third pillar for revenue generation.

Industry: In accordance with the industry transformation toward ICT convergence, ComSoc must expand its horizon to accommodate the entire ICT industry and must provide services to all those engineers involved in the converged ICT industry. Noting that a major reason for the existence of ComSoc is to serve the informational and professional development needs of industry-based members and companies, we should develop a credible value proposition for the ICT industry that convinces their executives that ComSoc is a business necessity. From the standards’ perspective, we should develop a rich set of standards-related programs in the form of forums, symposia, education, conferences, and publications that will tighten our relationship with industry enterprises and members. Therefore, it is desirable to set a comprehensive “Industry Outreach” program that is a cost-effective outreach

program of ComSoc’s products and services to industry and governmental communities.

Technical: The core value of ComSoc is the technical expertise of its volunteer participants in communications. Since the conventional communication industry has been transformed into the converged ICT industry, with a large number of non-communications people in different sectors joining the enlarged ICT world, it is important for ComSoc to recognize them as potential members and provide the services they need. In order to accommodate those potential members and to expand the technical horizon of the existing members, we should develop new services to meet their technical demand in the ICT environment. First of all, we should consider developing a one-stop ICT service system that provides all technical information they may be curious about and want to learn. Such a system should also be of high value to our existing members for learning more about the expanded ICT world. Therefore, it is desirable to build a “One-Stop ICT Service System (OSISS)” as the technical strategic project for new technical value creation.

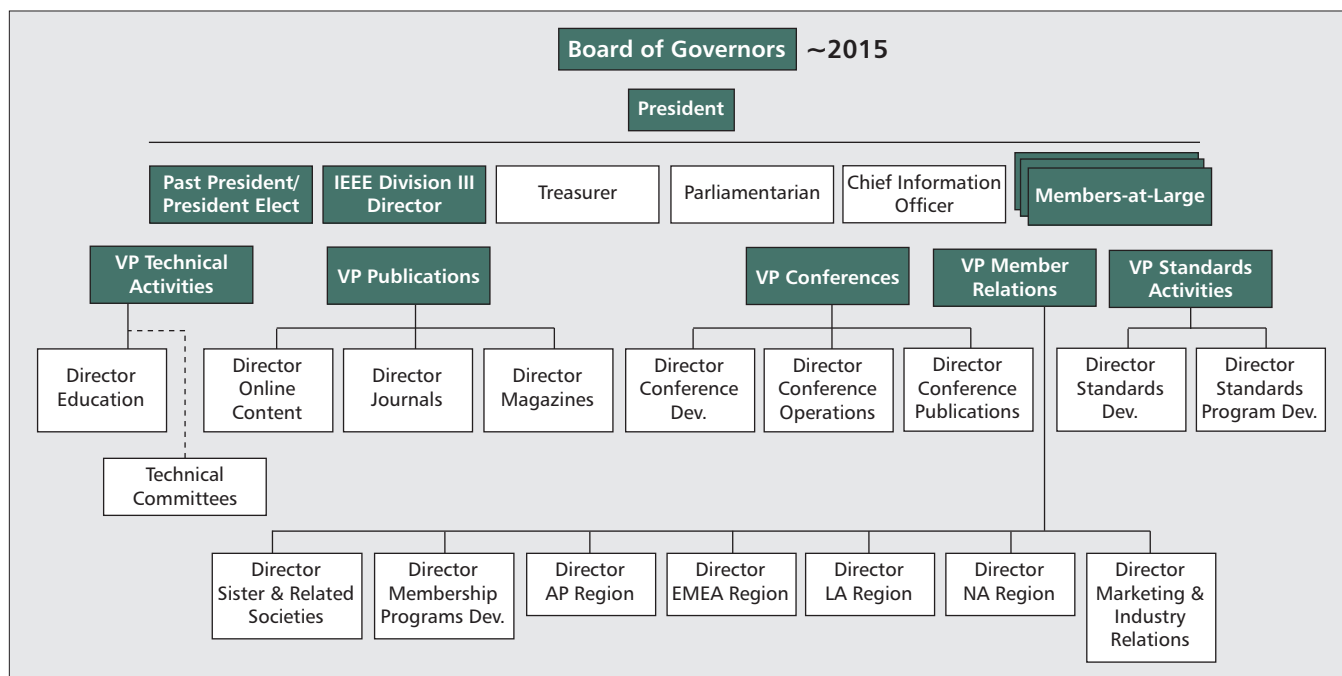
Member: A challenge that ComSoc faced over the last decade was the trend of declining membership and the difficulty of member retention. In order to reverse that trend, ComSoc has exerted a wide variety of efforts, but they were not fruitful because they were a kind of “blind” approach, not targeted at the interest of individual members. The time is now ripe to take advantage of today’s ICT technology and business practices to provide services to members at an individual level. If ComSoc changes the notion of membership service to “targeted service” such that ComSoc can care for members at an individual level, it will surely make a big difference in retaining membership. Therefore, we should take individual-level membership services as a strategic item and develop a “Member Activity and Service Supporting System (MASSS)” to operationalize the strategy. The goal of this system is to serve our members globally at an individual level, thereby creating new values to our members and contributing to the financial health of ComSoc as well.

After reporting the above four strategic goals to the BoG in June 2014, the SPC started developing the four business plans that can effectively realize those strategic items. The four groups extended their studies to develop concrete business plans that can realize their strategic goals and reported their results to the bi-monthly SPC teleconferences. By December 2014 the SPC completed developing the four business plans and Byeong reported them to the BoG at its December 2014 meeting where he requested that Sergio and the VPs take over implementation of the strategies. In January 2015 the SPC held a teleconference for a final review of the four business plans and prepared a final report. Byeong presented the final report in detail at the Management Retreat held in January 2015. After in-depth discussions at the Retreat, Sergio accepted the four business plans and directed three relevant VPs (VP-Technical Activities, VP-Standards Activities, and VP-Member Relations) to take the first steps toward implementing the four strategic business plans.

STRUCTURAL CHANGE AND BYLAWS REVISION

After submitting the final report describing the four strategic business plans, in February 2015 the SPC began to study how to change ComSoc’s structure to practically support the implementation of the business plans. It was the final step to complete the strategic plans for revitalizing ComSoc by reflecting the ComSoc 2020 Vision. We set a goal to secure one Director to take charge of each business plan and limit the structural changes at the Director level. In addition, we

THE PRESIDENT'S PAGE



Current ComSoc governing structure.

chose to meet the goal by redefining the responsibilities of some Directors without increasing the total number of Directors. We assumed that the structural change is a Stage 1 change, which would be followed by a subsequent structural change after a few years when the education and other businesses grow up to their target levels. The resulting structural changes are as described below.

First, we created the *Director-Educational Services* by redefining the existing Director-Education and Training. This Director (specifically, the Board chaired by this Director) is responsible for overseeing the Society's education and training activities, including administration of the Society's programs on continuing education, incorporating tutorials, short courses, lectures, and so on. In particular, it is responsible for developing and maintaining continuing professional education and training programs, or the CPET program.

Second, we created the *Director-Industry Outreach* by redefining the existing Director-Marketing and Industry Relation. This Director (specifically, the Board chaired by this Director) is responsible for assuring a comprehensive and cost-effective industry outreach program of the Society's products and services to industry and governmental communities. It is also responsible for developing liaisons with communications and networking related enterprises to promote ComSoc products and services and to attract industry and government leaders into ComSoc's volunteer community.

Third, we created the *Director-Technical Services* by taking the Director slot made available by removing the Director-Conference Publications and transferring its responsibility to the Director-Conference Operation. This Director (specifically, the Board chaired by this Director) is responsible for overseeing the promotion of the technical communities of the Society and their activities, including the promotion of technical content and development of educational content. In particular, it is responsible for developing and providing one-stop ICT (information-communications technology) services, or the OSISS service.

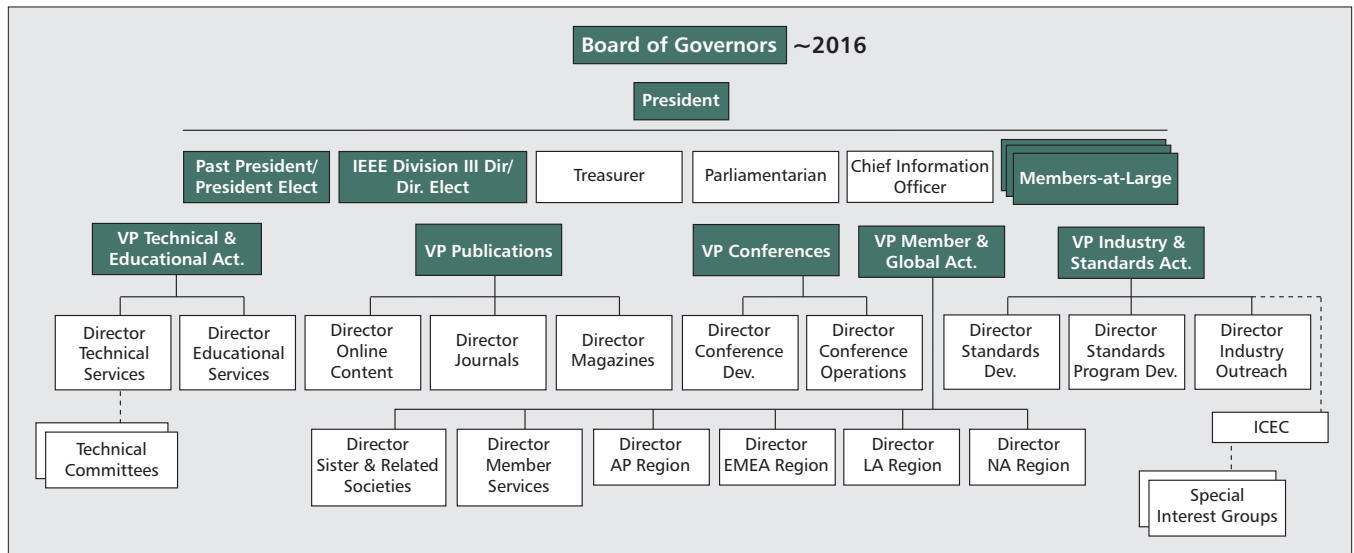
Fourth, we created the *Director-Member Services* by redefin-

ing the existing Director-Membership Programs Development. This Director (specifically, the Board chaired by this Director) is responsible for overseeing all services and programs addressed to members and Chapters, and oriented to membership retention and development in the four ComSoc geographical regions. In particular, it is responsible for developing and providing individual-level membership services globally, or the MASSS service.

In addition, we modified the responsibility of the *Director-Sister and Related Societies* such that it can more proactively stretch out ComSoc's activities globally by working more closely with ComSoc's and other Societies' chapters as well as related non-professional societies around the world. This Director (specifically, the Board chaired by this Director) is thus responsible for enhancing Society activities with our sister and related societies (SRS) by developing new programs with SRS and cooperating with SRS in offering Society/SRS products and services globally, and strengthening the Society's global and professional reach.

To augment the above Director-level structural changes, we changed the title of some Vice President positions as well as some reporting lines. First, we changed VP-Technical Activities to *VP-Technical & Educational Activities (TEA)* with the Director-Educational Services and the Director-Technical Services reporting to it. Second, we changed VP-Member Relations to *VP-Member & Global Activities (MGA)* with the Director-Member Services, the Director-Sister and Related Societies, and the four Regional Directors reporting to it. Third, we changed VP-Standards Activities to *VP-Industry & Standards Activities (VP-ISA)* with the Director-Industry Outreach, the Director-Standards Development, and the Director-Standards Programs Development reporting to it. We also changed two reporting lines. First, Technical Committees (TC), which previously reported to the VP-Technical Activities, will report to the Director-Technical Services. Second, the Industry Content & Exhibition Committee (ICEC), which previously reported to the President, will report to the VP-ISA, and the Special Interest Groups will report to the ICEC.

THE PRESIDENT'S PAGE



New ComSoc governing structure, effective 2016.

After completing the Stage 1 structural changes through repeated email and teleconference discussions, the SPC started revising ComSoc's Bylaws to reflect the changes and provide a means for formal approval by the BoG. The SPC first prepared a draft of the Bylaws revision and sent it to the Governance Committee (GC), chaired by Paul Cotae and vice-chaired by Stefano Galli. The GC performed an extensive and detailed review of the draft and returned a revised draft to the SPC. In June 2015 the SPC opened a final review session through teleconference and then finalized the revision of the Bylaws. The BoG approved the proposed structural changes and Bylaws revision at the June 2015 meeting, and thus the Stage 1 ComSoc structural changes became officially effective. In parallel, the Executive Director of ComSoc, Susan Brooks, devised a new internal organization of ComSoc staff to complement and match the society's new structure.

IMPLEMENTATION OF BUSINESS PLANS

Implementation was started in 2015 when the SPC completed creating four strategic business plans at the end of 2014. After the Management Retreat of January 2015, the implementation of the four business plans was handed over from Byeong to Sergio, with pertinent VPs taking charge of their initial-stage implementation. Specifically, the VP-Technical Activities, Khaled Letaief, took the responsibility of the initial trials of the "CEPT" program and the initial development of the "OSISS"; the VP-Standards Activities, Rob Fish, took responsibility of the detailed planning of the "Industry Outreach" program and its initial-stage implementation; and VP-Member Relations, Stefano Bregni, took responsibility of the initial development of the "MASSS". The VPs were supposed to start developing the programs and systems as much as they can and then relay the work, together with their experiences, to their successor VPs in 2016 so the new VPs can continue implementing the four business plans without any disconnection or discontinuity.

The ComSoc structural changes and Bylaws revision approved at the BoG meeting of June 2015 had the effect of putting two wings on the shoulders of the 2016-2017 VPs in charge of those four business plans to lead ComSoc's journey toward ComSoc 2020. Under the VP-TEA, two directors, the Director-Education Services and Director-Technical Services, are dedicated to implementing the business plans for CEPT and OSISS, respectively; under the VP-ISA, the Director-

Industry Outreach will take charge of the Industry Outreach program; and under the VP-Members and Global Activities, the Director-Member Services will take charge of the MASSS project. Therefore, we are now fully equipped to implement the four business plans from an organizational perspective.

However, the structural changes and Bylaws revision do not truly guarantee a successful implementation of the four business plans. The Directors in charge must have the leadership skills and willingness to follow through. Success will heavily depend on appointing the right people in 2016. They must be well qualified and experienced and, more importantly, they must be resolved to dedicate themselves to successfully fulfill their mission. By appointing such individuals, our 2016-2017 President, Harvey Freeman, will play a key role in assuring that we will have a successful journey toward ComSoc 2020.

If the implementation efforts go smoothly with the four business plans fully developed and implemented, ComSoc will have positioned itself to see a much more promising future by 2020 or by 2025 at the latest. The CPET program will make ComSoc recognized by professional engineers as the center for continuing education and training and, at the same time, will help ComSoc by creating additional revenue of about five million dollars each year. The Industry Outreach program, if it is well developed in harmony with the standards programs, will make ComSoc recognized by industry people as "the Society for industry," thereby attracting a large number of industry members to join and engage in ComSoc activities, and at the same time will help create additional revenue of about three million dollars each year. The OSISS project will make ComSoc recognized by all ICT people and all those who are interested in ComSoc as "the Center for ICT" and will thereby attract many new people to join ComSoc and engage in its activities. The developmental process of the OSISS can also provide a valuable opportunity for ComSoc to closely collaborate with the Computer Society and other related Societies. The MASSS project will make ComSoc recognized by ComSoc members as the ideal technical home to stay and develop their careers by actively participating in ComSoc's activities.

GOING FORWARD

From the view of ComSoc's operations, the four business plans contribute to the very fundamental elements of revenue and membership. The CPET and the Industry Outreach pro-

THE PRESIDENT'S PAGE

Four strategic business plans toward ComSoc 2020

	Strategic goals	Business plans	Directors in charge	VPs in charge
1	Education Programs, Content, and Services	CPET: Continuing Professional Education & Training	Director-Educational Services	VP-Technical & Educational Activities
2	Standards, Industry Programs and Services	Industry Outreach	Director-Industry Outreach	VP-Industry & Standards Activities
3	New Technical Values Creation	OSISS: One-Stop ICT Service System	Director-Technical Services	VP-Technical & Educational Activities
4	Serving Members Globally at Individual Level	MASSS: Member Activity and Service Support System	Director-Member Services	VP-Member & Global Activities

grams will contribute to creating additional revenues. If everything goes according to plan, these two programs will contribute toward creating new revenues of about eight million dollars, thereby growing ComSoc from a 17 million dollar operation to a 25 million dollar operation. On the other side, the OSISS and MASSS projects, as well as the Industry Outreach program, will contribute to increasing ComSoc membership. Specifically, the OSISS and the Industry Outreach will contribute to recruiting new members, and the MASSS will contribute to retaining existing members. If everything is well managed, ComSoc should again see its membership exceed the 58,000 of 2000, and even reach well over 60,000. Such a bright future of 2020, based on the healthy foundation of revenue and membership, will stimulate various new Society activities that nicely match the converged ICT communications environment, thereby enabling ComSoc to better serve humanity in this new ICT era.

Our journey toward ComSoc 2020 is a long journey to rejuvenate the Society in the new ICT convergence era. During the three decades of industry transformation, conventional circuit-mode communications has been completely replaced with IP packet-mode communications, and communications itself has blurred the line between computing and communications. This has created the era of a converged ICT industry, where we stand today. Going forward, ComSoc will embrace the opportunities created by this new environment, expanding its industry partnerships beyond the conventional telecom companies such as AT&T, Verizon, Telecom Italia, NTT, KT, and others, to also include Microsoft, Amazon, Apple, Google, Facebook, and others. At the same time, it will expand its membership base beyond communications engineers to include all those people participating in ICT. We look to you to be part of reaching the destination of our journey toward ComSoc 2020.

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

www@comsoc.org "About Us" (bottom of page)

NEWLY APPROVED AMENDMENTS TO THE IEEE COMSOC CONSTITUTION

In June 2015, the ComSoc Board of Governors approved a revision of ComSoc Constitution. The IEEE approved the revision in July 2015, and ComSoc Membership approval is currently pending.

The main changes in this revision can be summarized as follows:

- Alignment with IEEE governing documents.
- Alignment with ComSoc Bylaws, which were also revised in June 2015 – see this issue page 16.

- Various clarifications and language improvements.

The revised Constitution and the previous one can be found here:

<http://www.comsoc.org/about/documents/constitution>

Objections to proposed changes must be emailed to Susan Brooks, ComSoc Executive Director, at s.m.brooks@comsoc.org by 15 October 2015.

IEEE COMMUNICATIONS SOCIETY CONSTITUTION

(IEEE APPROVAL: JULY 2015)

(COMSOC MEMBERSHIP APPROVAL: PENDING)

Table of Contents – Articles

1. Name, Purposes and Scope
2. Bylaws
3. Policies and Procedures
4. Membership
5. Organization
6. Finances
7. Member Services
8. Amendments

Article 1 - Name, Purposes, and Scope

- 1.1 Name – The name of this organization is the IEEE Communications Society, hereinafter referred to as “the Society.” It is organized within the Institute of Electrical and Electronics Engineers, Inc. hereinafter called “the IEEE.”
- 1.2 Purposes – The purposes of the Society are: (a) scientific and educational - directed toward the advancement of the theory, practice and application of communications engineering and related arts and sciences; (b) professional - directed toward promotion of high professional standards, development of competency and advancement of the standing of members of the profession it serves.

The Society promotes cooperation and exchange of information among its members and those of other organized bodies within and outside the IEEE. Means to these ends may include, but are not limited to, the holding of meetings for the presentation and discussion of papers, the publication of journals, sponsorship of tutorial seminars and workshops, stimulation of research, the education of members, establishment of standards, and providing for the technical and professional needs of its members via organized efforts.

- 1.3 Scope – The IEEE Communications Society embraces the science, technology, applications and standards for information organization, collection and transfer using electronic, optical and wireless channels and networks, including but not limited to:
 - Systems and network architecture, control and management;
 - Protocols, software and middleware;
 - Quality of service, reliability and security;
 - Modulation, detection, coding, and signaling;
 - Switching and routing;
 - Mobile and portable communications;
 - Terminals and other end devices;
 - Networks for content distribution and distributed computing; and
 - Communications-based distributed resources control.
- 1.4 Authority – Society organization and operations are in accordance with the IEEE Constitution and Bylaws, IEEE Policies, and IEEE Operations Manuals of Major Boards reporting to the IEEE Board of Directors.

SOCIETY NEWS

Article 2 - Bylaws

- 2.1 Bylaws are rules and regulations adopted by the Society for governing its members and for the overall management of its affairs. They provide guidance to govern all phases of the organization, management and activities, as outlined in the Constitution. Bylaws may not be in conflict with the Constitution. Bylaws are approved and amended by the Board of Governors and may be changed as Society interests evolve.

Article 3 - Policies and Procedures

- 3.1 Policies and Procedures provide more detailed statements about specific policies, objectives, and procedures than are contained in the Constitution or Bylaws.

Article 4 - Membership

- 4.1 IEEE members of any grade shall become Society members upon application and payment of the Society membership fee. The membership fee and the cost of publications and other considerations members receive for the membership fee, are to be set as part of the Society's annual budget.
- 4.2 Grades of membership for the Society are the same as those of IEEE plus Affiliates. When used in the Bylaws, the term "member" includes all grades of membership.
- 4.3 Individuals who are not members of the IEEE may become Affiliate members of the Society upon (a) meeting the requirements established in the IEEE Bylaws for Affiliate membership; (b) by making proper application for Affiliate membership; and, (c) by making appropriate payment for Affiliate membership. Any other requirements for Affiliate membership shall be as established in the IEEE Bylaws. Affiliate members have the rights and privileges of a Society member. An Affiliate member may serve in an appointed capacity that does not require IEEE membership as a prerequisite.

Article 5 - Organization

- 5.1 Board of Governors – The Society is governed by an administrative body called the Board of Governors (BoG).
- 5.1.1 Elected officers – The BoG includes officers elected by the Society membership, such as, but not restricted to:
- President, immediate Past President, and President-elect – The President is the chief executive officer of the Society.
 - Vice Presidents – Each chairs a Council responsible for a key area of interest to the Society and reports to the BoG and OpCom.
 - Members-at-Large – They are elected for staggered multiyear terms. The operations of the Society are assessed periodically, and the number of Members-at-Large is adjusted in accordance with membership needs and growth.
 - Society membership representative(s) to IEEE Board of Directors
 - All elected officers have full voting privileges in handling Society affairs.
- 5.1.2 Appointed officers – The BoG includes a group of officers appointed by the President. As specified in the bylaws, such appointments may be made upon recommendation of the appropriate Vice President or not, and may require BoG approval or not. These officers assist in managing Society activities and are supervised by the President and appropriate Vice President. Appointed officers serve concurrently with the nominal terms of President and Vice Presidents, and include:
- Treasurer responsible for Society financial affairs.
 - Directors, each of whom chairs a major Board responsible for an area important to the Society.
 - Other officers with special titles and duties.
 - Appointed officers may propose resolutions but do not have voting privileges in the BoG. If appointed Officers are serving concurrently as elected Officers, then they have voting privileges in the BoG.
- 5.2 Operating Committee – Management of Society affairs between regular and special meetings of the BoG is delegated to an Operating Committee (OpCom). Actions of the OpCom shall be ratified at the next BoG meeting, except those actions taken in areas already delegated.
- 5.3 Councils – Councils are responsible for the policies of their Boards and Standing Committees and they oversee the operations of their Boards and Committees, address issues common to all their Boards/Committees, address issues that could not be resolved at the Board/Committee level, and escalate such issues to the BoG if they cannot be resolved at the Council level.
- 5.4 Boards – Boards are the major operational entities of the Society and are organized under Councils. Boards have scopes aligned with the scope of their Councils and, within their scope, they can decide policies and make operational decisions as allowed by their Councils. Major Boards are established and dissolved through resolutions approved by the BoG.
- 5.5 Committees – The Standing Committees are established and dissolved through resolutions approved by the BoG. Ad Hoc Committees may be established by the President or by the BoG and dissolved through resolutions approved by the BoG.
- 5.6 Professional Staff – The staff consists of paid professional employees of IEEE who support the activities of the Society. The staff is managed by a Society Executive Director (hereafter referred to as Executive Director) who also serves as BoG Secretary.

SOCIETY NEWS

Article 6 - Finances

- 6.1 Assets – All funds and property held by or for the Society are vested in the IEEE.
- 6.2 Revenues – Basic revenues consist of fees or assessments that are levied on members of the Society for membership dues, covering publications supplied and services rendered to all members. Other revenues may be raised from the sale of Society publications, advertising, expositions, contributions, and from other sources consistent with IEEE regulations. Proposed new income sources require the approval of IEEE.
- 6.3 Conference Registrations – Suitable registration charges may be collected by the Society from members and non-members attending Society meetings, symposia, conferences and conventions, consistent with IEEE policy and regulations.
- 6.4 Budget – An annual budget shall be prepared and approved by the BoG and IEEE in advance of each fiscal year. Any changes to the budget, or expenditures in excess of budgeted amounts or for unbudgeted items, require advance approval by the BoG before commitment and/or payment.
- 6.5 Debts – Neither the Society nor any officer or representative thereof has any authorization to contract debts for, pledge the credit of, or in any way bind the IEEE without prior approval by IEEE.

Article 7 - Member Services

- 7.1 Meetings and Conferences – Principal Society meetings are conferences, workshops, symposia and conventions, held either alone or in cooperation with other IEEE units and/or other professional or technical organizations.
 - 7.1.1 Organization – Meetings are organized according to IEEE regulations. All meetings are open on an equal basis to all IEEE and Society members. Registration fees at Society meetings may differ for the various grades of IEEE and Society members and be higher for non-members.
 - 7.1.2 Papers Selection – At meetings sponsored by the Society or in which it participates, the methods of selection of papers for presentation, and the publications procedures are consistent with IEEE regulations. The Society, through the Committees within the Technical Activities Council, offers guidance in the solicitation and review of papers and in organizing and moderating sessions at meetings.
- 7.2 Publications – The Society, subject to the editorial and fiscal policies of the IEEE, publishes magazines, transactions, journals and other technical materials, such as leading-edge technical articles, tutorials, conference papers, etc. Fees charged for such publications may be higher for non-member subscribers and purchasers than for Society members.
- 7.3 Education – Principal educational activities include basic and continuing education programs.
- 7.4 Standards – The Society sponsors standards development in accordance with the process defined and approved by IEEE Standards Associations. It also organizes standards-related activities that comply with applicable IEEE/ComSoc and/or IEEE-SA policies.

Article 8 - Amendments

- 8.1 Constitution
 - 8.1.1 Origin – Amendments to this Constitution may be initiated by:
 - Proposal approved by the BoG.
 - Petition submitted to the President by a minimum of 100 Members.
 - 8.1.2 Procedure on Proposals – Proposed amendments to the Constitution require two-thirds majority vote of all the voting members of the BoG. Amendments are subject to the approval of the IEEE Technical Activities Board (TAB). After approval by TAB, the proposed amendment shall be published in the Society magazine, or directly mailed to the membership. The amendment becomes effective unless one percent (or more) of the membership objects in writing to the designated IEEE office within 60 days.
 - 8.1.3 Procedure on Petitions – When a petition for a proposed amendment is submitted, the BoG shall prepare a summary statement and a recommendation for or against adopting the amendment. Summary statement and recommendation require a two-thirds majority vote of all the voting members of the BoG. The petition, summary statement, and recommendation shall be subject to approval by the IEEE TAB. After approval by TAB, the proposed amendment shall be published in the Society magazine, or directly mailed to the membership. The amendment becomes effective unless one percent (or more) of the membership objects in writing to the designated IEEE office within 60 days.
 - 8.1.4 Objections – If one percent objects, a ballot with the proposed amendment shall be mailed to all voting members of the Society. A return date of at least 60 days shall be allowed. Proposed amendments require a two-thirds majority of the returned ballots for approval.
 - 8.1.5 Amendments – Changes shall become effective 60 days after all necessary approvals and notifications.
- 8.2 Bylaws – Revisions and proposed amendments to the Bylaws shall be approved by a two-thirds vote of BoG members in attendance, with a quorum present. After approval, the amendment shall be published in the Society magazine or directly mailed to the membership.

NEWLY APPROVED AMENDMENTS TO THE IEEE COMSOC BYLAWS

In June 2015, the ComSoc Board of Governors approved a major revision of ComSoc Bylaws. The IEEE approved the revision in July 2015, and the revised Bylaws are now in force.

The main changes in this revision can be summarized as follows:

- Alignment with IEEE governing documents.
- Alignment with ComSoc Constitution, which also was revised in June 2015 – see this issue page 13.

- A change in the organization of the Board of Governors, including changes to the description of various Boards and Committees.

- Various clarifications and language improvements.

The revised Bylaws now in force and the previous ones can be found here:

<http://www.comsoc.org/about/documents/bylaws>

IEEE COMMUNICATIONS SOCIETY BYLAWS (JULY 2015)

TABLE OF CONTENTS – ARTICLES

1. Objectives
2. Membership
3. Officers and Operations
4. Councils
5. Technical Committees and Special Interest Groups
6. Boards
7. Standing Committees
8. Ad Hoc Committees
9. Society Representatives
10. Budget and Finance

ARTICLE 1 – OBJECTIVES

- 1.1 Objectives – The objectives of the Society are to provide to its members and the global community of communications professionals the services outlined in Clauses 1.2, 1.3, and 1.4.
- 1.2 Technical Information
 - Creation by research and innovation by the Communications Society community
 - Identification and promotion of hot topics
 - Dissemination worldwide by publications, presentations, and electronic media
 - Exchange by Chapter activities, workshops, discussions, mutual assessments, general networking on technical subjects, and other means of professional communication
 - Facilitation of standards activities
- 1.3 Education (basic and continuing)
 - Tutorials, short courses, lecture programs
 - Chapter support and other delivery mechanisms
- 1.4 Professional Services
 - Personal career growth by providing technical and personal development information
 - Job opportunity benefits through inter-personal networking and facilitation of interactions among members
 - IEEE programs

ARTICLE 2 – MEMBERSHIP

- 2.1 Society membership eligibility and grades are defined in the Constitution.
- 2.2 Unless otherwise stated, IEEE Members (including Graduate Student Members) are entitled to all rights and privileges of the Society except that Student Members are not entitled to hold office or vote.

- 2.3 A Society member who is delinquent in paying Society dues shall be dropped from membership according to IEEE procedures. A former member may reinstate membership upon payment of current dues.

ARTICLE 3 – OFFICERS AND OPERATIONS

- 3.1 All officers (except the Society Executive Director) who are members of the Board of Governors (BoG), Councils, Boards, and Standing and Ad Hoc Committees or are Technical Committee Chairs and Society Representatives shall be Members of the Society. The President-Elect and Vice Presidents (VPs) shall be Senior Members or Fellows of IEEE. Officers who do not meet these member requirements shall not be eligible for Society officer positions as defined in these Bylaws.

3.2 Elected Officers

- 3.2.1 President-Elect, Vice President-Technical and Educational Activities (VP-TEA), Vice President-Publications (VP-PUB), Vice President-Conferences (VP-CON), Vice President-Member and Global Activities (VP-MGA), Vice President-Industry and Standards Activities (VP-ISA), IEEE Division III Director-Elect, and Members-at-Large of the BoG are elected by direct vote of the voting Members of the Society (see Article 2.2).

- President-Elect shall be elected in even-numbered years, and Vice Presidents in odd-numbered years.
- One-third of the total (12) Members-at-Large shall be elected annually.
- When an elected officer is elected to another Society position, except President-Elect, during his or her term, he/she shall resign from the former position, upon taking office.
- If the President-Elect holds more than one voting position in the BoG or Operating Committee (OpCom), he/she will have no more than one vote.
- When an elected officer is elected to the position of President-Elect during his/her term, he/she may continue holding the earlier position through the conclusion of its term or upon entering the position of President, whichever comes first.

3.2.2 Terms of Office

- President-Elect shall serve a one-year term the year following his/her election (odd-numbered) and begin a two-year term as President the following year (even-numbered), and then continue for a one-year term (odd-numbered) as Past President.
- Vice Presidents shall serve two-year terms beginning in the even-numbered year following their election.
- Members-at-Large shall serve a three year term beginning the year following their election.

3.2.3 Eligibility for Re-election

- The President-Elect shall not be re-elected President-Elect for more than one term, consecutive or otherwise.
- Vice Presidents may be re-elected to the same office for a second consecutive two-year term, but are further ineligible for that office until the lapse of one year.
- A member shall be ineligible for a Vice President-level position after being elected for a total of any five vice-presidential terms, consecutive or otherwise.
- Members-at-Large may be re-elected as Members-at-Large for a second consecutive term, but are further ineligible for that office until the lapse of one year.

3.2.4 Absence or Incapacity of:

- President – Duties shall be performed by the President-Elect (odd-numbered years)/Past President (even-numbered years) and then by the Vice President-Technical and Educational Activities, Vice President-Publications, Vice President-Conferences, Vice President-Member and Global Activities, and Vice President-Industry and Standards Activities, in that order.
- President-Elect – The term shall be filled by the Past President, who shall continue in that capacity until a special election is held and a new President-Elect is chosen.
- Vice President – Individuals shall be identified from the appropriate candidate group slate, in the sequence of the number of votes received, and the individual receiving the most number of votes shall be automatically appointed to serve the remainder of the elected term.
- Member-at-Large – Individuals shall be identified from the same regional slate as the candidate being replaced, in the sequence of the number of votes received, and the individual receiving the most number of votes shall be automatically appointed to serve the remainder of the elected term.

If none of these individuals can serve, the vacancy shall be filled by action of the BoG upon proposal by the President; a person filling a position in this manner shall serve the remainder of the elected term.

- 3.2.5 Vacancies and Removal from Office. A Society Officer elected by the voting members of the Society may be removed from office, with or without cause, by a vote of the voting members of the Society within thirty days following either:

- The affirmative vote of two-thirds of the votes of the members of the Board of Governors present at the time of the vote, provided a quorum is present, on a motion to remove such individual, or
- Receipt by IEEE of a petition signed by at least 10% of the total number of voting members of the Society moving for the removal of such individual.

A ballot on such motion shall be submitted to the voting member of the Society. If a majority of the ballots cast by the voting members for or against such motion are to remove such individual, the individual shall be removed from such positions.

The vacant office shall be filled as defined in the governing documents of the Society.

3.3 Appointed Officers

- 3.3.1 The appointment of Treasurer, CIO, Parliamentarian, and all Directors shall be proposed by the President-Elect in consultation with the VP-Elect (if any) with whom the position is associated, and approved by the outgoing BoG at the last meeting of an odd-numbered year.
- 3.3.2 Unless indicated otherwise in the Bylaws, the Chairs of Standing Committees shall be proposed by the President-Elect in consultation with the VP-Elect (if any) with whom the position is associated, and approved by the outgoing BoG at the last meeting of an odd-numbered year.
- 3.3.3 The appointment of Chairs of Ad Hoc Committees rests with the authority that established the Ad Hoc Committee.
- 3.3.4 All Society appointed Officers shall serve for the nominal term of the President.
- 3.3.5 The Society appointed Officers are:
 - Treasurer – Is responsible for assuring sound financial practices, establishing prudent budgetary policies, overseeing preparation and presentation of the Society’s budget and working with IEEE on financial matters.
 - Chief Information Officer (CIO) – Oversees cost-effective planning, acquisition, maintenance and use of the Society’s information systems and networking, databases and telecommunications services.
 - Directors – Each chairs a Board and serves in the Council under which the Board is aligned.
 - Regional Directors shall be appointed by the President from lists containing at least two candidates from each region submitted by the respective Regional Board before December 15 of odd-numbered years. If the respective Board does not submit its list by this deadline the President shall make the appointment in consultation with the incoming Vice President- Member and Global Activities (VP-MGA)
 - Parliamentarian – Advises the President on rules of order and proper procedures during BoG meetings. The President may, in case of conflict, request a ruling on procedures from the Parliamentarian.
 - Standing and Ad Hoc Committee Chairs

3.4 President, President-Elect, and Past President-Responsibilities.

- 3.4.1 In the odd-numbered years, the President-Elect shall assist the President in discharging the responsibilities of that office. In the even-numbered years, the Past President shall assist the President in taking on the responsibilities of that office.
- 3.4.2 During his/her term, the President-Elect shall start selecting officers to appoint upon taking office as President in consultation with the appropriate VP-Elect.
- 3.4.3 The President shall inform the BoG of the roster of all Boards and Standing Committees as soon as they are finalized and by the first BoG meeting of his term.
- 3.4.4 The President is the highest ranking volunteer officer of the Society. He/she is responsible for leading the implementation of strategic actions and directions set by the ComSoc BoG. The President or his/her delegate represents the Society in negotiations with Sister Societies and other similar organizations.
- 3.4.5 The President shall oversee and coordinate handling of ethics and conduct issues involving members, including author misconduct, at the Society level. The President shall be assisted by the Executive Director and by volunteers with experience in such matters, as needed.

3.5 Vice Presidents – Responsibilities

- 3.5.1 Vice Presidents are accountable to the President for their areas of responsibility.
- 3.5.2 Vice President – Technical and Educational Activities is responsible for all technical activities and educational services within the Society. Reporting to this Vice President are:
 - Director – Educational Services
 - Director – Technical Services
 - Chair – Awards Committee
 - Chair – Communications History
 - Chair – Emerging Technologies Committee
 - Chair – Fellow Evaluation Committee
 - Chair – Technical Committees Recertification Committee
 - Chair – GLOBECOM/ICC Technical Content Committee
- 3.5.3 Vice President – Publications is responsible for all activities of the Society related to print and electronic products, such as journals, magazines and online offerings. Reporting to this Vice President are:
 - Director – Journals
 - Director – Magazines
 - Director – Online Content
- 3.5.4 Vice President – Conferences is responsible for all aspects of technical conferences, workshops, and professional meetings, including conference publications. Reporting to this Vice President are:
 - Director – Conference Development

- Director – Conference Operations
- Chair – GLOBECOM/ICC Management & Strategy Committee
- Chair – GLOBECOM/ICC Technical Content Committee

3.5.5 Vice President – Member and Global Activities is responsible for: a) all activities, services and programs associated with members and chapters, and oriented to membership retention, development and marketing in the four regions; b) all activities related to the organization and management of chapters; and c) relations with other IEEE and professional societies worldwide.

Reporting to this Vice President are:

- Director – Member Services
- Director – Sister & Related Societies
- Director – Asia/Pacific Region
- Director – Europe, Middle-East & Africa Region
- Director – Latin America Region
- Director – North America Region
- Chair – Women in Communications Engineering Standing Committee

3.5.6 Vice President – Industry and Standards Activities shall be responsible for overseeing all Society activities and programs related to all standards activities and industry services within the Society, including: (i) fostering technical activities related to relevant current standards development and industry services; (ii) identifying opportunities and fostering ComSoc’s engagement in new and/or existing standards development projects that are under development by different standards development organizations worldwide; (iii) increasing the visibility of ComSoc industry and standards initiatives within IEEE, the wider international standards community, and the broad international community of communications technologists; (iv) using ComSoc industry and standards activities to forge closer ties with ComSoc’s other departments and activities; (v) maintaining a close and informed relationship with the IEEE-SA; (vi) management within ComSoc, according to IEEE governing documents, ComSoc-sponsored IEEE Standards Association (IEEE-SA) projects, and (vii) fostering and implementing activities that are of interest to industry and government, including practitioners, managers, executives, young professionals and other industry professionals. The Vice President – Industry and Standards Activities shall be the official ComSoc liaison to the IEEE Standards Association Board of Governors.

Reporting to this Vice President are:

- Director – Standards Development
- Director – Standardization Programs Development
- Director – Industry Outreach
- Chair – Industry Content and Exhibition Committee (ICEC)

The Vice President – Industry and Standards Activities shall be an ex-officio voting member of the Standards Development Board, the Standardization Programs Development Board, the Industry Outreach Board, and the ICEC.

3.6 Board of Governors (BoG)

3.6.1 Officers on the Board of Governors:

- Elected (Voting) Officers:
 - President
 - President-Elect (odd-numbered years), immediate Past President (even-numbered years)
 - Vice Presidents
 - Members-at-Large
 - IEEE Division III Director
 - IEEE Division III Director-Elect (odd-numbered years)
- Appointed (Non-voting) Officers:
 - Directors
 - Chief Information Officer
 - Parliamentarian
 - Treasurer
 - Executive Director

3.6.2 The BoG shall hold at least two formal in-person meetings annually. Special BoG meetings may be held at the request of the President or any four members of the BoG. A majority of the voting members of the BoG constitutes a quorum. When a quorum is present, a majority vote is necessary to transact business. Proxy voting is not allowed.

- The presiding officer of the BoG shall have no vote on the BoG except if the vote is by secret ballot or unless the Chair’s vote can change the outcome of the vote.
- The vote of a majority of the votes cast of the members present and entitled to vote at the time of the vote, provided a quorum is present, shall be the act of the Society.

3.6.3 Each year, the President, Vice Presidents, and Executive Director shall submit the coming year’s Operating Plans to the BoG. Progress on these plans shall be reviewed throughout the year by the BoG.

SOCIETY NEWS

- 3.6.4 The BoG may meet and act upon the vote of its members in person, by any means of telecommunications, or by combination thereof. The normal voting requirements shall apply when action is taken whereby all persons participating in the meeting can hear each other and view or access presentations at the same time-
- 3.6.5 For meetings with in-person and remote participants, remote participants who either cannot hear other participants or who are not heard by other meeting participants do not meet the requirements for meeting attendance and, therefore, are not included in quorum calculations or allowed to vote.
- 3.6.6 Business may be conducted by means other than formally held meetings when the matter can be adequately handled via letter, electronic ballot, electronic mail interchange, etc., referencing Society policies and procedures. When transacting business without a meeting, a majority vote of all the BoG members eligible to vote is required for actions so taken.
- 3.7 Operating Committee (OpCom)
- 3.7.1 Between formal and special BoG meetings, business shall be managed by the Operating Committee (OpCom). OpCom is a subset of the BoG and its members are:
- President
 - President-Elect (odd-numbered years), immediate Past President (even-numbered years)
 - Vice Presidents
 - Members-at-Large (three – one from each annually elected group and appointed by the President)
 - IEEE Division III Director/ Director-Elect
 - Directors
 - Chief Information Officer
 - Treasurer
 - Parliamentarian
 - Executive Director
- 3.7.2 Only OpCom members who are voting members of the Board of Governors are eligible to vote.
- 3.7.3 OpCom shall meet twice annually, in person or by other means. Additional meetings may be held at the request of the President or any four members of OpCom. A majority of the voting members of OpCom constitutes a quorum. When a quorum is present, a majority vote is necessary to transact business. Proxy voting is not allowed. Actions of OpCom shall be submitted to the BoG for ratification in a consent agenda or further consideration in its next meeting.
- 3.7.4 All OpCom members are expected to attend OpCom meetings, except for:
- Directors
 - Chief Information Officer
- The President shall determine which of the above members shall be invited to a particular OpCom meeting.
- 3.7.5 Actions from a duly called OpCom meeting are not effective until ratified by the BoG.
- 3.8 Operations
- 3.8.1 Minutes of each BoG and OpCom meeting shall be distributed to the BoG within 30 days of the meeting. For executive sessions, only motions passed shall be included in the BoG and OpCom minutes. Brief executive session minutes shall be kept on file in the office of the ComSoc Executive Director.
- 3.8.2 Members of the BoG and OpCom shall receive notice of their formal meetings no fewer than 21 days prior to the scheduled meeting start date.
- 3.8.3 If a quorum is not present at a duly called BoG or OpCom meeting, actions may be formulated but are not effective until ratified by letter, electronic mail, or conference call. A majority vote of that specific body is required for ratification. Following ratification, approved decisions shall be recorded in the minutes of the meeting that did not have a quorum.
- 3.8.4 Business action that is formulated outside of a duly called meeting may be conducted by letter, electronic mail, or conference call, referencing Society policies and procedures. When transacting business without a meeting, a majority vote of all the assembly members eligible to vote is required for actions so taken. Approved decisions shall be confirmed promptly in writing or by electronic transmission and recorded in the minutes of the next meeting.
- 3.8.5 Business at Society meetings shall be conducted according to Robert's Rules of Order (latest revision) unless other rules and procedures are specified in the Not-for-Profit Corporation Law of the State of New York, the IEEE Certificate of Incorporation and IEEE governing documents.
- 3.8.6 Councils, Boards, and Committees shall hold meetings at the request of the relevant Chair with sufficient frequency to transact Society business with reasonable dispatch.
- 3.8.7 A Management Retreat may be held annually at the discretion of the President.
- 3.8.8 BoG members shall adhere to decisions of the BoG, unless such decisions violate IEEE or Society constitutions, bylaws or policies.

- 3.8.9 The Ombudsman shall be the first point of contact for reporting a dispute or complaint related to Society activities and/or volunteers. The Ombudsman shall investigate, provide direction to the appropriate IEEE resources if necessary, and/or otherwise help settle these disputes at an appropriate level within the Society. Nominations & Elections Committee shall nominate two candidates for the position of Ombudsman who are not currently on the Board and have not been on the Board for at least two years. The BoG shall then select one of the two candidates to serve for a two-year term beginning the second year of the President's term.
- 3.8.10 Constitution, Bylaws, and Policies and Procedures of the Society shall be in accordance with the IEEE Governing documents.
- 3.8.11 Proposed amendments to Society Governance documents should be reviewed by the Governance Committee prior to approval.
- 3.9 Professional Staff
 - 3.9.1 Subject to compliance with all applicable IEEE Bylaws and Policies, the Society may create an Executive Office supported by IEEE staff. The Society's Executive Office functions to coordinate and carry-out the day-to-day operations, policies and procedures concerning all aspects of the Society's business. The Office also maintains corporate memory and provides ongoing and ad hoc management reports/documents. In addition, the Society's Executive Office serves as one of the Society's primary points of contact for both members and IEEE staff.
 - 3.9.2 Subject to compliance with all applicable IEEE Bylaws and Policies, the Society may determine the budget for the Executive Office. The staff is hired by the IEEE and all conditions of employment shall be based upon IEEE Bylaws, staff policies and practices and all applicable laws and regulations. Office organization, job descriptions, IEEE staff policies and employment practices are available from the IEEE Human Resources Department.
 - 3.9.3 The Executive Director is the most senior position on the IEEE staff that supports the Society, and as such, he/she manages and develops, personally and through subordinate management staff, the paid IEEE staff members that support the Society's operations and activities. The Executive director supports the Society President, officers and volunteer leadership to achieve the Society goals. This Executive director reports through the Managing Director, Technical Activities, to the IEEE Executive Director.
 - 3.9.4 The Executive Director serves as BoG/OpCom secretary, assisted by staff members where needed.

ARTICLE 4 – COUNCILS

- 4.1 Councils are chaired by Vice Presidents to address Technical and Educational Activities, Publications, Conferences, Member and Global Activities, and Industry and Standards Activities. Directors reporting to and Standing Committees aligned under a Vice President serve on his/her Council and all are voting members together with the Chair. Vice Presidents may appoint a vice chair and a secretary (voting positions), and additional non-voting members as needed, and shall designate a member to serve as chair pro tempore in their absence.
- 4.2 Council policies and procedures are approved by the Council and with consent of the BoG.
- 4.3 Technical and Educational Activities Council (TEA-C) – This Council is responsible for the educational and technical interests of the Society, encompassing the broad range of communications and communications-related technical areas.
- 4.4 Publications Council (PUB-C) – This Council is responsible for the needs of the Society and Society Members related to print and electronic projects, such as journals, magazines, and online offerings, not including conference publications.
- 4.5 Conferences Council (CON-C) – This Council is responsible for the needs of the Society and Society Members related to technical conferences, workshops, and professional meetings.
- 4.6 Member and Global Activities Council (MGA-C) – This Council is responsible for all Society activities and programs related to members, chapters, membership development, sister and related societies, and Society regions.
- 4.7 Industry and Standards Activities Council (ISA-C) – This Council is responsible for the needs of the Society and Society members related to industry and standards. Additional (non-voting) members of the Council include the Chairs of any Standards Committees that are established and report to the Industry and Standards Activities Council Directors.

ARTICLE 5 – TECHNICAL COMMITTEES AND SPECIAL INTEREST GROUPS

- 5.1 Technical Committees (TCs)
 - 5.1.1 Technical Committees are established to promote and achieve the technical objectives of the Society and report to the Director – Technical Services.
 - 5.1.2 Technical Committees may be created, merged, or dissolved by resolution of the BoG. A resolution to create or substantially change a Technical Committee may be submitted by the VP-TEA and shall include the name, scope, tentative program for the first year, and approximate numbers of interested and potential members.
 - 5.1.3 New Technical Committees may also be proposed by petition of 25 Society Members; petitions shall be submitted to the BoG by the VP-TEA and include information detailed in 5.1.2.
 - 5.1.4 The Chair of a new Technical Committee is appointed for two years by the Director - Technical Services with the approval of the VP-TEA. During this period, a mentor is assigned to the committee by the Director of Technical Services. Subsequently, the Chair shall be elected by members of the Technical Committee.
 - 5.1.5 Elections for Technical Committee Chairs are held every two years for a two-year term. A Chair cannot serve more than two consecutive terms of office.

SOCIETY NEWS

- 5.1.6 Technical Committees shall have Policies and Procedures (P&Ps) which shall include officer positions and election procedures. P&Ps shall conform to the template specified in the Society P&Ps which may be modified with the approval of the Director of Technical Services.
- 5.1.7 Each Technical Committee shall have a technical scope that may be modified when appropriate, upon approval of the TEA-C and the BoG.
- 5.2 Special Interest Groups (SIGs)
 - 5.2.1 Special Interest Groups are established to cover substantial and diverse topical areas of current industry interests and report to the Chair of the Industry Content and Exhibition Committee (ICEC).
 - 5.2.2 Special Interest Groups may be created, merged, or dissolved by the Chair of ICEC.
 - 5.2.3 The creation, operation, and dissolution of SIGs are governed by the charter of the ICEC.
 - 5.2.4 The charter of the ICEC and any changes to it shall be approved by the VP-ISA and the BoG.
 - 5.2.5 Special Interest Groups shall have Policies and Procedures (P&Ps) which shall include officer positions and election procedures. P&Ps shall conform to the template specified in the Society P&Ps which may be modified by a Special Interest Group with the approval of the Chair of ICEC.
 - 5.2.6 Each Special Interest Group shall have a scope that may be modified, upon approval of the Chair of ICEC.

ARTICLE 6 – BOARDS

- 6.1 Boards are the operational and strategic entities of their respective Councils and are chaired by Directors. The Boards and the Council under which they are organized are listed below:
 - Conferences Council (CON-C)
 - Conference Development Board
 - Conference Operations Board
 - Member and Global Activities Council (MGA-C)
 - Member Services Board
 - Sister & Related Societies Board
 - AP Region Board
 - EMEA Region Board
 - LA Region Board
 - NA Region Board
 - Publications Council (PUB-C)
 - Journals Board
 - Magazines Board
 - Online Content Board
 - Industry and Standards Activities Council (ISA-C)
 - Industry Outreach Board
 - Standards Development Board
 - Standardization Programs Development Board
 - Technical and Educational Activities Council (TEA-C)
 - Educational Services Board
 - Technical Services Board
- 6.2 Directors are responsible for appointing members to their Boards, with the approval of the appropriate Vice President.

All Board members, including the Director, serve two-year terms concurrent with the nominal duration of the presidential term.

Each Director may appoint a secretary and additional non-voting members as needed. In addition to the voting members specified in the “Board Descriptions” clause, each Board may approve voting rights for additional members.
- 6.3 Policies and Procedures for each Board shall be developed by the Board and approved by the appropriate Council under which the Board is aligned. In the case that the council does not approve the P&Ps and a compromise cannot be found, the board may request the BoG to resolve the matter and approve the P&Ps. An exception exists for the Standards Development Board Policies and Procedures and those of the Standards Committees reporting to the Standards Development Board which shall be approved by the IEEE-SA Standards Board.
- 6.4 Board Descriptions
 - 6.4.1 Conference Development – This Board is responsible for the strategic planning, technical scope, and growth of all ComSoc financially-sponsored conferences (defined as portfolio conferences).

Members include representatives from the TEA-C and at least four Members-at-Large with at least one having served as the technical program chair and at least one as the general chair of a major conference. All are voting members, in addition to the Director.

SOCIETY NEWS

- 6.4.2 Conference Operations – This Board is responsible for the oversight and management of the operational, publications, and financial aspects of all ComSoc conferences.
- Members include the ComSoc Treasurer and at least five Members-at-Large with at least two having served as the General Chair and at least one as the Technical Program Chair of a major ComSoc conference. All are voting members, in addition to the Director.
- 6.4.3 Educational Services – This Board is responsible for the oversight of all Society education and training activities, including administration of the Society’s programs on continuing education, incorporating tutorials, short courses, lectures, etc. In particular, this Board is responsible for developing and maintaining continuing professional education and training programs.
- Members include representatives from the Technical Services, Conference Development, Conference Operations, Industry Outreach and Member Services Boards, ICEC, and at least two Members-at-Large. All are voting members, in addition to the Director.
- 6.4.4 Industry Outreach – This Board is responsible for assuring a comprehensive and cost-effective outreach program of Society products and services to industry and governmental communities. It is also responsible for developing liaisons with communications and networking related enterprises to promote ComSoc products and services and to attract industry and government leaders into ComSoc’s volunteer community.
- Members include a representative from each of the following: the Technical Services Board, the Educational Services Board, the PUB-C, the CON-C, the Standardization Programs Development Board, the MGA-C, and the ICEC. In addition, up to four members at large may be appointed by the Director to represent external industrial and governmental interests. All are voting members, in addition to the Director.
- 6.4.5 Journals – This Board is responsible for the oversight of Society journals. Voting members are the Editors-in-Chief of Society journals for which ComSoc is the Managing Partner and/or has a majority financial stake, and two Members-at-Large, in addition to the Director. Additional members can be appointed, including Liaison Editors to other IEEE journals. All of these are voting members, in addition to the Director.
- 6.4.6 Magazines – This Board is responsible for the oversight of Society magazines. Voting members consist of the Editors-in-Chief of Society magazines and two Members-at-Large, in addition to the Director. Additional members can be appointed, including Liaison Editors to other IEEE magazines. All of these are voting members, in addition to the Director.
- 6.4.7 Member Services – This Board is responsible for the oversight of all services and programs addressed to members and chapters, and oriented to membership retention and development in the four regions. In particular, this Board is responsible for developing and providing individual-level membership services globally.
- Members include the four Regional Directors plus one member per Region selected by the Director from a list of candidates, consisting of at least two names per region, submitted by each Regional Director. All are voting members, in addition to the Director.
- 6.4.8 Online Content – This Board is responsible for initiating, assessing and overseeing Society online content. It supports technical committee activities; online services; as well as publications, conferences, and education products and services.
- Members include representatives from the Conference Development, Conference Operations, Educational Services, Journals, and Magazines Boards; the TEA, MGA, and ISA Councils; the CIO; and up to three additional members. All are voting members, in addition to the Director.
- 6.4.9 Regional Boards – These Boards are responsible for stimulating, coordinating and promoting the activities of ComSoc members and chapters throughout the IEEE regions. The four regions, each with its own Board, are:
- Asia/Pacific (AP)
 - Europe, Middle-East & Africa (EMEA)
 - Latin America (LA)
 - North America (NA)
- Each Board shall have a minimum of five voting members, in addition to the Director.
- 6.4.10 Sister & Related Societies – This Board is responsible for enhancing Society activities with our sister and related societies (SRS) by developing new programs with SRS and cooperating with SRS in offering Society/SRS products and services globally, and strengthening the Society’s global and professional reach. This includes establishing and maintaining Society relationships on an international, regional, national or local scale with SRS.
- Where appropriate, enhancing Society SRS activities shall be accomplished through collaboration with IEEE sections/chapters, including Society and non-Society chapters.
- Membership includes up to three representatives from the MGA-C, up to three members from selected Sister and Related Societies, and up to three Members-at-Large. All are voting members, in addition to the Director.
- 6.4.11 Standards Development – This Board is responsible for the promotion and advancement of communications standards.
- It consists of eight members in addition to the Director, all of whom are voting members. The Director shall select the members in accordance with the priorities listed in the next paragraph.

The Director shall give priority to serve on the Standards Development Board to: Standards Committee Chairs, ComSoc appointed Chairs or Co-Chairs for joint Standards Committees, Working Group Chairs who are directly sponsored by the Standards Development Board, ComSoc-appointed Working Group Chairs or Co-Chairs for jointly sponsored Working Groups and volunteers in ComSoc Technical Committees.

The Director shall be the official ComSoc liaison to the IEEE Standards Association Standards Board (SASB).

- 6.4.12 Standardization Programs Development – This Board is responsible for launching pre-and post-standardization technical activities, not restricted to those standards being developed by the IEEE. These would include, but not be limited to Research Groups that lead to the discovery of standardization opportunities and, for completed standards, creation of follow-up programs, such as compliance testing, standards education, workshops, conferences, and publications on technical issues that are relevant to standards.

The Board shall consist of up to eight members in addition to the Director, all of whom are voting members. The Director of the Standardization Programs Development Board serves as the liaison to the IEEE-SA Industry Connections Program.

- 6.4.13 Technical Services – This Board is responsible for the oversight and promotion of the technical communities of the Society and their activities including the promotion of technical content and development of educational content. In particular, this Board is responsible for developing and providing one-stop ICT (information-communications technology) services. This Board is also responsible for managing the Distinguished Lecturers selection process.

Members include Technical Committee Chairs, Chair of the Emerging Technologies Committee, Chair of the Technical Committees Recertification Committee, and a representative of the Educational Services Board. All are voting members, in addition to the Director.

- 6.5 Other Boards may be created by a two-thirds majority vote of the BoG. The scope, responsibility, and Policies and Procedures shall be defined before incorporating the new Board into the Bylaws.

ARTICLE 7 – STANDING COMMITTEES

- 7.1 Standing Committees of the Society and their alignments under Society officers are:

• Finance	President
• Governance	President
• IEEE/ComSoc Coordination	President
• Nominations & Elections	President
• Operations & Facilities	President
• Strategic Planning	President
• GLOBECOM/ICC Management & Strategy	VP-CON
• GLOBECOM/ICC Technical Content	VP-CON/VP-TEA
• Women in Communications Engineering	VP-MGA
• Awards	VP-TEA
• Communications History	VP-TEA
• Distinguished Lecturers' Selection	VP-TEA
• Emerging Technologies	VP-TEA
• Fellow Evaluation	VP-TEA
• Technical Committees Recertification	VP-TEA
• Industry Content & Exhibition	VP-ISA

- 7.2 Unless otherwise specified in the Standing Committee description, Standing Committee Chairs shall be responsible for appointing Standing Committee members with the approval of the Society Officer under which the Standing Committee is aligned.

Unless otherwise specified in the Standing Committee description, Standing Committee members, including the Chair, are voting members and shall serve a two-year term concurrent with the nominal duration of the presidential term.

The Chair may appoint a secretary and additional non-voting members as needed. In addition to the voting members specified in the “Standing Committee Descriptions” clause, each Standing Committee may approve voting rights for additional members.

Standing Committees may have additional or fewer members as determined by the President on an exception basis and for a designated term.

- 7.3 A Standing Committee is aligned under the Council chaired by the Vice President to whom the Standing Committee Chair reports. Policies and Procedures of Standing Committees aligned under a Council shall be developed by the Standing Committee and approved by the Council under which the Standing Committee is aligned. In the case that the council does not approve the P&Ps and a compromise cannot be found, the Standing Committee may request the BoG to resolve the conflict and approve the P&Ps. Policies and Procedures for Standing Committees whose Chair reports to the President shall be developed by the Standing Committee and approved by the BoG.
- 7.4 Standing Committee Descriptions
- 7.4.1 Awards – This committee is responsible for all major awards and recognitions made or proposed by the Society. It consists of not less than twelve (12) members who shall serve for a three-year term. One-third of the members are appointed each year. Committee members may not provide nominations or reference letters while in office, nor participate in deliberations on awards or recognitions for which they may be under consideration.
- 7.4.2 Communications History – This committee is responsible for identifying, placing in electronic archives, and raising public awareness through all appropriate steps on the most important facts/person/achievements of communications history in particular, as well as telecommunication milestones in general. The committee consists of three members who shall serve a three-year term, with one member appointed each year.
- 7.4.3 Distinguished Lecturers Selection – This committee is responsible for establishing selection criteria and for the appointment of lecturers.
- The ex-officio Chair of this Standing Committee is the Vice Chair of the TEA-C. Voting members of this committee consist of the Chair, the VP-TEA, VP-MGA, the Director-Member Services, and the Chair of the Emerging Technologies Committee.
- 7.4.4 Emerging Technologies – This committee is responsible for identifying, describing, and nurturing new technology directions, recommending new programs, and nurturing potential technical committees for formal proposal via the VP-TEA.
- The Chair is appointed by the President from the members of the Strategic Planning Committee with the recommendation of the VP-TEA. Standing Committee members shall include at least one more member from the Strategic Planning Committee. The committee shall have six members appointed for three years with one-third appointed each year. In addition, the Editor-in-Chief of *IEEE Communications Magazine* and the Editor-in-Chief of *IEEE Journal of Selected Areas in Communications* are ex-officio voting members of the committee.
- 7.4.5 Fellow Evaluation – This committee is responsible for the Society’s evaluation of Fellow nominations being considered by the IEEE Fellow Committee. It consists of a Chair and nine members that shall serve a three-year term with one-third of the members being appointed each year. Chair and members must be IEEE Fellows and Members of the Society.
- 7.4.6 Finance – This committee is responsible for facilitating the Society’s budget process and for managing and providing direction in all aspects of Society financial matters. The committee meets twice a year at ICC/GLOBECOM. The committee is chaired by the Treasurer and includes the President, Past or President-Elect, Vice Presidents, CIO, Executive Director, and a representative from each MAL class.
- 7.4.7 GLOBECOM/ICC Management and Strategy (GIMS) – This committee is responsible for the successful conduct, strategic evolution, and policies of the IEEE Global Communications Conference (GLOBECOM) and the IEEE International Conference on Communications (ICC). The committee structure and operation is detailed in the GIMS charter. The voting members of the GIMS committee are: a Chair, three or four Members-at-Large, three past members of an ICC or GLOBECOM Organizing Committee, and the GITC committee Chair.
- 7.4.8 GLOBECOM/ICC Technical Content (GITC) – This committee is responsible for providing strategic vision and management of the technical content of GLOBECOM and ICC to guarantee timeliness and the highest level of quality. The committee structure and operation is detailed in the GITC charter. The voting members of the GITC committee are: a Chair, two to four members identified by the TEA-C, three past ICC or GLOBECOM Technical Program Chairs, and the GIMS committee Chair.
- 7.4.9 Governance – This committee is responsible for all matters related to Society Governance, including but not limited to: reviewing any proposed amendment to Society Governance documents (Constitution, Bylaws, Policies and Procedures) prior to its discussion in the BoG; crafting amendments to Society Governance documents that result from actions of the BoG; establishing Society-wide Governance best practices and overseeing their application across all Councils, Boards, and Committees; upon request or when needed proposing changes to existing Society Governance documents with the goal of keeping them current and consistent; and serving as an interpretive body on Governance issues. Voting membership: Chair appointed by the President, the Parliamentarian, and up to three additional members appointed by the President upon recommendation of the Committee Chair and shall include one previous ComSoc President and one sitting BoG Member at Large. The committee Chair may appoint up to three additional voting committee members. The Chair and committee members (except the Parliamentarian) shall serve three-year terms with one reappointment allowed. Terms of all members (except the Parliamentarian) shall be staggered so that no more than half of the members’ terms expire every two years; when necessary, such staggering shall be created by appointing members to terms shorter than three years, as indicated at the time of their appointment.
- 7.4.10 IEEE/ComSoc Coordination – This committee is responsible for the Society’s internal coordination and cooperation with IEEE entities and for enhancing Society relations with IEEE governance and staff. Voting members consist of the President, VP-MGA, Past President/President Elect, IEEE Division III Director, IEEE Division III Past Director/Director-Elect and Director- Sister & Related Societies.

7.4.11 Industry Content & Exhibition – This committee is responsible for developing and promoting a strategic vision and oversight for organizing and promoting internal ComSoc communities that are attractive to members from industry, government, or other non-academic sectors. This includes processes to assure the quality and value of content in industry oriented conferences, events and education. It is within the overall objective and mission of the ICEC to increase industry participation in ComSoc events.

Voting members of this committee are the Chair and 4-to-6 Members at Large.

7.4.12 Nominations & Elections – This committee is responsible for identifying candidates to fill elected Society office positions, and for the development, implementation and supervision of election procedures. It consists of nine voting members, each appointed by the President, upon recommendation by the Chair, and approved by the BoG for a three-year term, one-third of the members being appointed each year. At least one annual appointee must not be a member of the BoG. The tenth voting member is the current IEEE Division III Director. Immediately after being elected and officially announced, the President-Elect becomes an ex-officio non-voting committee member.

The committee is chaired by Past Presidents. Each President spends the first year after his/her term of office, e.g., as Past President, as an ex officio voting committee member (the eleventh voting member in odd years). This Past President assumes the Chair the following year (an odd-numbered year) and serves for two years. If the Past President is unable to serve, the President shall appoint a Chair for a similar period, with the approval of the BoG.

The Chair shall not be eligible to be elected to the BoG during his/her term of service. A committee member may be nominated for a position only if (i) the nomination is not made by a member of the committee, and (ii) the member resigns from the committee prior to its first meeting of the year in which the nomination shall be made. Meeting of this committee shall always be held in Executive Session. The number of signatures required for a petition candidate to appear on a ComSoc ballot shall be equal to the number required by the IEEE Bylaw on Nominations & Appointments. All candidates for appearance on the ballot must be identified at least 28 days prior to the date of the election. The operation of the committee is detailed in its policies and procedures.

7.4.13 Operations & Facilities – This committee is responsible for supporting the President in making recommendations to the BoG on operations, facilities and related capital expenses. It is chaired by the President and includes the following voting members who serve for the duration of the President's term:

- CIO
- Treasurer
- Four members appointed by the President from among volunteer BoG officers
- One member appointed by the President from among volunteers who are not on the BoG

The Executive Director is an ex-officio non-voting member of this committee. The Committee may approve the participation of invited experts as required by the agenda items. Meetings of this committee shall always be held in Executive Session.

7.4.14 Strategic Planning – This committee is responsible for preparing a long-term strategic plan to guide the direction and future of the Society and for preparing short-term plans to direct specific areas, as appropriate.

Voting members of the committee shall be the Chair, the Vice Presidents (or a representative named by the VP), and up to four Members-at-Large, all appointed by the President.

7.4.15 Technical Committees Recertification – This committee recommends the establishment of new Technical Committees and reviews current committees to determine whether they are fulfilling their responsibilities.

The committee is chaired by the VP-TEA and includes as voting members the TEA-C Vice Chair and six members appointed by the Chair for a three-year term. One-third of the members are appointed each year by the Chair from among Members-at-Large of the BoG. Committee members who are also officers of Technical Committees under review shall excuse themselves from deliberations related to their Technical Committee.

7.4.16 Women in Communications Engineering – This committee is responsible for encouraging the participation and membership of women communications engineers in the Society. The Chair is appointed by the President to a two-year term, and the appointment can be extended of one additional term. The Chair reports to the Vice President-Member and Global Activities (VP-MGA). Voting members are appointed by the Chair to a two-year term upon approval of the Vice President-Member and Global Activities (VP-MGA), and their appointment can be extended of one additional term. The committee will meet at least once a year at ICC or GLOBECOM. It will provide an annual written report, which will be distributed to the Society President, Vice Presidents, and Technical Committee Chairs prior to each ICC.

7.5 Other Standing Committees may be created by a two-thirds majority vote of the BoG. The scope, responsibilities, and Policies and Procedures shall be defined before incorporating the new Standing Committee into the Bylaws.

ARTICLE 8 – AD HOC COMMITTEES

8.1 Ad Hoc Committees may be established by the President or by the BoG to address broad technical or operational issues within the Society or IEEE.

- 8.2 Voting members shall be appointed by the authority that established the Ad Hoc Committee. The President shall report to the BoG the composition, mission, and duration of the Ad Hoc Committees established under his/her authority.
- 8.3 OpCom shall review Ad Hoc Committees annually and recommend to the BoG whether they should continue, disband, or be elevated to Standing Committees.
- 8.4 Policies and Procedures of Ad Hoc Committees shall be developed by the Ad Hoc Committee and approved by the authority that established the Ad Hoc Committee.

ARTICLE 9 – SOCIETY REPRESENTATIVES

- 9.1 Representatives are responsible for representing the Society in other IEEE organizations and other non-IEEE Societies and organizations. They are appointed by the President for terms as required by the other organization, in consultation with the appropriate Vice President.

ARTICLE 10 – BUDGET AND FINANCE

- 10.1 Officers shall prepare budgets for the coming calendar year in the first half of each year to be approved by the BoG at its mid-year meeting. Actuals shall be reviewed throughout the year, and a forecast reported at each meeting.
- 10.2 Dues and fees are set by the BoG in accordance with IEEE and Society guidelines and are based upon proposals by the Treasurer to the BoG. Billing and receipt of annual dues are part of the IEEE dues billing process.
- 10.3 Budget
 - 10.3.1 Each year the Society produces a budget which shall be approved by the BoG.
 - 10.3.2 The Treasurer is responsible for the development of the Society annual budget and submitting to IEEE Technical Activities for consolidation with other societies, and ultimately to the IEEE for their consolidated budget. The Treasurer monitors revenues and expenses, providing interim reports on budgets, forecast, actuals at each BoG and OpCom meeting. A complete financial report, including actual versus budget, net assets, and reserves is presented by the Treasurer annually.
- 10.4 Finance
 - 10.4.1 The Treasurer has oversight responsibility for all Society financial matters.
 - 10.4.2 Funds shall be handled as designated by the Treasurer and shall be deposited with IEEE or with external financial institutions, as approved by the BoG and/or IEEE Board of Directors.
 - 10.4.3 The Treasurer, or Society Executive Director, or their designee shall follow orderly procedures for disbursement of funds, providing sufficient checks and balances and appropriate record keeping. A budgeted expenditure requires no further approval beyond approval of the Treasurer.
 - 10.4.4 The Treasurer shall periodically review the Society finances and recommend adjustments needed to insure financial stability of the Society.
 - 10.4.5 The Treasurer shall cooperate with Society and IEEE officials to accomplish financial audits when requested. The results of these audits shall be presented to the BoG.

THE FIRST IEEE COMMUNICATIONS SOCIETY SUMMER SCHOOL TRENTO, ITALY, JULY 6-9, 2015

BY FABRIZIO GRANELLI, DISI – UNIVERSITY OF TRENTO

The city of Trento, the educational, scientific, financial and political centre of Northern Italy, served as the ideal launching point for the first IEEE Communications Society (ComSoc) Summer School. With a population of more than 100,000 residing in this noted university city, Trento is recognized throughout the country for its high quality of life, prosperous business opportunities, advanced research centers and renowned international cultural institutions. It also boasts a rich heritage of cultural and social lifestyles; enriched international conferences, meetings and exhibitions; traditional costume festivals and market fairs; historical events; and global festivals highlighted with theatre, music and dance.

Held July 6-9, the first IEEE ComSoc Summer School provided participants with top-level lectures on hot topics in communications as well as a myriad of networking opportunities. Targeting IEEE ComSoc PhD student members, the event 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 very happy and honored that Trento was chosen for the first edition of this event," says Prof. Fabrizio Granelli, the venue's local organizer. "We did our best to select the most promising students and then provide them with an exciting and unprecedented program in addition to the world communications authorities as speakers. We were deeply impressed by the positive and constructive attitude of the participants, who were all well prepared and interactive."

Selected from more than 100 applicants during the Spring of 2015, the Summer School hosted 43 participants from worldwide locations at the University of Trento, one of the youngest Italian universities (founded in the 1960s) and in conjunction with the Dept. of Information Engineering and Computer Science, the 2nd highest ranked ICT department in Italy.

On July 6, the conference opened with the introductions of the Rector of the University of Trento, Prof. Paolo Collini; Head of the Dept. of Information Engineering and Computer Science (host and co-founder of the initiative), Prof. Gian Pietro Picco; the ComSoc VP Member Relations, Prof. Stefano Bregni; and the Head of the ComSoc Education Board Training Working Group, Prof. Fabrizio Granelli.

Afterwards, the event's program began with an interactive lecture on "Collaborative Near-Capacity Wireless System Design" by Prof. Lajos Hanzo, University of Southampton, UK. Aimed at future wireless communication systems, this talk first discussed the limitations of MIMOs relying on co-located array-elements and then showed how single-antenna-aided cooperative mobiles can circumvent these limitations through the formation of MIMOs with distributed elements - a concept also referred to as Virtual



Participants of ComSoc's first Summer School program.

Antenna Arrays (VAA). The corresponding amplify-forward and decode-forward protocols as well as their hybrids were studied in relation to channel coding, which has to be specifically designed for the VAAs in order to prevent avalanche-like error-propagation. Hence, sophisticated three-stage-concatenated iterative channel coding schemes were identified, arguing that in the absence of accurate channel information at the relays the best way forward might be to use multiple-symbol differential detection. Indeed, it is rather unrealistic to expect that an altruistically relaying handset would also accurately estimate the source-relay channel for the sake of high-integrity coherent detection. EXIT-chart-aided designs were then introduced for creating near-capacity solutions in addition to a range of future research directions and open problems.

On July 7, Prof. Giuseppe Bianchi from the University of Rome – Tor Vergata, Italy 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 OpenState, a very recent proposal devised to permit platform-agnostic programmability of stateful tasks (formally described in terms of abstract state machines such as Mealy Machines). It also focused on architecture and implementation issues (including backward compatibility with OpenFlow commodity Hardware), as well as application examples (forwarding consistency, MAC learning, reverse path forwarding, fault recovery, etc).

On Tuesday afternoon, the participants had the chance to present their ongoing and future plans during a 2-hours poster session. The session was open to all students of the local doctoral school as well as to the professors and researchers of the ICT department.

The following day, Prof. Andrea Goldsmith from Stanford

University, USA, presented a four-hour lecture on “The Next Wave in Wireless Communications.” It started by discussing the fundamentals of wireless communications, i.e. channel models and impairments, capacity adaptive techniques, diversity, MIMO, OFDM, Spread Spectrum techniques and proceeded with an extensive discussion on multiuser systems, a survey of current cellular systems and a review of future expectations.

The speech then focused on ad hoc wireless networks, which are expected to gain additional relevance in scenarios of device-to-device communications and Internet of Things. However, the discussion underlined the need for a more flexible and efficient access technology for the scarce communication resources, paving the way for 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 on 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 the communications theory (e.g. neural science).

On the last day, Prof. Nelson L.S. da Fonseca of the State

University of Campinas, SP, Brazil, presented a seminar on “Networking for Big Data,” which included a discussion on its ecosystem, perspective in Society and Science and processing capabilities. Prof. Fonseca then spoke about network virtualization as a fundamental pillar to networking support for the Big Data area as well as research opportunities not only on networking for Big Data but also on Big Data for networking. Also provided were lab exercises for students to practice on their way back home.

During the course of the four-day event, the IEEE ComSoc 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.

“It was a challenging and exciting event. We are very proud of the outcome, and hope it will become an integral component of the yearly events organized by the IEEE ComSoc for its members,” concludes Prof. Granelli.

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

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 detailed structured data needed by the orchestrators and controllers in the ecosystem.

This Feature Topic addresses industries’ standards 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 approaches to API development is also invited. Solicited topics include (but are not limited to):

- Information modeling
- Data modeling
- Transforming information models to data models
- Service development life cycle 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 *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 *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/com-mag-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

IEEE ICC 2015 SETS ATTENDANCE RECORD IN LONDON, UNITED KINGDOM

NEARLY 3,000 INDUSTRY PROFESSIONALS, RESEARCHERS & SCIENTISTS ATTEND

MORE THAN 1,800 PRESENTATIONS DEDICATED TO LATEST “SMART CITY & SMART WORLD” ADVANCEMENTS

The 2015 IEEE International Conference on Communications (IEEE ICC 2015) recently concluded its global annual event with a record-breaking 2947 attendees participating in 1800 presentations dedicated to nearly every area of the communications spectrum.

“This year’s conference was a phenomenal success, which is a tribute to the efforts of our many organizers, patrons and staff,” says Executive Chair Professor Jiangzhou Wang of the University of Kent, United Kingdom. “Once again, everyone joined together to create an unsurpassed educational and networking platform designed to facilitate innovations worldwide as well as professional relationships and career opportunities that will span a lifetime.

“This was all performed within walking distance of some of the world’s most iconic attractions and London’s Tech City, the fastest growing technology cluster in Europe. As one of the most technologically advanced cities in the world, London was also the ideal backdrop for spurring the introduction of new revolutionary technologies and our next wave of telecom leaders. Set among a dense population of hi-tech offices, smart startups and long-established international Internet giants, we set our highest attendance record ever for an IEEE ICC event. This was also after receiving the most paper submissions (3,342) in the conference’s 50-year history. We are deeply appreciative for everyone’s contributions to this landmark achievement and unparalleled support of the communications industry and the professionals working to change lives every day.”

Held from June 8–12, 2015 at the ExCel London Convention Centre in London, IEEE ICC 2015 began Monday, June 8 with the first of two full days of tutorials and workshops addressing topics like “Device-to-Device Communication for Cellular and Wireless Networks,” “Greening Cloud Networks,” “Dynamic Social Networks” and “Energy Harvesting and Energy Cooperation in Wireless Communications.” For instance, after a brief introduction to 5G small cells, the Workshop on “Backhauling 5G Small Cells: Challenges and Potential Approaches” illustrated full-duplexing and massive MIMO approaches to multi-tier RAN architecture challenges, while the “Next Generation Backhaul/Fronthaul Networking and Communications” workshop accented the newest mmWave techniques for increasing the density and capacity of fifth generation (5G) cellular access systems.

Later that night, the conference officially commenced with the Annual Welcome Reception and grand opening of the conference exhibit hall. There, numerous industry giants such as Qualcomm, Huawei, Nokia, GENBAND, Keysight, National Instruments, InterDigital, P.I. Works, IEEE Big Data, Airvana, SoliD, Imperial College London, IEEE Green ICT, DIGILE, Fore-Mont, Anritsu, Artech House, Telefonica, Nutaq, Rodhe & Schwartz, Springer, Cambridge University Press, Wiley, and Cambridge demonstrated their wide breadth of products and services for the hundreds of attendees touring the banquet floor.

After brief introductions, the IEEE ICC 2015 keynote



agenda highlighted Tuesday morning with the presentations of Dr. Paul E. Jacobs, Executive Chairman of Qualcomm Incorporated, who discussed the “Mobile-Powered Future” and the creation of intelligently connected mobile ecosystems, and Professor H. Vincent Poor of Princeton

University, who spoke on “Smart Grid: The Role of the Information Sciences” and the integration of renewable energy sources.

During his speech, Dr. Jacobs highlighted his view for “transforming the edge of the Internet” with solutions that overcome the 1000× data challenge, create the connectivity fabric for everyone, and bring cognitive technology to life. This included aggregating the licensed and unlicensed spectrum to deliver greater performance, while utilizing LTE-U/LAA to deliver new levels of on-device intelligence and integration. He also foresaw the inevitability of “making cities smarter” with emergency alerts, vehicle to cloud communications, road sensors and wireless vehicle charging through devices operating with a “digital 6th sense,” cognitive connectivity and seamless and preemptive security.

Dr. Poor then began his address by describing the Smart Grid as the “The Internet of Energy” filled with pervasive control, self-monitoring and healing features, two-way communication, automated maintenance, and increased consumer choices. For this, he outlined a “cyber-physical approach” using information science and techniques like game and information theories as well as statistical inference. As an example, he highlighted a model for energy trading that includes plug-in electric cars exchanging energy with the main grid and the resulting decline in utility prices as more “players enter the game.”

At the Annual Awards Luncheon that afternoon, IEEE ComSoc President Sergio Bendetto honored this year’s IEEE Fellows, who were:

- Jean Armstrong — For contributions to the theory and application of orthogonal frequency-division multiplexing in wireless and optical communications
- Azzedine Boukerche — For contributions to communica-



2015 IEEE Fellows were honored at the annual Awards Ceremony.

tion protocols for distributed mobile computing and wireless sensor networks

- Iain Collings — For contributions to multiple user and multiple antenna wireless communication systems
- Xiqi Gao — For contributions to broadband wireless communications and multirate signal processing
- Monisha Ghosh — For contributions to cognitive radio and signal processing for communication systems
- Ekram Hossain — For contributions to spectrum management and resource allocation in cognitive and cellular radio networks
- Markus Rupp — For contributions to adaptive filters and communication technologies
- Aylin Yener — For contributions to wireless communication theory and wireless information security
- Amitava Ghosh — For leadership in cellular communication system standardization

After Paul E. Jacobs was presented the Distinguished Industry Leader Award for his leadership in the field of mobile communications and his contribution to the commercialization of mobile technology breakthroughs that have significantly contributed to the growth of the industry, IEEE ComSoc Awards Chair Lajos Hanzo announced the 2015 Prize Paper Award recipients as:

- David Gesbert, Stephen Hanly, Howard Huang, Shlomo Shamai Shitz, Osvaldo Simeone, and Wei Yu, who received the Best Tutorial Paper Award for “Multi-Cell MIMO Cooperative Networks: A New Look at Interference,” *IEEE Journal on Selected Areas in Communications*, vol. 28, no. 9, December 2010, pp. 1380–1408.

- Jakob Hoydis, Stephan ten Brink, and Mérouane Debbah, who were given the Leonard G. Abraham Prize for “Massive MIMO in the UL/DL of Cellular Networks: How Many Antennas Do We Need?” *IEEE Journal on Selected Areas in Communications*, vol. 31, no. 2, February 2013, pp. 160–171.

- Mehdi Bennis, Meryem Simsek, Andreas Czylwik, Walid Saad, Stefan Valentin, and Merouane Debbah, who were presented the Fred W. Ellersick Prize for “When Cellular Meets WiFi in Wireless Small Cell Networks,” *IEEE Communications Magazine*, vol. 51, no. 6, June 2013, pp. 44–50.

- Rui Zhang and Chin Keong Ho, who received the Marconi Prize Paper Award in Wireless Communications for “MIMO Broadcasting for Simultaneous Wireless Information and Power Transfer,” *IEEE Transactions on Wireless Communications*, vol. 12, no. 5, May 2013, pp. 1989–2001.

- Cem U. Saraydar, Narayan B. Mandayam, and David J. Goodman, who were awarded the Advances in Communication Recognition for “Efficient Power Control via Pricing in Wireless Data Networks,” *IEEE Transactions on Communications*, Vol. 50, No. 2, February 2002, pp. 291–303.

- David Hillerkuss, Rene Schmogrow, Matthias Meyer, Stefan Wolf, Meinert Jordan, Philipp Kleinow, Nicole Lindemann, Philipp C. Schindler, Argishti Melikyan, Xin Yang, Shalva Ben-Ezra, Bend Nebendahl, Michael Dreschmann, Joachim Meyer, Francesca Parmigiani, Periklis Petropoulos, Bojan Resan, Andreas Oehler, Kurt Weingarten, Lars Altenhain, Tobias Ellermeyer, Michael Moeller, Michael Huebner, Juergen Becker, Christian Koos, Wolfgang Freude, and Juerg Leuthold, who were given the IEEE Communications Society Charles Kao Award for Best Optical Communications & Networking Paper “Single-laser 32.5 Tbit/s Nyquist WDM Transmission,” *Journal of Optical Communications and Networking*, vol. 4, no. 10, October 2012, pp. 715–723.

- Ezio Biglieri and Marco Lops, who were honored with the *Journal of Communications and Networks (JCN)* Best Paper Award sponsored by the Korean Information & Com-



Attendees enjoyed the London Fanfare Band at the Welcome Reception.

munications Society (KICS) and technically cosponsored by IEEE Communications Society (ComSoc) for “Linear–Quadratic Detectors for Spectrum Sensing,” *Journal of Communications and Networks*, vol. 16, no. 5, October 2014, pp. 485–492.

- Fenyao Bao, Ing-Ray Chen, MoonJeong Chang, and Jin-Hee Cho, who received the William R. Bennett Prize for “Hierarchical Trust Management for Wireless Sensor Networks and Its Applications to Trust-Based Routing and Intrusion Detection,” *IEEE Transactions on Network and Service Management*, vol. 9, no. 2, June 2012, pp. 169–183.

- Hien Quoc Ngo, Erik G. Larsson, and Thomas L. Marzetta, who were given the Stephen O. Rice Prize for “Energy and Spectral Efficiency of Very Large Multiuser MIMO Systems,” *IEEE Transactions on Communications*, vol. 61, no. 4, April 2013, pp. 1436–1449.

- Tiangao Gou and Syed A. Jafar, who were presented the Heinrich Hertz Award for Best Communications Letter for “Optimal Use of Current and Outdated Channel State Information — Degrees of Freedom of the MISO BC with Mixed CSIT,” *IEEE Communications Letters*, vol. 16, no. 7, July 2012, pp. 1084–1087.

- Harpreet S. Dhillon and Radha Krishna Ganti, who were given the Young Author Best Paper Award for “Modeling and Analysis of K-Tier Downlink Heterogeneous Cellular Networks,” *IEEE Journal on Selected Areas in Communications*, vol. 30, no. 3, April 2012, pp. 550–560.

In addition, Tuesday initiated the introduction of a wide array of forums, panels, symposia, and demonstrations detailing the latest advances in communications technologies and regulatory policy innovations. This was earmarked by nearly 20 panel discussions chaired by representatives of Huawei, Cisco, Ericsson, Microsoft, InterDigital, Samsung, P.I. Works, the European Commission, and National Instruments, among others. These authorities highlighted the newest research and innovations in areas ranging from 5G Challenges and Opportunities, Ensuring the Long-Term Sustainability of the Internet, and Smart City and Sustainable Ecology to CAP Theorem Challenges and DevOps Approaches, Cellular IoT and Electronic Healthcare, IoT, and Telemedicine and Cloud: Standards, Challenges, and Opportunities.

For the first time, IEEE ICC 2015 offered “Presentations on the Podium” as brief discussions hosted by Anritsu, National Instruments, GENBAND, Keysight, Rohde & Schwarz, Solid Inc., Digile, InterDigital, P.I. Works, Nutaq, and Nokia. Dedicated to the introduction of new innovations



Daily keynote sessions captivated attendees.

and technologies, these short vignettes explored topics like understanding and improving the mobile device user experience, WebRTC un-wiring the communications network, rapid 5G systems prototyping, resolving in-building mobile coverage challenges with DAS, DIGILE IoT Program, oneTRANSPORT: Reshaping the Transport Sector with oneM2M, Next Generation Centralised SON, and Looking Ahead to 5G.

Other milestones included the conference's technical program that spanned Tuesday through Thursday and featured 1285 peer-reviewed papers delivered across 69 oral and interactive sessions. Among the hundreds of symposia topics were Big Data Security and Privacy, Adaptive Video Streaming, Compressed Sensing, Software, Mobile Social Networks, Energy Harvesting, CrowdSensing and Mobile IoT Solutions, Satellite Networking, e-Health, IPv6 and IoT Network Protocols, Elastic Optical Networks, and P2P and Opportunistic Communications.

Another prominent feature unveiled on Tuesday was the demonstrations of Nokia, Huawei, Keysight, InterDigital, GenBand, Anritsu, Digile, P.I. Works, Nutaq, and National Instruments on topics like the "Finnish Internet of Things," "Sub 6 GHz MIMO," "5G mmWave MIMO Channel Sounding," "WiFi Data Offloading," and "M2M Connected Vehicles."

Wednesday started with the keynotes of Dr. Wen Tong, CTO, Huawei Wireless, who focused on "5G to Embrace the Vertical Industries," and Professor Alwyn Seeds of University College London, who addressed "Wireless over Fibre Systems: From MHz to THz." First, Dr. Tong confronted the "Dawn of the Automated Society" and its role in transforming wireless to "redefine everything." As a result, by 2025 5G will surpass the employment of hundreds of billions of M2M devices and hundreds of millions of new real-time services operating with massive, mission-critical connectivity. Enabled by 5G ultra-narrowband waveform technologies, 5G will also be earmarked by the use of 1 MHz of spectrum for generating 900 billion sensor connections, a new class of industrial small cells enabling flexible design automation and a 5G air interface design offering "zero" latency across five million global macro sites.

Dr. Seeds followed this presentation by citing how the world's IP data traffic will exceed 130 exabytes a month by 2018, and mobile data usage is expected to grow three times faster than fixed IP traffic between 2013 and 2018. He then outlined multiple techniques for satisfying these rampant data demands, which include the use of multi-level modulation formats or MIMO at frequencies lower than 100 GHz, free space optical communications, and THz-over fiber systems.

According to Dr. Seeds, wireless over fiber is a proven method for getting base stations close to users at affordable

system costs, while "achieving ultra-high bit rate transmissions without the fog outages that affect free-space optical systems." Another benefit is the transfer of signals over optical fiber, enabling integration with fiber access systems. However, challenges that must still be overcome to ensure efficient and reliable THz wireless communications include the systemic introduction and implementation of high power photomixers, low noise amplifiers for receivers and active array antennas for mobile device tracking. Other areas concern the advance of techniques that will achieve increased transmitter output power; target 10 mW to 100 mW.

Immediately following the keynote session, attendees were feted at the Chief Technology Officer Forum highlighted by the viewpoints of EE CTO Fotis Karonis; Telefónica UK CTO Brendan O'Reilly; Huawei wireless CTO Wen Tong; and Qualcomm wireless standards EVP Ed Tiedemann. Confronted with the title question, "How Far Can We Evolve Mobile Networks — What Is Next?," each of the participants urged far greater collaboration among the industry's players to ensure 5G not only delivers higher quality service, but also overcomes significant data traffic among their mobile networks.

Moderated by Three UK CTO Bryn Jones, O'Reilly was insistent about the need for ubiquitous coverage, citing the 2012 London Olympics as an example of the cooperation necessary to accomplish instantaneous communications on a global scale. Others highlighted better interference management as a necessity for improving spectral efficiency as Dr. Tong stressed the need to "get back to basics" and prioritize spectral efficiency. Also lamented was the reluctance of consumers to pay more for more services, the rising cost of infrastructure and dwindling margins, as well as the continued failure of vendors to either fully understand industry challenges or generate solutions.

The final day of the conference's comprehensive technical agenda began on Thursday with the addresses of Professor Xiaohu You of Southeast University, Professor Lajos Hanzo of the University of Southampton, and Dr. Jürgen Schindler, head of the 5G Business Program at Nokia. Discussing "5G Mobile Communications in China" and the nation's strategies for creating massive cooperative cloud radio capabilities, Dr. You outlined the Chinese National High-Tech R&D Program encompassing 863 open platform and enabling technology projects. This includes the development of 1 Gb/s experiences, seamless wide-area coverage, 100 percent reliability, and ultra-low cost and power consumption. Highlighted were also the presentation of a 5G wireless technology roadmap consisting of new air interfaces (AI) and 5G evolution AI exploiting low-frequency bands below 6 GHz and high-frequency bands within the 6–100 GHz range by 2019.

Professor Hanzo then followed by detailing "A Stroll with Shannon to Next-Generation Plaza: Large-Scale MIMOs, Single versus Multiple RF Chains and All That..." During his talk, Dr. Hanzo "strolled" with Shannon to review numerous present-day scenarios exploring questions like: "Would the field of wireless have developed equally bandwidth consciously?" "What if governments had not imposed frequency-license fees?" "What about green radio and the tactile Internet?" The relevancy of each query was then explored in relation to current mm-wave and optical wireless, frequency-reuse factor, SNR pathloss and fading signal interference, and MIMO capacity applications with the goal of conceiving cooperative massive MIMO-aided unlicensed and optical wireless HetNets.

After Professor Hanzo's address, IEEE ICC 2015 Technical Program Chairs Nathan Gomes and Athanassios Manikas recognized the Best Paper Award winners for this year's



Attendees engaged one another during the interactive sessions.

event. Those cited for these honors were:

- Yezi Huang, Thomas Magesacher, and Eduardo Medeiros (Lund University, Sweden), Chenguang Lu and Per-Erik Eriksson (Ericsson Research, Sweden), and Per Ödling (Lund University, Sweden) for their paper “Mitigating Disorderly Leaving Events in G.fast.”
- Limei Guo (Central South University, USA), Hsiao-Chun Wu (Louisiana State University, USA), Yiyan Wu (Communications Research Centre, Canada), and Xian Liu (University of Arkansas, Little Rock, USA) for their entry on “Optimal Total-Downlink-Transmitting-Power and Subchannel Allocation for Green Cellular Networks.”
- Ting Wu (NYU Polytechnic School of Engineering, USA), Theodore Rappaport (New York University and NYU WIRELESS, USA), and Christopher Collins (New York University, USA) for “The Human Body and Millimeter-Wave Wireless Communication Systems: Interactions and Implications.”
- Pierre Coucheney, Kinda Khawam (Université de Versailles, France), and Johanne Cohen (LRI-CNRS & PRISM-CNRS, France) for “Multi-Armed Bandit for Distributed Inter-Cell Interference Coordination.”
- Ying Cui (Shanghai Jiaotong University, China), Muriel Médard (MIT, USA), Edmund Yeh (Northeastern University, USA), Douglas Leith (Hamilton, Ireland), and Ken R. Duffy (National University of Ireland Maynooth, Ireland) for “Optimization-Based Linear Network Coding for General Connections of Continuous Flows.”
- Jingya Li (Chalmers University of Technology, Sweden), Emil Björnson (Linköping University, Sweden), Tommy Svensson and Thomas Eriksson (Chalmers University of Technology, Sweden), and Mérouane Debbah (Supelec, France) for “Optimal Design of Energy-Efficient HetNets: Joint Pre-coding and Load Balancing.”
- Xiaomin Chen and Admela Jukan (Technische Universität Carolo-Wilhelmina zu Braunschweig, Germany), and Muriel Médard (MIT, USA) for “Linear Network Coding and Parallel Transmission Increase Fault Tolerance and Optical Reach.”
- Sai Qian Zhang (University of Toronto, Canada), Qi Zhang (University of Toronto and University of Waterloo,

Canada), and Hadi Bannazadeh and Alberto Leon-Garcia (University of Toronto, Canada) for “Network Function Virtualization Enabled Multicast Routing on SDN.”

- Sang-Woon Jeon (Andong National University, Korea), SongNam Hong (Ericsson Research, USA), Mingyue Ji (University of Southern California, USA), and Giuseppe Caire (Technische Universität Berlin, Germany) for “Caching in Wireless Multihop Device-to-Device Networks.”
- Shoujiang Ma, Daoyun Hu, Shengru Li, Nana Xue, Suoheng Li, Yan Shao, and Zuqing Zhu (University of Science and Technology of China, China) for “QoS-Aware Flexible Traffic Engineering with OpenFlow-Assisted Agile IP-Forwarding Interchanging.”
- Lili Wei and Geng Wu (Intel Corporation, USA), and Rose Qingyang Hu (Utah State University, USA) for “Sum-Capacity Optimal Spread-Spectrum Data Hiding in Video Streams.”

• Yun Liao, Tianyu Wang, Kaigui Bian, and Lingyang Song (Peking University, China) and Zhu Han (University of Houston, USA) for “Decentralized Dynamic Spectrum Access in Full-Duplex Cognitive Radio Networks.”

The event’s keynote program concluded with Dr. Jürgen Schindler, head of the 5G Business Program at Nokia, who spoke about “5G, Expanding the Human Possibilities of Technology” and “expanding the human possibilities of cellular technology.” In this design, everyone is an innovator operating within an ecosystem characterized by zero road fatalities, high industrial productivity, and healthier people living within safe connected homes. As explained by Dr. Schindler, a nationwide spectrum globally harmonized to enable guaranteed QoS and global roaming will be key to creating this scenario highlighted by “5G any to any connectivity,” simultaneous and native HetNet multi-connectivity, and session-on-demand resource efficiency for the sporadic data transmission of low-cost and low ARPU devices.

The final day of IEEE ICC 2015 then featured another full day of tutorials and workshops on Friday, June 12. Among the many subjects highlighted were Optical Wireless Communications, The Path Towards 5G, Cloud Radio Access Networks, Android Security, Cognitive Radios and Networks for Spectrum Coexistence, Green Communications and Networks with Energy Harvesting, Smart Grids and Renewable Energies, Dependable Vehicular Communications (DVC), and Security and Privacy for Internet of Things and Cyber-Physical Systems.

For additional information on IEEE ICC 2015, please visit <http://www.ieee-icc.org/2015>. Visitors are also invited to network with colleagues and peers as well as share their professional experiences through the conference Facebook, LinkedIn, and Twitter pages.

In addition, the planning for IEEE ICC 2016 to be held May 23–27 in Kuala Lumpur is currently underway. Expected to host thousands of global communications experts at the globally benchmarked Kuala Lumpur Convention Centre, the Call for Papers deadline has been set as October 16 for original submissions to this premier global venue located near some of the world’s largest ICT industries, and amidst this rising hub of economic and social innovations. All interested professionals are invited to visit <http://www.ieee-icc.org/2016> for ongoing conference updates and detailed submission information.

CONFERENCE CALENDAR

Updated on the Communications Society's Web Site
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>

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 ICOS 2015 — IEEE Int'l. Conference on Space Optical Systems and Applications, 27–28 Oct.

New Orleans, LA
<http://icos2015.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-iccc.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/>

IEEE NFV-SDN 2015 — IEEE Conference on Network Function Virtualization and Software Defined Networks, 18–21 Nov.

San Francisco, CA
<http://www.ieee-nfvdsn.org/>

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.dspcs-witsp.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.iitm.ac.in/coconet2015/index.html>

–Communications Society portfolio events appear in bold colored print.

–Communications Society technically co-sponsored conferences appear in black italic 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; fax: + (212) 705-8996. Items submitted for publication will be included on a space-available basis.



August 2015

ISSN 2374-1082

REGIONAL REPORT

IEEE Young Professionals, Lahore Section Win World Class Award "Hall of Fame 2015"

By Usman Munawar , IEEE Lahore Section, Pakistan

The IEEE Young Professionals Lahore Section, after winning the IEEE Hall of Fame Award 2011, the IEEE YP Region 10 Award 2012, and the MGA Individual Award 2012, has won the IEEE Hall of Fame Award 2015. Today our Young Professionals, Lahore Section list of active members has reached more than 50. More professionals are coming in to join our cause of serving the society. IEEE YP strives to bring reforms to the engineering sector in our country with valuable efforts and approaches in the engineering profession. The Hall of Fame 2015 Award is a tribute to our seniors for such achievements for Pakistan and Region 10.

SUMMARY OF THE EVENTS HELD IN 2014

IEEE PSYWC 2014: IEEE Pakistan Student, YP (Young Professionals), and WIE (Women In Engineering) Congress (PSYWC) was the largest event for the students and young professionals in the three IEEE sections (Lahore, Islamabad, and Karachi). The three-day event united 450 students, young professionals, and women in engineering, computer science, and allied fields to explore ideas, develop skills, and discuss issues in their profession. Students had a chance to interact with highly professional academia and industry personnel.

IEEE Technovate 2014: IEEE Technovate 2014 was a collaboration of IEEE Young Professionals Lahore, IEEE Lahore Section, with the IEEE team IEEE from Egypt that treads on the path paved by the Egyptian Engineering Day (EED) program, which is celebrated annually in Egypt (Region 8). The aim of this event was to select the best projects of Pakistan so those involved could be sponsored to fly to Egypt to compete on the next Regional level with the best projects of the entire region on the prestigious Egypt Engineering Day (EED).

Brainiac 2014: BRAINIAC 2014 was organized in collaboration with ICOSST 2014, a perfect mix of technical and non-technical competitions, seminars, conferences, and workshops. This three-day exciting venture consisted of a Business Idea competition, an Android App Development Competition, Extreme Programming, a Web Design Competition, a Pseudo Code Competition, and Circuit Mania.

Formation and promotion of new chapters, affinity groups, and student branches: A new milestone was achieved by establishing an IEEE PES Student Chapter for the first time in Pakistan by IEEE YP. The PES Young Professionals Committee is now dedicated to better serving student members of PES and assisting PES Young Professional members in their transition to corporate life. In addition to this, the IEEE YP Lahore Section played its vital role in introducing three technical Chapters (IEEE PES, IEEE IAS, and IEEE ComSoc), an Affinity Group (IEEE Consultant Network),



Members of the IEEE Young Professionals Lahore Section with the 2015 IEEE Hall of Fame Award.

five Student Branches (IEEE IIUI, IEEE BZU, IEEE SCET, IEEE UOG, and a fifth in process, IEEE CIIT Sahiwal, and seven Technical Society Student Chapters (four IAS Student Chapters (UET, LCWU, IUB and GCUF), and three PES Chapters (NFC, CIIT, and LCWU). IEEE Young professionals kept its eyes on the promotion of IEEE Societies to empower and update knowledge of IEEE members in the section, which has proven very valuable for engineers and society itself.

IEEE Day Celebrations UMT: An IEEE Student Branch was established at the University of Management and Technology on this special occasion. The event proved to be a greatly successful membership campaign. Various IEEE volunteers and student representatives highlighted the advantages of being an IEEE member and how IEEE provides countless opportunities in technology.

Major Achievements: Adding to the list of mega events are the 8th IEEE ICOSST, IEEE STEP, IEEE SPAC, IEEE NWPE, IEEE Fetex 2015, IEEE Final Year Project Exhibition 2014, and 6TH IEEE AEPEX 2014, and industrial visits (PTCL and KICS).

Recipe for IEEE YP Success: IEEE Young Professionals of Pakistan is an international community of enthusiastic, dynamic, and innovative members and volunteers. IEEE is committed to helping young professionals and graduates in order to develop their skill set. IEEE YP members are volunteers who help young graduates and professionals in recognizing their career path and goals and providing them with consultation and career counseling. IEEE YP plans and manages all the events with one common and main objective of professional development in Pakistan. IEEE YP held more than 20 Board meetings and general meetings in order to brainstorm events and out of the box ideas and to bring them to reality.

IEEE YP Team: IEEE YP includes more than 20 members who have been designated with multiple projects. Our Executive Committee includes the following professionals

- Dr. Amjad Hussain, Mentor
- Mr. Maroof Raza, Mentor
- Mr. Ijlal Haider, Mentor
- Usman Muhammad Ali, Chair, YP; Lecturer, University of Sargodha
- Usman Munawar, Secretary, YP; Research Office, University of Engineering and Technology, Lahore, Pakistan

(Continued on Newsletter page 4)

Distinguished Lecturer Tour of Nei Kato in China, May 2015

By Waheed ur Rehman, Beijing ComSoc Chapter Secretary, Weixiao Meng, Harbin ComSoc Chapter Chair, Liu Jia Jia, Xian ComSoc Chapter, China

Prof. Nei Kato was invited by the Xian, Beijing, and Harbin ComSoc Chapters to conduct a ComSoc Distinguished Lecturer Tour in May 2015. The first lecture was held at Xidian University on 15 May, which was hosted by the Xian ComSoc Chapter. The lecture topic was "Device-to-Device (D2D): Research Trends and Future Perspectives." There were more than 40 attendees including faculty members and graduate students from Xidian University. The attendees were very interested in the lecture topic and they actively participated in the Q&A session, asking questions ranging from current research progress of D2D to future open problems. They also raised practical issues regarding the achievable performance and implementation details of the "Relay-by-Smartphone" D2D prototype.

Prof. Kato delivered his second lecture at Beijing University of Posts and Telecommunications (BUPT) on 18 May 2015. It was hosted by the Beijing ComSoc Chapter. There were approximately 50 attendees, including university students, academic researchers, and public participants from ICT related industries. Prof. Kato started the lecture with background information about D2D and why we need it. The technical presentation covered a range of interesting details about D2D and the recent research development in this field. Prof. Kato also showed a short video about interesting research conducted in his lab, where the students showed how to send a message up to 2.5 km without any infrastructure. There were many different insights and raised curiosities in the participants' minds, which was then reflected during the Q & A session. The lecture was about an hour. Participants were very interested in the topic and contents. After the lecture, Prof. Kato had additional discussions with the other professors at BUPT, and he visited the related labs.



After the lecture at the Xi'an Chapter.



Prof. Nei Kato gave a technical presentation at the Harbin Institute of Technology.

After the lectures in Beijing, Prof. Kato traveled to Harbin to conduct his technical talk at Harbin Institute of Technology on 21 May 2015. The lecture was hosted by the Harbin ComSoc Chapter. Prof. Kato analyzed the background and significance of D2D, proposed that in the case of spectrum scarcity of resources, it is possible to increase the cellular communication system spectral efficiency and reduce terminal transmit power for D2D to a cer-

(Continued on Newsletter page 4)

Distinguished Lecturer Tour of Tom Hou to Beijing and Nanjing, China

By Waheed ur Rehman, IEEE Beijing ComSoc Chapter Secretary, and Prof. Lianfeng Shen, IEEE Nanjing ComSoc Chapter Chair, China

Dr. Tom Hou, IEEE Fellow, was invited by the IEEE Beijing Communications Society Chapter to give a 90-minutes lecture on "Advances in Wireless Networking for Cyber Physical Systems" on 6 May 2015 (Wednesday). The lecture was held at Beijing University of Posts and Telecommunications (BUPT) from 3 pm to 4.30 pm. There were approximately 30 attendees, including university students, academic researchers, and public participants from ICT related industries.

The lecture began at the preliminary level for the need and significance of cyber and physical systems (CPS). Dr. Hou's lecture focused on two main areas. First, he discussed the throughput in multi-hop wireless networks, which is a fundamental problem in CPS. Second, he considered the energy problem in sensor networks, which is another fundamental problem. Finally, Dr. Hou talked about the new frontiers in sharing radio spectrum and co-existence between primary and secondary networks beyond dynamic spectrum access.



Lecture at National Mobile Communications Research Lab, Southeast University.

Participants were very interested in the topic and contents, and they actively asked their questions during the Q&A session, with the questions covering both the basic and professional aspects of the topic. Due to the shortage of time, several participants had further exchanges with Dr. Hou after the session ended.

It was a pleasure for the Beijing ComSoc Chapter to host Dr. Hou, and the lecture brought many significant ideas to the research students.

After the technical talks in Beijing, Dr. Hou was invited to give

(Continued on Newsletter page 4)

European Wireless 2015 in Budapest

By Hassan Charaf, Frank Fitzek, Leonardo Badia, Marcos Katz and László Lengyel, Hungary

The 21st European Wireless (EW) Conference took place in Budapest, Hungary, 20-22 May 2015, and was organized by the Budapest University of Technology and Economics (BME) with the cooperation of VDE-ITG, represented by Klaus D. Kohrt, Germany, and HTE, represented by Peter Nagy, Hungary, both of them Sister Societies of ComSoc. BME is the most significant University of Technology in Hungary and is also one of the oldest Institutes of Technology in the world, having been founded in 1782. The 2015 edition of EW aimed at addressing the key theme of "5G and Beyond," with a focus on the latest trends, developments, and of course the future applications on top of the mobile and wireless communications. On the first day six tutorials, organized by Sergio Palazzo and Morten V. Pedersen, were offered in three tracks covering two different topics on 5G, Device-to-Device, IoT, network coding, and network programming for Android platforms.

The Technical Program of the conference consisted of 86 papers from more than 30 countries selected after a detailed review process, and was enriched by three keynote talks by Eric Dahlmann, Muriel Médard, and Tommaso Melodia. Eric



Thanks for Muriel Médard after the keynote.

Dahlmann from Ericsson explained the future of 5G and the evolution from previous cellular systems. Muriel Médard spoke about the interaction between network coding and the physical layer. Tommaso Melodia gave an interesting speech entitled "Toward Ultrasonic Networking for Implantable Intra-body Networks."

EW2015 featured two workshops (organized by Christian Wietfeld, László Lengyel, and Péter Ekler), with the topics "Device-to-Device Communication for 5G Systems (WD2DC)" and "Coding Techniques for 5G Networks."

Continuing with the tradition of the European Wireless conference, the scope of EW15 was broad, including, in addition to the subjects of the main conference theme, the most relevant topics in modern communications. This was one of the reasons why most conference attendees, more than 100 colleagues, were present continuously during the conference.

Accepted papers (all available on IEEE Xplore) were organized into 16 technical sessions.

The next EW conference will be held in Oulu, Finland on 18-20 May 2016. The theme for European Wireless 2016 is "5G: Applications, Businesses and Technologies." More information can be found at <http://ew2016.european-wireless.org/>.



The organizing Committee during the Conference dinner.

Novel In-Band Full-Duplex Communication Prototype for 5G Systems, Finland, May 2015

By Mikko Valkama, Professor at Tampere University of Technology

As 5G radio system research is rapidly developing and is attracting significant global interest, a team of researchers at Tampere University of Technology (TUT), led by Prof. Mikko Valkama, has been developing and building an advanced radio transceiver concept for improving the spectral efficiency of future 5G networks. 5G networks are aiming at very large data rates, very high capacity, very low latency, and very fast response time. The developed prototype is based on the principle of wireless in-band full-duplex communications, which means that transmission and reception are performed at the same time on the same center-frequency. Thus, all of the temporal and spectral resources can be used for both transmitting and receiving data, unlike in the traditional frequency-division or time-division duplexing based schemes. This can as much as double the spectral efficiency, which is a highly tempting feature for the upcoming high-data-rate 5G networks. However, the prevalent technical challenge in



Mikko Valkama (left), Shilpa Talwar, Yang-Seok Choi, Timo Huusari and Dani Korpi demonstrated in-band full-duplex communications at the 2015 Mobile World Congress.

such an in-band full-duplex transceiver is the own transmit signal, which is a powerful source of self-interference for the sensitive receiver circuitry. Thereby, without highly efficient attenuation of the own transmit signal in different parts of the receiver chain and antenna interface, in-band full-duplex communication is not possible.

Recently, in close collaboration with Intel Wireless Labs, the researchers at TUT have managed to greatly alleviate this problem by using various advanced RF/analog and digital self-interference cancellation techniques. With the current version of the

(Continued on Newsletter page 4)

DLT/NEI KATO/Continued from page 2

tain extent. He described the mechanism of DTN proposed by his group. It was suitable for the case when irregular nodes moved in a large scale with faster velocity or nodes were sparse, which could solve the routing failure of MANET caused by nodes' changes and movements to some extent. He also demonstrated a D2D application in the emergency network using a MANET and DTN's mixed strategy that was very effective. During the lecture, professors and students discussed the trends of D2D such as "the security issues of D2D," "the shortest path problem of D2D" and "the validity period of DTN information," among others. Prof. Kato also talked about the research activities on wireless networks in Tohoku University, providing valuable information to the local professors and students. There were approximately 100 attendees at the lecture.

Prof. Kato also attended a seminar organized by Prof. Weixiao Meng's group (Monday Seminar). Prof. Kato patiently provided guidance and directions to students in their research domain, and also had a discussion with the students and other professors, which lead to new ideas in their research work. The participants were impressed with Prof. Kato's rigorous learning style and pleasant personality.

DLT/TOM HOU/Continued from page 2

an IEEE Distinguished Lecture in the China Network Valley, which is supported by IEEE ComSoc Nanjing, on 11 May 2015. This talk was hosted by Prof. Xiaohu You from Southeast University. More than 150 attendees included professors and students from the IEEE ComSoc Nanjing Chapter.

During the lecture, Dr. Hou mainly discussed MIMO's related computational model in multi-hop wireless networks. First, Dr. Hou presented a clear classification of the existing MIMO model, which contains a matrix-based model and DOF-based model. The DOF-based model is divided into a conservative model and an optimistic model. Second, Prof. Hou focused on some new advances related to MIMO's capability in a multi-hop network environment. He presented his research on an optimal computational model, a new enabling technology to be used in the physical layer and the network layer, which combines the advantages

of the conservative model and the optimistic model. Finally, Prof. Tom Hou proposed that most of these research advances were interdisciplinary and required a synergistic exploration of knowledge bases from multiple technical domains for future research.

In the question and answer session, Prof. Hou gave detailed and professional answers to the teachers' questions concerning MIMO modeling, and the students' questions about opportunities for member recruitment. This gave audiences a better understanding of wireless networking and also stimulated their interests in exploring this area.

PROTOTYPE FOR 5G/Continued from page 3

prototype, the own transmit signal can be attenuated by as much as 110 dB, in spite of adopting low-cost RF components, already allowing for true full-duplex operation with low-to-medium transmit power levels. This was also demonstrated at the 2015 Mobile World Congress in early March, where Prof. Valkama and researchers Timo Huisari and Dani Korpi from TUT, together with Dr. Yang-Seok Choi and Dr. Shilpa Talwar from Intel, showcased their full-duplex prototype. The demonstration attracted great interest from the industry and research community alike. Thus, it is clear that these recent advances in full-duplex device implementations show great promise for improving the spectral efficiency of the next generation wireless networks.

In addition to successful research into in-band full-duplex transceivers, Prof. Valkama's group at TUT has also made innovative progress in other fields related to 5G networks, such as 5G radio interface and radio system design, radio localization, transmitter linearization, and cognitive radio. Moreover, his group has been given recognition and praise for their work from the Finnish Funding Agency for Innovation (Tekes) and Academy of Finland.

For further information about the research and demonstrations, please contact the authors (mikko.e.valkama@tut.fi, dani.korpi@tut.fi). Further technical information related to the full-duplex radio research is available at <http://www.tut.fi/full-duplex/>.

HALL OF FAME 2015/Continued from page 1

- Sahar Sultan, Vice Chair, YP; Lecturer, Lahore College for Women, University of Pakistan
 - Faiza Nasir, Treasurer, YP; Zong CMPak
- Below are project managers working in IEEE Young Professionals:
- Ms. Anum Tariq Khan, Project Coordinator Trainings, IEEE beENG Program
 - Mr. Rana Hassan, Project Coordinator Industry, IEEE beENG Program
 - Mr. Ehtasham Khan, Media Manager, YP
 - Mr. Muhammad Haris, Project Manager, IEEE Mentoring Program
 - Ms. Rabail Khizar Malik, Program Drive Lead/Project Manager, IEEE beENG Program
 - Mr. Saad Zafar, Project Manager, Graduate Assistance Program
 - Ms. Javeria Sohail, Documentation & Reporting Manager, IEEE YP Programs
 - Mr. Muhammad, Research and Development Program PM
 - Mr. Usman Rana, Website Manager, IEEE YP Programs
 - Mr. Bilal Javed, SAC Representative, IEEE YP
 - Mr. Saad Rind, Student Coordinator, IEEE YP
 - Mr. Ateeq Azam, Research and Development Program
 - Mr. Hamza Timahim, Project Coordinator, IEEE beENG Program

**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, 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
FAWZI BEHMANN (FAWZL.BEHMANN@GMAIL.COM)
NICOLAE OACA (NICOLAE_OACA@YAHOO.COM)
EWELL TAN, SINGAPORE (EWELL.TAN@IEEE.ORG)



A publication of the
IEEE Communications Society

www.comsoc.org/gcn
ISSN 2374-1082

Immerse yourself in the future.

Get an eye-opening look at tomorrow's mobile world, straight from the visionaries lighting the way.

CTIA Super Mobility 2015 takes you beyond the horizon.

Join a prominent roster of tech luminaries, business leaders and mobile pioneers at CTIA Super Mobility 2015 as they share their visions for tomorrow's mobile technology.

From the latest consumer trends to the next-gen innovations shaping our connected world, the world's foremost authorities map out the advances shaping the future of the mobile industry.



CTIA Super Mobility 2015 Keynote Speakers



Glenn Lurie
President & CEO
AT&T Mobility



Ron Smith
CTIA Chairman
President & CEO



Tom Wheeler
Chairman
Federal Communications
Commission



Bob Pittman
Chairman & CEO



Keynote Host:
Meredith Attwell Baker
President & CEO

CTIA
The Wireless Association®



Richard Maltzberger
Chief Development
Officer & President of
Lowe's International



Marcelo Claure
President & CEO



Robin Thurston
Chief Digital Officer



Marni Walden
EVP & President of
Product Innovation
and New Businesses



CTIA Super Mobility 2015

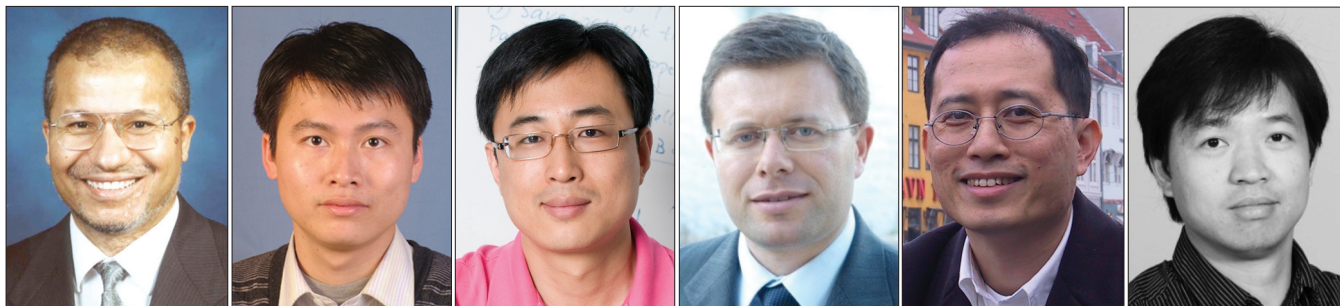
In partnership with Microsoft

SEPTEMBER 9 • 10 • 11 | LAS VEGAS | SANDS EXPO

REGISTER TODAY @ www.CTIASuperMobility2015.com

Use code **IEEE20** to save an additional 20% off the standard rate.

SECURITY AND PRIVACY IN EMERGING NETWORKS: PART II



Mohsen Guizani

Daojing He

Kui Ren

Joel J. P. Rodrigues

Sammy Chan

Yan Zhang

This is the second part of the “Security and Privacy in Emerging Networks” Feature Topic. In Part I, which was published in April 2015, we selected those contributions that dealt with the theory behind the security and privacy of such networks. In Part II, we present articles that overview new security schemes for emerging networks such as vehicular, biomedical, underwater, crowdsourcing, and mobile networks. We feel that even though these emerging networks have attracted many research efforts lately, the security and privacy aspects have not been investigated well. Thus, it is important to provide ways to protect such networks from various security and privacy attacks. The aim of this FT is to promote further research interests in security and privacy in emerging networks by providing a vehicle for researchers and practitioners to discuss research challenges and open issues, and disseminate their latest research results. This can pave the way to implementing emerging networks with the necessary protection from major vulnerabilities. We received a large number of submissions but were obliged to accept only the best 13 papers. Part I was composed of six contributions that dealt with the theory of security/privacy threats, while this issue (Part II) is composed of seven articles addressing security challenges in a specific set of emerging networks.

With wireless technology available in all new vehicles, it is expected that a large amount of information will be exchanged between vehicles and/or between vehicles and roadside units (RSUs). Therefore, malicious attacks (whether intentional or not) may inject untrustworthy information into the network and cause havoc for drivers. This could lead to fatal accidents and loss of lives. Through the first article, “Toward a Trustworthy Vehicular Social Network,” Q. Yang and H. Wang propose a social network approach to study trustworthy information sharing in vehicular networks. They first cover the research progress in measuring direct trust and modeling indirect trust in online social networks. They conclude with a discussion of how to apply those schemes to vehicular social networks and identify some research challenges.

As cloud-assisted wireless body sensor networks (WBSNs) are becoming increasingly popular in healthcare

applications, the security and privacy threats targeting WBSNs deserve more attention. The second article, “Verifiable, Privacy-Assured, and Accurate Biomedical Signal Collection for Cloud-Assisted Wireless Body Sensor Networks,” by C-M. Yu *et al.* focus on data privacy and data completeness where the authors propose a verifiable, privacy-assured, and accurate data collection scheme for cloud-assisted WBSNs. Through both simulation and prototype implementation, they show that the proposed scheme is energy-efficient and effective in protecting data privacy and completeness.

In addition, the use of underwater acoustic sensor networks (UASNs) has increased recently. However, most efforts in this area have not taken network security for UASNs seriously. Typically, UASNs are vulnerable to malicious attacks due to the unique characteristics of an underwater acoustic communication channel (e.g., low communication bandwidth, long propagation delays, and high bit error rates). In addition, the significant differences between UASNs and terrestrial wireless sensor networks (TWSNs) need special attention in the development of secure communication mechanisms for underwater sensor nodes. G. Han *et al.* address these issues in their contribution, “Secure Communication for Underwater Acoustic Sensor Networks.” They present a survey of emerging topics arising from secure communications in UASNs. Then they propose a number of open research problems that, once resolved, could lead to providing secure and efficient methods for UASNs.

As we all know, mobile devices are capable of initiating sophisticated cyber-attacks, especially when they coordinate together to form what is referred to as a mobile distributed botnet (MobiBot). MobiBots leverage the absence of basic mobile operating system security mechanisms and the advantages of device-to-device (D2D) communication in masking malicious code propagation, which make them a serious security threat to any machine/network. In the next article, “From Botnets to MobiBots: A Novel Malicious Communication Paradigm for Mobile Botnets,” the authors investigate the potential and impact of large-scale infection and coordination of neighboring devices. They highlight how mobile devices can leverage short-range

wireless technologies in attacks against other mobile devices that come within proximity. Later, they quantitatively measure the infection and propagation rates within MobiBots using short-range wireless technology such as Bluetooth.

On the other hand, the proliferation of mobile devices such as smartphones has enabled participatory sensing systems that collect data from users through their mobile devices and infer useful information from it. However, users have concerns regarding possible privacy leakage from their data. “Privacy-Preserving Participatory Sensing” by Q. Li and G. Cao addresses how to simultaneously protect privacy and provide incentives for participatory sensing. They review previous approaches, discuss their limitations, and propose two new types of participatory sensing systems with improved privacy protection.

A mobile crowdsourcing network (MCN) is a new promising network architecture that applies the principles of crowdsourcing to perform tasks using powerful mobile devices. However, it also raises some critical security and privacy issues that may prevent the applications and/or implementation of MCNs. The article “Security and Privacy in Mobile Crowdsourcing Networks: Challenges and Opportunities” by K. Yang *et al.* investigates some of these issues in order to achieve better understanding of these critical security and privacy challenges. They propose a general architecture for a mobile crowdsourcing network comprising of both crowdsourcing sensing and crowdsourcing computing. Then they discuss several critical security and privacy challenges that capture the essential characteristics of MCNs. They go on to formulate some research problems leading to possible research directions hoping to bring attention to further investigation into security and privacy solutions in MCNs.

The final article in this FT involves space information networks using satellites and high-altitude platform stations. Space information networks are able to enhance detection and transmission capabilities compared to current single Earth observation satellites. Although many attempts have been carried out concerning the space network architecture and protocols, the security issues have not been well investigated. C. Jiang *et al.* focus on the security problems in the space information networks from the perspectives of secure handoff, secure transmission control, secure key management, and secure routing. In their article, “Security in Space Information Networks,” they review the challenges and open problems, and provide some solutions regarding the security issues on space information networks.

We are confident that these articles will add value to your research activities and give an overall direction for those researchers interested in this topic.

The Guest Editors would like to thank the previous

Editor-in-Chief (Sean Moore) and the current Editor-in-Chief (Osman Gebizlioglu) for their guidance, feedback, and encouragement along the way. We are very grateful to them for allowing us to schedule two issues of the FT due to the large number of submissions received from highly qualified researchers. We also thank the IEEE Communications Magazine Publications Staff for their patience and hard work in making this issue a reality.

BIOGRAPHIES

MOHSEN GUIZANI [S'85, M'89, SM'99, F'09] (mguizani@ieee.org) is currently a professor and the associate vice president of graduate studies at Qatar University. He previously served as chair of the Computer Science Department at Western Michigan University (2002–2006) and chair of the Computer Science Department at the University of West Florida (1999–2002). He received his B.S., M.S., and Ph.D. degrees in electrical and computer engineering from Syracuse University, New York. His research interests include wireless communications and mobile computing, cloud computing, cyber security, and smart grid. He is the author of nine books and more than 400 publications in refereed journals and conferences. He served as an IEEE Computer Society Distinguished Speaker from 2003 to 2005. He is a member of the IEEE Communications and IEEE Computer Societies and ASEE, and is a Senior Member of ACM.

DAOJING HE (hedaojinghit@gmail.com, djhe@sei.ecnu.edu.cn) received his B.Eng. (2007) and M. Eng. (2009) degrees from Harbin Institute of Technology, China, and his Ph.D. degree (2012) from Zhejiang University, China. He is currently a professor in the Software Engineering Institute, East China Normal University. His research interests include network and systems security. He is an Associate Editor or on the Editorial Boards of a number of international journals such as *IEEE Communications Magazine*.

KUI REN (kuiren@buffalo.edu) is an associate professor at the State University of New York at Buffalo. His research interests span cloud and outsourcing security, and wireless and wearable security. His research has been supported by NSF, DoE, AFRL, MSR, and Amazon. He was a recipient of NSF CAREER Award in 2011 and Sigma Xi/IIT Research Excellence Award in 2012. He is an Associate Editor for IEEE TMC, TIFS, TSG, and so on. He is a Distinguished Lecturer of IEEE.

JOEL J. P. C. RODRIGUES [S'01, M'06, SM'06] (joeljr@ieee.org) is a professor in the Department of Informatics of the University of Beira Interior, Covilhã, Portugal, and a researcher at the Instituto de Telecomunicações, Portugal. He is the leader of NetGNA Research Group (<http://netgna.it.ubi.pt>), Chair of the IEEE ComSoc TC on eHealth, Past Chair of the IEEE ComSoc TC on Communications Software, and a Steering Committee member of the IEEE Life Sciences Technical Community. He is the Editor-in-Chief of three international journals, and a co-author over 400 papers, two books, and three patents. He is the recipient of several Outstanding Leadership and Outstanding Service Awards by IEEE Communications Society and several best paper awards.

SAMMY CHAN [S'87, M'89] (eeschan@cityu.edu.hk) received his B.E. and M.Eng.Sc. degrees in electrical engineering from the University of Melbourne, Australia, in 1988 and 1990, respectively, and a Ph.D. degree in communication engineering from the Royal Melbourne Institute of Technology, Australia, in 1995. He is an associate professor in the Department of Electronic Engineering, City University of Hong Kong.

YAN ZHANG (yanzhang@simula.no) received a Ph.D. degree from Nanyang Technological University, Singapore. Since August 2006, he has been working with Simula Research Laboratory, Norway. He is currently head of the Department of Networks at Simula Research Laboratory, and an adjunct associate professor at the Department of Informatics, University of Oslo, Norway. He is a Regional Editor, Associate Editor, on the Editorial Board, or a Guest Editor of a number of international journals. His recent research interests include wireless networks, cyber physical systems, and smart grid communications.

Toward Trustworthy Vehicular Social Networks

Qing Yang and Honggang Wang

ABSTRACT

Wireless vehicular networks offer the promise of connectivity to vehicles that could provide a myriad of safety and driving-enhancing services to drivers and passengers. With wireless technology available in each car, it is expected that huge amounts of information will be exchanged between vehicles or between vehicles and roadside infrastructure. Due to defective sensors, software viruses, or even malicious intent, legitimate vehicles might inject untrustworthy information into the network. Besides relying on the public key infrastructure, this article proposes a social network approach to study trustworthy information sharing in a vehicular network. We first cover recent research progress in measuring direct trust and modeling indirect trust in online social networks, and then discuss how to apply them to vehicular social networks despite several pressing research challenges.

INTRODUCTION

Emerging wireless technologies enable vehicles to connect to each other to form a vehicular network through wireless channels and share traffic- or entertainment-related information to provide improved safety and pleasure to drivers and passengers. Besides existing wireless technologies (e.g., Bluetooth), various connectivity solutions such as dedicated short-range communication (DSRC), cellular networks, and WiFi are being bundled with original equipment manufacturer (OEM) made cars. Potential applications of networking vehicles include enhanced driving safety, smart roadside information systems, and environment-friendly transportation. In fact, global connected car market shipments are expected to reach 59.86 million units and are likely to reach \$98.42 billion by 2018 [1]. The existing and expected consumer demands and mandates are the major drivers for the connected vehicle market.

Various types of information are exchanged between cars in a vehicular network, including traffic jams, road construction, incidents/crashes, road conditions, and weather alerts, so it is important for a vehicle/driver to distinguish trustworthy from untrustworthy data. Vehicles

sharing factual information with others are considered to be trustful, while those sending false information are distrustful. Currently, most research on trustworthy information sharing in vehicular networks rely on public key infrastructure (PKI). Although PKI builds the first line of defense, it is possible for legitimate vehicles to send untrustworthy information due to defective sensors or computer viruses, or even for malicious reasons. Untrustworthy information sent by distrustful vehicles has the potential to become the most harmful data within a vehicular network; for example, a driver might report wrong parking information to ensure that he can park in the desired parking lot. A vehicular network, on the other hand, is also a mobile social network where vehicles meet each other to establish friendship-like relations and thus are embedded in a social structure. Therefore, an interesting question arises: *Is it possible to achieve a trustworthy vehicular network by applying the research findings about trust in social networks to vehicular networks?*

VEHICULAR SOCIAL NETWORK

The term vehicular social network (VSN) was coined in [2]. A VSN connects vehicles that are physically close to each other and enables them to take advantage of their proximity to form a tightly coupled, ad hoc, virtual world. This definition only considers the social connections between vehicles that are in each other's communication range. However, a VSN could be defined in a broader sense as a network of physically or virtually connected vehicles that are interested in sharing information for a common purpose or benefit. In such a network, physical and virtual connections between vehicles can be built via DSRC, cellular networks, and cloud [3].

Recently, the development of VSNs has gained momentum, coming largely from various applications of VSNs in people's daily lives. Figure 1 categorizes existing VSN applications based on the physical and social distances between drivers. A driver can interact with others with different levels of familiarity, varying from family members to acquaintances (or even

Qing Yang is with Montana State University.

Honggang Wang is with the University of Massachusetts, Dartmouth.

strangers). These people could be within the driver's car, near the car, or far away. With the help of a VSN, the driver can efficiently share information with others, including, but not limited to:

- Having fun with *CarPlay* in a car (with the family)
- Sharing a restaurant review with friends via *FourSquare*
- Tracking locations of family members by *Life360*
- Scheduling a carpool with co-workers through *KarPooler*
- Sharing her location in real time with acquaintances through *Glympse*
- Sharing a ride and splitting the cost with another person who requests a ride along a similar route by *UberPool*
- Cooperatively driving with nearby cars
- Sharing traffic-related information to strangers via *Waze*

Some of the above-mentioned systems (e.g., *FourSquare*) can actually provide services to drivers at various social and physical distances.

While the authors believe VSN will become popular within families and between close friends, it will mature when drivers are able to share information with strangers because more drivers are participating in VSNs. Therefore, this article focuses on the technical challenges of realizing trustworthy information sharing between stranger vehicles in VSNs.

SOCIAL CONNECTIONS BETWEEN VEHICLES

To understand how and why stranger vehicles are socially connected in a VSN, we first revisit the life cycle of relationships among humans in traditional social networks. The lifespan of such a relationship can be divided into three stages:

- A weak connection is created.
- A weak connection is cultivated and becomes strong.
- A strong connection is maintained, and a weak connection is terminated.

At each stage, a connection can be created, cultivated, and maintained only when two people both have the desire to further their relationship.

In a VSN, it is evident that drivers desire to share information with each other, which has been demonstrated in *Waze.com*, one of the world's largest community-based traffic and navigation systems. However, it is unclear whether this desire leads to a social network among vehicles. Therefore, it is necessary to investigate the life cycle of relations in traditional social networks and compare them to the counterparts in a VSN.

In the first stage, weak connections in traditional social networks are created from communities or existing social networks. If two persons share a common community (e.g., a dancing club), the likelihood that they create a connection becomes significantly large. Another way of building new connections is through a person's existing social network; for example, your friend might introduce you to one of her friends. In VSNs, a car could connect with another if they encounter (within the communication range) others on the road and have common interest in

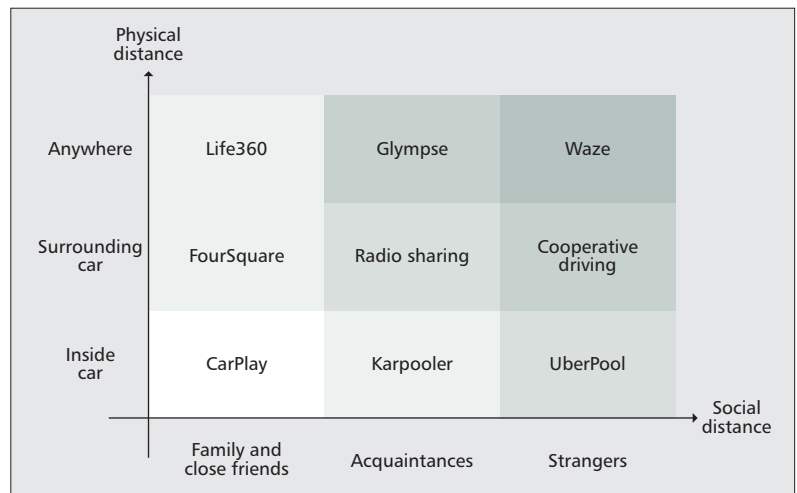


Figure 1. Applications of VSNs categorized by physical and social distances.

the shared information. While vehicles are exchanging messages, their local social connections could also be shared, so it is highly possible for a vehicle to create virtual connections to its "friends' friends."

In the second stage, a weak connection becomes a strong one when people in traditional social networks have the chance to interact with each other and cultivate their relationship. In VSNs, two vehicles frequently encountering and sharing information with each other (e.g., parking in the same parking lot) have the chance to make a weak connection strong.

In the third stage, strong connections in traditional social networks are maintained, while weak connections are terminated. In the context of VSNs, strong connections are maintained if and only if two vehicles keep sharing mutually beneficial information.

CONSTRUCTION OF VEHICULAR SOCIAL NETWORKS

A VSN can be constructed in a centralized or distributed way. In the centralized solution, a social connection between two vehicles is uploaded to the cloud [3] via existing cellular technologies (e.g., third generation, 3G, or Long Term Evolution, LTE). After collecting these social connections, a social network among vehicles could be built and downloaded by vehicles. In the distributed solution, each vehicle shares its social connections with nearby vehicles via DSRC technology. By exchanging social connections with others, a vehicle can incrementally construct its vehicular social network. Although the centralized one provides a reliable and real-time solution, the distributed one is cheaper as it has no dependence on infrastructures. We believe these two solutions are complementary to each other; that is, when cloud service is not available or too expensive, a distributed solution becomes a better candidate, and vice versa.

A social connection in a VSN initiates when two vehicles encounter and share information with each other. Such connection information will be either uploaded to the cloud [4] or saved by these two vehicles. In this article, we focus on

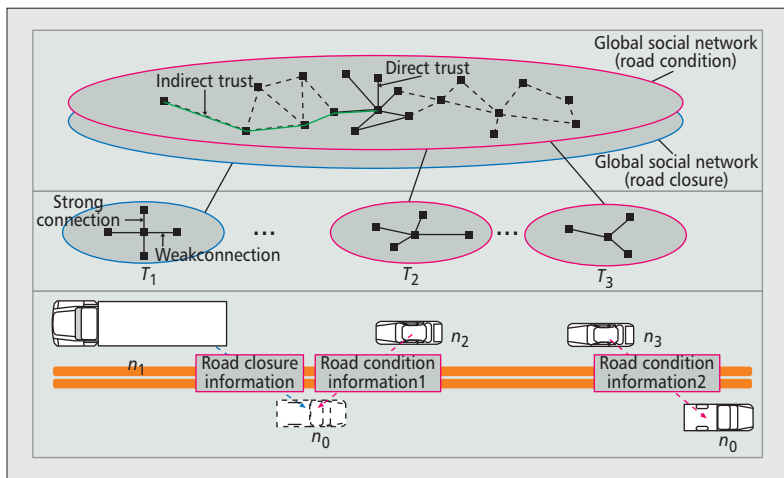


Figure 2. Illustration of a distributed vehicular social network construction.

constructing VSNs in a distributed manner as it is a cheaper solution requiring no infrastructures. In the distributed solution, each vehicle maintains a local social connection tree (SCT). The root of a local SCT is the vehicle that maintains this SCT; the edges in the SCT are the (direct or indirect) social connections of the vehicle.

As shown in Fig. 2, a vehicle (e.g., n_0) keeps tracking the interactions between itself and nearby vehicles. If it receives a message (e.g., road condition) from vehicle n_2 , it first checks whether n_2 is in its local SCT. If not, a node n_2 is created and inserted into its local SCT. Otherwise, the weight of the edge between n_0 and n_2 is updated based on the new interaction. More positive interactions lead to a stronger (heavier) connection (edge).

When two vehicles encounter each other on the road, they could also share their local SCTs. With the SCTs received from others, a vehicle can incrementally build its global social network. Knowing more vehicles with common interests and receiving more relevant information motivates vehicles to share their local SCTs with others. Both local SCTs and global social networks are built in a distributed manner, so the constructing process scales in large vehicular networks.

If a social network structure exists in a vehicular network, it is possible to study the trustworthiness among vehicles by a social network approach. For example, if two vehicles have a weak connection, it probably means that the trust between them is low. At least, there is a lack of positive evidence of a trustful connection between them. Later on, if they exchange lots of useful information, a weak connection may become a strong one, indicating high trustworthiness between them.

The trust computed from vehicle interactions is considered direct trust as vehicles are directly connected to each other. Another type of trust is called indirect trust, which describes the trust relations of vehicles that never physically meet. For example, a vehicle may have virtual connections to vehicles recommended by its “friends.”

Last but not least, trust is context-specific, so

different local SCTs and global social networks are needed for different contexts. For instance, when vehicle n_0 receives road closure information from truck n_1 , the connection between n_0 and n_1 is created in n_0 's local SCT dedicated to road closure. Later on, when it receives local SCTs about road conditions from n_2 and n_3 , these SCTs are integrated into n_0 's global road condition social network. How to classify messages into different categories is a trade-off where more categories (contexts) increase the processing time and storage, while fewer categories yield inaccurate trust estimations.

TRUSTWORTHY VEHICULAR SOCIAL NETWORKS

Thanks to the technological advances of wireless networks, drivers in populated locales are connected to form a tightly coupled and ad hoc mobile VSN. Within such a mobile social network, trustworthiness between temporarily connected vehicles could be mined from their interactions and applied within a VSN to achieve trustworthy information sharing among vehicles [5]. In this section, we start with the current research on direct and indirect trust in online social networks (OSNs) and then extend them to VSNs. Finally, we provide a discussion on what the authors believe are the most important research challenges that lie ahead.

DIRECT TRUST MEASUREMENT

Several works about direct trust measurement in online social networks exist. Researchers have focused on measuring trustworthiness in OSNs based on users' similarity [6] and interactions [7].

Direct Trust in OSNs: Strong correlation between trust and users' interest similarity was found in [6], based on the data obtained from two real-world online communities: All Consuming and FilmTrust. For each dataset, a user's interest profile is constructed from the ratings he makes on corresponding items (e.g., books or movies). Such interest profiles are then used to compute the similarities between different users. The trust level between any two users is then measured from the profile similarity. In other words, the more similar the users, the more likely they trust each other.

Besides user interest similarity, direct trust can also be mined from user interactions. In [7], the authors asked 35 participants to rate the trustworthiness of their Facebook friends, which served as the ground truth. Seventy-four Facebook variables were then collected for each participant and used to compute the trustworthiness of his friends. Finally, the authors model tie/connection strength between users as a linear combination of these predictive variables. They discover that user interaction data can be used to distinguish strong and weak ties with more than 85 percent accuracy.

Direct Trust in VSNs: Previous study on direct trust in OSNs reveals that it is possible to compute or mine direct trust from driver interest similarities and interactions in VSNs. To apply

the existing approaches, however, we need to make some modifications. For example, the user interest similarity between vehicles should be re-defined because the only product in VSNs are messages exchanged between vehicles. Information shared between vehicles should be classified into different categories as in [6], and a driver's interest profile vector should be built from ratings of all categories. The rating of a category can be modeled as the user's estimate of the trustworthiness of received information in that category. To determine whether a received message is trustworthy or not, each vehicle must have an information discrimination system. Vehicles could rely on their own sensors to evaluate the trustworthiness of received messages. In addition, they can upload received messages to the cloud to conduct information discrimination.

If information sharing between vehicles is considered as interactions between drivers, mining trust from interacting data becomes a suitable solution to obtaining direct trust between vehicles. However, messages are exchanged only when two vehicles with a common interest encounter each other, so research is needed to study how such interactions can reflect trust between vehicles. Besides, compared to OSNs, the volume of vehicular interaction data could be extremely large because vehicles can potentially exchange messages with all neighboring vehicles. Therefore, a resource-aware information discrimination scheme is needed so that only the most relevant information is shared between vehicles with a common or similar interest.

INDIRECT TRUST INFERENCE

Due to the propagative nature of trust, focus on indirect trust inference mainly lies in modeling trust propagation along trust relations between people who do not have direct trust connections [8]. To model and determine the trust between two users having no direct interaction in an OSN, there are mainly two method types in the literature: topology-based and evidence-based.

Indirect Trust in OSNs: A good example of topology-based method is proposed in [9], which leverages the truth that a disproportionately small "cut" exists between Sybil, users with multiple fake identities, and honest nodes to distinguish distrustful from trustful users. Since this approach is designed to identify Sybil users, indirect trust is considered a binary value indicating whether or not a user is trustful.

The authors discover that a distrustful user may create many fake identities but could build a limited number of connections (or edges) to legitimate users. By looking at the social connections between users, a community composed of distrustful nodes can be identified and eliminated from trustful ones. Topology-based indirect trust inference is good for some applications (e.g., Sybil node detection), but it has limitations in computing non-binary trustworthiness, which could be addressed by evidence-based approaches [10].

To understand indirect trust in online social networks, Jøsang proposed the seminal work of modeling trust by the subjective logic model

[10]. Subjective logic is a type of probabilistic logic based on the Dempster-Shafer belief theory, explicitly taking uncertainty and belief ownership into account while computing trustworthiness. It treats trustworthiness as opinions and introduces an algebra for opinion operations, such as discounting and consensus operations. An opinion in subjective logic contains three components: belief, distrust, and uncertainty, which reflect the subjectivity and uncertainty existing in users' assessment of others' trustworthiness.

Indirect Trust in VSNs: While direct trust could be obtained from vehicles' similarity and interactions, indirect trust is derived from other vehicles' recommendations. The indirect trust from a truster vehicle to a trusted vehicle highly depends on how strongly they are connected, which is affected by many factors such as the social distance, connection strengths, and social network topology between them. If the trustworthiness of all connections in a vehicle's global social network are known, the problem of computing indirect trustworthiness in the VSN can be defined as follows.

Given the global social network $G(V,E)$ of a certain vehicle u , \forall vehicle v such that $e(u, v)$ not in E , and \exists at least one path from u to v , how to compute the trustworthiness of v to u , that is, how u should trust a stranger v based on her social connections.

To solve this problem, the evidence-based approach is more appropriate than the topology-based ones because subjective logic is able to compute non-binary trust values, which enables the comparison between trustful vehicles, for example, to identify the most trustful vehicle. Moreover, non-binary trust values also mean more accurate trust evaluations.

Applying subjective logic to compute indirect trust in VSNs, however, faces a few challenges. First, subjective logic defines trust as an opinion vector containing belief, distrust, and uncertainty. Most existing direct trust datasets only provide the values of belief; for example, the trust level of a relation is 0.7. That the other 0.3, however, could mean either distrust, uncertainty, or both. Second, subjective logic cannot handle complex topologies. Although some approximation solutions are proposed to select the strongest path while computing indirect trust in a complex network, the selection procedure causes information loss in trustworthiness computation. Third, applying subjective logic on a graph requires finding all possible paths (social connections) between two vehicles, which is an NP-hard problem. Therefore, procedures are needed to control the size of local SCTs on vehicles, to trim a global social network by eliminating less trustful edges, and to design an approximation algorithm to find all paths between the truster and trusted vehicles.

In summary, a vehicle could evaluate the direct trust of another vehicle based on their past interactions, and infer indirect trust by investigating how closely they are socially connected. This approach is different from traditional ones that focus on securing or authenticating vehicles based on the public key

Drivers in populated locales are connected to form a tightly coupled and ad hoc mobile VSN. Within such a mobile social network, trustworthiness between temporarily connected vehicles could be mined from their interactions and applied within a VSN to achieve trustworthy information sharing among vehicles.

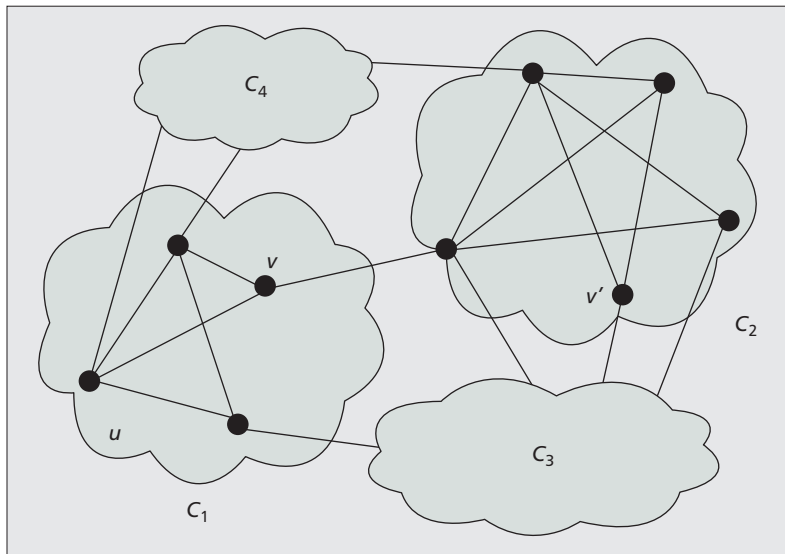


Figure 3. A trust social network is split into different clusters (or clubs).

infrastructure (PKI). While PKI builds the first line of defense, it guarantees only the identification of legitimate vehicles but not their trustworthiness.

RESEARCH CHALLENGES

There are many challenges in deploying a trustworthy vehicular social network, such as increasing the penetration rate of equipped vehicles, constructing adequate roadside infrastructure, and establishing new policies and laws. In this article, however, we only focus on the technical issues relevant to designing a trustworthy vehicular network. To realize a trustworthy VSN, the following research challenges need to be addressed:

- Information classification and discrimination
- Resource-aware information discrimination
- An appropriate indirect trust model
- An efficient algorithm for computing indirect trust

INFORMATION CLASSIFICATION AND DISCRIMINATION

To achieve context-aware trust evaluation, messages exchanged among vehicles must be classified into different categories and discriminated based on their trustworthiness. On each vehicle, a trust evaluation system will be developed to affirm the information conveyed in received messages. Vehicles could rely on their own onboard sensors and roadside infrastructures to obtain the ground truth, and then discriminate the trustworthiness of received messages. Such a system will enable a vehicle to distinguish positive from negative and uncertain messages. With the subjective logic model [10], trustworthiness between vehicles could be modeled based on the amount and nature (positive, negative, or uncertain) of their interactions.

Furthermore, this trust evaluation system should encompass the multifaceted (e.g., dynamic, asymmetric, and subjective) aspects of trust.

Analyzing messages exchanged among vehicles might cause the privacy leakage of drivers, so privacy protection mechanisms are also needed.

RESOURCE-AWARE INFORMATION DISCRIMINATION

With the basic trust evaluation system set up, it is essential to study resource-aware information discrimination due to the massive amount of messages exchanged between vehicles. For example, a probabilistic information discrimination mechanism might help a vehicle spend its resources on discriminating the messages with the greatest application relevance. In addition, a piggybacked notification scheme could reduce the communication overhead by sending the notification of untrustworthy vehicles to others. With this scheme, vehicles that are unable to identify untrustworthy messages will learn of such information from their trustworthy “friends.”

INDIRECT TRUST MODEL

A VSN can be considered a dynamic graph where nodes represent vehicles and time-varying edges indicate connections between temporarily connected vehicles. Given such a dynamic graph, the following research questions need to be studied. Does indirect trust exist between vehicles? Is it possible to compute the indirect trust between vehicles from their social connections? How can we appropriately model trust propagation within VSNs? Could multiple trust relations be combined to form a new one (and how)? What is the best way to compute the expected trustworthiness of a social connection?

EFFICIENT ALGORITHM ON COMPUTING TRUST

Assuming the indirect trust model is in place, an efficient indirect trust computation algorithm is needed. This is because finding all possible paths between two vehicles in a global social network is NP-hard. Therefore, a global social network must be preprocessed (e.g., divided into disjoint clusters with small diameters), so searching all paths only occurs within a cluster. Existing clustering algorithms, such as the k -means algorithm and spectral analysis, can be applied to split a global social network into several smaller clusters.

As shown in Fig. 3, a global social network is divided into four clusters, and only a sub-graph (e.g., cluster C_1) is used to compute the trustworthiness between nodes u and v . If these two nodes are not in the same cluster (e.g., u and v'), the original network (composed of clusters C_1 – C_4) will be used. Alternatively, we can dynamically adjust the algorithm’s searching depth; that is, if the computed indirect trust between the trustor and trustee is accurate, there is no need to search deeper to find all possible paths.

CONCLUSION

This article has provided a social network approach to study trustworthy vehicular networks, that is, measuring direct trust from past

interacting data and infer indirect trust from social recommendations among vehicles. The evolution of social connections between vehicles and the construction of a vehicular social network are first investigated, followed by an overview of the progress in research of direct and indirect trust in online social networks. Although there are similarities between a vehicular social network and its analogs, online social networks, there appears to be a divide between these two fields.

Leveraging results for online social networks, many new research opportunities exist in trustworthy VSNs, for example, VSN construction protocol, message classification and discrimination, trust information discrimination, privacy protection, direct trust measurement, indirect trust inference, trust computing algorithm, simulation and experiment platforms, and a VSN dataset containing trust information. With the increasing deployment of connected vehicles, the authors expect to see more interdisciplinary research efforts devoted to studying trustworthy VSNs.

REFERENCES

- [1] "Connected Car Market (2013-2018)," 2013, online; accessed 19 Sept. 2014.
- [2] S. Smaldone *et al.*, "Roadspeak: Enabling Voice Chat on Roadways Using Vehicular Social Networks," *Proc. 1st Wksp. Social Network Systems*, 2008, pp. 43–48.
- [3] R. Yu *et al.*, "Toward Cloud-Based Vehicular Networks with Efficient Resource Management," *IEEE Network*, vol. 27, no. 5, Sept. 2013, pp. 48–55.
- [4] G. Nan *et al.*, "Distributed Resource Allocation in Cloud-Based Wireless Multimedia Social Networks," *IEEE Network*, vol. 28, no. 4, July 2014, pp. 74–80.

- [5] D. Huang, X. Hong, and M. Gerla, "Situation-Aware Trust Architecture for Vehicular Networks," *IEEE Commun. Mag.*, vol. 48, no. 11, Nov. 2010, pp. 128–35.
- [6] C.-N. Ziegler and J. Golbeck, "Investigating Interactions of Trust and Interest Similarity," *Decision Support Systems*, vol. 43, no. 2, 2007, pp. 460–75.
- [7] E. Gilbert and K. Karahalios, "Predicting Tie Strength with Social Media," *Proc. CHI '09*, 2009, pp. 211–20.
- [8] G. Liu *et al.*, "Assessment of Multi-Hop Interpersonal Trust in Social Networks by Three-Valued Subjective Logic," *Proc. IEEE INFOCOM*, Apr. 2014, pp. 1698–1706.
- [9] H. Yu *et al.*, "Sybilguard: Defending against Sybil Attacks via Social Networks," *IEEE/ACM Trans. Network*, vol. 16, no. 3, 2008, pp. 576–89.
- [10] A. Jsang, "A Logic for Uncertain Probabilities," *Int'l. J. Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 9, no. 03, 2001, pp. 279–311.

BIOGRAPHIES

QING YANG, Ph.D., is a RightNow Technologies assistant professor in the Department of Computer Science, Montana State University. He received B.S. and M.S. degrees in computer science from Nankai University and Harbin Institute of Technology, China, in 2003 and 2005, respectively. He received a Ph.D degree in computer science from Auburn University in 2011. His research interests lie in the areas of wireless vehicular networks, network security, and trust in online social networks.

HONGGANG WANG received his Ph.D. in computer eEngineering from the University of Nebraska-Lincoln in 2009. He is an assistant professor at the University of Massachusetts, Dartmouth and is an affiliated faculty member of the Advanced Telecommunications Engineering Laboratory at the University of Nebraska-Lincoln. His research interests include wireless healthcare, body area networks, multimedia sensor networks, mobile multimedia and cyber security, wireless networks, and cyber-physical systems. His research is funded by the National Science Foundation, Department of Transportation, and the University of Massachusetts President's Office.

A piggybacked notification scheme could reduce the communication overhead by sending the notification of untrustworthy vehicles to others. With this scheme, vehicles that are unable to identify untrustworthy messages will learn of such information from their trustworthy "friends."

Verifiable, Privacy-Assured, and Accurate Signal Collection for Cloud-Assisted Wireless Sensor Networks

Chia-Mu Yu, Chi-Yuan Chen, and Han-Chieh Chao

ABSTRACT

As cloud-assisted WBSNs have become increasingly popular in healthcare applications, security and privacy threats deserve much more attention. In this article, focusing on data privacy and data completeness issues, we propose a VPAA scheme for cloud-assisted WBSNs. Through both simulation and prototype implementation, we confirm that VPAA is energy-efficient and effective in protecting data privacy and completeness.

WSNs AND WBSNs

Wireless sensor networks (WSNs) have found many applications due to intensive studies in the past. In particular, wireless body sensor networks (WBSNs) have become increasingly popular in healthcare applications [1], and various kinds of biosensors have been developed and used. This trend is also boosted by the decreasing cost of advanced sensors. Thus, we can already find many commercial solutions for healthcare monitoring and biomedical signal extraction on the market. For example, some of the solutions may help detect heart attacks, while some other systems are able to reduce the risk of asthma by evaluating body condition and sending an alert to the patient. The above functionalities are all contributed by biosensors continuously monitoring the biomedical signals of a patient.

A WBSN is composed of a number of wirelessly connected and miniaturized biosensors that can extract such signals. A WBSN alone is not useful; instead, a WBSN is usually connected to a user device or a remote controller unit (RCU). The RCU might be a physician or data processing server in a remote clinic or hospital. Depending on the required functionalities, the patient can purchase sensors from different manufacturers for practical use. Therefore, the RCU usually offers an interface (e.g., a website) through which the patient can register new sensors to the RCU and obtain credentials from the RCU. A potential application is one in which wearable, portable, embeddable, and even implantable sensors are with a patient, collecting

biomedical signals. Then the collected signal readings are forwarded to user devices such as laptop and smartphone, allowing the patient to know the corresponding health condition. In addition, the collected signal readings can be forwarded to a physician in a remote clinic or hospital, allowing remote diagnosis and treatment. The collected signal readings can also be forwarded to an automatic processing server that detects a physiological condition in a real-time fashion and responds to emergency events like heart attack and asthma.

Since the communication capability of sensors is usually limited, a wireless hub with Internet connection is needed to relay the signal readings to the RCU. More specifically, IEEE 802.15.6 is a promising solution for sensor communications, where a simple star topology is used with a wireless hub at the center and various sensors around the wireless hub. The data flow, in essence, is that sensors collect and forward the signal readings to the wireless hub, which in turn forwards the received signal readings to the RCU.

CLOUD-ASSISTED WBSNs

In the above architecture, the RCU in fact bears a huge amount of storage, computation, and communication overhead, because all of the signal readings generated by sensors flood into the RCU. The reasons are as follows. First, the RCU needs to have a huge amount of disk space to store the received signal readings, causing the storage burden. Second, since an RCU may offer more than one service on the WBSN, the incoming readings may require different processing techniques or even be ignored. Nevertheless, the service categorization of incoming readings needs to be performed even if no further processing is required. This leads to additional computation burden on the RCU. Third, for a similar reason as above, possibly only partial signal readings are sufficient for the RCU to make a correct decision for diagnosis and treatment. In this case, redundant communications, in fact, can be avoided. The above concerns are realistic, especially because the amount of data transmit-

Chia-Mu Yu is with Yuan Ze University and the Innovation Center for Big Data and Digital Convergence.

Chi-Yuan Chen is with National Ilan University.

Han-Chieh Chao is with National Ilan University and the School of Information Science and Engineering, Fujian University of Technology.

ted by a sensor of WBSNs monitoring physiological signals can easily reach nearly 2.7 GB. If a higher data rate is employed, this number can even reach up to 31 GB per day [2].

As the emerging cloud computing technologies evolve, the use of the cloud in different applications to facilitate the processing work becomes promising. Here, the cloud is introduced as an online storage with the processing capability to reduce the burdens on the RCU. Two representatives of such cloud services are Google Health and Microsoft HealthVault. After the placement of the cloud between sensors and the RCU, the data flow is changed accordingly. More specifically, sensors collect and forward the signal readings to the wireless hub, which in turn forward them to the cloud. The cloud keeps the signal readings in its permanent storage. Depending on the requests from the RCU, the cloud transmits the requested signal readings to the RCU.

An illustration of cloud-assisted WBSNs is shown in Fig. 1, where three patients are equipped with biosensors. The biosensors transmits the biomedical signals to the cloud, which will answer the queries made by the RCU. For example, the biosensors keep sending heartbeat signals to the cloud. The RCU is now interested in patient 3's heartbeat information, and therefore sends a query to the cloud. The cloud responds with patient 3's heartbeat information.

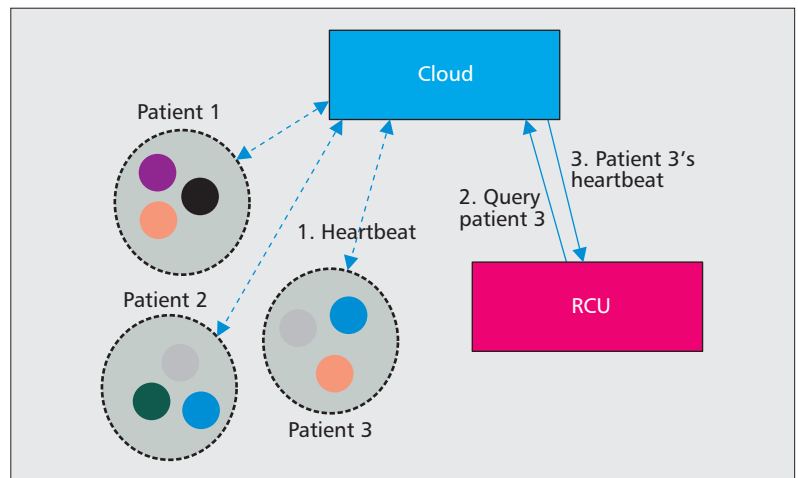


Figure 1. A conceptual illustration of cloud-assisted WBSNs.

assisted WBSN, the cloud, due to hardware error, software misconfiguration, or even malicious manipulation, may report to the RCU only an incomplete set of signal readings. The use of cryptographic techniques cannot solve this problem, because such misbehavior can only be regarded as a particular type of packet loss that cryptographic tools do not help.

SECURITY AND PRIVACY CONCERNS OF CLOUD-ASSISTED WBSNS

Although WBSNs have attracted many research efforts, the security and privacy issues of cloud-assisted WBSNs need more attention [3]. In this article, we consider the task of collecting signal readings from sensors in WBSNs. In other words, an RCU may be in charge of thousands of patients, each with a number of sensors extracting biomedical signals. The RCU may be interested in retrieving the signal readings from a particular sensor or particular types of sensors. The RCU may issue a query to the cloud, asking for the corresponding signal readings stored in the cloud. We believe that this functionality is one of the most popular types of queries, and is needed in almost all of the applications. However, achieving the goal is not trivial, especially after the cloud is introduced as a proxy to cache the signal readings.

Confidentiality and authenticity are issues with which all of systems with data forwarding must be concerned. Thus, they are also of great importance in the context of WBSNs, since the biomedical signals extracted by sensors can be used to identify individuals, resulting in confidentiality breaches, and the signal readings are also susceptible to modification, resulting in authenticity breaches. While data privacy and authenticity can be guaranteed by using off-the-shelf solutions such as AES and SHA-1, the requirement for data completeness is more subtle. Data completeness in the context of data collection in WBSNs states that a complete set of signal readings for a given sensor should be received by the RCU. Nevertheless, in a cloud-

DESIGN CHALLENGE OF PRIVACY-ASSURED SIGNAL COLLECTION

To handle the security and privacy concerns of cloud-assisted WBSNs, we face the following design challenges:

- The bit rate and delay requirements of WBSNs are important factors in designing protocols over the WBSNs. The bit rate and delay requirements may vary depending on the underlying applications. For example, applications such as cardiac pacemakers, implanted defibrillators, and neurostimulators require bit rates up to 400 kb/s. Other applications, such as swallowable camera pills, may require a bit rate of approximately 1 Mb/s and < 250 ms delay. EMG also requires a bit rate of approximately 1.5 Mb/s and < 250 ms delay. As a consequence, if the privacy-assured mechanism in use is too heavyweight and results in significantly reduced throughput or increased delay, it is unacceptable from the application point of view.

- The wearable, portable, embeddable, and even implantable sensors in WBSNs are all assumed to be extremely resource-constrained, mainly because of their small or even tiny size. The storage, processing, and communication capabilities are thus all very restricted and cannot afford the execution of complicated protocols. In addition, since sensors are usually battery-powered and a complicated protocol may consume energy in either computation or communication, a simple design of the privacy-assured mechanism is preferred;

- Although the resources for the design are very limited, a proper design of a simple privacy-assured mechanism can be done by leveraging a particular feature of WBSNs; that is, the RCU

The encrypted signal readings in both solutions are explicitly transmitted, incurring unavoidable communication overhead for the sensors. In addition, since AES has fixed output length and high-entropy output, no particular feature in encryption outputs can be utilized to achieve compression of sensor readings.

may assume to be resource-abundant. In particular, although sensors are usually battery-powered and resource-constrained, the RCU can be expected to have a PC-level computation capability without the concern of power failure. With the consideration of this feature, in the protocol design, the computation burden needs to be shifted from the sensor side to the RCU side in order to minimize the delay and energy consumption of sensors. On the other hand, as the radio circuit usually plays a major role in consuming energy, the communication task also has to be reduced by compressing the signal readings so as to further minimize energy consumption and prolong battery lifetime.

•The compression may reduce the volume of data to be transmitted. However, considering that the recovered signals need to be exactly the same as the uncompressed ones, the choice and design of compression mechanism could be difficult. Note that such a strict requirement on perfect signal recovery is needed because even an inaccurate signal reading might cause improper diagnosis/treatment and have serious consequences. Lossless compression techniques can apply to biomedical signals, but usually the compression ratio is not satisfiable, and the compression algorithm is complex. A technique called compressed sensing developed in the signal processing community is also adopted to achieve data volume reduction. Nevertheless, data reduction and recovery heavily rely on prior knowledge of signal sparsity and its signal dictionary, resulting in its inapplicability in certain types of biomedical signals.

EXISTING SOLUTIONS TO ACHIEVE PRIVACY-ASSURED SIGNAL COLLECTION

Many research efforts have been devoted to the issues of data privacy and completeness in generic two-tiered wireless sensor networks (WSNs). Since WBSNs and WSNs share similar data flows, the solutions in WSNs are supposed to be candidates for WBSNs. Nevertheless, they do not all directly apply to a healthcare monitoring system. More specifically, due to the nature of a bucket scheme, the solutions in [4–6] actually have privacy leakage. The adversary can narrow down the possibility, with an educated guess on signal readings. Due to the use of the neighborhood chain technique, the additional communication overhead in [7] grows linearly with the number of signal readings, which may be a huge communication burden on the sensor. The watermark technique [8], a technique similar to the neighborhood chain, also causes similar communication overhead.

A recent development, CDG, on efficient data collection in sensor networks, can be found in [9]. In [9], an advanced signal processing technique, compressed sensing, is used in data collection with the guarantee of reducing the overall communication cost significantly, especially in large-scale sensor networks. However, CDG also applies only to generic sensor networks.

Some of the other aspects of security and pri-

vacy issues in cloud-assisted WBSNs are discussed in [10–12]. Nevertheless, here we focus only on the signal collection problem in cloud-assisted WBSNs.

NAÏVE METHOD

Due to the signal sampling in sensor hardware, each sensor s_i is assumed to have β signal readings and the corresponding AES encryption outputs. A straightforward solution, S_1 , is for s_i to calculate and send out β encrypted signals and the corresponding cryptographic hashes during the sensor registration, and for these to be uniquely shared by s_i and the RCU. Another straightforward solution, S_2 , is for s_i to calculate a single cryptographic hash of all β encrypted signals, and then send out β encrypted signals and the calculated hash. Nevertheless, the encrypted signal readings in both solutions are explicitly transmitted, incurring unavoidable communication overhead for the sensors. In addition, since AES has fixed output length and high-entropy output, no particular feature in encryption outputs can be utilized to achieve compression of sensor readings.

RANDOMIZED AND DISTRIBUTED ORDER-PRESERVING ENCRYPTION FOR PRIVACY ASSURANCE

Order-preserving symmetric encryption (OPE) [13] is a deterministic encryption scheme over numerical values. If the plaintexts x_1 and x_2 satisfy $x_1 < x_2$, it can ensure that $\mathcal{E}_K(x_1) < \mathcal{E}_K(x_2)$, where $\mathcal{E}_K(\cdot)$ denotes the OPE function with key K .

A simple OPE was presented in [13]; given that y numbers, $x_1 < \dots < x_y$, are the possible plaintexts, we generate an array of y uniformly random numbers $k_1 < \dots < k_y$ as the key K . The encryption (decryption) is accomplished by searching in K for the ciphertext (plaintext) corresponding to the plaintext (ciphertext); for example, $\mathcal{E}_K(x_i) = k_i$, $1 \leq i \leq y$.

In OPE, its key also acts as possible ciphertexts. OPE reveals nothing but the numerical order of plaintexts because all of the ciphertexts k_1, \dots, k_y are distributed uniformly over a specific range. Despite the leakage of numerical order, OPE is in fact provably secure [14]. In spite of the drawback of large key size, we keep such a simple form of OPE in mind for the ease of presentation. More sophisticated OPE schemes can be found in [13, 14].

OPE has been applied widely to encrypted database retrieval, where the data are generated from a single authority. However, this is not the case in the WBSN setting. In addition, because the number of possible sensor readings may be limited and known from hardware specification, the relation between plaintexts and ciphertexts could be revealed. For example, if only 20 kinds of possible outputs can be generated by particular sensors, the adversary can practically derive the OPE key by investigating the numerical order of the intercepted ciphertexts.

Our solution is a novel use of OPE, called

rdOPE, which is randomized OPE, involving random encryption over distributed sources with a limited input value range. The technical challenge of rdOPE design is to maintain the numerical orders of encryptions from different sensors that use different OPEs. Nevertheless, since the possible mapping between plaintexts and ciphertexts are fixed by RCU in advance, the ciphertexts can be determined in the WBSN setting such that the numerical orders of ciphertexts in different sensors can be preserved.

More specifically, rdOPE for n sensors with r possible sensor readings is defined as an encryption scheme $\langle \mathcal{E}, \mathcal{D}, k^{(i)}, h_{rdOPE}(\cdot), n, r, b, c \rangle$ such that

$$\mathcal{E}_{k^{(i)}}(x_1) < \mathcal{E}_{k^{(j)}}(x_2) \text{ if } x_1 < x_2, 1 \leq i, j \leq n, \quad (1)$$

where $k^{(i)}$ and $k^{(j)}$ denote the rdOPE keys possessed by s_i and s_j , respectively, and the value ranges of the hash output $h_{rdOPE}(\cdot)$ and encryption function output $\mathcal{E}_{k^{(i)}}(\cdot)$ are $[1, c]$ and $[1, b]$, respectively. Two rdOPE design examples are shown in Figs. 2a and 2b.

An instance of rdOPE key construction (also called rdOPE table construction) works as follows. At first, rcn possibly distinct numbers, $k_1 \leq \dots \leq k_{rcn}$, are chosen randomly from $[1, b]$ by RCU. The numbers k_1, \dots, k_{rcn} are partitioned into r groups, g_1, \dots, g_r , where $g_{\hat{i}}$ consists of $k_{1+(\hat{i}-1)cn}, \dots, k_{\hat{i}cn}$, $1 \leq \hat{i} \leq r$. RCU randomly samples c numbers from $g_{\hat{i}}$ without replacement, and then stores them in s_i , $1 \leq i \leq n$. The c numbers from $g_{\hat{i}}$ are the possible ciphertexts of the plaintext input, \hat{i} . As a result, the rdOPE key $k^{(i)}$ for s_i is a $c \times r$ array containing s_i 's possible ciphertexts. In the above rdOPE key construction, if k_1, \dots, k_{rcn} are selected such that $k_{cn} \neq k_{cn+1}, k_{2cn} \neq k_{2cn+1}, \dots$, and $k_{(r-1)cn} \neq k_{(r-1)cn+1}$, the constraint of Eq. 1 can always be fulfilled based on the partitioning rule of g_1, \dots, g_r .

Let $k_{j,v}^{(i)}$ be the v th possible ciphertext corresponding to the j th plaintext in s_i . For example, $k_{3,2}^{(1)} = 8$ and $k_{5,1}^{(2)} = 14$ in Fig. 2b. rdOPE works as follows. When having a sensor reading x_j , the sensor s_i simply computes $\mathcal{E}_{k^{(i)}}(x_j) = k_{j,v}^{(i)}$ with $v = h_{rdOPE}(x_j || k_i)$. Once $k_{j,v}^{(i)}$ is received by RCU, the decryption $\mathcal{D}_{k^{(i)}}(k_{j,v}^{(i)})$ can be accomplished by searching in $k^{(i)}$ for the plaintext corresponding to $k_{j,v}^{(i)}$.

rdOPE in fact offers only one way of using OPE in a distributed system, and the distribution of $k^{(i)}$ on different sensors are still all uniform. Thus, the security of rdOPE can be guaranteed to be not less than that of OPE in the case of $c \geq 2$, while rdOPE degenerates to the simple OPE presented in [13] in the extreme case of $c = 1$.

Because of the c choices of encryption output of each plaintext, even if the number of possible inputs is limited, it is more difficult for the adversary to infer the plaintext by correlating the eavesdropped ciphertexts to the possible plaintexts.

The rdOPE scheme is also robust to bit error because all of the entries in $k^{(i)}$ are chosen manually by the RCU; therefore, when an RCU receives a number that claims to be an encryption of rdOPE, it is easy for the RCU to verify such a claim by checking whether the received number appears in $k^{(i)}$.

		Possible plaintext inputs							Possible plaintext inputs				
		1	2	3	4	5			1	2	3	4	5
Sensor ID	S_1	1	4	7	10	13	Sensor ID	S_1	1, 3	4, 5	7, 8	10, 11	13, 14
	S_2	2	5	8	11	14		S_2	2, 3	5, 6	8, 9	11, 12	14, 15
	S_3	3	6	9	12	15		S_3	3, 2	4, 6	7, 8	10, 12	13, 15

Figure 2. Examples of rdOPE: a) the case of $c = 1$; b) the case of $c = 2$.

There are two possible concerns regarding implementing rdOPE in sensor networks:

- The additional computation burden for RCU to calculate the rdOPE table. This concern involves the computation of rdOPE keys. Since the amount of computation linearly grows with the number of rdOPE keys, the effort of calculating an rdOPE table is affordable under the usual assumption of a powerful RCU.

- The additional space requirement for each sensor to store the corresponding rows of the rdOPE table. This concerns the storage overhead. The rdOPE table is of size $r \times c$. When the sensor readings are two-byte integers, the table size is as much as $2^{16}c$ bits. In the case of $c = 4$, this results in an additional 2^{18} -bit space requirement. As the current generation of sensor nodes usually has near or even more than hundreds of kilobytes, for ordinary sensors the effort of storing the rdOPE table can be deemed affordable as well.

VERIFIABLE, PRIVACY-ASSURED, AND ACCURATE SIGNAL COLLECTION

The cloud may report to an RCU incomplete biomedical signals for several reasons. For example, the cloud may be malicious, colluding with the adversary to deliberately ignore specific biomedical signals in an attempt to let important people miss proper diagnosis and treatment. On the other hand, the cloud can be benign but accidentally report incomplete data due to an unexpected misconfiguration/error. Here, inspired by our previous work [15], we propose a dummy reading-based compression framework, where a *virtual line segment approach* is proposed to reduce the communication overhead and at the same time keep the benefit of detecting data incompleteness, at the expense of increased computation burden on the RCU side. In particular, we propose verifiable, privacy-assured, and accurate signal collection (VPAA) by integrating the rdOPE scheme with the virtual line segment approach. In this way, the data privacy of VPAA is guaranteed by rdOPE, while data completeness is achieved by the virtual line segment containing both genuine and dummy readings.

VPAA works as follows. Let η be a system parameter denoting the difference between the maximum and minimum encrypted readings within an epoch. Note that the parameter η will be different with different types of sensors. Let $e_{i,j}$ be the j th signal reading of s_i . Then s_i con-

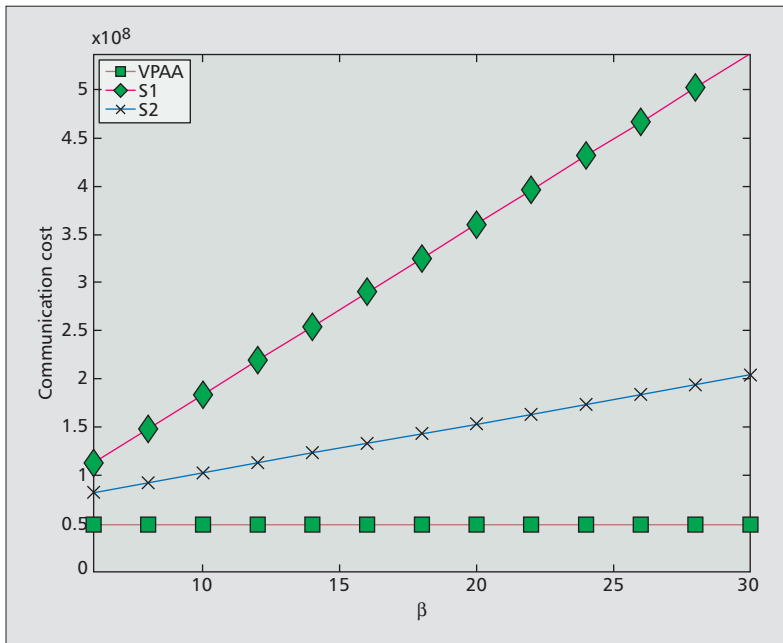


Figure 3. The communication cost of straightforward solution S_1 , straightforward solution S_2 , and the proposed VPAA scheme with different β s.

constructs a virtual line segment $\mathcal{L}_i = \langle \mathcal{L}_{i,L}, \mathcal{L}_{i,U} \rangle$ with $\mathcal{L}_{i,L}$ selected uniformly at random from $[e_{i,1} - \eta, e_{i,1}]$ and $\mathcal{L}_{i,U}$ selected uniformly at random from $[e_{i,\beta}, e_{i,\beta} + \eta]$, where $\mathcal{L}_{i,L}$ and $\mathcal{L}_{i,U}$ are two endpoints of \mathcal{L} . Obviously, \mathcal{L}_i contains all of the encrypted outputs. The purpose of \mathcal{L}_i is to anonymize $e_{i,1}, \dots, e_{i,\beta}$; the elements in \mathcal{L}_i other than the genuine encrypted readings can be thought of as dummy readings. The primary purpose of using dummy readings in such an implicit form is to anonymize the genuine readings such that the adversary cannot know which are genuine and therefore has no clue which to erase.

The use of such a virtual line segment achieves the compression task; its advantage is that there is no need to store and send each element in \mathcal{L}_i explicitly. In other words, since both genuine and dummy readings are compressed to a virtual line segment only, two numeric values, $\mathcal{L}_{i,L}$ and $\mathcal{L}_{i,U}$, are sufficient to represent each element in \mathcal{L}_i . Nonetheless, we still need additional information enabling the RCU to identify which elements are genuine. To achieve this goal, cryptographic keyed hashes of genuine readings are associated with the virtual line segment. As a result, s_i submits to the cloud $i, \mathcal{L}_i, \mathbb{H}^{(i)}, h_{\tilde{k}_i}(i | \mathcal{L}_i | \mathbb{H}^{(i)})$, where $\mathbb{H}^{(i)} = \{h_{\tilde{k}_i}(e_{i,1}), \dots, h_{\tilde{k}_i}(e_{i,\beta})\}$ denotes a set of hashes calculated from genuine individual readings. Although no dummy readings are used explicitly, the use of the virtual line segment containing dummy readings and the cryptographic keyed hash still prevent the genuine readings from intentional or accidental removal from the received signal readings, with the additional benefit of significantly reducing the number of readings that need to be transmitted.

When an RCU needs to acquire the biomedical signals from s_i , it issues a query to the cloud. The cloud submits $i, \mathcal{L}_i, \mathbb{H}^{(i)}, h_{\tilde{k}_i}(i | \mathcal{L}_i | \mathbb{H}^{(i)})$ to the RCU.

We can see that a number of unnecessary

readings included in the virtual line segment will also be returned. Nonetheless, the main focus of VPAA is the successful compression of individual sampled signal readings at the cost of the increased computation burden on the RCU. Hence, the RCU may extract the genuine signal readings by calculating the hash of each element in the received \mathcal{L}_i and checking whether the corresponding hash appears in $\mathbb{H}^{(i)}$.

A concrete example showing how VPAA works can be as follows. Assume that the sensor s_4 has $e_{4,1} = 2, e_{4,2} = 3, e_{4,3} = 4$, and $\eta = 5$. The virtual line segment \mathcal{L}_4 could be $\mathcal{L}_4 = [1, 9]$. It follows that the sensor submits to the cloud 4, $[1, 9], \mathbb{H}^{(4)}, h_{\tilde{k}_4}(4 | \mathcal{L}_4 | \mathbb{H}^{(4)})$, where $\mathbb{H}^{(4)} = \{h_{\tilde{k}_4}(e_{4,1}), \dots, h_{\tilde{k}_4}(e_{4,\beta})\}$. When the RCU receives 4, $[1, 9], \mathbb{H}^{(4)}, h_{\tilde{k}_4}(4 | \mathcal{L}_4 | \mathbb{H}^{(4)})$ from the cloud, it checks whether the calculated hash of each element in $[1, 9]$ can find a match in the received $\mathbb{H}^{(4)}$. The element belongs to genuine readings if so, and is a dummy one otherwise.

The privacy guarantee can be achieved by the use of rDOPE. In addition, the virtual line segment contains all of the signal readings. The use of cryptographic keyed hash ensures the integrity of the virtual line segment. Thus, while the readings are removed from the virtual line segment, the detection probability is close to 1.

In VPAA, the encrypted readings are replaced by a virtual line segment. s_i only needs to submit the parameters for representing the line segment and the necessary verification materials to the RCU. Hence, the communication cost $\ell_{id} + 2\ell_d + (\beta + 1)\ell_h$, where ℓ_{id} , ℓ_d , and ℓ_h denote the numbers of bits required for representing the ID, signal reading, and hash, respectively.

The simulation result for communication cost with different β s is shown in Fig. 3. In Fig. 3, we can see that due to the compression benefit of using the virtual line segment approach, the communication cost is minimal compared to the straightforward solutions S_1 and S_2 , which do not have signal reading compression. An obvious observation can also be made: as the number β of signal readings sensed increases, more communications can be saved in VPAA since only a fixed number of transmitted bits are required.

DISCUSSION

One may argue that the combined use of AES encryption and the virtual line segment approach can also achieve the same guarantee of data privacy and completeness. From the functionality point of view, the combined use of those two does offer the same functionality. Nevertheless, AES has the fixed 128-bit output length and high-entropy output. Therefore, an encryption output can be thought of as an almost uniformly random point over a range $[0, 2^{128} - 1]$. In this sense, the virtual line segment could be very lengthy. Thus, although AES is applicable in our WBSN setting, rDOPE is still necessary in the VPAA design because of its order preservation characteristic that leads to a shorter virtual line segment.

Actually, VPAA was implemented on TelosB motes on top of TinyOS (CPU: TI MSP430F1611; ROM: 48 KB + 256 B; RAM: 10

	ROM	RAM	CPU
VPAA	14448 bytes	640 bytes	1333.753 mJ

Table 1. Summary of prototype implementation.

KB; radio chipset: ChipCon CC2420). Our program code was also run on TOSSIM in TinyOS 1.1.15 to evaluate the energy consumption. In our setting, together with the AES encryption function in a CC2420 chipset, CBC-MAC mode is used to implement the hash function with $\beta = 10$ readings.

CONCLUSION

For the signal collection problem in cloud-assisted WBSNs, a particular application of WSNs, we propose a verifiable, privacy-assured, and accurate data collection (VPAA) scheme. VPAA has two salient security guarantees: data privacy and data completeness. The main focus of VPAA is its ability to significantly reduce the communication cost of sensors at the cost of increased computation cost on the powerful RCU side.

A possible future research direction is to reduce the computation overhead on the RCU sides. In particular, in our VPAA scheme, although the communication cost has been reduced via the proposed virtual line segment approach, the computation overhead on the RCU side turns out to be increased. This is because the RCU needs to check which readings are genuine and which are dummies. The anonymization level is proportional to the increased length of the virtual line segment, which is also proportional to computation time for the RCU. Therefore, a possible future research direction is to develop a technique to simultaneously reduce both communication and computation overhead.

ACKNOWLEDGMENT

This research was partly funded by the Ministry of Science and Technology of the R.O.C. under grants MOST 102-2218-E-155-006-MY2, 103-2221-E-197-018 and 101-2221-E-197-008-MY3, National Nature Science Foundation of China under No. 61170296 & 60873241, and the Program for New Century Excellent Talents in University under Grant No. 291184.

REFERENCES

- [1] J. Ko *et al.*, "Wireless Sensor Networks for Healthcare," *Proc. IEEE*, vol. 98, no. 11, 2010, pp. 1947–60.
- [2] M. Shoaib and H. Garudadri, "Digital Pacer Detection In Diagnostic Grade ECG," *IEEE Int'l. Conf. E-health, Networking, Application & Services*, 2011.
- [3] A. Abbas and S. U. Khan, "A Review on the State-of-the-Art Privacy-Assured Approaches in the E-Health Clouds," *IEEE J. Biomed. Health Informatics*, vol. 18, no. 4, July 2014, pp. 1431–41.
- [4] B. Sheng and Q. Li, "Verifiable Privacy-Assured Sensor Network Storage for Range Query," *IEEE Trans. Mobile Computing*, vol. 10, no. 9, Sept. 2011, pp. 1312–26.
- [5] J. Shi, R. Zhang, Y. Zhang, "A Spatiotemporal Approach for Secure Range Queries in Tiered Sensor Networks," *IEEE Trans. Wireless Commun.*, vol. 10, no. 1, Jan. 2011, pp. 264–73.

- [6] C.-M. Yu *et al.*, "Practical and Secure Multidimensional Queries in Tiered Sensor Networks," *IEEE Trans. Info. Forensics Security*, vol. 6, no. 2, Jun. 2011, pp. 241–55.
- [7] F. Chen and A. X. Liu, "Privacy- and Integrity-Preserving Range Queries in Sensor Networks," *IEEE/ACM Trans. Networking*, vol. 20, no. 6, Dec. 2012, pp. 1774–87.
- [8] Y. Yi *et al.*, "A Digital Watermarking Approach to Secure and Precise Range Query Processing in Sensor Networks," *IEEE INFOCOM*, 2013.
- [9] C. Luo, F. Wu, J. Sun, and C. W. Chen, "Compressive Data Gathering for Large-Scale Wireless Sensor Networks," *ACM Int'l. Conf. Mobile Computing and Networking*, 2009.
- [10] D. He *et al.*, "Lightweight and Confidential Data Discovery and Dissemination for Wireless Body Area Networks," *IEEE J. Biomed. Health Info.*, vol. 18, no. 2, March 2014, pp. 440–48.
- [11] D. He *et al.*, "Secure and Lightweight Network Admission and Transmission Protocol for Body Sensor Networks," *IEEE J. Biomed. Health Informatics (J-BHI)*, vol. 17, no. 3, May 2013, pp. 664–74.
- [12] Y. Tong *et al.*, "Cloud-Assisted Mobile-Access of Health Data With Privacy and Auditability," *IEEE J. Biomed. Health Informatics*, vol. 18, no. 2, Mar 2014, pp. 419–29.
- [13] R. Agrawal *et al.*, "Order Preserving Encryption for Numeric Data," *ACM Int'l. Conf. Management of Data*, 2004.
- [14] A. Boldyreva *et al.*, "Order-Preserving Symmetric Encryption," *Annual Int'l. Conf. Theory and Applications of Cryptographic Techniques*, 2009.
- [15] C.-M. Yu *et al.*, "Top-k Query Result Completeness Verification in Tiered Sensor Networks," *IEEE Trans. Info. Forensics and Security (T-IFS)*, vol. 9, no. 1, Jan 2014, pp. 109–24.

BIOGRAPHIES

CHIA-MU YU (chiamuyu@saturn.yzu.edu.tw) received his Ph.D degree from the Computer Science Group, Department of Electrical Engineering, National Taiwan University in 2012. He was a research assistant with the Institute of Information Science, Academia Sinica from 2005 to 2010, a visiting scholar at Harvard University from September 2010 to September 2011, a visiting scholar at Imperial College London from January 2012 to September 2012, and a postdoc researcher at IBM Thomas J. Watson Research Center from October 2012 to July 2013. He is currently an assistant professor with the Department of Computer Science and Engineering, Yuan Ze University. Since 2014, he has served as an Associate Editor of three journals: *IEEE Access*, *Security and Communication Networks*, and *Gate to Multimedia Processing*. His research interests include cloud storage security, sensor network security, and smart grid security.

CHI-YUAN CHEN [M] (chiyuan.chen@ieee.org) received his Ph.D. degree in electrical engineering from National Dong Hwa University in 2014. He is currently an assistant professor with the Department of Computer Science and Information Engineering of National Ilan University. He has been the Associate Editor-in-Chief for the *Journal of Internet Technology* since February 2014. His research interests include mobile communication, network security, and quantum computing. He is a member of ACM.

HAN-CHIEH CHAO [SM] (hcc@niu.edu.tw) is a joint appointed Chair Professor with the Department of Electronic Engineering and the Institute of Computer Science & Information Engineering, where also serves as the president of National Ilan University, I-Lan, Taiwan, R.O.C. He served as director of the Computer Center for the Ministry of Education from September 2008 to July 2010. He is the Editor-in-Chief of *IET Networks*, the *Journal of Internet Technology*, the *International Journal of Internet Protocol Technology*, and the *International Journal of Ad Hoc and Ubiquitous Computing*. He has served as a Guest Editor for *ACM Mobile Networking and Applications*, the *IEEE Journal on Selected Areas in Communications*, *IEEE Communications Magazine*, *Computer Communications*, *IEEE Proceedings on Communications*, the *Computer Journal*, *Telecommunication Systems*, *Wireless Personal Communications*, and *Wireless Communications & Mobile Computing*. He is a Fellow of IET (IEE). He is a Chartered Fellow of the British Computer Society.

The anonymization level is proportional to the increased length of the virtual line segment, which is also proportional to computation time for the RCU. Therefore, a possible future research direction is to develop a technique to simultaneously reduce both communication and computation overhead.

Secure Communication for Underwater Acoustic Sensor Networks

Guangjie Han, Jinfang Jiang, Ning Sun, and Lei Shu

ABSTRACT

UASNs are widely used in many applications, and many studies have been conducted. However, most current research projects have not taken network security into consideration, despite the fact that a UASN is typically vulnerable to malicious attacks due to the unique characteristics of an underwater acoustic communication channel (e.g., low communication bandwidth, long propagation delays, and high bit error rates). In addition, the significant differences between UASNs and terrestrial wireless sensor networks entail the urgent and rapid development of secure communication mechanisms for underwater sensor nodes. For the above mentioned reasons, this article aims to present a somewhat comprehensive survey of the emerging topics arising from secure communications in UASNs, which naturally lead to a great number of open research issues outlined afterward.

INTRODUCTION

Over the past few years, there has been a rapidly growing body of research on underwater acoustic sensor networks (UASNs) due to their wide applications in many underwater scenarios, including oceanographic data collection, assisted navigation, disaster prevention, and so on. In most of the above applications, UASNs are deployed in unattended and even hostile environments, so secure communication among sensor nodes must be considered to ensure the efficiency of UASNs. Recent research on UASNs has mainly focused on network construction or protocol management, while few efforts have been made in network security. Thus, in this article, we study the security issues for UASNs.

Recently, considerable research has been done on secure communication in terrestrial wireless sensor networks (TWSNs). However, the unique characteristics of UASNs make it impossible to directly use these secure communication mechanisms for TWSNs. The main differences between UASNs and TWSNs are listed in Table 1. High-frequency radio signals attenuate rapidly underwater. Therefore, radio communication, which is extensively used in TWSNs, can-

not work efficiently underwater. Instead, acoustic communication emerges as a better choice for underwater communications. Acoustic communication has several handicaps, such as low bandwidth, high propagation delay and high bit error rate. Moreover, the acoustic channel has low link quality mostly due to the multi-path propagation and time variability of the medium. In addition to the above mentioned points, in UASNs, all sensor nodes freely drift with ocean currents, which results in dynamic 3D network configurations. Furthermore, the production cost of underwater sensor nodes is often more than that of regular terrestrial sensor nodes. Therefore, a dense deployment is not an affordable choice. The large-scale and sparse structure of UASNs makes it easy to attack and hard to defend.

The contributions of this article are listed as follows:

- The secure communication protocols in each layer for UASNs are carefully analyzed in the article. To the best of our knowledge, this is the first work to study secure communication for UASNs layer by layer.
- The performance of the existing secure communication protocols is highlighted in terms of their strengths and weaknesses. In addition, the malicious attacks and their countermeasures in each layer are identified.

The rest of the article is organized as follows. The following section provides a detailed survey of the secure communication protocols. Their advantages and limitations are then compared. Furthermore, the malicious attacks and related countermeasures are summarized. Finally, open research problems are outlined, and conclusions are made in the final section.

SECURE COMMUNICATION PROTOCOLS FOR UASNs

Generally, malicious attack attempts on UASNs can be classified into two types: attacks on sensor nodes and attacks on network protocols. The first kind of malicious attack is the most effective method to damage UASNs. However, this method is not practical in real

Guangjie Han, Jinfang Jiang, and Ning Sun are with Hohai University.

Lei Shu is with Guangdong University of Petrochemical Technology.

applications. Since underwater sensor nodes in UASNs are always sparsely deployed, it is hard to simultaneously destroy several nodes. Therefore, the first kind of malicious attack is not very destructive unless the destroyed/damaged nodes are key nodes such as the sink node. Relatively, attacks on network protocols, especially communication protocols, occur more frequently in any kind of UASN. Once the communication protocols are broken, the whole network is useless. Therefore, in this article, we focus on secure communication protocols for UASNs in a bottom-up-layer manner. As shown in Fig. 1, the bottom-up network architecture of UASNs consists of five layers: physical layer, data link layer, network layer, transport layer, and application layer.

SECURE COMMUNICATION IN THE ACOUSTIC PHYSICAL LAYER

Compared to radio communications in TWSNs, underwater communications over acoustic channels are much more vulnerable to malicious attacks. The most common malicious attack in the physical layer is a jamming attack, which is a popular type of denial of service (DoS) attack. In a jamming attack, a malicious node that tries to disrupt the communications between sensor nodes is referred to as a jammer. A jammer interferes with the physical channels of normal sensor nodes by sending many useless signals on the same frequency band. Since the frequency bands of acoustic communication are narrow (from Hertz to kilohertz), UASNs are very vulnerable to jamming attacks.

Many existing solutions for jamming detection in TWSNs are not suitable for UASNs. For example, in TWSNs, other signals (e.g., infrared waves) can be alternatively used to avoid radio frequency jamming. This method cannot work well in UASNs since infrared waves are severely attenuated underwater. Therefore, in [1], S. Misra *et al.* proposed an UnderWater Jamming Detection Protocol (UWJDP) to detect and mitigate jamming attacks. It is assumed that a jamming attack is done by injecting a lot of malicious packets into the network at a high rate. There are three phases in UWJDP, neighbor discovery, jamming detection, and jammed area mapping. In the jamming detection phase, packet send ratio (PSR), packet delivery ratio (PDR), and energy consumption amount (ECA), are used to detect jamming. UWJDP considers that jamming always decreases PSR and PDR, while increasing ECA. However, PSR and PDR may also decrease due to the interruption of the communication link. In this case, the sender consumes higher ECA for packet retransmission. Hence, we argue that UWJDP is not capable of accurately detecting jamming. Furthermore, this scheme relies on full-dimensional location information of sensor nodes, which is difficult to obtain due to the dynamic underwater environment. Obtaining location information requires special equipment, such as GPS, or algorithms, which introduces high communication overheads and additional energy consumption.

In [2], the effects of jamming attacks on UASNs are studied based on real-world field

Features	UASNs	TWSNs
Architecture	Mostly 3D	Mostly 2D
Topology	Highly dynamic due to water current	Static or slightly dynamic
Node movement	Move with water current	Staying static
Deployment	Sparse deployment	Dense deployment
Communication medium	Acoustic or optical signals	Radio communication is extensively used
Speed of medium	Acoustic velocity in water is about 1500 m/s	The speed of radio frequency in the air is 3×10^8 m/s
Propagation delay	High due to the low-speed acoustic communication	Low
Frequency	Low frequency (Hz, kHz) because signal with high frequency is quickly absorbed in water	High frequency (MHz, GHz)
Bandwidth	Short communication distance has higher bandwidth	Bandwidth does not change with different distances
Link quality	High bit error rate (BER) and packet loss rate (PLR)	Relatively better

Table 1. The differences between UASNs and TWSNs.

tests, where malicious nodes are categorized into two categories: dummy (signal) jammer and smart (deceptive) jammer. The first kind of jammer does not know anything about the network, and uses noise to jam the acoustic communication channel. The second kind of jammer knows the network protocol and can pretend to be a legitimate node to launch malicious attacks. Network performance under jammer attacks is simulated, and the simulation results show that jamming attacks on UASNs can be launched easily, which may drastically degrade the performance of the network. The work in [2] provides new sights for researchers to further study jamming attacks.

In addition to the malicious jamming schemes mentioned above, there is also friendly jamming, which can be used to improve network security. For example, in [3], the authors proposed a secure underwater communication scheme based on cooperative friendly jamming, which is Jamming through Analog Network Coding (J-ANC). Unlike conventional cooperative jamming secure schemes, which employ artificial noise as a jamming source, J-ANC utilizes the same spreading code used by a legitimate link. The packet transmitted by the friendly jammer is known to the legitimate node, but not to the eavesdropper. Therefore, after obtaining the interfering packet, the legitimate node can correctly decode the received packet, while the eavesdropper cannot.

From the above-mentioned recent works, we can figure out that the security research in the

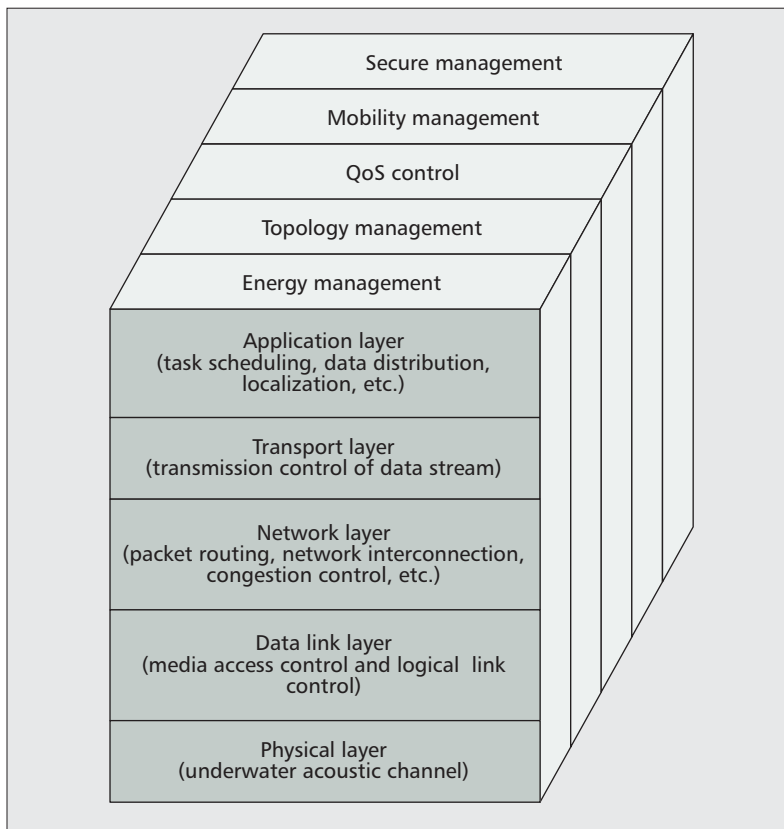


Figure 1. The network architecture of UASNs.

acoustic physical layer is still in its infancy. For jamming attacks, the open issues can be concluded to be:

- Jamming attacks need detailed classification for further study. For example, a jamming attack can be classified into three types: continuous, pulsed, and reactive jamming.
- In continuous jamming attacks, the jammer continually sends many useless packets and thus exhausts its power in a short time. Therefore, legitimate nodes can switch to sleep mode and wake up until the attack is over.
- In pulsed jamming attacks, the jammer can alternate between sleep and work mode. In this case, the jammer can save its energy consumption and randomly corrupt the network. Therefore, legitimate nodes can communicate with each other in the interval of the jammer's sleep time and report the attack to the base station.
- In reactive jamming attacks, the jammer can switch to sleep or work mode according to legitimate nodes' working conditions. If a transmission is detected, the jammer starts disturbing. Otherwise, it will turn to idle. In this case, legitimate nodes can use multi-frequency acoustic channels to transmit packets.

SECURE MAC PROTOCOLS IN THE DATA LINK LAYER

In the data link layer, an efficient and secure medium access control (MAC) protocol enables multiple sensor nodes to share a common wireless medium in an efficient and fair way. In the

MAC layer, precise and secure time synchronization is essential to adjust sensor nodes' sleep-wakeup mode. Therefore, secure time synchronization has been carefully studied. For example, in [4], a water-quality monitoring sensor network (WATER) is proposed, where a lightweight secure time synchronization mechanism is studied. The proposed correlation-based security WATER time synchronization mechanism (WATERSync) can efficiently detect outlier timestamp data. First, based on depth information, it is found that the acoustic propagation delays between two neighbor sensor nodes fit a normal distribution. That is, the timestamps between two neighbor nodes are correlated with each other. According to the coefficient, an outlier timestamp data can be detected. If sufficient outlier data are detected, the corresponding sensor node is considered to be malicious. However, WATERSync does not work well in a dynamic underwater environment. First, in order to detect malicious nodes, each node needs many timestamp data from neighbor nodes. Sensor nodes are sparsely deployed, and packet loss rates are relatively high in UASNs; therefore, the correlation judgment will lose its efficiency without enough timestamp data. In addition, due to long propagation delay and node mobility, the proposed WATERSync scheme is not suitable for real-time secure time synchronization.

In order to improve the performance of WATERSync, the authors in [5] propose a two-step security UASN synchronization model: correlation test, and a statistical reputation and trust model. In [4], only depth information is used; therefore, WATERSync can only achieve vertical clock synchronization. While the proposed scheme in [5], which is called secure vertical and horizontal synchronization (SVHS), can ensure both horizontal and vertical time synchronization.

In [6], a cluster-based secure synchronization (CLUSS) protocol is proposed for UASNs. CLUSS consists of three phases: an authentication phase, inter-cluster synchronization phase, and intra-cluster synchronization phase. In the authentication phase, all sensor nodes need to be authenticated to each other. First, cluster heads are authenticated to the beacon nodes on the water surface, and then the ordinary nodes are authenticated to their cluster heads. In the inter-cluster synchronization phase, cluster heads synchronize themselves by using beacons, while ordinary nodes synchronize with corresponding cluster heads in the intra-cluster synchronization phase. In order to improve accuracy of time synchronization, the propagation delays of downlink and uplink caused by node movement are carefully analyzed. In addition, part of the inter-cluster synchronization phase and intra-cluster synchronization phase can be concurrently executed to save energy consumption of sensor nodes and reduce the number of transmitted packets. Compared to traditional time synchronization protocols, CLUSS can significantly reduce energy consumption as well as synchronization errors. However, due to the variable and long propagation delays of acoustic communication, the time needed for sensor nodes' synchronization should be investigated to improve

the real-time efficiency of the secure time synchronization mechanism. Furthermore, since UASNs typically have limited resources, the secure time synchronization schemes should be designed with lower computational complexity and lower communication overhead.

EFFICIENT ACOUSTIC ROUTING PROTOCOLS IN THE NETWORK LAYER

As in TWSNs, secure routing at the network layer is a major concern in UASNs. The unique characteristics of the acoustic communication channel present a significant challenge to designing an efficient secure routing protocol. Another challenge for routing protocol in UASNs is handling node mobility. Secure routing protocols for TWSNs are mainly designed for stationary scenes. In UASNs, however, most underwater nodes are mobile, and the network topology changes dramatically with ocean current. Thus, the existing routing algorithms for TWSNs are unsuitable for UASNs. Although many routing protocols have been proposed for UASNs, few of them have been designed with consideration of security.

In the network layer, malicious attacks against routing protocols mainly include the following types: flooding attacks, Sybil attacks, wormhole attacks, sinkhole attacks, and black hole attacks. In [7], a new mechanism, distributed visualization of wormhole (Dis-VoW), is proposed for UASNs to detect wormhole attacks. In Dis-VoW, each sensor node first reconstructs its local network layout using the multidimensional scaling method. Then, based on the distortions in edge lengths and angles among neighbor sensor nodes, the wormhole nodes can be detected. However, Dis-VoW is not feasible in large-scale UASNs or high-density UASNs, since the reconstruction of network layout consumes a large amount of energy, especially for large-scale or high-density underwater environments.

In [8], the authors present a suite of novel protocols to enable wormhole-resilient secure neighbor discovery (WSND) in UASNs. The protocols are based on the direction of arrival (DoA) estimation of acoustic signals. The DoA method is a basic functionality that is easily available in underwater environments. Simulation results demonstrate that the proposed WSND protocol can efficiently thwart a wormhole attack. In addition, compared to some traditional routing protocols, which usually make hard requirements on accurate time synchronization, WSND is much easier to implement in UASNs.

In [9], a security suite is presented for UASNs that is composed of both fixed and mobile underwater nodes. The security suite consists of a secure routing protocol and a set of cryptographic primitives (SRCP). In order to protect confidentiality and integrity of acoustic communication, both the unique characteristics and constraints of the acoustic channel are taken into account to design the security suite. To the best of our knowledge, this is the first work that provides an efficient solution to protect the integrity and confidentiality of UASNs.

Based on the above mentioned secure routing protocols, it can be concluded that:

- It is necessary to develop novel secure routing protocols against wormhole attacks. For example, in [7], if the wormhole node can falsify the timestamp information, the Dis-VoW cannot efficiently detect the wormhole attack. In [8], the WSND protocol is also heavily affected by the distance-measuring errors between neighbor sensor nodes.
- In the existing secure routing protocols, only wormhole attacks are studied. The other insider attacks (Sybil attacks, sinkhole attacks, flooding attacks, etc.), which significantly threaten the security of UASNs, need further investigation.
- Besides the insider attacks, the outside intruders also need to be investigated.
- The development of reputation-based or trust-based routing protocols is suggested, where the behaviors of neighbor nodes can be detected and analyzed. Sensor nodes with good behavior are assigned with high trust values, while sensor nodes with bad behavior are assigned with low trust values. In this case, malicious nodes or selfish nodes that do not cooperate in routing can be efficiently eliminated.

RELIABLE DATA TRANSFER PROTOCOLS IN THE TRANSPORT LAYER

In the Internet, User Datagram Protocol (UDP) and Transmission Control Protocol (TCP) are the two most commonly used protocols in the transport layer. UDP is usually used for real-time service, while TCP is used to ensure reliable data transmission. Thus, a TCP-based approach is much more suitable for UASNs. However, secure TCP has not been investigated in UASNs.

There are always two kinds of methods to ensure secure data transmission: authentication and data encryption. In [10], several of the most popular digital signature schemes for end-to-end authentication in UASNs are evaluated in terms of energy consumption. After authentication, a secure symmetric key is always used to encrypt all the transmitted data. It is generally known that authentication and data encryption are strongly related to each other. First, the secret keys established in the process of data encryption can be used to perform sensor nodes' authentication. Conversely, once the sensor nodes verify each other's authenticity, they can also establish one or more secret keys to securely encrypt the transmitted data. Duo to the unique characteristics of the underwater acoustic channel, the conventional solutions for data encryption should be updated for UASNs. In [11], a key generation system (KGS) is presented which is robust to variable underwater environments. Unlike the conventional solutions, it exploits reciprocity, deep fades, a randomness extractor, and robust secure fuzzy information to generate the key. The key is generated based on the characteristics of the underwater channel. Therefore, the KGS is secure against adversaries who do not know the location of the deep fades but know the number of deep fades.

In the data link layer, an efficient and secure MAC protocol enables multiple sensor nodes to share a common wireless medium in an efficient and fair way. In the MAC layer, precise and secure time synchronization is essential to adjust sensor nodes' sleep-wakeup mode. Therefore, secure time synchronization has been carefully studied.

SECURE APPLICATIONS IN THE APPLICATION LAYER

The security problem in the application layer is always related to the practical applications. Take secure localization; for example, it is generally known that localization is a very important issue in UASNs. Many applications (e.g., underwater environment monitoring) need sensor nodes' location information. Therefore, it is urgent to research the secure localization problem.

In [12], a secure localization algorithm based on a trust mechanism (SLTM) is proposed for UASNs. In order to reduce the influence of malicious location information, a trust model is built based on Beta distribution to detect illegal beacon nodes. Meanwhile, in order to reduce the influence of the instability of the underwater acoustic channel on the trust evaluation process, a trust filter mechanism (TFM) is adopted to

calculate the trust values of beacon nodes. Simulation results show that the SLTM can recognize malicious beacon nodes quickly. The accuracy and security of localization are greatly improved. However, an SLTM is not applicable in a real underwater environment since an SLTM does not take node mobility into consideration. In addition, the localization of sensor nodes in the presence of malicious attacks such as Sybil or wormhole attacks is not considered. The detection rate of malicious nodes and the false detection rate of normal nodes are not investigated.

Against malicious attacks, a reputation or trust model was recently suggested as an effective security mechanism. In TWSNs, many trust models have been proposed [13]. However, the unique characteristics of UASNs make it impossible to directly use these trust models. In addition, when trust models defend against malicious nodes, they may also be attacked by adversaries.

Protocols	Methodology	Anti-attack	Anti-mobility	Energy efficiency	Advantages	Limitations
UWJDP [1]	PSR, PDR ECA	Jamming attack	No	No	Efficiently detect jamming attack	High communication overheads
J-ANC [3]	Friendly jamming	Eavesdropper	No	No	Low computational complexity	High communication overheads
WATER-Sync [4]	A normal distribution of propagation delays	Resist insider attacks	No	No	Depth information taken into consideration	Can only achieve vertical clock synchronization
SVHS [5]	Correlation test and statistical trust model	Resist insider attacks	No	No	Ensure both horizontal and vertical time synchronization	The trust model is not practical in real underwater environments
CLUSS [6]	Authentication	—	No	No	The first time synchronization for cluster-based UASNs	High computational complexity and communication overheads
Dis-VoW [7]	Visualizing the distortions in edge lengths and angles among neighbor sensor nodes	Wormhole attacks	No	No	Suitable for large scale UASNs	Is not suitable for mobile network environment
WSND [8]	DoA estimation of acoustic signals	Wormhole attacks	No	Yes	Low computational complexity	Cannot detect adjacent wormhole nodes
SRCP [9]	Cryptography	Data modification attack	Yes	No	Protecting confidentiality and integrity of packets	High computational complexity
KGS [11]	Key generation	Robust against adversaries who do not know the location of the deep fades but know the number of them	No	No	Protect data integrity, privacy and confidentiality	High computational complexity
SLTM [12]	Trust filter mechanism	Malicious anchor nodes	No	No	Improve localization accuracy with malicious anchor nodes	Cannot handle node mobility

Table 2. Comparison of secure communication protocols.

Currently, few trust models have been designed with consideration of security.

DISCUSSION AND OPEN RESEARCH ISSUES

As shown in Table 2, the performance of the secure communication protocols is compared in terms of the following six aspects: methodology, anti-attack, anti-mobility, energy efficiency, advantages, and limitations. It can be concluded that the security research of UASNs is still in the initial stage. In the design of communication protocols, only limited kinds of malicious attacks (e.g., jamming attacks, wormhole attacks) are taken into consideration. In addition, most studies do not take node mobility into consideration, and their computational complexity is relatively high, which ultimately introduces high communication overheads and energy consumption. In order to obtain an efficient and secure communication protocol, the following basic requirements should be satisfied.

Security: As a data-centric sensor network, the designed secure communication protocol for UASNs should first satisfy the security requirement. That is to say, whenever the network is attacked, the designed protocol needs to ensure that the critical data information is correct and available to authorized users, while not accessible to unauthorized users.

Robustness: On one hand, malicious nodes need to be efficiently and promptly detected or even eliminated from the network. On the other hand, the network should stay workable and operational under any kind of hostile attack.

Energy efficiency: For UASNs, energy efficiency becomes a significant metric since UASNs are energy limited. An energy-efficient secure communication protocol means that the protocol spends the least possible energy on reliable data transmission. Moreover, the secure communication protocol should be robust against nodes' mobility. However, nearly no current works consider how to overcome the above-mentioned problems.

Lightweight Pprotocol: A secure communication protocol of UASNs should be as simple as possible (i.e., without dependence on a special software or hardware), since underwater sensor nodes are resource-limited with low power, limited memory space, and communication bandwidth.

Based on the above analysis, the malicious attacks in each layer and the proper countermeasures are summarized in Table 3.

CONCLUSIONS

In this article, we have focused on the unique characteristics of the acoustic communication channel, and studied possible attacks and countermeasures of communication protocols in UASNs. First, we present a detailed survey of secure communication protocols for UASNs. The secure communication protocols in the five layers are analyzed in a bottom-up manner. Furthermore, we compare the protocols by discussing their methodology, energy efficiency,

Layers	Attacks	Countermeasures
Physical layer	Tampering	The perception mechanism for physical damage, encryption algorithm, etc.
	Jamming	Sleep-wakeup model, multi-frequency communication, using different transmission priority, etc.
Data link layer	Collision	Forward error correction (FEC) code
	Exhaustion	Limiting the transmission speed and retransmission times of packets
	Unfairness	Avoiding the use of long packets, redistributing transmission priority of packets, etc.
Network layer	DoS attacks	Detection of energy consumption
	Selective forwarding	Multi-path routing, reputation and trust model, etc.
	Sybil	Identity authentication of sensor nodes
	Wormhole	Construction of network topology
	Sinkhole	Traffic monitoring, identity authentication, multi-path routing, etc.
Transport layer	Flooding	Limit the broadcast range of sensor nodes

Table 3. Malicious attacks and their countermeasures in each layer.

robustness against attacks and mobility, advantages, and limitations. In summary, it is found that security research for UASNs is still in its infancy. The field of UASNs is rapidly growing, but many challenges remain wide open for future investigation. Below, we list some representative issues:

- In order to improve the security of communication, a reputation or trust model can be considered. By analyzing sensor nodes' communication behavior, rather than identifying the type of malicious behavior, we assign different trust values to different nodes to promote sensor nodes' collaboration.

- Based on the survey of the secure communication protocols, it can be concluded that most current research is only suitable for small and static UASNs. However, in real applications, sensor nodes freely drift with ocean current. Considering the mobility of sensor nodes can improve the performance of communication protocols.

- To prolong network lifetime, only a fraction of sensor nodes need to participate in secure communication. In UASNs, adjacent sensor nodes share common sensing tasks. This implies that not all the sensor nodes are necessary to transmit data packets during the whole system lifetime. Therefore, if we can schedule sensor nodes to work alternatively, the network lifetime can be prolonged.

- Due to the long propagation delay of the acoustic channel and the mobility of underwater nodes, a UASN can be seen as a delay-tolerant network. In this regard, how to design secure

As a data-centric sensor network, the designed secure communication protocol for UASNs should first satisfy the security requirement. That is to say, whenever the network is attacked, the designed protocol needs to ensure that the critical data information is correct and available to authorized users, while not accessible to unauthorized users.

communication protocols in a delay-tolerant UASN is another issue that needs further attention.

• In current research works, malicious attacks are separately considered in each layer of the communication protocol. However, a cross-layer design is suggested as a better choice. For example, as analyzed before, when a pulsed jamming attack happens in the physical layer, the sleep scheduling scheme in the MAC layer can be adopted to cause legitimate nodes to communicate with each other in the interval of the jammer's sleep time to improve the network throughput.

• Secure communication protocols are always designed based on interactions among sensor nodes to detect malicious nodes and ensure the integrity of sensed data. In addition to integrity and correctness, the privacy and confidentiality of data packets should also be protected.

ACKNOWLEDGMENT

The work is supported by "Jiangsu Province Ordinary University Graduate Innovation Project, No. CXZZ13_02," "Qing Lan Project," "National Natural Science Foundation of China, Grant No. 61401107," "2014 Guangdong Province Outstanding Young Professor Project," and "Natural Science Foundation of Jiangsu Province of China, No. BK20131137."

REFERENCES

- [1] S. Misra *et al.*, "Jamming in Underwater Sensor Networks: Detection and Mitigation," *IET Commun.*, vol. 6, no. 14, Sep. 2012, pp. 2178–88.
- [2] M. Zuba *et al.*, "Vulnerabilities of Underwater Acoustic Networks to Denial-of-Service Jamming Attacks," *Security Commun. Net.*, Feb. 2012, pp. 1–11.
- [3] H. Kulhandjian, T. Melodia, and D. Koutsonikolas, "Securing Underwater Acoustic Communications through Analog Network Coding," *Proc. SECON*, June 2014, pp. 1–9.
- [4] F. Hu, S. Wilson and Y. Xiao, "Correlation-Based Security in Time Synchronization of Sensor Networks," *Proc. WCMC.*, Mar. 2008, pp. 2525–30.
- [5] F. Hu *et al.*, "Vertical and Horizontal Synchronization Services with Outlier Detection in Underwater Acoustic Networks," *Wireless Commun. Mob. Com.*, vol. 8, no. 9, Nov. 2008, pp. 1165–81.
- [6] M. Xu *et al.*, "A Cluster-Based Secure Synchronization Protocol for Underwater Wireless Sensor Networks," *Int. J. Distrib. Sensor Networks*, vol. 2014, Apr. 2014, pp. 1–13.

- [7] W. Wang *et al.*, "Visualisation of Wormholes in Underwater Sensor Networks: A Distributed Approach," *Int'l. J. Security Net.*, vol. 3, no. 1, Jan. 2008, pp. 10–23.
- [8] R. Zhang and Y. Zhang, "Wormhole-Resilient Secure Neighbor Discovery in Underwater Acoustic Networks," *Proc. 29th IEEE INFOCOM.*, Mar. 2010, pp. 2633–41.
- [9] G. Dini and A.L. Duca, "A Secure Communication Suite for Underwater Acoustic Sensor Networks," *Sensors-Basel*, vol. 12, no. 11, Nov. 2012, pp. 15,133–58.
- [10] E. Souza *et al.*, "End-to-End Authentication in Underwater Sensor Networks," *Proc. ISCC.*, July 2013, pp. 299–304.
- [11] Y. Liu, J. Jing, and J. Yang, "Secure Underwater Acoustic Communication Based on a Robust Key Generation Scheme," *Proc. 9th ICSP*, Oct. 2008, pp. 1838–41.
- [12] Y. Zhang *et al.*, "Node Secure Localization Algorithm in Underwater Acoustic Sensor Network Based on Trust Mechanism," *J. Computer Applications*, vol. 33, no. 5, May. 2013, pp. 1208–11.
- [13] G. Han *et al.*, "Management and Applications of Trust in Wireless Sensor Networks: A Survey," *J. Comp. Sys. Sci.*, vol. 80, no. 3, May. 2014, pp. 602–17.

BIOGRAPHIES

GUANGJIE HAN [M'02] (hanguangjie@gmail.com) is currently a professor with the Department of Information & Communication Systems at Hohai University, China. He finished his work as a postdoctoral researcher with the Department of Computer Science at Chonnam National University, Korea, in 2008. He has served as an Editor of *IJAHUC*, *KSII*, and *JIT*. His current research interests are security and trust management, localization and tracking, and routing for sensor networks.

JINFANG JIANG (jiangjinfang1989@gmail.com) is currently pursuing her Ph.D degree from the Department of Information & Communication Systems at Hohai University. She received her B.S. degree in information and communication engineering from Hohai University in 2009. Her current research interests are security, trust management, routing, and localization for wireless sensor networks.

NING SUN (sunn2001@hotmail.com) received her Ph.D. degree in computer science from Chungbuk National University, Korea, in February 2013. Since then, she has been a lecturer in the Department of IOT Engineering at Hohai University. Her research interests are in the design and evaluation of network architectures and protocols. She is currently investigating wireless sensor networks, the Internet of Things, and network security.

LEI SHU [M'07] (lei.shu@ieee.org) received his Ph.D. degree from the National University of Ireland, Galway, in 2010. In October 2012, he joined Guangdong University of Petrochemical Technology, China, as a full professor. Meanwhile, he is also the vice-director of the Guangdong Provincial Key Laboratory of Petrochemical Equipment Fault Diagnosis, China. He is the founder of the Industrial Security and Wireless Sensor Networks Lab. His research interests include wireless sensor networks, middleware, security, and fault diagnosis.

From Botnets to MobiBots: A Novel Malicious Communication Paradigm for Mobile Botnets

Abderrahmen Mtibaa, Khaled A. Harras, and Hussein Alnuweiri

ABSTRACT

Cyber security is moving from traditional infrastructure to sophisticated mobile infrastructureless threats. A major concern is that such imminent transition is happening at a rate far exceeding the evolution of security solutions. In fact, the transformation of mobile devices into highly capable computing platforms makes the possibility of security attacks originating from within the mobile network a reality. Today, mobile devices are capable of initiating sophisticated cyberattacks, especially when they coordinate together to form what we call a MobiBot. MobiBots differ from classical botnets in that they exploit mobile operating system vulnerabilities and the advantages of device-to-device communication to mask malicious code propagation.

INTRODUCTION

Rising trends in mobile device adoption is paving the way for the proliferation of various mobile security threats. Trend Micro, in their security report, shows that the number of malicious applications doubled in just six months to reach more than 700,000 malwares in June 2013 [1]. Malwares represent an extreme security concern, especially when they help form a coordinated security threat such as mobile botnets, which have become an emerging threat for users and network operators. In 2012, an Android-based botnet was formed when victims were convinced to install the Madden game NFL 12, which hosted a malware that used Internet relay chat (IRC) for its command and control (C&C) channel [2]. In April 2014, an Android botnet targeted financial institutions in the Middle East by infecting more than 2700 phones and intercepting more than 28,000 SMSs [3]. These growing attacks exhibit similar behaviors as their predecessors by using the same C&C channels and the same network architecture connecting infected bots to their corresponding botmaster that manages the botnet.

The growth in mobile computation, sensor, and communication capabilities makes us anticipate an imminent move toward sophisticated mobile infrastructureless-based threats. Mobile

infrastructureless services have become a necessity, pushing researchers and companies to introduce novel device-to-device communication such as Bluetooth, WiFi-Direct, LTE Direct, and fifth generation (5G) [4, 5]. In October 2014, Hong Kong protesters massively used an infrastructureless application, FireChat, to organize themselves away from the “government-controlled” networks [6]. In addition, these devices are now capable of initiating sophisticated cyberattacks without the need to pass through a powerful centralized entity. We therefore anticipate extremely challenging future security attacks initiated by infected mobile devices that coordinate together, forming what we call a mobile distributed botnet (MobiBot). Leveraging the advantages of device-to-device (D2D) communication in masking malicious code propagation, as well as the absence of basic mobile operating system security mechanisms, make MobiBots a serious security threat to any system or network.

In this article, we quantitatively assess the potential risk for botnet and MobiBot security while emphasizing the common characteristics and differences between them. We investigate the potential for and impact of large-scale infection and coordination of mobile devices. We highlight how mobile devices can leverage short-range wireless technologies in attacks against other mobile devices that come within proximity. Similar to classical botnets, a MobiBot life cycle consists of three main phases; infection, communication, and attack. We quantitatively measure the infection and propagation rates within MobiBots. We show that MobiBot infection can infect a 96-node network in only a few minutes, which can scale up to 10,000-node networks.

CLASSIC BOTNETS

A botnet is a network of infected machines that we call bots, controlled by a malicious machine called a botmaster. The botmaster aims to control a network of bots in order to initiate malicious attacks against one or many victim machines. Bots and the botmaster communicate via exchanging C&C messages. In most cases, as shown in Fig. 1a, the botmaster creates a set of

*Abderrahmen Mtibaa and
Hussein Alnuweiri are with
Texas A&M University.*

*Khaled A. Harras is with
Carnegie Mellon University.*

A botnet is a network of infected machines that we call bots, controlled by a malicious machine called a botmaster. The botmaster aims to control a network of bots in order to initiate malicious attacks against one or many victim machines.

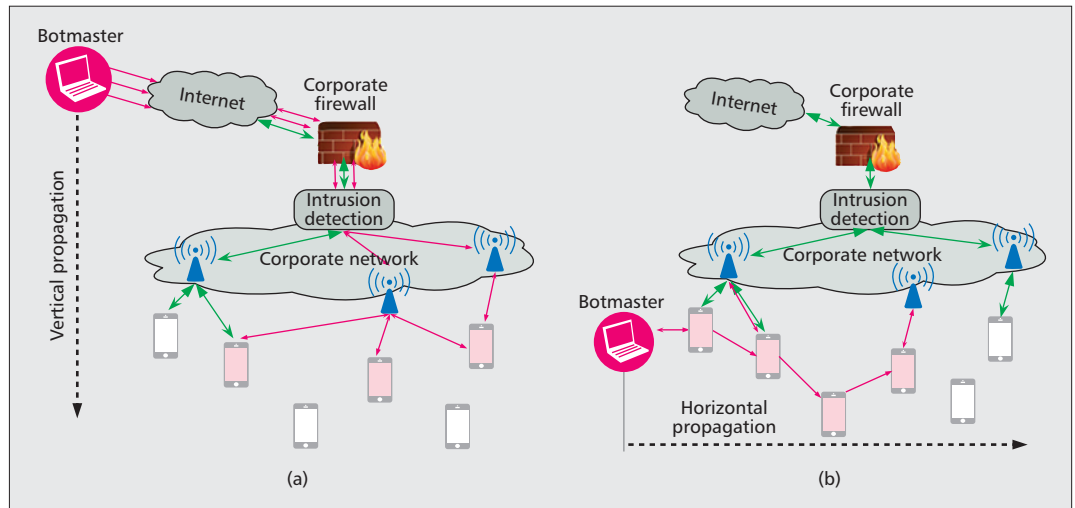


Figure 1. Comparison between a) classical botnet and b) MobiBot architectures (green arrows: non-malicious traffic, red arrow: malicious traffic): a) classical botnet “vertical” architecture; b) MobiBot’s “horizontal” architecture.

C&C servers to play an intermediate role between the botmaster and its bots. Servers can be rented by the botmaster to malicious clients to perform malicious acts such as sending spam and distributed denial of service (DDoS) attacks.

BOTNET ARCHITECTURE

A botnet architecture can be centralized, distributed, or hybrid. In a centralized architecture, bots contact the C&C server in order to receive information from the botmaster. In general, little time is spent in the transmission of a message from the botmaster to all its bots, and this represents one of the major advantages of this scheme. However, the C&C server constitutes a single point of failure. Thus, if the server shuts down, the entire network is dismantled. In a distributed architecture, all bots act simultaneously as servers and clients (i.e., peers). This philosophy avoids the existence of a single point of failure. However, the time required for a message to reach all nodes is much larger. Finally, hybrid botnets combine the advantages of the two previous architectures while introducing multiple management difficulties. In this case, one or more distributed networks, each with one or more centralized servers, exist. The disconnection of one of these servers implies, in the worst case, the fall of one of the distributed networks, allowing the rest to continue normal operation.

Regardless of the botnet architecture, the life cycle of any classical botnet consists of three main phases: infection, communication, and attack phases.

BOTNET INFECTION PHASE

Attacker View: In this phase, the botmaster designs a botnet in order to recruit as many machines as possible. Recruitment starts by infecting victim machines (called bots) using different techniques such as email messages, spam containing attachments, links that trigger installation of malicious programs, malware applications, or social networks.

Defender View: Apart from desperately look-

ing for solutions to demotivate potential botnet developers from undertaking the task of designing a botnet, providing techniques that stop or reduce the infection rate is one of the first steps toward limiting a botnet’s impact. Solutions consist of securing the end hosts by installing and updating anti-virus software, and designing networks resilient to botnet infection by implementing intrusion detection techniques and firewalls that limit the infection of most end host machines. Network-based solutions range from intercepting malware prior to infection, to monitoring and identifying the infection before securing it.

BOTNET COMMUNICATION PHASE

Attacker View: Once the victim machine is infected, it becomes a bot registered to its botnet network. Registration helps the C&C servers and/or the botmaster gathering information about the new bot and its capabilities. The C&C communication framework channel is used to communicate malicious attack orders such as downloading and executing software updates.

Defender View: This phase represents the first attempt to combat a botnet defined as a collection of individual infected machines. Communication defense consists of detecting and identifying malicious communication between the infected bots and their C&C servers, and isolating the botnet by reporting and shutting down the detected C&C server. Researchers have proposed many detection techniques based on honeynets [7] and/or intrusion detection [8].

Honeynets consist of making vulnerable machines available for eventual infection. The infected machine spies on C&C communication messages and collects information about the botmaster or its servers, and the botnet architecture/characteristics. Today, botmasters and botnet developers are aware about such defense techniques, and have started designing architectures and ways to avoid honeynets [7].

Intrusion detection systems monitor traffic coming in or out of an enterprise network or an

Internet service provider (ISP). Such systems search for network traffic anomalies, such as high network delays, high traffic throughput, and suspicious port usage [8].

BOTNET ATTACK PHASE

Attacker View: Upon receiving C&C messages, the botnet enters the attack phase where a group of recruited bots are asked to initiate a specific malicious activity such as DDoS, spam, click fraud, and phishing attacks [8]. The ultimate goal of a botnet is to successfully and stealthily execute the planned attack. Upon success, the botmaster tries to reuse the same bots for future attacks while infecting new bots, continuously increasing its size and its revenue [9].

Defender View: While the earlier defense stage the better efficiency in thwarting a botnet's success, defense mechanisms at the attack phase represent the last chance to ultimately stop malicious attacks [8].

Both enterprise and ISPs have to implement anti-hiding mechanisms that reduce the impact of the attack by increasing its detectability. Anti-hiding mechanisms increase the probability of botnet detection and isolation.

MOBIBOT SECURITY

In typical computer communication networks, attacks are assumed to originate outside the network. Hence, traffic is analyzed at the edge of the networks to detect anomalies, and firewalls are placed at the periphery of the network to protect against intrusions and attacks. While this has served the community well, we anticipate that infections of laptops and mobile handsets enable penetration of these protection mechanisms, and hence attacks can come from a device inside the network. With the proliferation of botnets and bring your own device (BYOD), we anticipate that it will be difficult to avoid insider attacks. This assumption changes the types of protection we have to design compared to protecting infrastructure networks. It also changes the landscape for potential solutions for discovered problems.

MOBIBOT ARCHITECTURE

A MobiBot is a mobile botnet consisting exclusively of a set of infected mobile bots capable of initiating sophisticated cyberattacks while coordinating together via D2D short-range wireless communication such as Bluetooth [10]. A MobiBot is originated and controlled by one or many botmasters. The botmasters' goal is mainly infecting as many mobile devices while ensuring that all malicious communication within a MobiBot is masked. They will then coordinate attacks originated from the infected bots toward a specific target or simply steal valuable information from the infected mobile devices.

Within a MobiBot, infected bots communicate mainly via mobile-to-mobile wireless communication. We call this communication paradigm *horizontal* (Fig. 1b) as opposed to the classical *vertical* communication between bots and their botmaster in classical botnet scenarios (Fig. 1a).

We consider a malicious MobiBot scenario depicted in Fig. 1b. We investigate the basic sce-

nario where one or possibly several malicious nodes called botmasters need to attack a given target such as a corporate network, steal a top secret file, and so on. The botmasters start by infecting one or more initial bots that will help disseminate malicious code toward nodes in order to form a *MobiBot*. Botmasters will then verify that they have infected enough bots before initiating their attacks. The verification process is challenging since the size of the network, N , is unknown and hard to predict by the botmasters. We investigate the time required to infect a given network. This time should be long enough to infect as many bots as possible in order to inflict significant damage. However, very long waiting times increases the probability of identifying and isolating the attack by users or administrators.

MOBIBOT INFECTION PHASE

The infection phase in the MobiBot life cycle is similar to that of the classical botnet. In addition to the botnet infection techniques, a MobiBot leverages D2D communication to recruit new bots. This involves utilizing short-range wireless technologies such as Bluetooth, WiFi Direct, and Zigbee, or long -range wireless technologies such as Long Term Evolution (LTE) Direct and 5G.

Attacker View: With respect to the botnet developer, MobiBots offer additional infection opportunities via D2D wireless communication that help recruit more bots, thus increasing the impact of the attack and improving its success probability.

To ensure D2D infection, users' devices must have their short-range wireless interfaces turned on. For instance, mobile users keep their Bluetooth interfaces off when they are not using them (other short-range technologies have similar characteristics). In addition, Bluetooth requires an explicit communication setup between the sender and the receiver. Botmasters are therefore forced to design their infection code to infect/recruit bots while turning on their wireless interfaces and modifying the Bluetooth security level to accept communications from other bots. Moreover, turning on wireless interfaces in a mobile device can be suspicious. Therefore, this operation should be transparent or short in time because, for instance, keeping a Bluetooth interface on for a long period of time will increase the probability of detecting such abnormal behavior of the mobile device, thus suspecting an infection. On the other hand, enabling Bluetooth for short periods of time may not be enough to transfer messages between infected bots. In addition, the Bluetooth interface should be on at the same time. This requires synchronization between devices in a fully distributed and disconnected network, which represents a very challenging problem.

Defender View: As shown in the previous section, infection defense is very important because it represents the first step toward limiting the success of eventual malicious attacks. Similar techniques as those discussed in the previous section apply here as well. However, any background thread trying to exchange messages via D2D wireless communication is suspicious, and should be detected and isolated.

With the proliferation of botnets and bring your own device (BYOD), we anticipate that it will be difficult to avoid insider attacks. This assumption changes the types of protection we have to design compared to protecting infrastructure networks.

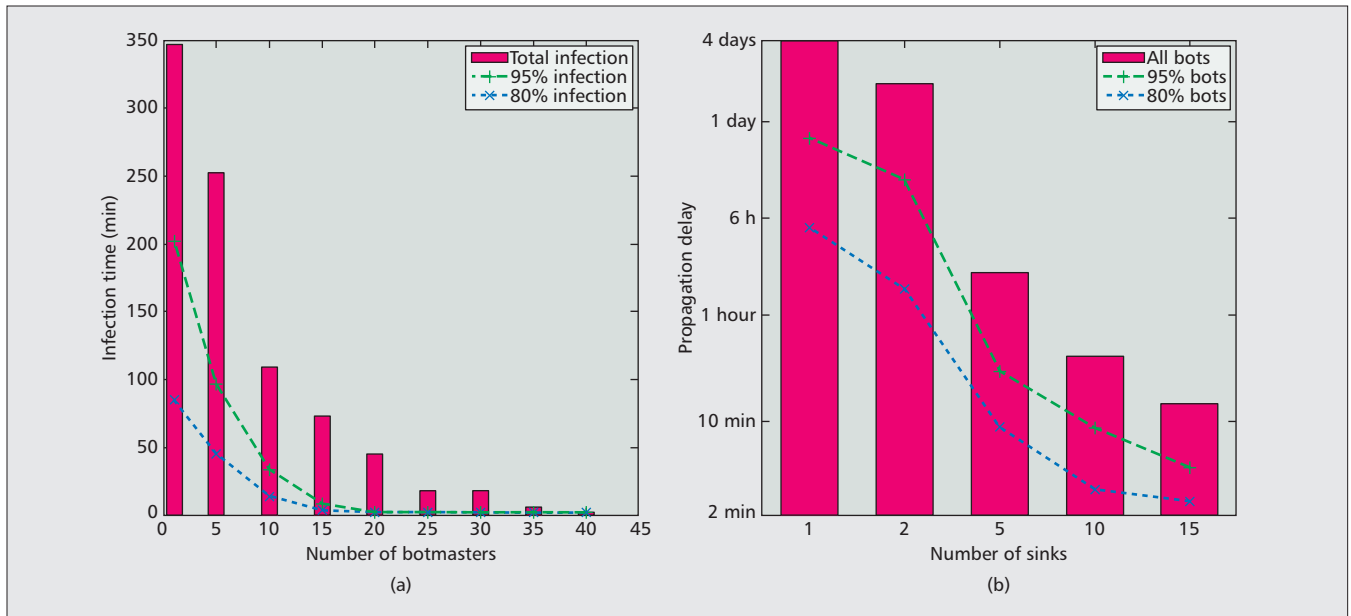


Figure 2. MobiBot infection rate and propagation rate in the Infocom '06 dataset: a) infection time (Infocom06 dataset); b) epidemic attack response time vs. number of sinks (data size = 7 MB).

Therefore, any application that requires enabling D2D wireless communication should be considered suspicious by antivirus programs. Antivirus programs are essential for today's mobile devices. These programs are as important as the ones utilized for desktops and laptops. The goal is to provide antivirus systems that are suitable for mobile devices with limited resource consumption (e.g., battery life).

MOBIBOT COMMUNICATION PHASE

Attacker View: A MobiBot has a unique architecture that relies mainly on D2D wireless communication. We leverage the fact that a mobile communications device, within its intrinsic motion patterns, makes frequent contact with entities that are capable of exchanging data/codes. Contact with some of these entities can be intermittent, limited in duration, and sometimes unpredictable. The challenge facing the botmaster is how to establish a two-way routing mechanism that helps C&C control messages travel between the botmaster and its bots. Botmasters need to send such messages to control the bots, send them updates, or coordinate an attack. On the other hand, bots require some-time to send data (e.g., data stealing attack) back to the botmaster.

A solution to this problem resides in electing a subset of bots that act as sinks to collect data and send messages for coordination.¹ Therefore, the botmaster contacts these sinks opportunistically to collect the stolen data and/or disseminate updates and messages through these sinks. We investigate the required time to disseminate a malicious code and infect 100, 95, and 80 percent of the nodes in the network using D2D opportunistic communication. We consider the mobile network environments given by the IEEE Infocom '06 real-world datasets [11]. This dataset consists of Bluetooth encountering logs between participants at the IEEE Infocom '06 confer-

ence. We assume that command messages used during the infection phase are small enough to be carried out with any contact log in our datasets. We vary the number of sinks used to disseminate C&C messages and measure the infection rate with respect to the number of sinks used.

While infection of the total number of nodes in the network takes more than 5 h, this time is reduced by almost 50 percent if the botmaster wants to infect only 95 percent of the nodes, and reduced by almost 75 percent if the botmaster wants to infect only 80 percent of the nodes (Fig. 2a). These results are verified using different network environments [10].

In order to increase the delivery success rate of both C&C messages and data messages, message replication can be utilized. Sinks create multiple replicas of the same message, which increases its probability to reach more bots with shorter delay. This can also be achieved by sending the same messages from different sinks located in different geographic zones. However, while redundancy increases the delivery ratio, it also increases the network overhead, which leads to reducing the stealthiness of the MobiBot.

Defender View: Mobile-to-mobile communication defense is very challenging. State-of-the-art solutions rely on traditional vertical communication paradigms. They utilize methods and techniques to identify the C&C message exchanges between an infected bot and its botnet servers. However, the MobiBot architecture does not rely on such vertical architecture. MobiBot nodes communicate together in a purely distributed manner, which makes traditional detection techniques fail.

Inspired by the success of the HoneyPot solution, a mobile HoneyPot technique can be utilized. It consists of exposing a few mobile devices located in several key zones to potential infection via mobile-to-mobile communication (e.g.,

¹ We note that, as opposed to sensor network definition of sinks, sinks in this article are bots that push control messages and malicious code in addition to gathering data from all bots

HoneyBot nodes should keep all their wireless interfaces on all the time and accept all incoming messages). However, as opposed to traditional HoneyPot nodes that communicate directly with the botnet C&C servers, HoneyBots only communicate with simple bots. The challenge, therefore, is identifying the elected sinks or the botmasters.

MOBIBOT ATTACK PHASE

Attacker View: A MobiBot, in addition to its stealthiness compared to classical botnet networks, takes advantage of node mobility, which increases the infection rate of the network and allows malicious code to reach very sensitive targets and regions that are normally well secured against network attacks. In fact, mobility allows the botmaster to control machines that are physically situated near sensitive targets. In addition, MobiBots are designed to implement very targeted attacks in a very short period of time. Infection and propagation delays do not exceed a few hours to one day, as shown above.

In our risk assessment, we consider two basic attack scenarios; targeted and epidemic attacks.

The Targeted Attack: It consists of initiating an attack from multiple bots in order to steal a particular data file from a given source such as a confidential file located in a corporate network or/and a server. This attack is successful when the botmaster gets the requested file from the first bot.

The Epidemic Attack: It consists of stealing confidential data from all infected nodes; for instance, collecting all localization data or emails/SMS periodically from all bots. This attack requires the transfer of many data files from multiple bots toward a few botmasters. Data is forwarded using D2D communications through opportunistic contacts.

We measure the propagation rate, defined as the average waiting time taken by the botmaster to acquire the requested files from the infected bots. It is the time difference between sending files from the bots until their reception at the botmaster. We also measure the impact of the number of sinks on the propagation time of an epidemic attack when 1, 2, 5, 10, or 15 sinks were used, as plotted in Fig. 2b. We use a 7 MB file that requires up to four days to propagate when only one sink is used. However, allowing 15 sinks reduces this time to only 10 min. Recruiting multiple sinks is therefore needed to perform epidemic attacks.

Defender View: A MobiBot attack defense consists of tracking and isolating infected bots as well as the sinks used by the botmaster. However, since MobiBot's tracking and isolation is very challenging, we consider MobiBot prevention techniques that hinder mobile infection by leveraging different trust levels to enforce confidence values to a given application. We investigate the impact of existing trust-based prevention techniques on network performance. Our results show that prevention techniques are very costly [10]. We show that while MobiBots are very hard to detect and isolate compared to common botnet networks, even a basic prevention technique makes a network perform at only 60 percent of its capacity [10].

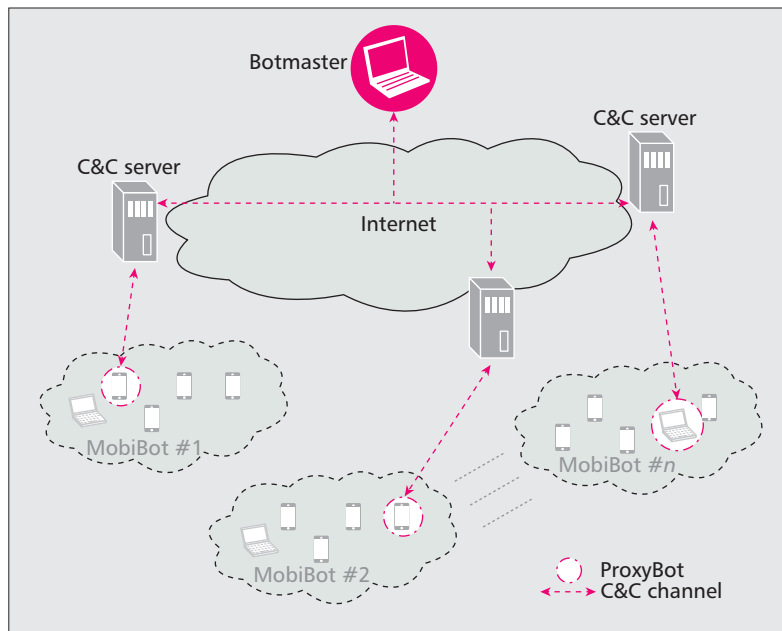


Figure 3. Hybrid MobiBot architecture.

MOBIBOT SCALABILITY

A MobiBot's horizontal communication architecture provides malicious users with attractive stealth tools and mechanisms to attack personal or corporate users. However, the D2D communication paradigm of MobiBots may not be as scalable as classical botnets. In fact, a MobiBot will never be able to recruit and manage a big number of bots as in classical botnets (e.g., hundreds of thousands of bots). We believe that the scalability of MobiBots is an important metric that determines its imminent deployment.

HYBRID ARCHITECTURE

The MobioBot hybrid architecture consists of n MobiBots, where $i = 1, \dots, n$. Connecting multiple MobiBots located in separate geographic areas allows the botmaster to control a larger number of bots simultaneously. We assume that bots are originally assigned to a given MobiBot based on their distinct location (e.g., city or country). Communications across MobiBots happen whenever bots from a given $MobiBot_i$ come in contact with other bots belonging to a different $MobiBot_j$, or through a fixed infrastructure-based server controlled by the botmaster, as shown in Fig. 3. Communication with these C&C servers is also intermittent. In fact, these servers do not maintain any connection to one or multiple bots. However, they opportunistically communicate with bots that are connected to the Internet.

We consider a hybrid MobiBot network consisting of a single botmaster controlling k C&C servers (Fig. 3). Each C&C server controls 1 to m mobile MobiBots via their elected bot acting as a sink. We call this elected bot the *ProxyBot*. Each $MobiBot_i$ can elect 1 to h $ProxyBot_i^j$, where $j = 1, \dots, h$. In addition, a given $ProxyBot_i^j$ can communicate with multiple C&C servers.

Whenever two mobile bots come within communication range, we call such event a contact,

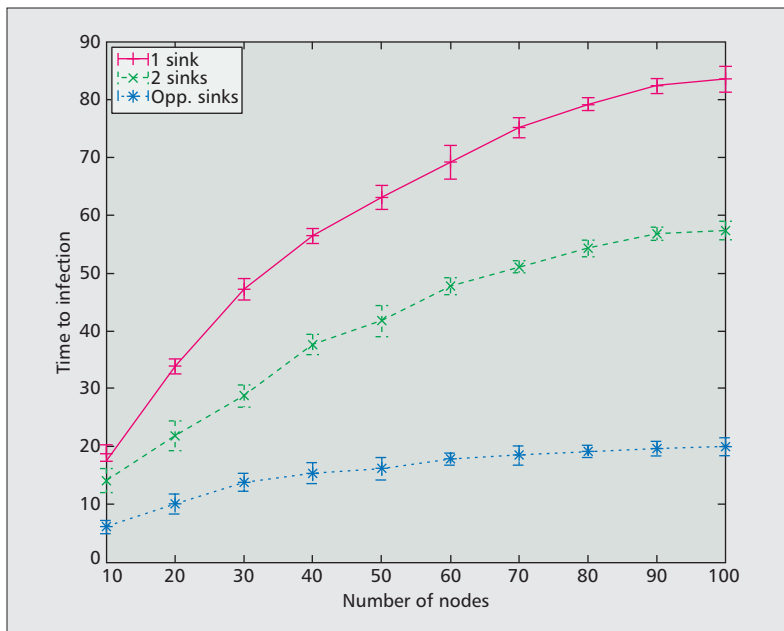


Figure 4. Efficiency vs. stealthiness.

during which mobile bots can exchange control or/and data messages. Such communication is allowed not only when the encountered bots belong to the same $MobiBot_i$, but also when they belong to two different $MobiBots$.

SCALABILITY STUDY

Our architecture introduces multiple challenges, such as how the system selects a sink, and how many sinks the system needs in order to satisfy contradictory metrics, efficiency, and stealthiness. In order to deal with such challenges and more, we adopt a data-driven approach that highlights these issues and the characteristics of large-scale $MobiBots$.

Due to the lack of large-scale experimental data sets, we artificially create a $MobiBot$ network that consists of 100 to 10,000 bots. We call this dataset $MBots14$. $MBots14$ consists of m disjoint communities. Each community is represented by a real-world mobility trace. Traces are generated by modifying the two original traces: INFOCOM '06 [12], and Sigcomm '09 [13]. We incorporate traces for a duration of three days, interpolate virtual infrastructure server nodes, and then generate opportunistic contacts between these servers and random nodes in each dataset. In each of the above traces, we randomly select from 50 to 100 percent of the nodes and/or 50 to 100 percent of the contacts in the original dataset.

We first focus on the characteristics of a single $MobiBot$ network. We select one, two, and multiple sinks per $MobiBot$ cluster. One and two sinks are selected randomly from the set of infected nodes in the dataset. However, we implement an algorithm that selects nodes which are opportunistically connected to the C&C servers. In other words, each node can act as a sink if it is able to establish a direct connection with the server for no less than δt (i.e., in our experiment we fix $\delta t = 10$ min).

Figure 4 compares the performance of the

$MobiBot$ when one, two, and multiple sinks are deployed. We plot, with 98 percent confidence intervals, the average infection time out of 50 runs where sinks are randomly selected. We show that with only two sinks, nodes can be infected in less than 1 h. Allowing two sinks helps improve the infection time by up to 20 percent. Note that we deployed the sinks randomly. We believe that better deployment of sinks that can be placed in opposite geographic locations or with the highest density can boost performance considerably. The goal of this plot is not to show the optimal infection rate but highlight the effectiveness of $MobiBot$ infection with only a few sinks. We also show that opportunistic communication with the server improves the infection rate considerably and infects all nodes in less than 20 min. However, opportunistic communication with the server introduces large overhead traffic, resulting in sending redundant data. Moreover, setting up multiple connections between infected nodes and the server reduces its stealthiness. Our experiments show that the opportunistic method introduces up to 650 percent overhead messages compared to those generated by the 2-sink method, and an average of 138 connections (i.e., connections are intermittent in time) with the server while the sink methods maintains only 2 connections.

We then consider our $MBots14$ dataset to compare the infection rate when the number of infected bots scale up to 10,000 bots. We show in Fig. 5a that the infection rate does not increase exponentially as we scale the number of bots. We also compare the infection rate using the hybrid and classical horizontal architectures. We use the SanFrancisco dataset [14] where taxis move in three different areas in San Francisco. We show that when taxis are concentrated in the bay area (100 to 150 cabs) the infection rate is faster than that of $MBots14$. This is mainly because cabs are faster than pedestrians, and this helps accelerate the infection of all nodes. However, when the number of nodes scales, the horizontal architecture fails to infect all nodes within 5 to 6 h, while the hybrid architecture helps stabilize the infection rate at no more than 1 h in large-scale networks.

We also investigate the overhead messages required to infect a large-scale network in Fig. 5b. Overhead messages include exchanged messages between any two encountered nodes. We show that both architectures require the same number of messages to infect up to 40 percent of the nodes. However, the horizontal architecture fails to infect the rest of the nodes with low overhead messages. Reaching the unpopular areas of San Francisco takes time and costs overhead that can be translated to energy consumption as well. This same phenomenon is also shown with the hybrid architecture, but the number of messages exchanged does not exceed 220 percent of the total number of infected nodes.

CONCLUSION

This article highlights the potential capabilities of mobile devices in initiating sophisticated cyberattacks while forming what we call $MobiBots$. $MobiBots$ leverage device-to-device com-

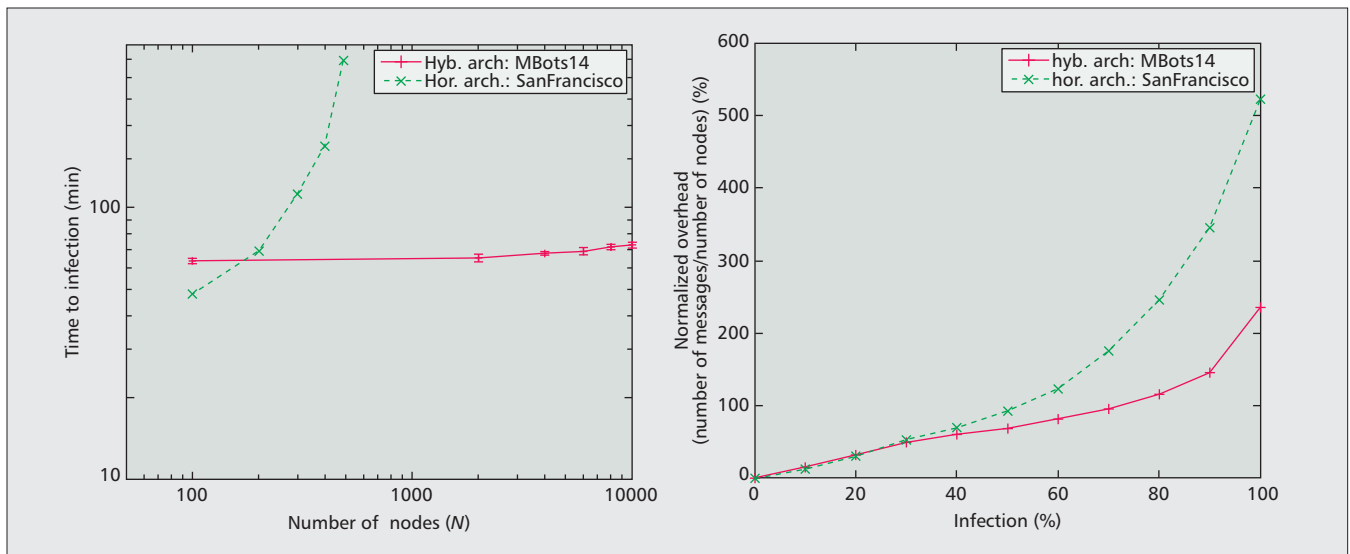


Figure 5. MobiBot scalability analysis with/without hybrid architecture.

munication to mask malicious infection and propagation in order to become an extremely challenging security threat. In this article, we present a risk assessment of MobiBots showing that they are able to efficiently infect and stealthily propagate data.

ACKNOWLEDGMENT

This publication was made possible by NPRP grant # 5-648-2-264 from the Qatar National Research Fund, a member of Qatar Foundation. The statements made herein are solely the responsibility of the authors.

REFERENCES

- [1] T. Micro, Mobile Security, <http://www.trendmicro.com/cloud-content/us/pdfs/security-intelligence/reports/rpt-2q-2013-trendlabs-security-roundup.pdf>, 2013.
- [2] D. Maslennikov, IRC Bot for Android, <http://securelist.com/blog/virus-watch/32310/irc-bot-for-android-14/>, 2012.
- [3] W. F. 2.0, Android Botnet Targets Middle East Banks, <http://krebsonsecurity.com/2014/04/android-botnet-targets-middle-east-banks/>, April 2014.
- [4] A. Mtibaa et al., "Towards Resource Sharing in Mobile Device Clouds: Power Balancing Across Mobile Devices," *Proc. 2nd ACM SIGCOMM Wksp. Mobile Cloud Computing*, 2013, pp. 51–56.
- [5] A. Mtibaa, A. Fahim, and K. A. Harras, "Towards Computational Offloading in Mobile Device Clouds," *Proc. IEEE 5th Int'l. Conf. Cloud Computing Technology and Science*, Dec. 2013.
- [6] T. McKay, "Here's the Ingenious Way Protesters in Hong Kong Are Organizing Themselves," <http://mic.com/articles/100148/here-s-the-ingenious-way-protesters-in-hong-kong-are-organizing-themselves>, Oct. 2014.
- [7] P. Wang et al., "HoneyPot Detection in Advanced Botnet Attacks," *Int'l. J. Info. Comp. Security*, vol. 4, no. 1, Feb. 2010, pp. 30–51.
- [8] A. Karim et al., "Botnet Detection Techniques: Review, Future Trends and Issues," *J. Zhejiang Univ. SCIENCE C*, 2014.
- [9] S. S. Silva et al., "Botnets: A Survey," *Computer Networks*, vol. 57, no. 2, 2013, pp. 378–403.
- [10] A. Mtibaa, K. Harras, and H. Alnuweiri, "Malicious Attacks in Mobile Device Clouds: A Data Driven Risk Assessment," *Proc. IEEE 23rd Int'l. Conf. Comp. Commun. and Networks*, 2014.

- [11] A. Chaintreau et al., "The Diameter of Opportunistic Mobile Networks," *Proc. ACM CoNext*, 2007.
- [12] A. Chaintreau et al., "Impact of Human Mobility on Opportunistic Forwarding Algorithms," *IEEE Trans. Mobile Computing*, vol. 6, no. 6, 2007, pp. 606–20.
- [13] A.-K. Pietilainen, CRAWDAD dataset thlab/sigcomm2009 (v. 2012-07-15), <http://crawdad.org/thlab/sigcomm2009/>, July 2012.
- [14] M. Piorkowski, N. Sarafijanovic-Djukic, and M. Grossglauser, CRAWDAD data set epfl/mobility (v. 2009-02-24), <http://crawdad.org/epfl/mobility/>, Feb. 2009.

BIOGRAPHIES

ABDERRAHMEN MTIBAA (amtibaa@tamu.edu), Ph.D., is currently an assistant research scientist at Texas A&M University in Qatar. Prior to that, he was a postdoctoral research associate in the School of Computer Science at Carnegie Mellon University in Qatar. Graduated on June 2010 from the University of Paris VI and Technicolor Paris Research Lab. His current areas of interest include mobile cloud computing, mobile security, mobile opportunistic networks/DTN, wireless and Ad-hoc networks, mobility models, protocol design, routing/forwarding, network communities and social networking. He has two US patents and more than 40 publications in numerous international prestigious journals, conferences, and workshops.

KHALED A. HARRAS [M], Ph.D., is currently an Associate Teaching Professor at Carnegie Mellon University. He is the founder and director of the Networking Systems Lab, and the Computer Science Program Director at CMU's campus in Qatar. His main research interests include delay and disruption tolerant networks, specifically protocol and architectural challenges in extreme networking environments, visual and wireless sensor networks, social pervasive systems, computational offloading, and multi-interface networking and communication. He has published more than 70 papers, workshops, and journals in top international venues. He is a member of the ACM.

HUSSEIN ALNUWEIRI, Ph.D., is a Professor in the Department of Electrical & Computer Engineering at Texas A&M University in Qatar. He has a long record of industrial collaborations with several major companies worldwide. He is also an inventor, and holds three US patents and one international patent. He has authored or co-authored over 150 refereed journal and conference papers in various areas of computer and communications research such as mobile Internet technologies, wireless protocols, routing and information dissemination algorithms for opportunistic networking, and quality-of-service provisioning and resource allocation in wireless networks.

Privacy-Preserving Participatory Sensing

Qinghua Li and Guohong Cao

ABSTRACT

The proliferation of mobile devices such as smartphones has enabled participatory sensing systems that collect data from users through their mobile devices and infer useful information from the data. However, users have concerns regarding possible privacy leakage from their data and lack incentives to contribute their data. To effectively motivate users to participate, both privacy and incentive issues need to be addressed. In this article, we address how to simultaneously protect privacy and provide incentives for participatory sensing. We review previous approaches, discuss their limitations, and propose new approaches for two types of participatory sensing systems.

INTRODUCTION

Smartphones have become an essential part of our daily life. Besides providing voice and data communication, smartphones are also acquiring richer functionality through various sensors. For example, the iPhone 5 includes eight different sensors: accelerometer, GPS, ambient light, dual microphones, proximity sensor, dual cameras, compass, and gyroscope. These sensors are very useful for gathering data about people and their environments. For example, GPS enables sensing of locations, microphones can record sounds of the surroundings, and an accelerometer enables sensing of users' movement and activities. Smartphone users are also considered as special types of sensors that generate human input data to be used in many surveys. Recently, these sensors have been used for *participatory sensing* in which the sensing data on multiple phones are collected by remote data collectors to support many interesting applications, including tracking the spread of disease across a city, building a noise map, monitoring traffic conditions, and so on.

In spite of many useful applications, there are two obstacles that hinder the large-scale deployment of participatory sensing applications [1, 2]. First, participatory sensing poses serious threats to user privacy. The data from mobile devices may be exploited to obtain private information about users, including their locations, health condition, lifestyle, religious activities, and so on. For example, an applica-

tion that tries to build a noise map for a city may request that each user continuously upload his or her current location and the noise level at this location. However, from this data, the data collector can infer where the user has been (e.g., if he or she has gone to a hospital or church) and possibly infer the user's activities. Second, users lack incentives to join in participatory sensing. To participate, a user has to trigger various sensors to measure data (e.g., to obtain GPS locations), which may consume much of his/her smartphone's power. Also, the user needs to upload data to the server, which may consume some data of a wireless quota (e.g., when the data is photo/video). Moreover, the user may have to be in a specific location to sense the required data. To motivate participation, both obstacles should be addressed.

Based on what kind of data is needed by the collector, participatory sensing systems can be divided into two categories. In one category, raw data is needed by the collector to do data mining; hence, mobile devices directly submit their sensed data (e.g., GPS coordinates, accelerometer reading, and user input). In the other category, the collector is interested in the aggregation statistics of a group of mobile devices' data instead of each individual's raw data. In many monitoring applications, such aggregation is periodically done to continuously identify interesting phenomena and track important patterns [3, 4]. For example, the average amount of daily exercise (which can be measured by motion sensors) people get can help infer public health conditions. The average level of air pollution and pollen concentration can help people plan their outdoor activities. In such systems, it is not necessary for a mobile device to submit its raw data. Instead, it can perturb the data value in some way as long as the collector can recover the correct aggregate statistics. For convenience, we call these two categories *raw-data-based sensing* and *aggregate-data-based sensing*, respectively. Due to the difference in the nature of data collected, they can be dealt with differently.

In this article, we address how to simultaneously protect privacy and provide incentives for raw-data-based and aggregate-data-based sensing systems. We introduce the challenges in providing privacy and incentive simultaneously,

Qinghua Li is with the University of Arkansas.

Guohong Cao is with Pennsylvania State University.

review previous approaches, discuss their limitations, and propose new approaches to address those challenges.

PARTICIPATORY SENSING SYSTEMS AND CHALLENGES

PARTICIPATORY SENSING

In a participatory sensing system, there are many mobile nodes and a remote data collector. Mobile nodes are mobile devices carried by users or mounted in vehicles, and they produce sensing data such as locations, pictures, and sound. The collector is interested in collecting different data from mobile nodes (we use node and user interchangeably). Nodes communicate with the collector through cellular networks or WiFi. To facilitate data collection, the collector publishes sensing tasks to nodes, and the latter submit data reports according to the specification of sensing tasks. To incentivize participation, the collector remits credits to nodes for their contributed data. The credits can be converted to monetary rewards in the real world, or used to purchase service from the collector. A sensing task typically specifies the type of sensor readings needed (e.g., GPS coordinates), the conditions of sensing (e.g., in Manhattan during rush hour), the number of data reports needed from each node, the number of credits paid to each node, creation time, and expiration time.

REQUIREMENTS FOR PRIVACY

The collector is untrusted. It is eager to infer the private information of nodes from the system. To protect privacy, the following properties are required:

- Given a readable data report, the collector does not know which node has submitted it. This is required since a data report may contain private information such as the node's location.
- Given a sensing task, the collector does not know if a given node has accepted the task. Sometimes, accepting a task itself causes privacy leakage. For instance, a task asks for the current noise level at Central Park. If the collector knows that Alice has accepted this task, even if it does not know which report is submitted by Alice, it can infer that most likely Alice is around Central Park now.
- The collector cannot link multiple readable reports submitted by the same node. If the collector can do so, and each of the linked reports includes when and where this report is generated, the collector can construct the node's movement path, which in turn helps the collector infer who the node is.
- The collector cannot link multiple sensing tasks accepted by the same node.

REQUIREMENTS FOR INCENTIVE

Nodes are greedy and want to earn as many credits as they can. For this purpose, they may abuse the system in many ways:

- A node may want to earn credits from a task without submitting data for it.

- A node may submit more reports for a task than allowed by the collector to earn more credits.
- Suppose there is a set of credentials for each particular task. Then a node may try to use the credentials for one task to earn credits from another task, because the latter task is paid at a higher rate.
- A node may spend a single credit token twice, causing the double-spending problem.
- A malicious node may compromise other nodes, steal their credentials, and use them to earn credits.

Such misbehavior must be mitigated.

RESOURCE CONSTRAINTS OF MOBILE DEVICES

Besides the requirements for privacy and incentive, another challenge is that mobile devices usually have limited power supply, computational resources, and bandwidth. Hence, a solution must be efficient in energy, computation, and communication. Another constraint is that peer-to-peer communications among mobile devices are expensive. Therefore, a solution should avoid or minimize the use of communications between nodes.

APPROACHES FOR RAW-DATA-BASED SENSING

PREVIOUS APPROACHES AND THEIR LIMITATIONS

Privacy in participatory sensing has been addressed by a lot of work. The simplest solution is to suppress or anonymize a node's identity during tasking and reporting. However, the protection provided by this solution is weak. The problem is that if multiple reports can be linked as being from the same node (e.g., through analyzing the submission times of those reports or through the same pseudonym included in those reports), the node's identity might be inferred [5].

Even if the node's identifier is suppressed and timing analysis is mitigated, if it uses the same IP or medium access control (MAC) address when submitting multiple reports, those reports can also be linked easily. To address these problems, Mix networks (e.g., Mixmaster and Tor) can be used when a node submits data reports. Mix networks enable hard-to-trace communications by sending a message through a chain of Mix nodes (proxy servers on the Internet). Each Mix node receives messages from multiple previous hops, randomly shuffles them, and forwards them to the next hop. In this way, the collector does not know the real source of a request or report. Dynamically changing the IP and MAC addresses also helps protect privacy.

Combining these techniques, AnonySense [6], a general-purpose privacy-preserving framework of participatory sensing, has been proposed. It delivers sensing tasks to anonymous nodes, and collects unlinkable data reports from anonymous nodes. DeCristofaro *et al.* [7] aim at cryptographic treatment of privacy and propose a different framework. In this framework, external entities can query specific users' data, and the

To facilitate data collection, the collector publishes sensing tasks to nodes, and the latter submit data reports according to the specification of sensing tasks. To incentivize participation, the collector remits credits to nodes for their contributed data.

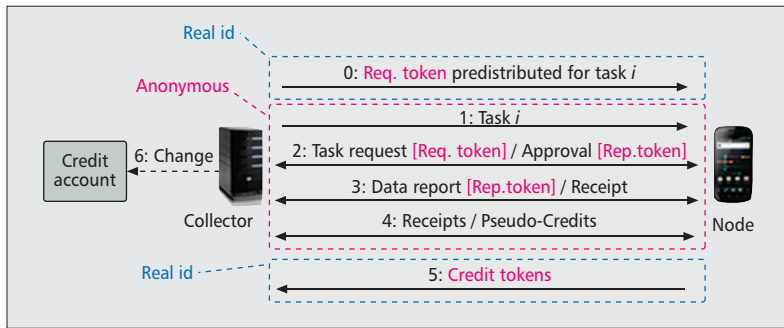


Figure 1. Providing privacy-preserving incentives [2].

framework can hide which mobile node matches a query.

Note that even if tasking and reporting are anonymous, the reported data itself may still cause privacy leakage. For instance, if a report includes a user's home as the location, the collector can link the report to the user. This problem has been addressed by many other techniques that are orthogonal to the approaches discussed in this article.

Several incentive schemes have been designed for participatory sensing to promote participation by paying credits to users. In [8], gaming and reverse auction theories are used to optimize the payment for a task. For example, each user bids for a sensing task with an expected payment, and the collector preferably assigns the task to those users with low bids to minimize the total payment. An optimal solution is found in which each user bids according to a certain strategy.

However, previous approaches address privacy and incentive *separately*, and thus cannot effectively promote user participation. It is challenging to address these two issues simultaneously. One option is to use blind signature to implement privacy-preserving credits, since it has been widely used for anonymous electronic payment. However, blind signature cannot *directly* solve the problem, since a malicious user that has compromised other users can steal and spend their credits without being detected. Other anonymous credential systems cannot be directly used for similar reasons. Privacy-preserving mechanism design and auctions can protect participants' types and valuations of a good, but they cannot be applied in mobile sensing to protect users' interest in sensing tasks.

It is natural to consider the combination of a privacy protection scheme and a credit-based incentive scheme to address both issues, but as pointed out in [2], the problem is not as simple as it appears. Specifically, when anonymity is provided as the protection for privacy, a greedy user can use different anonymous identifiers to submit unlimited data reports for the same sensing task (which may not be desirable) and earn unlimited credits. Also, a malicious node can compromise other nodes and use their credentials to earn credits without being detected. To address these problems, new designs are required to simultaneously address privacy and incentive.

We proposed a solution for providing privacy-aware incentives [1, 2]. Besides satisfying the privacy requirements, the solution ensures that, for a task paid at the rate of c credits per node ($1 \leq c \leq C$, where C is a system parameter), one node cannot earn more than c credits from this task.

The approach relies on a set of tokens (including task *request token*, *report token*, *receipt*, and *credit token*) to achieve the goals of incentive and privacy. Each node predetermines the tokens it will use to process each future task, and commits to the collector that it will really use them. Later, when a node processes a sensing task using the predetermined tokens, the collector will verify that those tokens have been appropriately committed. The design of protocol and commitment guarantees that no node can misuse tokens to earn more credits than it should. To protect privacy, tokens and their commitments are constructed and used in a privacy-preserving way.

Basic Protocol: A request token for each future task is pre-distributed to each node before the task is created. At random intervals, each node, say Alice, anonymously retrieves tasks from the collector. If Alice wants to submit data for a task, she anonymously sends a request to the collector. This request contains her request token for this task. If the collector approves this request, it sends back an approval message and issues a number of report tokens to Alice. At this time, the task is assigned to Alice. Alice collects data in the way specified by the task. Then she anonymously submits each data report attached to a report token, and receives a receipt for each report. An independent communication session is used to submit each report. After Alice submits all reports for a task, she anonymously submits the receipts of this task to the collector to redeem credits. The collector issues Alice some pseudo-credits, which are transformed into credit tokens by Alice. For each credit token, Alice waits a random time and then deposits the token with the collector. The collector maintains a credit account for each node in the system, and updates Alice's credit account accordingly. The collector ensures that each credit token can only be spent once, which solves the double-spending problem. Alice uses her real ID to obtain the request token. In this way, the collector can make sure that only one request token for each task is issued to Alice, such that Alice cannot request a task multiple times. When Alice deposits credit tokens, she also uses her real ID.

Credit Token: Since the user uses her real identity to deposit a credit token, and the collector knows the data report from which the corresponding pseudo-credit is earned, it is necessary to break the link between the credit token and the pseudo-credit such that the collector cannot know the data report from which the credit is earned. Blind signature is used to address this challenge. A blind signature scheme [9] enables a user to obtain a signature from a signer on a message without exposing the message content to the signer. In this approach, a credit token

consists of a random token identifier m chosen by the node and its blind signature σ signed by the collector. To obtain a credit token, the user chooses m , blinds it with a random blinding factor, and sends the blinded message m' to the collector. The collector signs on m' using a standard algorithm (e.g., Rivest, Shamir, and Adleman, RSA), and passes the signature σ' (i.e., pseudo-credit) back to the user. The user removes the blinding factor from σ' and obtains a valid signature σ on m . This process ensures that the collector cannot link $\langle m, \sigma \rangle$ to m' or σ' .

Request Token: Alice gets the request token when she communicates with the collector using her real ID, but when Alice uses a request token to request a task, it must be ensured that the collector does not know Alice is the requester. To achieve this goal, partially blind signature (PBS) [10] is used to construct request tokens. PBS is similar to blind signature, except that it allows the signer to put some common (not identifiable) information into the signature. A request token consists of a random token identifier τ , task index i (common information), and the collector's PBS over them. To get the token, Alice chooses τ , and then Alice and the collector run a PBS protocol during which Alice obtains the PBS. The property of PBS ensures that the collector cannot link the PBS or request token to Alice. Besides, the signature binds the token to a specific task, which means Alice cannot use it to process other tasks.

Report Token and Receipt: Multiple reports submitted by the same node should not be linkable. This means that the report tokens obtained in the same session should not be linkable to each other. Since a node's receipts for the same task are submitted together, we must also break the link between a report and its receipt. To achieve these goals, the report token and receipt are also constructed using PBS schemes in a similar way as a request token.

Mitigating Token Misuse: Tasks are indexed as 1, 2, 3, ... in the order of their creation time, and grouped into *task windows* of size W . The first (second) W tasks belong to the first (second) window, and so on. Tokens are managed based on task windows. Let us consider just one window without loss of generality.

Before any task in this window is created, for each task i in the window, each node predetermines the C credit token identifiers m_1, m_2, \dots, m_C it will use for this task. The node commits to the collector that it will use these credit tokens for this task. To do so, it builds an extended Merkle tree [2] over m_1, m_2, \dots, m_C , and then obtains a PBS from the collector over the root τ of the tree and task index i . Here, $\langle \tau, i, PBS \rangle$ is a commitment, and it is also used as the node's request token for task i . Given τ , the identifiers of the receipts the node will use in this task are also determined. The node also binds m_1, \dots, m_C to its own real ID using another extended Merkle tree rooted at α .

When a node submits the receipts of task i to redeem c credits, the collector verifies the relation between identifiers of the receipts and τ .

	Type I $n = 1$ $c = 1$	Type II $n = 1$ $c = 256$	Type III $n = 256$ $c = 256$	Type IV $n = 256$ $c = 1$
Node (Nexus S)	90 ms	95 ms	116 ms	112 ms
Collector (laptop)*	10 ms	14 ms	37 ms	33 ms

*The time needed to process a task for each node.

Table 1. The average running time of processing a task on a Nexus S phone and a laptop [2].

Then, through a blind signature scheme, the node obtains c blind signatures for m_1, \dots, m_c from the collector, and they form the c credit tokens. To prevent the node from abusing the unused identifiers m_{c+1}, \dots, m_C , the node needs to reveal these unused identifiers to the collector, as well as the proof that they are included in the extended Merkle tree rooted at τ . Each node that does not submit data for a task also reveals its credit token identifiers for this task. When a node deposits a credit token $\langle m, SIG(m) \rangle$, it also sends the appropriate elements of the extended Merkle tree rooted at α to prove that m has been bound to its real ID in the token distribution phase.

This process ensures that, given a request token (τ) or receipt (β), the credit tokens (m) which can be obtained and who can deposit these tokens are determined. This means that even if an attacker can compromise other nodes, it cannot use their tokens/receipts to obtain credit tokens that can be deposited to its own account.

This scheme has very low computation cost, as shown in Table 1.

APPROACHES TO AGGREGATE-DATA-BASED SENSING

For aggregate-data-based sensing, when it is important to protect what sensing tasks a node has taken, we can first use the approach for raw-data-based sensing to collect the participants' data and then compute the aggregate. In the following, we consider the case where it is not sensitive if a node has taken a certain task, but it is important to protect the content of each node's data (e.g., the daily amount of exercise). In this case, if we can design an *aggregator oblivious* aggregation scheme (such a scheme ensures that the aggregator learns the desired aggregate statistics but nothing else, e.g., any individual node's data value), nodes can use their real ID to communicate with the collector (i.e., aggregator) without privacy leakage. Then credits can be directly paid to nodes, and the requirements on incentive can easily be satisfied. In the following, we focus on how to do aggregator oblivious data aggregation.

PREVIOUS APPROACHES AND THEIR LIMITATIONS

Homomorphic encryption is widely used in privacy-preserving data aggregation. Homomorphic encryption allows certain computations over ciphertexts to be done in such a way that when

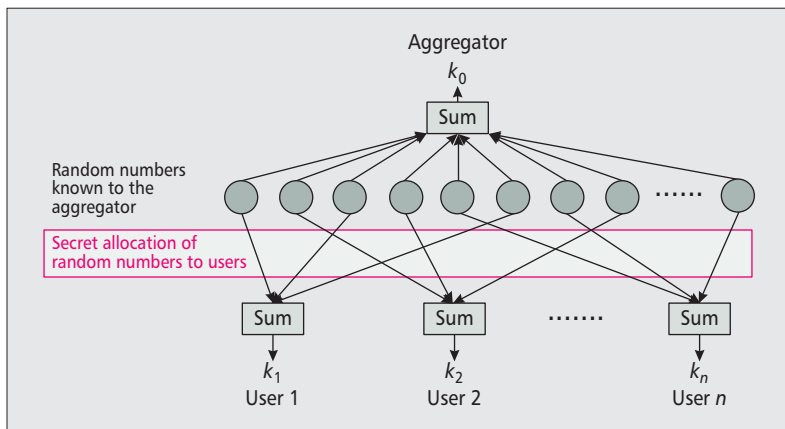


Figure 2. The intuition behind the basic encryption method for sum aggregation [3]. The aggregator computes the sum of a set of random numbers as the decryption key. These numbers are secretly allocated to the users, and each user computes the sum of its allocated numbers as the encryption key. The aggregator does not know which random numbers are allocated to each user, and thus does not know any user's key.

the encrypted result is decrypted, it matches the result of computations over plaintext. For example, suppose c_1 and c_2 are the ciphertexts of m_1 and m_2 . Then one may decrypt the sum of ciphertext $c_1 + c_2$ to get the sum of plaintext $m_1 + m_2$. To use such encryption algorithms in aggregation, user A can encrypt his data and pass the result to user B . User B encrypts her own data, adds the ciphertext to A 's ciphertext, and passes the result to user C . C does similar operations as B . Finally, the ciphertext is passed to the aggregator, which simply decrypts it to get the sum of all users' data. However, in most previous homomorphic encryption-based approaches, all users encrypt data with the same key, which is also known to the aggregator. If the communications between users can be intercepted, the aggregator will be able to decrypt each user's data value.

A similar but more complex approach is secure multi-party computation (SMC). It allows multiple parties to jointly compute a function (e.g., sum) over their inputs but keep these inputs private. However, most SMC algorithms are interactive and require parties to communicate with each other. This is impractical in participatory sensing.

Several recent constructions use the Paillier cryptosystem. It is a homomorphic encryption scheme, but users use different keys to encrypt their data, and the aggregator cannot decrypt any individual's data. The sum aggregation scheme in [11] divides the decryption key into portions and distributes them to nodes. The aggregator collects all nodes' ciphertexts, multiplies them together, and sends the aggregated ciphertext back to all nodes. Each node decrypts a share of the sum aggregate and sends it to the aggregator, which then obtains the final sum. However, this scheme requires multiple rounds of interaction between nodes and the aggregator. Other schemes based on the Paillier cryptosystem rely on inter-node communications, which is not practical in participatory sensing.

To remove multi-round communications and

inter-node communications, Shi *et al.* [12] proposed a sum aggregation scheme assuming that the decisional Diffie-Hellman problem is hard. In this approach, each node i ($i = 1, \dots, n$) gets an encryption key k_i , and the collector gets a decryption key k_0 , which satisfies that their sum is zero. Each node i encrypts its data x_i by computing $g^{x_i}H(t)^{k_i}$ and sends the ciphertext to the aggregator. The aggregator computes the product of all ciphertexts and multiplies it by $H(t)^{k_0}$, deriving $g_{\sum x_i}$. Then it solves the discrete log by trying every possible value of the sum until finding a matching one. This approach has a simple communication model, but the decryption needs to traverse the possible plaintext space of sum, which is inefficient for large plaintext space.

These approaches heavily rely on computationally expensive public key cryptography and hence are too expensive for resource-constrained mobile devices. Moreover, none of the previous schemes on stream data aggregation considers the Max/Min aggregate (i.e., the maximum/minimum value), which is also important.

EFFICIENT DATA AGGREGATION SCHEMES

To address the limitations, we proposed an efficient stream aggregation scheme for sum based on lightweight symmetric key cryptography, and extended it to support Max/Min [3, 4]. It is aggregator oblivious.

Based on a semantically secure additive homomorphic encryption algorithm [13], our basic encryption method for sum aggregation works as follows (see the basic idea in Fig. 2). Suppose each user's data value is an integer in range $[0, \Delta]$. Let n denote the number of users. In the initial setup phase, a key dealer generates a pool of nc different secrets, and divides them into n random disjoint subsets, with c secrets in each subset. It assigns all secrets to the aggregator, and assigns one subset to each user. Let \mathcal{S} denote the set of all secrets, and \mathcal{S}_i denote the subset assigned to user i . If no trusted key dealer is available, the setup phase can also be done through a standard secure multi-party protocol. In aggregation period t , user i computes $f(s, t)$ for each secret s in \mathcal{S}_i . Here f is a pseudorandom function that can be implemented using HMAC. Then user i sets his encryption key k_i^t as the sum of $f(s, t)$ modulo a large integer M . He encrypts his data x_i by computing the sum of x_i and k_i^t modulo M , and sends the ciphertext to the aggregator. The aggregator computes $f(s, t)$ for each secret in \mathcal{S} , and generates its decryption key k_0^t as the sum of $f(s, t)$ modulo M . It decrypts the sum of users' data by adding their ciphertexts and subtracting k_0^t from the result. Since \mathcal{S} is the union of all users' secret set \mathcal{S}_i , the correct sum is derived.

The confidentiality of each user's data relies on the fact that the aggregator does not know the subset of secrets assigned to any specific user. If sufficient secrets are assigned to each user, it will be computationally infeasible for the aggregator to guess the secrets of any specific user in a brute force way. With standard combinatorial techniques, we can derive the number of secrets needed by each user to achieve a required security level (e.g., 80-bit). In most practical settings, this number is very small (less than 10), which means very low computation

cost at mobile devices. The computation cost of the aggregator can also be significantly reduced [3]. Measurements on Nexus S phones and a Windows laptop showed that this scheme runs at least one order of magnitude faster than existing schemes (Table 2). The scheme can easily achieve differential privacy [14] (which is secure against adversaries with arbitrary auxiliary information) with very low noise in the sum [15].

Dealing with Dynamic Joins and Leaves: When a user leaves, her secrets should be removed from the aggregator’s secret pool or redistributed to other users, such that the aggregator can still get the correct sum. Some of the secrets may be redistributed to users compromised by the aggregator, and the aggregator knows that these secrets belonged to the leaving user. This reduces the aggregator’s uncertainty about the secrets used by the leaving user and other users, and decreases the security level of all users. To maintain the required security level, when a user leaves, a new set of secrets should be distributed to all remaining users, which means high communication cost. Similar issues exist when a user joins.

Redundancy in Security: We proposed to address this problem through maintaining redundancy in security [4]. The basic idea is to let each user maintain higher than l -bit security by using more secrets than required by l -bit security. When a users leaves, this user’s secrets together with part of several other users’ secrets are removed from the aggregator’s secret pool, such that the aggregator does not know which of them belong to the leaving user. Those affected users are notified to remove the corresponding secrets from their secret pool. This is carefully done such that each remaining user still maintains sufficient secrets for the required security level. Join is dealt with similarly.

Overlapped Grouping: Intuitively, grouping can be used to efficiently deal with dynamic joins and leaves. In a naïve solution, the key dealer divides users into small disjoint groups, and applies the basic encryption method to each group independently. The aggregator can decrypt the sum of each group, and add them together to obtain the sum of all users’ data. When a user joins or leaves a group, only the users in this group and the aggregator are redistributed secrets, which means low communication cost. However, since the aggregator knows the sum of each group, the solution is not aggregator oblivious. Also, if differential privacy is needed, sufficient noise should be added to the sum of each group, which means a large noise accumulated in the final sum. To address these problems, an overlapped grouping technique is proposed in [15] that guarantees aggregator obliviousness.

CONCLUSIONS AND OPEN CHALLENGES

Although participatory sensing has many useful applications, privacy concerns and lack of incentives prevent its large-scale deployment. In this

	n	10^3	10^4	10^5	10^6
Encryption (Nexus S)	Our scheme	2.4 ms	2.1 ms	1.9 ms	1.4 ms
	EXP	90 ms	90 ms	90 ms	90 ms
Decryption (laptop)	Our scheme	24 μ s	18 μ s	15 μ s	12 μ s
	EXP($\Delta = 10^2$)	1.8 s	5.6 s	18 s	56 s
	EXP($\Delta = 10^3$)	5.6 s	18 s	56 s	177 s
	EXP($\Delta = 10^4$)	18 s	56 s	177 s	560 s
	EXP($\Delta = 10^5$)	56 s	177 s	560 s	1770 s

Table 2. The running time of our Sum protocol and EXP [12] on a Nexus S phone and a laptop [4].

article, we address how to protect privacy and provide incentives to users simultaneously. Focusing on raw-data-based sensing and aggregate-data-based sensing, we review previous approaches and point out their limitations. We also propose novel privacy-aware incentive schemes for raw-data-based sensing and efficient approaches for aggregate-data-based sensing.

Open Challenges: The data collector relies on truthful sensing data collected from nodes to identify important patterns and derive useful statistics. However, dishonest nodes may manipulate their sensing data for benefits. A selfish node may report faked data without doing the actual sensing; a malicious node may manipulate its data values to “pollute” the aggregate data collection. To ensure the usefulness of collected data, it is crucial to mitigate data forgery attacks. How to thwart data forgery under the framework of privacy-aware incentive is an open challenge.

For raw-data-based sensing, one new attack on privacy is a credit-based inference attack. Specifically, the collector may infer if a node has submitted a data report for a task from the number of credits the node has earned. For example, suppose there are three tasks paid at a rate of 1, 2, and 5 credits. If a node has earned 3 credits, the collector can infer that it has taken the first two tasks. For another example, suppose the collector has published 10 tasks, 5 of which require a data report from Central Park, and each task is paid for one credit. If a node Bob has earned 6 credits, the collector can infer that Bob has submitted data for at least one of those 5 tasks. Thus, the collector knows that Bob has been to Central Park. How to address such attacks is also an open challenge.

Existing game-theory-based incentive schemes for participatory sensing have studied methods for setting the number of credits paid to each reporting node. How to integrate them into the framework of privacy-preserving incentive is another open problem.

ACKNOWLEDGMENT

This work was supported in part by the U.S. National Science Foundation (NSF) under grant number CNS-1320278.

Existing game-theory-based incentive schemes for participatory sensing have studied methods for setting the number of credits paid to each reporting node. How to integrate them into the framework of privacy-preserving incentive is another open problem.

REFERENCES

- [1] Q. Li and G. Cao, "Providing Privacy-Aware Incentives for Mobile Sensing," *Proc. IEEE PerCom*, 2013.
- [2] Q. Li and G. Cao, "Providing Efficient Privacy-Aware Incentives for Mobile Sensing," *Proc. ICDCS*, 2014.
- [3] Q. Li and G. Cao, "Efficient and Privacy-Preserving Data Aggregation in Mobile Sensing," *Proc. IEEE ICNP*, 2012.
- [4] Q. Li, G. Cao, and T. F. Porta, "Efficient and Privacy-Aware Data Aggregation in Mobile Sensing," *IEEE Trans. Dependable Secure Computing*, vol. 11, no. 2, 2014, pp. 115–29.
- [5] J. Krumm, "Inference Attacks on Location Tracks," *Proc. 5th Int'l. Conf. Pervasive Computing*, 2007, pp. 127–43.
- [6] C. Cornelius *et al.*, "Anonymsense: Privacy-Aware People-Centric Sensing," *Proc. ACM MobiSys*, 2008, pp. 211–24.
- [7] E. De Cristofaro and R. Di Pietro, "Preserving Query Privacy in Urban Sensing Systems," *Proc. ICDCN*, 2012, pp. 218–33.
- [8] D. Yang *et al.*, "Crowdsourcing to Smartphones: Incentive Mechanism Design for Mobile Phone Sensing," *Proc. ACM MobiCom*, 2012.
- [9] D. Chaum, "Blind Signatures For Untraceable Payments," *Proc. CRYPTO '82*, 1982.
- [10] M. Abe and T. Okamoto, "Provably Secure Partially Blind Signatures," *Proc. CRYPTO*, 2000, pp. 271–86.
- [11] V. Rastogi and S. Nath, "Differentially Private Aggregation of Distributed Time-Series with Transformation And Encryption," *Proc. ACM SIGMOD*, 2010.
- [12] E. Shi *et al.*, "Privacy-preserving Aggregation of Time-Series Data," *Proc. Network and Distrib. Sys. Security Symp.*, 2011.
- [13] C. Castelluccia *et al.*, "Efficient and Provably Secure Aggregation of Encrypted Data in Wireless Sensor Networks," *ACM Trans. Sensor Networks*, vol. 5, no. 36, 2009, pp. 20:1–20:3.
- [14] C. Dwork *et al.*, "Calibrating Noise to Sensitivity in Private Data Analysis," *TCC*, 2006.
- [15] Q. Li and G. Cao, "Efficient Privacy-Preserving Stream Aggregation in Mobile Sensing with Low Aggregation Error," *Proc. 13th Privacy Enhancing Technologies Symp.*, 2013.

BIOGRAPHIES

QINGHUA LI [M] received a B.E. degree from Xian Jiaotong University, an M.S. degree from Tsinghua University, and a Ph.D. degree from Pennsylvania State University. In 2013, he joined the University of Arkansas, where he is currently an assistant professor in the Department of Computer Science and Computer Engineering. His research interests are security and privacy in networked and mobile systems including mobile sensing, smart grid, and mobile cloud computing.

GUOHONG CAO [F] received a B.S. degree in computer science from Xian Jiaotong University and a Ph.D. degree in computer science from Ohio State University in 1999. Since then, he has been with the Department of Computer Science and Engineering at Pennsylvania State University, where he is currently a professor. His research interests include wireless networks, wireless security, smartphones, vehicular networks, wireless sensor networks, and distributed fault-tolerant computing. He has served on the Editorial Boards of *IEEE Transactions on Mobile Computing*, *IEEE Transactions on Wireless Communications*, and *IEEE Transactions on Vehicular Technology*, and has served on the organizing and technical program committees of many conferences, including the TPC Chair/Co-Chair of IEEE SRDS '09, MASS '10, and INFOCOM '13. He was a recipient of the NSF CAREER award in 2001.

Security and Privacy in Mobile Crowdsourcing Networks: Challenges and Opportunities

Kan Yang, Kuan Zhang, Ju Ren, and Xuemin (Sherman) Shen

ABSTRACT

The mobile crowdsourcing network (MCN) is a promising network architecture that applies the principles of crowdsourcing to perform tasks with human involvement and powerful mobile devices. However, it also raises some critical security and privacy issues that impede the application of MCNs. In this article, in order to better understand these critical security and privacy challenges, we first propose a general architecture for a mobile crowdsourcing network comprising both crowdsourcing sensing and crowdsourcing computing. After that, we set forth several critical security and privacy challenges that essentially capture the characteristics of MCNs. We also formulate some research problems leading to possible research directions. We expect this work will bring more attention to further investigation on security and privacy solutions for mobile crowdsourcing networks.

INTRODUCTION

With the rapid advances in mobile and communication technologies, most mobile devices today are equipped with powerful processors, various sensors, large memories, fast wireless communication modules, and so on. Because of these sophisticated components, mobile devices have become important tools to sense, communicate, and compute data. For instance, smartphones can be used to collect video/image data (with cameras), acoustic data (with microphone), location information (with GPS), and some other useful contextual information (with gyroscopes and accelerometers). They can also transmit data via cellular networks, WiFi, Bluetooth, NFC, and so on. According to the forecasts from Canalys, worldwide mobile device shipments, including notebook PCs, tablet PCs, smart phones, and phones, will reach 2.6 billion units by 2016.

Mobile devices (e.g., smartphones, tablets, wearable devices) are usually carried by humans, so they are more applicable in a crowdsourcing environment. Crowdsourcing is

the combination of two words, “crowd” and “outsourcing”, coined by Jeff Howe in 2006 [1], and defined as the “*act of taking a job traditionally performed by a designated agent (usually an employee) and outsourcing it to an undefined, generally large group of people in the form of an open call.*” Although there are many advantages of crowdsourcing, the most attractive one should be that it can bring massive intelligence to solve problems at an affordable price. Some tasks that are difficult for computers or individuals can be solved efficiently by crowdsourcing to a massive group of people, including image tagging, audio translation, and so on.

The mobile crowdsourcing network (MCN) is an emerging network paradigm that captures the advantages of powerful mobile devices and crowdsourcing. It applies the principles of crowdsourcing to perform tasks with human involvement and powerful mobile devices. In MCNs, crowdsourcing is involved in both data collection and data processing: *crowdsourcing sensing* [2–4] and *crowdsourcing computing* [5, 6]. Human mobility offers unprecedented opportunities for both data sensing and transmission with mobile devices. Mobile devices can sense the surroundings wherever their holders arrive, and the storage capability of a mobile device enables it to transmit data in a store-carry-and-forward way. Moreover, human capabilities also offer intelligent human computation with their devices. Different crowdsourcing applications may utilize different human capabilities, including human perception (understanding, feeling, intuition), intelligence, cognition, knowledge, visual recognition, common sense, experiences, and so on.

Due to human involvement and crowdsourcing, several challenging security and privacy concerns are raised in MCNs. For example, some sensed data may contain location information, which may implicitly reveal a mobile user’s movement. Compared to traditional networks, security and privacy issues in MCNs are more critical and challenging due to the following characteristics:

The authors are with the University of Waterloo.

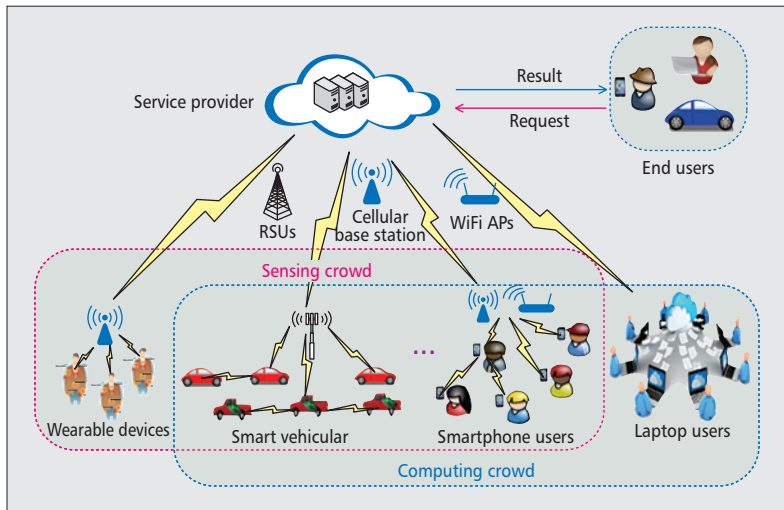


Figure 1. General architecture of mobile crowdsourcing networks.

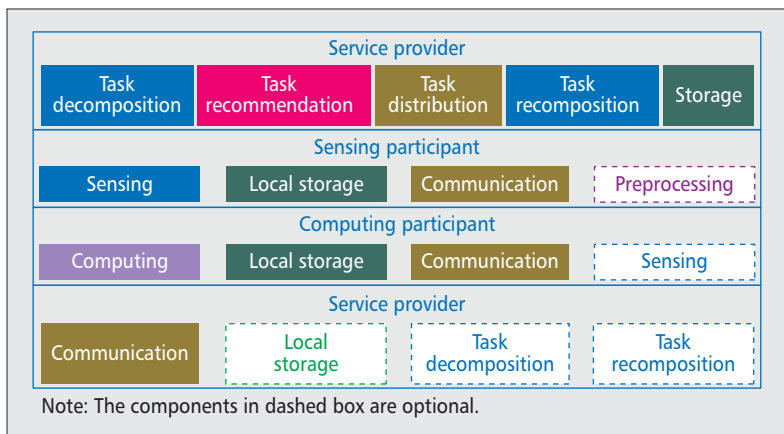


Figure 2. Components of each entity in mobile crowdsourcing networks.

Human involvement: The human is involved in both crowdsourcing sensing and crowdsourcing computing. The sensed data may contain not only sensitive information of mobile devices, but also private information of crowdsourcing participants. Moreover, mobile devices may be controlled by malicious holders to launch attacks.

Task crowdsourcing: Task crowdsourcing can raise big security concerns, especially when crowdsourcing tasks themselves contain sensitive information. When crowdsourcing tasks to a dynamic group of people, it is more difficult to protect the private information than only outsourcing tasks to a single server, as the size of the group cannot be predetermined.

Dynamic topology: Mobile users may accept crowdsourced tasks based on their interests, locations, or device conditions (residual battery, available sensors, etc.). The network topology may change over time due to human mobility and dynamic user join/leave, which may also increase the difficulty of exploring security and privacy solutions.

Heterogeneity: Various communication networks may be involved in MCNs, such as wireless sensor networks, cellular networks, WiFi, Bluetooth, and vehicular ad hoc networks

(VANETs). Besides, there are also a diversity of mobile devices participating in MCNs, which may produce heterogeneous data.

When participating in crowdsourcing sensing or crowdsourcing computing, mobile users consume their own resources (e.g., battery, cellular data, memory) and may suffer potential security and privacy threats. Although some incentive mechanisms [7, 8] are proposed to provide participants with enough rewards, if security and privacy cannot be guaranteed, many mobile users are still not willing to participate in and contribute to MCNs.

Aiming to address security and privacy challenges in MCNs well, in this article, we first describe several critical security and privacy challenges that capture the characteristics of MCNs. Then we point out some research problems that may lead to some possible research directions. We expect that this work will promote further investigation on security and privacy solutions for mobile crowdsourcing networks.

SYSTEM ARCHITECTURE OF MCNS

Figure 1 illustrates a general architecture of MCNs, which includes four basic types of entities: service provider, end users, sensing crowd, and computing crowd.

SERVICE PROVIDER

The service provider is a crowdsourcing platform that provides crowdsourcing services to both end users and public crowds. Generally, the service provider accepts service requests from end users, and partitions these tasks into several small tasks that can be crowdsourced. It then publishes these crowdsourced tasks on its service board and waits for the crowds to finish them. Upon receiving results from those crowdsourced tasks, the service provider performs the final process and sends the final results back to end users. In some scenarios, the service provider is only responsible for publishing decomposed crowdsourcing tasks from end users if end users choose to decompose and recompose tasks by themselves.

As illustrated in Fig. 2, the service provider is equipped with several fundamental components, including the task decomposition component, the task recomposition component, the task recommendation component, the task distribution component, and the data storage component. The task requested by the end user is first decomposed by the task decomposition component and then distributed to the crowds through the task distribution component, as shown in Fig. 3. The task recomposition component is responsible for performing the final process of the results from those crowdsourced tasks. The task recommendation component enables mobile users to submit subscription trapdoors (which are constructed under subscription policies defined by mobile users) to indicate their preferences on crowdsourced tasks. Once there is a crowdsourced task, the task recommendation component will check whether this task matches these subscription trapdoors. If the task matches the subscription trapdoor of a mobile user, the service provider will send an alert to this user.

There are two basic types of crowdsourced tasks: sensing and computing. Sensing tasks are designed to collect data from a crowd of mobile users who carry sensor-enabled mobile devices. Sensing tasks return the sensed data from sensing participants to the service provider, which may be stored in storage systems managed by the service provider or sent back to end users, depending on different applications. Computing tasks are designed to crowdsource the computation to a multitude of participants with their mobile computing devices, such as mobile phones, tablets, and smart cars. Usually, computing tasks also require human intelligence and capabilities, which are denoted as human computation [9].

The service provider is usually honest but curious in the sense that it may refer some personal information from sensed data (e.g., location information, identity, interests) and may also be interested in published tasks, computing results, as well as final results sent back to end users.

END USERS

End users are the customers who purchase or rent crowdsourcing services at certain costs. They send service requests to the service provider and receive results from it. As illustrated in Fig. 2, devices of end users should have basic communication components.

Under some circumstances, end users may decompose and recombine tasks by themselves and only rely on the service provider to publish their crowdsourced tasks. Together with decomposed crowdsourcing tasks, they may also provide some input data from their local storage systems or other purchased databases. In this setting, devices of end users may also be equipped with local storage components, task decomposition components, and task recombination components, which are described by dashed boxes in Fig. 2.

Similar to the service provider, end users are usually assumed to be honest but curious when they request crowdsourcing services from the service provider. Sometimes, the end user may also participate in crowdsourced tasks published by other end users. In this setting, the security assumption of the end user should be the same as the following assumption of crowdsourcing participants.

SENSING CROWD

The sensing crowd is a crowd of mobile users who accept and participate in crowdsourced sensing tasks. In order to perform sensing tasks, as shown in Fig. 2, the sensing devices of crowdsourcing participants should have sensing components, local storage components, and communication components. The sensing components may be cameras, microphones, GPS, gyroscopes, accelerometers, and so on, and the communication component may support cellular networks, WiFi, Bluetooth, NFC, and others. Besides, some of sensing devices should also be equipped with preprocessing components if they need to conduct some preprocessing on the sensed data.

There are a large number of crowdsourced sensing tasks published every day, but mobile

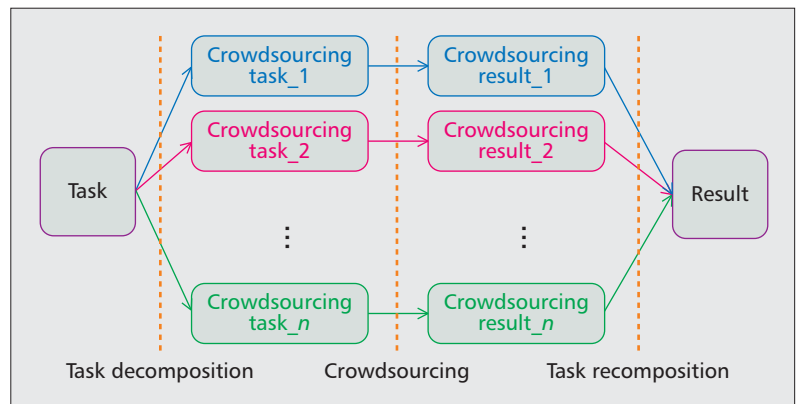


Figure 3. Task decomposition, crowdsourcing, and recombination.

users may be only interested in some of tasks. In order to filter sensing tasks, mobile users can provide *sensing subscriptions*, which indicate their preferences on sensing tasks. The preferences may be affected by conditions of mobile devices, and the interests and activities of mobile users. The sensing crowd is not always trusted. For instance, some sensing participants are malicious in the sense that they may report invalid data to the service provider or launch a distributed denial of service (DDoS) attack by accepting all the tasks without giving any results back.

COMPUTING CROWD

The computing crowd is a crowd of users who accept and participate in crowdsourced computing tasks. Based on the input data, computing tasks can be divided into two types.

Sensing-Based Computing Tasks: Sensing-based computing tasks take as inputs the current sensed data (e.g., the current GPS information) and the data provided by the service provider or end users (if applicable). These tasks are usually designed for people with sensor-enabled mobile devices. Different from crowdsourcing sensing, the sensed data here are directly used as computation inputs instead of being sent to the service provider.

Pure Computing Tasks: Pure computing tasks only take the data provided by the service provider (if applicable) as input. Thus, pure computing tasks can be accepted by people with any computing devices, such as mobile phones, laptops, and tablets.

As shown in Fig. 2, devices of computing participants should be equipped with computing components, local storage components, and communication components. Sensing components are also required if they accept sensing-based computing tasks. Actually, any mobile user with proper devices can be part of either sensing or computing crowds. Similar to a sensing crowd, a computing crowd may also provide *computing subscriptions*, which indicate their preferences on computing tasks. However, computing participants cannot be fully trusted either, as they may cheat and send back wrong results. Moreover, they may also be interested in the input data provided by the service provider or end users.

Due to the advantages of human involvement and powerful mobile devices, the architecture of

Application examples	Descriptions
Air pollution	Detect air pollution emitted by factories, cars, and farms.
Water quality	Monitor the water quality and study its eligibility for drinking.
Levels	Measurement of the energy radiated by cell stations and WiFi routers.
Smart navigation	Plan route according to weather conditions, accidents, and traffic jams.
Smart parking	Monitor parking space availability in the city and recommend with charges.
Smart traffic light	Control traffic lights according to traffic load and emerging events.
Health monitoring	Monitor health status from heart rate, electrocardiography, blood pressure, etc.
Disease diagnosis	Diagnose the disease from personal health parameters, and other cases.
Food recommendation	Recommend food or drinks according to personal health conditions.

Table 1. Application examples of mobile crowdsourcing networks.

MCNs can be well applied in many applications. Table 1 provides some mobile crowdsourcing application examples.

SECURITY AND PRIVACY CHALLENGES IN MCNs

Inspired by the characteristics of MCNs, we point out some security and privacy challenges, including privacy threats, reliability threats, and availability threats.

PRIVACY THREATS

In MCNs, the privacy may be leaked out from both the data and the task. Here, we discuss two privacy threats: *privacy threats from data* and *privacy threats from tasks*.

Privacy Threats from Data: We first describe privacy threats by analyzing data flows in MCNs. The data flows in sensing tasks are different from the ones in computing tasks in MCNs: In sensing tasks, the data are first sensed by the sensing crowd, transmitted to the service provider, then stored in the storage system managed by the service provider or sent back to end users. In computing tasks, the input data may come from various sources, for example, the storage system managed by the service provider, local storage systems from end users, cloud storage systems purchased by end users/the service provider, sensor-enabled mobile devices, and so on. The output data of computing tasks are first sent back to the service provider and finally to the end user after the task recomposition. Based on different data flows in MCNs, the privacy threats can be summarized as follows.

- **Privacy of Sensed Data:** The sensed data may contain sensitive information of sensing participants, such as identities, location information, biometric information, and so on. For example, the location information can be easily obtained either from GPS receivers embedded in mobile devices or triangulation based on WiFi or cellular networks. Moreover, some environmental data (e.g., precise air temperature, the light, the noise, etc.) may also reveal the location information. The disclosure of location information may leak the privacy of participants, such as home and workplace locations, routines, habits, etc.

- **Privacy of Computing Inputs:** The input data of crowdsourced computing tasks may also contain sensitive information, such as business financial records, proprietary research data, or personal health information. When sending the data to the computing crowd, it may leak out the private information of data contributors, data owners or other related people.

- **Privacy of Computing Results:** The output results of crowdsourced computing tasks may be sensitive or private. End users do not want the service provider to know the contents of the results or obtain some sensitive information from the output results. In some scenarios, even computing participants are not allowed to know the contents of the computing results.

Privacy Threats from Tasks: Besides the privacy leakage from data, the task itself may also reveal some private information of both end users and crowdsourcing participants, denoted as *Task Privacy of End Users* and *Task Privacy of Participants*.

- **Task Privacy of End Users:** The content of the task may reveal sensitive information of end users to the service provider. For example, if an end user publishes crowdsourcing tasks that can only be accepted by psychologists, the service provider may infer that this end user may suffer from some psychological diseases.

- **Task Privacy of Participants:** Some tasks may also leak out private information about crowdsourcing participants. For example, if a crowdsourcing participant accepts a temperature measurement task at a particular location X at time Y, it may reveal that this participant will be at location X at time Y. However, participants may not want to leak their current location when they are tasked. In this example, the participant cannot hide the exact location information and the time when executing the task. One countermeasure is to hide their identities when taking the task and reporting results, such that the service provider only knows that some participants are at location X at time Y, but it does not know exactly who.

RELIABILITY THREATS

In MCNs, any one with a mobile device or a computing device can accept crowdsourced tasks. Due to task crowdsourcing and human involvement, it is difficult to guarantee every participant to provide reliable data or computing results.

Reliability of Sensed Data: Some malicious sensing participants may report incorrect or invalid sensed data, which is referred to as a pollution attack. Although pollution attacks exist in traditional sensor networks, it is much more

challenging to detect and resist pollution attacks in MCNs due to the following reasons:

- Any adversary can be a participant in MCNs and provide pollution data.
- Powerful mobile devices can be configured by the adversary to jam specific pollution data.
- The anonymous mechanism of privacy protection also increases the difficulty of pollution attack detection.

Reliability of Computing Results: The adversary may also be a computing participant; thus, not all computing participants are honest. Malicious computing participants may provide invalid or incorrect results in order to save their computing resources. Sometimes, honest computing participants may also provide wrong results due to many reasons, for example, misunderstanding directions, or making mistakes due to personal bias or lack of experience. Because end users can only receive final results from the service provider, they can only verify whether or not final results are correct. However, it is difficult for them to identify the dishonest computing participants who provide the wrong results when final results are incorrect. Thus, the verification of the crowdsourcing result should also be performed by the service provider, which is a challenging issue as the service provider does not know the contents of the results due to privacy requirements.

Reliability of Transmission: The sensed data may be transmitted to the service provider via various channels, including 3G/4G, WiFi, VANET, or relayed by other devices. The reliability of transmission channels in these networks has been well studied in the literature. However, in MCNs, data may be tampered with by intermediated devices due to human involvement in data transmission. For example, some malicious participants may selectively discard certain data packages or modify data content.

AVAILABILITY THREATS

Extensive studies have been done on some denial of service (DoS) issues for traditional networks in the literature, such as network congestion due to message floods. In MCNs, however, several unique DDoS issues are raised due to task crowdsourcing:

- *DDoS by malicious participants:* Some malicious participants may accept all the crowdsourced tasks but refuse to give valid results or even ignore the tasks. This may cause the DDoS where valid crowdsourcing participants cannot get any tasks since the total number of crowdsourced tasks are usually fixed.
- *DDoS by honest but selfish participants:* A DDoS attack may also happen even when all the participants are honest. For example, a participant who is selfish and hopes to receive more rewards may accept all the computing tasks and complete them over a long time period due to its limited computation capabilities. This will also eliminate the advantage of diverse contributors in crowdsourcing as the results are from the same person.

It is challenging to prevent these DDoS attacks in MCNs because of the characteristics

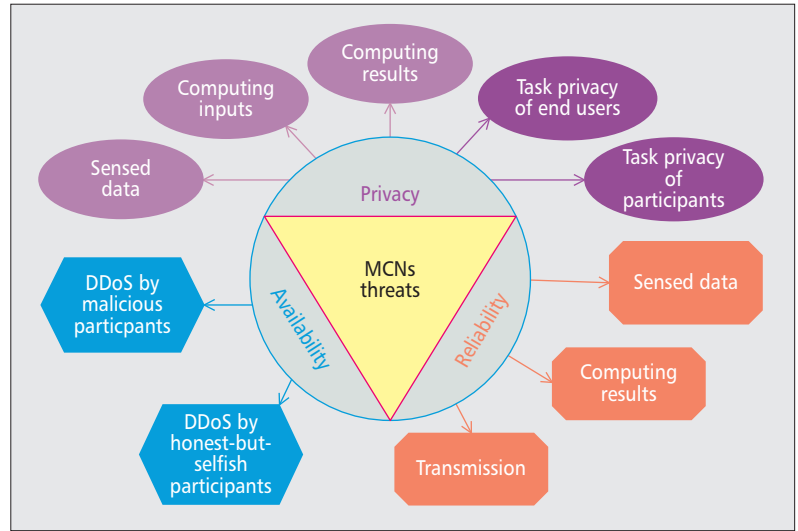


Figure 4. Security and privacy threats in MCNs.

of MCNs, including task crowdsourcing, human involvement, and dynamic topology. To deal with these availability threats, a series of methods may be required. First, all the users should be authenticated before joining the mobile crowdsourcing network. Second, a reputation mechanism is necessary to evaluate the reliability of crowdsourcing participants and detect unreliable crowdsourcing participants. Third, novel incentive mechanisms may also be required to provide fairness in MCNs. Figure 4 briefly summarizes the security and privacy threats that capture the characteristics of MCNs.

SECURITY AND PRIVACY OPPORTUNITIES IN MCNs

There are many potential research problems from privacy threats, reliability threats, and availability threats. Some traditional user privacy concerns have been discussed in [10]. In this section, we formulate several new research problems that may provide some possible research directions.

AUTHENTICATION OF CROWDSOURCING PARTICIPANTS

In order to cope with availability and reliability threats, it is necessary to authenticate crowdsourcing participants before they join the network. Normally, the authentication is conducted by the service provider. However, in some applications, the service provider is just a platform for task recommendation and distribution, and does not have the capability or privilege to authenticate crowdsourcing participants. Under this circumstance, crowdsourcing participants should be authenticated by end users. However, due to the large number of crowdsourcing participants, it is not efficient to let end users perform authentication. Therefore, a distributed authentication mechanism is desired to provide efficient and reliable authentication service for crowdsourcing participants.

In MCNs, the tags should be generated by end users, and the subscription policy should be generated by mobile users. Therefore, efficient methods are desired to protect both task privacy and subscription privacy of task recommendation in MCNs.

One possible approach is to let the existing crowdsourcing participants authenticate new mobile users. However, some of the existing crowdsourcing participants may be malicious and collude to authenticate invalid mobile users, which means that they may let unauthorized users join the system. To avoid collusion authentication, we can utilize a threshold-based group authentication method. In a (t, n) threshold group authentication scheme, at least t existing crowdsourcing participants can act as a group authenticator, while any less than t existing crowdsourcing participants cannot authenticate new mobile users successfully. Thus, this can tolerate at most $t - 1$ malicious existing crowdsourcing participants in MCNs. However, how to efficiently initialize the system remains a challenging issue in a threshold group authentication scheme.

Another challenging problem in a (t, n) threshold group authentication scheme is the equality of authentication privilege. In other words, the newly joined crowdsourcing participant should also be able to provide group authentication to other new mobile users, together with existing crowdsourcing participants in the system. Specifically, suppose the authentication secret s of a registered mobile user is shared among n existing crowdsourcing participants. When a new mobile user becomes a crowdsourcing participant, this secret s should be shared among the new $n + 1$ crowdsourcing participants. Here comes another critical problem: how to reshare the secret from n crowdsourcing participants to the new $n + 1$ crowdsourcing participants. To solve this problem, we can refer to the method in [11], which lets each existing participant further share its secret piece into secret sectors. The new crowdsourcing participants can then reconstruct the secret sectors and obtain the new shared secret piece. However, this may incur heavy communication overhead among the crowdsourcing participants. Therefore, how to reduce the communication overhead for distributed threshold-based group authentication becomes a promising research direction.

PRIVACY-PRESERVING TASK RECOMMENDATION

The task recommendation component provides a platform for task publication and subscription. Mobile users subscribe the tasks of their interests from the publishers by submitting subscription trapdoors. Due to the honest but curious service provider, the privacy issues become much more critical in task publication and subscription. There are two major concerns:

- *Task privacy.* When publishing crowdsourced tasks, end users do not want the service provider and other unauthorized participants to access their published tasks (including the corresponding input data).
- *Subscription privacy.* Mobile users also do not want the service provider to know what types of tasks that are of interest to them. Moreover, end users may hope to define the access policy of crowdsourced tasks themselves, and mobile users also hope to define the subscription policy themselves.

To protect task privacy, end users usually encrypt tasks and the input data before publishing to the honest but curious service provider, such

that the service provider cannot know the contents of tasks and their data without decryption keys. The content of the subscription trapdoor, such as interests and status of mobile devices, should also be encrypted to prevent the service provider from knowing the type of tasks in which mobile users are interested. However, it is difficult to encrypt tasks or trapdoors in MCNs because of the large number of crowdsourcing participants. For example, traditional public key encryption methods require a publisher to encrypt tasks with different keys for different users. It may produce many copies of encrypted tasks in the system, the number of which is proportional to the number of mobile users. Moreover, the publisher needs to know public keys of mobile users beforehand, which is impossible in MCNs. Toward symmetric key encryption methods, the publisher needs to be always online to distribute keys. Similar problems hold for the encryption of trapdoors.

Fortunately, attribute-based encryption (ABE) [12] is a new encryption technique that only produces one copy of the ciphertext. Specifically, the task is encrypted under an access policy defined by the encryptor. The secret key of the decryptor is associated with credentials/attributes of the decryptor. The ciphertext can be decrypted only when attributes of the decryptor can satisfy the access policy in the ciphertext. Thus, we can apply attribute-based access control schemes [13, 14] to protect task privacy.

To protect subscription privacy and allow subscribers to define subscription policy, a straightforward method is to construct a subscription trapdoor by using ABE with another set of parameters. However, this requires the authority responsible for attribute management and key generation in an ABE system to generate tags for each task or subscription policy for each query. In MCNs, the tags should be generated by end users, and the subscription policy should be generated by mobile users. Therefore, efficient methods are desired to protect both task privacy and subscription privacy of task recommendation in MCNs.

PRIVACY-PRESERVING VERIFIABLE COMPUTATION OUTSOURCING

To protect the privacy of input data and output results of computing tasks, a straightforward method is to encrypt the input data and do the computation directly on encrypted data by using homomorphic encryption methods such that the results, when decrypted, match computation carried on unencrypted data. However, the high overhead of homomorphic encryption makes it far from applicable in practice. Thus, new methods are required to achieve computation outsourcing in MCNs.

Toward the reliability of the results from crowdsourced computing tasks, both end users and the service provider should be able to verify the correctness of computing results. But end users may not have sufficient computation resources. Thus, the ultimate goal is to design efficient verifiable computation outsourcing protocols that can minimize the computational overhead of end users. Of course, the overhead incurred by the correctness verification of the

computation should be substantially smaller than running the computation itself. Another ambitious goal is to design protocols that minimize the communication overhead between end users/service provider and crowdsourcing participants.

However, some crowdsourced computing tasks contain human computation associated with human intelligence. In this case, computing results may be subjective to individual sensitivity and experience, and thus these results cannot be easily verified by verifiable computing protocols. Sometimes, crowdsourcing participants may also act honestly, but misunderstand task directions or make mistakes due to personal bias or lack of experience. Other methods may also be applied to ensure the reliability of results, such as building a reputation system, redundancy tasks, and statistical filtering. When end users are malicious, neither crowdsourcing participants nor end users can give unbiased verification results. In this context, similar to third party data integrity checking [15], a trusted third party may be employed to verify computing results.

CONCLUSION

The mobile crowdsourcing network is a promising paradigm in a ubiquitous computing era. In this article, we have proposed a general architecture of MCNs containing crowdsourcing sensing and crowdsourcing computing. We have outlined several critical security and privacy threats that capture the characteristics of MCNs. We have also formulated several research problems that lead to future research directions on security and privacy solutions for MCNs.

ACKNOWLEDGMENT

This research has been supported by a research grant from the Natural Science and Engineering Research Council (NSERC), Canada.

REFERENCES

- [1] J. Howe, "The Rise of Crowdsourcing," *Wired*, vol. 14, no. 6, 2006, pp. 1–4.
- [2] R. K. Ganti et al., "Mobile Crowdsensing: Current State and Future Challenges," *IEEE Commun. Mag.*, Nov. 2011, vol. 49, no. 11, pp. 32–39.
- [3] D. Christin et al., "A Survey on Privacy In Mobile Participatory Sensing Applications," *J. Systems and Software*, vol. 84, no. 11, 2011, pp. 1928–46.
- [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] M. Conti et al., "From Opportunistic Networks to Opportunistic Computing," *IEEE Commun. Mag.*, vol. 48, no. 9, 2010, pp. 126–39.
- [6] K. Parshotam, "Crowd Computing: A Literature Review and Definition," *Proc. South African Institute for Computer Scientists and Information Technologists Conference (SAICSIT '13)*; ACM, 2013, pp. 121–30.
- [7] A. Singla and A. Krause, "Truthful Incentives in Crowdsourcing Tasks Using Regret Minimization Mechanisms," *Proc. 22nd Int'l. Conf. World Wide Web*, 2013, pp. 1167–78.
- [8] H. Zhou et al., "Consub: Incentive-Based Content Subscribing in Selfish Opportunistic Mobile Networks," *IEEE JSAC*, no. 99, 2013, pp. 1–11.

- [9] A. J. Quinn and B. B. Bederson, "Human Computation: A Survey and Taxonomy of a Growing Field," *Proc. SIGCHI Conf. Human Factors in Computing Sys.*, ACM, 2011, pp. 1403–12.
- [10] Y. Wang, Y. Huang, and C. Louis, "Respecting user Privacy in Mobile Crowdsourcing," *Science*, vol. 2, no. 2, 2013, pp. pp–50.
- [11] K. Yang et al., "Threshold Key Redistribution for Dynamic Change of Authentication Group in Wireless Mesh Networks," *IEEE GLOBECOM '10*, 2010, pp. 1–5.
- [12] V. Goyal et al., "Attribute-Based Encryption for Fine-Grained Access Control of Encrypted Data," *Proc. 13th ACM Conf. Computer and Commun. Security*, New York, NY, USA: ACM, 2006, pp. 89–98.
- [13] R. Lu, X. Lin, and X. Shen, "SPOC: A Secure and Privacy-Preserving Opportunistic Computing Framework for Mobilehealthcare Emergency," *IEEE Trans. Parallel Distrib. Sys.*, vol. 24, no. 3, 2013, pp. 614–24.
- [14] K. Yang et al., "DAC-MACS: Effective Data Access Control for Multiauthority Cloud Storage Systems," *IEEE Trans. Info. Forensics Security*, vol. 8, no. 11, 2013, pp. 1790–1801.
- [15] K. Yang and X. Jia, "Data Storage Auditing Service in Cloud Computing: Challenges, Methods and Opportunities," *World Wide Web*, vol. 15, no. 4, 2012, pp. 409–28.

BIOGRAPHIES

KAN YANG (kan.yang@uwaterloo.ca) received his B.Eng. degree in information security from the University of Science and Technology of China in 2008 and his Ph.D. degree in computer science from the City University of Hong Kong in August 2013. He is currently a postdoctoral fellow with the Broadband Communications Research (BBCR) group in the Department of Electrical and Computer Engineering at the University of Waterloo, Canada. His research interests include cloud security and privacy, big data security and privacy, mobile security and privacy, big data mining, applied cryptography, wireless communication and networks, and distributed systems.

KUAN ZHANG (k52zhang@bbcr.uwaterloo.ca) received his B.Sc. degree in electrical and computer engineering and M.Sc. degree in computer science from Northeastern University, China, in 2009 and 2011, respectively. He is currently working toward a Ph.D. degree with the BBGR Group, Department of Electrical and Computer Engineering, University of Waterloo. His research interests include security and privacy for mobile social networks.

JU REN (ren_ju@csu.edu.cn) received his B.Sc. and M.Sc. degrees in computer science from Central South University, China, in 2009 and 2012, respectively. He is currently a Ph.D. candidate in the Department of Computer Science at Central South University, China. From August 2013 to the present, he is also a visiting Ph.D. student in the Department of Electrical and Computer Engineering, University of Waterloo. His research interests include wireless sensor networks, mobile sensing/computing, and cloud computing.

XUEMIN (SHERMAN) SHEN (xshen@bbcr.uwaterloo.ca) is a professor and University Research Chair, Department of Electrical and Computer Engineering, University of Waterloo. He was the associate chair for Graduate Studies from 2004 to 2008. His research focuses on resource management in interconnected wireless/wired networks, wireless network security, social networks, smart grid, and vehicular ad hoc and sensor networks. He served as the Technical Program Committee Chair/Co-Chair for IEEE INFOCOM '14, IEEE VTC '10-Fall, Symposia Chair for IEEE ICC '10, Tutorial Chair for IEEE VTC '11-Spring and IEEE ICC '08, and Technical Program Committee Chair for IEEE GLOBECOM '07. He also serves or has served as Editor-in-Chief of *IEEE Network*, *Peer-to-Peer Networking and Application*, and *IET Communications*. He is a registered Professional Engineer of Ontario, Canada, an Engineering Institute of Canada Fellow, a Canadian Academy of Engineering Fellow, and a Distinguished Lecturer of the IEEE Vehicular Technology and Communications Societies.

Some crowdsourced computing tasks contain human computation associated with human intelligence. In this case, computing results may be subjective to individual sensitivity and experience, and thus these results cannot be easily verified by verifiable computing protocols.

Security in Space Information Networks

Chunxiao Jiang, Xuexia Wang, Jian Wang, Hsiao-Hwa Chen, and Yong Ren

ABSTRACT

Space information networks were proposed to broaden the observation area and realize continuous information acquisition using satellites and high altitude platform stations. Space information networks are able to enhance detection and transmission capabilities compared to the current single Earth observation satellite. Although lots of technical work has been done concerning the space network architecture and protocols, the security issues have not been investigated well. In this article, we focus on the security problems in space information networks from four perspectives, that is, secure handoff, secure transmission control, key management, and secure routing. Existing works, together with their challenges and open problems, are discussed, and our proposed scheme is introduced. Overall, this article aims to help readers understand the motivation, problem formulation, and solutions regarding security issues on space information networks.

INTRODUCTION

Currently, Earth observations are based mainly on single satellites or satellite constellations, leading to a dilemma that continuous information capture and transmission are difficult to accomplish. Meanwhile, satellite ground stations cannot be deployed around the world in every country. This exacerbates the problem due to the fact that ground stations cannot receive data from satellites unless satellites come into certain detection ranges (i.e., the limitation of view angle of satellites above the horizon). In addition, the lack of collaborative management and assignment of satellite resources reduces the efficiency of information acquisition and transmission. Recently, rapid development in space technologies has catalyzed the emerging space information networks [1], which are composed of geosynchronous Earth orbit (GEO) satellites, medium Earth orbit (MEO) satellites, low Earth orbit (LEO) satellites, high-altitude platform stations (HAPSs), and so on. With the help of space information network technology, the coverage area and effectiveness of emergency monitoring can be improved significantly. As an example, shown in Fig. 1, an event occurring on the other side of the earth can be detected and transmitted back to a ground station located in

China through space information networks. Therefore, space information networks are expected to play a key role in disaster prevention and relief, emergency rescue, global location and navigation, space tracking and control, and so on.

As space information networks become more and more important in both civilian and military applications, the security requirements also become more demanding. In order to protect space communications and applications, some key secure technologies, including secure handoff, secure transmission control, key management, and secure routing, should be integrated into the space network architecture. Space information networks have several unique characteristic features: lack of fixed communication infrastructure, extremely long and variable propagation latency, asymmetrical forward and reverse link capacities, high link bit error rate, intermittent link connectivity, heterogeneous network integration, and high reliability and security requirements. The aforementioned characteristics make research on space information networks rather challenging. In recent years, researchers have done a lot of work regarding space network structure, network convergence, network management, routing, mobility management, and transmission control. However, the security related issues have not been investigated well.

In this article, we focus on the security issues in space information networks, where we survey the literature and propose our schemes, enabling readers to understand the motivation, problems, and solutions of secure space information networks. Specifically, we elaborate on the key technical issues, corresponding solutions, potential research directions, and challenges from four aspects, including secure handoff, secure transmission control, key management, and secure routing. The rest of this article is outlined as follows. In the next section, we review secure handoff technologies in space information networks with related works and open problems in the literature. Then we discuss the secure transmission control issues and their challenges. Following that, the issues of key management in space information networks are addressed in four different categories: centralized, distributed, topology-based, and preconfigured. Finally, related works on secure routing are surveyed together with our proposed intrusion-detection-based secure routing scheme.

Chunxiao Jiang, Xuexia Wang, Jian Wang, and Yong Ren are with Tsinghua University.

Hsiao-Hwa Chen is with National Cheng Kung University.

SECURE HANDOFF IN SPACE INFORMATION NETWORKS

Due to the rapid change in their relative positions among different layers of satellites, near-space platforms, and ground terminals, handoff needs to take place frequently, during which signaling information is exchanged among different nodes. For example, as shown in Fig. 2, when an aircraft is switching its connection from one LEO satellite to a GEO satellite, handoff information involves the previously accessed satellite and the newly accessed satellite, as well as the network control center. In such a case, when handoff happens in space information networks, signaling information might be eavesdropped, falsified, or fabricated. In order to guarantee basic security requirements (e.g., authenticity, confidentiality, and integrity), the signaling message should be meticulously protected during the handoff processes.

In a wireless heterogeneous network, handoff technology can be classified into three aspects, including horizontal handoff, vertical handoff, and diagonal handoff [2]. Horizontal handoff happens in a homogeneous network using the same access technology, such as the handoff among different base stations in cellular networks. In LEO satellite networks, due to the dynamics of satellites and their corresponding coverage areas, the switch of ground terminals from one LEO satellite to others belongs to horizontal handoff, as shown in Fig. 2. In contrast, vertical handoff refers to the traversal between two different access technologies, which is a characteristic feature of heterogeneous networks, for example, the handoff between cellular base station and WLAN access. From the perspective of the space-air-ground integrated network, ground terminals or aircraft equipped with multiple access technologies can switch between satellites and near-space information networks in a vertical manner. Diagonal handoff is the combination of horizontal and vertical handoffs, which is relatively rare in space information networks. A handoff is said to be diagonal when a node traverses those networks that use a common underlying technology (as in IEEE 802 family networks), and it allows a user to continue its applications with the required quality of service (QoS) from one network to another.

In the literature, the handoff schemes in satellite networks were surveyed in [3], especially the ground terminals' horizontal handoff in LEO satellite networks. The handoff schemes were classified into link-layer and network-layer schemes, where the link-layer handover schemes were further divided into three subcategories (i.e., spotbeam handover schemes, satellite handover schemes, and inter-satellite link, ISL, handover schemes), and the network-layer handover schemes were also further classified depending on their connection switch strategies. The survey in [3] focused mainly on the fault-tolerant and integrity performance of various handoff schemes, while the authenticity and confidentiality issues were not thoroughly investigated. Another work in [4] summarized the secure handoff optimization schemes for wireless het-

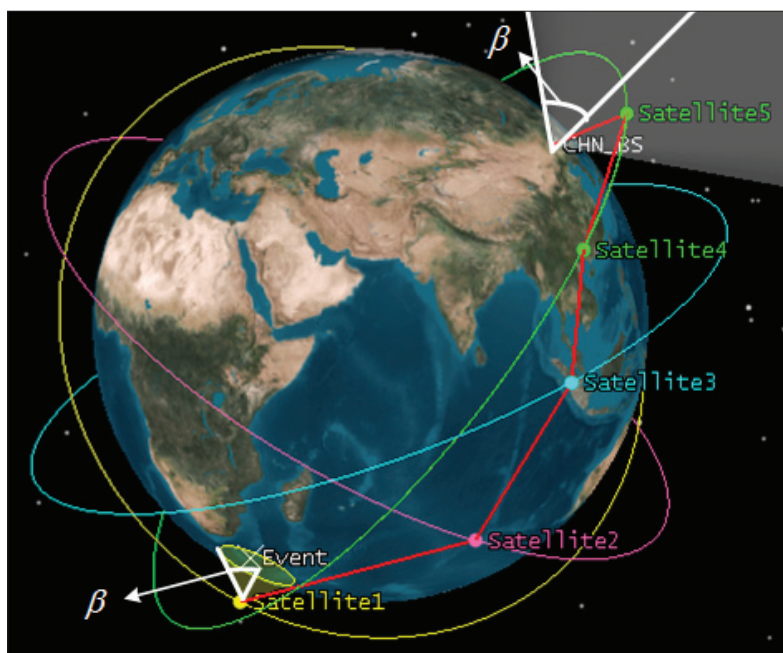


Figure 1. Example of event detection using space information networks.

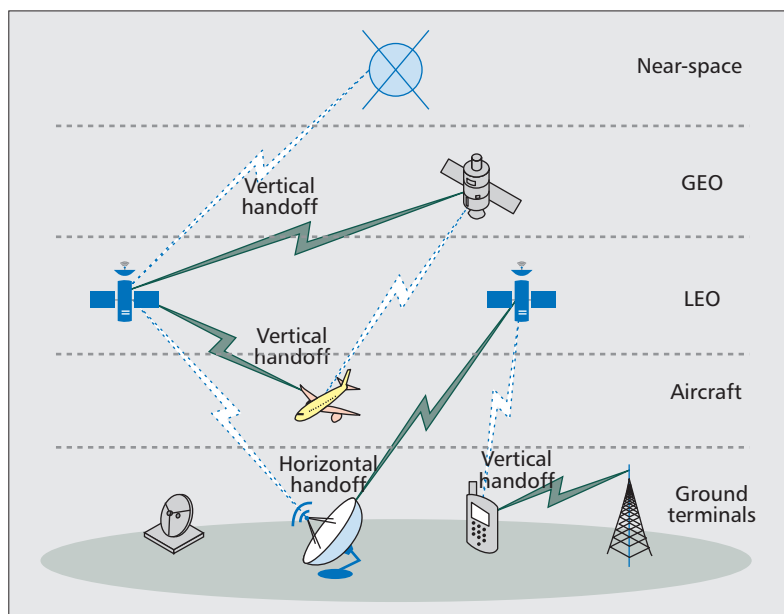


Figure 2. Handoff scenarios in space information networks.

erogeneous networks, including authentication, authorization, and accounting (AAA) context transfer, optimized integrated registration procedure of mobile IP and SIP with AAA operations (OIRPMSA), media-independent pre-authentication (MPA), shadow registration, and so on. Meanwhile, those schemes were compared in terms of public vs. secret key, mutual authentication, privacy, and non-repudiation, which assumed trust between domains. Recently, He *et al.*, proposed a secure and efficient handover authentication scheme based on bilinear pairing functions in [5]. To prevent denial-of-service (DoS) attacks, the authors also proposed a polynomial-based lightweight verification scheme.

There are several open issues regarding

Similar to terrestrial networks, key management is also of importance in space information networks. All cryptography-based security technologies are associated with key management, which penetrates into communication and security protocols at every layer.

secure handoff in space information networks. First, the implementation of vertical handoff is challenging due to high dynamics and the extremely long latency of satellites. Specifically, when an aircraft is switching from a MEO satellite to a near-space network, a dilemma would occur due to the high speed of the aircraft vs. the long-delay signaling information exchange. Moreover, a secure vertical handoff with an additional authentication or encryption phase is even more difficult. Second, for the secure message transmission, the Internet Engineering Task Force (IETF) recommended to use an authentication header (AH) protocol or the encapsulating security payload (ESP) of IPsec, as well as Internet key exchange (IKE). However, the key agreement through IKE has to rely on a pre-shared key or public key authentication. As a result, IPsec cannot be directly employed in space information networks, which require a more efficient method to reduce latency and signaling cost during handoff processes. Third, the issue of how to control the ratio of encryption overhead and the amount of information to be transmitted should also be considered. Due to relatively high bit error rate (BER) and long link delay, the space information network cannot afford a complicated encryption scheme. In such a case, the trade-off between encryption cost and security performance becomes a remarkable problem. Finally, when it comes to simultaneous authentication of multiple satellites, novel access authentication schemes for congestion avoidance are also required.

SECURE TRANSMISSION CONTROL IN SPACE INFORMATION NETWORKS

Space information networks are characterized by long round-trip time (RTT), asymmetrical broadband, and high dynamic topology, all of which make the traditional terrestrial Internet protocols incompetent. Especially, the pervasive and everlasting TCP is confronted with challenges in the implementation of space information networks due to the fact that TCP was originally designed for the terrestrial Internet. First, TCP features the use of window-based transmission control algorithms. Based on the sliding window flow control mechanisms, TCP can manage the amount of data a source can transmit by adjusting the window size according to the acknowledgment information from the destination. Second, it is well known that TCP cannot distinguish data losses caused by link failure from those caused by network congestion. These issues may not cause obvious performance degradation in the terrestrial Internet. However, when it comes to extreme space communication environments, involving, for example, long link delays, noisy channels, and asymmetric channel rates, the aforementioned problems would become obvious and even serious.

In the literature, Wang *et al.*, surveyed the reliable data transport protocols for space information networks in [6]. Seventeen different types of data transport protocols available for space information networks were divided into three categories:

- Protocols involving only changes to TCP, such as Stream Control Transmission Protocol (SCTP), Satellite Transport Protocol (STP), TCP-Peach, and TCP-Planet
- Protocols involving changes to TCP and/or network infrastructure, for example, explicit congestion control (XCP), explicit and fair window adjustment (REFWA), and performance enhancing transport architecture (PETRA)
- Delay-tolerant networking (DTN) and other protocols, such as Bundle Protocol (BP), Consultative Committee for Space Data Systems (CCSDS) File Delivery Protocol (CFDP), and Licklider Transmission Protocol (LTP)

One can see that most of the existing works on transport protocols in space information networks are amended or extended from the traditional TCP. As discussed earlier, in order to adapt to the special working environment of space information networks, the traditional TCP has to mutate to overcome the aforementioned problems. A more recent work on variant-TCP design for satellite networks can be found in [7], where an improved TCP-friendly rate control (TFRC) protocol was proposed. The authors designed a loss differentiation algorithm to differentiate packet losses and avoid TFRC's misclassifications, which depends on the queuing delay calculated based on the RTT measurement and makes the loss event rate calculation of the TFRC more accurate.

As far as securing data transmission in space information networks, most existing works focus on two aspects: the space communications protocol specifications (SCPS) series proposed by the CCSDS, and the satellite IP series with the corresponding amendments. They address issues in improvement of transmission performance and reliability, and the secure transmission schemes relying mostly on network layer security technologies, for example, IPsec and SCPS-SP (security protocol) in CCSDS. In the typical transport layer security protocol transport layer security its predecessor, secure sockets layer (TLS/SSL), certification-based public key systems were widely adopted, such as Rivest-Shamir-Adleman (RSA) and elliptic curve cryptography (ECC). However, TLS requires public key transmission and verification between clients and servers, resulting in long handshake latency. Although some researchers proposed to incorporate ID-based cryptography (IBC) into the hand-shake protocol to shorten handshake duration [8], such a protocol suffers from high computational and communication costs, which cannot be applied directly to space information networks. In [9], the authors proposed a TCP header compression method to reduce the header overhead and improve the performance of satellite-based virtual private networks (VPNs), while at the same time guaranteeing a given encryption performance. The proposed solution generated a cryptographic hash of flow identification information and facilitated TCP performance enhancing proxy (PEP) agent and IPsec coexistence.

From the above literature review, we can see that secure transmission control in space networks is still a widely open problem. First, for

Kinds	1-affects- <i>n</i>	Advantages	Disadvantages
Centralized	No	<ul style="list-style-type: none"> • Real time, high efficiency • Large-scale networks 	<ul style="list-style-type: none"> • High computational cost • Not suitable for dynamic networks
Distributed	Several	Low computational cost	<ul style="list-style-type: none"> • Only for small-scale network utility • Additional message exchange cost
Topology-based	Medium	<ul style="list-style-type: none"> • High efficiency • Large-scale networks 	<ul style="list-style-type: none"> • Vulnerable cluster head or root node • Not suitable for dynamic networks
Preconfiguration	Several	<ul style="list-style-type: none"> • Mobile/dynamic networks • Short key-agreeing time 	Complex rekeying and updated key issue

Table 1. Comparison of different kinds of key management schemes.

basic transmission control, how to mitigate high BER, long round-trip delay, asymmetric channels, and intermittent connection in space information networks is one major issue. For instance, typical BER in a satellite communication link can be 10^{-2} ; the propagation delay from a ground terminal to a LEO spacecraft is about 25 ms, 110 ms to a medium earth orbit (MEO) satellite, and 250 ms to geostationary earth orbit (GEO) satellite. The ratio of downlink bandwidth and uplink bandwidth can even be roughly 1000:1. Second, when the security requirement is added, authentication and encryption have to be taken into account. This would inevitably incur more time and communication costs, making secure transmission control even more challenging. How to trade off between security performance and additional security costs is another serious issue. Finally, some traditional TCP vulnerability problems (e.g., TCP sequence number prediction/emulation) may also appear in space information networks, and thus need to be redesigned considering the unique features of space information networks.

KEY MANAGEMENT IN SPACE INFORMATION NETWORKS

Similar to terrestrial networks, key management is also of importance in space information networks. All cryptology-based security technologies are associated with key management, which penetrates into communication and security protocols at every layer. For example, in secure handoff and secure routing, public key and symmetric key cryptosystems have been widely adopted. Contrary to terrestrial wireless networks, space information network nodes spread over an extremely large geographic area, leading to difficulty in building a powerful online key management center. Meanwhile, the communication costs among different space information network nodes are extremely high due to the complex space communication environment. Nevertheless, the space information network nodes (e.g., satellites or near-space platforms) possess relatively high-performance hardware and high computational capacity. Therefore, traditional key management solutions cannot be applied directly to space information networks. Researchers have to redesign appropriate

encryption, authentication, and integrity protection mechanisms based on the characteristics of space information networks. Up to now, most of the existing key management works for space information networks focused on satellite networks, while the problems in satellite and near-space hybrid multicast communication networks are still not well investigated.

In the literature, the existing key management schemes can be classified into four categories, as follows.

Centralized Key Management: An online central server with powerful capability is in charge of the key management task of the whole network, including generating and distributing keys.

Distributed Key Management: All nodes contribute their share to generate a common key with several rounds of message exchanges in a public channel.

Topology-Based Key Management: The schemes are optimized according to a special network topology structure.

Pre-Configuration Key Management: The network nodes can get the keys in advance from an off-line key management server.

In Table 1, we summarize the advantages and disadvantages, comparing those key management schemes where the 1-affects-*n* problem means when one node joins/leaves the network, all other nodes need to update their shared key. For example, the work in [10] proposed and analyzed an interworking key management solution between multilayer IPsec and logical key hierarchy (LKH), which belongs to the aforementioned third category (topology-based). The authors showed for LKH how user preregistration and periodic admission reduces this initialization cost, and how the optimum outdegree of a hierarchical tree varies with the expected user volatility and rekey factor. Another example is given in [11], where a secure identity-based authenticated key-exchange (AKE) protocol was proposed, which belongs to the aforementioned second category (distributed). As opposed to the other identity-based AKE protocols, the proposed protocol in [11] requires only two rounds of message sending in one session and three exponentiations per node in each session, achieving a relatively high communication and computation efficiency.

The challenges of key management in space

The nodes in space information network are usually satellites or aircrafts with much higher computational capability than that in terrestrial networks. How to utilize such properties to design key management schemes for a space information network is a worthy future research topic.

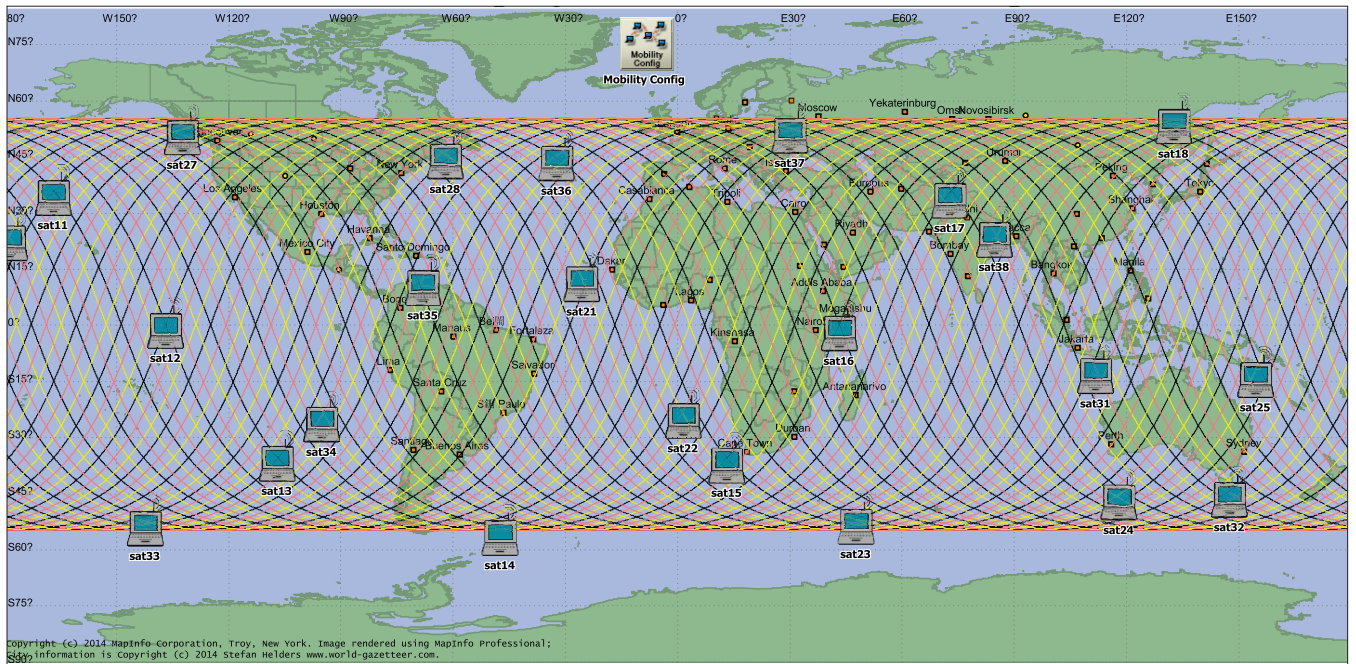


Figure 3. Simulation scenario with Walker 24/3/2 constellation.

information networks lie in several aspects. First, space information network nodes are distributed over an extremely vast open area, making it impossible to build a powerful online key management center, as in centralized schemes. Second, the space network environment is fundamentally different from the terrestrial environment; for example, the channel is easily interfered with, so the BERs of a space channel are relatively high, and communication load and time delay are strictly restricted when designing security protocols. In other words, reducing the round complexity and improving rekeying efficiency are important targets in designing a key management scheme for space information networks. Third, the high dynamics of space information network nodes may lead to frequent joining/leaving phenomenon, and the 1-affects- n problem should be emphasized. Meanwhile, one special characteristic feature of space network dynamics is that the orbits are predefined, and the movement trend can be accurately predicted. Such regularity may help design customized key management schemes, that is, using the topology rules as in the topology-based schemes. Finally, the nodes in space information network are usually satellites or aircrafts with much higher computational capability than that in terrestrial networks. How to utilize such properties to design key management schemes for a space information network is a worthy future research topic.

SECURE ROUTING IN SPACE INFORMATION NETWORKS

In satellite networks, any two satellites within line of sight are connected to each other via ISLs. There are two kinds of ISLs: the links between neighboring satellites in the same orbital plane, intra-plane ISLs, and the links

between neighboring satellites in different orbital planes, inter-plane ISLs. This raises the issue of routing in satellite networks. When it comes to the broader concept of space networks, different network entities in different layers have to rely on a trustworthy routing strategy in order to be connected. The objectives of routing strategy design fall into several aspects as follows:

- Guarantee an acceptable block probability due to link handovers.
- Minimize the length of the paths and the number of hops to decrease propagation delay and end-to-end delay.
- Do load balancing (i.e., prevent congestion of some links while others are idle).
- Perform efficient multicast over space information networks.
- Ensure robustness pertaining to different attacks.

In the literature, satellite routing protocols have been intensively studied, while space information network routing has not been well investigated, let alone the related security issues. The authors in [12] conducted a thorough survey on the routing strategies in satellite networks, with an emphasis on the routing issues in a dynamic topology, reducing link handovers and rerouting issues, path minimization algorithms, traffic-based routing, routing from a space-ground integration point of view, multicast routing, and load balancing algorithms. Later, a multiservice on-demand routing for LEO satellite networks was proposed in [13]. The authors divided typical real-time network services into three classes: voice, best-effort traffic that has no specific QoS constraints, and bandwidth-sensitive services such as video. The proposed protocol adjusts the routing procedure to the QoS requirements of different traffic classes in two steps: route request area formation and path discovery. Dynamic routing with security considerations was proposed in [14]. Different from most of

the existing works focusing on cryptography algorithms and secure system infrastructures, the authors randomized delivery paths for data transmission to achieve security. The underlying idea is to minimize the path similarity without introducing any extra control messages, and thus reduce the probability of eavesdropping consecutive packets over a specific link. A recent work in [15] designed a so-called intelligent routing strategy based on traffic lights (the congestion status). The idea supporting this work is that a packet may adjust the route dynamically according to the real-time color of traffic lights at each intermediate node, that is, a combination of preliminary planning and real-time adjustment.

As aforementioned, little effort has been made regarding routing in space information networks and its security considerations. Tackling this problem, we propose a secure routing protocol based on intrusion detection for space information networks. As a proactive security protection technology, intrusion detection is not only able to take control over illegal system attackers, but can also collect abnormal behavior information such as vulnerability holes and DoS attacks. A secure space network routing protocol works mainly against network attack behaviors, including false distance vector attacks, false destination sequence attacks, routing table overflow attacks, nodes' selfish behavior, and so on. In order to deal with these security threats, we design a space network security mechanism as follows. Detect the malicious attacks using intrusion detection system operating over the routing protocols, then degrade its creditworthiness via a reputation system. A normal node would be labeled as a malicious node and then isolated once its credit is below some threshold. If malicious nodes are detected, activate the intrusion response system to reconstruct the route. Considering the dynamic topology changes of space information networks, it may be difficult to discover and affirm malicious behaviors via a single-node intrusion detection system. In order to enhance the accuracy of attack detection, we propose to use multiple nodes to identify malicious acts, called a joint intrusion detection system. When a single node detects malicious acts but cannot confirm whether it is a malicious node, it would broadcast single intrusion detection messages, which can stimulate the neighboring nodes into a joint intrusion detection phase to make the decision cooperatively.

We simulate the intrusion-detection-based Optimized Link State Routing (OLSR) protocol via the Satellite Tool Kit (STK) and OPNET. As shown in Fig. 3, the simulation scenario consists of a Walker 24/3/2 constellation with 24 satellites distributed on a circular orbit at an altitude of 11,096 km in three planes inclined at 55°, where 2 defines the relative spacing between the satellites in adjacent planes. The intrusion satellites are assumed to be able to falsify one-hop neighbor node information in the Hello packet. The top of Fig. 4 shows that the network throughput with the proposed intrusion-detection-based routing scheme is apparently higher than that without the scheme, which verifies the effective-

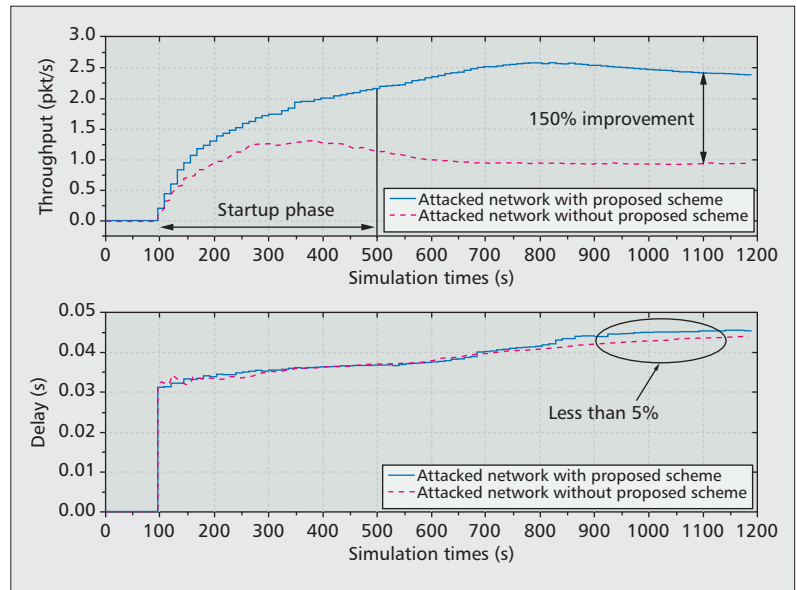


Figure 4. Performance comparison between attacked networks with and without the proposed scheme.

ness of the proposed scheme. Meanwhile, we also investigate the impact on the network delay performance, as shown in the bottom of Fig. 4. It can be seen that the delay increases slightly when the proposed scheme is adopted, which is negligible compared to the remarkable throughput enhancement.

CONCLUSION

In this article, the network security issues in space information networks are investigated. We study the secure handoff and transmission control issues in space information networks via surveying the existing work and open problems. Meanwhile, the key management issues in space information networks are also investigated in four different categories, including centralized, distributed, topology-based, and preconfigured. In addition, secure routing issues are introduced and our proposed scheme is presented. In summary, we expect that this article gives an overview of the fundamental issues and key techniques with regard to the security problems in space information networks.

ACKNOWLEDGMENT

This research was supported by NSFC China under projects 61371079, 61273214, 61271267 and 91338203, and by a Postdoctoral Science Foundation funded project.

REFERENCES

- [1] J. Mukherjee and B. Ramamurthy, "Communication Technologies and Architectures for Space Network and Interplanetary Internet," *IEEE Commun. Surveys and Tutorials*, vol. 15, no. 2, 2013, pp. 881–97.
- [2] L. Giupponi and C. Ibars, "Enabling Vertical Handover Decisions in Heterogeneous Wireless Networks: A State-of-the-Art and A Classification," *IEEE Commun. Surveys and Tutorials*, vol. 16, no. 2, 2014, pp. 776–811.
- [3] P.K. Chowdhury, M. Atiquzzaman, and W. Ivancic, "Handover Schemes in Satellite Networks: State-of-the-Art and Future Research Directions," *IEEE Commun. Surveys and Tutorials*, vol. 8, no. 4, 2006, pp. 2–14.

- [4] G. Karopoulos, G. Kambourakis, and S. Gritzalis, "Survey of Secure Handoff Optimization Schemes for Multimedia Services Over All-IP Wireless Heterogeneous Networks," *IEEE Commun. Surveys and Tutorials*, vol. 9, no. 3, 2007, pp. 18–28.
- [5] D. He, C. Chen, S. Chan, and J. Bu, "Secure and Efficient Handover Authentication Based on Bilinear Pairing Functions," *IEEE Trans. Wireless Commun.*, vol. 11, no. 1, 2012, pp. 48–53.
- [6] R. Wang et al., "Protocols for Reliable Data Transport in Space Internet," *IEEE Commun. Surveys and Tutorials*, vol. 11, no. 2, 2009, pp. 21–32.
- [7] Y. Sun, Z. Ji, and H. Wang, "TFRC-Satellite: A TFRC Variant with a Loss Differentiation Algorithm for Satellite Networks," *IEEE Trans. Aerospace Electronics Sys.*, vol. 49, no. 2, 2013, pp. 716–24.
- [8] C. Peng, Q. Zhang, and C. Tang, "Improved TLS Handshake Protocols Using Identitybased Cryptography," *IEEE IECC*, 2009, pp. 135–39.
- [9] A. Parichehreh and B. Eliasi, "VPN over Satellite: Performance Improving of E2E Secured TCP Flows," *IFIP WOCN*, 2008, pp. 1–4.
- [10] M. P. Howarth et al., "Dynamics of Key Management in Secure Satellite Multicast," *IEEE JSAC*, vol. 22, no. 2, 2004, pp. 308–19.
- [11] Y. Zhong and J. Ma, "A Highly Secure Identity-Based Authenticated Key-Exchange Protocol for Satellite Communication," *J. Commun. Net.*, vol. 12, no. 6, 2010, pp. 592–99.
- [12] F. Alagoz, O. Korcak, and A. Jamalipour, "Exploring the Routing Strategies in Next-Generation Satellite Networks," *IEEE Wireless Commun.*, vol. 14, no. 3, 2007, pp. 79–88.
- [13] S. Karapantazis, E. Papapetrou, and F.-N. Pavlidou, "Multiservice On-Demand Routing in LEO Satellite Networks," *IEEE Trans. Wireless Commun.*, vol. 8, no. 1, 2009, pp. 107–12.
- [14] C.-F. Kuo, A.-C. Pang, and S.-K. Chan, "Dynamic Routing with Security Considerations," *IEEE Trans. Parallel Distrib. Sys.*, vol. 20, no. 1, 2009, pp. 48–58.
- [15] G. Song et al., "TLR: A Traffic-Light-Based Intelligent Routing Strategy for N GEO Satellite IP Networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 6, 2014, pp. 3380–93.

BIOGRAPHIES

CHUNXIAO JIANG [S'09, M'13] (chx.jiang@gmail.com) received a B.S. (Hons.) degree in information engineering from Beihang University, Beijing, China, in 2008, and a Ph.D. (Hons.) degree from Tsinghua University (THU), Beijing, in 2013. From 2011 to 2013, he visited the Signals and Information Group at the Department of Electrical and Computer Engineering, University of Maryland at College Park, with Prof. K. J. Ray Liu. He currently holds a postdoctoral position with the Department of Electrical Engineering, THU, with Prof. Y. Ren. His research interests include the applications of game theory and queuing theory in wireless communication and networking, and social networks. He was a recipient of the Best Paper Award from IEEE GLOBECOM in 2013, the Beijing Distinguished Graduated Student Award, the Chinese National Fellow-

ship, and the Tsinghua Outstanding Distinguished Doctoral Dissertation in 2013.

XUEXIA WANG (wangxuexia@tsinghua.edu.cn) received her B.S. M.S. and Ph.D. degrees in electronic engineering from Harbin Institute of Technology, China, in 2002, 2004, and 2008, respectively. She is currently working as a postdoctoral with the Department of Electrical Engineering, THU, with Prof. Y. Ren. Her research interests include space-based information networks and secure routing protocol design.

JIAN WANG (jian-wang@tsinghua.edu.cn) received a Ph.D. degree in electronic engineering from THU in 2006. In 2006, he joined the faculty of THU, where he is currently an associate professor with the Department of Electronic Engineering. His research interests are in the areas of information security, signal processing in the encrypted domain, and cognitive networks.

HSIAO-HWA CHEN [S'89, M'91, SM'00, F'10] (hshwchen@ieee.org) is currently a Distinguished Professor in the Department of Engineering Science, National Cheng Kung University, Taiwan. He obtained his B.Sc. and M.Sc. degrees from Zhejiang University, China, and a Ph.D. degree from the University of Oulu, Finland, in 1982, 1985, and 1991, respectively. He has authored or co-authored over 400 technical papers in major international journals and conferences, six books, and more than 10 book chapters in the areas of communications. He has served as general Chair, TPC Chair, and Symposium Chair for many international conferences. He has served or is serving as an Editor or/and Guest Editor for numerous technical journals. He is the founding Editor-in-Chief of Wiley's *Security and Communication Networks Journal* (www.interscience.wiley.com/journal/security). He was the recipient of the best paper award at IEEE WCNC 2008 and a recipient of the IEEE Radio Communications Committee Outstanding Service Award in 2008. Currently, he is also serving as Editor-in-Chief for *IEEE Wireless Communications*. He is a Fellow of IET and an elected Member at Large of IEEE ComSoc.

YONG REN (reny@tsinghua.edu.cn) received his B.S. M.S., and Ph.D. degrees in electronic engineering from Harbin Institute of Technology in 1984, 1987, and 1994, respectively. He worked as a postdoctoral researcher at the Department of Electronics Engineering, THU, from 1995 to 1997. He is a professor of the Department of Electronics Engineering and the director of the Complexity Engineered Systems Lab of THU. He holds 12 patents, and has authored or co-authored more than 100 technical papers on the behavior of computer networks, P2P networks, and cognitive networks. He has served as a reviewer of *IEICE Transactions on Communications*, *Digital Signal Processing*, *Chinese Physics Letters*, *Chinese Journal of Electronics*, *Chinese Journal of Computer Science & Technology*, and *Chinese Journal of Aeronautics*, among others. His current research interests include complex systems theory and its applications to the optimization and information sharing of the Internet, space-based information networks, Internet of Things and ubiquitous networks, cognitive networks, and cyber-physical systems.

AUGUST 2015



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:

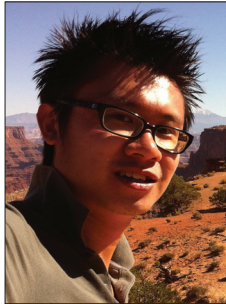


Limited Time Only at >> ww.comsoc.org/freetutorials



For this and other sponsor opportunities,
please contact Mark David // 732-465-6473 // m.david@ieee.org

ENERGY HARVESTING COMMUNICATIONS: PART III



Chau Yuen



Maged Elkashlan



Yi Qian



Trung Q. Duong



Lei Shu



Frank Schmidt

Over the last decade, energy harvesting has emerged as a promising approach to enable self-sufficient and self-sustaining operation for devices in energy-constrained networks by scavenging energy from the ambient environment to power up devices.

In particular, for wireless sensor networks, if the sensors, which spread throughout homes or factories, in buildings or even outdoors, are powered by energy harvesting, there are no batteries to replace and no labor costs associated with replacing them. For a cellular network, energy harvesting can be used to provide power in many elements of a telecom network, saving considerable costs in electricity supply, and providing low-maintenance monitoring. As another important focus, RF energy is currently broadcast from billions of radio transmitters around the world. The ability to harvest RF energy, from ambient or dedicated sources, enables wireless charging of low-power devices and has significant benefits to product design, usability, and reliability.

This Feature Topic focuses on energy harvesting related issues in communications, through presenting a holistic view of research challenges and opportunities in the emerging area of energy harvesting communications. We hope this Feature Topic is able to help readers to obtain better understanding of some key issues in energy harvesting, and drive more research interest.

This is Part III of this Feature Topic, starting with “A Hierarchical Packet Forwarding Mechanism for an Energy Harvesting Wireless Sensor Network” by Dapeng Wu *et al.*, which investigates a dynamic gradient-aware hierarchical packet forwarding mechanism based on the relative positions of the sensor nodes.

The article “Wireless Information and Power Transfer: From Scientific Hypothesis to Engineering Practice,” by Rong Zhang *et al.*, surveys the state-of-the-art findings on simultaneous wireless information and power transfer (SWIPT), including design challenges, potential solutions, and research ideas.

The article “Delay-Sensitive Dynamic Resource Control for Energy Harvesting Wireless Systems with Finite Energy Storage,” written by Fan Zhang *et al.*, studies optimization

of the delay performance of energy harvesting wireless systems with finite energy storage, covering the resource allocation to the channel fading information, data queue length, and energy queue.

The article “Toward Secure Energy Harvesting Cooperative Networks,” by Jiawen Kang *et al.*, addresses the security issues and solutions in energy harvesting cooperative networks, where the vulnerability of energy cooperation models in energy attacks is analyzed, and several energy defense solutions are proposed.

BIOGRAPHIES

CHAU YUEN (yuenchau@sutd.edu.sg) received his B. Eng and Ph.D. degrees from Nanyang Technological University, Singapore, in 2000 and 2004, respectively. He was a postdoctoral fellow at Lucent Technologies Bell Labs, Murray Hill, New Jersey, during 2005. He was a visiting assistant professor at Hong Kong Polytechnic University in 2008. During the period of 2006–2010, he worked at the Institute for Infocomm Research, Singapore, as a senior research engineer. He joined Singapore University of Technology and Design as an assistant professor in June 2010. He serves as an Associate Editor for *IEEE Transactions on Vehicular Technology* and was awarded Top Associate Editor for three consecutive years. In 2012, he received the IEEE Asia-Pacific Outstanding Young Researcher Award. He has held positions on several conference organizing committees, and is on Technical Program Committees of various international conferences.

MAGED ELKASHLAN received his Ph.D. degree in electrical engineering from the University of British Columbia, Canada, in 2006. From 2006 to 2007, he was with the Laboratory for Advanced Networking at the University of British Columbia. From 2007 to 2011, he was with the Wireless and Networking Technologies Laboratory at the Commonwealth Scientific and Industrial Research Organization (CSIRO), Australia. He also held an adjunct appointment at the University of Technology Sydney, Australia, between 2008 and 2011. In 2011, he joined the School of Electronic Engineering and Computer Science at Queen Mary, University of London, United Kingdom, as an assistant professor. His research interests include millimeter wave communications, energy harvesting, cognitive radio, and wireless security. He currently serves as an Editor for *IEEE Transactions on Wireless Communications*, *IEEE Transactions on Vehicular Technology*, and *IEEE Communications Letters*. He received Best Paper awards at IEEE ICC '14, International Conference on Communications and Networking in China in 2014, and IEEE VTC-Spring 2013. He received the Exemplary Reviewer Certificate of *IEEE Communications Letters* in 2012.

YI QIAN [M'95, SM'07] is an associate professor in the Department of Electrical and Computer Engineering, University of Nebraska-Lincoln (UNL). Prior to joining UNL, he worked in the telecommunications industry, academia, and the government. Some of his previous professional positions include serving as a senior member of scientific staff and technical advisor at Nortel Networks, a senior systems engineer and technical advisor at sev-

eral startup companies, an assistant professor at the University of Puerto Rico at Mayaguez, and a senior researcher at the National Institute of Standards and Technology. His research interests include information assurance and network security, network design, network modeling, simulation and performance analysis for next generation wireless networks, wireless ad hoc and sensor networks, vehicular networks, smart grid communication networks, broadband satellite networks, optical networks, high-speed networks, and the Internet. He has a successful track record in leading research teams and publishing research results in leading scientific journals and conferences. Several of his recent journal articles on wireless network design and wireless network security are among the most accessed papers in the IEEE Digital Library. He is the current Chair of the Communications and Information Security Technical Committee in the IEEE Communications Society. He is an IEEE Distinguished Lecturer.

TRUNG Q. DUONG received his Ph.D. degree in telecommunications systems from Blekinge Institute of Technology (BTH), Sweden, in 2012, and then continued working at BTH as a project manager. Since 2013, he has joined Queen's University Belfast, United Kingdom, as a lecturer (assistant professor). He held visiting positions at Polytechnic Institute of New York University and Singapore University of Technology and Design in 2009 and 2011, respectively. His current research interests include cooperative communications, cognitive radio networks, green communications, physical layer security, massive MIMO, cross-layer design, mmWave communications, and localization for radios and networks. He has been a TPC chair for several IEEE international conferences and workshops, including most recently the IEEE GLOBECOM '13 Workshop on Trusted Communications with Physical Layer Security. He currently serves as an Editor for *IEEE Communications Letters* and *Wiley Transactions on Emerging Telecommunications Technologies*. He served as Lead Guest Editor of the Special Issue on Location Awareness for Radios and Networks of the *IEEE Journal on Selected Areas in Communications*, Lead Guest Editor of the Special Issue on Secure Physical Layer Communications of *IET Communications*, Guest Editor of the Special Issue on Green Media: Toward Bringing the Gap between Wireless and Visual Networks of *IEEE Wireless Communications*, Guest Editor of the Special Issue on Millimeter Wave Communications for 5G of *IEEE Communications Magazine*, Guest Editor of the Special Issue on Cooperative Cognitive Networks of the *EURASIP Journal on Wireless Communications and Networking*, and Guest Editor of the Special Issue on Security Challenges and Issues in Cognitive Radio Networks of the *EURASIP Journal on Advances Signal Processing*. He was awarded the Best Paper Award at IEEE VTC-Spring '13 and the Exemplary Reviewer Certificate of *IEEE Communications Letters* in 2012.

LEI SHU [M] received his B.Sc. degree in computer science from South Central University for Nationalities, China, in 2002, his M.Sc. degree in computer engineering from Kyung Hee University, Korea, in 2005, and his Ph.D. degree from the Digital Enterprise Research Institute, National University of

Ireland, Galway, Ireland, in 2010. Until March 2012, he was a specially assigned researcher in the Department of Multimedia Engineering, Graduate School of Information Science and Technology, Osaka University, Japan. He is a member of IEEE IES, IEEE ComSoc, EAI, and ACM. In October 2012, he joined Guangdong University of Petrochemical Technology, China, as a full professor. In 2013, he started to serve Dalian University of Technology as a Ph.D. supervisor in the College of Software, Beijing University of Posts and Telecommunications as a Master's supervisor in information and communication engineering, Wuhan University as a Master's supervisor in the College of Computer Science, guest professor at Tianjin University of Science and Technology, and a guest researcher at Guangzhou Institute of Advanced Technology, Chinese Academy of Sciences. Meanwhile, he is also working as vice-director of the Guangdong Provincial Key Laboratory of Petrochemical Equipment Fault Diagnosis, China. He is the founder of the Industrial Security and Wireless Sensor Networks Lab. His research interests include wireless sensor networks, multimedia communication, middleware, fault diagnosis, and security. He has published over 230 papers in related conferences, journals, and books in the area of sensor networks. Currently, his H-index is 21 in Google Citation. Total citations of his papers by other people are more than 1600. He developed an open source wireless sensor networks simulator, NetTopo, to evaluate and demonstrate algorithms. NetTopo has been downloaded more than 3420 times over the past three years, and is widely used by international researchers and students. He was awarded the MASS 2009 IEEE TCs Travel Grant and the Outstanding Leadership Award of EUC 2009 as Publicity Chair, GLOBECOM 2010, and ICC 2013, the ComManTel 2014 Best Paper Award, and the Outstanding Service Award of IUCC 2012 and ComcomAP 2014. He also received a few more awards from the Chinese government: Top Level Talents in "Sailing Plan" of Guangdong Province, China, and Outstanding Young Professor of Guangdong Province, China. He has been serving as Editor-in-Chief for *EAI Endorsed Transactions on Industrial Networks and Intelligent Systems*, and Associate Editor for *IEEE Access*, *ACM/Springer Wireless Networks*, *Journal of Network and Computer Applications*, *Transactions on Emerging Telecommunications Technology*, and several other publications. He has served as Co-Chair for many international conferences. He has obtained more than 4 million RMB in research grants since October 2012.

FRANK SCHMIDT is a pioneer in energy harvesting and the visionary in the management team of EnOcean. As chief technology officer he is responsible for the overall technical orientation, patent related activities, as well as the relationship management with educational, research and scientific organizations. Before joining EnOcean he was at the Central Research Department of Siemens AG, where he created self-powered wireless sensor technology as early as 1995. He has been granted more than 40 patents for his energy harvesting inventions and is the author of numerous technical publications in this field. He is a physicist and studied at the Technical University of Chemnitz, Germany.

A Hierarchical Packet Forwarding Mechanism for Energy Harvesting Wireless Sensor Networks

Dapeng Wu, Jing He, Honggang Wang, Chonggang Wang, and Ruyan Wang

ABSTRACT

Energy harvesting technologies have gained widespread attention for their perpetual energy supply for sensor nodes. However, the energy resources are still insufficient while the harvesting module is added on the node. To prolong the network lifetime and meet the demand of a green wireless network, a dynamic gradient-aware hierarchical packet forwarding mechanism is designed. According to the relative positions of nodes, gradient-aware clusters are established. Consequently, considering the energy conversion efficiency and relative distance, cluster heads are selected reasonably. Furthermore, by exploiting the available energy and the number of cluster members, packets can be forwarded to the sink in an energy-efficient manner. Results show that the network lifetime can be noticeably improved.

INTRODUCTION

The wide deployment of wireless sensor networks (WSNs) enables a great variety of remarkable applications including remote environment monitoring, healthcare, air quality monitoring, and so on. However, the performance of WSNs is restrained by limited node processing capacity and battery capacity. Currently, several approaches have been introduced to reduce the energy consumption and correspondingly maximize the lifetime of WSNs, such as data aggregation, green routing, and sleep scheduling mechanisms. Among them, energy harvesting (EH) technology [1] is very promising due to the unlimited energy supply provided by power sources such as solar power and thermal power. By adding an EH module to the existing node structure, a novel EH-WSN was proposed in [2], which substantially changed the design of the network by eliminating the major constraint: the battery capacity.

EH sensor nodes can effectively exploit ambient energy and convert it into electricity to power the sensors, which potentially prolongs the network lifetime. Although ambient energy can be either directly consumed or stored in a rechargeable battery for later use, the availability

of ambient energy is still limited by different environmental factors, such as light intensity for solar power and temperature for thermal power. Therefore, the harvested energy is constantly changing and relatively limited in a given period, which makes the trade-off between the harvested energy and consumed energy the major factor affecting the application performance. To effectively manage the harvested energy, there are still several unsolved problems, including packet forwarding mechanisms customized for EH-WSNs, which optimize the energy consumption and maximize the network lifetime.

Existing energy-efficient packet forwarding strategies can be classified into three categories [3]: data-centric forwarding, location-aware forwarding, and hierarchical forwarding. Hierarchical forwarding establishes a clustered topology to balance the energy consumption and the load of nodes; at the same time, scalability can be guaranteed due to the simplified topology maintenance. However, cluster heads (CHs) should be appropriately selected to be uniformly distributed in the sensing field to guarantee coverage with minimized cluster establishment and maintenance overheads. Furthermore, the clusters are established where CHs can help aggregate the packets within their clusters and then forward the aggregated packets according to given routing rules. In practice, the energy consumption of the backbone network consisting of CHs is imbalanced, because CHs closer to the sink will spend more energy on relaying packets from farther clusters. To address this so-called “hot spot” problem and avoid the network partitioning caused by it, the authors in [4] provided a clustering method to adaptively control the sizes of clusters according to their distances from the sink. The probability of nodes becoming CHs is determined according to their positions. However, this clustering method requires the lifetimes of all sensor nodes in the CH selection procedure, and massive resources will be consumed by exchanging the status information in a distributed WSN. Reference [5] proposed an energy-efficient hierarchical routing algorithm that utilizes a special packet header to dynamically update the node residual energy. When the

Dapeng Wu, Jing He, and Ruyan Wang are with Chongqing University of Posts and Telecommunications.

Honggang Wang is with the University of Massachusetts Dartmouth.

Chonggang Wang is with Interdigital Communications, USA.

residual energy of all nodes is smaller than the threshold, the clusters are forced to re-establish. Two sleep scheduling mechanisms for geographic routing in WSNs were proposed in [6] to effectively transmit data on paths with shorter average distances, which dynamically adjusts the sleep schedule according to the requirements of the geographic routing and fully considers the mobility of nodes and duty cycle. Reference [7] provided a holistic overview of EH communication systems and proposed a mathematical model that enables perpetual applications in remote and unattended areas. On the other hand, the proposed mathematical model is customizable according to given applications, where the EH rate changes rapidly under different scenarios, such as light intensity and sunshine duration. By estimating the future energy consumption trend based on the monitored energy harvesting status, a real-time adaptive energy management strategy is proposed in [8] to observably improve the channel throughput and guarantee reliable data transmission. Besides, an energy storage method is proposed to mainly use the harvested energy instead of the stored energy and thus reduce battery loss. An energy-loss-aware routing scheme for EH-WSNs is proposed in [9] to achieve a relatively high degree of the global residual energy by considering the energy loss in battery charging and minimizing the energy loss caused by battery overcharge. The authors in [10] proposed an adaptive energy-aware clustering packet forwarding mechanism for EH-WSNs. The node energy state is taken into account in the CH selection procedure, and the cluster establishment can be adjusted according to the network deployment environment. But it may have a long transmission delay when the available time of solar energy is very short, because each node harvests energy from the environment without a backup chemical battery. By employing EH nodes as the dedicated relays for CHs, a cluster-based packet forwarding mechanism for EH-WSNs is introduced in [11] to effectively prolong the network lifetime.

Taking the features of EH sensor nodes into account, a gradient-based energy-efficient clustering (GEEC) packet forwarding mechanism is presented in this article. The proposed mechanism constructs a gradient model according to the hop counts from nodes to the sink, and CHs are selected in a distributed manner according to the distance between a CH candidate and the center line of its circular ring (for simplicity, we use the term ring hereafter). Besides, the node available energy is evaluated based on the EH rate. Then unequal clusters are constructed, and the inter-cluster routing paths are established and updated according to the available energy, the number of non-CH nodes, and the relative node positions, which finally fulfill the objectives of balanced network energy consumption, reduced network overhead, and prolonged network lifetime.

With the proposed dynamic unequal clustering and inter-cluster packet forwarding methods, the collected sensing data can be more effectively uploaded to the sink. Therefore, the overall packet transmission energy consumption is reduced, and the node failures when the harvest-

ed energy is insufficient are avoided, which ensures the connectivity of the network and prolongs the network lifetime.

The main contributions of this article are as follows:

First, a dynamic gradient-aware hierarchical data forwarding mechanism is proposed for EH-WSNs to establish the gradient model according to the relative node positions and further construct the network topology in an unequal clustering manner.

Second, an energy-balanced CH selection method is proposed to comprehensively consider the residual energy level and relative distances, and thus determine the CH competition value. By selecting the node with strong competitiveness as the CH, energy consumption is notably balanced among clusters.

Last, an energy-efficient data forwarding mechanism is proposed according to the designed cost function based on the residual energy, number of cluster members, and relative position relationship between nodes. When entering the stable data transmission stage, low-gradient relays with small cost function values are selected to forward packets, which dramatically reduces the transmission cost and enhances the energy efficiency.

NODE AND NETWORK MODELS

Usually, an EH sensor node consists of five modules: the EH module, energy storage module, processing unit, sensing unit, and wireless transceiver, respectively. The EH module is responsible for converting ambient energy to electricity, and we employ a photovoltaic panel in this article. The super capacitor is adopted as the energy storage module, which stores the harvested energy to later power the processing, sensing units, and wireless transceiver. Other modules serve the same functions as conventional WSN nodes. The node structure is shown in Fig. 1.

The major difference between the EH sensor node and a conventional WSN node is the attachment of an EH module, which substantially prolongs the node lifetime. In a practical application, different energy conversion techniques are demanded for different ambient energy. For instance, the size of the photovoltaic panel is a crucial design problem in terms of the cost efficiency, environmental factors, available time of solar energy; and the capacity of the rechargeable battery is another important design problem to reduce the network cost; at the same time, the available energy of a given node is limited and quantifiable until the next energy harvesting round. Even when engaging the EH module, ambient energy is unavailable during the night due to the lack of light. Therefore, to ensure stable energy consumption and reliable packet transmission, and to avoid network partitioning, the appropriate use of available energy is still important for EH-WSNs.

Generally, a clustered EH-WSN can be modeled as follows. Assuming that the radius of the sensing field is R , the sensing nodes are uniformly distributed in a fixed manner, and the sink is at the center of the sensing field. Nodes within

Even when engaging the EH module, ambient energy is unavailable during the night due to the lack of light. Therefore, to ensure stable energy consumption and reliable packet transmission, and avoid network partitioning, the appropriate use of available energy is still important for EH-WSNs.

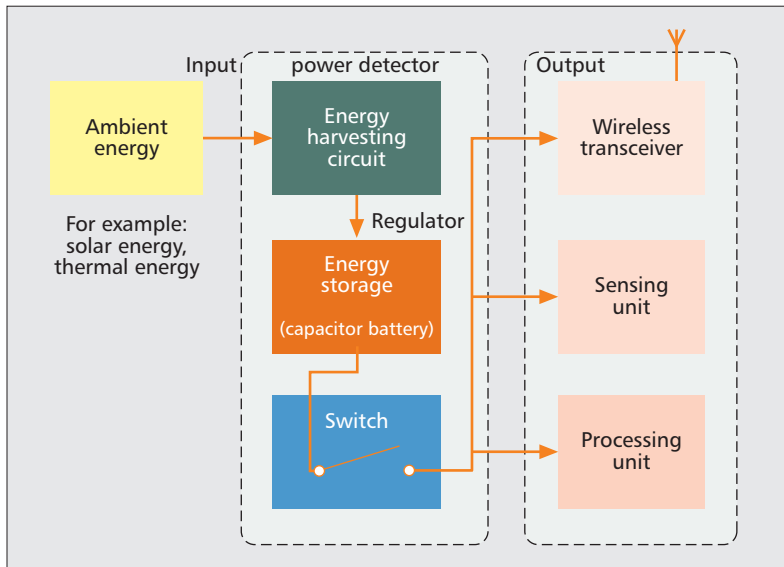


Figure 1. The architecture of EH sensor nodes.

the communication range of the sink can directly communicate with it, whereas other nodes form clusters in a self-organized manner and transmit packets via CHs. The node transmission power is adjustable according to the transmission distance to reduce energy consumption. To enhance the scalability of the proposed mechanism, we assume that nodes are not equipped with Global Positioning System (GPS) modules, so the location information is unknown and estimated by the signal strength. To reduce the computation complexity, we employ the free space channel model (d^2 power loss) and multipath fading channel model (d^4 power loss) [12] to estimate the transmission energy consumption. According to the fundamental theory of data aggregation, packets sensed by a group of neighbor nodes are highly space-correlated, and transmitting all of them causes redundancy. Thus, a CH first aggregates the data sensed by its cluster members, and the energy consumption of the aggregation is defined as E_{DA} .

GEEC MECHANISM

To reasonably exploit the available energy of EH sensor nodes and avoid energy depletion before the next energy harvesting round, an energy-efficient packet forwarding mechanism is designed, which includes four phases: gradient establishment, CH selection, inter-cluster routing construction, and stable packet forwarding. Initially, packets consisting of hop counts are broadcast in the network to establish the gradient, where this parameter reflects the distance from a node to the sink, and it can be utilized to select the energy efficient route. Second, combining the residual energy with the relative distances between nodes, the unequal clusters are constructed based on the established gradient. Furthermore, packet forwarding paths are set up according to the gradients, which enable the shortest transmission path of packets from CHs to the sink. In the stable packet transmission procedure, cluster members directly communicate with their CHs,

and then CHs forward the packets along the established paths, where the relays are selected according to their available energy, positions, and number of cluster members. As a result, the packet forwarding procedure can be achieved in an energy-efficient manner.

GRADIENT ESTABLISHMENT

To reserve the energy resources, the gradient model is established according to the minimum hop count from a given node to sink. First, the minimum hop (MH) to the sink from each node needs to be initialized. To distinguish the sink, its gradient is set to 0, whereas the gradients of other nodes are set to infinity. The gradient initialization is triggered by the sink, and the sink broadcasts the initialization packet to all its neighbor nodes, which contains a hop count (HC) whose value is 0. After a sensor node receives the initialization packet successfully, the node first checks whether its gradient is infinity. Subsequently, the node will update the gradient (HC value plus one) and rebroadcast the updated initialization packet to all the neighbor nodes. The above procedure is repeated until the gradients of all nodes are updated.

After the gradient model is established, all sensor nodes obtain their shortest paths to the sink. The sensing field is divided into K concentric circles with linearly increasing radii. In particular, the sensor nodes within the i th ring have the same gradient. According to the outer and inner radii of the predefined rings, the center lines of a ring can be determined to further calculate the distance from node m to the center line, which can later be adopted to evaluate the node competitiveness to reasonably select CHs.

DYNAMICAL CLUSTER RADIUS CALCULATION

As mentioned above, the availability of ambient energy changes constantly, and the network lifetime should be maximized with the harvested energy. In a hierarchical WSN, there are two major types of packet forwarding: intra-cluster and inter-cluster. Obviously, the energy consumption of intra-cluster forwarding is determined by the packet length and the number of cluster members, whereas the energy consumption of inter-cluster communication varies greatly, because CHs closer to the sink have to forward more packets for other CHs. Although dynamic CH selection can balance the energy consumption of clusters closer to the sink by alternating CHs periodically, it cannot balance the energy consumption of all CHs and between all clusters, which may lead to network partitioning. Therefore, to maximize the network lifetime and ensure balanced energy consumption among different rings, the unequal clustering mechanism is adopted in this article. Thus, the heavy loads on CHs closer to the sink and the network partitioning caused by the energy depletion of these nodes can be avoided. Based on the established gradient model, clusters are constructed and cluster sizes are determined, while the clusters within the same ring are of the same size.

According to the energy consumption model of wireless communications [12], the average energy consumption of CHs in the i th ring can be evaluated. As described before, the energy

consumption of all CHs in the network should be balanced in order to maximize network lifetime. Therefore, the corresponding cluster radius is calculated to achieve the global energy balance by equalizing the energy consumptions among K rings.

According to the gradient model, packets are forwarded from outer rings to inner rings until they arrive at the sink. Due to the uniform distribution of sensor nodes, the average transmission distances are equal for CHs from adjacent rings.

The packets to be forwarded by CHs in the i th ring consist of two parts: the aggregated packet of the intra-cluster packets, and the packets received from CHs of the $(i + 1)$ th ring. Therefore, according to the above factors, the average number of packets to be forwarded per second of CHs in the i th ring can be calculated. According to our unequal clustering mechanism, the radius of the cluster in the i th ring, r_i , can be calculated based on the radius of the outermost ring and ring number i .

GRADIENT PACKET TRANSMISSION

Clusters in different rings are constructed in a distributed manner according to the earlier established gradient, and then relay nodes are selected from CHs to reasonably build forwarding paths. As can be seen, during the CH selection procedure, it is crucial to ensure that the selected CHs cover the whole sensing field. Furthermore, for inter-cluster packet transmission, it is necessary to select the shortest forwarding path that also guarantees the balanced energy consumption, because the inter-cluster packet forwarding consumes the major part of the network energy. The gradient packet transmission consists of three phases: CH selection, inter-cluster routing construction, and stable packet transmission.

Cluster Head Selection — As mentioned above, dynamic CH selection should be designed to balance the loads of each ring and equalize the energy consumption rates of nodes in the same ring. According to the existing EH technologies, EH efficiency is affected by multiple factors, such as light intensity, temperature, weather condition, and node location, and the EH rate is technically a nonlinear random process. Thus, EH rate prediction and available energy estimation are very difficult. However, the prediction of the EH rate in the next round can effectively improve the accuracy of the available energy estimation. While there is little existing research about EH prediction, it still remains an obstacle to the large-scale application of EH sensor nodes. Because the EH rate sequence is highly self-correlated, periodic, and nonlinear, and the EH rate changes constantly due to multiple impact factors, a back propagating neural network [13] is employed in this article to obtain the available energy of a given node as the basis to select relay nodes more reasonably. The back propagating neural network is capable of self-learning and error correction, and has strong approximation capability when handling nonlinear data by obtaining the network weights and structure through training and learning. The

actual EH rates are input into the m -input neurons, and the hidden layer contains $2m + 1$ neurons. Last, the predicted EH rates are exported to the output neurons.

By adopting the above back propagating neural network prediction model, the EH rate of the next round can be predicted according to the historical EH rates. Besides, for the reason that the intra-cluster communication cost is minimized when the CH is fixed at the cluster center [14], the above mentioned available energy should be combined with the distance from the given node to the ring center line to determine the node competitiveness for the CH selection, and the competitiveness of node m can be obtained by weighting the two parameters. Apparently, the value domains of these two parameters are different. Thus, they are normalized by a log function. Therefore, nodes with more available energy and closer to the center line are more competitive and more likely to become CHs.

The detailed CH selection process is as follows.

First, every node becomes a CH candidate with a preset threshold T , which can be adjusted to control the proportion of candidates. Because only CH candidates participate in CH selection, other nodes staying awake in the working mode will inevitably consume energy and thus reduce the network lifetime. To guarantee network connectivity and reduce the energy consumption, a sleep scheduling algorithm is applied to some sensors, where nodes except CH candidates enter the sleeping mode. After the CH selection process is completed, CHs send wakeup messages to those sleeping nodes for cluster formation. Finally, the introduced sleep scheduling algorithm involving the node mobility consideration can flexibly suit the dynamic topology changes, which provides favorable scalability.

Second, according to the gradient model, the cluster radii of the candidates are calculated. Every candidate then broadcasts a contention message (CM), where the broadcast radius R should be restricted and less than the node distance threshold to avoid the d^4 power loss. When candidate s receives a CM from candidate g , candidate g will be added into the neighbor candidate set S_{CH} of candidate s if they are in the competing range of each other and share the same gradient. After S_{CH} is constructed, candidates with no neighbor candidates become CHs. Otherwise, the candidate with the largest competitiveness value wins the competition, and broadcasts a CH selection completed message (CHSCM) to candidates in its S_{CH} . Any candidate receiving the CHSCM will quit the competition and broadcast a quit message (QM) to its S_{CH} , whereas any candidate receiving the QM will remove the sender from its S_{CH} , which ensures that there is only one CH within its competing range.

Last, CHs will wake the sleeping nodes by broadcasting a CH announcement message (CHAM). While the broadcast radius is the cluster radius, the receiving node s will add the sender to its neighbor CH set if they are in the same ring, and then select its CH from the neighbor CH set. Eventually, node s sends a corre-

To reasonably exploit the available energy of EH sensor nodes and avoid energy depletion before the next energy harvesting round, an energy-efficient packet forwarding mechanism is designed, which includes four phases: gradient establishment, CH selection, inter-cluster routing construction, and stable packet forwarding.

sponding message to join the cluster, and thus clusters are constructed.

Based on the pre-established gradient model and calculated dynamic cluster radius, the above-mentioned cluster establishment process describes the unequal clustering method and cluster formation principle. The unequal cluster structure based on the gradient model is shown in Fig. 2.

Inter-Cluster Routing Construction — As can be seen, the original intention of the unequal cluster structure is to balance the heavy inter-cluster communication load, where the communication cost is mainly determined by the

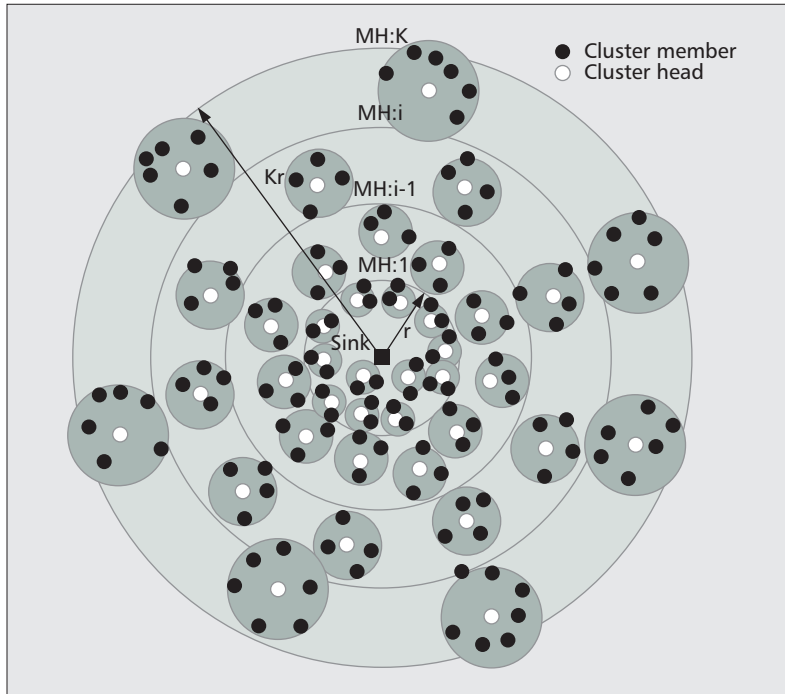


Figure 2. Gradient-aware unequal cluster structure.

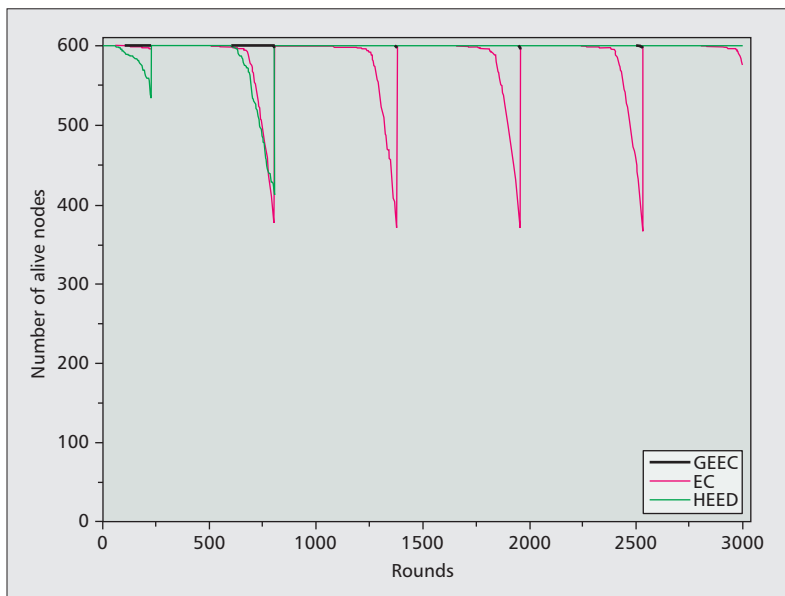


Figure 3. The number of alive nodes vs. operation rounds.

communication radius. Therefore, the single-hop communication radius should be restricted by the node distance threshold to avoid the above mentioned d^4 power loss.

In addition, according to our established network structure, CHs in the first ring directly communicate with the sink, and the forwarding paths for other CHs should be selected to minimize the energy consumption. Based on the assigned gradient for each node, the loop-free path can be established easily. Before the inter-cluster packet forwarding, non-CH nodes enter sleeping mode, and all CHs except the ones in the outermost ring broadcast a node status message (NSM) with radius R . When receiving this message from CH n , CH m will compare their gradient values to determine the next hop relationship.

To further balance the energy consumption of CHs from different rings, the next hop node table is established, which consists of the available energy, MH, distance to sink, cluster member number of the next hop, and distance from CH m to the next hop. Whenever CH m receives an NSM, and MH in this message is smaller, the table is updated by adding the sender to the next hop set until the inter-cluster route structure is completed. Every next hop set includes one or more next hop nodes. In particular, CHs in the first ring only have the sink as their next hop node. As a result, inter-cluster routing can be achieved in a dynamic manner by selecting the next hop node adaptively.

Stable Packet Transmission— All sensor nodes periodically send packets to their CHs, whereas CHs aggregate received packets into a single packet. As can be seen, CHs in the i th ring will forward the packets received from their cluster members and from CHs in the $(i+1)$ th ring to CHs in the $(i-1)$ th ring. By repeating the hop-by-hop transmission, packets can be forwarded to the sink in the direction of descending gradients. Because the gradient of a CH equals its minimum hop count to the sink, minimum hop routing is achieved. When selecting relay nodes from the next hop set, the available energy, relative distance, and number of cluster members are comprehensively considered to guarantee both the energy efficiency and balanced energy consumption. Because relay nodes have to aggregate packets from their cluster members, forwarded packets received from other CHs that consume more energy resources CHs with more available energy should be preferentially selected as relay nodes. On the other hand, with fewer cluster members, the intra-cluster communication traffic is lighter. Therefore, CHs with fewer members should be selected as relay nodes to reserve more energy for the inter-cluster forwarding procedure. Furthermore, the shortest path from a CH to the sink should be the straight line connecting them, so the relay nodes are expected to be near this straight line to minimize the transmission distance and reduce the transmission cost. In the procedure of stable inter-cluster packet transmission, we define a cost function by weighting the above three factors. In each packet forwarding round, CH m will choose the node with the smallest cost value as the relay node, and packets will eventually be

forwarded to the sink by repeating this procedure. In terms of the weighting coefficients, the transmission distance greatly affects the transmission cost, whereas fewer cluster members lead to more available energy for inter-cluster communication.

NUMERICAL RESULTS

In this section, Matlab is adopted to validate the proposed packet forwarding mechanism in EH-WSNs, and its performance is compared to HEED [15] and EC [4] mechanisms. Eventually, the performance metrics are the network lifetime, total number of successfully delivered packets, and packet energy consumption, where the network lifetime is defined as the lifetime of the first energy-exhausted node, and the graphic results also show the numerical trend of alive nodes.

The sensing field size is 500 m × 500 m, where 600 sensor nodes are uniformly deployed, and the sink is at the center of the sensing field. The initial energy of nodes is set to 0.5 J, and the simulation duration is calculated by rounds. Other parameters are set as follows: the energy consumption per bit of the transmitting circuit is 50 nJ/bit, the power amplifier parameters for free space and multipath fading channels are 10 pJ/bit/m² and 0.0013 pJ/bit/m⁴, respectively, the energy consumption of the data aggregation is 5 nJ/bit/signal, the communication radius of the sink is 40 m, the total ring number is 5, the probability of becoming a CH candidate is 0.5, the data packet length is 4000 bits, and the control packet length is 200 bits.

NETWORK LIFETIME

First, the energy efficiency of three algorithms is verified according to the above mentioned network lifetime, where the numerical trend of alive nodes is shown in Fig. 3. The network lifetimes of GEEC, EC, and HEED are 802, 104, and 60 rounds, respectively. Obviously, the trend of GEEC is stable when compared to the other two mechanisms, and the number of alive nodes is the largest. Due to the combination of unequal clustering and energy-efficient inter-cluster multi-hop routing, GEEC effectively balances the energy consumption of clusters far or close to the sink, which optimizes the energy consumption of inter-cluster communication and the number of alive nodes. As shown in Fig. 3, the graphic trend of the alive node number is zigzag due to the changing illumination intensity. When the illumination intensity is weak or in darkness, the available energy for nodes is limited and exhaustible, and nodes fail if their batteries run out.

TOTAL NUMBER OF SUCCESSFULLY DELIVERED PACKETS

The total numbers of successfully delivered packets under different node numbers are depicted in Fig. 4, where the number of operation rounds is 3000. Apparently, this is because the CHs selected by GEEC provide better coverage and can balance the energy consumption, which allows nodes continuous energy supply before the arrival of the next EH round. Because of the corresponding continuous packet trans-

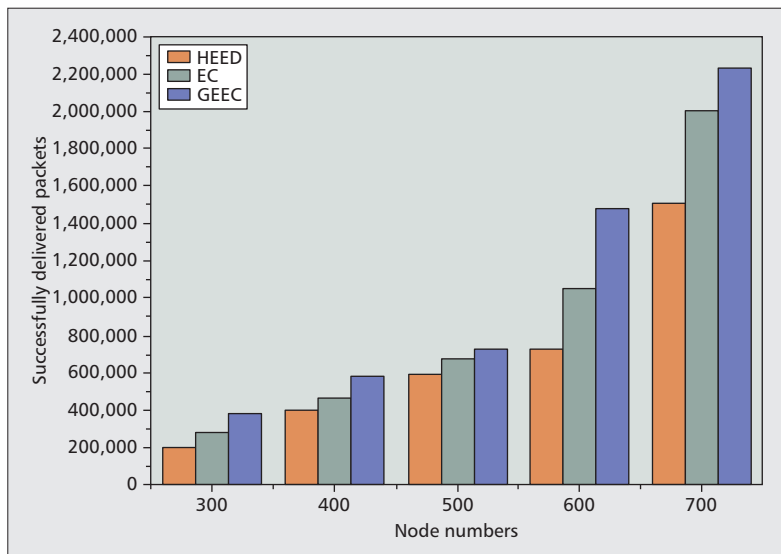


Figure 4. The successfully delivered packets in 3000 rounds vs. node numbers.

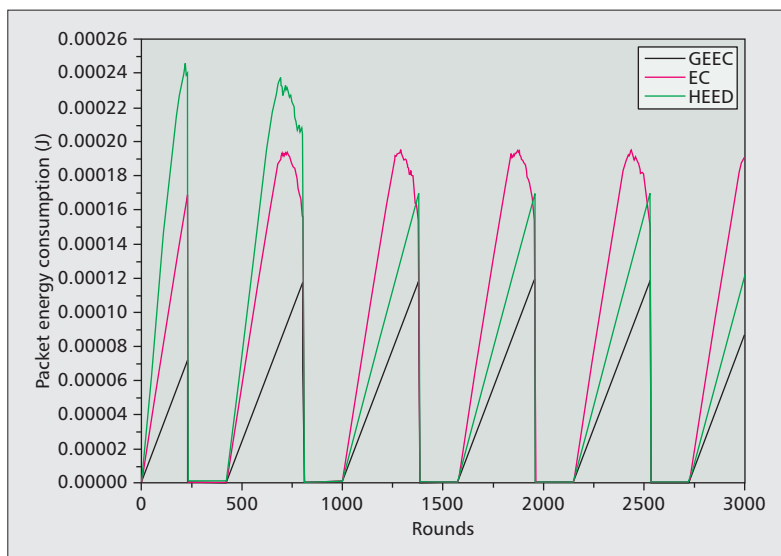


Figure 5. Packet energy consumption vs. operation rounds.

missions, more packets are successfully delivered by GEEC.

PACKET ENERGY CONSUMPTION

The energy consumption for a single packet is depicted in Fig. 5. As can be seen, compared to EC and HEED, GEEC is more energy-efficient due to lower energy consumption per packet, and the packet energy consumption for GEEC is slowly increasing because GEEC solves the “hot-spot” problem.

CONCLUSION

Energy harvesting technology has greatly boosted the development of WSNs by providing theoretically unlimited energy supplies for longer network lifetime. However, the harvested energy is restricted by the device cost and strongly dependent on parameters such as light intensity and sunshine duration in many practical applica-

Due to the development of mobile sensor nodes in practice, further research is planned on scalability, applicability, and efficiency of data forwarding and sleep scheduling algorithms, by which EH technology can be introduced to WSNs with mobile sensors, further optimizing energy efficiency and network lifetime.

tions, which motivate us in designing an energy-efficient packet forwarding mechanism. The proposed GEEC designs the intra-cluster and inter-cluster routing based on the harvested energy to improve the performance of WSNs. The effectiveness of the proposed approach has been proved through our simulations.

In future work, we plan to thoroughly study the EH rate prediction technology for a more precise node capability estimation. Due to the development of mobile sensor nodes in practice, further research is planned on scalability, applicability, and efficiency of data forwarding and sleep scheduling algorithms, by which EH technology can be introduced to WSNs with mobile sensors, further optimizing energy efficiency and network lifetime.

ACKNOWLEDGMENTS

This work is partially supported by the National Natural Science Foundation of China (61371097), Chongqing Natural Science Foundation (Grant No. CSTC2013JJB40006), and Youth Talents Training Project of Chongqing Science & Technology Commission (cstc2014kjrc-qnrc40001).

REFERENCES

- [1] C. Ho and R. Zhang, "Optimal Energy Allocation for Wireless Communications with Energy Harvesting Constraints," *IEEE Trans. Signal Processing*, vol. 60, no. 9, May 2012, pp. 4808–18.
- [2] R. Vullers *et al.*, "Energy Harvesting for Autonomous Wireless Sensor Networks," *IEEE Solid-State Circuits Mag.*, vol. 2, no. 2, Spring 2010, pp. 29–38.
- [3] Y. Qu *et al.*, "Towards a Practical Energy Conservation Mechanism with Assistance of Resourceful Mules," *IEEE Internet of Things J.*, Nov 2014.
- [4] D. Wei *et al.*, "An Energy-Efficient Clustering Solution for Wireless Sensor Networks," *IEEE Trans. Wireless Commun.*, vol. 10, no. 11, Sept. 2011, pp. 3973–83.
- [5] T. Du *et al.*, "An Energy Efficiency Semi-Static Routing Algorithm for WSNs Based on HAC Clustering Method," *Info. Fusion*, vol. 21, Jan. 2015, pp. 18–29.
- [6] C. Zhu *et al.*, "Sleep Scheduling for Geographic Routing in Duty-Cycled Mobile Sensor Networks," *IEEE Trans. Industrial Electronics*, vol. 61, no. 11, 2014, pp. 6346–55.
- [7] D. Gunduz *et al.*, "Designing Intelligent Energy Harvesting Communication Systems," *IEEE Commun. Mag.*, vol. 52, no. 1, Mar. 2014, pp. 210–16.
- [8] S. Peng and P. Chor, "Throughput Optimal Energy Neutral Management for Energy Harvesting Wireless Sensor Networks," *IEEE WCNC*, Shanghai China, April 2012, pp. 2347–51.
- [9] G. Martinez, S. Li, and C. Zhou, "Wastage-Aware Routing in Energy-Harvesting Wireless Sensor Networks," *IEEE Sensors J.*, vol. 14, no. 9, April 2014, pp. 2967–74.
- [10] J. Meng, X. Zhang and Y. Dong, *et al.* "Adaptive Energy-Harvesting Aware Clustering Routing Protocol for Wireless Sensor Networks," *IEEE CHINACOM*, Kunming, China, Aug 2012, pp. 742–47.
- [11] P. Zhang, G. Xiao, and H. Tan, "Clustering Algorithms for Maximizing the Lifetime of Wireless Sensor Networks with Energy-Harvesting Sensors," *Computer Networks*, vol. 57, no.14, Oct 2013, pp. 2689–704.

- [12] W. Heinzelman, A. Chandrakasan, and H. Balakrishnan, "An Application-Specific Protocol Architecture for Wireless Microsensor Networks," *IEEE Trans. Wireless Commun.*, vol. 1, no. 4, Oct 2002, pp. 660–70.
- [13] G. Fernando *et al.*, "Interval Type-2 Fuzzy Weight Adjustment for Back Propagation Neural Networks with Application in Time Series Prediction," *Info. Sciences*, vol. 260, no. 1, 2014, pp. 1–14.
- [14] J. Wang *et al.*, "PWDGR: Pair-Wise Directional Geographical Routing Based on Wireless Sensor Network," *IEEE Internet of Things J.*, Nov 2014.
- [15] O. Younis and S. Fahmy, "HEED: A Hybrid, Energy-Efficient, Distributed Clustering Approach for Ad Hoc Sensor Networks," *IEEE Trans. Mobile Computing*, vol. 3, no. 4, Oct 2014, pp. 366–79.

BIOGRAPHIES

DAPENG WU (wudp@cqupt.edu.cn) (wudp@cqupt.edu.cn) received his M.S. degree in communication and information system in June 2006 from Chongqing University of Posts and Telecommunications, and his Ph.D. degree from Beijing University of Posts and Telecommunications in 2009. His research interests include ubiquitous networks, IP QoS architecture, network reliability, and performance evaluation in communication systems.

JING HE received her B.S. degree of communication engineering in June 2008 from Yangtze University, College of Technology and Engineering. She is now a Master's degree candidate in information and communication engineering at Chongqing University of Posts and Telecommunications. Her research interests include wireless sensor networks, routing protocol, and energy efficiency.

HONGGANG WANG (hwang1@umassd.edu) worked for Bell Labs Lucent Technologies China from 2001 to 2004 as a member of technical staff. He received his Ph.D. in computer engineering from the University of Nebraska-Lincoln in 2009. He is an assistant professor at the University of Massachusetts (UMass) Dartmouth and is an affiliated faculty member with the Advanced Telecommunications Engineering Laboratory at the University of Nebraska-Lincoln. He is also a faculty member of the Biomedical Engineering and Biotechnology Ph.D. program (BMEBT) at UMass Dartmouth. His research interests include wireless health, body area networks (BANS), cyber security, mobile multimedia and cloud, wireless networks and cyber-physical systems, and big data in mHealth.

CHONGGANG WANG (cgwang@ieee.org) received his Ph.D. degree from Beijing University of Posts and Telecommunications in 2002. He is currently with InterDigital Communications. His R&D focuses on the Internet of Things (IoT), machine-to-machine (M2M) communications, future Internet, and mobile networks, including technology development and standardization. Before joining InterDigital in 2009, he performed R&D at NEC Laboratories America, AT&T Labs Research, the University of Arkansas — Fayetteville, and Hong Kong University of Science and Technology.

RUYAN WANG received his Ph.D. degree in 2007 from the University of Electronic and Science Technology of China (UESTC) and his M.S. degree from Chongqing University of Posts and Telecommunications (CQUPT), China, in 1997. Since December 2002, he has been a professor with the Special Research Centre for Optical Internet and Wireless Information Networks at CQUPT. His research interests include network performance analysis and multimedia information processing.

Wireless Information and Power Transfer: From Scientific Hypothesis to Engineering Practice

Rong Zhang, Robert G. Maunder, and Lajos Hanzo

ABSTRACT

Recently, there has been substantial research interest in the subject of simultaneous wireless information and power transfer (SWIPT) due to its cross-disciplinary appeal and its wide-ranging application potential, which motivates this overview. More explicitly, we provide a brief survey of the state of the art and introduce several practical transceiver architectures that may facilitate its implementation. Moreover, the most important link-level as well as system-level design aspects are elaborated on, along with a variety of potential solutions and research ideas. We envision that the dual interpretation of RF signals creates new opportunities as well as challenges requiring substantial research, innovation and engineering efforts.

INTRODUCTION

Scientifically Oriented Background: In thermal and statistical physics, the earliest and most famous experiment regarding information and energy was conceived by Maxwell in 1867, which is referred to as “Maxwell’s Demon” [1], where the second law of thermo-dynamics was hypothetically violated by the bold hypothesis of information to energy conversion. This stimulated further intriguing research in the mid-20th century as to whether information processing itself dissipates energy, which subsequently led to “Landauer’s principle” [2] suggesting that thermo-dynamically reversible manipulations of information in computing, measurement, and communications do not necessarily dissipate energy, since no energy is required to perform mathematical calculations. Despite the fact that information to energy conversion is elusive, it was suggested from a fundamental perspective that the separate treatment of information and energy may have to be challenged in the practical design of engineering systems.

Naturally, information is carried by attaching itself to a physical medium, such as waves or particles. In molecular and nano-communications, information is delivered by conveying encoded particles from the source to the destina-

tion. Similarly, in optical communications, information is delivered by photons having information-dependent intensities, which may be detected by a photon counting process. Given the nature of the process, the system is capable of providing a heating/illumination/propulsion function. Both of the above examples suggest that the underlying matter that carries information can be effectively reused for diverse applications. The explicit concept of transporting both information and energy simultaneously was raised by the authors of [3], which was recently further extended to the wireless regime [4, 5] since the underlying electro-magnetic (EM) wave can carry both information and energy. Thus, it is desirable that a mobile device is free from being tethered in any way.

Engineering-Oriented Background: Before delving into the topic of simultaneous wireless information and power transfer (SWIPT) [4, 5], a brief introduction of wireless power transfer (WPT) is warranted [6]. Generally speaking, WPT can be carried out in two basic ways, based on either *near-field EM induction* in the form of inductive coupling and resonant coupling relying on coils or *far-field EM radiation* using microwave frequencies by relying on so-called rectennas. When compared to near-field EM induction, which only tolerates a small misalignment between the transmitter and receiver, far-field microwave power transfer (MPT) supports a wider coverage area and thus may be considered to be more suitable for employment in SWIPT systems.

The earliest experiments on MPT were conducted by Tesla with the ultimate goal of creating a worldwide wireless power distribution system. In the mid-20th century, MPT was conceived for high-power applications in the mega/kilowatt range. The need for power increased substantially due to the development of electronic devices in the late 20th century. In particular, research has been focused on the design of compact and efficient rectennas conceived for applications in the milli/microwatt range, where the startup company Powercast¹ has reported that microwatt-scale power was

The authors are with the University of Southampton.

The financial support of the EPSRC under the India-UK Advanced Technology Centre (IUCATC), that of the EU under the Concerto project as well as that of the European Research Council’s (ERC) Advanced Fellow Grant is gratefully acknowledged.

¹ <http://www.powercast-co.com/>

Due to the fact that EM radiation is restricted by both health and safety regulations, it remains an open challenge at the current state of the art to power a mobile phone, typically requiring around a few hundred milliwatts of power.

transmitted over a distance of a few meters at transmission power of 23 dBm at a frequency around 900 MHz.

As is widely known across the wireless community, another innovative and revolutionary exploitation of EM wave propagation over the last century was wireless information transfer (WIT), inspired by the late 19th century radio experiments conducted by Marconi, which eventually led to the mid 20th century tactical military use of radar in the Second World War and to the pervasive commercial revolution of the mobile industry. It is thus the maturing WPT and WIT fields that make SWIPT an interesting emerging research topic.

Application-Oriented Background: Compared to the pervasive WIT, WPT is less well exploited at the time of writing, which is mainly due to the high attenuation of RF signals over distance. Thanks to recent advances in both antenna technologies and power electronics, significantly increased power can now be transferred to wireless devices, which are themselves becoming increasingly more energy-efficient. Thus, SWIPT becomes especially compelling in scenarios where charging the battery incurs either high risk or high cost, as in medical implants, in-building sensors, and so on. To elaborate a little further, conventional energy harvesting relies on ambient energy sources (wind, vibrations, heat, etc.), and at the current state of the art they are only capable of supporting low-rate communications. By contrast, the unique characteristic of a SWIPT system is that it is capable of operating in an environment with insufficient ambient energy sources, while delivering a controllable amount of wireless information and energy concurrently, hence supporting low-cost sustainable operations.

Due to the fact that EM radiation is restricted by both health and safety regulations, it remains an open challenge at the current state of the art to power a mobile phone, typically requiring around a few hundred milliwatts of power. Nonetheless, this technology is becoming more appealing for employment in low-power sensors. Consequently, it is promising for the family of near-future wireless systems with the goal of connecting billions of low-power devices globally. In this scenario, known as the Internet of Things (IoT), SWIPT is expected to become a pervasive enabler, which is “immortal” from an energy replenishment perspective. To be more specific, low-power sensor networks, wireless body area networks, and long-distance RF identification (RFID) networks would directly benefit from SWIPT techniques.

Organization — Based on the above background, we provide an overview of the emerging topic of SWIPT, commencing from a crisp literature review to elaborating on practical architectures as well as on their important design aspects and their potential solutions. Our overview article is thus organized as follows. We commence by surveying the existing literature on SWIPT, with an emphasis on its theoretical aspects. We then summarize the existing practical architectures facilitating SWIPT, where three potential archi-

tectures are introduced. This is followed by our elaborations on the key link-level components, ranging from multi-antenna-aided techniques to both coding and modulation. These link-level components are the fundamental building blocks of system-level studies, as discussed, where we explicitly reveal the associated pros and cons of interference. Finally, we conclude our discourse.

STATE OF THE ART

Let us commence with a crisp survey of the SWIPT literature, followed by the summary of practical transceiver architectures. Due to the limited number of references allowed, we apologize that only a few of the seminal SWIPT references are cited; a comprehensive survey will follow in our future work.

THEORETICAL STUDIES

There is a paucity of literature on the current theoretical research advances in SWIPT, which concentrate almost entirely on the *energy vs. capacity* trade-off. The pioneering work of [3] revealed that a nontrivial trade-off exists between the available energy and the achievable capacity for typical channel models, such as the binary symmetric channel (BSC) and the additive white Gaussian noise (AWGN) channel, indicating that maximizing the information rate is to a degree coupled with maximizing the energy transfer. Indeed, when a wideband fading channel is considered, the authors of [4] confirmed this fundamental trade-off between energy and capacity, where the classic water-filling-based power allocation applied across the entire frequency band maximized the attainable information rate, while energy transfer was maximized by transmitting at a single frequency using the total available power, thus reducing the information rate. Similar conclusions were also drawn concerning this fundamental trade-off for narrowband multiple-input multiple-output (MIMO) channels [5], where water-filling-based power allocation spanning all eigenvalues of the channel matrix maximizes the information rate. By contrast, the energy transfer is maximized by concentrating all available power in the specific direction corresponding to the maximum eigenvalue. Further research considered various other important channel models, such as the families of interference channels, multiple access channels, unicast/multicast channels, and secrecy channels.

Alongside those theoretical results, practical receiver architectures facilitating SWIPT were proposed in [7]. In addition to the widely investigated *one-way* SWIPT focusing on receiver architectures, the scope of SWIPT may also be further expanded to *two-way* SWIPT [8], where a pair of nodes interactively communicate and exchange power. A related topic is the so-called wireless-powered communications [9], where the transmitter conveys power to the receiver, which is then converted to DC power and reused for the destination's information transmission in the reverse direction. Finally, a range of other treatises considered diverse energy-transfer-aided systems, such as multi-carrier systems, relay-assisted arrangements, cognitive-radio settings, and beamforming-aided systems.

PRACTICAL ARCHITECTURES

Most of the above information theoretical insights were drawn under the assumption that the information can be retrieved at the same time the power is received. However, practical transceivers may not readily support this operation. From an implementational point of view, there are four types of transceiver architectures that may be capable of supporting a SWIPT system, as shown in Fig. 1.

Building Blocks: Before proceeding to our detailed elaborations, we first introduce the concept of power transfer tunnel and information transfer tunnel, as seen in Fig. 1a. The *power transfer tunnel* typically consists of a rectifier used for converting the received RF power to DC power, which is followed by a (multi-stage) DC to DC booster. Subsequently the power is stored in the battery. Note that from a pure power transfer point of view, the antennas and rectifier are jointly known as a rectenna. On the other hand, the *information transfer tunnel* typically consists of a baseband (BB) digital signal processing (DSP) module invoked after the RF-to-BB down conversion process at the front-end. More explicitly, there are various RF-to-BB designs, which typically include filters, mixers, amplifiers, analog-to-digital (A/D) converters, and so on.

Different EM Waves: The first type of architecture supports a SWIPT system by employing two different EM waves for information transfer and power transfer, respectively [10]. The most straightforward option is the *parallel independent* architecture seen in Fig. 1a, where two transmitter and receiver pairs are applied in parallel using two well isolated EM waves, naturally one for power transfer and one for information transfer. The two interfaces may be independently operated or coordinated by a controller. By contrast, an amalgamated option is constituted by the *parallel combined* architecture seen in Fig. 1b, where the power transfer process is incorporated into the information transfer process by mixing the low-frequency carrier used for power transfer with the high-frequency carrier invoked for information transfer. Following the down-conversion operation of the superheterodyne receiver, the low-frequency carrier shifted to DC by the rectifier subsequently enters the power transfer tunnel, while the high-frequency carrier shifted to the intermediate frequency (IF) is further processed before entering the information transfer tunnel.

Identical EM Waves: The second type of architecture facilitates SWIPT by relying on a single EM wave for information transfer and power transfer. Practical receivers facilitating SWIPT may be operated in two different modes [5], on either a *time-switching* basis or a *power-splitting* basis, as seen in Fig. 1c. To be more specific, in time-switching mode, the receiver alternatively and opportunistically activates the information transfer tunnel and the power transfer tunnel. On the other hand, in the power-splitting mode, a certain portion of the received power is used

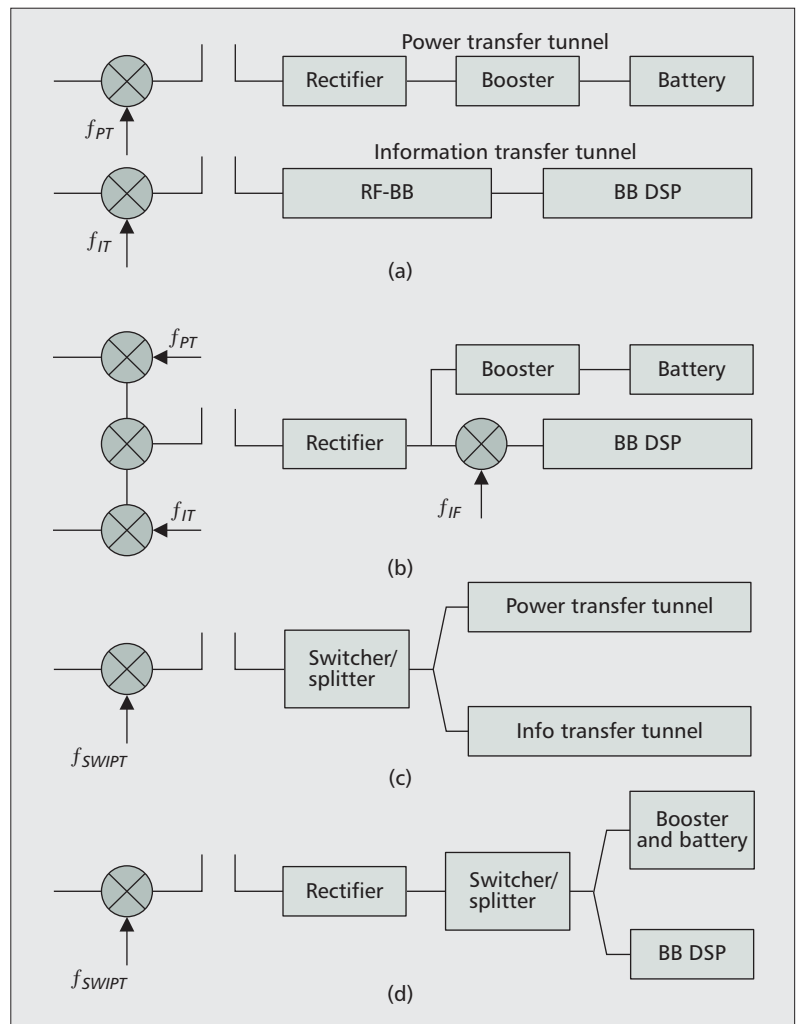


Figure 1. Transceiver architectures conceived for SWIPT: a) parallel independent scheme; b) parallel combined scheme; c) time switching/power splitting scheme; d) integrated scheme.

for powering the receiver, while the remaining received power is used for retrieving information. Note that the time-switching receiver may be viewed as a special case of the power-splitting receiver when the splitting ratio is either one or zero. When channel information is available at the transmitter, the two operating modes may be further optimally configured.

Varshney's Concept: Apart from the above-mentioned architectures, the more beneficial fully integrated information and power transfer philosophy can be used as the ultimate objective in the spirit of Varshney's seminal concept [7], which proposed that energy and information transfer should be innately inter-linked. In this light, the *integrated* architecture of [7] was proposed, which is seen in Fig. 1d, where the rectifier employed for power transfer at the receiver also acts as the front-end for RF-to-BB down-conversion. The resultant DC signal is then passed to both the energy storage and the information retrieval blocks in either a switching or splitting approach. In order to enable SWIPT using an integrated receiver, a specifically designed modulation-specific energy is required,

The channel coding design in SWIPT must also strike a desirable energy efficiency balance between the transmission energy efficiency and processing energy efficiency. This is because channel codes that facilitate high transmission energy efficiency are typically complex and hence suffer from low processing energy efficiency.

which relies on Varshney's concept stating that "information is patterned matter energy." Finally, it is also worth noting that the underlying difference between the approaches of Figs. 1c and 1d is that the former requires both a rectifier and RF-to-BB down-conversion, while the latter only needs a rectifier.

DESIGN ASPECTS AND POTENTIAL SOLUTIONS

Let us now focus our attention on exploiting the identical EM waves and hence the architectures discussed earlier, while considering the important aspects in designing the SWIPT system from both a link-level as well as a system-level perspective, along with a variety of promising solutions and potential research ideas.

LINK-LEVEL DESIGN ASPECTS

Several important building blocks deserve special attention in the link-level design of SWIPT. These include efficient and compact design of rectennas, smart design of battery management, and powerful as well as robust design of BB signal processing algorithms. Let us now elaborate on the last point from a communications perspective in the context of a SWIPT system, while referring readers to [11] for cross-disciplinary details on rectenna design and battery management. More particularly, when taking into account the characteristics of a SWIPT system, three physical layer components have to be discussed in more detail, with the design aspects of multi-antenna-aided techniques being the most critical ones for the practical exploitation of SWIPT.

Multiple Antennas: Key to Practical Applications — There are several constraints that fundamentally limit the development of SWIPT systems. First, due to the fact that EM radiation is restricted by both health and safety regulations, the power supplied by the source is typically limited. For example, a macro base station emits 46 dBm power, while an indoor access point has a transmit power of 23 dBm. Second, the hostile wireless propagation, including path loss, shadowing, and multipath fading, substantially reduces the average received RF power. Third, a state-of-the-art rectenna exhibits a conversion efficiency of about 50 percent reported by Powercast, which can only be activated above a certain RF input power level, typically above -15 dBm. As a result, the combined effect of the above facts suggests that a SWIPT system is limited to a very *short range*.

To extend the attainable range of SWIPT, multiple antennas are necessary, since they are capable of providing a larger antenna aperture and a higher antenna gain. Practically, to accommodate a large number of antennas in a compact shirt-pocket-sized communicator, a higher carrier frequency is found to be beneficial in a SWIPT system, such as the 5.8 GHz band and higher. When equipped with multiple antennas, two different signal processing operations are promising: analog domain *beamforming* using a phase shifter with/without complex weighting

and digital domain *precoding*, which may be flexibly designed for satisfying a predefined rate and/or power constraint.

Indeed, beamforming techniques may be deemed the key for SWIPT systems, and most of the open literature is focused on this topic. For example, Xu *et al.* in [12] considered the energy beamforming concept in the context of a multiple-input single-output (MISO) downlink system. Furthermore, Park and Clerckx in [13] designed a beamforming-aided solution for a K -user interference-limited system relying on multiple transmitter and receiver pairs, while Li *et al.* in [14] considered a collaborative energy beamforming design using multiple transmitter and receiver pairs. The beneficial impact of large-scale/massive MIMOs was studied in [15]. Finally, active research also addressed a range of practical aspects, such as only having partial channel state information (CSI) relying on a realistic CSI feedback design. Regrettably, these solutions cannot be cited here due to the limited number of allowable references.

Channel Coding: Key to Link Reliability — SWIPT would cause energy depletion that will result in *processing-induced errors* due to voltage variations in addition to *channel-induced errors* due to wireless propagation only found in conventional wireless communications. To elaborate a bit more, SWIPT receivers are usually considered to be passive or semi-passive, relying on the received power to drive the receiver circuits. In this scenario, the time-varying operating and channel conditions can result in fluctuations in the received voltage level from time to time. In this occurrence, the propagation delays of the electronic signals within the receiver circuits will be extended and may exceed the period of the clock that synchronizes the timing of the circuits. This will result in processing errors, corrupting the data represented by the electronic signals in the receiver circuits. Therefore, inherently *robust* channel coding design is required in order to mitigate both the channel-induced errors as well as the processing-induced errors, which may occur within any of the receiver circuits including the channel coding circuit itself.

The channel coding design in SWIPT must also strike a desirable *energy efficiency balance* between the transmission energy efficiency and processing energy efficiency [16]. This is because channel codes that facilitate high transmission energy efficiency are typically complex and hence suffer from low processing energy efficiency. In this context, Fig. 2 compares both the transmission-related and signal-processing-induced energy consumption of a range of different channel codes. In addition, another characteristic of SWIPT systems is their potential *discontinuous operation*, where discontinuous energy and information reception will be observed at the receiver. In this light, the channel coding design should have to tolerate this variability or discontinuity, which would require a powerful short-delay code design. To sum up, in order to holistically design a channel coding scheme for SWIPT systems, both the hardware as well as the algorithmic aspects of channel coding have to be revisited. It is thus envisioned

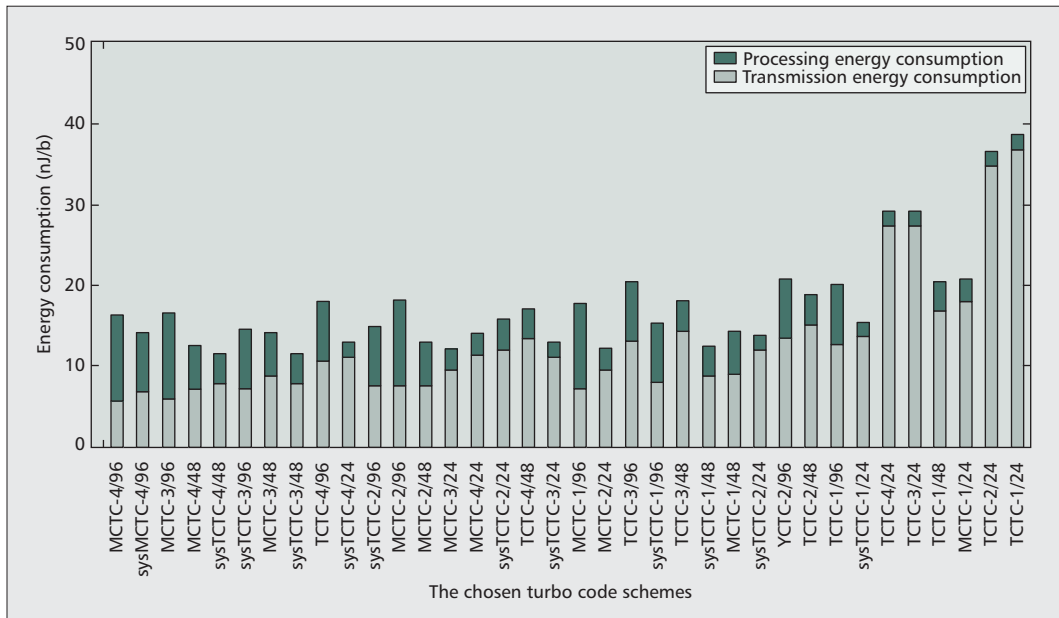


Figure 2. Transmission energy consumption and processing energy dissipation of different turbo codes at a bit error ratio of 10^{-5} . In this plot, we listed different twin-component and multi-component turbo codes in the format of TCTC-x/c and MCTC-x/c, where x is shorthand for a particular combination of coding parameters, and c is their fixed complexity. More explicitly, sufficiently high transmission energy reserves are required to maintain an acceptable signal-to-noise ratio (SNR), while the processing energy consumption encapsulates the energy dissipation imposed by four dominant turbo code components: the datapath, the controller, the memory, and the interleaver. This plot clearly exemplifies that jointly considering both power dissipation factors of the channel codes is indeed important in the design of passive/semi-passive receiver. Finally, note that the complexity is given as $2^m * N$, with m being the number of memory elements per component codes and N being the number of BCJR operations with sequential decoder activation order.

Another characteristic of SWIPT systems is their potential discontinuous operation, where discontinuous energy and information reception will be observed at the receiver. In this light, the channel coding design should have to tolerate this variability or discontinuity, which would require a powerful short-delay code design.

that delay-limited short code design is required for counteracting their potentially discontinuous operations when also taking their energy depletion problems into account.

Modulation: Key to Data Integrity — When the power-splitting receiver architecture of Fig. 1c is employed, any conventional modulation scheme may potentially be employed. By contrast, when Varshney’s principle is adopted and the integrated receiver architecture of Fig. 1d is employed, energy-conscious modulation schemes would be desired. This is because no explicit RF-to-BB down-conversion chain is used by the integrated receiver of Fig. 1d. On the other hand, the rectennas invoked for power conversion at the front-end of the integrated receiver extract only the energy. This is reminiscent of energy harvesting from light, where, for example, pillars of silicon are aligned in parallel to the incoming light in order to improve the achievable photon harvesting, and a radial PN junction is used to efficiently collect the carriers. Hence, the intensity modulation, which is routinely used in optical communications, is of great interest. A typical example is constituted by classic on-off keying, where a binary one carries power while a binary zero does not [8]. At the same time, the pattern of power transfer (the particular instance of the appearance of a one) conveys information. Hence, the duality of power and information transfer explicitly manifests itself. A related observation is that a variable-length constraint

channel coding scheme [17] that generates unequal number of ones and zeros may be found beneficial in satisfying the joint power transfer and information transfer requirements.

Inspired by spatial modulation, a novel generalized modulation scheme conceived for SWIPT was proposed by invoking multiple antennas [18], where information is carried not by the classic radio *waveforms* but by energy *patterns*, as seen in Fig. 3. More explicitly, the specific choice of the transmitted pattern energy embeds information into the pattern of the power delivered, which may assume:

- A position-based energy pattern, which is reminiscent of the pulse position modulation (PPM) concept, but invoked in the spatial domain.
- An intensity-based energy pattern, which is similar to pulse amplitude modulation (PAM), but exclusively relying on positive values

Note that other intensity/pulse based modulation schemes, such as pulse width modulation (PWM) and pulse interval modulation (PIM), may also be beneficially combined with the concept of energy-pattern-based transmission conceived for an integrated receiver. Finally, the rate vs. energy trade-off of the advocated energy-pattern-aided SWIPT is illustrated in Fig. 4.

SYSTEM-LEVEL DESIGN ASPECTS

The above-mentioned key components are critical enablers in the development of SWIPT at the

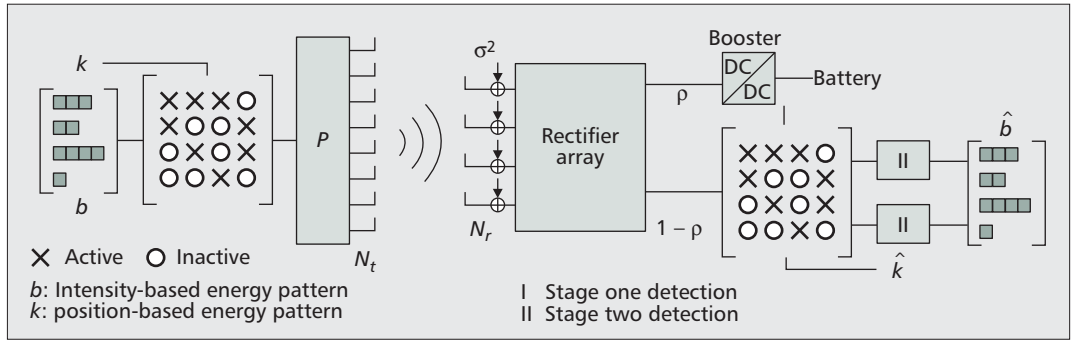


Figure 3. Energy pattern aided SWIPT with integrated receiver. In this plot, position-based energy patterns are used for activating a particular receive antenna index using an appropriate transmit precoding method. The intensity-based energy pattern is generated by mapping the positive PAM symbols onto the activated receive antennas. At the receiver, the rectifier array serves the dual purpose of converting RF power to DC power for power transfer and the RF signal to BB signal for information transfer, which is arranged in two stages for the sake of extracting the information embedded in the position-based and intensity-based energy patterns.

link level for transmission from a single point to a single point. Let us now focus our attention on the system-level design of SWIPT.

Beyond Point-to-Point Transmission — We commence our discussions with the widely considered scenario of a single-cell multi-user SWIPT system. A typical observation in this sce-

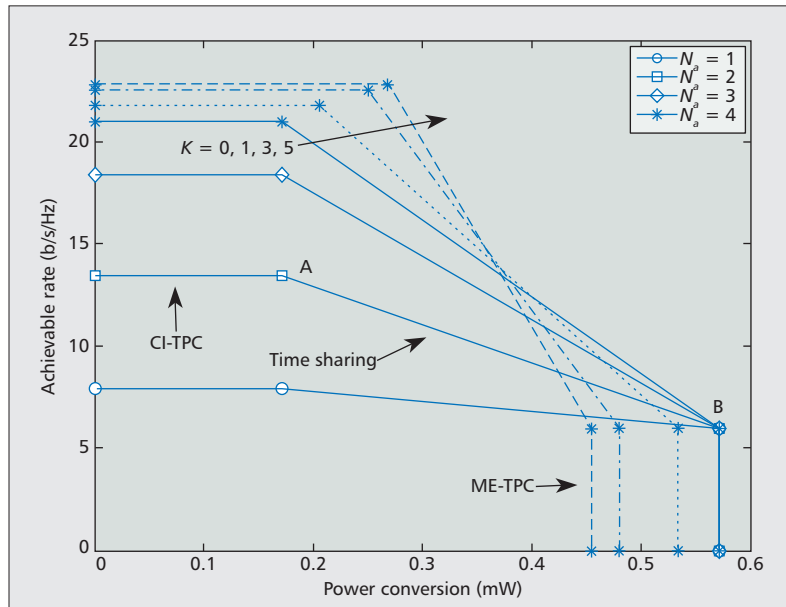


Figure 4. Trade-offs between the achievable rate and power conversion of an $\{N_t, N_r\} = \{8, 4\}$ energy-pattern-aided SWIPT system employing different numbers of activated receive antennas $N_a = [1, 2, 3, 4]$ for position-based energy pattern aided and $M = 64$ -ary intensity-based energy patterns assisted transmission. We used Rician fading associated with $K = [0, 1, 3, 5]$ at a transmit power set to $P_t = 30$ dBm with the path loss model of $P_r = P_t A_t A_r (d_0 \lambda_c)^{-\alpha_{pl}}$, where $d_0 = 5$ m is the distance between the TA and RA arrays. Furthermore, λ_c denotes the wavelength corresponding to a carrier frequency of 5.8 GHz, the path-loss exponent is set to $\alpha_{pl} = 4$ and the aperture of the TA and RA is represented by A_t and A_r , respectively, with the per-antenna aperture being set to 1 cm^2 . Moreover, both the channel inversion (CI)-based transmitter pre-coding (TPC)-aided arrangement and the maximum eigenvalue (ME)-based TPC-aided arrangement are characterized, where the lines connecting points A and B are achieved by time-sharing between the two TPC methods.

nario is the *coverage difference* between the power transfer zone and information transfer zone [9]. This is due to the different operational point of the power transfer functionality and the information transfer functionality, where the former typically operates above, say, -15 dBm, while the latter only has to maintain a certain average received signal-to-noise ratio (SNR). In a noise-limited system, the noise floor is typically lower than -100 dBm (including both the thermal noise and processing noise). Then the requirement of a 10 dB average received SNR would result into an operational point of -90 dBm. Hence again, this disparity translates into a substantially different radius for the power-transfer zone and for the information-transfer zone.

A plausible observation in this scenario is that there is a potentially imminent benefit due to having a *multi-user gain* when compared to the point-to-point case. This is because the existence of multiple users provides scheduling opportunities for ensuring that the power supplied at the source would not be wasted by the users experiencing severely faded channels. As a result, in a single-cell multi-user SWIPT system, a challenge arises when the transmitter has to optimally allocate its resources and then has to efficiently schedule its information as well as its power transfer action associated with the most appropriate multi-antenna signaling methods. This will become more insightful if a specific quality of service (QoS) metric, such as the delay, is also taken into account.

The Pros and Cons of Interference: When considering a multi-cell multi-user SWIPT system, a conflicting view of the interference emerges, since interference is desirable from a power transfer point of view, but it is unwanted from an information transfer point of view. Intelligently harnessing the interference would reshape the above-mentioned coverage gap between the power-transfer zone and the information-transfer zone. More explicitly, the *proactive* prevention of interference, with the aid of interference avoidance via (fractional) frequency reuse and interference averaging via frequency

hopping, may require a radically new design in the context of SWIPT. On the other hand, the *passive* retrospective reduction of interference, such as parallel or serial interference cancellation and interference suppression relying, for example, on interference rejection combining (IRC) may still be employed. As further advances, the interference may be beneficially turned into a precious source of designed signal energy when the concept of coordinated multi-point (CoMP) transmission is adopted in a multi-cell multi-user SWIPT system [14]. Finally, interference alignment may require further fundamental revisiting of the impact of the power transfer requirements.

A closely coupled topic with interference as the central focus lies in the cognitive model, where in addition to the SWIPT system established over the primary link, the battery of the passive or semi-passive primary receiver may also be charged by the transmissions of the secondary link. As a result, a conventionally harmful secondary transmission becomes potentially useful in the context of a SWIPT system. In this case, the underlying interference model is of paramount importance, where in a pair of primary and secondary links, four potential operation combinations are feasible at the respective receivers of the primary and secondary links, as detailed in [13]. This naturally necessitates further game theoretical and information theoretical insights. Last but not least, there has been emerging work on system-level SWIPT, which relied on powerful stochastic geometry to draw macroscopic impact [19].

CONCLUSIONS

We introduce the background rationale and motivation of SWIPT systems with the aid of a brief survey and provide a basic architectural summary. We also discuss its practicality by detailing several critical aspects ranging from link-level to system-level design. We envision that the dream “to facilitate and cheapen the transmission of intelligence” of Nikola Tesla will become a reality, with the added benefits of having a substantially reduced carbon footprint, provided that the research community improves both rectenna and battery design as well as communications and signal processing techniques.

REFERENCES

- [1] K. Maruyama, F. Nori, and V. Vedral, “The Physics of Maxwell’s Demon and Information,” *Reviews of Modern Physics*, vol. 81, Jan 2009, pp. 1–23.
- [2] R. Landauer, “Irreversibility and Heat Generation in the Computing Process,” *IBM J. Research and Development*, vol. 5, no. 3, 1961, pp. 183–91.
- [3] L. Varshney, “Transporting Information and Energy Simultaneously,” *2008 IEEE Int’l. Symp. Info. Theory*, 2008, pp. 1612–16.
- [4] P. Grover and A. Sahai, “Shannon Meets Tesla: Wireless Information and Power Transfer,” *2010 IEEE Int’l. Symp. Info. Theory*, 2010, pp. 2363–67.
- [5] R. Zhang and C. K. Ho, “MIMO Broadcasting for Simultaneous Wireless Information and Power Transfer,” *IEEE Trans. Wireless Commun.*, vol. 12, no. 5, May 2013, pp. 1989–2001.

- [6] J. Garnica, R. Chinga, and J. Lin, “Wireless Power Transmission: From Far-Field to Near-Field,” *Proc. IEEE*, vol. 101, no. 6, June 2013, pp. 1321–31.
- [7] X. Zhou, R. Zhang, and C. K. Ho, “Wireless Information and Power Transfer: Architecture Design and Rate-Energy Tradeoff,” *IEEE Trans. Commun.*, vol. 61, no. 11, Nov, 2013, pp. 4757–67.
- [8] P. Popovski, A. Fouladgar, and O. Simeone, “Interactive Joint Transfer of Energy and Information,” *IEEE Trans. Commun.*, vol. 61, no. 5, May 2013, pp. 2086–97.
- [9] H. Ju and R. Zhang, “Throughput Maximization for Wireless Powered Communication Networks,” *IEEE Trans. Wireless Commun.*, vol. 13, no. 1, Jan. 2014, pp. 418–28.
- [10] S. Yoshida *et al.*, “Experimental Demonstration Of Coexistence of Microwave Wireless Communication and Power Transfer Technologies for Battery-Free Sensor Network Systems,” *Int’l. J. Antennas and Propagation*, 2013.
- [11] Z. Popovic *et al.*, “Low-Power Far-Field Wireless Powering for Wireless Sensors,” *Proc. IEEE*, vol. 101, no. 6, June 2013, pp. 1397–1409.
- [12] J. Xu, L. Liu, and R. Zhang, “Multiuser MISO Beamforming for Simultaneous Wireless Information and Power Transfer,” *IEEE Trans. Sig. Proc.*, vol. 62, no. 18, Sept 2014, pp. 4798–4810.
- [13] J. Park and B. Clerckx, “Joint Wireless Information and Energy Transfer in a K-User MIMO Interference Channel,” *IEEE Trans. Wireless Commun.*, vol. 13, no. 10, Oct. 2014, pp. 5781–96.
- [14] S. Lee, L. Liu, and R. Zhang, “Collaborative Wireless Energy and Information Transfer in Interference Channel,” to appear, *IEEE Trans. Wireless Commun.*
- [15] X. Chen, X. Wang, and X. Chen, “Energy-Efficient Optimization for Wireless Information and Power Transfer in Large-Scale MIMO Systems Employing Energy Beamforming,” *IEEE Wireless Commun. Letters*, vol. 2, no. 6, Dec. 2013, pp. 667–70.
- [16] L. Li *et al.*, “Energy-Conscious Turbo Decoder Design: A Joint Signal Processing and Transmit Energy Reduction Approach,” *IEEE Trans. Vehic. Tech.*, vol. 62, no. 8, Oct. 2013, pp. 3627–38.
- [17] A. Fouladgar, O. Simeone, and E. Erkip, “Constrained Codes for Joint Energy and Information Transfer,” *IEEE Trans. Commun.*, vol. 62, no. 6, June 2014, pp. 2121–31.
- [18] R. Zhang, L.-L. Yang, and L. Hanzo, “Energy Pattern Aided Simultaneous Wireless Information and Power Transfer,” *IEEE JSAC*, vol. 33, no. 9, Sept. 2015, pp. 1–13, <http://eprints.soton.ac.uk/365524/>.
- [19] I. Krikidis, “Simultaneous Information and Energy Transfer in Large-Scale Networks With/Without Relaying,” *IEEE Trans. Commun.*, vol. 62, no. 3, Mar 2014, pp. 900–12.

BIOGRAPHIES

RONG ZHANG (rz@ecs.soton.ac.uk) received his Ph.D. in 2009 from the University of Southampton, United Kingdom. He was a research assistant at the Mobile Virtual Center of Excellence, United Kingdom, a research fellow at the University of Southampton, and a system algorithms expert for Huawei Sweden R&D. He is now a lecturer at the Southampton Wireless group of ECS, Southampton University. He has 40+ journals in prestigious publication avenues and many more in major conference proceedings. More details can be found at <http://www.ecs.soton.ac.uk/people/rz>

ROBERT G. MAUNDER (rm@ecs.soton.ac.uk) has studied at the University of Southampton since October 2000. He was awarded a first class honors B.Eng. in electronic engineering in July 2003, a Ph.D. in wireless communications in December 2007, and became an associate professor in February 2013. His research interests include joint source/channel coding, iterative decoding, irregular coding, and modulation techniques. More details can be found at <http://www.ecs.soton.ac.uk/people/rm>

LAJOS HANZO (lh@ecs.soton.ac.uk) is a Wolfson Fellow of the Royal Society, Fellow of the Royal Academy of Engineering (FREng), FIEEE, FIET, and a EURASIP Fellow. He has co-authored 20 IEEE Press-Wiley books, published 1400+ research entries at IEEE Xplore, organized and chaired major IEEE conferences, and has been awarded a number of distinctions. More details can be found at <http://www.mobile.ecs.soton.ac.uk>

We envision that the dream “to facilitate and cheapen the transmission of intelligence” by Nikola Tesla would become a reality, with the added benefits of having a substantially reduced carbon footprint, provided that the research community improves both the rectenna and battery design as well as the communications and signal processing techniques.

Delay-Sensitive Dynamic Resource Control for Energy Harvesting Wireless Systems with Finite Energy Storage

Fan Zhang and Vincent K. N. Lau

ABSTRACT

Energy harvesting technology has become a promising solution to enhance the energy efficiency and reduce carbon emissions in future wireless systems. In future wireless systems, most of the data throughput will come from *delay-sensitive* applications. To ensure a good experience for an end user, we target the optimization of the delay performance of an EH wireless system with finite energy storage. As such, it is necessary to adapt the resource allocation to the channel fading information, data queue length, and energy queue length information. The channel fading information provides the channel quality, the data queue length information provides the dynamic urgency of the transmitted data flows, and the energy queue length information provides the information on how much available energy is left in the energy buffer. Such a problem is quite challenging because it belongs to an infinite dimensional *stochastic optimization*. In this article, we review the existing works on the resource allocation problem in EH wireless systems. We also propose a low-complexity delay-sensitive resource control scheme and discuss valuable design insights.

INTRODUCTION

In recent years, energy harvesting (EH) technology has become a promising solution to the energy efficiency issue in future wireless systems because it not only prolongs the operation lifetime of battery-limited devices but also helps reduce greenhouse gas (such as carbon) emissions. Furthermore, EH technology is quite popular and is being intensively discussed for designing the future wireless systems, such as D2D communications, EH wireless sensor networks, and future cellular systems. Specifically, in EH wireless systems, the transmission nodes harvest energy from the ambient environment by means of EH devices, such as solar panels, wind turbines, and thermoelectric generators, and convert the harvested renewable energy into electricity [1]. However, since the renewable energy sources may appear to be random and bursty in

nature, *energy storage* is needed to buffer the random and bursty supply of renewable energy. Due to the high cost of large-capacity energy storage, in practice, EH devices usually have finite energy storage. Moreover, it is important to dimension the energy storage capacity so as to achieve both stability and good performance of the EH network. On the other hand, delay-sensitive applications such as video streaming and online gaming will take up a significant portion of the capacity demand in future wireless systems. Current wireless systems, such as WiFi and third generation (3G), cannot ensure a good experience for an end user with delay-sensitive applications. Therefore, it is also very important to take into account the delay requirements in designing resource control schemes to support delay-sensitive applications in EH wireless systems.

In this article, we study dynamic resource control in an EH wireless system with finite energy storage, as shown in Fig. 1, to support real-time delay-sensitive applications. The transmitter is solely powered by a solar panel that harvests energy from the surrounding environment. In order to have good delay performance, the dynamic resource control should be adaptive to the channel state information (CSI), data queue state information (DQSI), and energy queue state information (EQSI). Specifically, the CSI reveals the transmission opportunities of the time-varying wireless fading channel between the transmitter and receiver, the DQSI reveals the dynamic urgency of data flows, and the EQSI reveals information on the quantity of available renewable energy in the energy buffer. A control policy adaptive to CSI, DQSI, and EQSI is very challenging because the dynamic resource control problem belongs to an infinite dimensional *stochastic optimization*. In addition, the coupling between the data queue and energy queue in an EH wireless system further complicates the problem. It is well known that the Markov decision process (MDP) is widely used to deal with such a stochastic optimization [2]. However, classical value iteration algorithms (VIAs) [3] for solving the MDP only give numerical solutions, which suffer from slow convergence issues and lack of design insights. In this

The authors are with Hong Kong University of Science and Technology.

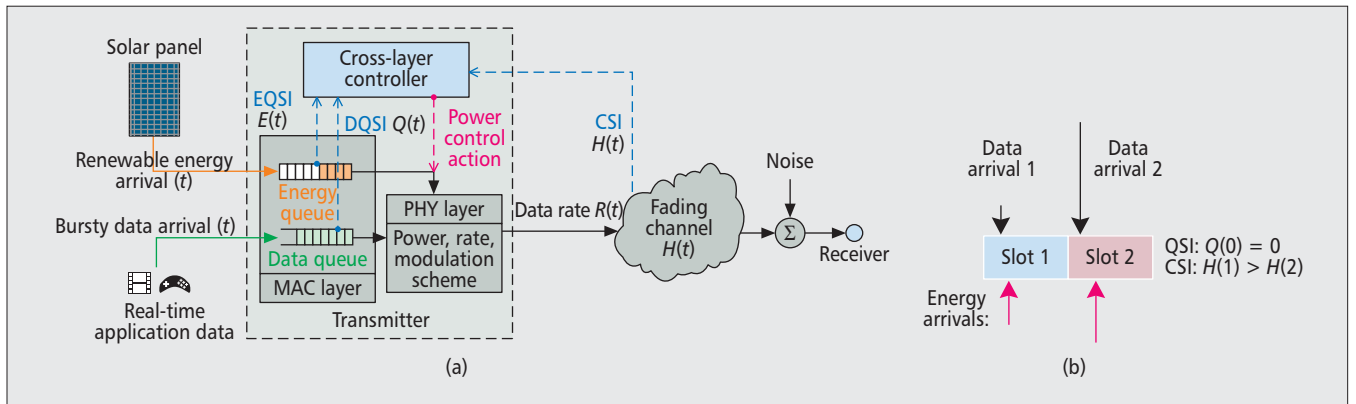


Figure 1. a) System model of a point-to-point EH wireless system with finite energy storage; b) an example of the bursty data arrivals and energy arrivals for two consecutive slots.

article, we provide an asymptotically optimal low-complexity delay-aware resource control scheme and obtain valuable design insights on the dynamic resource for EH wireless systems.

PRIOR WORKS ON RESOURCE CONTROL FOR EH WIRELESS SYSTEMS

In this section, we review some of the prior works on resource control for EH wireless systems. The differentiations among them are summarized below.

INFINITE/FINITE ENERGY STORAGE CAPACITY

Due to the random and bursty nature of renewable energy sources, energy storage is used to buffer the renewable energy so that wireless devices can have a sustained energy supply for delivering data to remote receivers. For example, in a one-day span, sunlight is weak during sunrise and sunset, intense at noon, and gone in the night. Therefore, it is desirable to harvest sufficient solar energy in the daytime and save it in storage to maintain a functional wireless system without sunlight at night. Some prior works assume that the energy storage in wireless devices has infinite capacity [4]. However, this is not practical due to the high cost of high-capacity renewable energy storage, and *energy storage with finite capacity* is a key cost component in the EH wireless systems. As such, it is necessary to appropriately choose the renewable energy buffer size and analyze the impact on how the finite energy storage affects system performance. There are some works that consider EH systems with finite energy storage [5–7], but the derived schemes are not good for delay-sensitive applications. Later in the article, we give the minimum energy buffer size to not only ensure the stability but also achieve good delay performance of an EH system by taking into account the random and bursty nature of renewable energy sources.

NON-CASUAL/CASUAL CONTROL

There are a lot of existing works on EH wireless systems in which it is presumed that the EH devices have non-causal knowledge of the time-

varying wireless channel, bursty data arrival, and renewable energy arrival profiles [4–6]. That is, the wireless device knows the future realizations of the channel conditions, data arrivals, and renewable energy arrivals. Such an assumption enables mathematical tractability of the associated resource allocation problem. However, this assumption is not realizable in practice due to the difficulty in predicting the random channel conditions and renewable energy source activities. In the article, we consider causal power control, which means that the control action depends on the instantaneous CSI, DQSI, and EQSI. In addition, [8, 9] consider causal transmission mode control based on the observed system state for EH wireless systems with infinite energy storage. They consider the minimizations of average power consumption and packet error rate, and assume that the information flow is delay-insensitive, which does not guarantee any delay performance.

CONTROL OBJECTIVES

For delay-sensitive applications, we need to target minimization of the end-to-end average delay performance of EH wireless systems. Specifically, we define delay as the average time between a data packet entering the data queue buffer at the transmitter to the time when the packet is received at the receiver. Different control objectives are considered in prior works on EH wireless systems. For example, some focus on maximization of the transmission data rate for a given deadline [5], while some consider minimization of the packet error rate or transmission complication time for a given data rate [6–9]. However, these formulations do not translate into delay minimization for delay-sensitive applications. We illustrate this using the toy example in Fig. 1b: we focus on two consecutive slots, where the CSI quality of the first slot is better than the second slot (i.e., $H(1) > H(2)$). The data arrival of the first slot is much smaller than the second slot, and there are some energy arrivals for both slots. Using the algorithms proposed in [5, 6], since $H(1) > H(2)$, more harvested energy is used to deliver data in slot 1 than slot 2. However, since data arrival 1 < data arrival 2, the average queue length (i.e., delay performance) will be larger. From the above

Due to the random and bursty nature of renewable energy sources, energy storage is used to buffer the renewable energy so that wireless devices can have a sustained energy supply for delivering data to remote receivers. For example, in a one-day span, sunlight is weak during sunrise and sunset, intense at noon, and gone in the night.

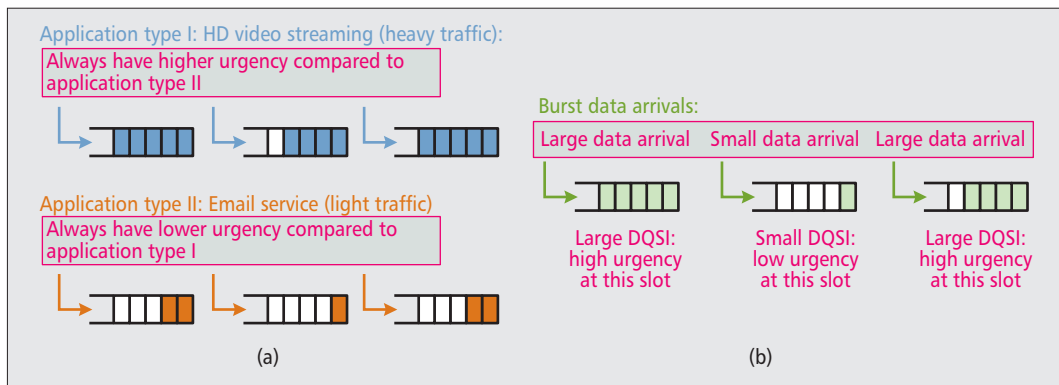


Figure 2. Illustrations of static urgency (priority) and dynamic urgency (priority).

example, it can be seen that throughput maximization for a given deadline in [5, 6] does not lead to smaller average delay performance because they have ignored the bursty (random) arrivals. Some existing works consider minimization of the average delay using the MDP approach for EH systems with finite energy storage [10]. However, the MDP problem therein is solved using numerical VIA, which gives no intuitional design insights. We propose a systematic low-complexity resource control scheme and give valuable control insights on how to dynamically control the renewable resource in EH wireless systems.

DELAY-SENSITIVE DYNAMIC RESOURCE CONTROL

MOTIVATION FOR CONSIDERING DELAY-SENSITIVE CONTROL

Since delay-sensitive applications will dominate the data stream in future wireless systems, it is necessary to design a dynamic resource control scheme that aims to provide good end-to-end delay performance. Control schemes aimed at physical layer objectives (e.g., throughput, packet error rate) cannot achieve good delay performance. Also, it is not optimal to just transmit all the data whenever there is sufficient energy. This is because one has to balance the current reward and future potential reward. We consider the following two examples with $\text{CSI} = \{\text{Good}, \text{Bad}\}$:

- **Good CSI and large DQSI:** If the current channel state is good and the data queue is long, it is wise to spend the power to empty the data queue.
- **Bad CSI and large DQSI:** If the channel state is bad but the data queue is long, it is not optimal to use all the energy to empty the entire data queue because you may not have power left for future time slots (where you may have a better CSI realization).

From this example, we can see that it is desirable to dynamically control the resource according to the instantaneous CSI, DQSI, and EQSI, where the CSI reflects the channel quality, the DQSI reflects the dynamic urgency of the data flow, and the EQSI reflects the energy availability conditions in the energy storage. We therefore propose a dynamic resource control scheme that can strike a balance among these factors.

WHY DOES DQSI CAPTURE DATA FLOW URGENCY?

Note that the urgency of a data flow dynamically changes (depending on the instantaneous queue length). Note that the average delay of a flow is given by average data queue length/average data arrival rate (according to Little's law). Hence, if the DQSI is large, it means that the average delay is likely to be large if you do not act on this. This is like going to a supermarket. If the shop manager sees a long queue at a particular counter, that counter needs to be served more urgently, he needs to act on it. In our case, if we see a large DQSI, it means the dynamic urgency or dynamic priority of the flow is higher, and even if we only have average CSI, we might still want to spend more power to transmit more data. On the other hand, if the DQSI is small (dynamic urgency is low), we probably do not want to transmit or transmit less. As such, it is not surprising that the DQSI indicates the *dynamic urgency* of a data flow. On the other hand, the urgency of a flow also depends on the application type, but that is static urgency or static priority. This is very different from the dynamic urgency of a flow as indicated by the instantaneous DQSI mentioned above. Figure 2 illustrates the static urgency (priority) w.r.t. application type and the dynamic urgency (priority) w.r.t. DQSI for a given application type.

SLOTTED EH WIRELESS SYSTEM MODEL

In this article, we consider a point-to-point EH wireless system with finite energy storage, as shown in Fig. 1. The transmitter delivers data to the remote receiver over the time-varying wireless channel $H(t)$. The transmitter is solely powered by the renewable energy source. For example, the transmitter is equipped with a solar panel that converts harvested solar energy to electricity for operation. The data to be delivered can be video streaming packets or game interactive data, and the remote end user is sensitive to delay of the arrival data flow. Specifically, if there is severe delay, the end user will experience playback interruption of high-quality video or freezing of an online game. In the transmitter, there are two buffers: the data queue buffer and energy queue buffer. The data queue buffer $Q(t)$ is for buffering the bursty traffic data flows toward the end user. The data queue is a

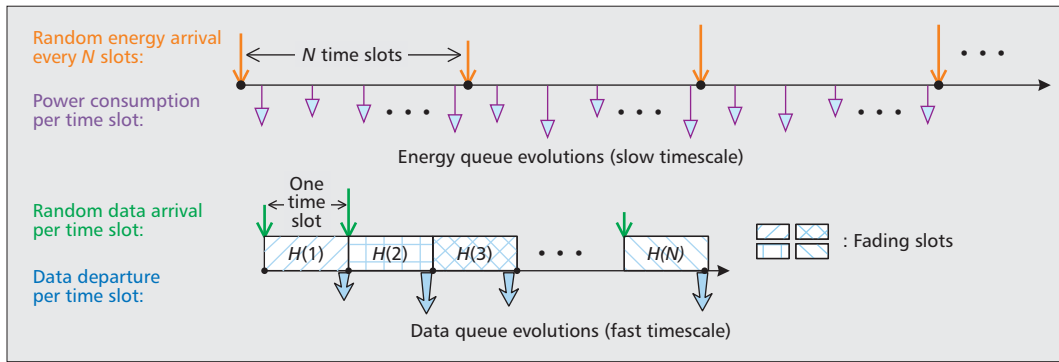


Figure 3. Illustration of evolutions for the data queue and energy queue in two timescales with slot duration. The energy queue evolves on a slow timescale, while the data queue evolves on a fast timescale.

controlled Markov chain with random data arrivals. The energy queue buffer $E(t)$ is for buffering the bursty renewable energy arrivals. The energy queue is also a controlled Markov chain with random energy arrivals. We consider that the data queue has infinite buffer size for simplicity. In practice, this may not be unreasonable because the cost of data storage is much cheaper than energy storage, and it is easy to equip wireless devices with large amounts of data storage. For example, a wireless sensor module (e.g., ConnectCore i.MX53 module) with a 1 GB external flash can store about 3×10^4 packets under a packet size of 1.5 kB in a Long Term Evolution (LTE) network for supporting video streaming. For a data network with a data rate of 10.8 Mb/s (slot duration = 10 ms), this buffer size is large enough for data storage. Therefore, it is very close to an infinite data queue. However, since high-capacity renewable energy storage is very expensive, it is presumed that renewable energy storage has finite capacity.

The EH wireless system works as a *slotted system* like most practical wireless systems. The realizations of the energy arrivals change every N consecutive time slots, and those of the data arrivals change once per time slot. Figure 3 illustrates the arrival processes of the energy and data arrivals, as well as the evolutions of the data and energy queues.

DYNAMIC RESOURCE CONTROL PROBLEM STATEMENT

For delay-sensitive applications, it is important to dynamically control the communication resource according to the CSI (captures the transmission opportunities), DQSI (captures the data urgency), and EQSI (captures the energy availability). Therefore, the global system state is characterized by the CSI, DQSI, and EQSI in the EH wireless system. The objective of the delay-sensitive control optimization is to find a power control policy¹ that minimizes the infinite horizon average delay as follows:

Delay-Sensitive Power Control Optimization:

$$\begin{array}{ll} \text{minimize} & \text{average delay} \\ \text{subject to} & \text{energy availability constraint} \end{array}$$

For delay-sensitive applications, it is important to dynamically control the communication resource according to the CSI, DQSI, and EQSI. Therefore, the global system state is characterized by the CSI, DQSI, and EQSI in the EH wireless system.

The control policy should satisfy the *energy availability constraint*. This means that the energy consumption at each time slot cannot exceed the current available energy in the renewable energy storage.

KEY TECHNICAL CHALLENGES

The delay-sensitive power control problem in the EH wireless system is an infinite dimensional MDP. There are several challenges associated with the MDP problem.

Challenges due to queue-dependent control: Control policies adaptive to the data and energy queue evolutions are quite challenging because the underlying problem embraces *information theory* (to model the channel dynamics) and *queueing theory* (to model the queue dynamics).

Complex coupling between the data queue and energy queue: The data rate of the transmitter in the EH wireless network depends on the current available energy in the energy storage. As such, the dynamics of the data and energy queues are coupled together. This further makes the MDP problem a *coupled multi-dimensional* stochastic optimization problem.

Challenges due to the random and bursty nature of the renewable energy source: In the previous literature on EH wireless systems, the bursty energy arrivals are modeled as independent and identically distributed (i.i.d.) random processes for analytical tractability. In practice, most of the renewable energy arrivals are not i.i.d., and such a *non-i.i.d. nature* will have a huge impact on the dimensioning of the energy storage capacity.

LOW-COMPLEXITY DELAY-AWARE RESOURCE CONTROL SCHEME

In this section, we propose an asymptotically optimal low-complexity power control solution for solving the delay-sensitive MDP problem.

DYNAMIC POWER CONTROL SOLUTION

The delay-sensitive power control problem in EH wireless systems is an infinite horizon average cost MDP. Using the divide-and-conquer principle, the infinite dimensional MDP is transformed into a per-stage optimization problem

¹ A power control policy is a mapping from the global system state to the power control action of the transmitter.

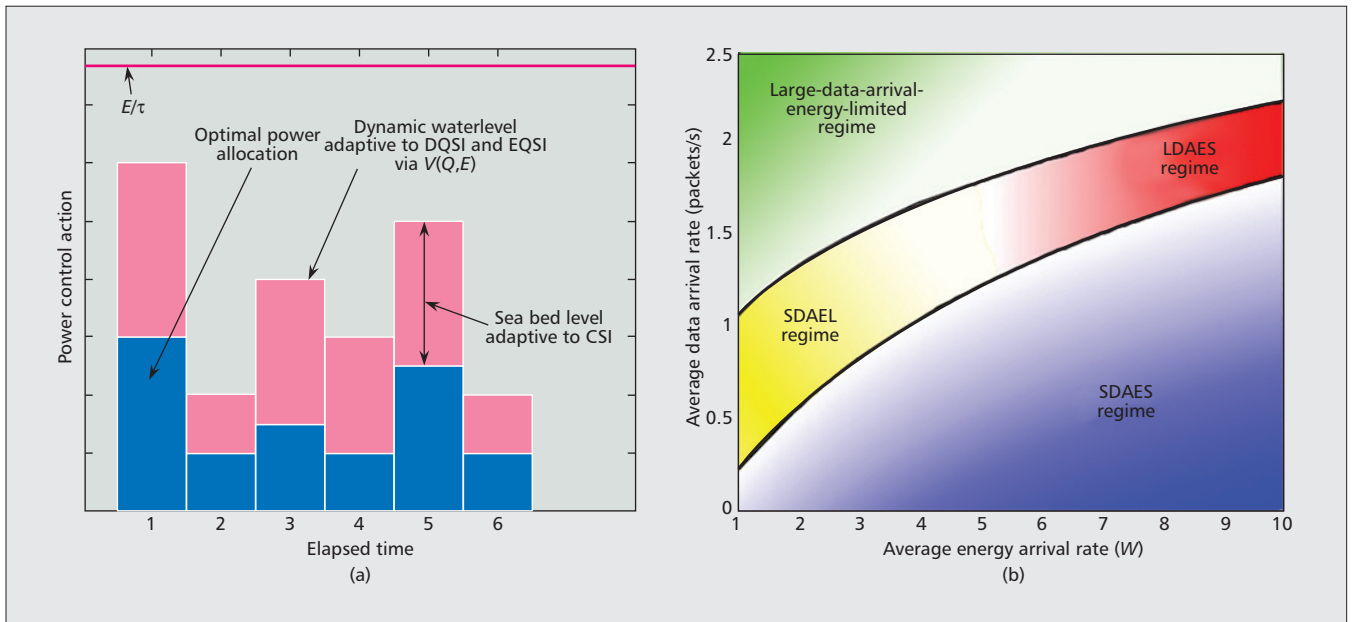


Figure 4. a) An example of the dynamic multi-level water-filling structure; b) asymptotic operating regimes. The two curves represent log-increasing functions w.r.t. the average energy arrival rate. The green area is less interesting because it corresponds to the heavy traffic regime for the data queue, and the delay will be large no matter what control policy is adopted (since the traffic loading of the system is at the instability margin).

given by the *Bellman equation* [3]. The optimal power control solution of the Bellman equation has the following structure:

$$\min \left\{ \left(\underbrace{f(V(Q,E))}_{\text{dynamic water level}} - \frac{1}{H^2} \right)^+ \frac{E}{\tau} \right\},$$

where $V(Q, E)$ is called the *priority function* that captures the dynamic priority of the data flow under different system state realizations w.r.t. DQSI and EQSI.

It can be observed that the optimal power control solution is very similar to the classical water-filling solution, but with a *dynamic multi-level water-filling* structure as shown in Fig. 4a. Furthermore, the optimal power control solution depends on the instantaneous CSI, DQSI, and EQSI via priority function $V(Q, E)$, which captures how the DQSI and EQSI affect the *overall priority* of the data flow. Furthermore, the optimal power control also depends on the current available energy E in the energy buffer.

CLOSED-FORM APPROXIMATION OF PRIORITY FUNCTION

The optimal power control solution depends on the priority function $V(Q, E)$. In the following, we focus on obtaining an *analytical expression* of the priority function. Note that obtaining the priority function is equivalent to solving a system of nonlinear fixed point equations (i.e., the Bellman equation). Classical VIA can only give numerical solutions, and suffers from slow convergence issues and lack of design insights. To overcome this challenge, we shall adopt a continuous time perturbation (CTP) approach [11] so

as to obtain closed-form solutions and low-complexity control schemes, and discuss valuable design insights.

Using the CTP approach, we obtain the closed-form approximate priority function under the following asymptotic regimes, illustrated in Fig. 4b.

Large-data-arrival-energy-sufficient (LDAES) regime: In this regime (red area in Fig. 4b), we have a large average data arrival rate and a large average energy arrival rate. This regime corresponds to the scenario where the wireless device has sufficient renewable energy supply for the energy queue to combat the heavy data traffic for the data queue.

Small-data-arrival-energy-limited (SDAEL) regime: In this regime (yellow area in Fig. 4b), we have a small average data arrival rate and a small average energy arrival rate. This regime corresponds to the scenario where we have insufficient energy supply for the energy queue, but the data traffic for the data queue is light.

Small-data-arrival-energy-sufficient (SDAES) regime: In this regime (blue area in Fig. 4b), we have a small average data arrival rate and a large average energy arrival rate. This regime corresponds to the scenario where we have a very sufficient renewable energy supply in the energy queue to keep the data queue stable and maintain a low data queue length.

The overall solution is asymptotically optimal when the slot duration is small. Please refer to [11] for the detailed derivations on this CTP approach. Furthermore, the small slot duration condition is easily satisfied in practical wireless systems. For example, in LTE, the physical layer is organized into radio frames (corresponding to the slots in our system), and the generic radio frame has a time duration of 10 ms.

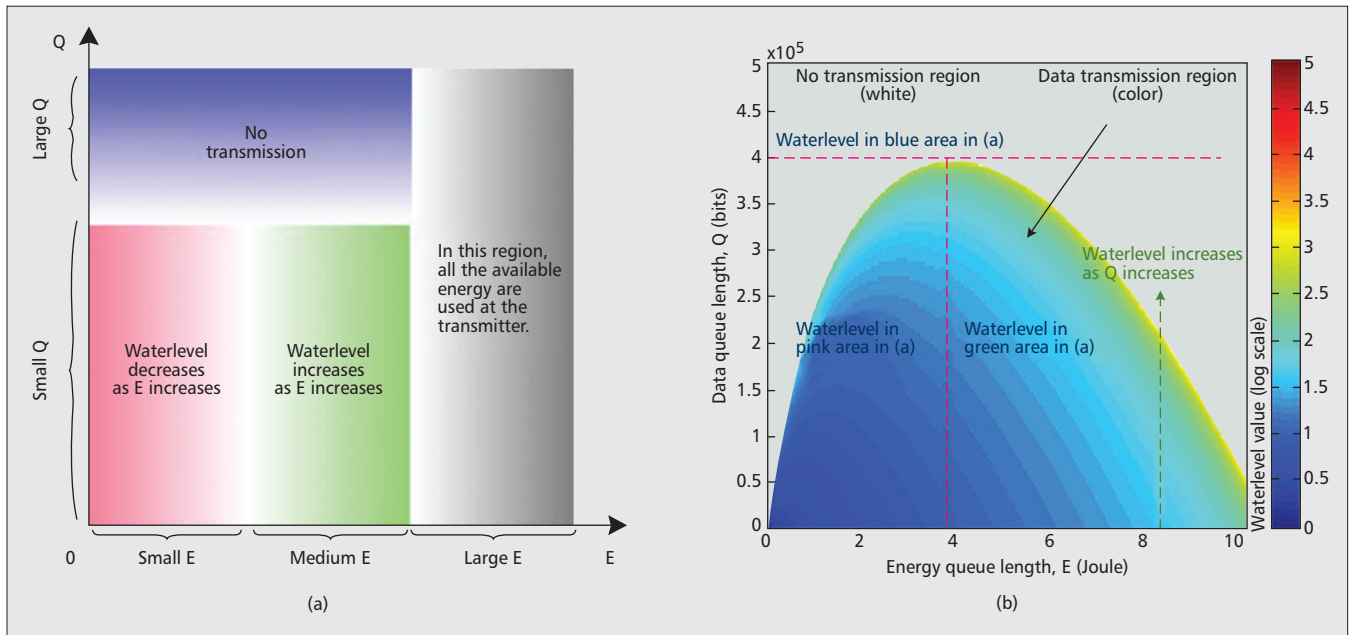


Figure 5. a) Decision region partitioning w.r.t. (Q, E) ; b) an example of the water level for LDAES and SDAEL regimes.

CONTROL INSIGHTS UNDER ASYMPTOTIC OPERATING REGIMES

Control insights under LDAES and SDAEL regimes: For both the LDAES and SDAEL regimes, the priority function has different closed-form properties in each of the four areas in Fig. 5a. Figure 5b illustrates an example of the water level vs. the DQSI and EQSI for small and medium E . Specifically:

- **Small and medium E , large Q (blue area of Fig. 5a):** We do not use any renewable energy to transmit data. The reason is that even though we can use the limited energy for data transmission, the data queue length will not decrease significantly, which contributes very little to the delay performance. Instead, if we do not use the energy at the current slot, we can save it and wait for the future good transmission opportunities.
- **Small E , small Q (pink area of Fig. 5a):** We use the available energy for transmission, and the water level is increasing w.r.t. Q , which is in accordance with the high urgency of the data flow. Furthermore, large E leads to a lower water level. This is reasonable because it is appropriate that for small E , we can save some energy in the current slot for better transmission opportunities in future slots.
- **Medium E , small Q (green area of Fig. 5a):** We use the available energy for transmission, and the water level is increasing w.r.t. Q . Furthermore, large E leads to a higher water level because we have sufficient available energy, and it is appropriate to use more power to decrease the data queue.
- **Large E (grey area of Fig. 5a):** The transmitter uses all the available energy to make room for future energy arrivals.

Solution and control insights under SDAES regime: For the SDAES regime, based on the priority function, the optimal power control

solution is to use all the available energy in the energy buffer at each time slot. This is reasonable because in this regime, there is plenty of renewable energy, and it is sufficient to use all the available energy to support the data traffic and maintain the data queue stability.

STABILITY CONDITIONS

In this part, we discuss the conditions that ensure the stability of EH wireless systems (i.e., the stability of the data queue).

Condition on the energy harvesting capability: We require that the average energy arrival rate be at least at exponential order of the average data arrival. This means that for given average data arrival rate, if the EH rate is too small, even if we use all the available energy in the energy buffer at each time slot, the data queue cannot be stabilized.

Condition on the energy storage capacity: The capacity of the energy storage should be at least at a similar order of $N \times$ average energy arrival per time slot. The condition gives a *first order design guideline* on the dimensioning of the energy storage capacity. This condition ensures that the energy storage at the transmitter has sufficient stored energy to support data transmission for N slots when energy arrivals are small.

PERFORMANCE EVALUATION

We compare our proposed low-complexity power control solution with the following baselines:

- **Baseline 1, greedy strategy (GS) [10]:** Full power is always used.
- **Baseline 2, CSI-only water-filling strategy (COWFS) [10]:** Classical CSI-dependent water-filling power control is used by maximizing the ergodic channel capacity.
- **Baseline 3, DQSI-weighted water-filling strategy (DWWFS) [12]:** CSI- and DQSI-dependent water-filling power control is

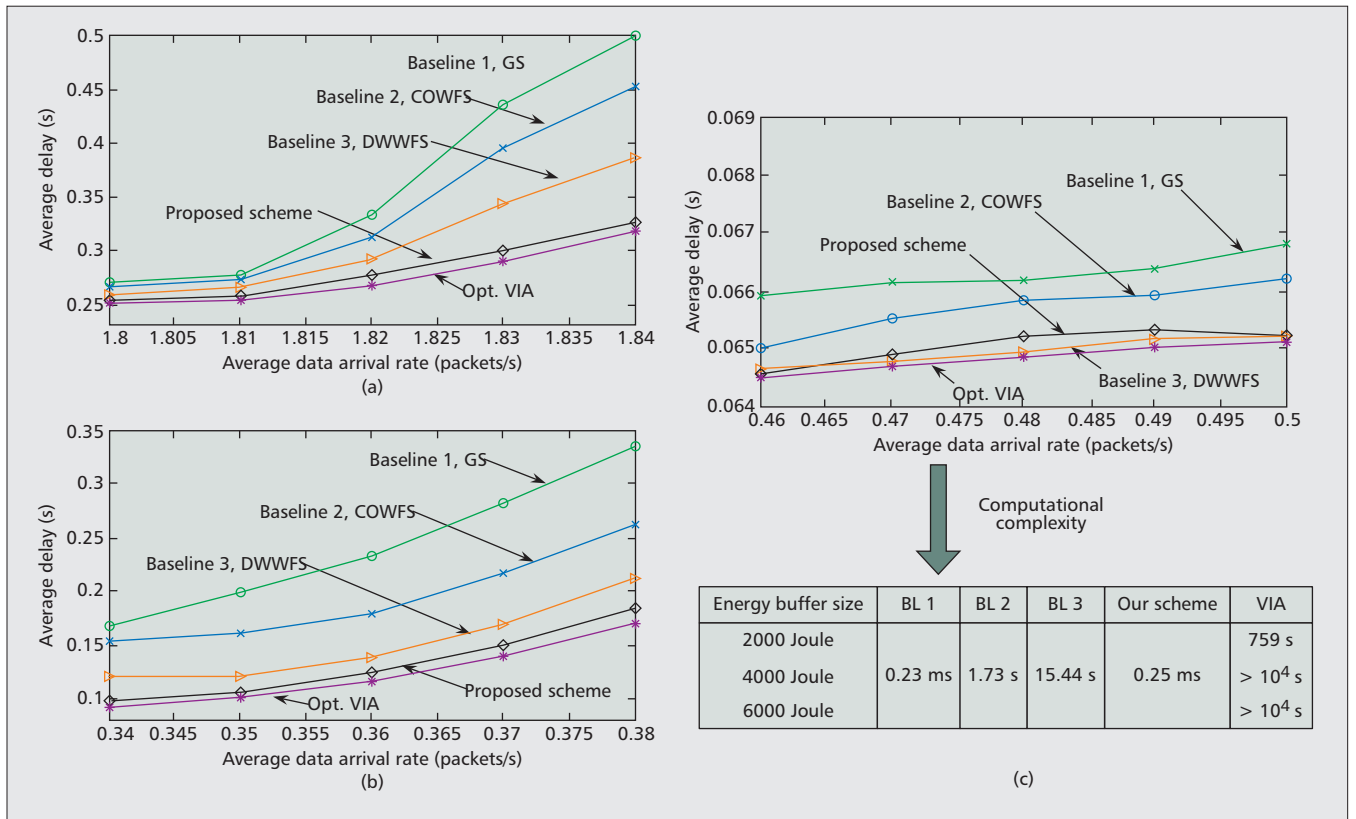


Figure 6. a) Delay performance under the LDAES regime; b) delay performance under the SDAEL regime; c) delay performance under the SDAES regime. Comparison of the MATLAB computational time under the SDAES regime.

used by maximizing the queue weighted ergodic capacity.

In the performance evaluation, the base station (BS) is equipped with 5-cascaded $2\text{ m} \times 2\text{ m}$ solar panels with energy harvesting performance $1\sim 10\text{ mW/cm}^2$. If the surrounding environment has sufficient sunlight, the EH performance is high. Otherwise, the EH performance is low [13]. The noise spectral density is -174 dBm/Hz , path loss is 160 dB , bandwidth is 20 MHz , and slot duration is 10 ms . The energy arrival rate at the BS changes every 5 min , and the renewable energy is stored in a 1.2 V 2000 mAh lithium-ion battery. Under this practical scenario, the BS can provide a data rate of 10.8 Mb/s for supporting delay-sensitive application of the mobile user.

Figure 6 illustrates the average delay vs. the average data arrival rate for the LDAES, SDAEL, and SDAES regimes. The proposed scheme achieves significant performance gain over all the baselines. The gain is contributed by the *DQSI- and EQSI-aware* dynamic water-filling structure. It can also be observed that the performance of the proposed low-complexity solution is very close to the optimal VIA [3]. Figure 6c also illustrates the comparison of the MATLAB computational time of the proposed solution, the baselines, and the brute-force VIA [3]. Note that the proposed scheme has similar complexity to baseline 1 due to the closed-form priority function. Therefore, our proposed scheme achieves significant performance gain with negligible computational cost.

CONCLUSION

Dynamic delay-sensitive power control for EH wireless systems with finite energy storage is challenging. This article surveys some of the existing works on resource control in EH wireless systems. To address the dynamic control challenge, we use a CTP approach to propose a low-complexity power control solution, which is adaptive to the CSI, DQSI, and EQSI. We also give valuable design insights on the dynamic control structures. Numerical results show that the proposed power control scheme has much better performance than the state-of-the-art baselines.

REFERENCES

- [1] D. Niyato *et al.*, "Wireless Sensor Networks with Energy Harvesting Technologies: A Game-Theoretic Approach to Optimal Energy Management," *IEEE Commun. Mag.*, vol. 14, no. 4, Sept. 2007, pp. 90–96.
- [2] Y. Cui *et al.*, "A Survey on Delay-Aware Resource Control for Wireless Systems — Large Deviation Theory, Stochastic Lyapunov Drift and Distributed Stochastic Learning," *IEEE Trans. Info. Theory*, vol. 58, no. 3, Mar. 2012, pp. 1677–700.
- [3] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, 3rd ed. Boston, MA: Athena Scientific, 2005.
- [4] J. Yang and S. Ulukus, "Optimal Packet Scheduling in an Energy Harvesting Communication System," *IEEE Trans. Commun.*, vol. 60, Jan. 2012, pp. 220–30.
- [5] K. Tutuncuoglu and A. Yener, "Sum-Rate Optimal Power Policies for Energy Harvesting Transmitters in An Interference Channel," *J. Commun. Net.*, vol. 14, no. 2, Apr. 2012, pp. 151–61.
- [6] K. Tutuncuoglu and A. Yener, "Optimum Transmission Policies for Battery Limited Energy Harvesting Nodes," *IEEE Trans. Wirelless Commun.*, vol. 11, no. 3, Mar. 2012, pp. 1180–89.

-
- [7] M. Gorlatova, A. Wallwater, and G. Zussman, "Networking Low-Power Energy Harvesting Devices: Measurements and Algorithms," *IEEE Trans. Mobile Comp.*, vol. 12, no. 9, Sept. 2013, pp. 1853–65.
- [8] A. Seyedi and B. Sikdar, "Energy Efficient Transmission Strategies for Body Sensor Networks with Energy Harvesting," *IEEE Trans. Commun.*, vol. 58, no. 7, July 2010, pp. 2116–26.
- [9] H. Li, N. Jaggi, and B. Sikdar, "Relay Scheduling for Cooperative Communications in Sensor Networks with Energy Harvesting," *IEEE Trans. Commun.*, vol. 10, no. 9, Sept. 2011, pp. 2918–28.
- [10] V. Sharma *et al.*, "Optimal Energy Management Policies for Energy Harvesting Sensor Nodes," *IEEE Trans. Wireless Commun.*, vol. 9, no. 4, Apr. 2010, pp. 1326–36.
- [11] F. Zhang and V. K. N. Lau, "Cross-Layer MIMO Transceiver Optimization for Multimedia Streaming In Interference Networks," *IEEE Trans. Signal Proc.*, vol. 62, no. 5, Mar. 2014, pp. 1235–44.
- [12] L. Huang and M. J. Neely, "Utility Optimal Scheduling in Energy Harvesting Networks," *IEEE/ACM Trans. Net.*, vol. 21, no. 4, Dec. 2011, pp. 1117–30.
- [13] C. Park and P. H. Chou, "Ambimax: Autonomous Energy Harvesting Platform for Multi-Supply Wireless Sensor Nodes," *Proc. IEEE SECON*, Sept. 2006, pp. 168–77.

BIOGRAPHIES

FAN ZHANG [SM'10] received his B.Eng. (First Class Hons) from Chu Kochen Honors College at Zhejiang University in 2010. He is currently pursuing a Ph.D. degree in the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology (HKUST). His research interests include cross-layer delay-sensitive resource allocation, and low-complexity stochastic optimization for wireless communication and control systems.

VINCENT K. N. LAU [F'11] obtained his B.Eng. (Distinction 1st Hons) from the University of Hong Kong in 1992 and his Ph.D. from Cambridge University in 1997. He joined Bell Labs from 1997 to 2004 and the Department of ECE, HKUST in 2004. He is currently a professor and the founding director of Huawei-HKUST Joint Innovation Lab at HKUST. His current research focus includes delay-optimal cross layer optimization, massive MIMO, interference mitigation and networked control systems.

Toward Secure Energy Harvesting Cooperative Networks

Jiawen Kang, Rong Yu, Sabita Maharjan, Yan Zhang, Xumin Huang, Shengli Xie, Hanna Bogucka, and Stein Gjessing

ABSTRACT

The concept of energy harvesting cooperative networks is an emerging technology that has very high potential for a large variety of applications. However, energy transfer capability may lead to unprecedented security challenges. In this article, we study energy security issues and the solutions in energy harvesting networks. We first identify typical energy related attacks and then propose defense solutions against these attacks. We also carry out security analysis and performance analysis to evaluate our proposed solutions. Simulation results have shown that the proposed defense solutions are effective and efficient.

INTRODUCTION

With the rapid development of wireless technologies, people with different mobile devices can receive various information and services. However, as power in devices is typically supplied by pre-charged batteries, the terminals have limited energy. Recently, a new concept called energy harvesting networks has emerged, which offers high potential to address this problem. In energy harvesting networks, wireless devices operate perpetually by using ambient energy, such as solar, vibration, and thermoelectric effects. This will significantly extend network lifetime and decrease maintenance costs. The intermittent nature of energy sources calls for energy management policies about data transmission.

Recently, several studies have investigated the energy management problems in different scenarios of energy-harvesting-based communication [1, 2]. In [1], the notion of energy cooperation is proposed, considering which nodes can share a portion of their harvested energy with others through a separate wireless energy transfer unit.

Unlike the classical definition of cooperation in information transmission, the concept of energy cooperation mainly focuses on the energy level. Since both energy and information cooperation can be exploited, an energy harvesting cooperative network is attracting much attention, especially in the context of optimizing transmission policies to use the available energy

[3]. However, the security issues in energy harvesting cooperative networks have been largely ignored. Existing work mainly focuses on security issues in energy harvesting wireless sensor networks. In [4], the authors present a taxonomy of attacks in energy harvesting wireless sensor networks. The study in [5] allows the system to change communication security settings at runtime with the goal of improving node lifetime while providing a suitable security level. In [6], the adaptability of some well-known cryptography algorithms to energy harvesting wireless sensor networks is discussed.

In this article, we focus on the unique security issues and challenges in energy harvesting cooperative networks. Unlike traditional wireless networks, nodes in energy harvesting cooperative communication not only transmit information, but also transmit energy to the cooperators. The major security issues about traditional information transmission, including key establishment, authentication, privacy, and robustness to denial-of-service attacks, have been thoroughly investigated [7]. The information can be encrypted to ensure confidentiality, and digital signatures can ensure information integrity and authentication. However, the security issues about energy are quite different. Energy cannot be encrypted to achieve confidentiality. There is no digital signature to guarantee energy integrity and authentication. Once an adversary receives energy from others, it can use the energy. Thus, the energy is in danger of being intercepted and consumed during cooperation. Therefore, energy cooperation transmission poses unexplored security challenges.

In this article, we first explore several typical energy attacks and challenges in energy harvesting cooperative networks. Then we propose the solutions to address these distinctive security issues. Finally, we conclude the article. The contributions of this article are as follows:

- We summarize the special vulnerabilities and security requirements of energy cooperation networks.
- We introduce four typical energy attacks in energy harvesting cooperative networks, and provide some examples to illustrate them.
- We address these unique security challenges via three fundamental defense solutions.

Jiawen Kang, Rong Yu, Xumin Huang, and Shengli Xie are with Guangdong University of Technology.

Sabita Maharjan and Yan Zhang are with Simula Research Laboratory.

Stein Gjessing is with the University of Oslo.

Hanna Bogucka is with Poznan University of Technology.

ENERGY COOPERATION AND ITS VULNERABILITIES IN ENERGY ATTACKS

ENERGY COOPERATION

Figure 1 illustrates an energy harvesting cooperative network, which consists of three types of communication nodes: *source*, *cooperator* and *destination*. All nodes in the model are able to harvest energy from surrounding environment. The amount of current energy in the battery of a node is called *energy state*. The source node (S) evaluates its current energy state, and operates collaboratively with the cooperator node i (C_i) to transmit information to the destination node (D). There are two main components in the communication nodes as follows.

- **Energy harvesting module.** Communication nodes, which can harvest energy from the environment, have the potential to operate beyond the timeframe limited by the finite capacity of their batteries. The harvested energy is stored in the energy buffer.
- **Information transmission module.** The information transmission module generates information, and uses the energy in the energy buffer to transmit information to the destination with the help of the transmitter.

In energy harvesting cooperative networks, there exists three energy cooperation models illustrated in Fig. 2:

- Cooperative relay with energy prepayment
- Cooperative relay with energy debt
- Direct transmission with energy debt [8]

Cooperative Relay with Energy Prepayment

— As shown in Fig. 2a, in cooperative relay with energy prepayment (CREP), the source node works with a cooperator to finish the information transmission. Here, the cooperator works as a *relay* to transmit information. After consulting with the relay in terms of demanded energy for transmission, the source will transfer the demanded energy and then transmit information to the relay when the source has sufficient energy. The relay transmits the information to the destination using transmission energy given by the source after receiving the energy and the information. During transmission, it does not consume any of its own energy. Meanwhile, the source sends a notice to the destination, which alerts the destination to get ready to receive information from the specific relay. The destination sends a confirm message back to the source after receiving the information from the specific relay, which indicates the end of cooperative transmission.

Cooperative Relay with Energy Debt

— As shown in Fig. 2b, in cooperative relay with energy debt (CRED), the source works with the cooperator to transmit information. In this case, the cooperator works as a relay but lends its energy to the source. There is a consultation of demanded energy for transmission as in CREP between the source and the cooperator. If the energy in the source is available to transmit the information, but not enough to provide the demanded energy to the relay, it first transmits

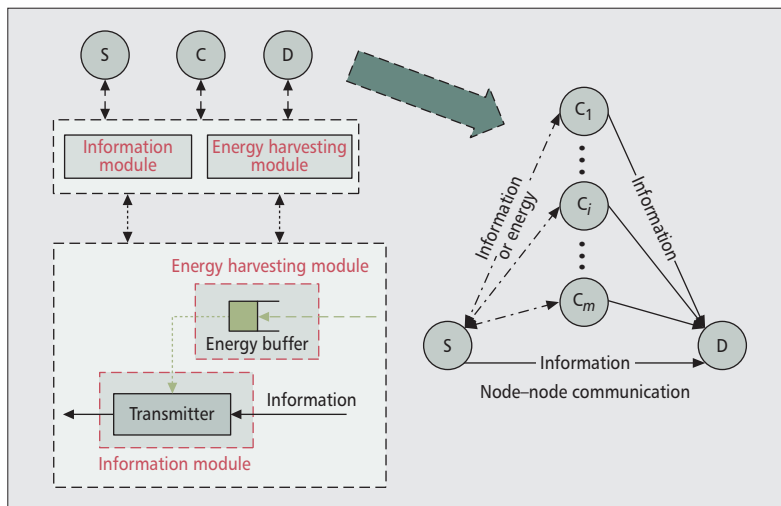


Figure 1. Network model.

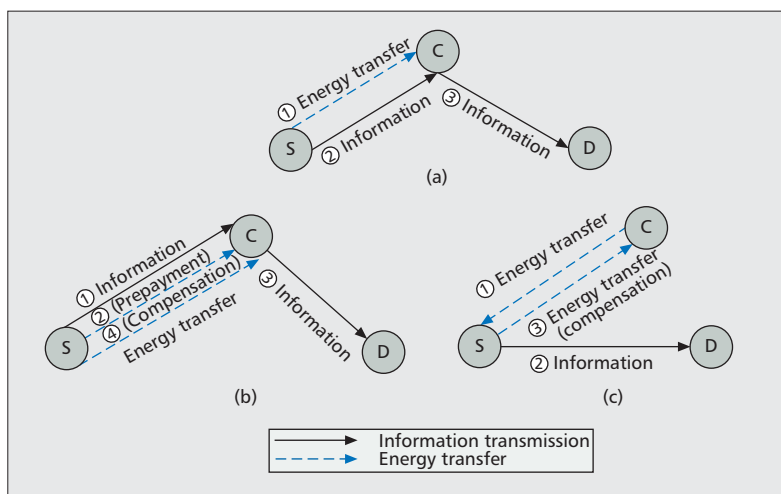


Figure 2. Energy cooperation models: a) cooperative relay with energy prepayment (CREP); b) cooperative relay with energy debt (CRED); c) direct transmission with energy debt (DTED).

the information and the available energy to the relay. Then the source requests the relay to lend a certain amount of energy for transmission, and they jointly generate an energy debt bill (EDB) for each other. The source stores the energy debt bill as an energy debt. The relay records the energy debt bill from the source and transmits the information to the destination via the energy from the source node and a portion of its own available energy in the buffer.

When the source node is idle, it harvests energy from surroundings. If there are several energy debt records at the source node, the source node compensates energy to the relay in the order of generated EDBs. In other words, the earliest energy debt bills should be the first compensated ones.

Direct Transmission with Energy Debt

— As shown in Fig. 2c, in direct transmission with energy debt (DTED), if a source is near the destination and does not have enough energy to transmit information, it will ask a nearby cooper-

Potential attacks against energy harvesting cooperative networks could be broadly considered from two different views. One is the attack against the energy, and another is against the information. As the communication nodes take different roles in different energy cooperation models, the behaviors of attacks vary widely according to different models.

ator for help. In this case, the cooperator works as an energy renter to rent a certain amount of energy to the source. They make a deal about the demanded energy, and then the source generates an EDB to the renter, which is the same as in CRED. The renter verifies the EDB from the source and transfers its stored energy by wireless link to support the source's transmission. In this way, the source transmits the information by itself instead of transmitting through a relay. The process of energy compensation is the same as that in CRED.

ENERGY COOPERATION VULNERABILITY

Energy has a prominent feature: once an adversary receives energy from others, it can directly use it. As a result, there are new security challenges during the process of cooperative energy transfer in energy harvesting cooperative networks. Energy attacks come from several vulnerabilities associated with energy cooperation as follows:

- *No encryption for energy.* Unlike traditional information transmission, energy cannot be encrypted to ensure that it only belongs to the entitled nodes. Any adversary can harvest it by illegal means and then take possession of it, which results in energy cooperation cheating and energy repudiation of reception.
- *No backup for energy.* The transmitted information can be backed up and recovered by duplication, while energy cannot be. This special property induces energy transmission cheating and repudiation of reception during transfer.
- *Unpublished energy state information.* The energy state information of the nodes is in danger of being attacked. This is because a malicious node may forge or alter its energy state information, leading to an energy state forgery attack.

These vulnerabilities may not only bring serious energy loss to the victim, but may also cause delay in information transmission or packets being dropped out. What is worse, it may lead to system crash because of highly redundant information.

SECURITY REQUIREMENTS FOR ENERGY COOPERATION

To preserve nodes against different types of adversaries and overcome the aforementioned vulnerabilities, the following security requirements should be satisfied.

- *Energy integrity.* During the energy transfer process, the energy should be divided into energy packet units to reduce energy loss or damage in harsh communication environments.
- *Energy non-frameability.* A trusted node cannot perform actions or, more generally, prove that a node performed the malicious actions when the node never did so.
- *Energy state credibility.* The node cannot tamper with or conceal its true energy state for illegal interests.
- *Energy usability.* As additional computation consumes additional energy, simplicity and reliability of management operations (e.g.,

encryption and decryption, maintenance) are also crucial. Traditional encryption algorithms need to be adjusted to fit the wireless network for reducing extra energy cost.

The security requirements of information in this article are similar to traditional information requirements, which include information confidentiality, authentication and integrity, availability, and so on. The requirement for information non-repudiation is related to a specific energy attack in our system. The special energy attack means that one node cannot collude with another node to frame a case against the third node.

TYPICAL ENERGY ATTACKS IN ENERGY HARVESTING COOPERATIVE NETWORKS

Potential attacks against energy harvesting cooperative networks may broadly be considered from two different views. One is the attack against energy, and another is against information. As the communication nodes take different roles in different energy cooperation models, the behaviors of attacks vary widely according to different models.

For energy attacks, we focus on three typical energy attacks. For information attacks, we only consider energy-related information attacks in this article as traditional information attacks have been extensively studied. As the impacts of energy attacks are similar in different cooperation models, we provide three kinds of typical attacks rather than describing cooperation-model-specific actions of attacks in this section. For each kind of attack, we provide some examples to illustrate related adversarial behaviors.

ENERGY STATE FORGERY

Energy state forgery means that an attacker counterfeits its energy state and sends the state information to a cooperator for illegal purposes during the nodes' communication. There are three main scenarios concerning this attack:

- Scenario one: Victim j sends a cooperation request to attacker i (a selfish cooperator). Attacker i , with sufficient energy, misrepresents that it has insufficient energy, and then refuses to help.
- Scenario two: An attacker i (a selfish source) with sufficient energy deceives that it has no or insufficient energy. i sends a cooperation request to ask for help. Victim j may be cheated and may establish a cooperation relationship with i .
- Scenario three: When victim j sends a cooperation request to attacker i (a selfish cooperator), i , without sufficient energy, feigns having sufficient energy. Then j establishes a cooperation relationship with i .

Next, we take some cases in different energy cooperation models as examples to illustrate this attack. In the CREP model shown in Fig. 2a, the source (attacker) with enough energy may misrepresent that its energy is not enough to execute CREP. Then it executes the CRED model. The source frequently sends requests to the cooperator (victim) for help.

Energy cooperation models	Attacker	Attack Types				Impact on Victim		Probability of occurrence
		Energy state forgery (ESF)	Energy cooperation cheating (ECC)	Energy repudiation of reception (ERR)	Energy-related information attack (ERIA)	Energy loss	Information lag or loss	
CREP	S	S→C ¹	—	S→C	—	ESF, ERR	ESF, ERR	Low
	C	C→S	C→S	C→S	C→S	ESF, ECC, ERR	ESF, ERR, ERIA	High
	D	—	—	—	—	—	—	—
	S+D	—	—	—	S+D→C ²	ERIA	—	Moderate
CRED	S	S→C	—	S→C	—	ESF, ERR	—	Moderate
	C	C→S	C→S	C→S	C→S	ECC, ER, ERIA	ESF, ERIA	High
	D	—	—	—	D→C	ERIA	—	Low
	S+D	—	—	S+D→C	S+D→C	ERR, ERIA	—	High
DTED	S	S→C	S→C	S→C	—	ESF, ECC, ERR	—	High
	C	C→S	—	C→S	—	ERR	ESF	Moderate
	D	—	—	—	D→S	ERIA	—	Low
	S+D	—	—	S+D→C	—	ERR	—	Low

¹ A→B: A attacks B; ² A+B→C: A colludes with B to attack C.

Table 1. Energy attacks and energy-related information attacks in different energy cooperation models.

In the CRED model shown in Fig. 2b, a cooperator without enough energy may pretend that it has sufficient energy to provide help to the source. Actually, the attacker is not able to transmit information to the destination.

For a selfish cooperator, in the DTED model of Fig. 2c, the cooperator misrepresents itself as idle now, and refuses to cooperate with the victim (source). The source thus gets no help from the cooperator, and the destination has to wait for a long time to receive the information.

ENERGY COOPERATION CHEATING

For energy cooperation cheating, an attacker (malicious source or cooperator) is motivated by its own benefit to send bogus messages or take fraudulent action to obtain trust from a victim. After establishing a cooperation relationship with the victim, the attacker requests that it transfers energy to the attacker, which leads to the loss of victim's energy. For example, in the DTED model, a victim (cooperator) loses energy because a malicious source broadcasts an energy request to the cooperator and hordes the energy by not transmitting the information to the destination. In the CREP model, the malicious cooperator (attacker) pretends to agree to information transmission, and then steals energy from the source, resulting in energy loss.

ENERGY REPUDIATION OF RECEPTION

Energy repudiation of reception refers to a situation where the attacker disavows reception of energy after receiving energy from the cooperator. An attacker may also claim that it has not received enough preconcerted energy from the cooperator when it borrows energy from others. Then the attacker refuses to return energy or provide help to the cooperator.

For example, in the DTED model, a cooperator (victim) may have transferred sufficient energy to a malicious source (attacker). But the source may claim that it has only received part of the preconcerted energy from the cooperator. The cooperator thus suffers from energy loss in this case. In the CRED model, after an attacker (malicious cooperator) transmits the information to the destination and gets energy compensation from the source, it may deny that it has any compensated energy, or has sufficient energy from the borrower.

ENERGY-RELATED INFORMATION ATTACKS

Next, we consider energy-related information attacks, which are special information attacks closely related to energy in our system.

As the process of energy transfer only occurs in the source or the cooperator, the energy attacks are initiated by one of them. However,

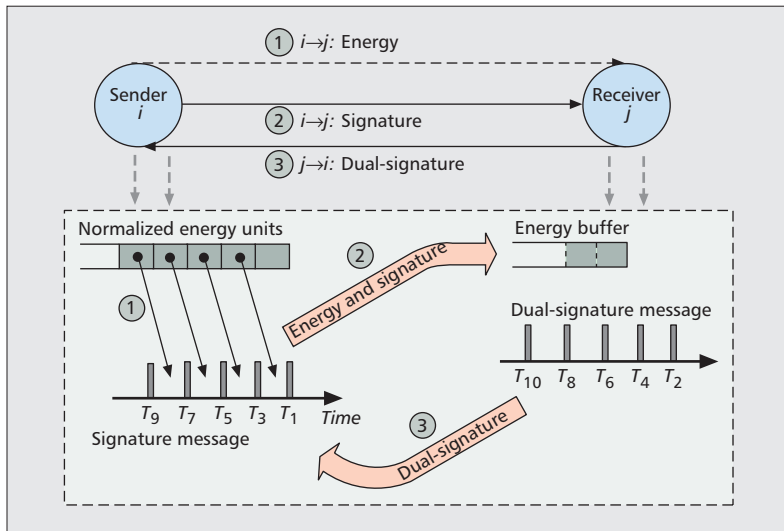


Figure 3. Energy dual-signature.

energy-related information attacks are initiated with nodes collusion. The most typical energy-related information attacks is the source-destination (S-D) information denial, which refers to the condition when the source and the destination collude to defame the cooperator. For example, in the CREP and CRED models, the source first colludes with the destination. Then they frame a case against the cooperator by claiming that the destination has not received any information from the cooperator. The cooperator has to transmit the information again, and thus suffers energy loss.

Table 1 shows various attacks with main properties in different energy cooperation models for energy harvesting cooperative networks.

ENERGY DEFENSE SOLUTIONS FOR ENERGY HARVESTING COOPERATIVE NETWORKS

SECURITY MODEL FOR DEFENSE

The security model for defense mainly consists of the aforementioned communication nodes, registration authority, and trust modules:

- **Registration authority:** The trusted registration authority provides authentication and authorization services to nodes. Each node registers to the registration authority. After authentication by the registration authority, nodes obtain public and private keys with corresponding certificates, which are stored in trust modules of the nodes.
- **Trust modules:** Each node is equipped with a trust module (i.e., built-in hardware and firmware), which has two functionality types:
 - Cryptographic operations and storage, and
 - Specific messages (e.g., current energy state) record and transmission, with contents that are trusted

Only the registration authority is allowed to read or modify the information stored in a trust module.

To defend potential energy attacks in energy harvesting cooperative networks, we devise three fundamental defense solutions based on the above security model:

- Periodic energy state report (ESR)
- Energy dual-signature
- Symmetric and asymmetric encryption

For potential energy attacks, these solutions work together to defend attacks and satisfy the security requirements discussed earlier.

Periodic Energy State Report — In this article, the ESR refers to the current energy state in the energy buffer of a node, which is recorded and transmitted to the registration authority by trust modules. The common format of a report includes a node's identity, geographic location, energy state in the energy buffer, percentage of energy change since last report, and timestamp and signature of the ESR. For privacy preservation, identity of the node can be protected by pseudonyms.

We consider the ESR to be periodically recorded and transmitted to the registration authority every few hundred milliseconds. In order to defend energy state forgery, a legitimate node can request that the registration authority verify the energy states of cooperation nodes. The registration authority gives back a reply (“true” or “false”) without specific state values to the applicant because of privacy preservation rules.

Energy Dual-Signature — In this article, energy dual-signature is used to validate and authenticate that a sender has transferred sufficient energy to a receiver with confirmation from the receiver. As shown in Fig. 3, there are three steps in the process of energy dual-signature. In the first step, the transferred energy of sender i (source or cooperator) is divided into a number of normalized energy units before transferring. Sender i generates a set of digital signatures matching every normalized energy unit. The digital signatures contain the following information: identity of the sender and receiver, event record of cooperation, timestamp, digital digest of the event record, total amount of transferred energy, and serial number of current transferring energy. Sender i arranges each energy unit and corresponding digital signature in an energy queue by time order. Then sender i transfers an energy unit and sends a corresponding digital signature message about this energy unit to receiver j at time T_1 . Receiver j verifies the signature message from i and confirms the transfer event after receiving the energy unit. Receiver j will generate a dual-signature message based on the signature of T_1 from i and send it back to i as a reply at T_2 . After sender i verifies the dual-signature, it continues transferring the next energy units to j at T_3 until the energy transfer is complete.

Symmetric and Asymmetric Encryption — To ensure the confidentiality of sensitive messages, the messages are both signed and encrypted. Each node has its own public and private key

pairs. However, the computing resource of nodes in energy harvesting cooperative networks is limited. The asymmetric encryption algorithms consume much more energy and computing resource than that of symmetric encryption. Thus, the symmetric encryption algorithms are better in most situations. However, we still need to use asymmetric encryption to ensure secure communication in the process of transmitting a symmetric key.

Once node i and node j both know the symmetric key, they can communicate with each other through well-known message authentication code (MAC or HMAC). But a drawback of symmetric encryption is that non-repudiation cannot be ensured well, although the likelihood of data being surreptitiously changed is extremely low. This is a compromise solution between efficiency and security. To achieve a higher level of security for sensitive messages, one can apply active security mechanisms or adopt public key infrastructure (PKI) encryption at the cost of losing a certain amount of efficiency.

For energy-relay information attacks in CREP or CRED, the key point is to prove that the cooperator has finished the information transmission. Therefore, we can provide the evidence via an information dual-signature, which is similar to the energy dual-signature. More specifically, the information is also divided into a number of normalized information packets. Then these information packets are transmitted by the same steps as normalized energy units shown in Fig. 3. The destination should give back the information dual-signature of each normalized information packet to the cooperator, which is the evidence that the destination has received the information packet. If an energy-related information attack is launched, with the help of the registration authority, the cooperators can defend it by making use of the information dual-signature.

ENERGY SECURITY ANALYSIS AND NUMERICAL RESULTS

ENERGY SECURITY ANALYSIS

For potential energy attacks in energy harvesting cooperative networks, the security requirements for energy cooperation are achieved by the above defense solutions. As shown in Fig. 3, the energy dual-signature includes digital signatures and energy normalization, which guarantees that energy is indivisible to satisfy energy integrity. By the help of dual-signature, energy transfer is undeniable. More specifically, non-frameability is achieved as follows:

- A node cannot pretend to be another node since a transmitted digital signature is signed by its private key, which is unique and verifiable.
- A receiver cannot launch energy repudiation of reception because the energy is bound with its signature from the sender and dual-signature signed by the anonymous key of the receiver. Likewise, a node cannot claim that a signature message was repeated as each signature message includes a timestamp.

- A node cannot cheat other nodes by issuing false energy state since its periodic energy state report provides true energy state information.

The periodic energy state report aims to verify the energy state of the nodes. The forged energy state information can easily be disproved by the registration authority, which stores all the periodically uploaded energy state information from trusted modules of the nodes. Consequently, the requirement of energy state credibility is achieved.

In order to preserve anonymity of the node and minimize the energy consumption for encryption, we adopt a symmetric and asymmetric encryption to ensure energy usability. The lifetime of a node is closely associated with its computation capability, memory, and battery power. An efficient encryption algorithm reduces not only delay but also energy consumption.

For energy-related information attacks, digital signature and information normalization ensure information integrity and non-frameability during the process of information dual-signature, which is similar to energy dual-signature. If an attacker is held accountable for its illegal actions, the registration authority will reveal the true identity of the attacker and put it on a blacklist. Then all the legitimate nodes will get the blacklist and refuse to communicate with the attacker.

NUMERICAL RESULTS

We carry out a simulation to evaluate the performance of the proposed energy defense solutions. We first compare the performance of energy loss under different attack strengths and arrival rates of information. The energy consumption of our security solutions is also evaluated. For the purpose of illustration, we simulate the performance among 100 pairs of nodes for 100 min. The communication range between the nodes is set to be 10 m. The arrival process of an energy packet is a sequence of independent identically distributed (iid) random variables, which takes values from the set $\sigma = \{0.25, 0.5, 0.75, 1\}$ mJ with equal probability [1]. The information packet arrival rate (i.e., μ) in the node's data buffer follows a Poisson process. In this article, the request and reply information packet transmitted between nodes is 10 bytes, which consumes 4.08 μ J energy. A signature is set to be 140 bytes (1 digital signature, 1 key, and 1 certificate), which is a typical range of elliptic curve cryptography. A node consumes 123.2 μ J energy to send a signature to others [9].

Figure 4a shows the performance of our energy defense solutions and that without security protection under different attack strengths (the ratio of attacker nodes among 100 pairs of nodes). These curves illustrate that the security protection performance is much better than that without security protection when the attack strength increases. This is due to the fact that most of the energy attacks are defended against by our security protection solutions. Comparing the rate of transmitting information arrival packets, the energy loss when $\mu = 2$ is 12 percent larger than when $\mu = 1$ in the scenario without

If an attacker is held accountable for its illegal actions, the registration authority will reveal the true identity of the attacker and put it on a blacklist. Then all the legitimate nodes will get the blacklist and refuse to communicate with the attacker.

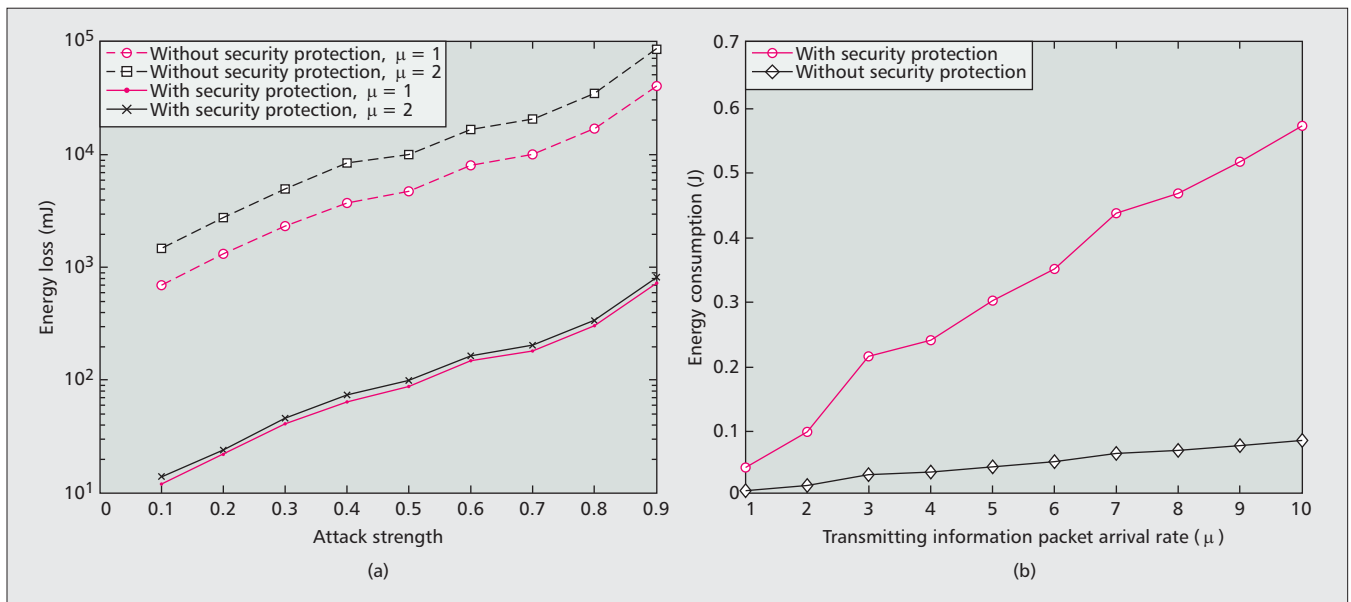


Figure 4. Performance comparison between scenarios with our security protection and without security protection: a) energy loss performance; b) energy consumption performance.

security protection. It implies that the higher the information packet arrival rate, the greater the energy loss. However, for the scenario with security protection, the energy loss is much less than that without security protection since our security solutions can resist most of attacks. Therefore, the performance is not sensitive to the changing rate of information arrival.

Figure 4b shows the energy consumption of our energy defense solutions under different transmitting information packet arrival rate. The curves illustrate that the higher packet arrival rate results in more energy consumption. Meanwhile, the case with security protection incurs more energy than that without protection. In addition, the energy protected by our defense approach (Fig. 4a) is 15 times the amount of the energy consumption (Fig. 4b) for our defense solution for $\mu = 1$. The results therefore indicate that our energy defense solutions are effective and feasible.

CONCLUSION

In this article, we have presented security issues and solutions in energy harvesting cooperative networks. We have performed security and performance analysis to evaluate our proposed solutions. Through numerical results, we have shown that our energy defense solutions are effective and efficient.

There are several interesting problems that can be further studied, such as optimal cooperator selection algorithm [10], denial of service attack, and energy state privacy preservation issues.

ACKNOWLEDGMENT

Rong Yu is the corresponding author for this article. The work is supported in part by programs of NSFC under Grant nos. 61370159, U1201253, and U1301255, Guangdong Province Natural Science Foundation under Grant no.

S2011030002886, Special-Support Project of Guangdong Province under grant. not. 2014TQ01X100, High Education Excellent Young Teacher Program of Guangdong Province under grant no. YQ2013057, Science and Technology Program of Guangzhou under grant no. 2014J2200097 (Zhujiang New Star Program), the projects 240079/F20 funded by the Research Council of Norway, and the European Commission FP7 Project CROWN (grant no. PIRSES-GA-2013-627490).

REFERENCES

- [1] B. Gurakan et al., "Energy Cooperation in Energy Harvesting Communications," *IEEE Trans. Commun.*, vol. 61, Dec. 2013, pp. 4884–98.
- [2] K. Tutuncuoğlu and A. Yener, "Optimum Transmission Policies for Battery Limited Energy Harvesting Nodes," *IEEE Trans. Wireless Commun.*, vol. 11, no. 3, 2012, pp. 1180–89.
- [3] D. Gunduz and B. Devillers, "Two-Hop Communication with Energy Harvesting," *2011 4th IEEE Int'l. Wksp. Computational Advances in Multi-Sensor Adaptive Processing*, 2011, pp. 201–04.
- [4] A. Di Mauro, D. Papini, and N. Dragoni, "Security Challenges for Energy-Harvesting Wireless Sensor Networks," *PECCS*, 2012, pp. 422–25.
- [5] A. V. Taddeo, M. Mura, and A. Ferrante, "Qos and Security in Energy-Harvesting Wireless Sensor Networks," *Proc. 2010 IEEE Int'l. Conf. Security and Cryptography*, 2010, pp. 1–10.
- [6] S. Pelissier et al., "Providing Security in Energy Harvesting Sensor Networks," *IEEE Consumer Commun. and Net. Conf.*, 2011, pp. 452–56.
- [7] D. Martins and H. Guyennet, "Wireless Sensor Network Attacks and Security Mechanisms: A Short Survey," *2010 13th IEEE Int'l. Conf. Network-Based Info. Systems*, 2010, pp. 313–20.
- [8] K. Ishibashi, H. Ochiai, and V. Tarokh, "Energy Harvesting Cooperative Communications," *2012 IEEE 23rd Int'l. Symp. Personal Indoor and Mobile Radio Commun.*, 2012, pp. 1819–23.
- [9] W. R. Heinzelman, A. Chandrakasan, and H. Balakrishnan, "Energy-Efficient Communication Protocol for Wireless Microsensor Networks," *Proc. 33rd Annual IEEE Hawaii Int'l. Conf. System Sciences*, 2000, p. 10.
- [10] B. Medepally and N. B. Mehta, "Voluntary Energy Harvesting Relays and Selection in Cooperative Wireless Networks," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, 2010, pp. 3543–53.

BIOGRAPHIES

JIAWEN KANG (kangjw@mail2.gdut.edu.cn) is now pursuing his Ph.D. degree at Guangdong University of Technology (GDUT), China. His research interests mainly focus on resource management, security, and privacy protection in wireless communications and networking.

RONG YU [M] (yurong@ieee.org) received his Ph.D. degree from Tsinghua University, China, in 2007. He is a full professor at GDUT. His research interest mainly focuses on wireless communications and networking, including cognitive radio, wireless sensor networks, and home networking. He is the co-inventor of over 10 patents, and an author or co-author of over 70 international journal and conference papers. He is currently serving as Deputy Secretary General of the Internet of Things (IoT) Industry Alliance, Guangdong, and Deputy Head of the IoT Engineering Center, Guangdong. He is a member of the Home Networking Standard Committee in China, where he leads the standardization work on three standards.

SABITA MAHARJAN [M] (sabita@simula.no) is currently a post-doctoral fellow at Simula Research Laboratory, Norway. She received her M.E. from Antenna and Propagation Lab, Tokyo Institute of Technology, Japan, in 2008, and her Ph.D. from the University of Oslo and Simula Research Laboratory, Norway, in 2013. Her research interests include wireless networks, smart grid communications, cyber-physical systems, machine-to-machine communications, software defined networks, and network security.

YAN ZHANG [SM] (yanzhang@ieee.org) received a Ph.D. degree from Nanyang Technological University, Singapore. Since August 2006, he has been working at Simula Research Laboratory. He is currently head of the Department of Networks; and an adjunct associate professor at the Department of Informatics, University of Oslo. He is a Regional Editor, Associate Editor, Guest Editor, or on the Editorial Boards of a number of international journals. His recent research interests include wireless networks, cyber physical systems, and smart grid communications.

XUMIN HUANG (huangxm@ieee.org) is now a master student of networked control systems at GDUT. His research interests mainly focus on network performance analysis, simulation, and enhancement in wireless communications and networking.

SHENGLI XIE [SM] (shlxie@gdut.edu.cn) received his M.S. degree in mathematics from Central China Normal University in 1992 and his Ph.D. degree in control theory and applications from SCUT in 1997. He is presently a full professor and head of the Institute of Intelligent Information Processing at GDUT. His research interests include wireless networks, automatic control, and blind signal processing. He is the author or coauthor of two books and more than 150 scientific papers in journals and conference proceedings. He received the second prize in China's State Natural Science Award in 2009 for his work on blind source separation and identification.

HANNA BOGUCKA [SM] (hbogucka@et.put.poznan.pl) received M.Sc. and Ph.D. degrees in telecommunications from Poznan University of Technology (PUT), Poland, in 1988 and 1995, respectively. Since 1988 she has been employed at PUT, currently in the Chair of Wireless Communications as a full professor and deputy dean for research on the Faculty of Electronics and Telecommunication. She has research interests in the area of wireless communications, cognitive radio, and green communication systems. She is the author of around 150 papers and five monographs and handbooks the area of radio communications.

STEIN GJESSING [SM] (steing@ifi.uio.no) is an adjunct researcher at Simula Research Laboratory and a professor of computer science in the Department of Informatics, University of Oslo. He received his Ph.D. degree from the University of Oslo in 1985. He acted as head of the Department of Informatics for four years, beginning in 1987. From February 1996 to October 2001 he was the Chairman of the national research program Distributed IT-System funded by the Research Council of Norway. He has participated in three European funded projects: Macrame, Arches, and Ascissa. His current research interests are routing, transport protocols, and wireless networks, including cognitive radio and smart grid applications.

ENERGY-EFFICIENT OPTICAL NETWORKS



Osman Gebizlioglu



Vijay Jain

The emerging trends of 2014 have continued to dominate the developments in the first half of 2015:

1. The packet optical transport system (P-OTS) concept supports the service provider goal to build a common network infrastructure by integrating Ethernet, IP/MPLS, and DWDM
2. There has been intensifying focus and attention on network functions virtualization (NFV) and software-defined networking (SDN) to provide service providers with the tools for more effective operation and management of communications networks in general, and optical communications networks in particular.

All major global service providers have announced their target timeframes for transitioning to all-IP network operations, while major industry standards organizations and fora have been striving to publish documents addressing use cases and plans for NFV and SDN. These important trends are expected to gain additional focus through the rest of 2015.

In this issue, we have selected four contributions that address high-data-rate coherent optical slot-switched networks, data center interconnection with elastic optical networks, optical interconnects for energy-efficient data center networks, and a new standard for energy efficiency in future optical access networks

In the first contribution, “High Data Rate Coherent Optical Slot-Switched Networks: A Practical and Technological Perspective,” Y. Pointurier, G. de Valicourt, Jesse E. Simsarian, Jürgen Gripp, and Francesco Vacondio present a review of node architectures for optical slot switching ring networks. Such networks can be used for metro and data center applications. Optical slot switching is a variation of optical packet switching, where all switched entities are slots of the same duration. Industrial research on optical packet switching has recently focused on optical slot switched (OSS) rings for use in metro networks. In addition to a review of node architectures, this contribution provides a comparison of key networking features, and a description of complex components and subsystems used in coherent slotted ring networks. Quality of service (QoS), latency and protection are discussed.

In the second contribution, “CSO: Cross Stratum Optimization for Optical as a Service,” H. Yang, J. Zhang, Y. Zhao, Y. Ji, J. Han, Y. Lin, and Y. Lee present data center interconnection with elastic optical networks as a promising solution to meet the high-bandwidth requirements for services. Many data center services require lower latency and

higher availability with the end-to-end guaranteed QoS, and this requires optimization through not only the data center networks, but also the optical transport networks that connect the data center networks. However, elastic optical networks and data centers are deployed separately. Since interworking among these separate systems requires the use of complex and inelastic interfaces, this is not an efficient solution for guaranteeing the required QoS, low latency, and effective protection. Thus, this article presents a novel cross stratum optimization (CSO) architecture in elastic data center optical interconnection to enable global optimization and control across an elastic optical transport network and the data center application layer. The authors describe CSO architecture and functional service provisioning schemes that are based on CSO. The optical as a service (OaaS) testbed deploying CSO with Open-Flow-enabled elastic optical node devices is presented and discussed with numerical results.

In the third contribution, “Optical Interconnects at the Top of the Rack for Energy-Efficient Data centers” J. Chen, Y. Gong, M. Fiorani, and S. Aleksic present three major types of passive optical interconnects and results of a performance assessment that covers the ability to host data center traffic, scalability, optical power budget, complexity of the required interface, and cost and energy consumption. Data center traffic volume has been steadily increasing with the growing use of cloud and multimedia services. This growth has been driving the demand for highly scalable, flexible, and energy-efficient networks inside data centers. Thus, optical networks in data centers have been recognized for energy efficiency and cost effectiveness. In order to take advantage of what optical networks offer, passive optical interconnects are considered for high-capacity data centers since no feasible solutions based on optical switching are available for handling high traffic volumes. The results have shown that the passive optical interconnects can achieve a significant reduction in power consumption and maintain cost at a similar level compared to electronic interconnects.

In the fourth contribution, “The Watchful Sleep Mode: A New Standard for Energy Efficiency in Future Access Networks” R. O. C. Hirafuji, K. B. da Cunhay, D. R. Campelo, A. R. Dhaini, and D. A. Khotimsky present the watchful sleep mode, a new mode that unifies the doze and cyclic sleep modes in a standardized PON energy efficiency mechanism into a single power management mode. The two standardized modes require signaling and result in energy waste. The new

mode saves energy by eliminating additional control signaling, and it has been included in international standards such as International Telecommunication Union — Telecommunication Standards Sector (ITU-T) 984 (G-PON) and ITU-T G.987 (XG-PON). In the new mode, an ONU periodically turns off its receiver and transmitter, as in the cyclic sleep mode, and performs infrequent bidirectional handshakes, as in the doze mode. The watchful sleep mode not only simplifies the implementation of the power management scheme at the ONU and OLT, but also combines the advantages of the cyclic sleep and doze modes, and outperforms both of them. More interestingly, a PON system supporting the watchful sleep mode can actually emulate the cyclic sleep or doze mode as a special case. It is also being considered for the NG-

PON2 standard (ITU-T G.989), which aims at standardizing TWDM-PON networks. The authors provide an overview of the main standardized (or considered for standardization) power saving techniques in PON systems. They also discuss its practical implementation details, and present a comparative study to highlight its advantages. Furthermore, they discuss the future outlook and how the new mode can be used in NG-PON2 systems.

In this second (and last) Optical Communications Series (OCS) issue of 2015, we thank all authors and reviewers for their valuable contributions to the OCS in past years and invite submissions on all aspects of optical communications technologies, with our best wishes to you for the rest of 2015 and beyond.

CALL FOR PAPERS

**IEEE TRANSACTIONS ON MOLECULAR, BIOLOGICAL, AND
MULTISCALE COMMUNICATIONS**

COMMUNICATIONS BEYOND CONVENTIONAL ELECTROMAGNETISM

This journal is devoted to the principles, design, and analysis of signaling and information systems that use physics beyond conventional electromagnetism, particularly for small-scale and multi-scale applications. This includes molecular, quantum, and other physical, chemical, and biological (and biologically inspired) techniques, as well as new signaling techniques at these scales.

As the boundaries between communication, sensing and control are blurred in these novel signaling systems, research contributions in a variety of areas are invited. Original research articles on one or more of the following topics are within the scope of the journal: mathematical modeling, information/communication-theoretic or network-theoretic analysis, networking, implementations and laboratory experiments, systems biology, data-starved or data-rich statistical analyses of biological systems, industrial applications, biological circuits, biosystems analysis and control, information/communication theory for analysis of biological systems, unconventional electromagnetism for small or multi-scale applications, and experiment-based studies on information processes or networks in biology. Contributions on related topics would also be considered for publication.

Editor-in-Chief

Urbashi Mitra, University of Southern California, USA

Associate Editor-in-Chief

Andrew W. Eckford, York University, Canada

Submit today!

<https://mc.manuscriptcentral.com/tmbmc>

EDITORIAL BOARD

Behnaam Aazhang, Rice University, USA
 Chan-Byoung Chae, Yonsei University, Korea
 Faramarz Fekri, Georgia Tech, USA
 Ananth Grama, Purdue University, USA
 Negar Kiyavash, University of Illinois, USA
 Vikram Krishnamurthy, University of British Columbia, Canada
 Tommaso Melodia, Northeastern University, USA
 Stefan Moser, ETH Zurich, Switzerland
 Tadashi Nakano, Osaka University, Japan
 Christopher Rozell, Georgia Tech, USA

High Data Rate Coherent Optical Slot Switched Networks: A Practical and Technological Perspective

Yvan Pointurier, Guilhem de Valicourt, Jesse E. Simsarian, Jürgen Gripp, and Francesco Vacondio

ABSTRACT

We review several node architectures for optical slot switching ring networks, which can be used in metropolitan or datacenter applications, and compare them for their networking aspects. The dimensioning, quality of service, latency, and protection issues are discussed for the different approaches. The main devices, i.e. fast wavelength-tunable laser, burst-mode coherent receiver (which is required to enable high data rate transmission at 100 Gb/s and above), and a slot blocker for improved wavelength usage efficiency are described, and available technologies for each key building block are reviewed.

INTRODUCTION

Optical packet switching refers to any technology that brings into optical communications the packet switching paradigm, which is well known in electronics. Optical packet switching generally relies on the following principles: opto-electronic conversions (and processing) occur only at source and destination nodes, while intermediate nodes switch data at a short (nanosecond to millisecond) granularity, using a so called all-optical or “transparent” switching fabric. The reduction in opto-electronic conversions translates into energy savings, and facilitates network upgrades, since transparent switching fabrics can often be made independent of the data rate of the signals used in the network. The switching of data with sub-wavelength granularity enables statistical multiplexing, i.e. the sharing of the transport medium between several data flows. Packet switching is also well adapted to the increasing amount of Internet protocol (IP) traffic (as opposed to voice circuits) carried in telecommunication networks. Optical slot switching is a version of optical packet switching, where all switched entities are slots of the same duration. Industrial research on optical packet switching has recently focused on optical slotted switched (OSS) rings [1–4], as depicted in Fig. 1, in the context of metropolitan-area (metro) networks.

Topologies more complex than rings, such as meshes and interconnected optical rings, can also be built.

Optical slot switching enables sub-wavelength networking directly at the optical layer, without resorting to costly intermediate opto-electronic conversions. The same transponder (a combination of an emitter and of a receiver) can be used to transmit (or receive) data to (from) several destinations (sources), as opposed to traditional circuit-switched networks, which establish a static connection between a source and a destination, whether this connection is fully used or not. Sub-wavelength networking is relevant for instance in next-generation optical mobile backhaul networks, where nearby base stations cooperate to cancel interferences (enabled by coordinated multipoint technology) and need direct, low-latency communication at data rates well below typical channel rates of 10 Gb/s or more [5]. In addition, optical slot switching networks enable highly dynamic adaptation to traffic variations. Those features are especially useful in metro networks, which rely today either on static circuit-switched networks or power-hungry all-electronic switches; and in datacenters, which again today rely on all-electronic switches that can easily adapt to the highly variable and unpredictable traffic, at the cost of high power consumption [6].

Even though the components and subsystems presented in this article are suitable for ring and mesh networks, we limit the architecture discussion to *slotted rings*. Optical packet switching in *meshed* networks introduces additional challenges, mostly related to ranging and scheduling of the optical packets. Asynchronous and variable-size IP packets that can be as small as 46 bytes are too short at high data rates (10 Gb/s in access networks, 100 Gb/s and above in core transport networks) to be efficiently switched with current optical technologies. To resolve this problem, packets are first aggregated at the network edges into “bursts” (here, slots) of data, which are sufficiently long to be compatible with the switching times of optical elements [7]. Contention, which occurs when several nodes want

Yvan Pointurier, Guilhem de Valicourt, Jesse E. Simsarian, Jürgen Gripp, are with Alcatel-Lucent Bell Labs.

Francesco Vacondio is with Alcatel-Lucent Bell Labs and Molini Industri-ali SpA.

to send packets to the same node simultaneously, can be resolved by scheduling when the slots are sent into the ring or mesh network [8], or by having the capability of electronic buffering at intermediate routing nodes [9].

In the past few years Alcatel-Lucent Bell Labs has proposed several node architectures for high speed OSS rings [1, 2] using coherent signaling for high data rate transmission of 100 Gb/s per channel and above, targeted to metro or datacenter networks; data transmission is all-optical, but control and scheduling are performed in the electronic domain. This paper reviews the formerly proposed node architectures and compares them for key networking aspects. Coherent OSS rings require complex components, which has so far hindered their development and deployment. We describe those components and subsystems and show how we resolved some of the key technological challenges. In addition, we show how the key components could be integrated thanks to technologies such as silicon photonics, so that node cost could be reduced and made compatible with the metro or datacenter segments. Finally, we draw brief conclusions.

NODE ARCHITECTURES

A typical OSS ring network is depicted in Fig. 1. Data is encapsulated in fixed-duration slots, converted to a wavelength division multiplexed (WDM) optical signal, sent onto a ring network, where it transits transparently through the intermediate nodes, and is converted back to the electronic domain at the destination only. In addition, a control channel that is demodulated and electronically processed at every node carries the headers of each synchronous WDM slot. Those headers contain routing information such as the source and destination node for each WDM slot. The control channel is also used to carry a clock to all nodes and ensure synchronous operation of the ring.

Figures 2a–c present three OSS node architectures (OSSv1 in Fig. 2a, OSSv2 in Fig. 2b, and OSSv3 in Fig. 2c), which all rely on the same basic key blocks:

- An optical blocker, implemented either as a static wavelength blocker (for instance, using a wavelength selective switch (WSS)) in OSSv1, or a dynamic slot blocker that can physically erase slots after they are received so as to reuse the fiber capacity in OSSv2 and OSSv3. Hence, with OSSv2 and OSSv3, the same wavelength can be shared by several transponders.
- “Burst-mode” receivers, which are capable

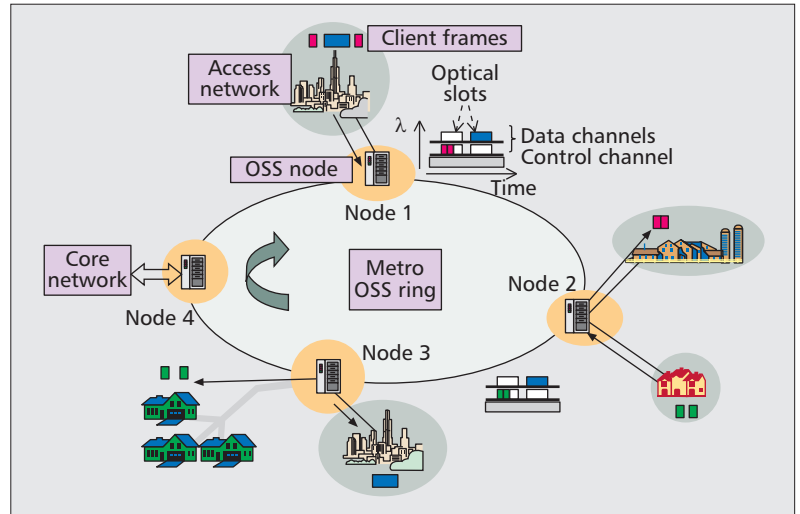


Figure 1. A sample optical slot switching network with two data channels and a control channel. Node 1 encapsulates a blue frame (coming from an access network) in a slot inserted on the first (top) data channel with destination node 3, and two red frames in a slot inserted on the second (bottom) data channel with destination node 2. Node 2 receives the slot containing the red frames and inserts a new slot containing green frames with destination node 3. Node 3 extracts the slots containing the blue and green frames.

of receiving bursts of data (duration: a few microseconds) separated by short guard intervals (duration: well below 1 microsecond)

- Fast tunable lasers, which can change their emission wavelength on a per-slot basis, are used as local oscillators to make the coherent receivers fast wavelength-tunable. Fast tunable lasers are also used in the *transmitters* of OSSv3.

In addition, amplification may be provided by standard optical amplifiers (EDFA). The utilization of specialized burst-mode EDFA with complex gain control electronics can be avoided by sending continuous streams of data in OSSv1 and keeping the guard intervals short (a few tens or hundreds of nanoseconds) in OSSv2 and OSSv3. Since no data can be transmitted during the guard intervals, slots have to be several microseconds long and may encapsulate several typical-sized (1500 bytes) client packets at data rates of 100 Gb/s and beyond. The guard interval has to be kept sufficiently short to avoid interaction with the EDFA automated gain control mechanisms and also to reduce wasted capacity, since no useful data is sent during the guard interval; at the same time the guard interval should be sufficiently long to allow the tuning of

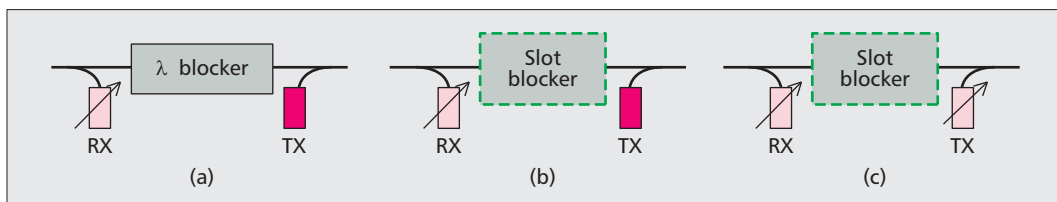


Figure 2. Optical slot switching node architectures: a) OSSv1; b) OSSv2; and c) OSSv3. An arrow indicates that the component contains a fast wavelength tunable laser, whereas the dark color TX block with no arrow is a fixed-wavelength laser.

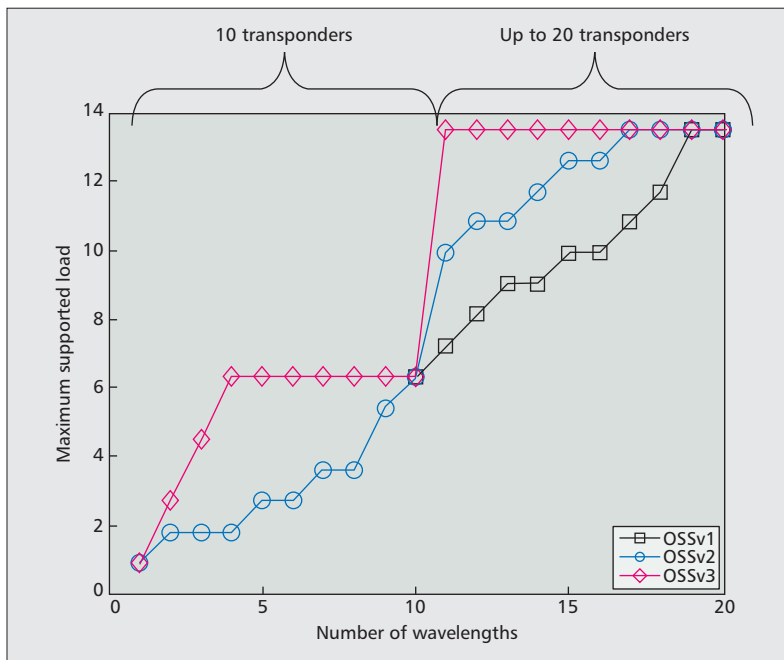


Figure 3. Comparison of the maximum supported load by the three proposed OSS nodes architectures in a 10-node unidirectional ring network, for a given amount of deployed resources (transponders and wavelengths). The supported load is normalized by the channel data rate. OSSv1 requires a minimum of 10 wavelengths for full connectivity.

the switching elements (fast tunable lasers, slot blocker) when they change state.

In OSSv1, “colored” emitters transmit on a predefined/static wavelength. Light from such an emitter enters the network after a wavelength blocker (Fig. 2a), which statically blocks incoming wavelength(s) to allow emission by the node at those wavelengths. In OSSv2, emitters are also colored but can be turned on or off on a slot-by-slot basis to allow the optical blocker, which is also reconfigurable at the slot granularity, to selectively erase slots. The node can re-use the capacity freed by the slot blocker in order to send new data onto the ring. OSSv3 uses both wavelength-tunable emitters and receivers (each leveraging a fast-tunable laser), in addition to a slot blocker. The OSS node architectures are increasingly flexible and dynamic, but also more complex and costly. The benefits of the flexibility in terms of maximum supported load in the network will be investigated.

NETWORKING ASPECTS

DIMENSIONING

In order to compare the three optical slot switching networks introduced in Fig. 2, we devised a dimensioning heuristic algorithm for each OSS variation. Given a traffic demand matrix, dimensioning gives the number of transponders and wavelengths required to carry the demand. Minimization of the number of transponders can be done by allocating as many transponders as required by the traffic matrix (this is the same number for all OSS flavors) and dedicating a separate wavelength for each transponder. Note that the utilization of slot blockers enables wavelength sharing by several transponders with

OSSv2 and OSSv3. Wavelength sharing decreases the number of required wavelengths, which in turn also reduces the network cost (since OSSv2 and OSSv3 require one gate per channel on every node), and a coarser transmission grid and hence less expensive devices. A heuristic algorithm is then used to minimize the number of deployed channels in the configurations with slot blockers (OSSv2 and OSSv3).

In Fig. 3 we show the maximum network load (the sum of all network demands; the load is normalized by the channel data rate) that is supported in a 10-node network, for each architecture, under a random uniform traffic demand, when the number of transponders and the number of wavelengths are fixed. OSSv1 requires one channel per node to carry traffic while OSSv2 and OSSv3 start carrying traffic when only a single wavelength is deployed. As more wavelengths are deployed, OSSv2 and OSSv3 carry up to twice as much traffic as OSSv1 for the same number of transponders and wavelengths.

RING SIZE

Although an OSS ring could be dimensioned for an arbitrary number of nodes, physical impairments sustained by signals crossing the nodes limit the cascadability of the OSS nodes. Unlike OSSv1 nodes, which use components that are similar to those found in circuit switched networks (splitters, WSS, EDFA) and hence are subject to the same cascadability limitations, OSSv2 and OSSv3 nodes include not only optical filters but also optical gates, which may distort signals. Optical gates such as (reflective) semiconductor optical amplifiers ((R)SOA) add optical noise and non-linear distortions. The largest count (typically several dozens) of cascable devices is obtained at the optimal trade-off between these impairments.

QoS SUPPORT

A modern metro ring network should be able to carry traffic with different classes of service (CoS), possibly under strict latency constraints, for instance to meet the Carrier Ethernet quality of service (QoS) specifications. In OSS, QoS support is achieved when client frames are encapsulated into slots, and when slots are inserted on the medium. In particular, CoS differentiation can be performed either during slot encapsulation (mixing frames with different classes in the same slot, but start filling slot with high priority frames) or at insertion time (prioritize adding slots containing more frames with higher CoS). Slot insertion on the fiber medium may be done according to a pre-defined schedule that is computed by a centralized controller, or according to reservations that are dynamically computed by each node in a distributed fashion, or simply opportunistically, i.e. by inserting data on the first available slot. Resource utilization of up to 80% with those methods was reported [10]. All methods are applicable to any of the OSS flavors. Observe, however, that OSSv3 is more flexible than OSSv1 and OSSv2 because all nodes can use any wavelength to send slots to any other node. Such flexibility facilitates decreased latency compared with OSSv1 and OSSv2.

PROTECTION

Protection in ring networks is typically achieved by using two counter-directional rings, with a wrapping or folding mechanism to isolate faults such as link or node failures. Such protection is applicable to OSS rings, which are compatible with both dedicated (resources reserved for both working and backup traffic) or shared (the same backup resources are shared to protect several working traffic streams) protection, or a combination of those to decrease network cost [11].

KEY DEVICES FOR AN OPTICAL SLOT SWITCHING NODE

SLOT BLOCKER

We show a slot blocker structure in Fig. 4a, consisting of a wavelength demultiplexer, one optical gate per wavelength that enables it to selectively erase any slot, and a wavelength multiplexer. The optical gate is one of the most crucial elements of the OSSv2/OSSv3 node as it is responsible for the (optical) quality of the transiting slots (when the gate is passing) and the inserted slots (when a gate is blocking, any remaining signal on the blocked wavelength may degrade the quality of a slot inserted further on the ring). Since a slot blocker structure contains several optical gates, it is especially important to minimize the cost and footprint of the optical gates without sacrificing functionality. Unfortunately, the cost of the slot blocker when implemented with discrete components has remained high for a long time. Recently, in order to decrease its cost, we proposed an implementation of the slot blocker using low-cost devices such as reflective semiconductor optical amplifiers (RSOA) [12] which have been developed for access networks: the packet reflective optical switch (PROS) shown in Fig. 4b [13]. Furthermore, the PROS requires only one fiber connector and is hence less costly than an integrated version of the slot blocker depicted in Fig. 4a, which requires two fiber connectors. The key specifications for the optical gates are their switching time (within the guard interval duration) and their extinction ratio, i.e. their ability to fully erase an incoming optical signal.

Active devices (in the standard InP platform) are excellent optical gate candidates due to their fast switching time, high extinction ratio, and high optical gain (SOA, RSOA, and EAM-SOA).

Considering the Silicon on Insulator (SOI) platform, various silicon-photonics devices have been proposed as building blocks for fast slot blocker structures such as Mach-Zehnder modulators (MZMs), ring resonator structures, or variable optical attenuators (VOAs). However, high insertion losses need to be compensated by optical amplifiers, and limited extinction ratio has been demonstrated [13].

Finally, hybrid III-V on silicon devices have been proposed [14] and may provide the benefits of SOAs as well as high integration capability.

FAST TUNABLE LASER

In circuit transmission, the major recent breakthrough to increase capacity has been the introduction of the digital coherent receiver, in

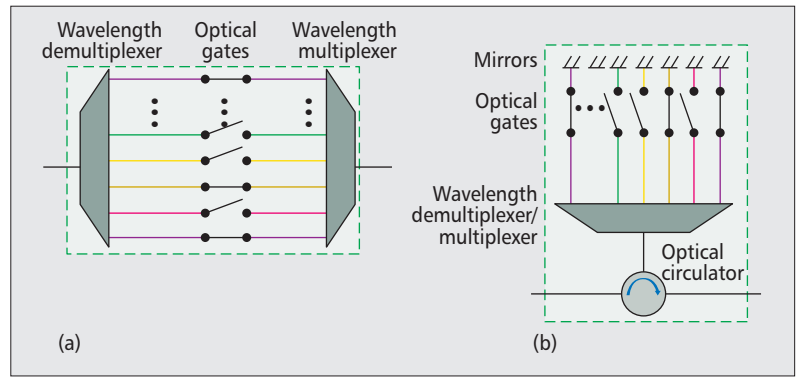


Figure 4. Slot blocker structures: a) basic and b) packet reflective optical switch “PROS”.

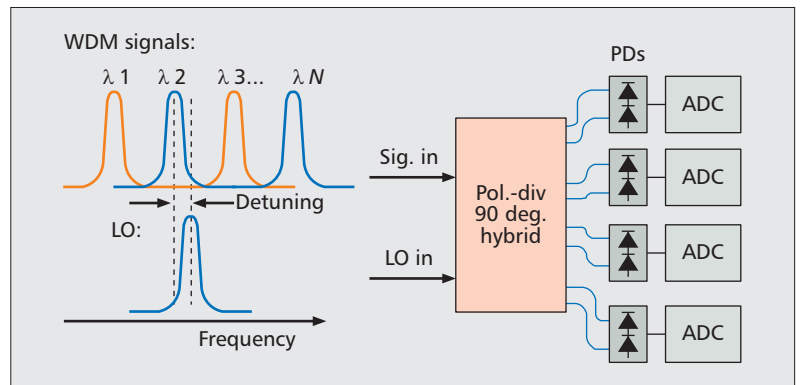


Figure 5. Comb of WDM signals and the local oscillator (LO) (left) at the inputs of the fast-tunable coherent receiver (right).

combination with polarization division multiplexing and high-order modulation formats.

We use coherent detection in the OSS network to take advantage of the high spectral efficiency modulation formats such as polarization division multiplexing quadrature phase shift keying (PDM-QPSK) that is presently used in commercial 100 Gb/s transmission systems. Figure 5 shows a diagram of the WDM signals and local oscillator wavelength that are received by a dual-polarization 90 degree optical hybrid. Whereas traditional WDM systems use an optical filter to select only one of the wavelengths before the receiver, the coherent receiver is able to select one of the wavelengths from the comb without optical filtering. This selectivity is possible since neighboring channels will have a larger optical beat frequency between the channel and local oscillator that can be removed by the low-pass analog filtering of the photodiodes (PDs) and analog-to-digital converters (ADCs) as well as any filtering that may be done by digital signal processing. This “colorless” capability of the receiver allows any node to receive optical slots from any other node by rapidly tuning the wavelength of the local oscillator.

The tuning of the laser needs to be fast since data cannot be transmitted or received during this time. Therefore, we use devices that tune electronically, such as the digital supermode distributed Bragg reflector (DS-DBR) laser [15]. The DS-DBR is monolithically-integrated on InP with multiple sections that operate with cur-

The frequency offset between transmit laser and local oscillator, the polarization state of the received signal, and the timing of the data symbols can change on a slot-by-slot basis. All these parameters need to be recovered rapidly, i.e. much faster than the guard interval duration which is typically well below a microsecond.

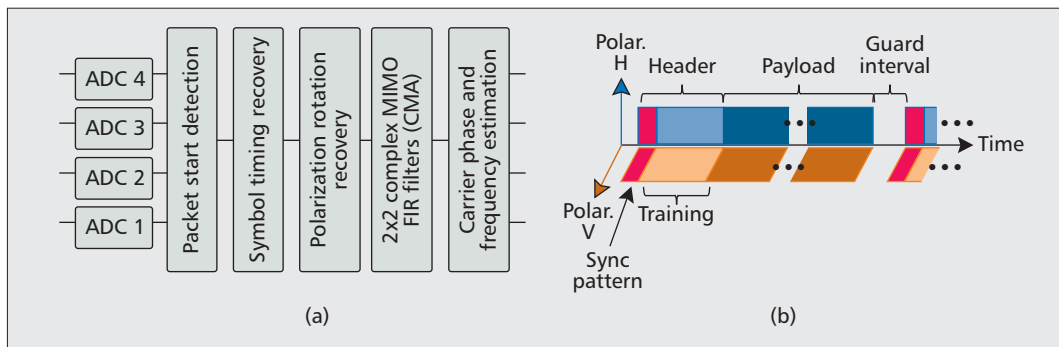


Figure 6. a) Digital signal processing blocks required in a burst-mode coherent receiver; b) structure of headers and payload for the packets.

rent injection. The gain and SOA sections produce and amplify light, respectively, and the front, phase, and rear sections allow for tuning across the wavelength band.

The phase noise of the transmitter and local oscillator lasers will negatively impact the performance of the network, especially coming from the local oscillator laser at the coherent receiver in the presence of fiber dispersion. The DS-DBR linewidth has contributions from low frequency noise (at frequencies $< \sim 100$ MHz), mainly coming from $1/f$ tuning-current noise, white noise at intermediate frequency (at frequencies from ~ 100 MHz to several GHz), and a relaxation oscillation at high frequency (at a frequency of several GHz). For high symbol rate systems, the intermediate and high frequency noise will cause the most serious impairments since this noise is more difficult to track with digital signal processing.

When designing the electronic drive circuitry of the laser, care must be taken to not introduce excessive phase noise on the laser from electronic noise on the driving currents while maintaining fast wavelength tunability. Using the DS-DBR laser as the local oscillator, we have demonstrated rapid and accurate wavelength switching, with frequency offset < 340 MHz after 80 ns, and a linewidth of 0.62 MHz after 120 ns [2]. To achieve these results, the phase section of the laser was used to compensate for frequency offsets and drifts caused by electronic carrier settling and thermal effects in the laser. With the phase compensation technique, rapid recovery of PDM-QPSK optical data with a bitrate of 112 Gb/s was achieved, even with up to 1280 km of fiber transmission.

BURST MODE RECEIVER

Coherent detection is now a commercial reality for circuit switching, but the use of a fast tunable laser as the local oscillator and burst-mode operation introduce a number of new challenges. The frequency offset between transmit laser and local oscillator, the polarization state of the received signal, and the timing of the data symbols can change on a slot-by-slot basis. All these parameters need to be recovered rapidly, i.e. much faster than the guard interval duration which is typically well below a microsecond.

We designed and implemented a (non real-time) fast slot receiver relying on the algorithmic blocks depicted in Fig. 6a and the slot format

depicted in Fig. 6b [16]. Independently of the modulation format, two PDM-QPSK headers are inserted at the slot beginning. The first 128 symbols are used to find the slot beginning, and 128 additional symbols that are used for channel estimation and initialize a 2×2 multiple-input multiple-output (MIMO) equalizer. We verified that, for slot lengths up to a few microseconds, the channel can be considered static over the payload, and therefore adaptive equalization is not required. After channel equalization, carrier frequency recovery is performed, and phase estimation is achieved with a maximum likelihood search. Finally, a decision is taken on the symbols, and the bit errors are counted on the payload. The algorithms are initialized on a slot-by-slot basis, and no specific block for chromatic dispersion compensation is used to cover the full range of distances discussed here. We showed in [16] that the receiver could decode slots independently modulated with various formats up to PDM-16QAM after transmission over 50 km of standard fiber *without amplification* in the context of datacenter, access, or small metro networks. Distances compatible with large metro or even core networks can be achieved with amplification.

CONCLUSION

We reviewed several optical slot switching ring networks for metro or datacenter applications, and their key building blocks: slot blockers for efficient wavelength utilization, and fast wavelength tunable lasers and burst-mode coherent receiver for high data rate transmission.

In order to ensure high system performance and data rate, it is essential that tunable lasers have sufficiently low phase noise. Suitable lasers are now commercially available, although the electronic control of these lasers must be enhanced to allow fast-switching capability. Burst-mode coherent reception is a novel challenge, and research breakthroughs were presented. We also described several technology platform options to implement a slot blocker. III-V material provides amplification and good optical performance and integration has been demonstrated at a research and commercial level, but the chip cost remains high. Silicon photonics also promises good optical performance, and integrated (but only passive) devices are currently commercially feasible. Hybrid III-

V/silicon technology has the largest potential, combining the benefits of silicon and InP into integrated active devices. First building blocks have been demonstrated, but higher levels of integration are not yet commercially feasible.

Technology for optical slot switching nodes with discrete components is already available and lab prototypes exist. However, further integration, along with the development of other non-hardware based blocks such as a novel control plane, will be needed before high-data rate, coherent optical slot switching reaches commercialization.

ACKNOWLEDGMENTS

This work was partly funded by the CELTIC+ SASER-SAVENET project.

REFERENCES

- [1] D. Chiaroni *et al.*, "Packet OADMs for the Next Generation of Ring Networks," *Bell Labs Technical J.*, vol. 14, no. 4, Winter 2010, pp. 265–83.
- [2] J. E. Simsarian *et al.*, "Fast-Tuning Coherent Burst-Mode Receiver for Metropolitan Networks," *IEEE Photon. Tech. Lett.*, vol. 26, no. 8, Apr. 2014, pp. 813–16.
- [3] Intune Networks, <http://www.intunenetworks.com/>.
- [4] S. Cao *et al.*, "A Novel Optical Burst Ring Network with Optical-Layer Aggregation and Flexible Bandwidth Provisioning," *Proc. OFC*, paper OThR5, 2012.
- [5] Q. Wei *et al.*, "Optical Mobile Network," *NTT Docomo Technical J.*, vol. 14, no. 2, Oct. 2012, pp. 43–54.
- [6] A. Greenberg *et al.*, "VL2: A Scalable and Flexible Data Center Network," *Proc. ACM SIGCOMM*, 2009, pp. 51–62.
- [7] C. Qiao and M. Yoo, "Optical Burst Switching (OBS) — A New Paradigm for an Optical Internet," *J. High Speed Networks*, vol. 8, no. 1, Mar. 1999, pp. 69–84.
- [8] I. Widjaja *et al.*, "Light Core and Intelligent Edge for a Flexible, Thin-Layered, and Cost-Effective Optical Transport Network," *IEEE Commun. Mag.*, vol. 41, no. 5, May 2003, pp. S30–S36.
- [9] R. Urata *et al.*, "First Demonstration of a Prototype Hybrid Optoelectronic Router," *Proc. ECOC*, paper PDP 3.2, 2009.
- [10] N. Benzaoui *et al.*, "Optical Slot Switching Latency in Mobile Backhaul Networks," *J. Lightw. Technol.*, vol. 33, no. 8, Apr. 15, 2015, pp. 1491–99.
- [11] A. Gravey *et al.*, "QoS of Optical Packet Metro Networks," *Proc. IEEE/OSA OFC*, paper We.1.C.5, Mar. 2014.
- [12] G. de Valicourt *et al.*, "Reflective Packet Add-Drop Multiplexer based on Modulation Format Agnostic and Low Cost Optical Gate," *Proc. ECOC*, paper We.2.B.5, 2012.
- [13] G. de Valicourt *et al.*, "Monolithic Integrated Silicon-based Slot-Blocker for Packet-Switched Networks," *Proc. ECOC*, paper We.3.5.5, 2014.

- [14] G. de Valicourt *et al.*, "A Next-Generation Optical Packet-Switching Node Based on Hybrid III-V/silicon Optical Gates," *IEEE Photonics Technol. Lett.*, vol. 26, no. 7, Apr. 1, 2014, pp. 678–81.
- [15] A. J. Ward *et al.*, "Widely Tunable DS-DBR Laser with Monolithically Integrated SOA: Design and Performance," *IEEE J. Sel. Topics Quantum Electron.*, vol. 11, no. 1, 2005.
- [16] F. Vacondio *et al.*, "Flexible TDMA Access Optical Networks Enabled by Burst-Mode Software Defined Coherent Transponders," *Proc. ECOC*, paper We.1.F.2, Sept. 2013.

BIOGRAPHIES

YVAN POINTURIER [S'02, M'06, SM'12] received a Ph.D. in electrical engineering in 2006 from the University of Virginia, USA. In 2009 he joined Alcatel-Lucent Bell Labs in France as a research engineer. He is the author or co-author of 16 European patents, more than 70 technical articles, and he is a co-recipient of a Best Paper Award at the IEEE ICC 2006 conference.

GUILHEM DE VALICOURT received the Ph.D. degree from Telecom ParisTech and III-V lab. In 2011 he joined Bell Labs Alcatel-Lucent in France and then in the U.S. in 2014. His main research interests include the study of advanced integrated photonics devices. He has authored or co-authored more than 70 scientific papers, two book chapters, and holds more than 15 patents. He received the 2012 Marconi Young Scholar Award and was a finalist for the 2012 Paris-Tech Ph.D. prize.

JESSE E. SIMSARIAN [M'04, SM'13] received the Ph.D. degree in physics from the State University of New York at Stony Brook. In 2000 he joined Bell Labs as a researcher in optical networking and switching. He has investigated software-defined transport networks, fast-switching coherent optical receivers, and scalable optical packet routers. He has authored or coauthored more than 70 publications. He is a senior member of the Optical Society of America.

JÜRGEN GRIPP [M'09] joined Bell-Labs in 1997 after receiving his Ph.D. in physics from the State University of New York at Stony Brook. Since joining Alcatel-Lucent Bell-Labs he has worked on optical Soliton transmission, OTDR-like dispersion mapping, fast wavelength switching tunable lasers, optical switch fabrics, high-capacity all-optical routers, and metropolitan-area networks. He is a member of the Optical Society of America.

FRANCESCO VACONDIO was a research engineer at Bell Labs' Centre de Villarceaux in Nozay, France. He received an M.S. degree in telecommunication engineering from the Università di Parma (Italy), and a Ph.D. degree in electrical engineering from Laval University in Canada. He joined Bell Labs in 2010, where he worked on the modeling of physical impairments and digital signal processing for dynamic and rapidly reconfigurable optical networks. He is now with Molini Industriali SpA (Italy).

Technology for optical slot switching nodes with discrete components is already available and lab prototypes exist. However, further integration, along with the development of other non-hardware based blocks such as a novel control plane, will be needed before high-data rate, coherent optical slot switching reaches commercialization.

CSO: Cross Stratum Optimization for Optical as a Service

Hui Yang, Jie Zhang, Yongli Zhao, Yuefeng Ji, Jianrui Han, Yi Lin, and Young Lee

ABSTRACT

Data center interconnection with elastic optical networks is a promising scenario to meet the high burstiness and high-bandwidth requirements of services. Many data center services require lower delay and higher availability with end-to-end guaranteed QoS, which involves both application and transport network resources. However, in the current mode of operation, the control of elastic optical networks and data centers is separately deployed. Enabling even limited interworking among these separated control systems requires the adoption of complex and inelastic interfaces among the various networks, and this solution is not efficient enough to provide the required QoS. In this article, we present a novel cross stratum optimization (CSO) architecture in elastic data center optical interconnection. The proposed architecture can allow global optimization and control across elastic optical transport network and data center application stratum heterogeneous resources to meet the QoS requirement with the objective of optical as a service (OaaS). The functional modules of CSO architecture, including the core elements of application and transport controllers, are described in detail. The cooperation procedure in CSO-based service provisioning and cross stratum service resilience modes is investigated. The overall feasibility and efficiency of the proposed architecture is also experimentally demonstrated on our OaaS testbed with four OpenFlow-enabled elastic optical nodes, and compared to MFA, ALB, and CSO-DGLB service provisioning schemes in terms of path setup/release/adjustment latency, blocking probability, and resource occupation rate. Numerical results are given and analyzed based on the testbed. Some future discussion and exploration issues are presented in the conclusion.

INTRODUCTION

Due to the rapid evolution of cloud computing and high-bit-rate data-center-supported applications, network operators are recurrently rethinking the way their networks are controlled to provide interconnection between data centers

and users. Since data center services are typically diverse in terms of required bandwidths and usage patterns, the network traffic supporting such services shows high burstiness and high-bandwidth characteristics. The traditional wavelength-division multiplexing optical transport network is inefficient to carry these applications due to the fixed International Telecommunication Union-Telecommunication Standards Sector (ITU-T) wavelength grids and spacing. To accommodate these services in flexible bandwidth connectivity channels, the architecture of an elastic optical network has been proposed and demonstrated [1], which is achieved by taking advantage of the orthogonal frequency-division multiplexing technology [2, 3]. The key idea of such a network is to allocate necessary spectrum resources with a fine tailored granularity for various user connection demands, and offer cost-effective and highly available connectivity channels [4]. Data center interconnection by elastic optical networks is a promising scenario to allocate spectrum resources for applications in a highly dynamic, tunable, and efficiently controlled manner.

Traditionally, the control of data centers and elastic optical networks is separately deployed in data center interconnect architecture [5]. The interaction through the user-network interface (UNI) among separated control systems of various vendors has been overly complex due to confidentiality and manageability [6]. Additionally, many delay-sensitive services require high-level end-to-end quality of service (QoS) guarantees, which make significant use of the optical network resources in the form of bandwidth consumption, and relate to application resources (e.g., computing and storage resource) in data centers. Depending on the technological heterogeneity and resource diversity, traditional data center interconnection does not provide a mechanism to exchange resource information across the boundaries of elastic optical network and data center application in independent operation scenario [7]. It is hard to expect that the current architecture will provide the integrated end-to-end dynamic connectivity and high-level performance requirements of these applications.

As a promising centralized control, software

Hui Yang, Jie Zhang,
Yongli Zhao, and Yuefeng
Ji are with Beijing University of Posts and Telecommunications.

Jianrui Han, Yi Lin, and
Young Lee are with
Huawei Technologies Co.,
Ltd.

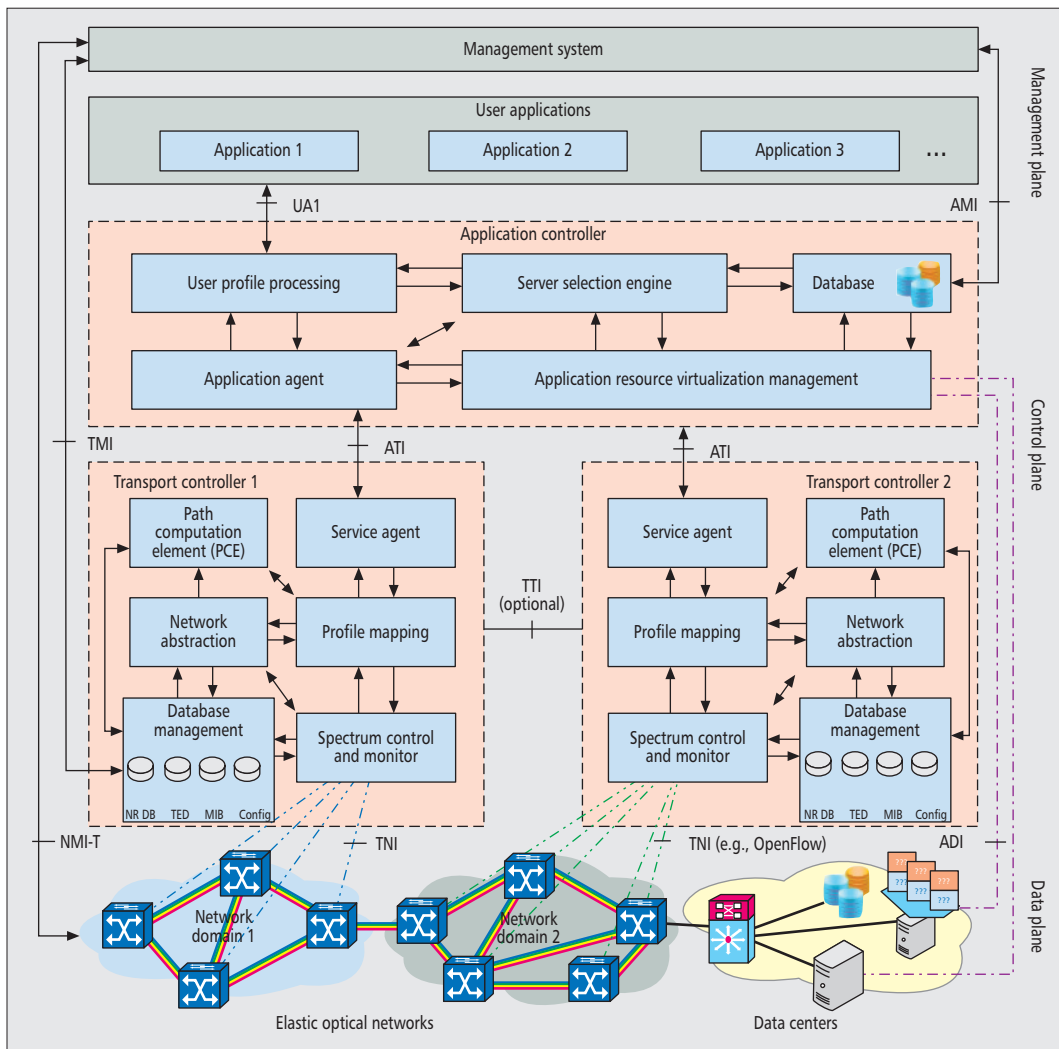


Figure 1. The CSO architecture.

The TCs sustain optical transport network stratum information abstracted from the physical network, and accordingly perform dynamic light-path provisioning in elastic optical networks. Six functional modules are included in a TC, which can control and manage the related network domain.

defined networking (SDN) enabled by OpenFlow protocol has gained popularity by supporting programmability of data center and network functionalities [8–10]. SDN can provide maximum flexibility for operators, unified control over various resources, and an abstraction mechanism for these resources [11, 12]. Based on the benefits of SDN, cross stratum optimization (CSO) is proposed as an architecture that allows global optimization and control across elastic optical network and data center application stratum heterogeneous resources to meet the QoS with a global view [13, 14]. The CSO can enable joint application and network stratum resources optimization, enhance responsiveness to end-to-end data center demands, and perform service resilience after a disaster via cooperative recovery techniques between application and network. This architecture is intended to enable optical as a service (OaaS) realizing both spectrum and application elasticity.

The rest of the article is organized as follows. We describe the CSO architecture and functional modules in detail. A cooperation procedure for CSO in two modes is proposed and discussed. Three service provisioning schemes are proposed based on CSO. The OaaS testbed

deploying CSO with OpenFlow-enabled elastic optical node devices is presented and discussed with numerical results. Finally, we conclude the article and discuss future research issues.

CSO ARCHITECTURE FOR DATA CENTER SERVICES

We present the CSO architecture in elastic data center optical interconnection to focus on the users' QoS requirement at two dimensionalities. First, heterogeneous services are replicated over multiple data centers so that a user request can be served from one of several potential data center application resources supporting the specified service [15], which is called *application elasticity*. Second, an elastic optical network can adjust the spectrum parameter in the physical layer to accommodate the service according to the transmission distance of the path. It enhances the spectrum utilization and realizes *spectrum elasticity*. We investigate the CSO to enhance QoS, and improve the transport network and data center application resources utilization, which consider the *application* and *spectrum elasticity* from two dimensions.

Message type	Process module
Network abstract topology request	Network abstraction
Path computation request	Path computation element
Path setup or adjust physical layer parameter request	Spectrum control and monitor
Computation result reply	Service agent

Table 1. Message types and corresponding process modules.

CSO ARCHITECTURE IN ELASTIC DATA CENTER OPTICAL INTERCONNECTION

The CSO architecture in elastic data center optical interconnection for data center services is illustrated in Fig. 1. The elastic optical networks are used to interconnect the distributed data centers in such architecture. It follows that the network architecture mainly consists of two strata: the optical transport network resources stratum (e.g., spectrum) and the application resources stratum (e.g., CPU and memory). Each resource stratum is software defined with OpenFlow, and controlled locally by a transport controller (TC) and an application controller (AC), respectively, in a unified manner. Multiple TCs are used to control multiple domain elastic optical networks in the CSO architecture, where each TC maintains the information of its own domain. To control the heterogeneous networks for data center interconnection with extended OpenFlow protocol (OFP), OpenFlow-enabled elastic optical device nodes with OFP agent software are required, which are referred to as a software defined optical transport network (SD-OTN) and regenerator (SD-R), as proposed in [15]. The motivations for CSO in elastic data center optical interconnection are twofold. First, the CSO architecture emphasizes the cooperation between TC and AC to select the optimal data center destination with application elasticity achieving joint and global optimization of transport network and application resources. Second, based on the distance of path and network resource, the architecture can control elastic optical networks and adjust the physical layer parameter (e.g., bandwidth and modulation format) to realize the software defined path (SDP) for service provisioning and further optimize resource utilization with spectrum elasticity.

TRANSPORT CONTROLLER

The TCs sustain optical transport network stratum information abstracted from the physical network, and accordingly perform dynamic light-path provisioning in elastic optical networks. Six functional modules are included in a TC, which can control and manage the related network domain. The corresponding functional modules are described as follows.

The service agent is the external communication module to interact with the AC and the profile mapping module, which certifies the request to ensure security control for elastic path computation.

The profile mapping is the core module of

the TC, which is responsible for message scheduling among the functional modules in a TC and delivering different messages to the corresponding process modules. It can translate the location of the data center server and the service type into connection and service parameters in a transport network (e.g., bandwidth, delay and jitter). The different types of messages and the required modules to handle a specified incoming message after parsing are illustrated in Table 1.

The network abstraction can abstract and manage the network topology from a physical layer network through PCE calculation, and provide abstracted resource information to the AC.

The path computation element (PCE) is capable of computing a network path or route based on a network graph and applying computational constraints [17]. If the destination node is located in another network domain, the inter-domain path computation should also be performed using PCE.

The spectrum control and monitor is responsible for monitoring physical layer network elements and controlling the spectrum bandwidth and modulation format in the underlying layer network. The end-to-end SDP can be provisioned by controlling all corresponding SD-OTNs along the computed path by using extended OFP. Note that the modulation format of service, such as quadrature phase shift keying (QPSK) and 16-quadrature amplitude modulation (16-QAM), is determined and adjusted based on the length of SDP. The SD-Rs are allocated and controlled via the OFP agent in the proposed architecture in order to increase the range of optical transmission, thus leading to potential transmission delay and control process latency.

The database management conserves the information of the SDP from the PCE and the abstract topology information from the network abstraction module.

APPLICATION CONTROLLER

The centralized AC is responsible for monitoring and maintaining application stratum resources in data centers. A typical AC consists of five essential modules, which are described as follows.

The user profile processing certifies the data center service requests including authentication, authorization, and accounting, and maps them into application request parameters (e.g., service type and application resource).

The application resource virtualization management monitors and maintains the virtual application resources obtained from data centers.

The server selection engine chooses the most optimal server or virtual machine for users according to application and elastic optical network resources utilization, allocates application resources, and determines the location of an application or where to migrate virtual machines.

The application agent is the communication module that interacts with TCs through the application-transport interface (ATI) and obtains network information from them periodically or based on event-based trigger.

The database stores the abstract application and network information from TCs and data centers.

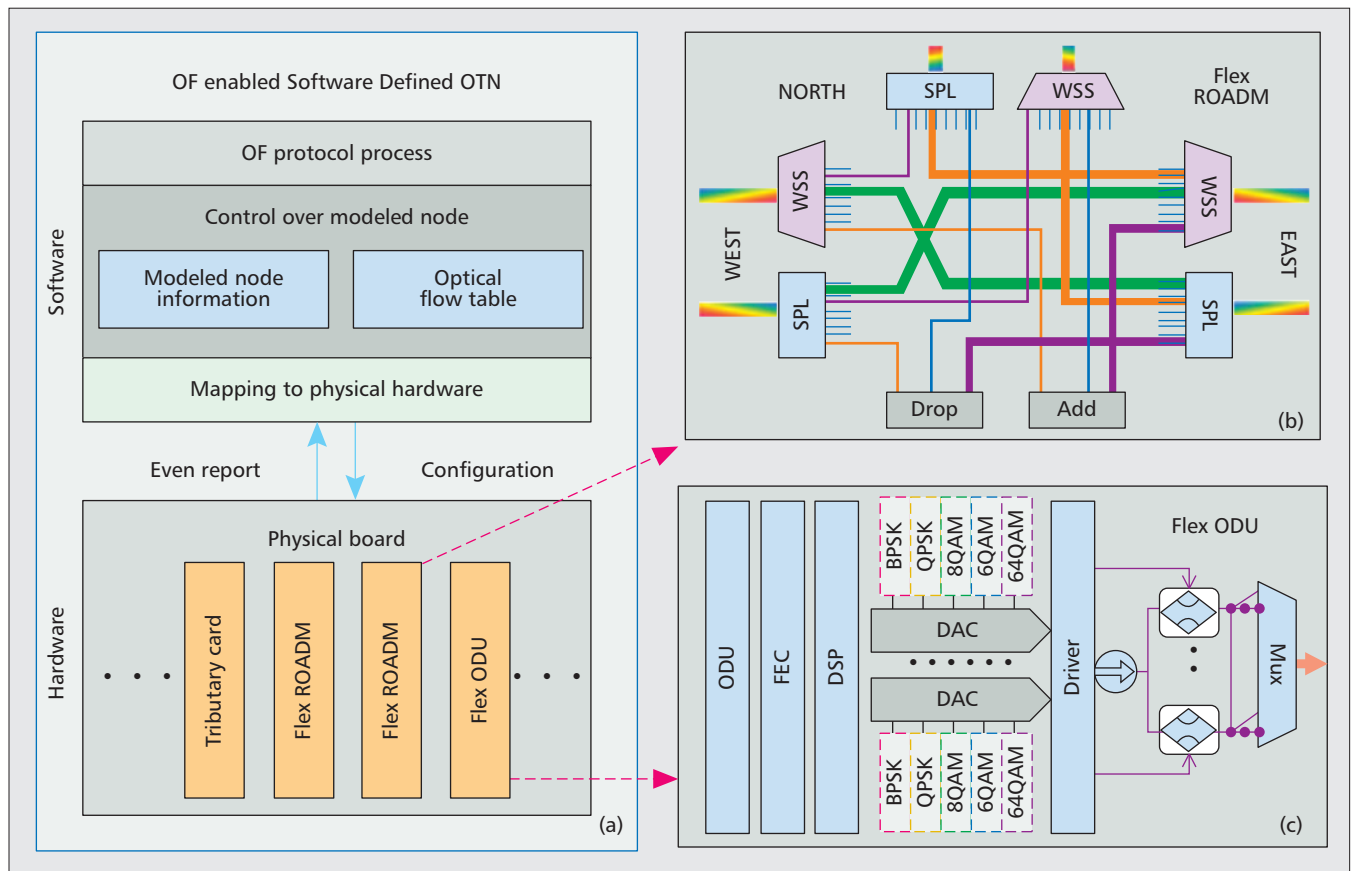


Figure 2. The functional models of a) OpenFlow-enabled SD-OTN; b) Flex ROADM; c) Flex ODU.

FUNCTIONAL MODEL OF SD-OTN

The functional models of SD-OTN are described in Fig. 2a. In SD-OTN, OpenFlow-enabled agent software is embedded to keep the communication between a TC and an optical node, which realizes the OpenFlow protocol process, controls modeled node and mapping to the physical hardware. Through the agent, the SD-OTN can maintain an optical flow table and modeled node information as software, and map the content to configure and control the physical hardware, and receive the report of spectrum control. The hardware of SD-OTN is composed of a series of physical boards, such as flexible reconfigurable optical add/drop multiplexer (ROADM), the optical demultiplexing unit (ODU) board, and the corresponding tributary card, due to the expansibility and convenience. The flexible ROADM and ODU boards meet the non-blocking and gridless requirement, and support various modulation format of the signal, as shown in Figs. 2b and 2c, respectively.

BUSINESS CONTEXT OF THE CSO

In the traditional communication ecosystem, the OTNs and data centers are usually operated by different business parties in the data center interconnection. Due to the inconvenience, and the development of user requirement and company scale, the single business party is forced to extend its types of services to adapt to the environment and promote competitiveness. Service providers can establish the OTN to interconnect their data centers, while the network operators

establish data centers to provide storage and computing resources for their clients. In such a scenario, the CSO and control of both OTNs and data centers can be performed by the single business party. Also, in the current standardization, many network vendors, operators, and service providers focus on the interfaces among the multi-domain controllers (e.g., in the Optical Networking Forum, ONF, and the Internet Engineering Task Force, IETF), including the southbound interface (e.g., OpenFlow) to the transport devices and northbound interface (e.g., the RESTful application programming interface, API) to the application. If the interfaces among the application and transport controllers are presented as the standard, the CSO can also be performed by different business parties in a traditional business environment.

COOPERATION PROCEDURE FOR CSO IN DIFFERENT SERVICE MODES

Cooperation between the AC and TC among different service modes is one of the key issues for CSO. It helps to achieve the service provisioning and resilience after a disaster for CSO. This section summarizes different cooperation procedures between two types of controllers.

CSO-BASED SERVICE PROVISIONING MODEL

It is necessary for the AC to collaborate with corresponding TCs when a data center service request needs cross application and network

Three service provisioning schemes are proposed based on the CSO architecture in the elastic data center optical inter-connection: modulation format adaptive adjustment (MFA), application load balancing (ALB), and CSO-enabled dynamic global load balancing (CSO-DGLB).

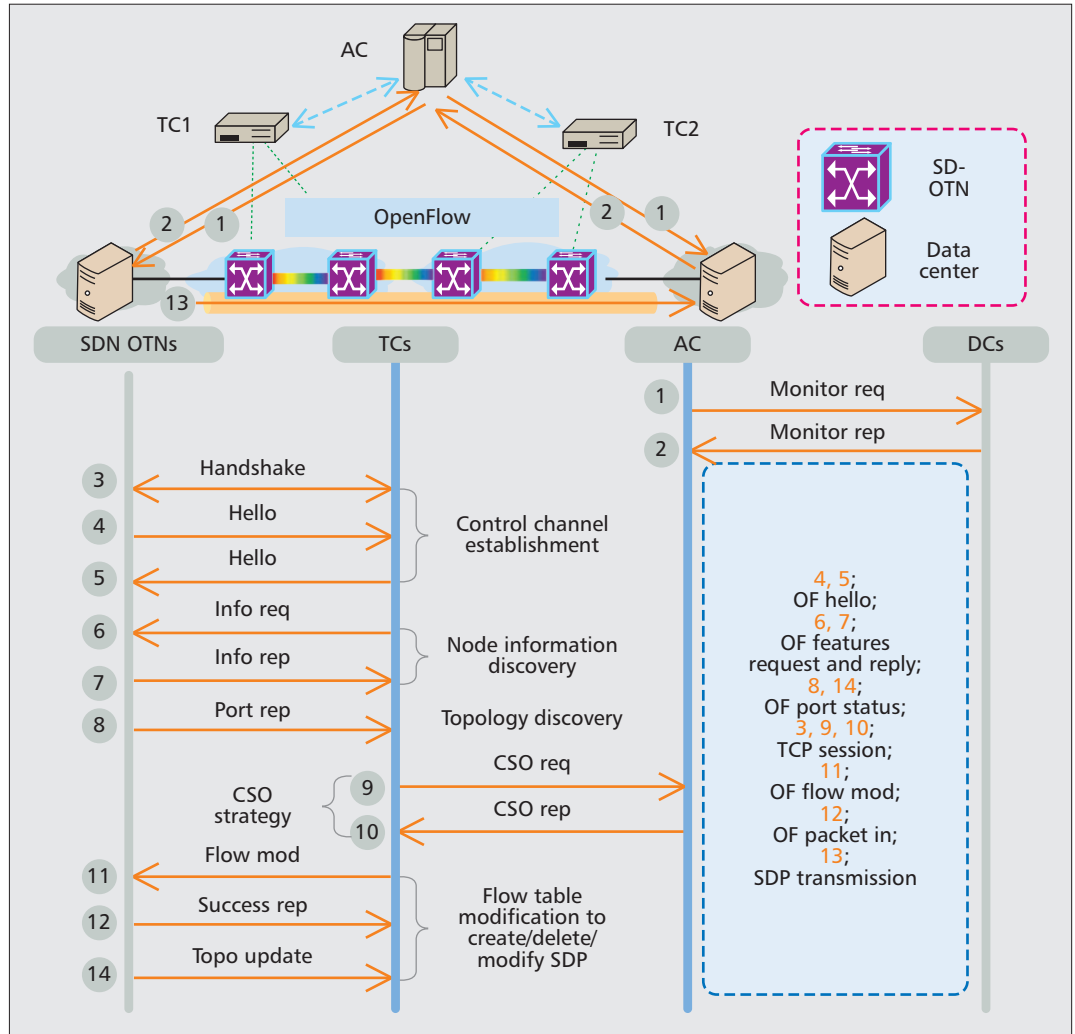


Figure 3. Cooperation procedure in the CSO-based service provisioning model.

strata in this model. The AC could be shared among the application and network strata and make the best use of the cross-stratum and cross-domain resources. In contrast, each TC usually runs to fulfill intra-domain routing and spectrum allocation and control in its own domain.

In Fig. 3, the AC is used to monitor the application resource utilization from data centers (steps 1–2). Once the SD-OTN first accesses the network, the handshake message of a TCP session between the TC and SD-OTN can be inter-worked to set up TCP connection preparing for the OpenFlow control (step 3). Each corresponding SD-OTN sends the OpenFlow hello message to the TC actively at the initial phase, while obtaining the acknowledgment hello message replying from TC (steps 4–5). TC sends the node monitor request to each SD-OTN in its own domain using *OFPT_FEATURES_REQUEST* message periodically (step 6). The corresponding SD-OTN responds the node and port information with *OFPT_FEATURES_REPLY* and *OFPT_PORT_STATUS* messages, respectively so that the TC can consolidate node information to discover the overall network node and topology information (steps 7–8). If a new request of the data center service arrives, the CSO-enabled scheme can

be completed in AC and choose the optimal destination node after receiving the CSO request from the TC of the source domain, and then AC sends the result and application resource to related TCs (steps 9–10). TCs receive the CSO reply and compute a path with optimal modulation format in the network according to the transmission distance and optical network bandwidth information, then proceed to set up/adjust an end-to-end SDP by controlling the corresponding SD-OTN using the *OFPT_FLOW_MOD* message (step 11). The elastic optical network can be needed to guarantee the end-to-end QoS. All corresponding SD-OTNs should report the setup status to TC through an *OFPT_PACKET_IN* message (step 12). When TCs obtain a setup success reply from the last SD-OTN, the data center service can be provisioned to the chosen destination node with the optimal modulation format for utilizing the application and spectrum resources effectively (step 13). After that, the network information occupancy in TCs can be updated to keep synchronization by receiving the update message from the corresponding SD-OTN via an *OFPT_PORT_STATUS* message (step 14). Note that these existing OpenFlow messages are reused to simplify our implementation.

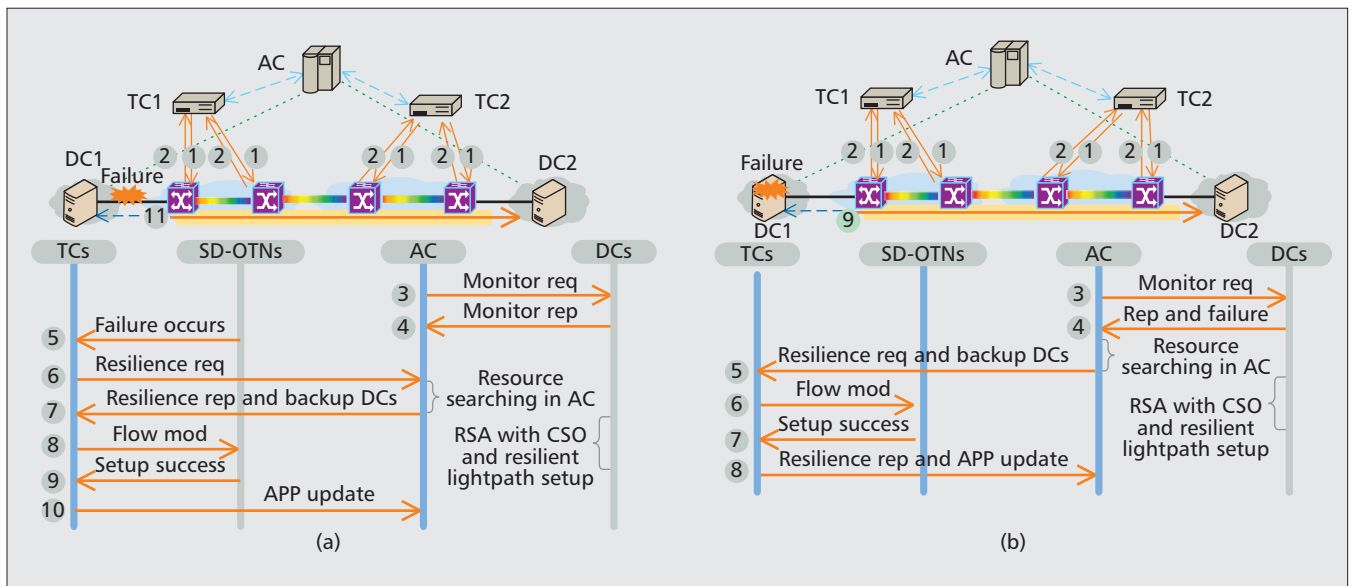


Figure 4. Cooperation procedure in cross stratum service resilience model when (a) inter-domain link failure occurs; b) data center domain failure occurs.

CROSS STRATUM SERVICE RESILIENCE MODEL

The cross stratum service resilience model in CSO architecture could ensure the high-performance QoS of user demands after a disaster. We discuss two typical failure cases in the analysis of cooperation procedure for simplicity (i.e., inter-domain link and data center domain failure), the latter of which simulates the regional disaster situation. In Fig. 4a, the TCs detect the status of optical nodes through interworking with SD-OTNs (steps 1–2), while the AC monitors the application resource information from data centers (steps 3–4). In case there is a failure in the inter-domain link, the service provided by the server in the same data center domain will fail to be offered to the user. Due to the automatic discovery mechanism based on the link manager protocol (LMP), the adjacent optical node finds the link failure and reports the information to the TC (step 5). With the analysis in the network stratum, the TC escalates the resilience of the data center application for a possible change of the resource origin and forwards the request to the AC in turn (step 6). The AC receives the resilience request and searches for alternative data center resources and locations in the same service type, then responds the resilience reply to the TC with the backup data center servers (step 7). The resilient path, considering the CSO of network and application resources, can be computed in TCs, then proceed to establish an end-to-end SDP by controlling the corresponding SD-OTNs (step 8). When TCs obtain a setup success reply from the last SD-OTN, the service can be recovered utilizing the cross stratum resources effectively. After that, the application usage in the AC can be updated to keep synchronization by receiving the update message from TCs (steps 9–11).

In Fig. 4b, the TCs monitor the network information of optical nodes by the interaction with SD-OTNs (steps 1–2). As the application stratum resources are monitored and maintained

in the AC (steps 3–4), in case of data center domain failure, the AC discovers it quickly and looks for alternative servers in the same host location. If alternative resources are only available in remote locations, the AC provides such information to the network stratum for possible connectivity change of the SDP in order to guarantee the QoS (step 5). After receiving the resilience request with application information, the TCs compute the lightpath with the CSO-enabled scheme and assign the spectrum at the corresponding SD-OTN (steps 6–7). Then the TCs update the data center application utilization to the AC after successfully provisioning a resilient lightpath (steps 8–9).

The control mechanism presented above is also applicable for VM migration. Essentially, the mechanism for VM migration is similar to the one for service resilience after a failure. The control patterns of VM migration and service resilience in case of data center failure are the same, because they need the destination node choice through both the network and application strata. Note that the signaling latency in the control plane is a relatively short portion compared to the whole latency of migration, which depends on the data volume of a VM and the distance between the source and destination of VM migration. The service resilience is triggered by link or data center failure. Different from service resilience, VM migration can be performed based on event-based trigger (e.g., maintenance management and data backup demand).

SERVICE PROVISIONING SCHEMES IN ELASTIC DATA CENTER OPTICAL INTERCONNECTION

Three service provisioning schemes are proposed based on the CSO architecture in the elastic data center optical interconnection: modulation format adaptive adjustment (MFA), application

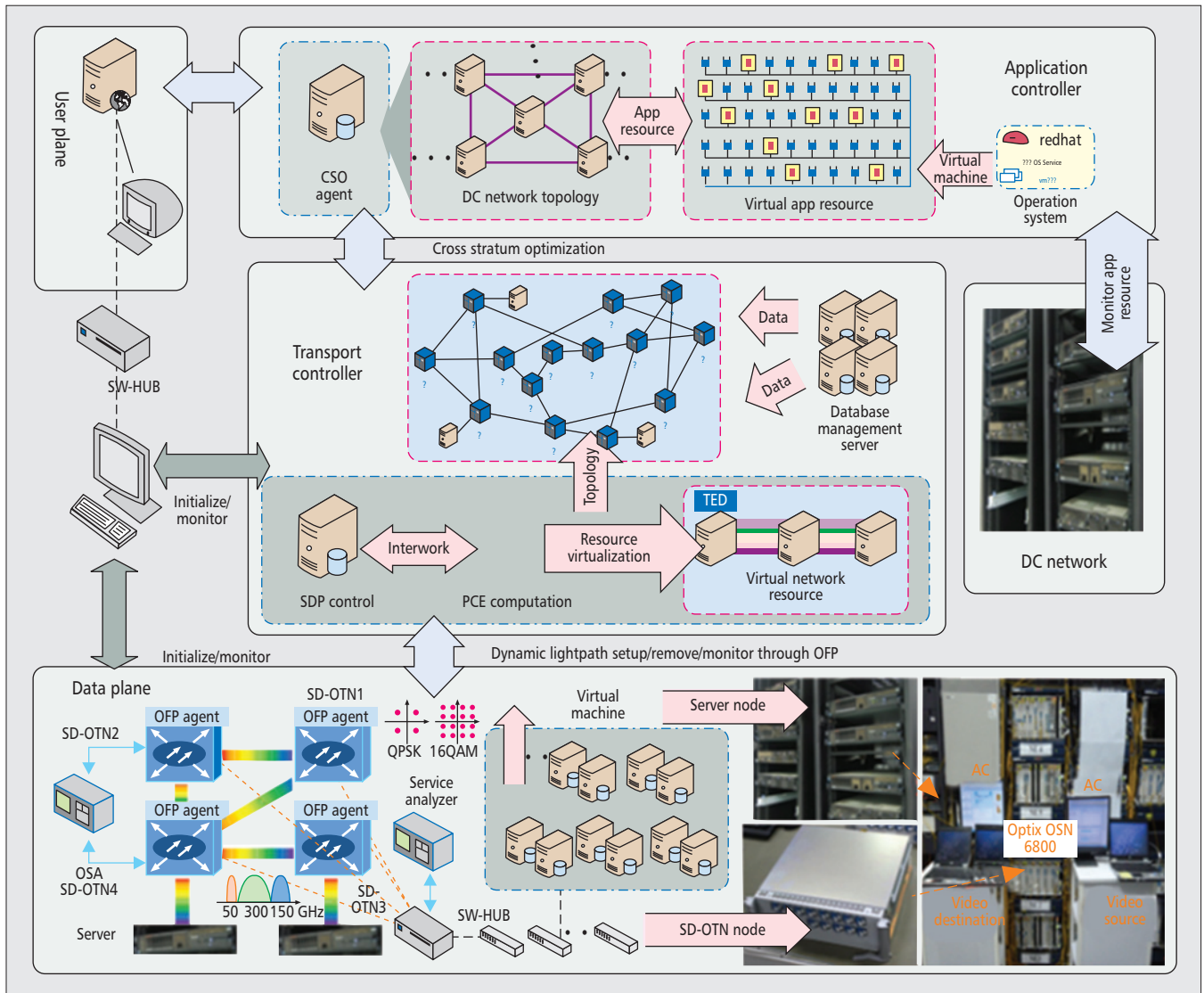


Figure 5. Experimental testbed and demonstrator setup.

load balancing (ALB), and CSO-enabled dynamic global load balancing (CSO-DGLB) schemes. The description is as follows.

In the MFA scheme, the AC selects the data center destination that has the path with the minimum distance from the source to destination. We use the hierarchical routing scheme to calculate the end-to-end path due to realization simplification. Each TC from a related domain can abstract the physical optical network and send the abstracted information to the AC. When a path computation request arrives at a TC, the TC will send the request to the AC if it finds that the destination is not in the local domain. Then the AC can calculate a TC sequence according to the abstracted inter-domain topology obtained from the TCs. After that, the AC will notify the involved TCs to compute the k -shortest paths in each domain. The end-to-end path computation can be chosen using Dijkstra's algorithm in the AC from the candidate paths, which gather all the computed paths from each TC. Furthermore, the modulation format of service could be chosen based on

the transmission distance of spectrum path and optical network bandwidth information. For the ALB scheme, according to the storage and CPU usage, the AC chooses the data center server node with the minimum application utilization as the destination to keep the balance of server load. The CSO-DGLB scheme could comprehensively consider the status of data center application resources (e.g., CPU and storage utilization) and elastic optical network resources (e.g., bandwidth and modulation format) for CSO. The AC selects the server node and data center location as the destination based on the application status collected from data centers and network condition provided by TCs dynamically. Receiving a service request and pairs of source and destination node from the AC, TCs can complete the end-to-end path computation in the connection and service parameter constraints, then select the appropriate modulation format according to the SDP distance, and perform spectrum assignment with first fit for the computed path and the lightpath provisioning by OFP.

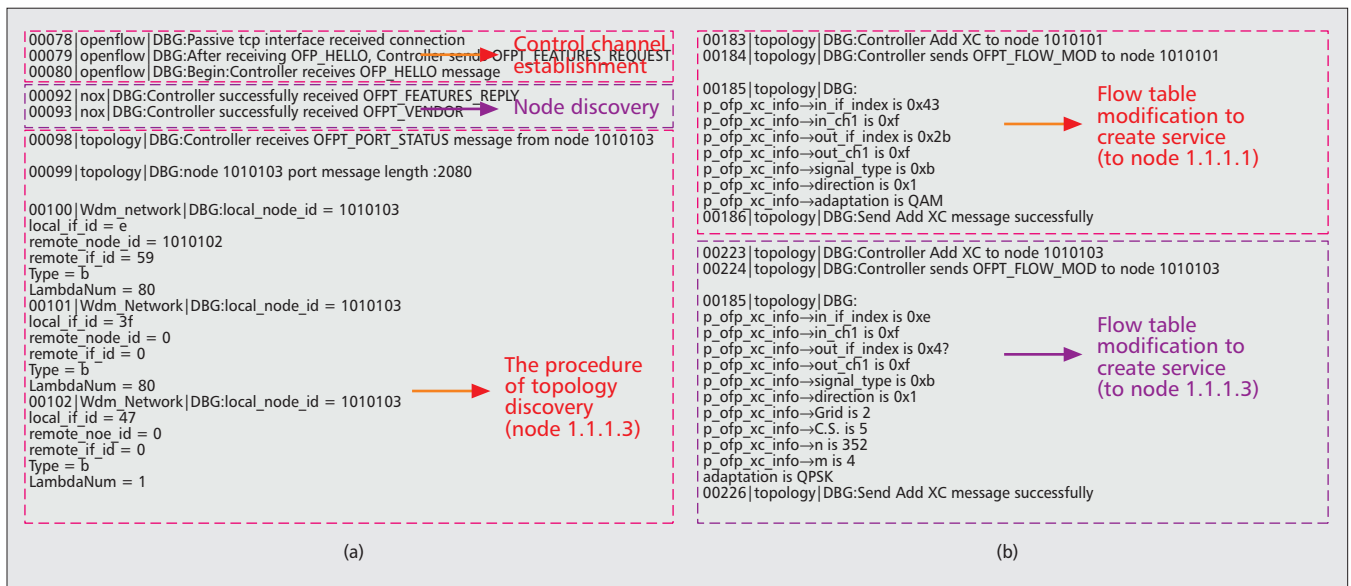


Figure 6. a) The capture of the OpenFlow messages for automatic network discovery; b) the capture of extended flow table message for SDP setup.

Items	Scheme processing	OFPT propagation	Device setup	Device release	Device adjust	Device hardware
Latency	~2.5 ms	~0.5 ms	~55 ms	~20 ms	~50 ms	~100 ms

Table 2. SDP setup/release/adjustment latency.

EXPERIMENTAL SETUP AND RESULTS DISCUSSION

The feasibility and efficiency of the proposed CSO architecture is evaluated on our OaaS testbed, which consists of control and data planes as shown in Fig. 5. In the data plane, four OpenFlow-enabled elastic optical nodes are equipped with Huawei Optix OSN 6800, which can switch or transport the elastic signal in optical networks. We develop a software OFP agent according to the API function to control the hardware of flex ROADM and ODU boards through OFP. Due to the short transmission distance, the SD-Rs are not deployed in the actual experiments. Data centers and the other nodes are implemented on an array of VMs in IBM X3650 servers. Since each VM has an operating system, its own CPU, and storage resource, it can be considered a real node. Therefore, system virtualization technology makes it easy to set up experiment topology based on NSFNet. Three data centers are deployed in various nodes of the experiment topology. In each data center, 10 VMs are used to accommodate and provide the data center application. For the control plane, TCs are assigned to support the architecture, and are deployed in three servers for elastic spectrum control, PCE computation, and resource abstraction, while the database management servers are responsible for maintaining a traffic engineering database, management information base, and configuration. Several controllers can realize the function of SDN, such

as NOX, Opendaylight, and Ryu. In our work, the SDN controller for the TC is implemented by NOX due to the realization simplification in the testbed. The NOX is deployed in each TC internally for spectrum control. Also, the PCE is integrated within the transport controller and interworks with NOX to realize the path computation and SDN-enabled spectrum control. The AC server is used as the CSO agent and monitors the application resources from the data center with the CSO scheme. The user plane is deployed in a server and deploys the service information generator to implement batch data center services for the experiments.

We have designed and experimentally verified SDP provisioning in elastic optical networks based on CSO for data center service. The destination data center is determined by the AC with the CSO scheme based on various application utilizations and current network resources, and the SDP for the service is set up from source to destination node. In addition, the spectrum bandwidth and corresponding modulation format can be tunable according to different SDP distances. The end-to-end SDP setup/release/adjustment latency is measured through dozens of experiments as well as timing performance of the AC and TC with SD-OTN nodes, which comprises the scheme processing time of the controller, OFP propagation latency, and software and hardware device handling times for the sake of observation and analysis, as shown in Table 2. The scheme processing time of AC and OFP propagation time occupy an extremely limited portion of the overall latency, and are

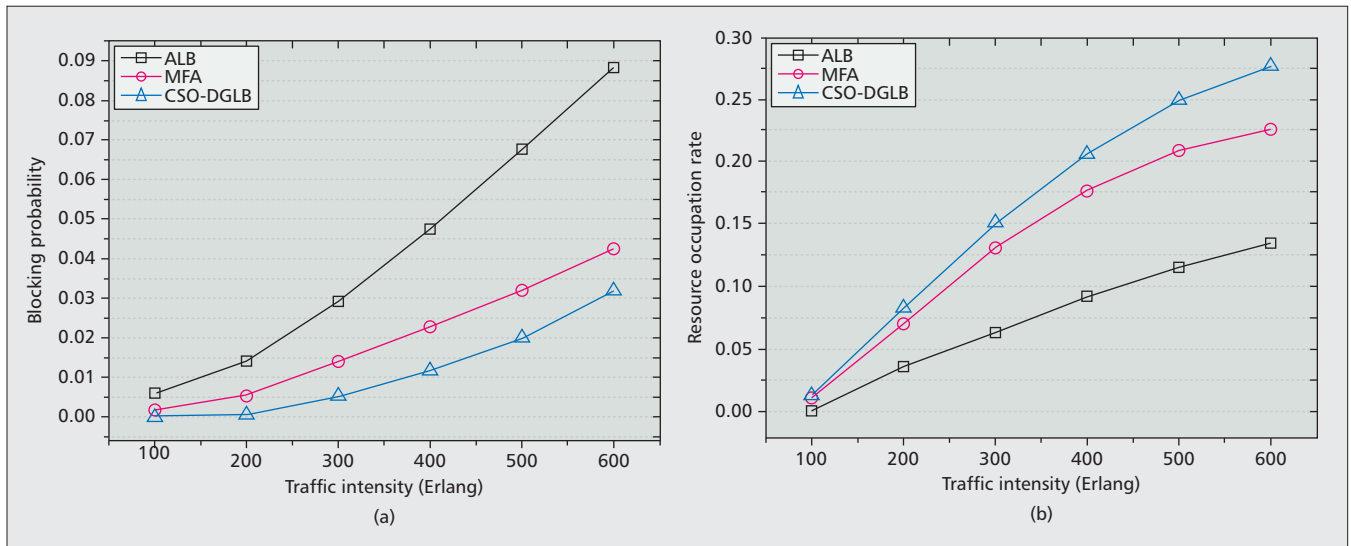


Figure 7. a) Blocking probability; b) resource occupation rate among the three schemes under a heavy traffic load scenario.

around 2.5 ms and 0.5 ms, respectively. The relatively key contributor is signal action time in the device including device setup time, release time, and adjustment time. The setup time (~ 55 ms) is a little higher than adjustment latency (~ 50 ms), since the equipment will respond a little slower with extra processing when a new SDP needs to be established. The release latency (~ 20 ms) is much lower than the latency of setup and adjustment, because the device does not need to reserve the spectrum resources when the established SDP will be released. The time of device hardware handling (~ 100 ms) is the main contributor to the overall latency on our testbed.

Figure 6 shows the capture of the OpenFlow message exchange for CSO through Wireshark deployed in the TC. When an SD-OTN connects to the network, the initial TCP session handshake is performed through the interface to establish the TCP connection, and then the SD-OTN sends the *OFPP_HELLO* message to the TC first. After receiving it, the TC responds with the *OFPP_HELLO* message to the corresponding SD-OTN, and the control channel between the TC and optical node is established. For the node information discovery, the TC periodically sends a monitor request to each SD-OTN using an *OFPT_FEATURES_REQUEST* message through extended OFP, while obtaining the status information of each one with an *OFPT_FEATURES_REPLY* message from them. After that, the overall network topology can be split jointly among TCs according to the neighboring port information through an *OFPT_PORT_STATUS* message. Since the all-optical nodes report the message simultaneously, Fig. 6a illustrates the capture of the OpenFlow message exchange of the SD-OTN node with 1.1.1.3 for the automatic network discovery procedure, which comprises control channel establishment, and node and topology discovery. After calculating the SDP with the CSO scheme, the TC proceeds to set up an end-to-end SDP by controlling all corresponding SD-OTNs along the computed path by using the *OFPT_FLOW_MOD* message. The

experimental results correspond to the procedures depicted in Fig. 3. Figure 6b shows a snapshot of the extended flow table modification message for SDP setup, which verifies the OPF extensions for CSO over elastic optical networks.

We also evaluate the performance of CSO under a heavy traffic load scenario, and compare the CSO-DGLB scheme with the MFA and ALB schemes through VMs. The traffic requests to the data center are established with spectrum randomly from 50 to 400 GHz, where the adjustable minimal frequency slot is 12.5 GHz. The needed application resource usage in data center is selected randomly from 1 to 0.1 percent for each service demand. They arrive at the network following a Poisson process, and results have been extracted through the generation of 100,000 demands per execution. Figure 7 compares the performances of three schemes in terms of blocking probability and resource occupation rate. The CSO-DGLB scheme reduces blocking probability more effectively than the MFA and ALB schemes, especially when the network is heavily loaded. The reason is that the CSO-DGLB scheme realizes global optimization considering both application and network resources integrally, and furthermore, on the basis of it, economizes the spectrum resource again by choosing a high-level modulation format according to the SDP distance. The MFA scheme outperforms the ALB in blocking probability. That is because the MFA scheme adjusts spectrum bandwidth through modulation format to save the doubled spectrum resource, and greatly increases the available resources for new arriving demands. In Fig. 7b, the CSO-DGLB scheme outperforms the other schemes in the resource occupation rate significantly. The main reason is that more resources can be occupied when the blocking probability is lower.

CONCLUSION

In order to meet the QoS requirements of data center services, this article presents a novel CSO architecture in elastic data center optical inter-

connection, which can allow global optimization and control across elastic optical transport network and data center application stratum heterogeneous resources to implement OaaS. The functional modules of the architecture and their cooperation procedure in CSO-based service provisioning and cross stratum service resilience modes are described and investigated. The performance of CSO is verified on our OaaS testbed for data center services. We evaluate its performance under heavy traffic load and compare it with the CSO-DGLB, MFA, and ALB schemes. Numerical results show that CSO with the CSO-DGLB scheme can utilize cross optical network and application stratum resources effectively, and dynamically adjust modulation format and bandwidth in elastic optical networks without increasing blocking probability.

Our future projects for CSO include several aspects. Multiple failures could be discussed for CSO architecture in the near future. The network function virtualization in the data center interconnect scenario with optical networks on our OaaS testbed should be studied in future work.

ACKNOWLEDGMENT

This work has been supported by NSFC project (61271189, 61201154), Ministry of Education-China Mobile Research Foundation (MCM20130132), the Fundamental Research Funds for the Central Universities (2015RC15), and Fund of State Key Laboratory of Information Photonics and Optical Communications (BUPT), P. R. China.

REFERENCES

[1] 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.

[2] W. Shieh, H. Bao, and Y. Tang, "Coherent Optical OFDM: Theory and Design," *Optics Express*, vol. 16, no. 2, Jan. 2008, pp. 841–59.

[3] S. Gringeri *et al.*, "Flexible Architectures for Optical Transport Nodes and Networks," *IEEE Commun. Mag.*, vol. 48, no. 7, July 2010, pp. 40–50.

[4] 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.

[5] E. Mannie, Ed., "Generalized Multi-protocol Label Switching (GMPLS) Architecture," IETF RFC 3945, Oct. 2004.

[6] S. Das, G. Parulkar, and N. McKeown, "Why OpenFlow/SDN Can Succeed Where GMPLS Failed," *Proc. ECOC*, paper Tu.1.D.1, Sept. 2012.

[7] H. Yang *et al.*, "Multi-Stratum Resource Integration for OpenFlow-based Data Center Interconnect [Invited]," *J. Opt. Commun. Net.*, vol. 5, no. 10, Oct. 2013, pp. A240–A248.

[8] M. Channegowda *et al.*, "Experimental Demonstration of an OpenFlow Based Software-Defined Optical Network Employing Packet, Fixed and Flexible DWDM Grid Technologies on an International Multi-domain Testbed," *Optics Express*, vol. 21, no. 5, Mar. 2013, pp. 5487–98.

[9] H. Yang *et al.*, "Global Resources Integrated Resilience for Software Defined Data Center Interconnection Based on IP Over Elastic Optical Network," *IEEE Commun. Letters*, vol. 18, no. 10, Oct. 2014, pp. 1735–38.

[10] R. Casellas *et al.*, "An Integrated Stateful PCE/OpenFlow Controller for the Control and Management of Flexi-Grid Optical Networks," *Proc. OFC/INFOEC*, paper OW4G.2, Mar. 2013.

[11] F. Paolucci *et al.*, "OpenFlow-Based Flexible Optical Networks with Enhanced Monitoring Functionalities," *Proc. ECOC*, paper Tu.1.D.5, Sept. 2012.

[12] L. Liu *et al.*, "OpenSlice: an OpenFlow-Based Control Plane for Spectrum Sliced Elastic Optical Path Networks," *Proc. ECOC*, paper Mo.2.D.3, Sept. 2012.

[13] H. Yang *et al.*, "Performance Evaluation of Time-Aware Enhanced Software Defined Networking (TeSDN) for Elastic Data Center Optical Interconnection," *Optics Express*, vol. 22, no. 15, July 2014, pp. 17630–43.

[14] H. Yang *et al.*, "Performance Evaluation of Multi-Stratum Resources Integrated Resilience for Software Defined Inter-Data Center Interconnect," *Optics Express*, vol. 23, no. 10, May 2015, pp. 13,384–98.

[15] J. Abley, A. Canada, and K. Lindqvist, "Operation of Anycast Services," IETF RFC 4786, Dec. 2006.

[16] 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/INFOEC*, paper PDP5B.1, Mar. 2013.

[17] A. Farrel, J. P. Vasseur, and J. Ash, "A Path Computation Element (PCE)-Based Architecture," IETF RFC 4655, Aug. 2006.

BIOGRAPHIES

HUI YANG (yanghui@bupt.edu.cn) is an assistant professor at the Institute of Information Photonics and Optical Communications at Beijing University of Posts and Telecommunications (BUPT). He received his Ph. D degree in communication and information systems from BUPT in 2014. His main research interests include network architecture, optical transport network, optical and wireless networks, software defined networks, and cross stratum optimization. He has authored or coauthored more than 40 papers in prestigious international journals and conferences.

JIE ZHANG is a professor and vice dean of the Institute of Information Photonics and Optical Communications at BUPT. He is sponsored for over 10 projects of the Chinese government. He has published eight books and more than 100 articles. Seven patents have also been granted. He has served as a TPC member for ACP '09, PS '09, and ONDM '10, among others. His research focuses on optical transport networks, packet transport networks, and more.

YONGLI ZHAO is currently a lecturer at 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 from BUPT in 2005 and 2010, respectively. More than 100 articles have been published. His research focuses on wavelength switched optical networks, optical transport networks, and packet transport networks.

YUEFENG JI is a professor and dean of the Institute of Information Photonics and Optical Communication at BUPT. He is also an expert in the Communication Technology Group of the National High-Tech Research and Development Program (863 Program) of China. His research interests are primarily in the areas of optical fiber communication and broadband information networking.

JIANRUI HAN is a system engineer in the Advanced Technology Department of the Wireline Network Business Unit and a leader in providing next generation telecommunications networks of Huawei Technologies, Ltd., Co. She joined Huawei in 2001. Her research interests are GMPLS control plane, wavelength switched optical networks, and transport software defined networks.

YI LIN is a research engineer at Huawei Technologies, Ltd., Co. He received his B.S. degree in electronic information science and technology in 2005 and his M.S. degree in radio physics in 2007 from Sun Yat-Sen University, and joined Huawei in 2007. His main research topic is intelligent control of transport networks, including ASON/GMPLS, PCE, transport SDN, and so on.

YOUNG LEE received his B.A. degree in applied mathematics from the University of California at Berkeley in 1986, his M.S. degree in operations research from Stanford University, California, in 1987, and his Ph.D. degree in decision sciences and engineering systems from Rensselaer Polytechnic Institute, Troy, New York, in 1996. He is currently a principal technologist at Huawei Technologies USA Research Center, Plano, Texas. He is leading optical transport control plane technology research and development.

Numerical results show that CSO with the CSO-DGLB scheme can utilize cross optical network and application stratum resources effectively, and dynamically adjust modulation format and bandwidth in elastic optical networks without increasing blocking probability.

Optical Interconnects at the Top of the Rack for Energy-Efficient Data Centers

Jiajia Chen, Yu Gong, Matteo Fiorani, and Slavisa Aleksic

ABSTRACT

The growing popularity of cloud and multi-media services is dramatically increasing the traffic volume that each data center needs to handle. This is driving the demand for highly scalable, flexible, and energy-efficient networks inside data centers, in particular for the edge tier, which requires a large number of interconnects and consumes the dominant part of the overall power. Optical fiber communication is widely recognized as the highest energy- and cost-efficient technique to offer ultra-large capacity for telecommunication networks. It has also been considered as a promising transmission technology for future data center applications. Taking into account the characteristics of the traffic generated by the servers, such as locality, multicast, dynamicity, and burstiness, the emphasis of the research on data center networks has to be put on architectures that leverage optical transport to the greatest possible extent. However, no feasible solution based on optical switching is available so far for handling the data center traffic at the edge tier. Therefore, apart from conventional optical switching, we investigate a completely different paradigm, passive optical interconnects, and aim to explore the possibility for optical interconnects at the top of the rack. In this article, we present three major types of passive optical interconnects and carry out a performance assessment with respect to the ability to host data center traffic, scalability, optical power budget, complexity of the required interface, cost, and energy consumption. Our results have verified that the investigated passive optical interconnects can achieve a significant reduction of power consumption and maintain cost at a similar level compared to its electronic counterpart. Furthermore, several research directions on passive optical interconnects have been pointed out for future green data centers.

Jiajia Chen and Matteo Fiorani are with KTH Royal Institute of Technology.

Yu Gong is with Zhejiang University.

Slavisa Aleksic is with Vienna University of Technology.

¹ Facebook Luleå Data Center, <https://www.facebook.com/notes/lule-percentageA5-data-center/lule-percentageA5-goes-live/474321655969861>.

INTRODUCTION

To serve the “networked society,” a rapidly growing amount of data is stored, processed, computed, transmitted, and instantly made available upon request. Users of these data range from large organizations (e.g., big enterprises, government, universities) to individual customers

who increasingly rely on the services offered by data centers. According to Cisco [1], the amount of traffic handled by data centers was already 2.6 zettabytes (1 zettabyte = 10^{21} bytes) per year in 2012 and is expected to reach 6.4 zettabytes per year in 2016, which represents a 25 percent compound annual growth rate (CAGR). It is anticipated that the required bandwidth for large-scale data centers will grow 20 times every 4 years (i.e., by 111 percent CAGR, with 1 Pb/s in 2012, 20 Pb/s in 2016, and 400 Pb/s in 2020) [2]. On the other hand, due to the thermal dissipation problem, the power consumption that can be afforded by the network equipment in data centers is only allowed to increase at a much lower rate (19 percent CAGR, e.g., 0.5 MW in 2012, 1 MW in 2016, and 2 MW in 2020, is acceptable for the network equipment within a data center [2]). Obviously, keeping business as usual cannot sustain the future data center traffic.

The consumers of energy in data centers are information technology (IT) equipment (e.g., servers, network equipment) and other supporting facilities (e.g., lighting, cooling). In order to identify how efficiently a data center uses its power, a measure called power usage effectiveness (PUE) is defined as the ratio of the total facility power to the IT equipment power. Much effort has been applied toward reducing PUE. For instance, a smart selection of data center location could greatly reduce the energy required for cooling and significantly improve PUE. It was very recently reported that Facebook carefully chose the location and launched an Arctic data center (consisting of three 28,000 m² buildings) in Sweden. By utilizing icy conditions in the Arctic Circle, the data center can reach a PUE around 1.07.¹ Such a low level of PUE implies that in modern data centers the major focus on energy savings should be moved to IT equipment. Currently, network equipment in a data center may take up to approximately 15 percent of the total energy, and this value is expected to grow in the future [3]. Thus, in order to sustainably handle ever increasing traffic demand, it is of extreme importance to address the energy consumption issue in data center networks (DCNs), which provide interconnections among different servers within a data center as well as interfaces to the Internet.

Optical fiber communication is by far the

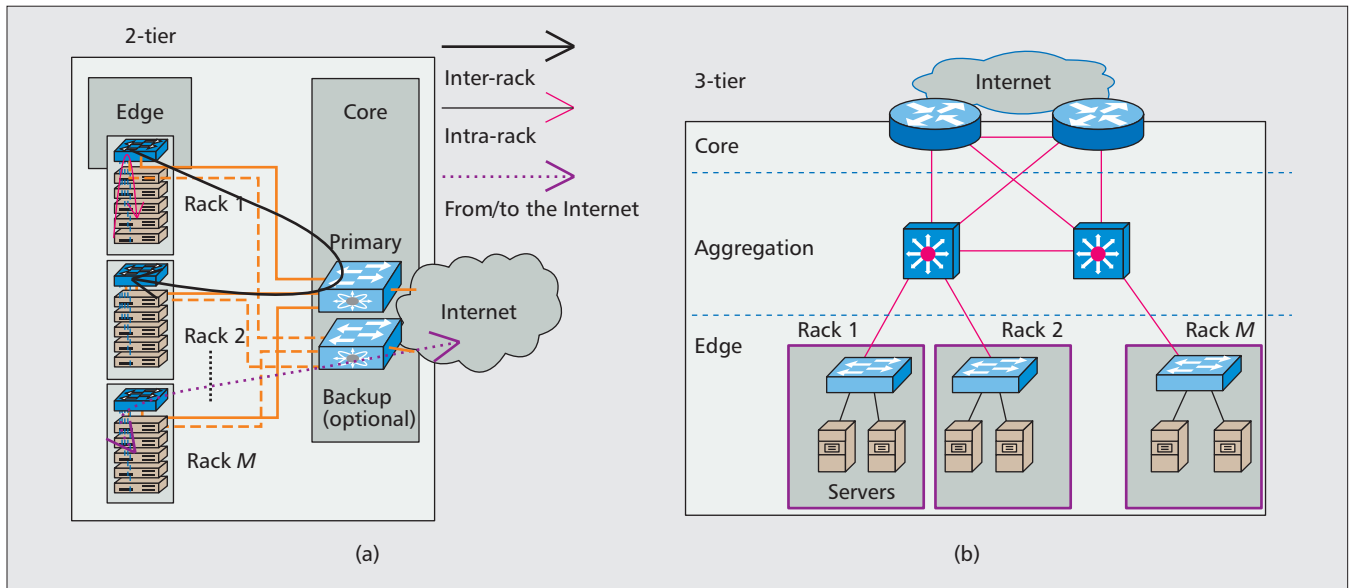


Figure 1. Data center network architectures: a) two-tier; b) three-tier.

least energy-consuming and least costly technique to offer ultra-high bandwidth for telecommunication [4]. It has also been considered as a promising transmission technology for data center applications [2]. However, in current DCNs, switching is still done in the electronic domain by the commodity switches, which consume an extensive amount of power and also cause a bottleneck for capacity upgrade. To reduce or eliminate the electronic components in DCNs, many architectures [2] based on optical switching have been proposed in the literature.

The existing solutions of optical switches for DCNs can be divided into two major categories [2]: hybrid electronic/optical switching and all-optical switching. For hybrid switching (e.g., in [5]), the electronic part deals with switching at the packet level, maintaining fine switching granularity while the optical circuit switching (OCS)-based part offers high capacity. However, scalability is its vital problem due to the lack of efficient solution to upgrade capacity of the electronic part. The scalability is not a problem in the second category, which is based on all-optical switching. This category can be further divided into two groups based on whether optical packet switching is employed or not. Schemes without optical packet switching (OPS) (e.g., in [6]) suffer from poor switching granularity, so the bandwidth utilization can be low, particularly if the traffic demands vary significantly. OPS could greatly enhance the switching granularity in DCNs, but this technology still has several fundamental problems, especially the lack of flexible optical buffering and signal processing technologies for contention resolution. Therefore, several proposed architectures involving OPS, e.g., [7] still need buffering and signal processing in the electronic domain. The extra optical-electrical (OE) and electrical-optical (EO) conversions significantly increase power consumption and introduce limitations for capacity upgrade.

Typically, DCNs include several tiers. For instance, two-tier DCNs (Fig. 1a) include two

stages, edge and core, while in three-tier architectures (Fig. 1b), an aggregation tier is introduced as a middle stage between the core and edge levels. For high scalability, DCNs can even have four tiers or more, where the aggregation tier can be further extended to more than one stage. The majority of the research efforts on optical interconnects so far (e.g., the aforementioned work [2, 5–7]) have been focusing on core/aggregation switches (i.e., switching among different racks). It should be noted that the power consumed by the switches at the edge tier (i.e., at the top of the rack, ToR), which interconnect the servers located in the same rack, is dominant because of a huge number of ToR switches. It can reach up to 90 percent of the total power consumed by all types of switches in DCNs [8]. Therefore, improving energy efficiency at the ToR level should be considered in the first place in order to decrease the overall DCN power consumption. However, there have been very few studies on this problem. With this in mind, in this article we concentrate on the edge tier of DCNs and explore whether optical interconnects are feasible at ToR in order to realize energy-efficient data center solutions. We investigate a few possibilities and carry out a performance assessment of different candidates for ToR interconnects.

The remainder of this article is organized as follows. We discuss the feasible options for optical interconnects at the edge tier of DCNs. First, traffic characteristics of the servers in data centers are summarized, which are of importance to be considered in the design of the interconnects at the ToR. Then we compare active vs. passive optical network solutions and attempt to identify the most promising ones for the edge tier of DCNs. We present several types of optical interconnects that utilize passive components to connect servers within a single rack. We carry out an evaluation of all the considered schemes. Finally, conclusions along with some points for future research directions are drawn.

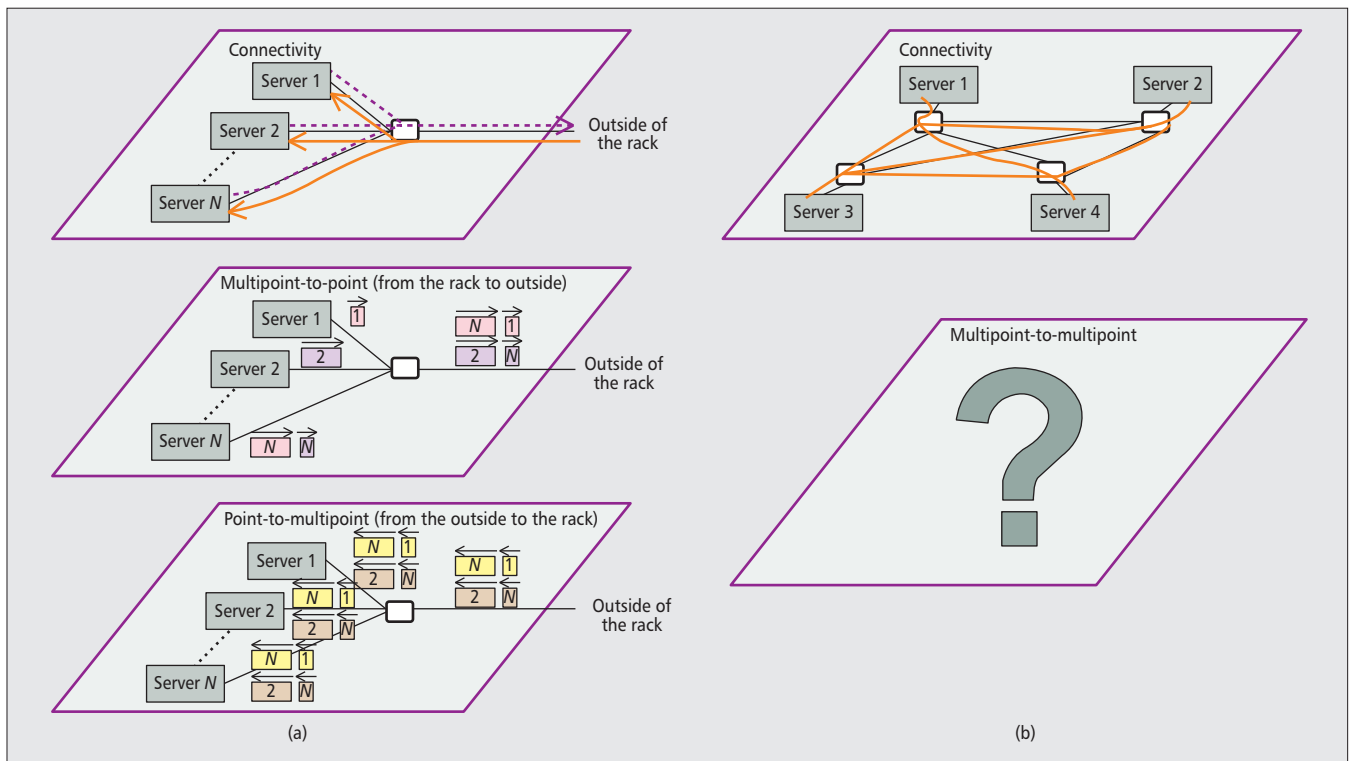


Figure 2. Two traffic scenarios at ToR in data centers: (a) traffic from/to outside of the rack; (b) intra-rack traffic.

DISCUSSION OF THE FEASIBILITY OF OPTICAL INTERCONNECTS AT THE TOR

In this section, the feasibility of optical interconnects at the ToR is discussed. First, we present the characteristics of the traffic generated by the servers; second, we highlight the advantages of passive over active optical interconnects at the ToR; finally, the main issues in designing passive optical interconnects for intra-rack communications are introduced.

TRAFFIC CHARACTERISTICS OF SERVERS IN DATA CENTERS

In order to determine whether and which optical technologies are feasible for interconnects at the ToR, it is of great importance to understand the characteristics of the traffic generated by the servers and its specific requirements. A number of key points from the literature [2, 9, 10] are summarized as follows:

Traffic Locality: The ratio of the intra-rack traffic to the total traffic can vary significantly depending on the application. For instance, in data centers that are used for cloud computing, the majority of the traffic (up to 80 percent) is intra-rack, while in the case of educational organizations and private enterprises, the portion of intra-rack traffic ranges from 10 to 40 percent [9]. Therefore, high flexibility is of prime concern for the design of interconnects at the ToR in order to meet the varied capacity requirements for intra-rack traffic generated by different applications.

Multicast Capability: Concurrent traffic flows

per server (around 10 in the majority of current data center servers [9]), which may have the same content but different destinations (e.g., for parallel computing purposes), can benefit from multicasting.

Variable Flow Capacity: Many applications (e.g., science computing) pose challenges in terms of elastic resource allocation, as the amount of required bandwidth per flow may greatly fluctuate from 1 to 40 Gb/s (or even higher) [10]. Although the average utilization is low, 1 Gb/s is insufficient for satisfying the peak traffic requirement and might not sustain future traffic demands.

Burstiness: The traffic inside the rack is very bursty [9]. Therefore, dynamic bandwidth allocation should be supported for intra-rack communication in order to improve utilization while offering sufficient capacity for the traffic peaks.

PASSIVE VS. ACTIVE

To satisfy all the aforementioned traffic requirements, it is extremely difficult to realize optical interconnects based on any of the existing optical switching technologies. OPS still has fundamental technical problems to achieve flexible buffering and signal processing in the optical domain. Therefore, extra OE and EO conversions are typically needed in OPS, which result in a significant increase in both power consumption and cost. Meanwhile, OCS has very coarse granularity, so it can hardly accommodate very busy traffic at the edge tier of DCNs. Most probably, these are the main reasons why no research progress has been reported so far on optical switching at the ToR.

Therefore, a completely different strategy has to be considered in contrast to optical switching

in order to explore a feasible approach for the optical interconnects at ToR. As a very successful solution for access networks, a passive optical network (PON) only uses passive components (e.g., splitters/combiners, arrayed waveguide gratings — AWGs) to interconnect fiber ports in the field. It has been proven that PON outperforms many other schemes such as point-to-point active optical network (AON) in terms of energy and cost efficiency [11]. Besides, it is also widely known that compared to AON, where an electronic switch is employed to interconnect different users, PON has higher reliability and is easier to maintain due to its passive nature. The products for 40 Gb/s PON have already been demonstrated by several vendors (e.g., Ericsson, Huawei), while PONs with higher speed (e.g., 100 Gb/s) are under development and expected to be ready in the near future. In this regard, passive solutions could be more appropriate for optical interconnect at the edge tier of DCNs than active-switching-based ones.

ISSUES IN PASSIVE OPTICAL INTERCONNECTS

In a PON, all users do not have direct connection to each other; instead, all of them need to first connect to a central node, that is, an optical line terminal (OLT). Therefore, it has a multipoint-to-point structure for upstream and a point-to-multipoint one for downstream. In contrast to PON, interconnects at the ToR are even more complicated. They are required to support two traffic scenarios (Fig. 2):

- Traffic to/from the outside of the rack (including inter-rack communication within the data center and that to connect to the outside of the data center)
- Intra-rack communication, where each server should be able to freely connect to all the others in the same rack

The same way in which a PON handles upstream and downstream traffic can be directly applied in the first scenario. However, the second scenario could be extremely difficult to handle because a multipoint-to-multipoint communication is required. Meanwhile, it should be noted that these two scenarios must be coordinated with each other as a common ToR structure should deal with both cases at the same time.

PASSIVE OPTICAL INTERCONNECTS AT THE TOR

In this section, we present different kinds of passive optical interconnects (POIs), which are considered as candidates for interconnects at the edge tier of DCNs to tackle both of the aforementioned traffic scenarios. They are categorized according to the types of passive components used for interconnection at the ToR. In all the schemes, each server includes an optical interface (OI) that can send the traffic to different destinations (i.e., to another server within the same rack or outside of the rack). A rack controller (in-band or out-of-band, such as the one presented in [14]) is needed in the control plane to solve the contention issue for all three presented types of schemes. A proper media access control (MAC) protocol and an

efficient resource allocation algorithm are required to sustain the traffic grooming (i.e., dealing with the packet layer). For the MAC protocol, there is a classical one widely considered in Ethernet networks, IEEE 802.3, carrier sense multiple access with collision detection (CSMA/CD). A similar concept could be applied to our architectures. On the other hand, new MAC protocols and tailored resource allocation strategies are required for high resource utilization and flexibility.

It should also be noted that all the presented POIs are self-contained. They could connect to any type of core/aggregation switch (e.g., the commodity switch or the ones proposed in [5–7]) through a proper interface (e.g., OE conversion). In order to achieve high energy efficiency, an all-optical data center network as proposed in [14] can be considered. Regarding the control plane, several tiers (e.g., proposed in [14]) could be applied; the core/aggregation switch controller is responsible for the resource allocated to inter-rack traffic as well as that to/from the Internet while the rack controller is responsible for the traffic within the rack.

SCHEME 1: $N \times N$ AWG-BASED POI

Figure 3a shows an example of the first type of POIs, which only use arrayed waveguide grating (AWG) to connect different ports within a rack [12]. Thanks to the cyclic property of the $N \times N$ AWG, a proper wavelength plan (Fig. 3b) can be made for each server to set up a connection to any other server as well as to the outside of the rack without any conflicts of spectrum. The cyclic property can guarantee that the signals carried by the same wavelength but injected into different input ports of the $N \times N$ AWG are distributed among the distinct output ports in a non-overlapping manner. Furthermore, a typical insertion loss of AWG is around 3–5 dB and not that related to the total number of ports (i.e., N). The major concern regarding scalability is related to the number of wavelengths needed for the whole rack. N ports of AWG correspond to at least N wavelengths. Meanwhile, its optical interfaces to the core/aggregation switch should have the ability to deal with multiple wavelengths (i.e., up to $N - 1$ channels). It should be also noted that connecting a pair of OI transceivers to both sides of the $N \times N$ AWG (Fig. 3a) inherently provides 1 + 1 or 1:1 protection, thereby leading to high reliability.

In the example shown in Fig. 3a, only one port at each side of the AWG is reserved for the traffic to/from the outside of the rack. For a general case, k ($1 \leq k < N$) ports at each side of the $N \times N$ AWG can be used for the traffic to/from the core/aggregation tier, while the remaining $N - k$ ports are used to connect servers within the rack. However, to fully utilize ports in the case of $k > 1$, multiple transceivers are required at each OI.

SCHEME 2: AWG + COUPLER-BASED POI

The second type of optical interconnects consider using both AWGs and couplers at the edge tier of DCNs. Here, a wavelength-independent coupler is used as either a combiner

In order to determine whether and which optical technologies are feasible for interconnects at the ToR, it is of great importance to understand the characteristics of the traffic generated by the servers and its specific requirements.

or splitter. Figure 3c shows an example of such a solution initially proposed in [13], where a wavelength label is assigned to each server (Fig. 3d). All signals sent from different servers in the same rack are first merged by the combiner and then routed to destination servers through the AWG. The total number of wavelengths needed for one rack is equal to the number of servers in the rack. Medium access control and bandwidth allocation need to be applied in order to make sure that the signals from different sources destined to the same server are received without interference. Furthermore, in this scheme the signals leaving the rack experience a lower insertion loss than the intra-rack traffic. The total insertion loss for intra-rack communication is highly dependent on the splitting ratio of the coupler, which becomes an important concern for scalability. For a rack with up to 64 servers, approximately 25 dB optical power budget is required for communication between each pair of servers within the rack [13]. A similar level of optical power budget is required for a typical PON configuration, and hence is still acceptable in a wide sense. Similar to scheme 1, this method also requires an interface to the core/aggregation switch that can handle multiple wavelengths from the same port.

SCHEME 3: COUPLER-BASED POI

The third considered scheme for optical interconnect only uses an optical coupler. Figure 3e shows an example of such a structure as proposed in [14], in which OIs are connected to N input ports of an $N \times 2$ coupler (here N represents the number of servers in one rack). By passing through a wavelength selective switch (WSS), channels $\lambda_1, \dots, \lambda_C$ belonging to the intra-rack communication are sent back to the coupler and broadcast to all the connected OIs in the same rack (Fig. 3f). In this way the multicast capability can easily be offered by the broadcast-and-select nature of the coupler. The channels assigned for the traffic to/from the outside of the rack, that is, channels $\lambda_{C+1}, \dots, \lambda_M$, are sent to (or received from) the outside of the rack. The WSS is able to support flexible channel allocation for the traffic within the rack and that from/to the core/aggregation switch. Additionally, using tunable transceivers in optical interfaces leads to even higher flexibility, but at the cost of increased complexity at both the component and control levels. Moreover, for intra-rack communication, the signals have to pass the coupler twice, which causes extremely high optical power loss, and hence can affect scalability significantly.

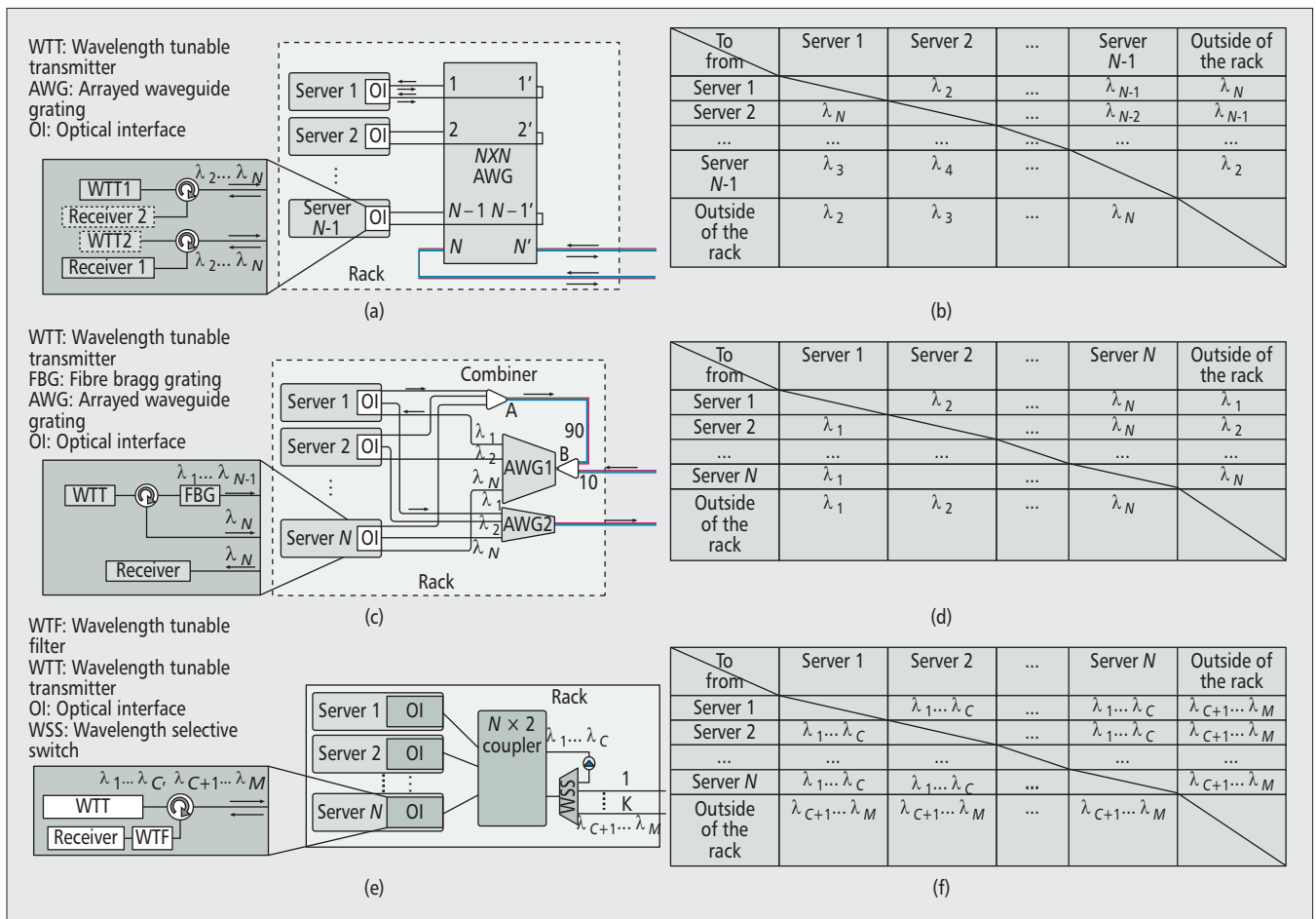


Figure 3. Block diagrams of different schemes of passive optical interconnects (POIs) and the corresponding wavelength plans: a) and b) for scheme 1 AWG-based POI; c) and d) for scheme 2 AWG+coupler-based POI; e) and f) for scheme 3 coupler-based POI.

PERFORMANCE COMPARISON

In this section, we carry out a performance assessment of the three aforementioned types of OIs at the ToR in terms of different aspects: ability to host data center traffic, scalability, optical power budget, and complexity of control protocol, as well as the required interface to the upper tier (e.g., core/aggregation switch). Besides, we also estimate their cost and energy consumption, and compare the presented OIs at the ToR with the conventional commodity switch.

Table 1 summarizes the main characteristics of the three considered schemes. Scheme 3, using an $N \times 2$ coupler, provides a number of benefits and clearly outperforms the other two schemes with respect to data center traffic requirements. It can easily support varied traffic locality for different applications by adapting the number of channels used for intra-rack connections. Due to the broadcast feature of the coupler, it can inherently support multicast in the optical domain. When employing wavelength tunable transponders (WTTs) and flex-grid WSSs [14], scheme 3 can offer great flexibility to vary the flow capacity. Bandwidth variable transponders (BVTs) provide the ability to tune the wavelength and change the number of spectral slots assigned to a flow. Thus, the capacity allocated to a flow can vary upon request from 1 to 100 Gb/s and beyond. On the other hand, due to the fact that the AWG is wavelength-sensitive and has a fixed grid, schemes 1 and 2 can hardly support elastic channel allocation. Besides, given that the OIs are equipped with burst-mode transceivers, scheme 3 is able to handle bursty traffic by assigning time slots with varied lengths. A similar approach can also be used in schemes 1 and 2. However, due to the fixed wavelength plan in schemes 1 and 2, the wavelength has to be changed if the destination is different, which may significantly affect bandwidth utilization because of the tuning time of the transmitter. On the other hand, the fixed wavelength plan makes design of the MAC protocol for schemes 1 and 2 relatively simple, as the MAC mechanism is only required to solve the conflicts in the time domain (i.e., to avoid signals from different destinations arriving at the same server simultaneously). For scheme 3, high flexibility is given at the expense of increased complexity of the control plane as both the spectrum and time domains should be taken into account. Regarding multicast capability, scheme 2 could also support it by reordering the AWG and the coupler (i.e., the signals pass the AWG and then are broadcast by the splitter). However, this feature is not offered by all the approaches belonging to this type of POI (e.g., the example shown in Fig. 3c cannot provide multicast).

The estimated power consumption and cost of the three considered options and a conventional ToR commodity switch are presented in Fig. 4 for a rack composed of up to 64 servers and varied capacities per server. In our analysis, we consider the power consumption and cost of all components of the interconnection network inside the rack, including network interfaces at servers and interconnects at ToR. The values and models used for the calculation of power

Features		Scheme 1	Scheme 2	Scheme 3
Flexibility to host traffic locality		--	--	++
Multicast capability (in optical domain)		--	O	++
Dynamic bandwidth allocation	Varied flow capacity	O	O	++
	Bursty traffic	O	O	++
Complexity of media access control protocol		+	+	O
Optical power budget		++	O	--
Scalability		+	+	-
Complexity of interface to core/aggregation tier		-	-	+
(Legend: ++ very good; + good; O medium; - poor; -- very poor)				

Table 1. Performance assessment.

consumption and cost are taken from [14, 15]. A summary of these values is reported in Table 2. We assume a basic cost unit (CU) as the cost of a 10 Gb/s electronic network interface card, which is estimated to be US\$150.00 [15].

Figure 4a shows that using POIs at the ToR makes it possible to achieve very large power savings compared to the conventional commodity-switch-based solution (power consumption reduction is attested to be at least by a factor of 6). The higher the data rate, the larger the reduction of power consumption can be achieved by POIs. On the other side, among the three presented types of OI architectures, scheme 1 with backup can provide relatively high reliability, but is the most power consuming due to using a pair of optical transceivers for protection. Scheme 3 is slightly more power consuming than scheme 2 and scheme 1 without backup, mainly because of the power consumption introduced by the wavelength tunable filter (WTF) and the flexgrid WSS. Note that due to the highly bursty nature of data center traffic, the power consumption of all the considered options can be further reduced by using sleep or low-power mode at optical interfaces. We include this as an important future research direction to explore different possibilities of green technologies tailored for the proposed POIs.

Figure 4b shows the cost results. With 10 Gb/s or higher guaranteed capacity, all three considered optical schemes at the ToR (without backup) are less expensive than the commodity-switch-based solution. This is because at high data rate, the electronic switch must be equipped with a large number of expensive network components for interconnection at the ToR. In contrast, the optical interconnects are based on a much simpler architecture, which makes use of mainly passive and low-cost components. Furthermore, these passive components are not sensitive to data rate, so their cost does not increase with capacity. Among all the optical architec-

We plan to develop a proper DBA that is able to provide sufficiently low latency and jitter together with high throughput. It should be noted that such an algorithm for the scheme purely based on a coupler needs to take care of both spectrum and time domains.

Component	Power (W)			Cost (CU)		
	10 G	40 G	100 G	10 G	40 G	100 G
Optical Interconnects at ToR						
Wavelength tunable transceiver (WTT)	1.5	2.5	8	1.3	12	30
Wavelength tunable filter (WTF)	1.5	1.5	1.5	3	3	3
Flexgrid wavelength selective switch (WSS)	15	15	15	8.3	8.3	8.3
Arrayed waveguide grating (AWG) per port	N/A	N/A	N/A	0.1	0.1	0.1
Fiber Bragg grating (FBG)	N/A	N/A	N/A	1	1	1
Circulator	N/A	N/A	N/A	0.7	0.7	0.7
Coupler	N/A	N/A	N/A	1.3	1.3	1.3
Isolator	N/A	N/A	N/A	3	3	3
Commodity-switch-based ToR solution						
Electronic ToR switch per port	9.4	15.6	50.1	3	9.2	23.1
Electronic network interface card (NIC)	1	2	4	1	3	7.5

Table 2. Power consumption and cost of the optical and electronic network equipment inside the rack [14, 15].

tures without any backup, scheme 3 is the most expensive one due to the relatively high cost of the WTF. At 40 and 100 Gb/s guaranteed capacities, the cost of the OIs is dominated by the WTTs, so the difference among the three architectures becomes less obvious. Scheme 1 with backup is always the most costly solution due to the fact that each server is equipped with two expensive optical interfaces, particularly for the cases with 100 Gb/s guaranteed capacity. On the other side, the other optical ToR architectures (scheme 1 without backup, scheme 2, and scheme 3) are still less expensive than the commodity-switch-based solution because of the simplified interconnect architecture at the ToR, which compensates the very high cost of the WTTs. It is expected that WTTs will become less expensive thanks to massive production, making the optical interconnects more beneficial from the economics point of view.

Figure 4c shows the power consumption of the overall data center network as the function of the total number of servers for three different configurations in order to verify the gain by using the proposed POI in terms of energy efficiency. The first configuration, referred to as an E/E network, is a conventional 3-tier architecture based on commodity switches. The second one, referred to as an O/E network, employs one of our presented POIs (i.e., scheme 3) at the edge tier, while the aggregation and core tiers are based on conventional electronic switches. Scheme 3 has been considered because it is the one that has the highest power consumption (excluding the backup), which can represent an upper bound. Finally, the last one, referred to as an O/O network, employs scheme 3 at the ToR and optical

switching architecture in the aggregation/core tier as proposed in [14]. We consider an oversubscription ratio of 10 (each server can transmit up to 100 Gb/s, but only 10 Gb/s is guaranteed when the data center is fully loaded), and the input data for core/aggregation tier refer to [14, 15]. Besides, the number of servers in the data center is varied from 24,000 to 72,000. The results show that the E/E network is the most power-consuming, regardless of the size of the data center. Replacing the conventional electronic ToR switches with the POIs can significantly decrease the power consumption. The best results are achieved by using the all-optical architecture (i.e., O/O network), which is able to reduce the overall power consumption by at least a factor of 10 with respect to the conventional solution that only uses commodity switches.

CONCLUSIONS AND FUTURE WORK

In this article, we introduce a completely new paradigm for designing intra-rack interconnects in data center networks that involves passive optical interconnects at the top of the rack. Optical interconnects at the ToR can significantly reduce the energy consumption of the edge tier, which is responsible for up to 90 percent of the total power consumed by the data center network.

We have categorized this novel paradigm for optical interconnects into three major types and performed a comparison taking into consideration different aspects such as flexibility, complexity, scalability, cost, power consumption, and multicasting capability. The results have clearly shown that passive optical interconnects can sig-

nificantly reduce power consumption while providing moderate cost saving compared to the existing solution based on the commodity electronic switch. It should be noted that the current energy consumption calculation does not consider any dynamic energy management scheme. Due to the highly bursty nature of data center traffic, there is a good opportunity to further reduce power consumption by using sleep or low-power mode at optical interfaces.

Passive optical interconnects for the ToR using AWG for multiplexing/demultiplexing can achieve good scalability but suffer in flexibility. The scheme purely based on a coupler can meet the requirements set by the host data center traffic thanks to its broadcast-and-select nature, but provides limited scalability due to high optical insertion loss. This trade-off between flexibility and scalability has to be further investigated in order to identify the best option of optical interconnects at the edge tier of data center networks.

Furthermore, proper MAC protocols and dynamic bandwidth allocation algorithms for passive optical interconnects at the ToR are greatly required in order to efficiently avoid conflicts at the receiver side. The proposed concept offers high energy efficiency, but it may affect network performance. Therefore, we plan to develop a proper DBA that is able to provide sufficiently low latency and jitter together with high throughput. It should be noted that such an algorithm for the scheme purely based on a coupler needs to take care of both spectrum and time domains.

ACKNOWLEDGMENT

The work described in this article was carried out with the support of the projects “Enabling Scalable and Sustainable Data Center Networks,” funded by Swedish Foundation of Strategic Research, and “Towards Flexible and Energy-Efficient Data Center Networks,” funded by Swedish Research Council.

REFERENCES

[1] Cisco Global Cloud Index: Forecast and Methodology, 2012–2017, White Paper, 2013.
 [2] C. Kachris, K. Kanonakis, and I. Tomkos, “Optical Interconnection Networks in Data Centers: Recent Trends and Future Challenges,” *IEEE Commun. Mag.*, vol. 14, no. 9, Sept. 2013, pp. 39–45.
 [3] P. Mahadevan et al., “On Energy Efficiency for Enterprise and Data Center Networks,” *IEEE Commun. Mag.*, vol. 49, no. 8, Aug. 2011, pp. 94–100.
 [4] D. C. Kilper et al., “Power Trends in Communication Networks,” *IEEE J. Selected Topics in Quantum Electronics*, vol. 17, no. 2, Mar.–Apr. 2011, pp. 275–84.
 [5] N. Farrington et al., “Helios: A Hybrid Electrical/Optical Switch Architecture for Modular Data Centers,” *ACM SIGCOMM*, 2010.
 [6] K. Chen et al., “OSA: An Optical Switching Architecture for Data Center Networks with Unprecedented Flexibility,” *IEEE/ACM Trans. Net.*, vol. 22, no. 2, Apr. 2014, pp. 498–511.
 [7] Y. Yin et al., “LIONS: An AWGR-Based Low-Latency Optical Switch for High-Performance Computing and Data Centers,” *IEEE J. Selected Topics in Quantum Electronics*, vol. 19, no. 2, Mar.–Apr. 2013, pp. 3600–3409.
 [8] R. Pries et al., “Power Consumption Analysis of Data Center Architectures,” *Green Commun. and Networking*, 2012.
 [9] S. Kandula et al., “The Nature of Data Center Traffic: Measurements & Analysis,” *ACM SIGCOMM Internet Measurement Conf.*, 2009.

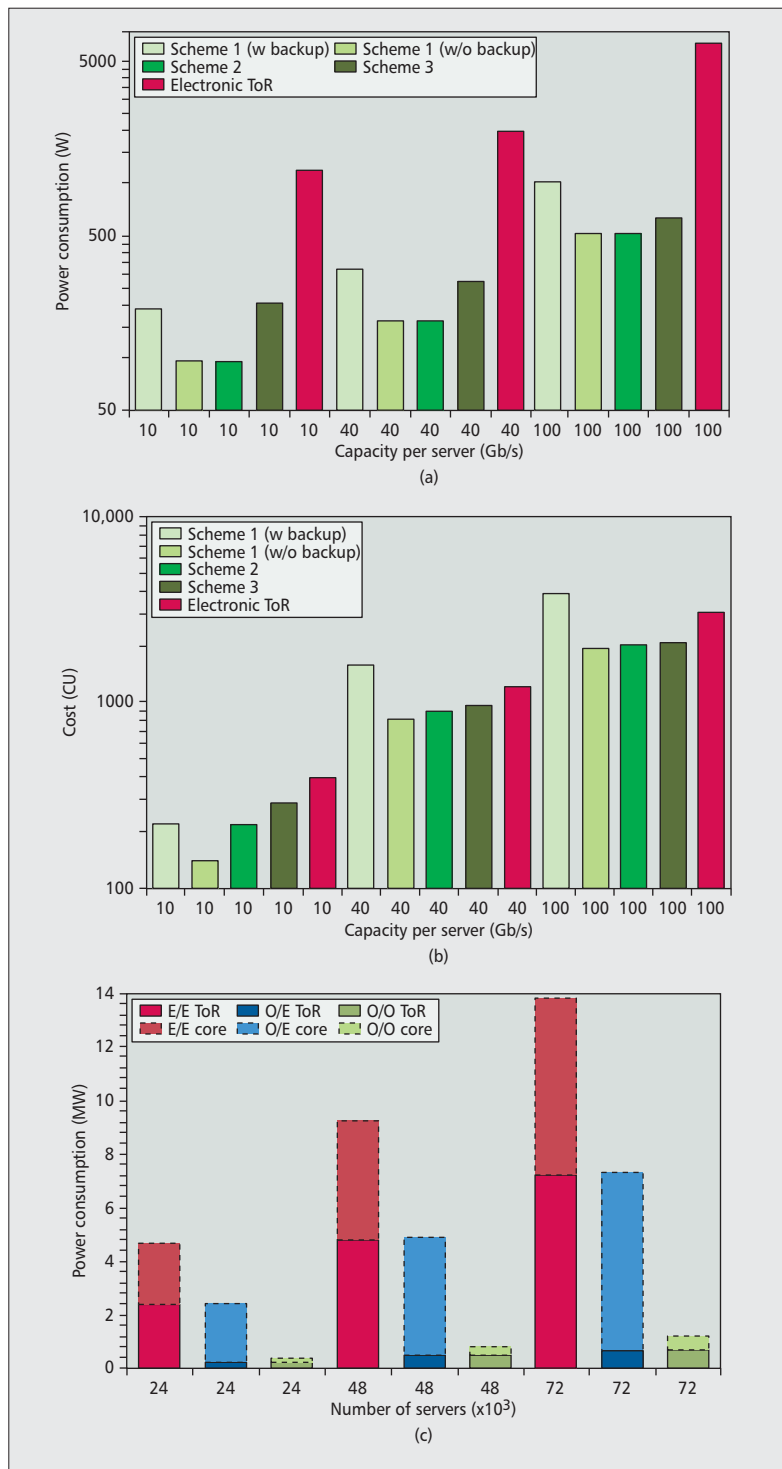


Figure 4. a) Power consumption; b) cost of the three presented schemes and a conventional commodity switch at ToR; c) total power consumed by the overall data center considering three different configurations.

[10] C. Develder et al., “Optical Networks for Grid and Cloud Computing Applications,” *Proc. IEEE*, vol. 100, no. 5, May 2012, pp. 1149–67.
 [11] D. Breuer et al., “Opportunities for Next-Generation Optical Access,” *IEEE Commun. Mag.*, vol. 49, no. 2, Feb. 2011, pp. S16–S24.
 [12] Y. Gong et al., “Highly Reliable Architecture Based on AWG for ToR in Datacenters,” filed patent, app. no. 201410207478.1, May 2014.
 [13] Y. Gong et al., “Passive Optical Interconnects at Top of the Rack for Data Center Networks,” *IEEE Optical Network Design and Modeling*, 2014

-
- [14] M. Fiorani *et al.*, "Energy-Efficient Elastic Optical Interconnect Architecture for Data Centers," *IEEE Commun. Letters*, vol. 18, no. 9, Sept. 2014, pp.1531–34.
- [15] L. Popa *et al.*, "A Cost Comparison of Datacenter Network Architectures," *ACM CoNext*, 2010.

BIOGRAPHIES

JIAJIA CHEN (jjajiac@kth.se) received a B.S. degree (2004) from Zhejiang University, China, and a Ph.D. degree (2009) from KTH Royal Institute of Technology, Sweden. She is working as an associate professor in the Optical Networks Lab (ONLab) at KTH. She is a co-author of over 100 publications in international journals and conferences in the area of optical networking. Her main research interests are optical transport and interconnect technology supporting future 5G and cloud environments. She has been involved in various European research projects like European FP7 projects IP-OASE and IP-DISCUS, and EIT-ICT projects. Moreover, she is principal investigator/co-principal investigator of several national research projects funded by the Swedish Foundation of Strategic Research (SSF) and the Swedish Research Council (VR).

YU GONG (yugong@coer.zju.edu.cn) received his Bachelor's degree (2013) from the Optical Engineering Department, Zhejiang University, China. He then joined the Centre for

Optical and Electromagnetic Research at Zhejiang University, where he is studying as a Ph.D. student. His research interests are passive optical network and optical interconnects for data centers.

MATTEO FIORANI (fiorani@kth.se) received a Master of Science in telecommunications (2010) and a Ph.D. in information and communication technologies (2014) from the University of Modena and Reggio Emilia, Italy. For one year he was a visiting student at Vienna University of Technology, Austria. He joined KTH Royal Institute of Technology in February 2014, where he is currently a postdoctoral researcher. He has co-authored 20 papers published in international journals and conference proceedings. His research interests include: optical interconnects for data centers and 5G transport technologies.

SLAVISA ALEKSIC [SM] (slavisa.aleksic@tuwien.ac.at) received M.Sc. (1999) and Ph.D. (2004) degrees in electrical engineering from Vienna University of Technology. His current research interests include communication networks, photonic networks, energy efficiency, optical and electrical signal processing, as well as high-speed media access control protocol design and implementation. He is an author or co-author of more than 110 scientific publications. He is a member of OVE, MRS, and IEICE. He has received several international awards, grants, and recognitions.

The Watchful Sleep Mode: A New Standard for Energy Efficiency in Future Access Networks

Raisa O. C. Hirafuji, Kelvin B. da Cunha, Divanilson R. Campelo, Ahmad R. Dhaini, and Denis A. Khotimsky

ABSTRACT

The continuously increasing consumption of power to access the Internet has been a major concern for network operators and equipment vendors. Passive optical network (PON) systems are widely seen as the future of broadband access. In 2010, ITU-T standardized a protocol-based PON energy efficiency mechanism that is comprised of two main modes, the doze mode and the cyclic sleep mode, which promise to save significant amounts of energy. However, the use of these two standardized alternative modes requires extra signaling and wastes energy. In this article we present the watchful sleep mode, a new mode that unifies the doze and cyclic sleep modes into a single power management mode. The new mode eliminates the extra control signaling and maximizes the amount of energy saved by keeping only the necessary hardware ON. Recently, the watchful sleep mode has been included in the ITU-T G.984 (G-PON) and ITU-T G.987 (XG-PON) recommendations. It is expected to be operated as the only power management mode in future PON systems.

INTRODUCTION

In the last decade energy efficiency has become a predominant theme in the design and operation of next-generation homes, electronics, software, and machinery. Recent studies show that by the year 2035, at the current power production rate, only 50 percent of the global power need will be satisfied [1]. Furthermore, to ensure human sustainability, global green house gas (GHG) emissions must be reduced by large quantities. This also applies to information and communication technologies (ICT), which are accountable for about 2 percent to 4 percent of the worldwide GHG emissions. More significantly, the power consumption of all ICT equipment, when active, accounts for roughly 40 percent to 60 percent of ICT's GHG emissions, and it is estimated that by the year 2020 these emissions will double if no action is taken [2].

Constituting the "bread-and-butter" of ICT,

the Internet is estimated to account for 1 percent to 2 percent of total electricity consumption in broadband-enabled countries; this is expected to grow even further with the data transmission rate trending higher [3]. More importantly, the Internet currently consumes about 10,000 times more energy than the minimum required to operate [1]. Meanwhile, the ever growing increase in Internet traffic (about 20 percent per year) and its impact on network energy consumption have become a growing concern for network operators and equipment vendors. The high electricity costs put pressure on the profitability of Internet service providers and operators, and the increased power consumption is causing heat dissipation problems in network gear. The scientific community reacted to the challenge by analyzing and experimenting with new techniques and tools for reducing the global Internet energy consumption [4].

Accounting for a large portion of modern Internet infrastructure, access networks in general, and optical access networks in particular, such as passive optical networks (which are gradually becoming the de-facto technology for next-generation broadband access networks) have been extensively studied, and many techniques have been proposed aiming at standardizing the energy efficiency mechanism in PON [5].

PON consists of a point-to-multipoint fiber infrastructure connecting a number of optical network units (ONUs) located at the customers' premises to an optical line terminal (OLT) residing at the service provider's central office using a shared fiber. Interested readers can refer to the works in [1] and [6] for a detailed overview of PON. PON systems have matured significantly in the last decade due to the proliferation of several delay-sensitive or bandwidth-intensive applications such as voice over Internet protocol (VoIP), high-definition video (HDTV), and video conferencing, all of which have stringent quality-of-service (QoS) requirements. Gigabit Ethernet PON (EPON) and its counterpart Gigabit PON (GPON) have been accepted as the standard for PON deployments. However, those killer bandwidth-hungry applications have

Raisa O. C. Hirafuji is with Universidade de Brasília (UnB).

Kelvin B. da Cunha and Divanilson R. Campelo are with Universidade Federal de Pernambuco (UFPE).

Ahmad R. Dhaini is with American University of Beirut.

Denis A. Khotimsky is with Verizon Communications, Inc.

made advanced versions of PON systems, such as 10 gigabit-capable PON (XG-PON), 10 gigabit Ethernet PON (10G-EPON), wavelength division multiplexing PON (WDM-PON), and time-wavelength division multiplexing PON (TWDM-PON) to become requisite technologies for the deployment of next-generation PON (NG-PON) [6].

Although it was shown in [3] that PON systems consume the least power among all the reported access network technologies, its evolution is expected to significantly increase its power consumption [6]. Therefore, reducing the energy consumption in PON is deemed necessary. To date, putting the ONU into low-energy mode has been considered the most cost-effective and promising method for saving energy in PON [1, 7–10]. Multiple proposals for a low-power mode have been put forward; these have been known as the *sleep* mode, which can either be *fast/cyclic* or *deep*, and the *doze* mode. The low-power modes are expected to decrease ONU energy consumption by almost 80 percent [8]. Interestingly though, the standard does not present a justification for proposing these different independently operated modes. More importantly, in order to achieve maximum energy saving, intelligent, and thus complex, arbitration of the different modes must be implemented at both the ONU and OLT.

In this article we present the *watchful sleep* mode, a new mode that unifies the doze and cyclic sleep modes into a single power management mode [11]. In the new mode, an ONU periodically turns off its receiver and transmitter, as in the cyclic sleep mode, and performs infrequent bidirectional handshakes, as in the doze mode. The watchful sleep mode not only simplifies the implementation of the power management scheme at the ONU and OLT, but also combines the advantages of the cyclic sleep and doze modes, and outperforms both of them. More interestingly, a PON system supporting the watchful sleep mode can actually emulate the cyclic sleep or doze mode as a special case.

Due to its effectiveness, the watchful sleep mode has been approved to be included in the ITU-T G.984 (G-PON) and ITU-T G.987 (XG-PONs) standards. It is also being considered for the NG-PON2 standard (ITU-T G.989), which aims at standardizing TWDM-PON networks.

The rest of the article is organized as follows. We present an overview of the main standardized (or considered for standardization) power saving techniques in PON systems. We describe the watchful sleep mode. We also discuss its practical implementation details, and present a comparative study to highlight its advantages. We discuss the future outlook and how the new mode can be operated under NG-PON2 systems. We summarize the article and present our conclusions.

PON ENERGY SAVING TECHNIQUES

In general, energy efficiency in PON systems can be achieved through improvements in components and module designs, one-sided power management techniques (i.e. either at the OLT or ONU), and protocol-based power manage-

ment mechanisms. As an example of improvements in component and module design, the use of vertical cavity surface emitting lasers (VCSELs)-based ONUs, as proposed in [12], can minimize the power consumed in both the active and power saving modes, and it can maximize the time the ONU spends in the power saving phase.

One-sided power management techniques are specific to either the OLT or ONU, thus they do not involve signaling between them. An example of an ONU-sided technique is *power shedding*, which is supported autonomously by the ONU [5]. With power shedding, the ONU turns off some of its components to reduce power consumption, while keeping the optical link in full function. An OLT-side power management technique has been proposed in [13], in which working ONUs of a TWDM-PON are organized to send and receive traffic over a minimum number of wavelengths, allowing the OLT to turn off some of its transceivers so as to save power.

Protocol-based power management techniques do require signaling between the OLT and ONU. Well known examples of such techniques are the doze and cyclic sleep modes, which have been standardized by the ITU-T [14]. In both modes, the ONU alternates between active and power saving phases, which are initiated and terminated through the exchange of signaling information between the OLT and ONU. During a power saving phase, the ONU swaps between two states: a full power state, referred to as the *aware state* for both modes; and a low power state. In doze mode, the low power state is known as the *listen state*, in which the ONU transmitter is OFF and the ONU receiver is ON. In the cyclic sleep mode, the low power state is known as the *sleep state*; here, both the ONU's transmitter and receiver are OFF.

To enable the doze and cyclic sleep modes in PON, the ONU and OLT implement the state machines illustrated in Fig. 1. Due to the centralized nature of PON systems, the OLT maintains a different state machine for each ONU. As shown in Fig. 1a, in the ONU state machine, the *active held* and *active free* states constitute the active phase of an ONU, while in the active held state, the ONU cannot enter a power saving phase. Meanwhile, in the active free state, the ONU can freely enter a power saving phase as soon as a local doze indication (LDI) occurs, which initiates the doze mode operation. Similarly, a local sleep indication (LSI) will initiate the cyclic sleep mode operation. The *doze aware*, *sleep aware*, *listen*, and *asleep* states comprise the power saving phase of the ONU. The transition from a power saving phase to an active phase is triggered by:

- A local wake-up indication (LWI).
- A message from the OLT.
- A forced wake up indication (FWI) bit in the bandwidth allocation map from the OLT.

The implementation of the LWI, LDI, and LSI is left to the ONU vendor and/or PON operator.

In the OLT state machine, the *awake forced* and *awake free* states correspond to the ONU in the active phase, whereas the *low power doze/sleep* and *alerted doze/sleep* states corre-

Due to its effectiveness, the watchful sleep mode has been approved to be included in the ITU-T G.984 (G-PON) and ITU-T G.987 (XG-PONs) standards. It is also being considered for the NG-PON2 standard (ITU-T G.989), which aims at standardizing TWDM-PON networks.

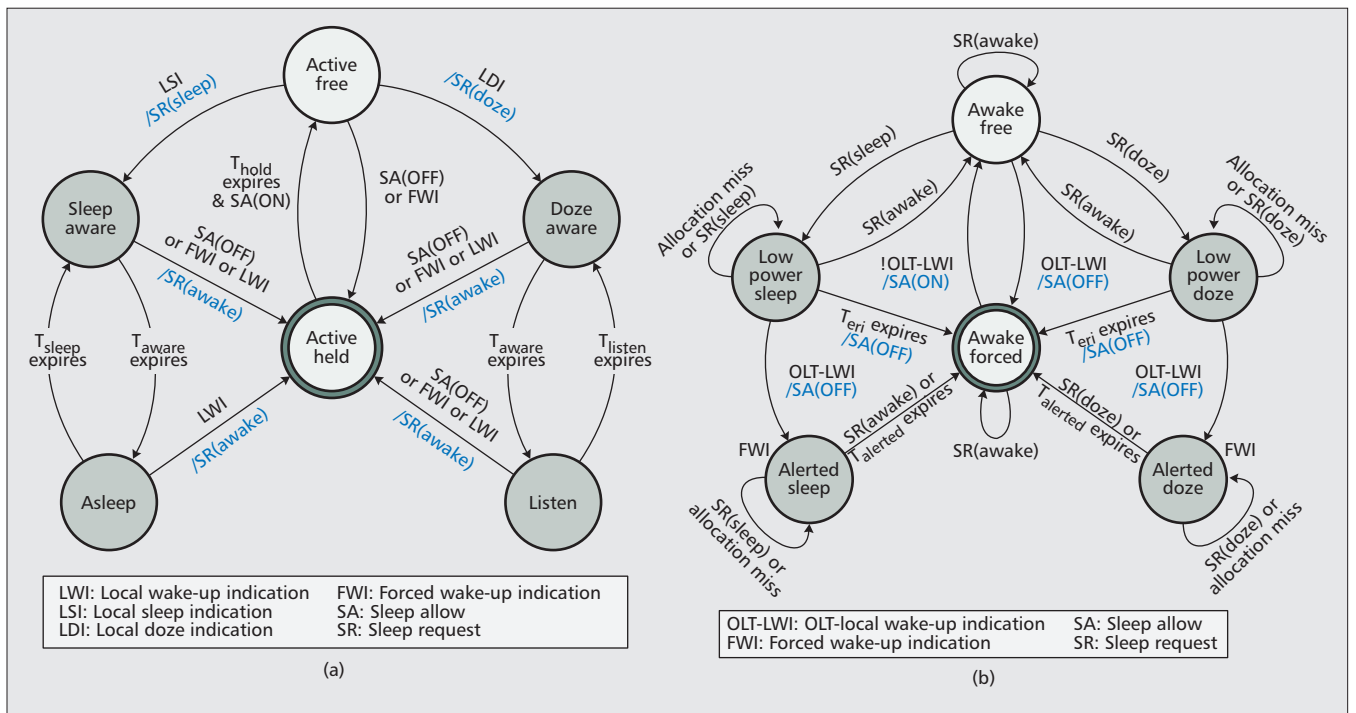


Figure 1. Cyclic sleep and doze modes state machines: a) ONU state machine; b) OLT state machine. The notation MESSAGE(ARG) represents the command and its corresponding argument as defined in the ITU-T standard. For example, SA(ON) is an OLT message that “allows” the ONU to start a power saving phase.

spond to the ONU in the power saving phase. The awake forced state indicates that the ONU cannot enter a power saving phase. In the awake free state, the ONU can enter a power saving mode, which is either the low power doze or low power sleep state. Finally, the *alerted doze* and *alerted sleep* states correspond to the period in which the OLT attempts to wake up the ONU. The transition from the awake forced state to the awake free state, which grants permission for the ONU to start a power saving phase, occurs upon the cessation of the LWI at the OLT, denoted by !OLT-LWI. The OLT stimulus OLT-LWI always triggers a transition to the awake forced or alerted doze/sleep states; this revokes the ONU permission to be in a power saving phase.

The sojourn time of the ONU in the aware, listen, and sleep states is represented by the notations T_{aware} , T_{listen} , and T_{sleep} , respectively, and they are set by the OLT according to the following criteria:

- To configure T_{aware} , the OLT should keep track of the time required to perform at least one handshake between the OLT and ONU; meanwhile the ONU would be in the active state. Therefore, the period T_{aware} is configured to be bigger than the duration of one handshake.

- To configure T_{listen} , the OLT should know how often the periodic bidirectional handshakes between the OLT and ONU occur. These may include, for example, an adjustment margin to account for keep-alive maintenance and ranging periods. Therefore, the period T_{listen} would be configured to be smaller than the maximum interval between these periodic handshakes. The ONU remains in the listen state for the period T_{listen} unless there is a local or an external wake-

up stimulus, which forces it to transition to the aware state.

- To configure T_{sleep} , the OLT must primarily take into consideration the tolerable latency of the external wake-up stimulus, that is, the elapsed period between the time to wake up the ONU and the time needed to restore bidirectional communication (i.e., between the OLT and ONU). Therefore, T_{sleep} shall be smaller than the tolerable latency of the external wake-up stimulus. The ONU remains in the sleep state for the period T_{sleep} if it is not interrupted by the arrival of a local wake-up stimulus.

The minimum ONU sojourn time in the active held state is denoted as T_{hold} (Fig. 1a). In the OLT state machine, the timing parameters T_{eri} and T_{alerted} denote the timeout for the low power doze/sleep state, and the timeout for the alerted doze/sleep state, respectively.

Even though the doze and cyclic sleep modes can improve the energy efficiency in PON, they impose several limitations. First, the OLT and ONU need to negotiate and agree upon the operated power saving mode (i.e., either the doze mode or the cyclic sleep mode). Second, in the Doze mode, the ONU keeps the receiver always ON even when there is no traffic destined for that ONU. Finally, in the cyclic sleep mode, in order to probe for external wake-up stimulus, the ONU must periodically turn ON both its receiver and transmitter, even though turning the receiver ON only is sufficient to achieve the foregoing goal.

As noted, the foregoing limitations may either downgrade the system performance (e.g., due to extra sophisticated signaling overhead), or waste energy due to keeping some of the devices unnecessarily ON. Furthermore, the ITU-T stan-

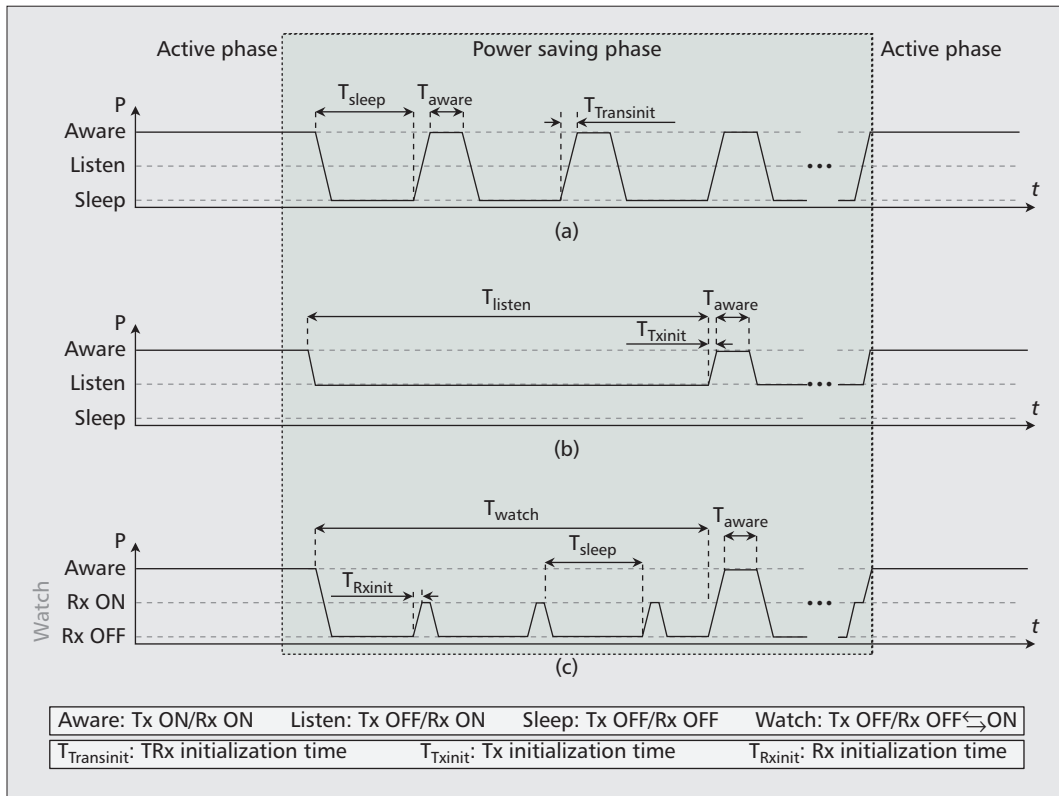


Figure 2. Power consumption in a) cyclic sleep mode; b) doze mode; c) watchful sleep mode.

Like the doze and cyclic sleep modes, the watchful sleep mode alternates between the active and power saving phases. In addition, during the power saving phase, it switches between the full power state and the low power state. The difference between the watchful sleep mode and the other two modes lies in the semantics of the low power state.

ard does not present a justification to maintain a separation of the two aforementioned power saving modes. Thus, the complexity of implementing and operating these two modes is not justified.

THE WATCHFUL SLEEP MODE

The *Watchful Sleep* mode is a new mode that unifies the doze and cyclic sleep modes into a single power saving mode. The new mode mitigates the drawbacks imposed by the operation of the doze and cyclic sleep modes by eliminating the need for negotiation of the employed mode between the OLT and ONU. It also maximizes the amount of energy saved in the ONU by keeping only the required devices ON, while the rest will be OFF.

Like the doze and cyclic sleep modes, the watchful sleep mode alternates between the active and power saving phases. In addition, during the power saving phase, it switches between the full power state and the low power state. The difference between the watchful sleep mode and the other two modes lies in the semantics of the low power state, as depicted in Fig. 2.

In the low power state of the watchful sleep mode, referred to as the *watch state*, the ONU maintains its transmitter OFF, but periodically turns the receiver ON for a short time to check for external wake-up stimuli. As a result, during the watch state, the ONU may be in two possible power consumption levels:

- The power level in which both the transmitter and the receiver are OFF.
- The power level in which the transmitter is OFF and the receiver is kept ON.

The ONU remains in the watch state for a period denoted T_{watch} , unless there is a local or an external wake-up stimulus that forces the ONU to transition to the aware state. The ONU sojourns with the receiver OFF for the period T_{sleep} before turning the receiver back ON to check for external wake-up stimuli; the receiver is kept ON until the ONU confirms that there is no external wake-up stimulus. This interval can be as small as one frame (e.g., 125 μ s in GPON) if the OLT sends a scheduling grant to the ONU in each frame. The OLT configures T_{aware} and T_{sleep} according to the same criteria used for the choice of these parameters in the cyclic sleep mode, whereas T_{watch} is configured according to the same criteria used to set T_{listen} in the doze mode.

Figure 3 illustrates the OLT and ONU state machines for employing the watchful sleep mode in PON. In comparison to the doze/cyclic sleep-based state machines shown in Fig. 1, the state machines for the watchful sleep mode have fewer states and signaling messages. Specifically, in the ONU state machine, the doze/sleep aware states have been merged into the *single* *wsleep* aware state, and the listen and asleep states have been merged into the single *watch* state. Similarly, in the OLT state machine, the low power doze/sleep states have been merged into the single *low power watch state*, and the alerted doze/sleep states have been merged into the single *alerted watch state*. Furthermore, the LSI and LWI stimuli have been merged into the local low power indication (LPI). Clearly, the new “simpler” state machines make the implementation and operation of the watchful sleep mode in

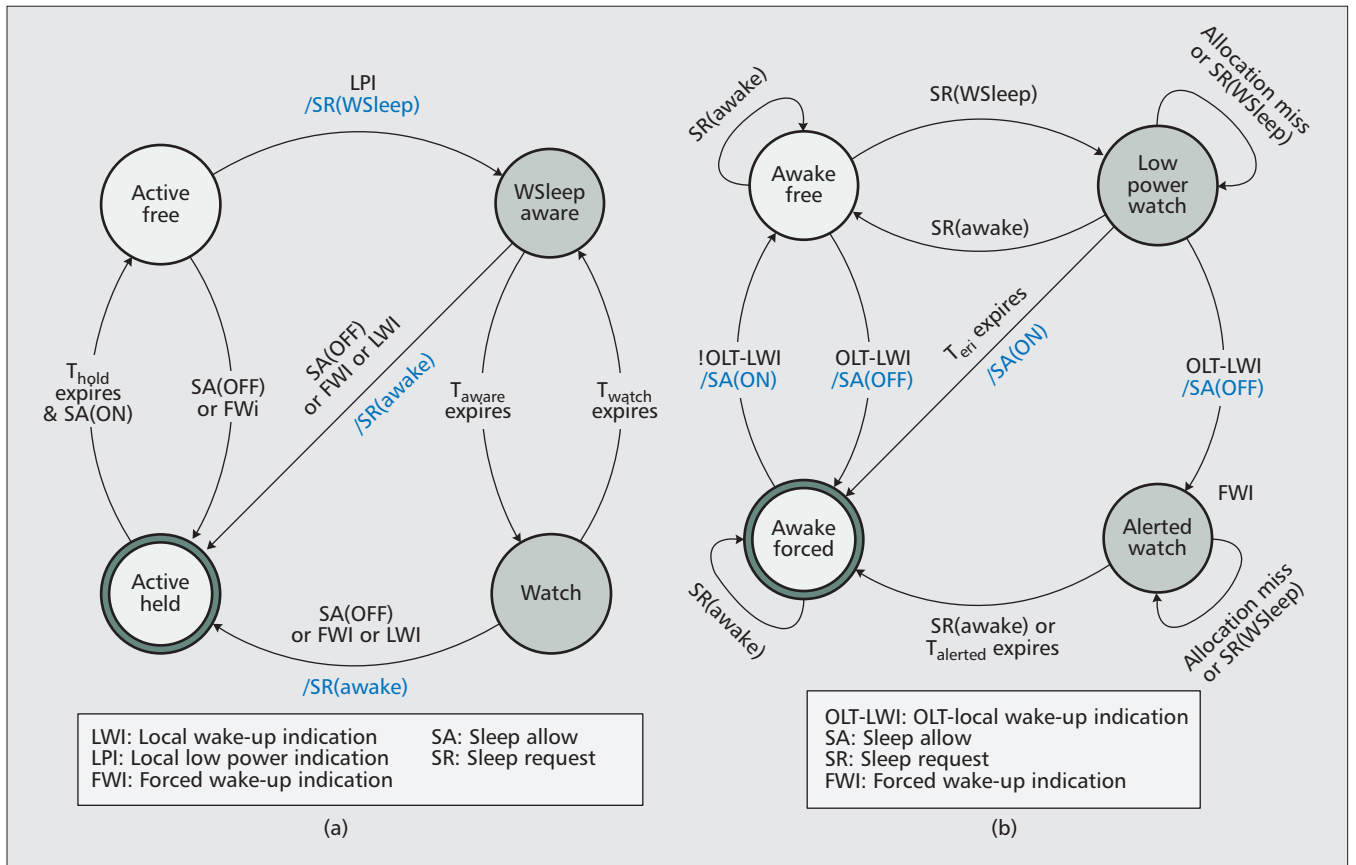


Figure 3. Watchful sleep mode state machines: a) ONU state machine; b) OLT state machine. The notation MESSAGE(ARG) represents the command and its corresponding argument as defined in the ITU-T standard. To adhere with the ITU-T standard, we define and implement similar commands. For example, SR(WSleep) is an ONU message sent to the OLT, expressing intent to start a power saving phase in the watchful sleep mode.

Parameter	Watchful sleep mode	Doze mode	Cyclic sleep mode
T_{aware}	5 ms	5 ms	5 ms
T_{sleep}	10 ms	—	10 ms
T_{listen}	—	10 s	—
T_{watch}	10 s	—	—
P_{Active}^1	100%	100%	100%
P_{Listen}^2	40%	40%	—
P_{Sleep}^3	5%	—	5%
T_{Txinit}^4	3 ms	3 ms	—
T_{Rxinit}^4	2 ms	—	—
$T_{\text{Transinit}}^4$	3 ms	—	3 ms

¹ P_{Active} : the power consumption when both the receiver and transmitter are ON.

² P_{Listen} : the power consumption when the receiver is ON and the transmitter is OFF.

³ P_{Sleep} : the power consumption when both the receiver and transmitter are OFF.

⁴ Refer to Fig. 2.

Table 1. Simulations parameters.

PON systems easier and more efficient. Furthermore, with the watchful sleep mode, the OLT and the ONU do not need to negotiate the power saving mode, since there is only one power management mode. More interestingly, the watchful sleep mode can emulate the doze and cyclic sleep modes in a very simple manner. To emulate the doze mode, the parameters will be simply configured as follows: $T_{\text{watch}} = T_{\text{listen}}$ and $T_{\text{sleep}} = 0$. To emulate the cyclic sleep mode, the OLT simply configures: $T_{\text{watch}} = T_{\text{sleep}}$.

COMPARATIVE SIMULATION STUDY

To highlight the advantages of the watchful sleep mode over the doze and cyclic sleep modes in terms of energy saving, we have conducted extensive simulations using OMNeT++. We consider an XG-PON system with 16 ONUs. The distance between each ONU and the OLT is set as 20 km, and a self-similar traffic generator is used to inject traffic in the network. The packet size is uniformly distributed between 64 and 1518 bytes. Table 1 shows the set of parameters used in the simulations for the doze, cyclic sleep, and watchful sleep modes. For the sake of generality, the power consumption of the XG-PON transceivers are represented by percentages in lieu of absolute values.

Most of the results in the literature pertaining to energy saving in PON consider the implementation of the wake up stimuli based on the

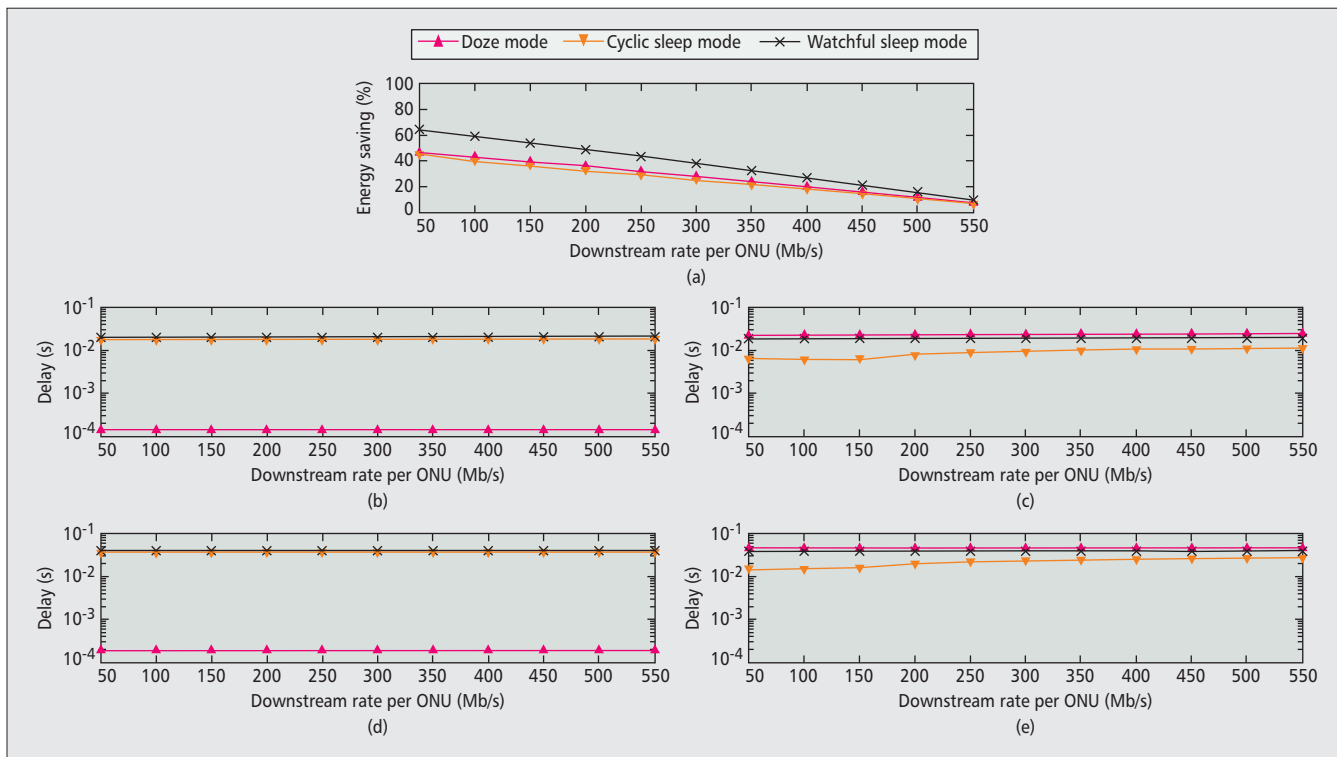


Figure 4. Simulation results with the watchful sleep mode under XG-PON: a) energy saving; b) mean downstream delay; c) mean upstream delay; d) maximum downstream delay; e) maximum upstream delay.

arrival of a user data frame [7, 11]. Consequently, the reported results do not exhibit any energy saving for high traffic loads, since a single frame arrival terminates the power saving phase. Conversely, in this work the wake up indicators are configured to enable energy saving at the ONU even in the presence of intense traffic, while maintaining a maximum end-to-end delay requirement of 56 ms (this value is considered to be acceptable for delay-sensitive applications such as voice over IP [9]).

The LWI indicator is configured to be activated 50 ms after a frame arrival, whereas the OLT-LWI indicator is configured to be activated 40 ms after a frame arrival. This difference of 10 ms between both indicators is due to the value of $T_{\text{sleep}} = 10$ ms in Table 1. Since, in the worst case scenario, the OLT might decide to wake the ONU as soon as the ONU switches to the Rx-OFF level in the aware state, a difference of $T_{\text{sleep}} (10$ ms) is necessary to maintain the maximum delay close to 50 ms. The indicators LDI, LSI, and LPI are configured to be activated once the ONU empties the queue. The !OLT-LWI stimulus is used by the OLT to trigger the ONU to enter the power saving mode when its queue is empty.

Figure 4 compares the performances of the doze, cyclic sleep, and watchful sleep modes in terms of the following: energy efficiency; mean downstream delay; mean upstream delay; maximum downstream delay; and maximum upstream delay. All the results are presented as a function of the packet arrival rate (measured in Mb/s) in the downstream direction. The packet arrival rate for the upstream direction is set to be always 1/4 of the downstream rate. The maximum delay

curves look constant due to the 40 ms (at the OLT) and 50 ms (at the ONU) extra waiting time in the queue. Besides, even without the power management modes, these curves would look constant due to the fact that these rates are not even close to the maximum capacity of XG-PON systems. It is worth noting that the watchful sleep mode not only outperforms the other two modes in terms of energy efficiency, but it also achieves about 50 percent energy saving for a downstream rate of 200 Mb/s, as shown in Fig. 4a. Furthermore, under the new mode, the maximum delay in the downstream and upstream directions is 55 ms. This demonstrates the ability of the watchful sleep mode to save more energy than the doze and cyclic sleep modes, without violating the delay requirements.

FUTURE OUTLOOK

THE WATCHFUL SLEEP MODE FOR NG-PON2

Currently, ITU-T is working toward the standardization of a 40-Gigabit-capable passive optical network (NG-PON2), which will be based on the TWDM-PON architecture [15]. Due to its effectiveness, the watchful sleep mode is also being considered to be operated as the default energy management mode in NG-PON2.

To anticipate the performance of the watchful sleep mode under NG-PON2, which may contribute to the foregoing decision of adopting the watchful sleep mode by ITU-T, we conducted a simulation study. Here we consider a TWDM PON with four wavelengths and 64 ONUs, with each wavelength serving 16 ONUs. Each wavelength pair has a downstream rate of 10 Gbps and an upstream rate of 2.5 Gbps. The

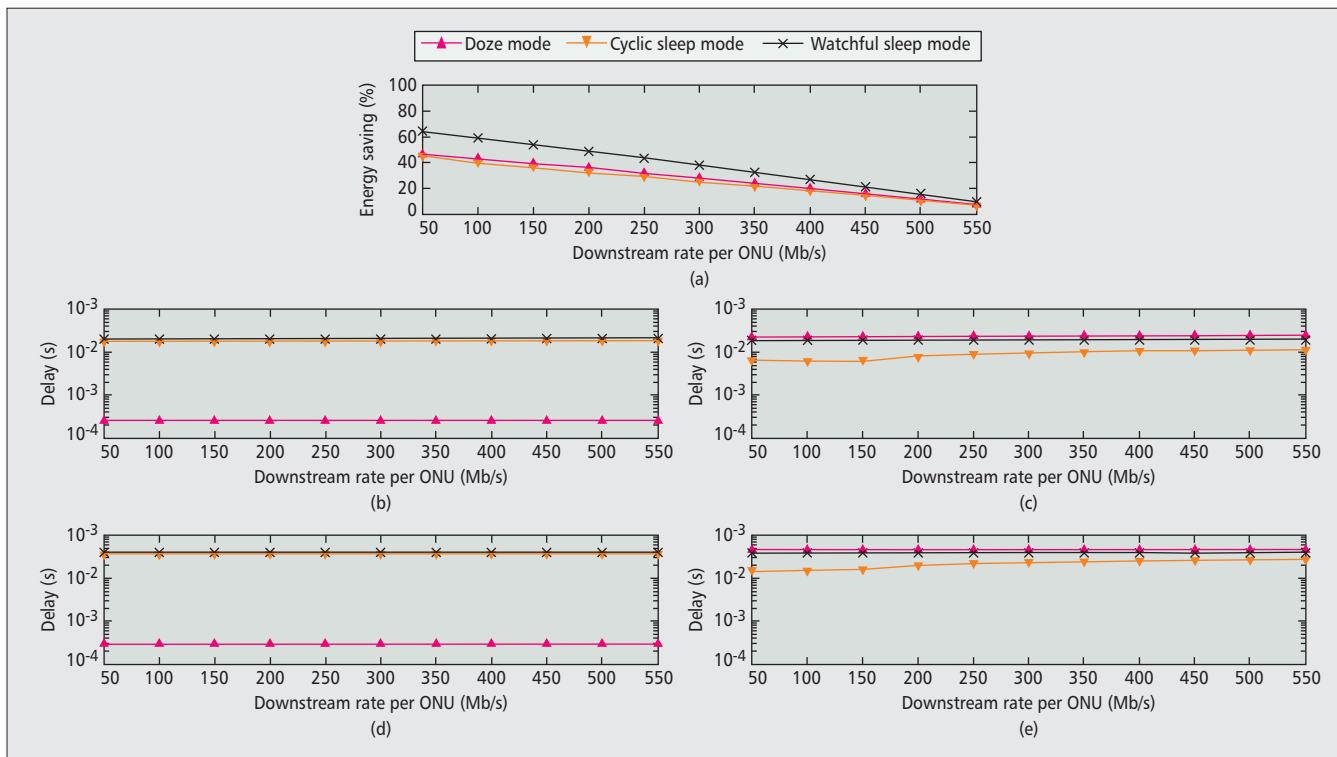


Figure 5. Simulation results with the watchful sleep mode under TWDM-PON: a) energy saving; b) mean downstream delay; c) mean upstream delay; d) maximum downstream delay; e) maximum upstream delay.

distance between the ONUs and the OLT is assumed to be 40 km, resulting in a propagation delay of 200 μ s. The same parameter set used for the XG-PON scenario is used here (Table 1).

Figure 5 compares the performance of the watchful sleep mode to that of the doze and cyclic sleep modes under the simulated TWDM-PON network. As noted, with the unified mode, the energy saved at lower rates is significantly larger than that at higher rates, due to the ONU spending more time in the power saving phase at lower rates. The results in Fig. 5 are quite similar to those presented in Fig. 4 since the TWDM-PON system is formed by four stacked XG-PON systems. In terms of delay, it can be seen in the figure that the delay requirement for delay-sensitive applications is not compromised. These results show that the watchful sleep mode can indeed be employed for NG-PON2, as it can effectively save more energy than the other two modes without impairing the services delivered to end users.

CONCLUSIONS

Energy efficiency is central to the design and operation of next-generation broadband access networks such as passive optical networks. Several effective techniques have been proposed to reduce the energy consumption of the network elements in PON systems, especially the customer premises equipment, without impairing the services delivered to end users.

In this article we present a novel power saving mode, namely the watchful sleep mode, which unifies two standardized alternative PON power saving modes (i.e., the doze and cyclic

sleep modes) into a single power management mode. The watchful sleep mode combines the advantages of the doze and cyclic sleep modes and outperforms them in terms of energy efficiency and simplicity of implementation. Due to its efficacy, the watchful sleep mode has been included in the ITU-T G.984 (G-PON) and G.987 (XG-PON) Recommendation suites, and is currently being considered for adoption in the ITU-T G.989 Recommendation series, which standardizes TWDM-PON systems.

ACKNOWLEDGMENTS

This work was supported in part by CNPq, CAPES, and FACEPE.

REFERENCES

- [1] A. R. Dhaini, P.-H. Ho, and G. Shen, "Toward Green Next-Generation Passive Optical Networks," *IEEE Commun. Mag.*, vol. 29, no. 11, Nov. 2011, pp. 94–101.
- [2] W. Vereecken *et al.*, "Power Consumption in Telecommunication Networks: Overview and Reduction Strategies," *IEEE Commun. Mag.*, vol. 29, no. 6, June 2011, pp. 62–69.
- [3] K. Hinton *et al.*, "Power Consumption and Energy Efficiency in the Internet," *IEEE Network*, vol. 25, no. 2, Mar./Apr. 2011, pp. 6–12.
- [4] "GreenTouch," available: <http://www.greentouch.org>
- [5] ITU-T, "GPON Power Conservation," available: <http://www.itu.int/rec/T-REC-G.SUP45-200905-I/en>.
- [6] J.-I. Kani, F. Bourgart, A. Cui, A. Rafel, and S. Rodrigues, "Next-Generation PON — Part 1: Technology Roadmap and General Requirements," *IEEE Commun. Mag.*, vol. 47, no. 11, Nov. 2009, pp. 43–49.
- [7] B. Skubic and D. Hood, "Evaluation of ONU Power Saving Modes for Gigabit-Capable Passive Optical Networks," *IEEE Network*, vol. 25, no. 2, Mar./Apr. 2011, pp. 20–24.
- [8] S. Wong *et al.*, "Sleep Mode for Energy Saving PONs: Advantages and Drawbacks," *Proc. IEEE GLOBECOM'09*, Hawaii, USA, Dec. 2009.

- [9] L. Valcarenghi *et al.*, "Energy Efficiency in Passive Optical Networks: Where, When, and How?," *IEEE Network*, vol. 26, no. 6, Nov./Dec. 2012, pp. 61–68.
- [10] A. R. Dhaini *et al.*, "Energy Efficiency in TDMA-based Next-Generation Passive Optical Access Networks," *IEEE/ACM Trans. Net.*, vol. 22, no. 3, May 2014, pp. 8504–63.
- [11] D. A. Khotimsky *et al.*, "Unifying Sleep and Doze Modes for Energy-Efficient PON Systems." *IEEE Commun. Letters*, vol. 18, no. 4, Apr. 2014, pp. 688–91.
- [12] E. Wong *et al.*, "Energy Efficiency of Optical Network Units with Vertical-Cavity Surface-Emitting Lasers," *Optics Express*, vol. 20, no. 14, 2012, pp. 14960–70.
- [13] H. Yang *et al.*, "ONU Migration in Dynamic Time and Wavelength Division Multiplexed Passive Optical Network (TWDM-PON)," *Optics Express*, vol. 21, no. 18, 2013, pp. 21491–99.
- [14] ITU-T, "10-Gigabit-Capable Passive Optical Networks (XG-PON): Transmission Convergence (TC) Layer Specification," ITU-T Recommendation G.987.3, Jan. 2014.
- [15] Y. Luo *et al.*, "Time- and Wavelength-Division Multiplexed Passive Optical Network (TWDM-PON) for Next-Generation PON Stage 2 (NG-PON2)," *IEEE/OSA J. Lightwave Technology*, vol. 31, no. 4, Feb. 2013, pp. 587–93.

BIOGRAPHIES

RAISA O. C. HIRAFUJI (raisahana@ieee.org) received her communication networks engineering degree from the Universidade de Brasília, Brazil, in 2012. She is currently working toward the M.Sc. degree in electrical engineering at the same university. Her current research interests are in next-generation passive optical networks and green communications.

KELVIN B. DA CUNHA (kbc@cin.ufpe.br) is currently an undergraduate student in computer engineering at Centro de Informática (CIn) of Universidade Federal de Pernambuco (UFPE), Recife, Brazil. His research interests include energy efficiency in next-generation optical access networks.

DIVANILSON R. CAMPELO (dcampelo@cin.ufpe.br) is an associate professor of computer engineering and computer science at Centro de Informática (CIn) of Universidade Federal de Pernambuco (UFPE), Recife, Brazil. He received his electrical engineer degree from UFPE in 1998 and his master's and Ph.D. degrees, also in electrical engineering, from the Universidade Estadual de Campinas (UNICAMP), Campinas, Brazil, in 2001 and 2006, respectively. In 2000–2001 he was a systems engineer at Nortel Networks, Sao Paulo, Brazil. In 2007–2008 he was an assistant professor at the Universidade Presbiteriana Mackenzie, Sao Paulo, Brazil. In 2008–2009 he was a visiting assistant professor at Stanford University, working in the Photonics and Networking

Research Laboratory (PNRL) led by Prof. Leonid Kazovsky. In 2010–2012 he was an assistant professor at the Universidade de Brasília, Brasília, Brazil. He is a member of the IEEE, IEEE ComSoc, ACM and ACM SIGCOMM. He has authored/coauthored a book chapter and more than 30 journals and conference papers. His current research interests include broadband access networks, energy efficiency in ICT, automotive networking, and software defined networking.

AHMAD R. DHAINI (ahmad.dhaini@aub.edu.lb) is a visiting assistant professor of computer science at the American University of Beirut (AUB), Lebanon. He received his B.Sc. in computer science from AUB in 2004; his M.Sc. degree in electrical and computer engineering from Concordia University, Montreal, Canada in 2006; and his Ph.D. degree in electrical and computer engineering from University of Waterloo, Canada in 2011. His master's dissertation was nominated for the best thesis award. In 2007–2008 he was a software analyst and consultant at TEKSystems, Montreal, Canada. In 2007–2008 he was a software designer at Ericsson, Montreal, Canada. In 2012–2014 he was as a postdoctoral scholar at Stanford University, working in the Photonics and Networking Research Laboratory (PNRL) led by Prof. Leonid Kazovsky, after being awarded the prestigious Natural Sciences and Engineering Research Council of Canada Postdoctoral Fellowship (NSERC PDF). He is an inventor with US patents. He has also authored/co-authored one book, one book chapter, and more than 30 highly cited research articles in top IEEE journals and conferences. He is a reviewer for NSF, NSERC, and several US universities' internal grants. He also serves as an editor for Springer's *Photonics Networks Communications*, and as a reviewer and technical program committee (TPC) member for several major IEEE journals and conferences. His research interests cover several themes of optical networks such as fiber-wireless (FiWi) broadband access networks, mission-critical networks, green communications, and software-defined networking.

DENIS A. KHOTIMSKY (dkhprim@gmail.com) received the Dipl.Eng. degree from the Moscow Aviation Institute in 1987 (electronic engineering) and the Ph.D. from the University of California in Santa Barbara in 1996 (computer science). He is currently a distinguished member of technical staff at Verizon Corporate Technologies in Waltham, Massachusetts, where he is focusing on PON-based access network planning. For the last eight years he has been involved in the standardization of passive optical network systems, working with FSAN, ITU-T SG15, and the Broadband Forum. He has served as an editor for several key ITU-T Recommendations related to G-PON, XG-PON1, and NG-PON2. In 2012 he proposed and subsequently defended the idea of watchful sleep mode in ITU-T.

Crowdsending Based Public Transport Information Service in Smart Cities

Károly Farkas, Gábor Fehér, András Benczúr, and Csaba Sidló

Editor's Note: The IEEE Communications Society has a Sister Society agreement with HTE (The Hungarian Association for Infocommunications). The terms of the agreement include republication of articles of HTE's *Infocommunications Journal* in IEEE Communications Society publications. The article below has already appeared in *Infocommunications Journal*. The citation is: *Infocommunications Journal*, no 4, pp. 13–20, Dec. 2014

Osman Gebizlioglu, Editor-in-Chief, *IEEE Communications Magazine*

ABSTRACT

Thanks to the development of technology and the emergence of intelligent services smart cities promise to their inhabitants enhanced perception of city life. For example, a live timetable service of public transportation can increase the efficiency of travel planning substantially. However, its implementation in a traditional way requires the deployment of some costly sensing and tracking infrastructure. Mobile crowdsensing is an alternative, when the crowd of passengers and their mobile devices are used to gather data for almost free of charge.

In this paper, we put the emphasis on the introduction of our crowdsensing based public transport information service, what we have been developing as a prototype smart city application. The front-end interface of this service is called TrafficInfo. It is a simple and easy-to-use Android application which visualizes real-time public transport information of the given city on Google Maps. The lively updates of transport schedule information relies on the automatic stop event detection of public transport vehicles. TrafficInfo is built upon our Extensible Messaging and Presence Protocol (XMPP) based communication framework what we designed to facilitate the development of crowd assisted smart city applications. The paper introduces shortly this framework, than describes TrafficInfo in detail together with the developed stop event detector.

INTRODUCTION

Services offered by smart cities aim to support the everyday life of inhabitants. Unfortunately, the traditional way of introducing a new service usually implies a huge investment to

deploy the necessary background infrastructure.

One of the most popular city services is public transportation. Maintaining and continuously improving such a service are imperative in modern cities. However, the implementation of even a simple feature which extends the basic service functions can be costly. For example, let's consider the replacement of static timetables with lively updated public transport information service. It requires the deployment of a vehicle tracking infrastructure consisting of among others GPS sensors, communication and back-end informatics systems and user interfaces, which can be an expensive investment.

An alternative approach to collect real-time tracking data is exploiting the power of the crowd via participatory sensing or often called mobile crowdsensing¹ [1], which does not call for such an investment. In this scenario (Fig. 1), the passengers' mobile devices and their built-in sensors, or the passengers themselves via reporting incidents, are used to generate the monitoring data for vehicle tracking and send instant route information to the service provider in real-time. The service provider then aggregates, cleans, analyzes the data gathered, and derives and disseminates the lively updates. The sensing task is carried out by the built-in and ubiquitous sensors of the smartphones either in participatory or opportunistic way depending on whether the user is involved or not in data collection. Every traveler can contribute to this data harvesting task. Thus, passengers waiting for a ride can report the line number with a timestamp of every arriving public transport vehicle at a stop during the waiting period. On the other hand, onboard passengers can be used to gather and report actual position information of the moving vehicle and detect halt events at the stops.

In this paper, we focus on the introduction of our crowdsensing based public transport information service, what we have been developing as a prototype smart city application. The front-end interface of this service, called TrafficInfo, is a simple and easy-to-use Android application which visualizes real-time public transport information of the given city on Google Maps. It is built upon our Extensible Messaging and Presence Protocol (XMPP) [2] based communication framework [3] what we designed to facilitate the development of crowd assisted smart city applications (we also introduce shortly this framework). Following the

K. Farkas and G. Fehér are with the Inter-University Centre for Telecommunications and Informatics, Debrecen, Hungary, and the Budapest University of Technology and Economics, Budapest, Hungary.

A. Benczúr and Cs. Sidló are with the University of Debrecen, Hungary and the Institute for Computer Science and Control, Hungarian Academy of Sciences (MTA SZTAKI), Budapest, Hungary

¹ We use the terms *crowdsensing*, *crowdsourcing* and *participatory sensing* interchangeably in this paper.

publish/subscribe (pub/sub) communication model the passengers subscribe in TrafficInfo, according to their interest, to traffic information channels dedicated to different public transport lines or stops. Hence, they are informed about the live public transport situation, such as the actual vehicle positions, deviation from the static timetable, crowdedness information, etc.

To motivate user participation in data collection we offer a day zero service to the passengers, which is a static public transportation timetable. It is built on the General Transit Feed Specification (GTFS, designed by Google) [4] based transit schedule data and provided by public transport operators. TrafficInfo basically presents this static timetable information to the users which is updated in real-time, if appropriate crowdsensed data is available. To this end, the application collects position data; the timestamped halt events of the public transport vehicles at the stops; and/or simple annotation data entered by the user, such as reports on crowdedness or damaged seat/window/lamp/etc. After analyzing the data gathered live updates are generated and TrafficInfo refreshes the static information with them.

The rest of the paper is structured as follows. After a quick overview of related work we introduce our generic framework to facilitate the development of crowdsourcing based services. We show our live public transport information service together with the developed stop event detector. Finally, we summarize our work with a short insight to our future plans.

RELATED WORK

In this section, we discuss the challenge of attracting users to participate in crowdsensing and review the relevant works in the field of crowd assisted transit tracking systems.

A crowdsourcing based service has to acquire the necessary information from its users who are producers and consumers at the same time. Therefore it is essential for the service provider to attract users. However, we face a vicious circle here. The users are joining the service if they can benefit from it and at the same time they contribute to keep running the service which can persuade others also to join. But how can the users be attracted if the service is not able to provide the expected service level due to the lack of contributors? This also means that the service cannot be widely spread without offering a minimum service level and until it has a sufficiently large user base.

Moovit² is a similar application to TrafficInfo which is meant to be a live transit app on the market providing real-time information about public transportation. It faces the above mentioned problem in many countries. Moovit has been successful only in those cities where it has already a mass of users, just like in Paris, and not successful in cities where its user base is low, e.g., in Budapest. In order to create a sufficiently large user base Moovit provides, besides live data, schedule based public transportation information as a day zero service, too. The source of this information is the company who operates the public transportation network. The best practice is for providing such information is using GTFS [4]. According to the GTFS developer page, cur-

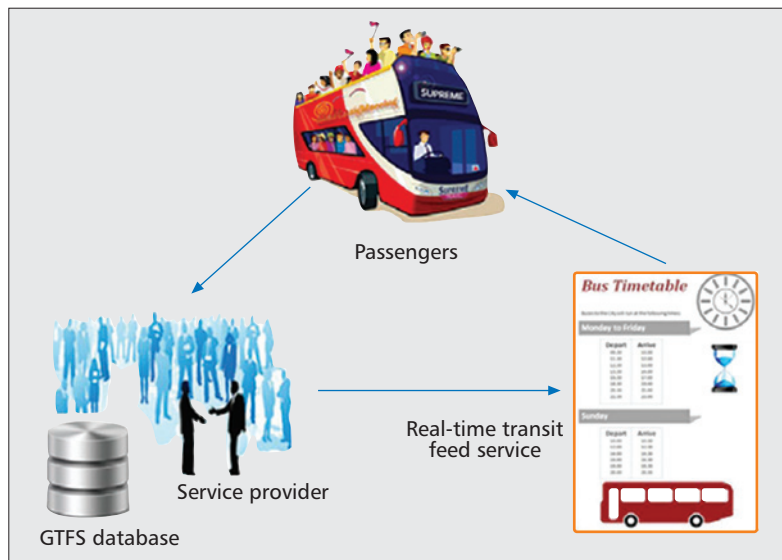


Figure 1. Real-time public transport information service based on mobile crowdsensing.

rently 263 public transportation companies provide transit feeds from all over the world. Moovit partially relies on GTFS and is available in 350 cities attracting more than 6.5 million users. We also adopted this solution in TrafficInfo.

Several other mobile crowdsensing based transit tracking ideas have been published recently. For instance, the authors in [5] propose a bus arrival time prediction system based on bus passengers' participatory sensing. The proposed system uses movement statuses, audio recordings and mobile celltower signals to identify the vehicle and its actual position. The authors in [6] propose a method for transit tracking using the collected data of the accelerometer and the GPS sensor on the users' smartphone. The authors in [7] use smartphone sensors data and machine learning techniques to detect motion type, e.g., traveling by train or by car. EasyTracker [8] provides a low cost solution for automatic real-time transit tracking and mapping based on GPS sensor data gathered from mobile phones which are placed in transit vehicles. It offers arrival time prediction, as well.

These approaches focus on the data (what to collect, how to collect, what to do with the data) to offer enriched services to the users. However, our focus is on how to introduce such enriched services incrementally, i.e., how can we create an architecture and service model, which allows incremental introduction of live updates from participatory users over static services that are available in competing approaches. Thus, our approach complements the above ones.

FRAMEWORK FOR CROWDSENSING BASED SMART CITY APPLICATIONS

In this section, we shortly describe our generic framework [3], which is based on the XMPP publish-subscribe architecture, to aid the development of crowdsensing based smart city applications. TrafficInfo is implemented on top of this framework.

² <http://www.moovitapp.com>

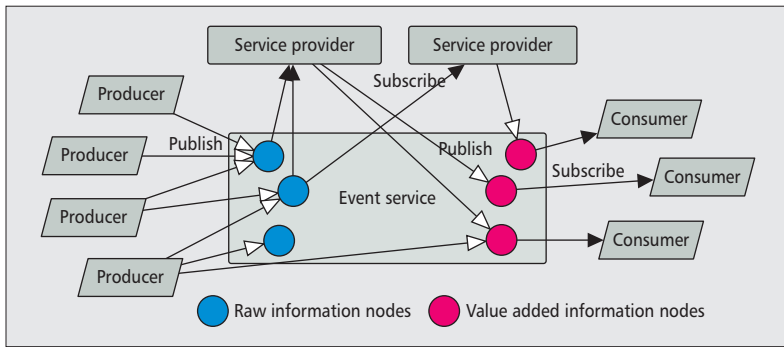


Figure 2. Crowdsensing model based on publish/subscribe communication.

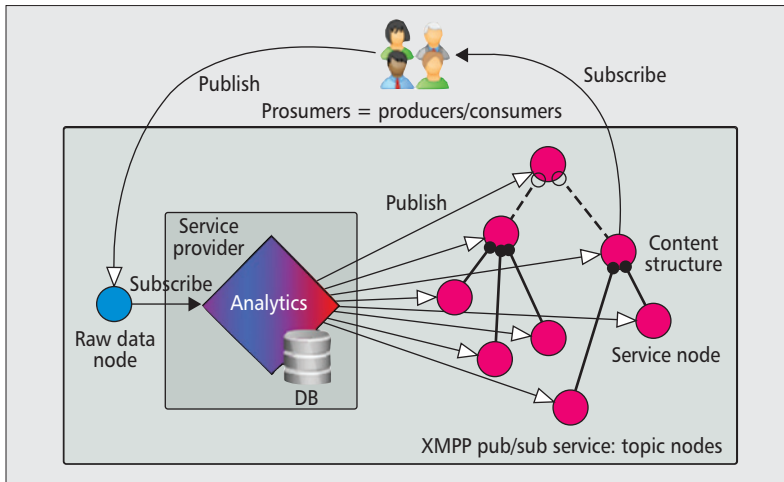


Figure 3. Mobile crowdsensing: the publish/subscribe value chain using XMPP.

COMMUNICATION MODEL

XMPP [2] is an open technology for real-time communication using Extensible Markup Language (XML) [9] message format. XMPP allows sending of small information pieces from one entity to another in quasi real-time. It has several extensions, like multi-party messaging [10] or the notification service [11]. The latter realizes a publish/subscribe (pub/sub) communication model [12], where publications sent to a node are automatically multicast to the subscribers of that node. This pub/sub communication scheme fits well with most of the mobile crowdsensing based applications. In these applications, the users' mobile devices are used to collect data about the environment (publish) and the users consume the services updated on the basis of the collected data (subscribe).

Hence, we use XMPP and its generic publish/subscribe communication model in our framework to implement interactions. In this model, we define three roles, like *Producer*, *Consumer* and *Service Provider* (Fig. 2). These entities interact with each other via the core service, which consists of event based pub/sub nodes.

Producer: The Producer acts as the original information source in our model producing raw data streams and plays a central role in data collection. He is the user who contributes his mobile's sensor data.

Consumer: The Consumer is the beneficiary of the provided service(s). He enjoys the value

of the collected, cleaned, analyzed, extended and disseminated information. We call the user as *Prosumer*, when he acts in the service as both Consumer and Producer at the same time.

Service Provider: The Service Provider introduces added value to the raw data collected by the crowd. Thus, he intercepts and extends the information flow between Producers and Consumers. A Service Provider can play several roles at the same time, as he collects (Consumer role), stores and analyzes Producers' data to offer (Service Provider role) value added service.

In our model, depicted in Fig. 2, Producers are the source of original data by sensing and monitoring their environment. They publish (marked by arrows with empty arrowhead) the collected information to event nodes (raw information nodes are marked by blue dots). On the other hand, Service Providers intercept the collected data by subscribing (marked by arrows with black arrowhead) to raw event nodes and receiving information in an asynchronous manner. They extend the crowdsensed data with their own information or extract cleaned-up information from the raw data to introduce added value to Consumers. Moreover, they publish their service to different content nodes. Consumers who are interested in the reception of the added value/service just subscribe to the appropriate content node(s) and collect the published information also in an asynchronous manner.

ARCHITECTURE

We can directly map this model to the XMPP publish/subscribe service model as follows (Fig. 3):

- Service Providers establish raw pub/sub data nodes, which gather Producers' data, for the services they offer.
- Consumers can freely publish their collected data to the corresponding nodes with appropriate node access rights, too. However, only the owner or other affiliated Consumers can retrieve this information.
- Producers can publish the collected data or their annotations to the raw data nodes at the XMPP server only if they have appropriate access rights.
- Service Providers collect the published data and introduce such a service structure for their added value via the pub/sub subscription service, which makes appropriate content filtering possible for their Consumers.
- Prosumers publish their sensor readings or annotations into and retrieve events from XMPP pub/sub nodes.
- Service Providers subscribed to raw pub/sub nodes collect, store, clean-up and analyze data and extract/derive new information introducing added value. This new information is published into pub/sub nodes on the other side following a suitable structure.

The pub/sub service node structure can benefit from the aggregation feature of XMPP via using collection nodes, where a collection node will see all the information received by its child nodes. Note, however, that the aggregation mechanism of an XMPP collection node is not appropriate to filter events. Hence, the Service Provider role has to be applied to implement scalable content aggregation. Figure 3 shows XMPP aggregations

as dark circles at the container node while empty circles with dashed lines represent only logical containment where intelligent aggregation is implemented through the service logic.

REAL-TIME PUBLIC TRANSPORT INFORMATION SERVICE

In this section, we shortly overview the architecture of our public transport information service, then describe TrafficInfo, its front-end Android interface together with our stop event detector.

SERVICE ARCHITECTURE

Our real-time public transport information service architecture has two main building blocks, such as our crowdsensing framework described earlier and the TrafficInfo application (Fig. 4). The framework can be divided into two parts, a standard XMPP server and a GTFS Emulator with an analytics module.

The XMPP server maps the public transport lines to a hierarchical pub/sub channel structure. We turned the GTFS database into an XMPP pub/sub node hierarchy. This node structure facilitates searching and selecting transit feeds according to user interest.

Transit information and real-time event updates are handled in the *Trip* nodes at the leaf level. The inner nodes in the node hierarchy contain only persistent data and references relevant to the trips. The users can access the transit data via two ways, based on *routes* or *stops*. When the user wants to see a given trip (vehicle) related traffic information the route based filtering is applied. On the other hand, when the forthcoming arrivals at a given stop (location) are of interest the stop based filtering is the appropriate access way.

The GTFS Emulator provides the static timetable information, if it is available, as the initial service. It basically uses the officially distributed GTFS database of the public transport operator of the given city. However, it also relies on another data source, which is OpenStreetMap (OSM), a crowdsourcing based mapping service [13]. In OSM maps, users have the possibility to define terminals, public transportation stops or even public transportation routes. Thus, the OSM based information is used to extend and clean the information coming from the GTFS source. The analytics module is in charge of the business logic offered by the service, e.g., deriving crowdedness information or estimating the time of arrivals at the stops from the data collected by the crowd.

TrafficInfo handles the subscription to the pub/sub channels, collects sensor readings, publishes events to and receives updates from the XMPP server, and visualizes the received information.

TRAFFICINFO FEATURES

TrafficInfo has three main features, but most of the users will benefit from its visualization capability that visualizes public transport vehicle movements on a city map.

1) *Visualization*: An example of this primary feature can be seen on Fig. 5a displaying trams 1, 4, 6 and buses 7 and 86 on the Budapest map in Hungary. The depicted vehicles can be fil-

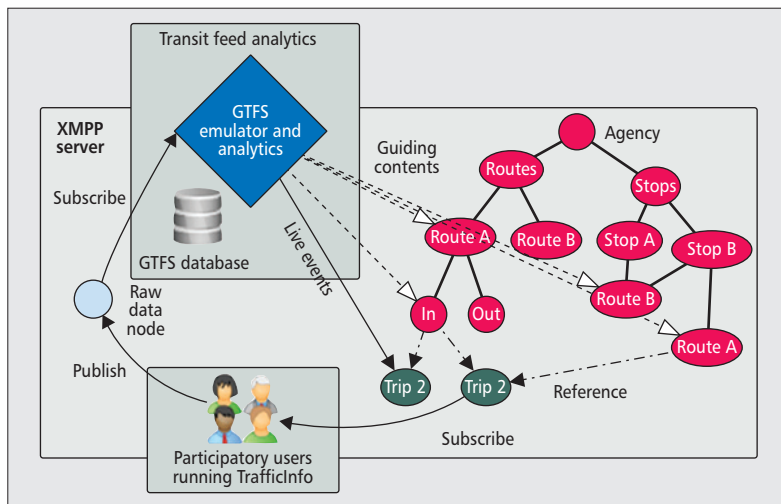


Figure 4. Real-time public transport information service architecture.

tered to given routes. The icon of a vehicle may reflect various attributes, such as the number, progress or crowdedness of the specific vehicle. Clicking on a vehicle's icon a popup shows all known information about that specific vehicle.

2) *Information Sharing*: The second feature is about information sharing. Passengers can share their observations regarding the vehicles they are currently riding. Figure 5b shows the feedback screen that is used to submit the observations. The feedback information is spread out using our crowdsourcing framework and displayed on the devices of other passengers, who might be interested in it. It is up to the user what information and when he wants to submit, but we are planning to provide incentives to use this feature frequently.

3) *Sensing*: The third feature is collecting smartphone sensor readings without user interaction, which is almost invisible for the user. User positions are reported periodically and are used to determine the vehicle's position the passenger is actually traveling on. In order to create the link between the passenger and the vehicle, we try to identify the movement of the user through his activities. To this end we are using various sensors, e.g., accelerometer, and try to deduct the timestamped stop events of the vehicles (our automatic stop event detection mechanism is described earlier). The duration between the detected stops coupled with GPS coordinates identifies the route segment, which the user actually rides.

Besides the GPS coordinates Google also provides location information on those areas, where there is no GPS signal. Usually this position is highly inaccurate, but the estimated accuracy is also provided. We also use the activity sensor, which guesses the actual activity of the user. Currently, the supported activities are: *in vehicle, on bicycle, on foot, running, still, tilting, walking and unknown*. Accuracy is provided here, as well.

The collected sensor readings, on one hand, are uploaded to the XMPP server, where the analytics module processes and shares them among parties who are subscribers of the relevant information; on the other hand, are used locally. For example, user activity is analyzed on the server side and it is used to create non real-time stop patterns through machine learning.

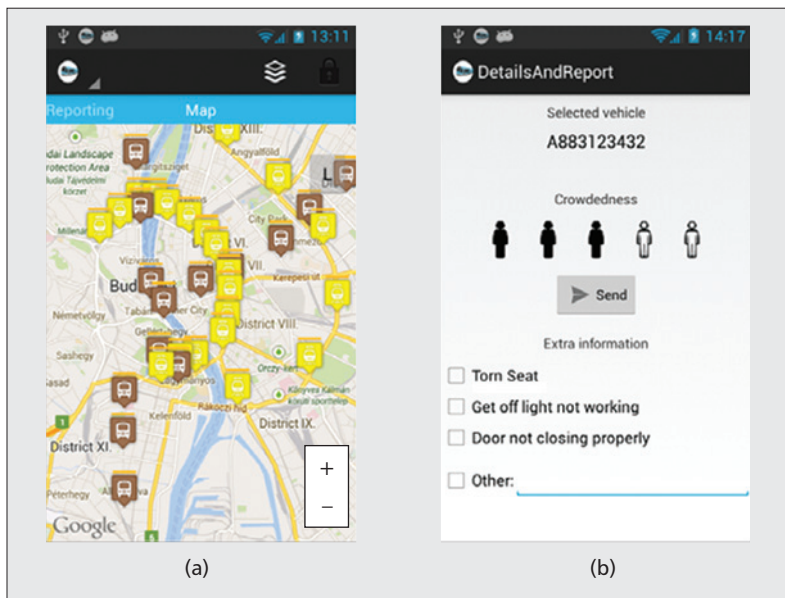


Figure 5. TrafficInfo screenshots: a) vehicle visualization; b) user feedback form.

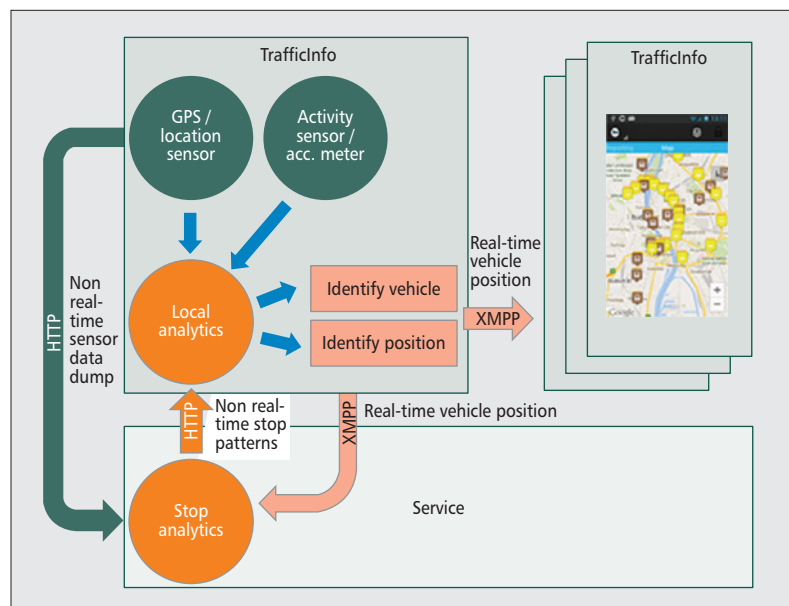


Figure 6. Sensor data flows in TrafficInfo.

These patterns are delivered back to the application, where further local analytics can use them for, e.g., identifying the vehicle and providing position information. These sensor data flows are depicted in Fig. 6. Note that, at the moment, stop events are detected locally on the device due to resource usage reasons and only the detected events with a timestamp are reported back to the server. Based on this information the server side analytics estimate the upcoming arrival times of the given vehicle and disseminate live timetable updates to the subscribers.

SERVICE LEVELS

Running TrafficInfo in a small city is different than in big cities, like Budapest. The cause of this difference is the unavailability of static pub-

lic transportation information in, e.g., GTFS format. If even the static public transportation schedule is not presented by the application, people will likely not use it. Furthermore, fewer users will generate less live traffic data which makes the whole application useless. Hence, it is clear that we should apply a different approach in cities where static public transportation information has not been available, yet.

1) *City Explorer*: In a new city, at the beginning we assume that there is zero knowledge in our system about the city's public transportation. The goal is to gather the relevant information in a fast and inexpensive way. When a reliable GTFS or OSM information base of the given city is available, we import this data into our databases. In other situations, we use crowdsourcing to gather this information. We assume that some users will install the TrafficInfo application either to contribute to city exploration or just simply for curiosity (or some incentive mechanism has to be introduced to grab users). We expect no other contributions than installing the application, carrying the smartphone during the day, traveling on public transportation and answering some simple questions asked by the application. The smartphones, using their built-in sensors, collect all the necessary data without user interaction. The questions are used to annotate the collected data.

Every day the captured data is uploaded in a batch to the server for analysis. At the same time, the application downloads information about what to ask on the following day(s).

Figure 8 depicts an uploaded activity log of a particular user. In this example, the information source is the Google activity recognition module mentioned above. The blue bars show the detected activity during the capture time. In addition, another sensor module recorded the motion, too. Its output is the red bars, recognizing only still or moving states. The height of the bars expresses the confidence of the recognition. Although the values represented with blue and red bars are coming from two different sensors, they usually have the same results for the still state. There are only a few differences, where the activity recognition shows unknown event, while the motion sensor signals still state. It is not displayed on the figure, but the GPS position and its accuracy are also logged for every event.

The captured logs are processed by the server during the night, when the users are typically inactive and the system tries to guess the public transport stops and routes in the city. The more users report the same information the higher the chance is to guess the transportation system correctly. A database stores all the possible stop locations together with their confidence. This database is then downloaded to the application which will ask simple questions to the users to identify stops. For instance, the application might ask: "Are you standing at a stop, waiting for public transport?" We expect simple answers for simple questions until we can construct the public transportation stop database. Routes are explored in a similar way. When the user travels between already known stops, we assume that there is a public transport route among these stops. The application might ask the user about the route type and the line number.

2) *Schedule (Re)Construction*: Once the public transportation stops and routes are explored in most parts of the city, we can assume with high confidence that more users join and use the application. Visualizing stops and routes aids users to get orientation. However, the exploration of the city is continuing, the sensor readings are always collected, but questions are asked only regarding to the partially explored areas.

When the number of users exceeds a certain level and the trips can be guessed, the automatic detection of the stop events comes into the picture. The detected events are reported to the server by the application. The server filters this data and analyzes the patterns of each transport line. As more stop events are captured the patterns are more complete and finally the public transportation schedule is constructed.

3) *Live Schedule*: TrafficInfo providing public transportation stops, routes and schedules is assumed to attract many users, similarly to those applications that are available in big cities based on GTFS data. One advantage of TrafficInfo is that it provides an alternative way to collect all necessary information from scratch which does not require the cooperation of the public transport operator company, rather relies on the power of the crowd.

When the number of users is high enough and (static) schedule information is available, the continuously collected position and stop event data is used to create and propagate lively updates. These updates refresh the timetable if necessary and reflect the actual public transport traffic conditions.

4) *Information Sharing on Public Transport Conditions*: On-line users are able to send and receive information about the vehicle's conditions they are actually riding. This requires user interaction on a voluntary basis as current sensors are not able to detect crowdedness, torn seats, bad drivers, etc. If the application has a wide user base we can always expect some volunteers to report on such conditions. The application provides easy to use forms to enter the relevant data.

5) *Additional Services*: When TrafficInfo is running in a full-fledged manner, it can cooperate with other services targeting public transportation. For example, a rendezvous service can be paired to the TrafficInfo application to organize dates on public transportation vehicles.

STOP EVENT DETECTION

One of the fundamental functions of TrafficInfo is to detect stop events of public transport vehicles. We implemented such a detector locally on the mobile device. The reason behind that is twofold. First, cheaper devices produce bogus raw GPS location data that, if directly transmitted to the XMPP server, would mislead the service. Second, raw logs are generated at a very high rate and it would cause a substantial burden to transmit the raw logged data to the server in real-time for further processing. Instead, only when stop events are detected a summary of information, e.g., the timestamp of the event and the time elapsed since the last stop event, will be transmitted.

To illustrate the challenge of stop event detection, we show the logged trajectory on tram routes 4 and 6 in Budapest from two devices, a Samsung Galaxy S3 and a Nexus4 smartphone, in Fig. 9

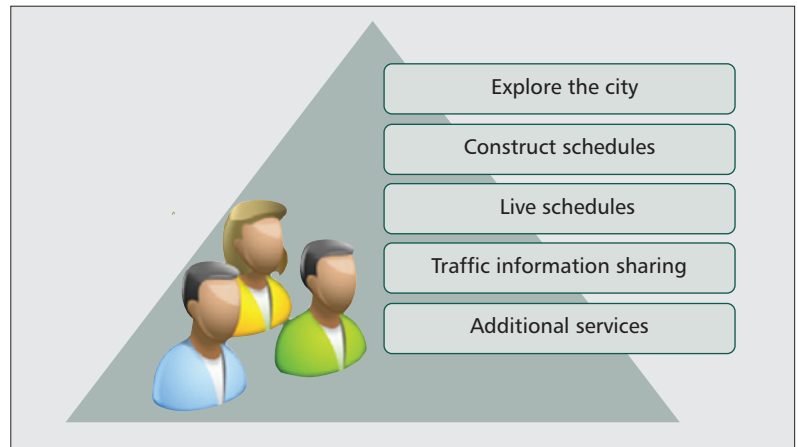


Figure 7. Service levels vs. user base.

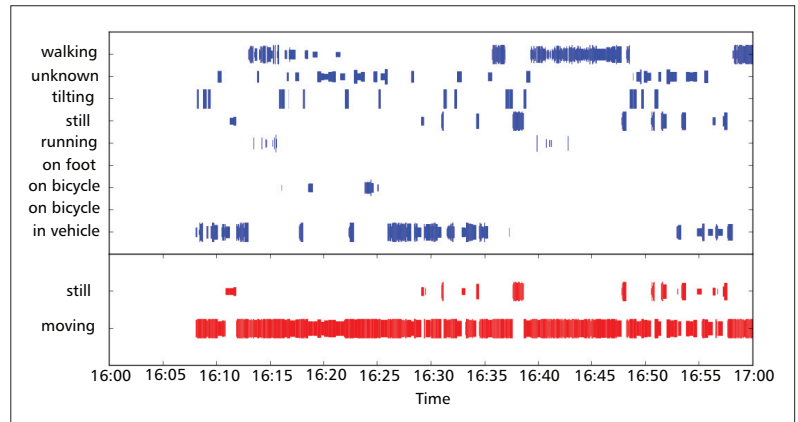


Figure 8. Captured sensor data during user activity.

and Fig. 10, respectively. In case of Nexus4 (Fig. 10), yellow dots indicate the predicted locations of the stop events. Note that Nexus4 with network information provides correct position data, similar in quality to the Galaxy S3 device. Unfortunately, we were not able to collect GPS position data from all the devices we used in our experiments even if the device was equipped with GPS sensor.

Our solution for stop event detection is based on features. Hence, we generated several features from the experimental usage logs collected during the testing period. The measurement object we used to collect context data is summarized in Table 1. It includes among others GPS, WiFi, network and acceleration sensor readings, etc.

The features we defined are the following:

- Latitude, Longitude: raw GPS data;
- AccAbsMax and AccAbsMin: maximum and minimum value of acceleration in the past 20 seconds;
- Last Annotation Time: in seconds, depending on the annotation type (Stopped at Station or Leaving Stop);
- Closest Station: distance calculated from raw GPS data;
- GPS Distance: distance traveled during the last 20 seconds based on raw GPS data.

We collected Android sensor and location data by using the Android Location API.³ The device can have multiple LocationProvider sub-

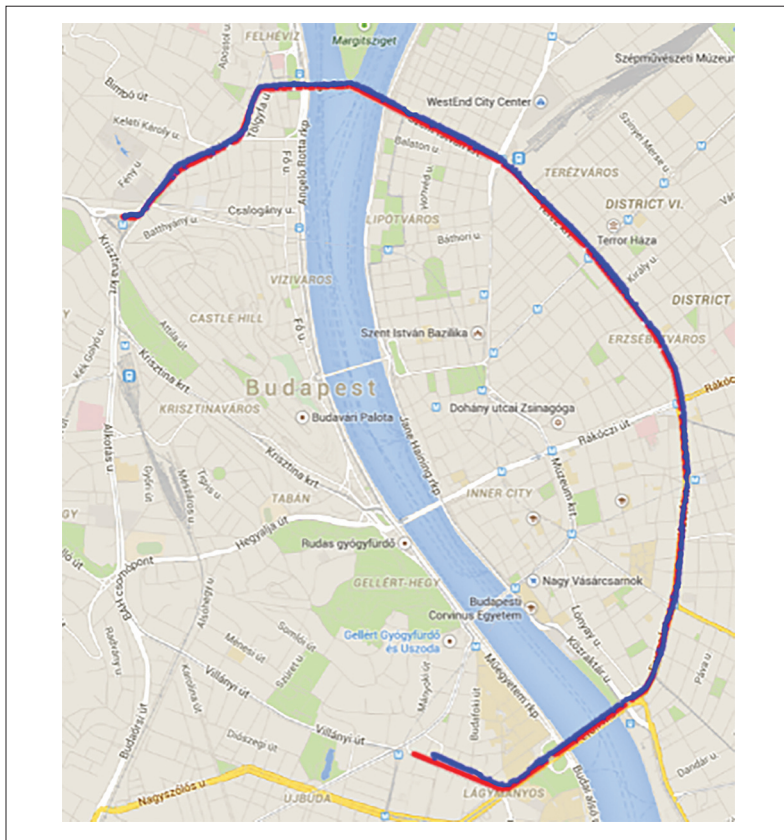


Figure 9. GPS position trajectory (blue) and the real tram route (red) as logged by a Samsung Galaxy S3 device.

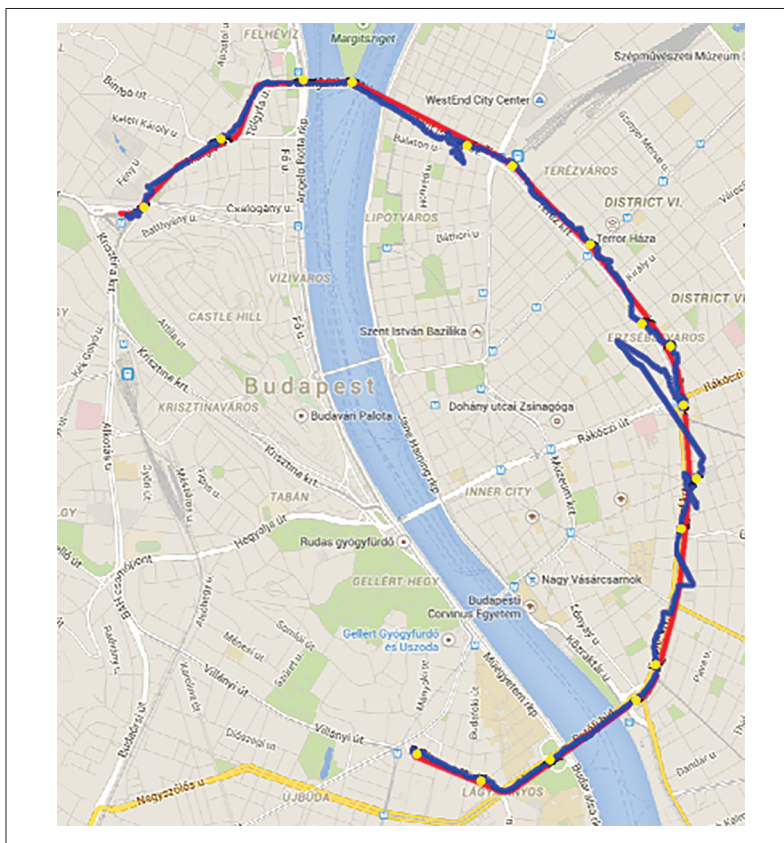


Figure 10. GPS position trajectory (blue), the real tram route (red) and stops (yellow dots) as logged and detected by a Nexus4 device.

classes based on network, GPS and passive sensors, and the location manager has a method to select the best provider. Accessing the sensors requires three level permissions: `ACCESS_FINE_LOCATION`, `ACCESS_COARSE_LOCATION`, `INTERNET`. The GPS sensor can be accessed by the NMEA listener.⁴ Accelerometer is accessible through the Google Location Services API, part of Google Play Services, a high level framework that automates the location provider choice.

For classification we used the J48 decision tree implementation of the Weka data mining tool.⁵ The final output of our detector is the detected stop event, including location and timestamp. With the combination of the defined features and models we could detect stop events with high accuracy within 13 seconds after the arrival at the station.

We measure the accuracy of the method by computing the precision, recall and AUC (Area Under the Curve) [14] of our classifiers in a 10-fold crossvalidation setting. We consider AUC as the main stable measure for classifier performance that does not depend on the decision threshold separating the predicted stop events. The best classifier reached precision 0.97, recall 0.95, F-measure 0.96. The corresponding best AUC was 0.86, which means that a random time point when the tram is at a stop is predicted 86% more likely a stop than another random time point when the tram is in between two stops. In general, an AUC between 0.8-0.9 is considered in the literature to be good to excellent.

SUMMARY

In this paper, we shortly introduced our XMPP based communication framework that we designed to facilitate the development of crowd assisted smart city applications. Then we presented our crowdsensing based real-time public transport information service, implemented on top of our framework, and its front-end Android application, called TrafficInfo, in detail together with our stop event detector. This detector was developed to automatically detect halt events of public transport vehicles at the stops.

As future work, we plan to develop TrafficInfo further and enhance the different services of all the introduced service levels. Moreover, we intend to recruit a noticeable user base and carry out field experiments with these real users. Their feedback is important to plan the directions for improvements.

ACKNOWLEDGMENT

The publication was supported by the TÁMOP-4.2.2.C-11/1/KONV-2012-0001 project. Károly Farkas has been partially supported by the Hungarian Academy of Sciences through the Bolyai János Research Fellowship.

³ <http://developer.android.com/guide/topics/sensors/index.html>

⁴ <https://developer.android.com/reference/android/location/GpsStatus.NmeaListener.html>

⁵ <http://www.cs.waikato.ac.nz/ml/weka/>

REFERENCES

- [1] R. Ganti, F. Ye, and H. Lei, "Mobile Crowdsensing: Current State and Future Challenges," *IEEE Commun. Mag.*, Nov. 2011, pp. 32–39.
- [2] P. Saint-Andre, "Extensible Messaging and Presence Protocol (XMPP): Core," RFC 6120 (Proposed Standard), Internet Engineering Task Force, Mar. 2011, <http://www.ietf.org/rfc/rfc6120.txt>.
- [3] R. L. Szabo and K. Farkas, "A Publish-Subscribe Scheme Based Open Architecture for Crowd-Sourcing," *Proc. 19th EUNICE Wksp. Advances in Communication Networking (EUNICE 2013)*, Springer, Aug. 2013, pp. 1–5.
- [4] Google Inc., General Transit Feed Specification Reference, <https://developers.google.com/transit/gtfs/reference/>.
- [5] P. Zhou et al., "How Long to Wait?: Predicting Bus Arrival Time with Mobile Phone based Participatory Sensing," *Proc. 10th Int'l. Conf. Mobile Systems, Applications, and Services (MobiSys 2012)*, June 2012.
- [6] A. Thiagarajan et al., "Cooperative Transit Tracking Using Smart-phones," *Proc. 8th ACM Conf. Embedded Networked Sensor Systems (SenSys 2010)*, 2010, Nov. 2010, pp. 85–98.
- [7] L. Bedogni, M. Di Felice, and L. Bononi, "By Train or by Car? Detecting the User's Motion Type Through Smartphone Sensors Data," *Proc. IFIP Wireless Days Conf. (WD 2012)*, 2012, pp. 1–6.
- [8] J. Biagioni et al., "EasyTracker: Automatic Transit Tracking, Mapping, and Arrival Time Prediction Using Smartphones," *Proc. 9th ACM Conf. Embedded Networked Sensor Systems (SenSys 2011)*, pp. 1–14.
- [9] T. Bray et al., Extensible Markup Language (XML) 1.0 (Fifth Edition), W3C, W3C Recommendation REC-xml-20081126, Nov. 2008, <http://www.w3.org/TR/2008/REC-xml-20081126/>.
- [10] P. Saint-Andre, "XEP-0045: Multi-User Chat," XMPP Standards Foundation, Standards Track XEP-0045, Feb. 2012, <http://xmpp.org/extensions/xep-0045.html>.
- [11] P. Millard, P. Saint-Andre and R. Meijer, "XEP-0060: Publish-Subscribe," XMPP Standards Foundation, Draft Standard XEP-0060, July 2010, <http://xmpp.org/extensions/xep-0060.html>.
- [12] P. Eugster et al., "The Many Faces of Publish/Subscribe," *ACM Comput. Surv.*, vol. 35, no. 2, June 2013, pp. 114–31.
- [13] M. Haklay and P. Weber, "OpenStreetMap: User-Generated Street Maps," *IEEE Pervasive Computing*, Oct. 2008, <http://dx.doi.org/10.1109/MPRV.2008.80>.
- [14] J. Fogarty et al., "Case Studies in the Use of ROC Curve Analysis for Sensor-based Estimates in Human Computer Interaction," *Proc. Graphics Interface 2005, ser. GI '05*, School of Computer Science, University of Waterloo, Waterloo, Ontario, Canada: Canadian Human-Computer Communications Society, 2005, <http://portal.acm.org/citation.cfm?id=1089508.1089530>, pp. 129–36.

BIOGRAPHIES

KÁROLY FARKAS (farkask@hit.bme.hu) received his Ph.D. degree in Computer Science in 2007 from ETH Zurich, Switzerland, and his M.Sc. degree in Computer Science in 1998 from the Budapest University of Technology and Economics (BME), Hungary. Currently he is working as an associate professor at BME. His research interests cover the field of communication networks, especially autonomic, self-organized, wireless and mobile ad hoc networks, and mobile crowdsourcing. He has published more than 70 scientific papers in different journals, conferences and workshops and he has given a plenty of regular and invited talks. In the years past, he supervised a number of student theses and coordinated or participated in several national and international research projects, such as CityCrowdSource of EIT ICTLabs. Moreover, he acted as program committee member, reviewer and organizer of numerous scientific conferences, thus he took the general co-chair role of the IEEE PerCom 2014 conference and the TPC co-chair role of the CROWDSENSING 2014 and the CASPer 2015 workshops. He is the coordinator of the local Cisco Networking Academy and was the founding initiator of the Cisco IPv6 Training Laboratory and the BME NetSkills Challenge student competition at BME. Between 2012–2015 he has been awarded the Bolyai János Research Fellowship of the Hungarian Academy of Sciences, and received the Gold Medal and Pollák-Virág awards of HTE (Scientific Association for Informatics Hungary) in 2015.

GÁBOR FEHÉR graduated in 1998 at the Budapest University of Technology and Economics on the Faculty of Electronic

Field description	Examples, possible values
Event type	Initialization, manual, sensor
Timestamp	Time when the event occurred
Track (tram, bus line)	Tram 6
Phone type	E.g., Nexus4 including IMEI
Acceleration	Absolute or axes X, Y, and Z
GSM signal strength	As defined in the relevant standard
Android GPS, network, and passive location accuracy, longitude and latitude values	Android GPS Location Provider data, accuracy radius with 68% confidence
CellID, WiFi MACID	LAC (location area code) and CID (Cell ID)
Vehicle number	ID of the transport vehicle
Direction	Onward or backward
Arrived at	Time of arrival at the stop
Manual input	–Stopped at station –Revoke stopped at station –Leaving stop –Revoke leaving stop –Stopped at traffic light –Revoke stopped at traffic light –Revoke last input

Table 1. Semantics of TrafficInfo measurement logs.

Engineering and Informatics. In 2004 he received a Ph.D. degree, the topic of his thesis was resource control in IP networks. Currently he is an associate professor at the same university. Besides giving lectures, he is also contributing to various national and international research projects. From 2004 he is continuously involved in three consecutive EU founded IST/ICT projects. He has teaching activity on the faculty's Smart City specialization working with microelectronics, sensor networks and smartphones. He and his students are working on more projects with crowdsourcing and crowdsensing, from the basic research up to the prototype applications.

ANDRÁS BENCZÜR received his Ph.D. at the Massachusetts Institute of Technology in 1997. Since then his interest turned to Data Science. He is the head of 30 doctoral students, post-docs and developers at the Institute for Computer Science and Control of the Hungarian Academy of Sciences (SZTAKI). He is site coordinator in the Hungarian FutureICT project, and cloud computing activity leader in the Budapest node of EIT ICTLabs. He serves on the program committees of leading conferences including WWW, WSDM, ECML/PKDD, he was Workshop Chair for WWW 2009 and main organizer of the ECML/PKDD Discovery Challenge 2010. In 2012 he was awarded the Momentum grant of the Hungarian Academy of Sciences for his research in Big Data.

CSABA SIDLÓ started working on data warehousing projects and application driven research problems of extremely large data sets in 2000. He joined the Institute for Computer Science and Control of the Hungarian Academy of Sciences (SZTAKI) in 2004; he is now head of the Big Data Business Intelligence research group. His main interest is Big Data analytics and business intelligence on scalable distributed architectures. His industrial projects include master data entity resolution, integration and analytics of log, web and location data. He is currently involved in several big data projects for web, telecom and sensor data analytics. He authored several research papers and book chapters, and has a Ph.D. in Informatics from Eötvös University, Hungary.

Network Distance Prediction for Enabling Service-Oriented Applications over Large-Scale Networks

Haojun Huang, Hao Yin, Geyong Min, Dapeng Oliver Wu, Yulei Wu, Tao Zuo, and Ke Li

ABSTRACT

Knowledge of end-to-end network distances is essential to many service-oriented applications such as distributed content delivery and overlay network multicast, in which the clients have the flexibility to select their servers from among a set of available ones based on network distance. However, due to the high cost of global measurements in large-scale networks, it is infeasible to actively probe end-to-end network distances for all pairs. In order to address this issue, network distance prediction has been proposed by measuring a few pairs and then predicting the other ones without direct measurements, or splicing the path segments between each pair via observation. It is considered important to improve network performance, and enables service-oriented applications over large-scale networks. In this article, we first illustrate the basic ideas behind network distance prediction, and then categorize the current research work based on different criteria. We illustrate how different protocols work, and discuss their merits and drawbacks. Finally, we summarize our findings, and point out potential issues and future directions for further research.

INTRODUCTION

Networks have become an indispensable part of our daily life as an information platform for communications. Over the past few decades, widespread distributed service-oriented applications, such as peer-to-peer file sharing services and overlay network multicast [1], have evolved considerably beyond the traditional client-server model, where an end user (client) only communicates with one server. In contrast, in the current service-oriented applications, end users have much flexibility in choosing their servers, with little or no information about the potential performance of different communication paths toward them. In order to overcome this issue, network distance prediction (NDP) has been proposed, which provides benefit for end users in selecting intelligent paths based on network

performance and constructs much more convenient networks. For instance, in content delivery networks, an end user can conveniently obtain its desired web resources from a particular site with knowledge of the predicted network distance.

Similar to [2–15], in this article the network distance between two nodes is defined as the communication delay or latency between them, in the form of either one-way delay or more often round-trip time (RTT). Obviously, it is infeasible to ceaselessly probe network distances among all pairwise nodes in large-scale networks because globally accurate measurements are difficult and costly to achieve and maintain. A natural idea is to probe a small set of pairs and then predict the distances between others without direct measurements, or splice the path segments between each pair; this is NDP. This understanding has motivated a great deal of research on developing NDP [1]. Figure 1 illustrates NDP operations by matrix factorization [9], with four landmarks and two ordinary hosts. Based on knowledge of inter-landmark distances, the distance matrix among landmarks L_1 , L_2 , L_3 , and L_4 can be factorized to the incoming and outgoing vectors for distance prediction among the other hosts. Given an ordinary host H_1 , which desires to know the distance to another host (e.g., H_2), it first measures its distance vectors to and from landmarks L_1 , L_2 , L_3 , and L_4 as $[0.5, 1.5, 1.5, 2.5]$. Then it can calculate its outgoing and incoming vectors as $Y_{H_1} = [1.5, 0, 1]$ and $X_{H_1} = [-1.5, 0, 1]$ (the calculation shown in [9]), respectively. Similarly, H_2 can obtain its distance vector as $[2.5, 1.5, 1.5, 0.5]$, and its outgoing and incoming vectors: $X_{H_2} = [-1.5, 0, -1]$ and $Y_{H_2} = [-1.5, 0, 1]$. If H_1 learns the incoming vector of H_2 , the distance between them can then be predicted as $X_{H_1} \cdot Y_{H_2} = 3.25$ with a tolerable predicted error of 8.3 percent, instead of relying on direct measurements. This means that as long as an acceptable predicted network distance can be obtained for host H_1 , the small measurement cost can be neglected and the remaining overhead is amortized over all distance predictions.

Haojun Huang is with Wuhan University and Tsinghua University.

Hao Yin is with Tsinghua University.

Tao Zuo is with Wuhan University of Science and Technology.

Geyong Min and Yulei Wu are with the University of Exeter.

Oliver Wu is with the University of Florida.

Ke Li is with Southwest Jiao Tong University.

NDP has been considered important to improve the performance of many service-oriented applications and bridge the gap between end users and large-scale networks, and thus has received increasing attention. However, the existing approaches have proven to be difficult to use in deployed applications [1–3, 5, 8]. In this article, we survey the various NDP approaches reported in the current literature, illustrated in Fig. 2, as a reference for further research. We analyze their emerging challenges and discuss future NDP developments.

The remainder of this article is organized as follows. The next section illustrates how different approaches work, and discusses their merits and drawbacks, followed by the existing evaluation metrics. Then the emerging challenges behind NDP are highlighted, and open issues and opportunities for further research are outlined. The final section draws our conclusions.

SURVEY ON NETWORK DISTANCE PREDICTION

Essentially, NDP resorts to the predicted distance without performing direct measurements, with the aid of infrastructures such as landmarks, tracers, Domain Name Service (DNS) servers, and routers, to represent the actual distance for a given host pairwise. It means that the predicted distances should be very close to the actual ones. Therefore, NDP makes sense if and only if an acceptable predicted accuracy, with quantified evaluation metrics investigated in the next section, can be guaranteed. The major challenges in designing NDP include symmetry, consistency, security, dynamics, cluster, and triangle inequality violations (TIVs), as elaborated in the following section. Currently, NDP approaches can be classified into three categories: coordinate-based approaches, path fitting approaches, and data-driven approaches. Figure 2 depicts most classic NDP approaches in these categories, and their historical development and evolution. In this section, we explain their operations with special focus on how they work, and discuss their merits and drawbacks.

COORDINATE-BASED APPROACHES

The basic ideal of such approaches is to design a finite-dimension virtual metric space and embed the hosts into that space with the constraint that errors between predicted distances and measured distances are minimized. The network distance between two reachable hosts is then predicted as the distance between their coordinates.

Figure 3 elaborates a matching between the large-scale networks such as the Internet and the metric space. The Internet network distances between four hosts are represented by green lines in Fig. 3a. This distance can be, for instance, the RTT. Figure 3b embeds the predicted network distances into a metric space. In such a space, the predicted distance is evaluated using the traditional distance metrics such as the relative error and stress [2, 9, 10]. In the following, we investigate the main coordinate-based NDP approaches.

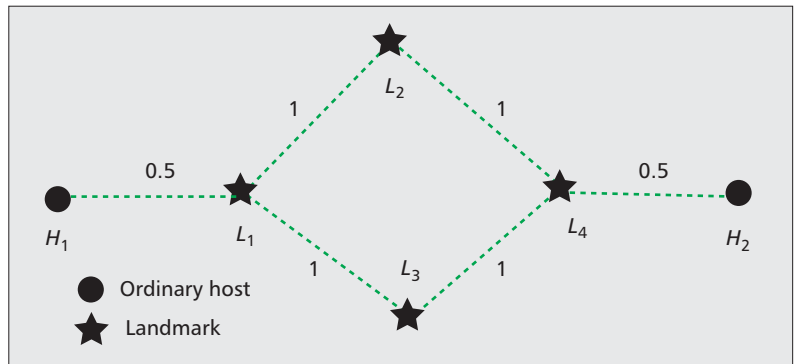


Figure 1. Network distance prediction derived from real measurements.

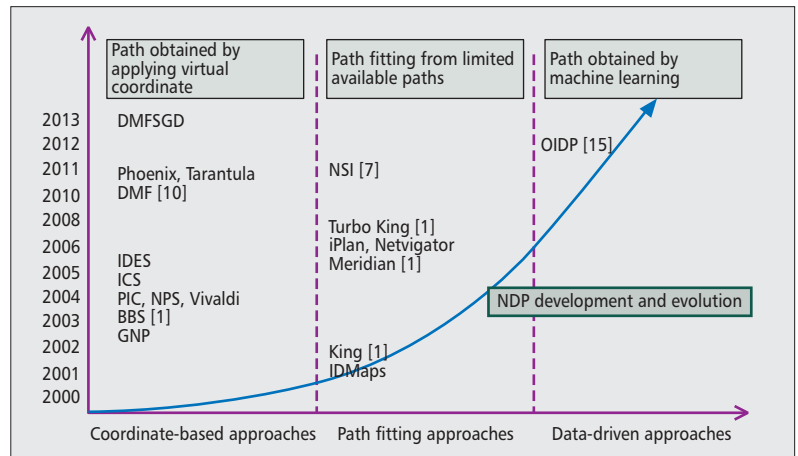


Figure 2. Network distance prediction categories and historical developments.

GNP: Global network positioning (GNP) [2] is the first network embedding system to predict network distance, which relies on a small number of fixed landmarks. It assigns the locations of all hosts in an n -dimensional Euclidean space, n being the number of landmarks. Given any two reachable hosts, GNP approximates the latency between them in the original networks as their distance in this space, with the assumption that the predicted distances among any three hosts satisfy the triangle inequality.

It starts by instructing the n landmarks to measure all the latencies among them. With this information, it computes all the landmark coordinates such that the distances among their coordinates in Euclidean space are as close as possible to their measured latencies in the original networks. Then the ordinary hosts calculate their own coordinates with respect to the landmarks. In this way, any network distance among pairwise hosts can be computed. A distinct drawback of GNP is that it is vulnerable to landmark failures since hosts join and leave networks frequently.

PIC: Practical Internet coordinates (PIC) [3] is the first security-aware mechanism to predict the Internet network distance. In PIC, an ordinary host selects any host with a coordinate that has already been given as a landmark. This is similar to GNP, but GNP selects a fixed set of

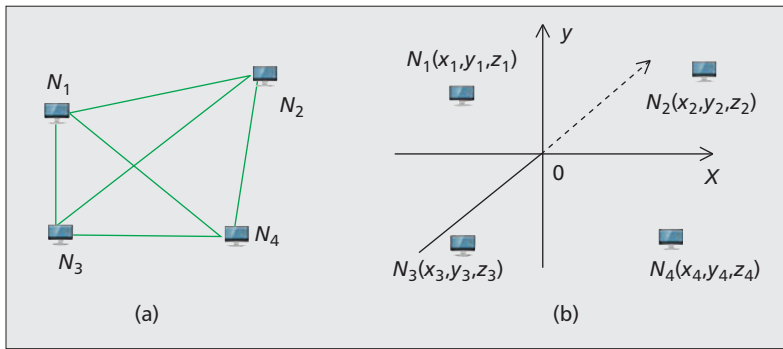


Figure 3. Correspondence between the Internet and the metric space. The network distances are represented in green lines in a) and metric distances in b).

landmarks for all of the ordinary nodes as the reference nodes. On the contrary, it needs to probe the network distance to each landmark, and uses an active node discovery protocol to discover some nearby hosts used as computing coordinates. In this way, it obtains the coordinates of all landmarks and then uses a multidimensional global optimization method (e.g., Simplex DownHill) to compute its coordinates such that the errors in the predicted distances are minimized. In order to prevent network attacks, PIC explores reference select point techniques based on triangle inequality to detect independent malicious participants.

NPS: The network positioning system (NPS) [6] is a hierarchical NDP system that maintains consistency while enabling decentralization to compute accurate and stable network locations. Essentially, it is a GNP extension aiming to overcome its landmarks' failure by selecting any ordinary host that knows its own position to serve as a landmark for other hosts by a membership server. The role of the membership server is to initialize configuration parameters of ordinary hosts, such as identifying the landmarks. In order to ensure the consistency of hosts, NPS imposes a hierarchical position dependency among the hosts, by eliminating any landmark called a malicious node that provides significant relative errors compared to the other landmarks.

ICS: An Internet coordinate system (ICS) [8] is a low-dimensional NDP system that retains as much topology information as possible. In ICS, the distances between a host and landmarks (called beacon nodes in ICS) are expressed as a distance vector, with a dimension equal to the number of landmarks. For any host, it does not need to measure the distance from itself to all the landmarks, but rather to a set of landmarks, and thus obtains its distance vector. The landmarks calculate their locations by multiplying the distance vector with a transformation matrix based on principal component analysis (PCA), which aims at distance dimension reduction by refining the variables that retain most of the original information. Such a linear transformation essentially extracts the topology information from delay measurements between landmarks and keeps it in a novel metric space.

IDES: The Internet distance estimation service (IDES) [9] is the first system based on matrix factorization to predict distances for large-scale networks, and has the same landmark-based architecture as GNP while enduring TIVs. Its essential idea is to approximate a large distance matrix with elements representing pairwise distances by the product of two smaller matrices. Such a model allows a representation of distances violating TIVs and asymmetric distances based on matrix factorization, and can be regarded as a form of dimensionality reduction with both singular value decomposition and non-negative matrix factorization algorithms. In particular, it addresses the questions of the impact of both landmark placement and measurement error on predicted performance.

Tarantula: Tarantula [4] is an alternative hierarchical NDP system, which focuses on mitigating the impact of TIVs. It dynamically divides whole networks into three clusters following the major actual clusters of the world networks — America, Asia-Pacific, and Europe — and integrates each two clusters into a subspace, meaning each host belongs to one cluster and two subspaces at the same time. It runs a Vivaldi [1] system on each cluster or subspace. Considering the inter-cluster links account for more TIVs than other links in the hierarchical NDP [1, 6], Tarantula uses the cluster-based system and subspace-based system to predict the intra-cluster and inter-cluster distances, respectively. In this way, it outperforms the existing hierarchical NDP greatly in predicted accuracy.

DMFSGD: Decentralized matrix factorization by stochastic gradient descent (DMFSGD) [10] is an IDES extension seeking to overcome the limitation generated by the failure of landmarks. Different from IDES, it only requires each host to measure local distances to and from a small set of neighbors and then predict the distances to the other hosts, without explicit matrix constructions and special hosts such as fixed landmarks and central servers. DMFSGD is simple, decentralized, and scalable. With these features, it addresses many practical issues in NDP such as measurement dynamics and network churn with time in large-scale networks.

Phoenix — Phoenix [5] is a weight-based NDP system by matrix factorization. Basically, it is an extension of DMFSGD and IEDS that aims to address the inaccurate coordinate impacts on distance prediction by introducing weights to reference coordinates. It assigns each landmark a graded weight in its coordinate based on its accuracy, and trusts the hosts with the higher weight values more than the others. Therefore, Phoenix can substantially mitigate the impact of error propagation and improve prediction accuracy over the existing approaches.

PATH FITTING APPROACHES

Different from the coordinate-based NDP approaches that treat the network as a black-box, the path fitting approaches utilize the internal network structure such as the network topology, DNS servers, tracers, and existing routing to predict the network distance by spli-

ing the path segments by observation or approximating the path among the nodes closer to the client-server as the network distance from the client to the server. Essentially, the path fitting approaches use direct approximative measurements [10] or reactive aggregations [11] instead of predictions to approximate the network distance for given client-servers.

Figure 4 illustrates how path fitting approaches predict network distance. Given nodes S and D , the network distance between them can be obtained by either splicing a short path segment from S to an intersection I_1 from which a path toward D has been observed, or approximating the distance between S and its nearest tracer (or DNS) T_1 , plus the distance between D and its nearest tracer T_2 , plus the distance between T_1 and T_2 , or even approximately equal to the distance from T_1 to T_2 . Three classic examples of those approaches are IDMaps [10], iPlane [11], and Netvigator [12].

IDMaps: The Internet distance map service (IDMaps) is the first Internet distance prediction system and is considered the predecessor of NDP approaches. It starts by building a simplified overlay topology map of the Internet performed by special hosts (called tracers in IDMaps) based on network measurements. Then it performs the shortest path routing on this map, such that the routing can be used as the predicted network distance for any two reachable hosts with valid IP addresses. Given two hosts, x and y , the predicted distance between them is expressed as the sum of the distance from x to its nearest tracer T_1 , the distance from y to its nearest tracer T_2 , and the shortest routing distance from T_1 to T_2 over this map. IDMaps provides general distance query such that service-oriented applications can easily obtain the network distances among the end users and services.

iPlane: The information plane (iPlane) is a scalable service that aims at providing the accurate Internet path performance predictions for application-level overlay networks. It continuously requires vantage hosts that locate in different geographic regions to map the Internet topology such that they obtain the observed paths with a rich set of link and router attributes such as latency, available bandwidth, and loss rate. With such information, it can predict paths for any two arbitrary reachable nodes by being combined with the measured performance of path segments in the Internet. In order to reduce measurement overhead, it clusters IP prefixes into border gateway protocol atoms such that it generates the target list. Compared to current NDP systems, iPlane not only provides latency prediction between two reachable nodes, but also automatically infers important network behavior information such as loss rate, capacity, available bandwidth, and isolated anomalies.

Netvigator: Network navigator (Netvigator) is an efficient NDP system that focuses on proximity estimation and distance prediction. Initially, it requires each host to send and receive probe packets to and from the landmarks and mile-

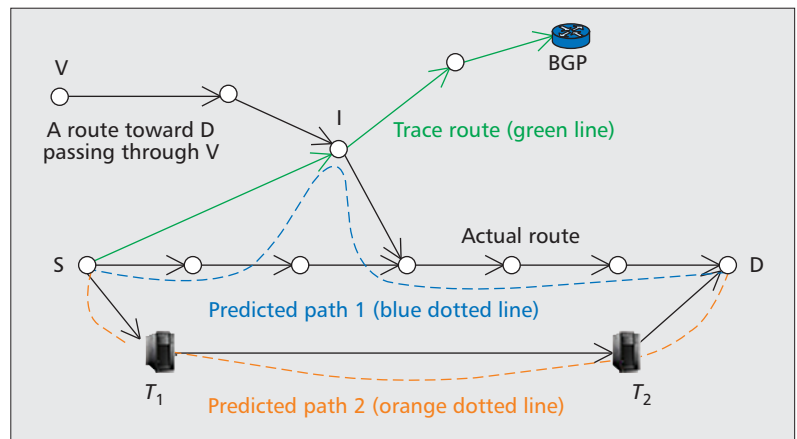


Figure 4. The predicted network distance derived from the path fitting.

stones (also called intermediate routers in Netvigator) if any, and then constructs a landmark vector including the distances to all the landmarks and the milestones through which the probe packets pass, and sends it to a central server with the global information table. Once receiving a query from a client host, the central server applies the clustering algorithm to identify k closest candidates for its proximity distance. In this way, the proximity distance from a client host to some server host can be predicted. A significant merit of Netvigator is that it can avoid false clustering by introducing the milestones used for refinement, and therefore achieves very accurate predictions of proximity estimation.

DATA-DRIVEN APPROACHES

Nowadays, we are in the age of transitioning from a hypothesis-driven world to a data-driven world brought by big data. Big data has changed our most basic understanding of how to comprehend network behavior and explore the Internet. For instance, assuming algorithms A and B can solve the same problem, the result of A is clearly better than B in the small-scale data world, that is, algorithm A can achieve better performance. However, such a situation may not exist as the amount of data increases. There are too many similar phenomena in current networks. These findings have brought computer science and its deuterogenic disciplines a landmark revelation: the data itself (rather than the approaches and models used by data analysis) can guarantee the validity of the data analysis results with the increase of the available data. Even though the approaches or models lack precision, it can make the conclusion closer to the truth as long as enough data is owned.

In the near future, our understanding of network designs will be driven more by the abundance of data rather than hypotheses, described as the fourth paradigm of scientific exploration. We argue that the data-driven idea will transfer the design philosophy of the future Internet in all aspects including architecture design, resource management, task scheduling, and network distance prediction [14]. However, the existing research has not taken full advantage of data-driven thought to steer the design of NDP. The quasi prototype of data-driven NDP

Approaches	Measurement overhead	Prediction prerequisites	Embedding model	Churn recovery	Infrastructure dependability	Scalability	Prediction accuracy
GNP	$O(m^2 + mn)$	A small fixed set of landmarks	Euclidean	No	No	Medium	General
PIC	$O(m^2 + mn)$	A small set of landmarks, P2P substrate	Euclidean	No	No	High	General
NPS	$O(m^2 + mn)$	A small set of landmarks	Euclidean	No	No	Medium	General
Tarantula	$O(m^2 + mn)$	Three clusters and subspaces	Euclidean	Yes	No	High	High
ICS	$O(m^2 + mn)$	A small set of landmarks	Lipschiz	No	No	High	General
IDES	$O(m^2 + mn)$	A small fixed set of landmarks	Matrix factorization	No	No	Medium	High
DMFSGD	$O(m^2 + mn)$	A small set of landmarks	Matrix factorization	Yes	No	High	High
Phoenix	$O(m^2 + mn)$	A small set of landmarks	Matrix factorization	Yes	No	High	High
IDMaps	$O(r^2 + k)$	A small fixed set of tracers	—	Yes	Yes	High	High
iPlane	$O(q^2 + n)$	A small set of vantage hosts mapping the Internet topology	—	Yes	Yes	High	High
Netvigatator	$O(mn + pn)$	A small set of landmarks and milestones	—	Yes	Yes	Medium	General

Note: m : the number of landmarks; n : the number of hosts; r : the number of tracers; k : the number of clustered address prefixes; p : the number of milestones; q : the number of vantage hosts.

Table 1. Characteristics of network distance predictions available in the current literature.

approaches can be found in [15], where an Internet NDP approach seeks to capture geographical characteristics between Internet host pairs by machine learning, instead of relying on direct measurements. Although without explicitly exposing data-driven thought, it still reveals a novel and efficient solution to NDP.

We think that NDP can be executed by employing data-driven approaches. Given client u and server v , the network distance between them, denoted by $P(u, v)$, can be obtained through machine learning by analyzing large amounts of user behavior data. Let $P(t_{i,j})$ be the path from node i to node j with the network distance of $t_{i,j}$. The network distance $P(u, v)$ can be expressed as

$$\begin{cases} P(u, v) = \sum P(t_{u, i_1}) + P(t_{i_1, i_2}) + \dots + P(t_{i_n, v}), \\ s. t. \min(t_{u, i_1} + t_{i_1, i_2} + \dots + t_{i_n, v}), t_{i_j, i_{j+1}} > 0. \end{cases} \quad (1)$$

This achievement depends on enough available data and our ability to harness big data. An inevitable trend is that many data-driven NDP approaches will be proposed in the future to provide intelligent server selection suggestions for clients. In order to facilitate the understanding of the existing NDP approaches comprehensively, Table 1 summarizes the above discussed approaches together with the performance criteria, including measurement overhead, prediction

prerequisites, embedding model, churn recovery, infrastructure dependability, scalability, and prediction accuracy.

EVALUATION METRICS

In order to quantify the magnitude of the differences between predicted distances and original distances, some evaluation metrics have been proposed. Let $d(i, j)$ denote the measured distance, $\bar{d}(i, j)$ the predicted distance computed from some function, ϕ the metric space, and $\phi(x)$ the coordinate of node x in ϕ . The current evaluation metrics can be summarized as follows.

Relative error: The relative error [2, 9–10], denoted by e_r , is defined as

$$e_r = \frac{d(i, j) - \bar{d}(i, j)}{\min(d(i, j), \bar{d}(i, j))},$$

or

$$\frac{|d(i, j) - \bar{d}(i, j)|}{\min(d(i, j), \bar{d}(i, j))}.$$

This metric was proposed to evaluate the accuracy of the distance prediction.

Stress: The stress [10] is given by

$$\sqrt[2]{\frac{\sum_{i,j} (d(i, j) - \bar{d}(i, j))^2}{\sum_{i \neq j} d(i, j)^2}},$$

which measures the overall fitness of the embedding and is used to illustrate the convergence of the proposed NDP schemes.

Median absolute estimation error (MAEE): The MAEE [8] is given by $\text{median}_{i,j} (|d(i,j) - \bar{d}(i,j)|)$. This metric is designed to evaluate the absolute prediction error between predicted distances and measured distances for any pair of reachable nodes.

Distortion: Let

$$r(\phi, x, y) = \frac{\phi(x) - \phi(y)}{d(x, y)},$$

$\text{exp}(\phi) = \max r(\phi, x, y)$, and $\text{con}(\phi) = \min r(\phi, x, y)$ as x and y range. The distortion of ϕ is defined as the ratio of

$$\frac{\text{exp}(\phi)}{\text{con}(\phi)}.$$

The distortion is a worst case measure of the quality of an embedding, and used to measure the worst case change in the relative distances of the embedding.

Local relative rank loss: Let $p(z) = \{(x, y) | x = y \text{ and } \text{swapped}(z, x, y)\}$, and

$$s = \frac{(|N|-1)(|N|-2)}{2}.$$

The local relative rank loss, denoted by $\text{rrl}(\phi, x)$, is defined as

$$\text{rrl}(\phi, x) = \frac{|p(z)|}{s},$$

where $|N|$ is the set of nodes, (x, y) are elements of $N \times N$, and $\text{swapped}(z, x, y)$ is true if the relative relationship of z to x and y is different in the original networks and the embedding space. This metric is designed to reflect the probability that the relationship between any two nodes in the original networks will have a different relative order in embedding space.

Closest neighbors loss: Given node x , the closest neighbors loss, denoted by $\text{cnl}(\phi, x)$, is defined as 0 if the nodes closest to x remain closest in the embedding space ϕ , and 1 otherwise. For n nodes, the average closest neighbors loss is defined by

$$\frac{\sum_{i=1}^n \text{cnl}(\phi, x_i)}{n}.$$

This metric is designed to reflect the average percentage of nodes with closest neighbors that are not preserved in the embedding space.

k -closest preservation: Given node x , let $\text{cn}(k, x)$ denote its k closest neighbors in the original networks, and $\text{cn}(\phi, k, x)$ its k closest neighbors in the metric space ϕ . The k -closest preservation is defined as

$$\frac{|\text{cn}(k, x) \cap \text{cn}(\phi, k, x)|}{k}.$$

This metric is used to reflect the ability of node x to keep the first k closest neighbors in the embedding space.

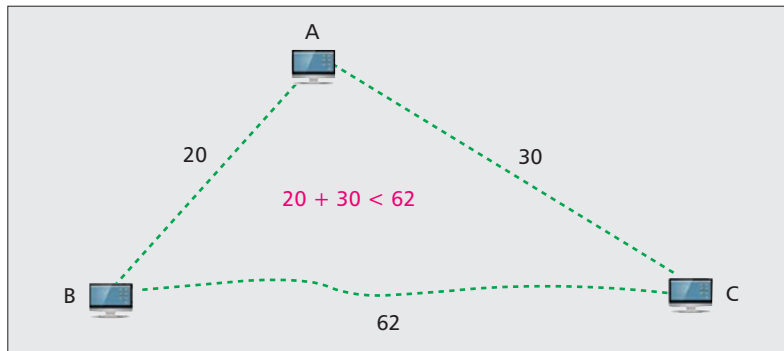


Figure 5. Triangle inequality violations derived from real measurements.

EMERGING CHALLENGES AND OPEN ISSUES

EMERGING CHALLENGES IN NDP

This subsection highlights the important challenges in NDP.

Symmetry: Given any nodes i and j , the distance between them is given by the function $d(i, j)$. The symmetry requires that for i and j , $d(i, j) = d(j, i)$. However, the network distance between any two reachable nodes is not necessarily symmetric due to the network structure and routing policy [1, 16]. Existing research studies [16] have measured the proportion of asymmetric routing: more than 20 percent of the links have borne the asymmetric flow/packet/byte, and more than 14 percent of the flows in the Internet have shown autonomous system asymmetry. These challenge accurate predictions of network distances.

TIV: Most NDP approaches, including [2–4] and others, are based on the embedding of host positions in a finite-dimensional space, commonly the Euclidean coordinate system, where Euclidean distance is used to form the desired estimate. However, in such a system the predicted distance always violates the triangle inequality due to routing policies or path inflation. Existing studies show that TIVs are widespread and persistent [6, 16]. A TIV occurs among a triple of nodes in the Internet when the latencies between them cannot form a valid triangle, and changes with time. Figure 5 illustrates such a scenario derived from real measurements [16]. A TIV represents, through a real network path, that there is a closer route to a host through an intermediate host than the direct route, but cannot hold for metric space. Thus, the TIV inevitably yields inaccurate predictions of network distances.

Consistency: Let $d(i, j)$ and $\bar{d}(i, j)$ denote the measured distance and the predicted distance from node i to node j , respectively. Given any reachable pairwise nodes i and j , and nodes u and v , consistency requires that the path among them $\bar{d}(i, j) > \bar{d}(u, v)$ if and only if $d(i, j) > d(u, v)$, and vice versa. In practice, there is a significant disparity in NDP in terms of the prediction accuracy for different distance ranges due to many factors such as dynamic topology/paths,

Strictly speaking, the network distance does not match the notion of metric space due to the inherent network characteristics such as violations of symmetry and TIVs. Thus, we can infer that any predicted distance obtained from NDP is an approximate network distance.

Approaches	Common challenges					
	Symmetry	TIV	Consistency	Security	Cluster	Dynamics
GNP	Low	Medium	Low	High	No	Medium
PIC	Low	Medium	Low	Low	Low	Low
NPS	Low	Medium	Low	Low	Low	Low
Tarantula	Low	Medium	Low	High	Low	Low
ICS	No	No	Low	Medium	No	Medium
IDES	No	No	Low	High	No	Medium
DMFSGD	No	No	Low	High	No	Low
Phoenix	No	No	Low	High	No	Low
IDMaps	No	No	No	Medium	No	Medium
iPlane	No	No	No	Medium	No	Medium
Netvigator	No	Medium	No	Medium	No	Medium

Table 2. Common challenges affecting network distance predictions.

time-varying traffic, congestion, and metric space [1], thus lowering the validity and scalability of NDP. How to achieve significant consistency in NDP is a common challenge for the current approaches.

Security: It has been shown that NDP methods are rather primitive [1] and cannot defend against all types of attacks including disorder attack, repulsion attack, and isolation attack. Such attacks make NDP very susceptible to the malicious nodes from either inside or outside of the network. The current solution is to test the TIV [8] or eliminate the landmarks that provide significant relative errors [6], which has been considered suboptimized to fulfill the security requirement. Security has become a great challenge in NDP.

Cluster: In order to achieve comparable prediction quality, the hierarchical NDP approaches divided the networks into several clusters in a distributed way. In such approaches, each host keeps different sets of neighbors and coordinates in different layers so that it can predict the intra-cluster distance and inter-cluster distance based on the local coordinates and global coordinates, respectively. Two significant issues of these approaches are that there are still no consensus on how to divide the networks into clusters and how many clusters the networks should be divided into [6, 8], which also challenges the predicted quality.

Dynamic: The NDP schemes should be resilient against the network dynamics mainly due to host dynamics and time-varying traffic, including host failures/joining, temporary network partitioning, and traffic change/overload. In particular, there are the distance predicted schemes, such as GNP

and IDES, suffering from landmark failures and overloading, and furthermore, all the predicted distances that have originated from these schemes are partially decided by direct measurements. However, such measurement results fluctuate frequently with time because of time-varying traffic, which has a slight impact on the direct measurements. Such dynamics prevent the system from building models of network distances and potentially degrade the predicted accuracy of NDP.

Table 2 summarizes how the main approaches proposed previously are influenced by each of these challenges.

OPEN ISSUES AND FUTURE DIRECTIONS IN NDP

In this subsection, we list some potential issues and future directions for further research on NDP.

Embedding Models: Strictly speaking, the network distance does not match the notion of metric space due to the inherent network characteristics such as violations of symmetry and TIVs. Thus, we can infer that any predicted distance obtained from NDP is an approximate network distance. Currently, there are many embedding models such as Euclidean space, hyperbolic space, and hybrid space used to predict network distance, while these models cannot very accurately draw out the characteristics of the real networks. For instance, the widespread TIVs in NDP cannot represent networks with complex routing policies such as suboptimal routing or asymmetric routing, and thus yield inaccurate predictions of network distances. In order to address such issues, we should focus on how to design embedding models imposing the network features.

QoS-Based Network Distance Prediction:

The NDP designed at the beginning mainly aims to provide predicted distances without performing direct measurements for end users to greatly benefit from intelligent path selection based on network performance. However, in most cases, they cannot provide quality of service (QoS)-based guaranteed distances except iPlane. For instance, the predicted distances fluctuate frequently and are not convincing in some ways. There is a performance gap between them and the ideal NDP. We argue that it is desirable to design QoS-based NDP approaches in the future. To achieve this goal, we can focus on time-varying measurements by exploiting differentiated services code point (DSCP) or type of service (ToS) bits to probe the network distance, and designing QoS-based predicted metrics to estimate the predicted accuracy.

Multi-Metric Network Distance Prediction:

Currently, much related work provides only a limited subset of the metrics of interest, commonly latency between a pair of nodes. In reality, however, latency is just one of many metrics, such as available bandwidth and packet loss rate, that affect the performance of service-oriented applications. Compared to the latency as an additional parameter, the available bandwidth is a concave parameter, and the packet loss rate is a multiplicative parameter. If we simply embed these parameters into Euclidean space like the current approaches, they cause prediction distance to be arbitrarily wrong. Therefore, designing multi-metric embedding models based on different performance metrics is an efficient approach to achieve the desired predicted performance and needs to be further investigated.

Predicted Errors: The NDP inevitably generates the predicted errors caused by various factors such as landmarks failure, embedding metric spaces, and evaluated metrics. The predicted errors directly determine the QoS for end users. In order to provide better service to end users, and bridge the gap between service-oriented applications and large-scale networks, we need to investigate the impact of the predicted errors on service-oriented applications and identify the reasons behind them, then design high-precision NDP. In addition, we should evaluate the predicted performance under the practical measurement platforms such as PlanetLab, DIMES, OneLab, and EmuLab using the multi-metrics, and create a system that leverages both theoretical approaches and actual distance prediction in the network, thus catalyzing the evolution of NDP into a service-oriented architecture.

Security-Aware Network Distance Prediction:

The most recent predicted mechanisms assume that the participating nodes can be trusted. Unfortunately, it has been proven that NDP methods are rather primitive and cannot prevent a variety of attacks such as disorder, repulsion, isolation, and system control attacks, providing a potentially attractive fertile ground for the disruption or collapse of the many applications and overlays that would use these services. To the best of our knowledge, currently there are only a

few simple methods in NDP to defend against malicious behaviors. For example, PIC uses a test based on the triangle inequality to detect malicious nodes, and NPS regards a landmark that provides significant relative errors compared to the reference nodes as a malicious node. Such security mechanisms designed for NDP approaches have shown that they are still susceptible to the intrusions.

Security-aware NDP studies are still in their infancy. How to prevent network attacks determines the QoS of distance prediction. It is desirable to design more security-aware NDP approaches and therefore prevent from various malicious behaviors.

CONCLUSIONS

The network distance prediction has been considered important to improve the performance of service-oriented applications and bridge the gap between end users and large-scale networks, and thus has received increasing attention. In this article, we have investigated the important existing NDP approaches, and categorize the current research work based on different criteria. We provide general information on the behaviors of NDP approaches, and discuss their merits and drawbacks. Finally, we point out potential issues and future directions for further research.

ACKNOWLEDGMENT

This work has been partially supported by the National Basic Research Program of China (No. 2011CB302601 and No. 2012CB315801), the National Natural Science Foundation of China (No. 61402343, No. 61222213 and No. 61170290), the EU FP7 CLIMBER Project (No. PIRSESGA-2012-318939), the Natural Science Foundation of Hubei Province (No. 2013CFB332), and the Jiangsu International Cooperation Program of Science and Technology (No. BZ2013018).

REFERENCES

- [1] B. Donnet, B. Gueye, and M. Kaafar, "A Survey on Network Coordinates Systems, Design, and Security," *IEEE Commun. Surveys and Tutorials*, 2010, vol. 12, no. 4, pp. 488–503.
- [2] T. S. E. Ng and H. Zhang, "Predicting Internet Network Distance with Coordinates-Based Approaches," *Proc. IEEE INFOCOM '02*, New York, NY, June 2002, pp.170–79.
- [3] M. Costa *et al.*, "PIC: Practical Internet Coordinates for Distance Estimation," *Proc. IEEE ICDCS 2004*, Tokyo, Japan, Mar. 2004, pp.178–87.
- [4] Z. Chen *et al.*, "Tarantula: Towards an Accurate Network Coordinate System by Handling Major Portion of TIVs," *Proc. IEEE GLOBECOM '11*, Houston, TX, Dec. 2011, pp. 5–9.
- [5] Y. Chen *et al.*, "Phoenix: A Weight-Based Network Coordinate System Using Matrix Factorization," *IEEE Trans. Network Services and Mgmt.*, 2011, vol. 8, no. 4, pp. 334–47.
- [6] T. S. E. Ng and H. Zhang, "A Network Positioning System for the Internet," *Proc. USENIX ATEC '04*, Boston, MA, 2004, pp.141–54.
- [7] H. Oktay *et al.*, "Distance Estimation for Very Large Networks Using MapReduce and Network Structure Indices," *Proc. WIN '11 Wksp.*, New York, NY, 2011, pp. 1–6.
- [8] H. Lim, J. Hou, and C. Choi, "Constructing an Internet Coordinate System Based on Delay Measurement," *IEEE/ACM Trans. Networking*, 2005, vol. 13, no. 3, pp. 513–25.
- [9] Y. Mao, L. Saul, and J. M. Smith, "IDES: An Internet Distance Estimation Service for Large Networks," *IEEE JSAC*, Dec. 2006, vol. 24, no. 12, pp. 2273–84.

Security-aware NDP studies are still in their infancy. How to prevent network attacks determines the QoS of distance prediction. It is desirable to design more security-aware NDP approaches and therefore prevent from various malicious behaviors.

- [10] Y. Liao *et al.*, "DMFSGD: A Decentralized Matrix Factorization Algorithm for Network Distance Prediction," *IEEE/ACM Trans. Networking*, 2013, vol. 21, no. 5, pp. 1511–24.
- [11] P. Francis *et al.*, "IDMaps: A Global Internet Host Distance Estimation Service," *IEEE/ACM Trans. Networking*, 2001, vol. 9, no. 5, pp. 525–40.
- [12] H. Madhyastha *et al.*, "iPlane: An Information Plane for Distributed Services," *Proc. USENIX OSDI '06*, Seattle, WA, Nov. 2006, pp. 367–80.
- [13] P. Sharma *et al.*, "Estimating Network Proximity and Latency," *ACM SIGCOMM Comp. Commun. Rev.*, 2006, vol. 36, no. 3, pp. 39–50.
- [14] H. Yin *et al.*, "Big Data: Transforming the Design Philosophy of Future Internet," *IEEE Network*, 2014, vol. 28, no. 4, pp. 14–19.
- [15] A. Jain and J. Pasquale, "Internet Distance Prediction Using Node-Pair Geography," *Proc. IEEE NCA '12*, Cambridge, MA, 2012, pp. 71–78.
- [16] Y. Zhu *et al.*, "Taming the Triangle Inequality Violations with Network Coordinate System on Real Internet," *Proc. ACM ReArch 2010*, Philadelphia, USA, 2010, pp. 1–6.

BIOGRAPHIES

HAOJUN HUANG (hhj0704@hotmail.com) is a lecturer within the College of Electronic Information at Wuhan University, and also a postdoctoral researcher in the Research Institute of Information Technology (RIIT) at Tsinghua University. He received his B.S. degree from the School of Computer Science and Technology of Wuhan University of Technology in 2005, and his Ph.D. degree from the School of Communication and Information Engineering from the University of Electronic Science and Technology of China in 2012. He has published 20 papers in international journals and conferences, such as *IEEE Communications Magazine* and *IEEE Network*. His research interests focus on big data, mobile Internet, wireless communication, and ad hoc networks.

HAO YIN (h-yin@tsinghua.edu.cn) is a professor in the Research Institute of Information Technology (RIIT) and National Laboratory for Information Science and Technology (TNList) at Tsinghua University. His research interests span broad aspects of multimedia networks, future networks, and big data-driven network science and engineering. Some of his research results have been widely used in industry and adopted by industry standards. He has published over 100 papers in refereed journals and conferences.

GEYONG MIN (g.min@exeter.ac.uk) is a professor of high-performance computing and networking in the Department of Mathematics and Computer Science within the College of Engineering, Mathematics and Physical Sciences at the University of Exeter, United Kingdom. He received his Ph.D. degree in computing science from the University of Glasgow, United Kingdom, in 2003, and his B.Sc. degree in computer science from Huazhong University of Science and Technology, China, in 1995. His research interests include future Internet, computer networks, wireless communications, multimedia systems, information security, high-performance computing, ubiquitous computing, modeling, and performance engineering.

DAPENG OLIVER WU (wu@ece.ufl.edu) is a professor at the Department of Electrical and Computer Engineering, University of Florida, Gainesville. He received his B.E. in electrical engineering from Huazhong University of Science and Technology in 1990, his M.E. in electrical engineering from Beijing University of Posts and Telecommunications, China, in 1997, and his Ph.D. in electrical and computer engineering from Carnegie Mellon University, Pittsburgh, Pennsylvania, in 2003. His research interests are in the areas of networking, communications, video coding, image processing, computer vision, signal processing, and machine learning.

YULEI WU (Y.L.Wu@exeter.ac.uk) is a lecturer in computer science at the University of Exeter. He received his Ph.D. degree in computing and mathematics and his B.Sc. (First Class Honors) degree in computer science from the University of Bradford, United Kingdom, in 2010 and 2006, respectively. His recent research focuses on future network architecture and protocols, wireless networks and mobile computing, cloud computing, and performance modeling and analysis.

TAO ZUO (zuotao@wust.edu.cn) is an associate professor in the College of Information Science and Engineering, Wuhan University of Science and Technology. He received his Ph.D. degree in electronics and information engineering from Wuhan University in 2011. His research interests include photoelectric tracking, predictive filtering, and robot network control.

KE LI (keli@swjtu.edu.cn) is a lecturer in the School of Information Science and Technology Southwest Jiao Tong University. Her major interests are network optimization, network design, and next generation networks. She has published over 10 papers in refereed journals and conferences.

CAINE: A Context-Aware Information-Centric Network Ecosystem

George Kamel, Ning Wang, Vassilios Vassilakis, Zhili Sun, Pirabakaran Navaratnam, Chonggang Wang, Lijun Dong, and Rahim Tafazolli

ABSTRACT

Information-centric networking (ICN) is an emerging networking paradigm that places *content identifiers* rather than *host identifiers* at the core of the mechanisms and protocols used to deliver content to end users. Such a paradigm allows routers enhanced with content-awareness to play a direct role in the routing and resolution of content requests from users, without any knowledge of the specific locations of hosted content. However, to facilitate good network traffic engineering and satisfactory user QoS, content routers need to exchange advanced network knowledge to assist them with their resolution decisions. In order to maintain the location-independency tenet of ICNs, such knowledge (known as *context information*) needs to be independent of the locations of servers. To this end, we propose CAINE — Context-Aware Information-centric Network Ecosystem — which enables context-based operations to be intrinsically supported by the underlying ICN routing and resolution functions. Our approach has been designed to maintain the location-independence philosophy of ICNs by associating context information directly to content rather than to the physical entities such as servers and network elements in the content ecosystem, while ensuring scalability. Through simulation, we show that based on such location-independent context information, CAINE is able to facilitate traffic engineering in the network, while not posing a significant control signalling burden on the network.

INTRODUCTION

Information-centric networking (ICN) is an emerging networking paradigm that places *content identifiers* rather than *host identifiers* at the core of the mechanisms and protocols used to deliver content to end-users. It essentially makes content *location independent*, thus inherently supporting features such as multicast and user mobility, and potentially enhancing content delivery performance for networks and end-users alike. For this reason, ICN is becoming increas-

ingly appealing to network operators and manufacturers, who are investing in ICN research to bring it to the level of maturity needed for wide-scale commercial rollout. However, this is still a long way off, with many challenges yet to address satisfactorily such as content naming, security, and routing and resolution system scalability [1].

ICNs play a direct role in the routing and resolution of content requests from users, supporting fine-grained content access and distribution, with the ability to handle both complexity and uncertainty. This is done without reliance on any dedicated domain name system (DNS) like entity sitting “outside” the network. To facilitate this role, content routers (CRs), i.e. routers enhanced with content awareness, need to possess advanced network knowledge, known as *context information*, such as content availability, content popularity, content server load, and end-to-end path conditions for content delivery. This context information may then be used by CRs to make routing and resolution decisions that fulfill the traffic engineering requirements of the network and the quality-of-service (QoS) requirements of users.

Although existing ICN schemes support context information dissemination and traffic engineering to some degree, this is typically done by explicitly associating context information to physical elements in the ecosystem, but not intrinsically to the content itself. As such the current practice of using context information in ICN environments still fails to support location independence in terms of context awareness. For example, in the data-oriented network architecture (DONA) [2] and content-ubiquitous resolution and delivery infrastructure for next-generation services (CURLING) [3] approaches in which resolution is carried out by dedicated handlers, only server load information is captured; path load information is not, as it is technically challenging/unsuitable to do so given their centralized architectures. In the named-data networking (NDN) scheme [4], content routers periodically flood user requests for chunks of different contents toward all potential content sources. The routers will then learn the best interface to use for all near-future content

George Kamel, Ning Wang, Vassilios Vassilakis, Zhili Sun, and Pirabakaran Navaratnam are with the University of Surrey.

Chonggang Wang and Lijun Dong are with InterDigital Communications LLC.

Rahim Tafazolli is with University of Surrey.

Our proposed scheme ensures that location-independent context information is exchanged efficiently among content routers, affording more frequent context information exchanges. This, in turn, reduces context information staleness and improves the performance of content delivery for both the network and the users.

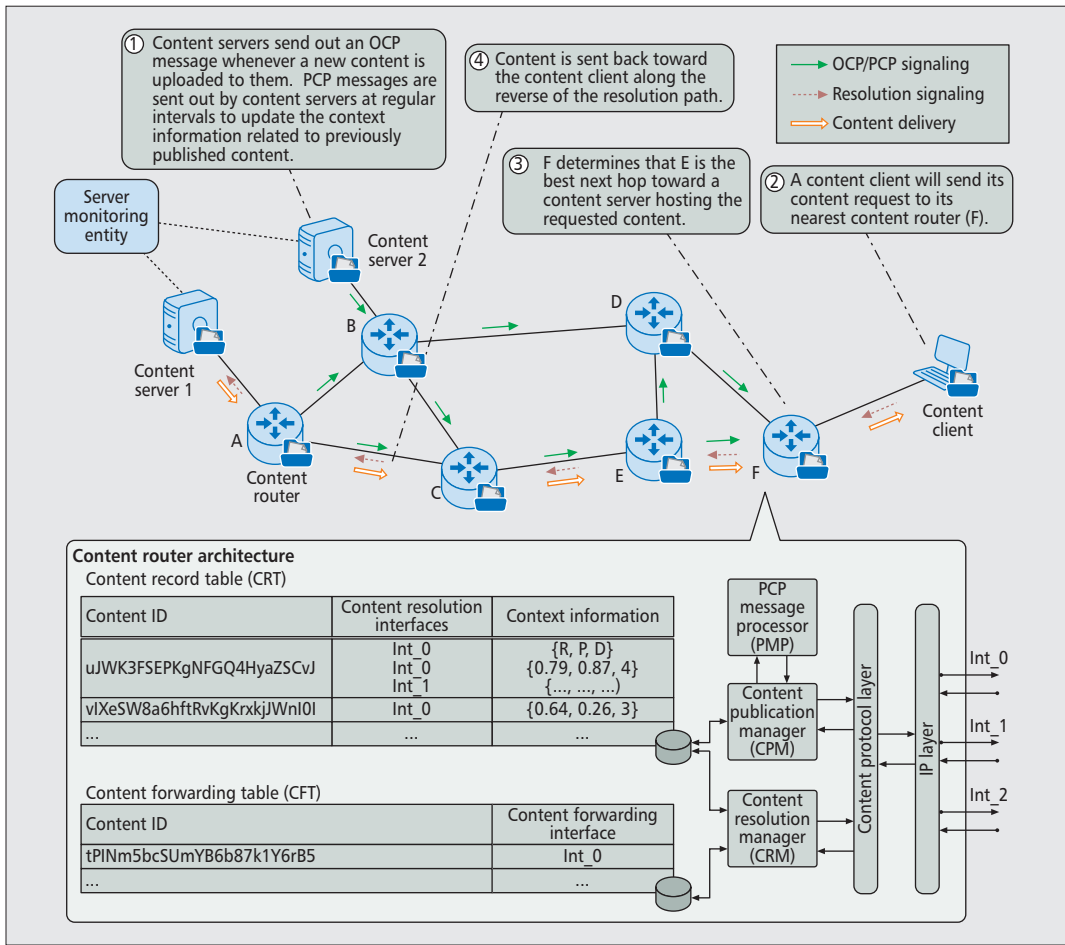


Figure 1. CAINE network and router architecture.

requests based on the time taken to receive the requested chunks. However, such an approach allows only the least delay path to be gauged without really obtaining concrete information about the distance to the servers, their loads, and the bandwidth available along each potential path.

In this article we introduce CAINE — Context-Aware Information-centric Network Ecosystem — which enables advanced context-based operations to be intrinsically supported by the underlying ICN routing and resolution functions. Our approach has been designed to maintain the location-independent philosophy of ICNs by associating context information directly to content rather than to the explicit servers hosting it, or indeed the network elements delivering it. Regardless of this, our proposed scheme ensures that location-independent context information is exchanged efficiently among CRs, affording more frequent context information exchanges. This, in turn, reduces context information staleness and improves the performance of content delivery for both the network and the users.

CAINE FRAMEWORK

In any information-centric network, CRs need to be able to handle content publication and resolution messages, manage content record entries and forwarding states, and correctly handle the transmission of content data itself toward the

clients. In addition to these basic ICN primitives, CAINE also needs to be able to handle the exchange and processing of context information which we integrate with content publication.

Figure 1 shows the architecture required to be implemented within each CR. Such an architecture enables the following three main functionalities related to the life cycle of a content.

Content publication (box 1): Content servers send out two types of content publication messages (elaborated upon later): an original-content publication (OCP) message whenever a new content is uploaded to a server, and a pseudo-content publication (PCP) message sent at regular intervals to update the context information related to content information a server has previously published.

Content resolution (boxes 2 and 3): Content clients submit content requests to their first-hop CR. This CR then determines the best next hop toward a content source hosting the requested content. The resolution routing decision is then carried out hop-by-hop at each CR until a content server is reached. The resolution decision made by each CR is based on the context knowledge each CR has accumulated from the PCP messages they have received.

Content delivery (box 4): Content delivery is made along the reverse of the resolution path based on forwarding states created during the associated content resolution phase.

The aim of the PCP is to disseminate knowledge about the current network conditions to the CRs, without revealing explicit information about servers' identities. This information empowers CRs to make decisions about routing content requests to the best available content source.

Each CR contains two main content management tables: a content record table (CRT) and a content forwarding table (CFT). The CRT contains information associating content IDs with one or more network interfaces through which content requests can be routed toward a content node holding the requested content. In addition, each content item may have network or server context information indicating some quality-of-service (QoS) metrics related to reaching the given content from the node at which the CRT resides. In the example given in Fig. 1, three types of context information are shown:

- *D*: the distance (number of CR-hops) from the content router to the content source.
- *P*: the bandwidth available on the path between the content router and the content source, which is given by the minimum link bandwidth along the path.
- *R*: the resources available at the content server, e.g, the number of additional connections it can support.

The presence of context information in the CRT departs from the approach taken by NDN [4], which specifies in their forwarding information base (FIB) (the equivalent of our CRT), only content ID (specifically, content prefix) and content resolution interface fields. A content publication manager (CPM) interfaces with the CRT to add or modify content records in response to different types of publication messages (expounded upon later) it receives from other content-aware nodes. In addition, a content resolution manager (CRM) also interfaces with the CRT in order to perform look-ups on the next-hop interface to which to forward content requests. The CPM communicates with a PCP message processor (PMP) to perform operations on pseudo-content publication messages and determine the content IDs to which the received messages pertain.

Finally, the CFT (equivalent to NDN's pending interest table (PIT)) contains forwarding states related to ongoing sessions. It maintains associations between content IDs and outgoing next-hop interface(s) through which to forward received content toward the relevant client(s). The CRT will interface with the CFT to install forwarding states in response to content requests it receives.

CAINE is designed to facilitate accurate decision-making during content resolution, so as to ensure that traffic load is well-balanced across the network. Contrary to many existing ICN approaches, we take a more distributed approach in which all content routers (CRs) within a domain are empowered with knowledge to help them make accurate content resolution decisions. We follow a gossip-style approach to content publication and resolution, essentially coupling together the physical signalling routes of content resolution and corresponding content delivery.

The main focus and novelty of our approach lies in the content publication process, which is performed in two stages and which together serves to facilitate that of *context-aware* content resolution. The following sections elaborate on the two content publication processes. In a later section we explain how the information dissemi-

nated during content publication is used to intelligently resolve content requests in such a way that ensures well-balanced network load.

ORIGINAL-CONTENT PUBLICATION

An original content publication (OCP) message is sent by a server whenever a new content is uploaded to or created at it to make CRs in the same domain aware of the *presence* of a new piece of content, as well as the direction toward it. An OCP message is encapsulated within an IP packet, and contains two fields: "Message Type," and "New Content ID." The "Message Type" field simply specifies that this message is an "OCP" message, whereas the "New Content ID" field contains the ID of the new content being published. A server identifier is not sent in the OCP, thus completely decoupling content identifiers from server identities.

When a CR receives an OCP message, it will first confirm receipt by sending an acknowledgement to the previous-hop CR. It will then proceed to create a new content record in the CRT, filling in the content ID and the content resolution interface, which is the interface through which the OCP was received; the context information is not filled until the next PCP message is sent by the server, since the aim of an OCP is to make content routers (CRs) aware of merely the presence of and direction toward the content in the network. To ensure that *all* CRs are made aware of the presence of the new content, a simple dissemination mechanism is employed in which each CR forwards the OCP messages it receives across all of its interfaces, except the one on which the message originally arrived. This allows the CRs to forward requests for content along the correct interface leading to the content source. In the case where multiple sources along multiple network interfaces exist for a particular requested content, the CR must make a decision on the interface to use, i.e. the direction in which to perform content resolution for the incoming request. To facilitate this decision, we propose the use of a special PCP detailed in the following section.

PSEUDO-CONTENT PUBLICATION

Pseudo-content publication messages are sent periodically to update CRs with the latest context information (server and network state) related to the contents hosted at a server. The aim of the PCP is to disseminate knowledge about the current network conditions to the CRs, without revealing explicit information about servers' identities. This information empowers CRs to make decisions about routing content requests to the best available content source.

The PCP message contains three fields: "Message Type," "Context Information," and "Bloom Filter." The "Message Type" simply indicates that the message is a "PCP." The "Context Information" field contains network and server context associated with the content hosted at that server. Specifically, this context information relates to the server resource availability, *R*, path bandwidth, *P*, and the distance, *D*, in CR-hops

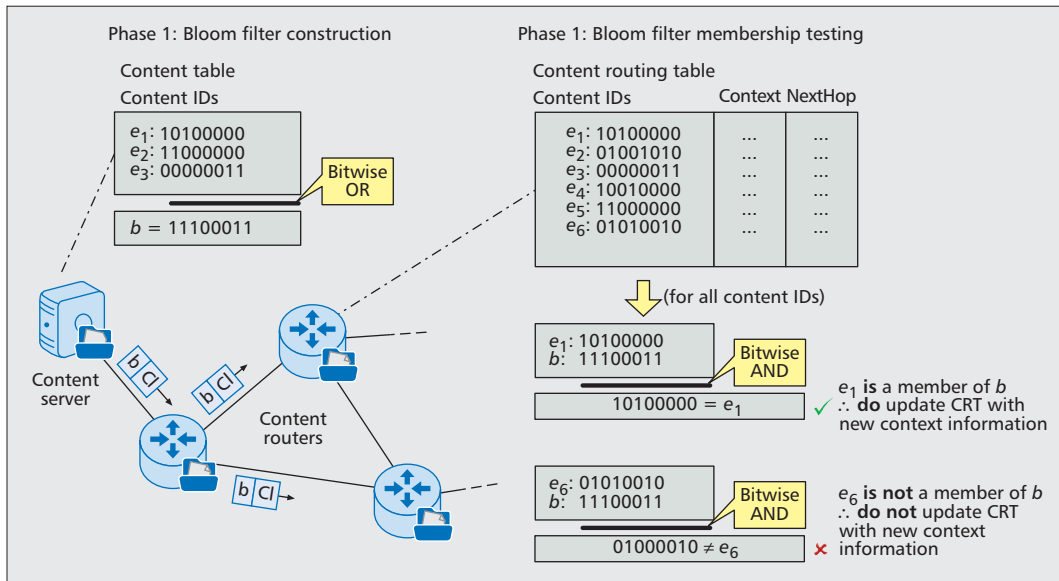


Figure 2. Bloom filter construction and element membership testing.

to the source hosting the requested content.

The last field in the PCP is a Bloom filter [5], which is a probabilistic data structure that allows for a set of elements to be represented by a single space-efficient bit string. Computationally-efficient logic-based set membership queries can then be performed on it to determine if an element is a member of the set it represents. In our case, the set of elements represented by the Bloom filter is the set of IDs of contents hosted at the server that generated the PCP. Set membership queries are performed by CRs using prior knowledge of the content IDs gained through OCP dissemination. Therefore, with both the OCP and PCP messages, the CRs can build up next-hop routing knowledge for each content, together with their associated context, without ever having exposed to the CRs any form of server identifier.

The rest of this section details the way in which the Bloom filters are constructed and elements tested for membership, the strategy for disseminating PCP messages throughout the network, and the way in which the PCP messages are processed at CRs.

BLOOM FILTER CONSTRUCTION

In order to produce the Bloom filter message, each element, e_i , $1 \leq i \leq n$, in the set, S , of contents hosted at a server is hashed k times using k independent hash functions. The resulting Bloom filter bit array representing S is formed of m bits, which is given by [6]

$$m \geq \frac{n \log_2(1/p_f)}{\ln(2)} \quad (1)$$

where p_f is the false-positive probability, i.e. the probability that a test for element membership of element $e_j \notin S$ is positive when it should be negative.

In Fig. 2 we show a simple example of the construction of a Bloom filter by a content server, and the testing of elements for membership by the CRs. In this example, the Bloom filter, B ,

is constructed by performing the bitwise OR operation on three elements, e_1 , e_2 , and e_3 . The content server then places the resulting Bloom filter, B , in a PCP packet together with the related context information, C . Each content router receiving the PCP checks, in turn, each content ID within its CRT for membership in B . This is achieved by performing a *bitwise AND* operation with B . In the example, the result of the bitwise AND operation of e_1 and B yields e_1 , indicating that e_1 is a member element of the Bloom filter. If some other element e_6 were to be tested for membership in B , the result would be negative, since the bitwise AND operation of e_6 with B does not yield e_6 .

From [6] we know that the optimal number of hash functions, k , is given by

$$k = \log_2 \left(\frac{1}{p_f} \right) \quad (2)$$

whereby *optimal* implies that m is minimized subject to meeting the target false-positive probability, p_f .

From Eq. 1 and Eq. 2 we can observe two key characteristics of Bloom filters that make them ideally suited to our application:

- The size of a Bloom filter is independent of the size of the elements, meaning it is possible to use very long content IDs without increasing the size of the Bloom filter. For instance, to create a Bloom filter for 50,000 content IDs, and given that the optimal number of hash functions is used, then to achieve a false-positive probability of 2 percent, approximately 8 bits per element will need to be used, giving a total Bloom filter size of approximately 50 kB. If the size of each content ID is assumed to be 256 bits, then without Bloom filters, conveying information about 50,000 IDs would require approximately 1.5 MB of space.

- The number of hash functions used has a bearing on the computational complexity of the Bloom filter, since the number determines the bits that need to be read to test for membership.

The last field in the PCP is a Bloom filter, which is a probabilistic data structure that allows for a set of elements to be represented by a single space-efficient bit string. Computationally-efficient logic-based set membership queries can then be performed on it to determine if an element is a member of the set it represents.

To minimize the effect of Bloom filter false-positives, the PCP message processing algorithm checks not only for membership of the content ID in the Bloom filter; it also checks for membership of the identifier of the network interface on which the PCP message was received in the given content ID's CRT interface list.

```

1:  $B \leftarrow \text{PCP.BloomFilter}$ 
2:  $C \leftarrow \text{PCP.ContextInfo}$ 
3:  $I \leftarrow \text{IP.GetRcvInterface(PCP)}$ 
4:  $C.D \leftarrow C.D + 1$  {update hop count}
5:  $C.P \leftarrow \min(C.P, I.BW)$  {update path bandwidth}
6: for each row in CRT do
7:    $E \leftarrow \text{row.ContentID}$  {content ID element to hash}
8:    $E_{\text{hashed}} \leftarrow E$  { initialize hashed element}
9:   for each HashFn in B.HashFns do
10:     $E_{\text{hashed}} \leftarrow \text{HashFn}(E)$ 
11:   end for
12:   check  $\leftarrow (B.BitString \mid E_{\text{hashed}})$ 
13:   if (check =  $E_{\text{hashed}}$ ) and
     ( $I \in \text{row.Interfaces}$ ) then
14:     CRT.UpdateContent(row.ID, I, C)
15:   end if
16: end for
17: PCP.ContextInfo  $\leftarrow C$ 

```

Figure 3. Pseudocode for processing a Publish primitive at a CR.

From Eq. 1 we can deduce that the optimal number of hashes grows only linearly with the number of bits per element, B , where B is given by the ratio m/n .

In the case where a false-positive occurs, we anticipate that this will not have an adverse effect on the performance of our proposed mechanism, for reasons that we will explain in the following sections.

PCP DISSEMINATION STRATEGY

Once a PCP message has been constructed by a content server, it needs to be disseminated to other CRs within the network. Since the PCP message is building up a distance-vector-type view of network context such as available path bandwidth and path length, each CR is required to forward the PCP message across all of its interfaces, as in the case of the dissemination of OCP messages. In the process of dissemination, context information on the network side will be added and updated within the PCP at each CR hop. This is done as an “offline” background process such that it does not interfere with the routing and resolution efficiency of CRs. To avoid potential routing loops, PCPs are forwarded using the *split-horizon rule* [7], i.e. CRs forward PCPs along all of its interfaces except the one on which it was originally received.

An important issue to consider is the frequency with which PCP messages are disseminated, as stale context information reduces the accuracy of content resolution decisions. Since PCP messages relate not only server resources but also network path bandwidth, the frequency with which servers send out PCP messages also needs to take into account the dynamics of link bandwidth availabilities. However, from experiments carried out based on real network traffic traces as well as real traces of user requests to YouTube servers [8], we found there is a strong correlation between the available server resources and the link conditions of the network. Therefore, such strong corre-

lation between the two metrics validates a mechanism to determine the most suitable update frequency based purely on the server resource availability.

Since our aim is not necessarily to achieve perfectly load-balanced servers, but rather to avoid servers from becoming overloaded, we propose a PCP dissemination frequency based on a non-linear set of triggers. With such triggers, PCP messages would be sent more often when server load is high, and less often when it is low. At times when the server load fluctuates little, hence not crossing any triggers, a PCP message may still be disseminated after a given time has elapsed from when the last one was sent. The purpose of these time-driven PCPs is to ensure that new contents recently published by a server using an OCP have some context attached to them.

PCP PROCESSING STRATEGY

When a CR receives a PCP, it updates the context information contained within the PCP, determines which entries in its CRT to update with the updated context information, and then forwards the PCP to the next-hop CR(s). The full PCP message processing algorithm is shown in Fig. 3.

In order to determine which entries in its CRT to update with the new context information, the CR checks in turn each content ID, e_i , $1 \leq i \leq N_{\text{CRT}}$, within its CRT for membership within the Bloom filter, B . This is done in the manner described previously and as illustrated in Fig. 2. If e_i is found to be a member of B and the ID of the network interface on which the PCP message was received matches one of the interface field entries in the CRT for the given content ID, then the context information related to that content ID and interface is updated with the new context information.

To minimize the effect of Bloom filter false-positives, the PCP message processing algorithm checks not only for membership of the content ID in the Bloom filter; it also checks for membership of the identifier of the network interface on which the PCP message was received in the given content ID's CRT interface list. Thus, if a false-positive does occur, it will affect only the *accuracy* of the context information of that particular content ID's interface information. As a result, the content resolution process (discussed in the next section) will still be able to route content requests toward *one* of the available sources of the requested content, although the routing decision may be suboptimal.

Once a CR has extracted the relevant information from the PCP message, the context information contained therein is updated before being forwarded to the next-hop CR(s). For example, the distance metric is incremented by one, whereas the path bandwidth is updated with the measured bandwidth of the link through which the PCP message was received, and is updated only if the locally measured link bandwidth is less than the overall path bandwidth.¹ The server resource availability information is not changed along the PCP dissemination paths.

¹ The bandwidth of a link is reverse estimated by a CR by passively measuring the rate of data it receives and subtracting this from the total supported data rate.

CONTENT RESOLUTION

When a content router receives a content request and there is more than one interface in its CRT attached to that content ID, the CR must make a decision about the “best” interface to use to forward the content request. To make this decision, the CR prioritizes the various context metrics, and performs tie-breaker tests on each priority metric in turn, as illustrated in Fig. 4. If the values of a given metric are equal, or lie within a certain pre-determined range from each other, those interfaces qualify for the next lower-priority round of selection. For this particular work, server resources, R , is given the highest priority, followed by path bandwidth, P , and then finally the distance, D . In order to allow all context metrics to be given consideration, we use non-linear ranges, such that higher QoS metric values have looser ranges, and vice versa.

In the example shown in Fig. 4, for a given content ID there are five interfaces from which to choose. In the first priority metric selection, the algorithm qualifies to the next round of selection the three interfaces toward servers having the highest resource availability and that are within the non-linear bounds. In the second priority, the algorithm qualifies out of the three interfaces the two having the highest path bandwidth to the content source and that are within the non-linear bounds. Finally, for the third priority metric, out of the two interfaces that qualified to the third stage of selection, the one with the least number of hops toward the content server is selected as the *best* interface.

Once a best next-hop along which to forward the content request is determined, the CR will install a forwarding state in the CFT to indicate the interface through which to send the content data toward the content client.

EVALUATION

The performance of CAINE was evaluated by means of computer simulation using the GÉANT topology as the reference topology [9], a pan-European point-of-presence (PoP) data network for the research and education community. Specifically, we used the topology of the year 2004, which consists of 23 PoP nodes, and 74 high capacity interconnecting network links. This choice of topology was influenced primarily by the availability of real path congestion measurements extending over four months conducted within the TOolbox for Traffic Engineering Methods (TOTEM) project [10], hence ensuring that the modelling is realistic. Five content servers with equal connection capacity were deployed at various locations within the network, each hosting 10,000 contents selected randomly from a pool of 25,000 possible contents, except for the 1,000 most popular contents, which were hosted by *all* servers.

All of the contents hosted were videos, the rates and durations of which were made to follow the measured trends reported by Cheng *et al.* in [11]. User request patterns were synthetically generated for a 24-hour period based on the characteristics of real YouTube request traces collected by Zink *et al.* [8]. This synthesis

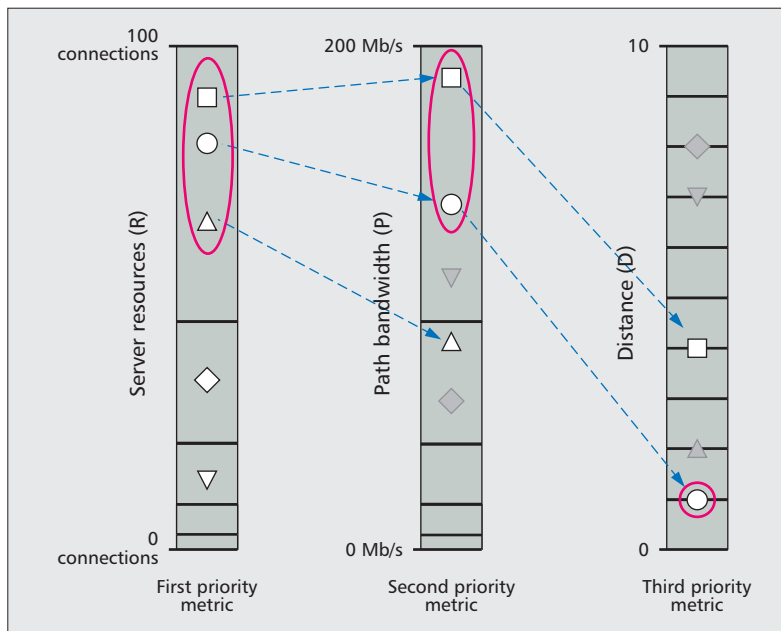


Figure 4. Example of metric prioritization for content resolution.

was achieved by observing the mean request rate, λ , in each 15-minute interval, Δ_t , and generating Poisson-distributed sets of content request requests with different factors, β , of the *set* of mean content request rates. The popularity of the contents followed a Zipf distribution with a shape parameter, α , of 1.0. Finally, the performance of the CAINE scheme is compared against a context-unaware scheme, i.e. one that is agnostic to the server resource and path bandwidth availabilities, and that forwards all content requests along the shortest path toward the nearest source that holds the requested object.

Figure 5 shows some key traffic engineering performance metrics of CAINE based on our simulation of content request events over a 24-hour period. To illustrate CAINE’s load-balancing efficacy, we captured in Fig. 5a the proportion of time for which the most heavily loaded content server was saturated during the 24-hour period. At low content request rate factors, β , content servers never reached saturation, but with increasing β , content servers suffer from significantly longer duration of saturation in the context-unaware case than in the context-aware case. At $\beta = 2.5$, the most heavily loaded server was saturated over 29 percent of the time under the context-unaware scheme, and only 9 percent of the time under the context-aware scheme. We can also glean some insight from the degree of server load imbalance across the five content servers, which we define to be the mean statistical range of server loads across the 24-hour period that was simulated. It was found that under low β , the context-aware approach had a load imbalance of 13 percent, which is marginally greater than that of the context-unaware approach. This is due to the fact that at low content request rates, server utilizations are relatively low, resulting in reduced frequency of PCP updates, and hence reduced freshness of context information and less optimal resolution decisions. However, with increasing values of β , as

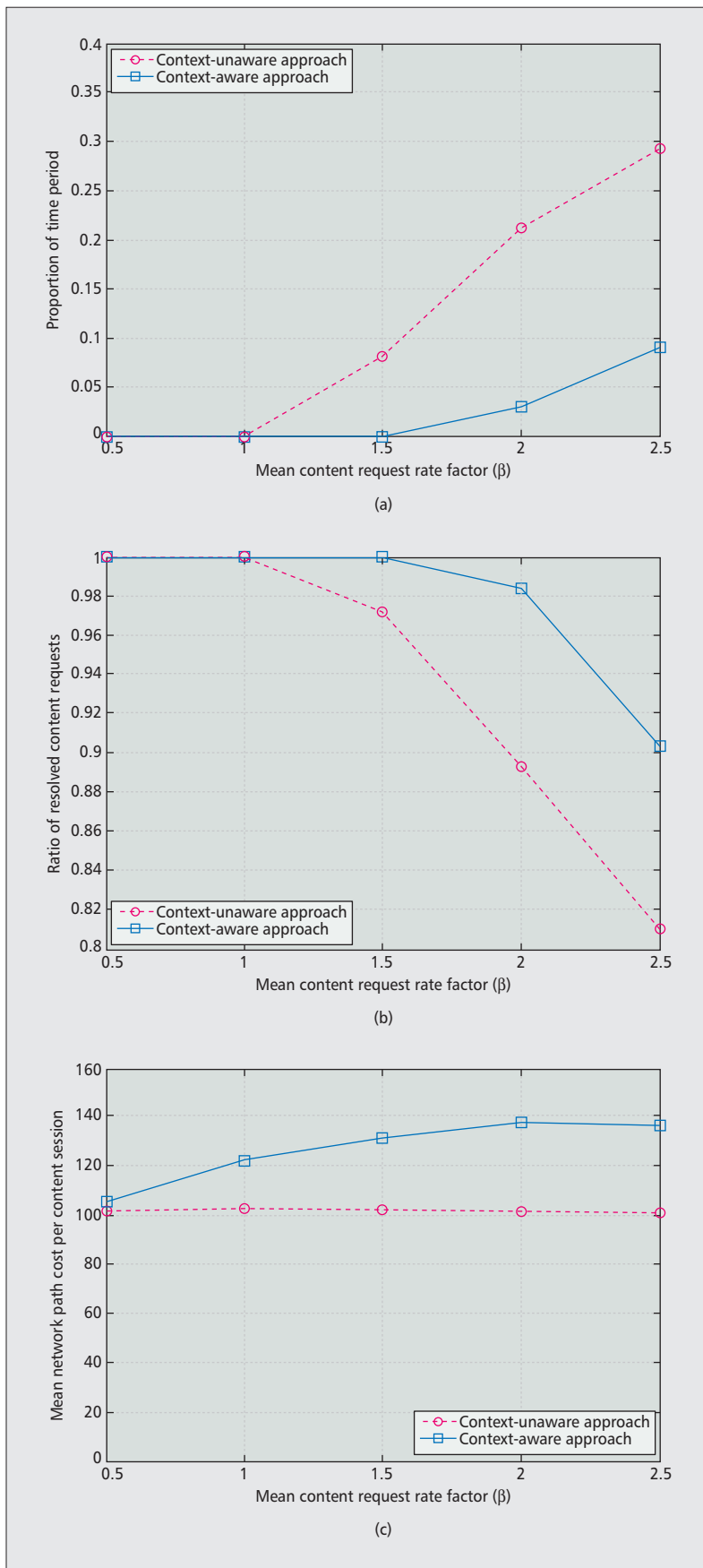


Figure 5. Simulation results: a) proportion of the 24-hour time period in which the most heavily loaded content server is saturated; b) average ratio of successfully resolved content requests across the 24-hour period; c) mean network path cost per admitted content session.

server utilizations reach more critical levels and PCP exchanges become more frequent, the load imbalance of the context-aware approach tended toward 10 percent, whereas that of the context-unaware approach increased to as much as 33 percent.

One of the benefits of CAINE's load balancing efficacy is its optimal network utilization, reflected by its higher ratio of successfully resolved content requests in comparison to the context-unaware approach, as shown in Fig. 5b. At high β , the resolved content request ratio of the context-aware scheme was as much as 11 percent more than the context-unaware scheme, due to CAINE's ability to use alternative paths to the shortest path, as well as alternative servers with more resources available. However, the tradeoff of CAINE's increased network utilization is higher network path cost per content session, which we computed based on the sum of the costs of the individual links traversed along a content delivery path. The cost of using a link is given by a piecewise function defined in [12], which is dependent on the link's utilization, and increases exponentially with it. Therefore, it essentially captures together both the link utilization and path length metrics. From Fig. 5c it can be seen that at the lowest simulated β , the network path cost of the context-aware scheme is only marginally greater than the context-unaware scheme, but can be as much as 35 percent more with higher β .

With regular PCP exchanges being at the heart of CAINE's operation, we looked at the mean frequency with which PCP messages are exchanged by servers at different values of β . The PCP messaging frequency was observed to increase exponentially with increasing β , from PCPs being exchanged every 14 minutes at $\beta = 0.5$, to every 20 seconds at $\beta = 2.5$. However, such a high rate of messages occurs only when the content servers approach saturation. Furthermore, such messages are approximately only 10 KB in size to advertise 10,000 contents, and are processed efficiently at CRs as a background ("offline") process.

CONCLUSION

To facilitate efficient use of network resources in ICNs, we have proposed a novel context-aware ICN-based scheme called CAINE in which location-independent context information is efficiently published to the network to facilitate better decisions during content resolution. Bloom filters are periodically constructed to efficiently convey the IDs of contents hosted at servers, to which up-to-date context information is then attached and disseminated to the content routers. Such a mechanism avoids the need to reveal to the network explicit condition information associated with physical elements within the ecosystem, thereby upholding the key ICN principle of location-independence.

Through simulation, we have shown that CAINE can achieve optimized network utilization and effective load-balancing between servers, particularly when their utilizations are at critical levels, albeit at the cost of increased network path cost. Furthermore, the use of Bloom

filters ensures that such messages do not pose a significant overhead to the network, from the perspectives of both transmission and content router processing.

REFERENCES

- [1] B. Ahlgren *et al.*, "A Survey of Information-Centric Networking," *IEEE Commun. Mag.*, vol. 50, no. 7, 2012, pp. 26–36.
- [2] T. Koponen *et al.*, "A Data-Oriented (and Beyond) Network Architecture," *SIGCOMM Comput. Commun. Rev.*, vol. 37, no. 4, Aug. 2007, pp. 181–92.
- [3] G. Pavlou *et al.*, "Internet-Scale Content Mediation in Information-Centric Networks," *Annals of Telecommun.*, vol. 68, no. 3–4, 2013, pp. 167–77.
- [4] V. Jacobson *et al.*, "Networking Named Content," *Proc. 5th Int'l. Conf. Emerging Networking Experiments and Technologies*, ser. CoNEXT '09, New York, NY, USA: ACM, 2009, pp. 1–12.
- [5] B. H. Bloom, "Space/Time Trade-Offs in Hash Coding with Allowable Errors," *Commun. ACM*, vol. 13, no. 7, July 1970, pp. 422–26.
- [6] A. Broder and M. Mitzenmacher, "Network Applications of Bloom Filters: A Survey," *Internet Mathematics*, 2002, pp. 636–46.
- [7] C. Hedrick, "Routing Information Protocol," RFC 1058 (Historic), Internet Engineering Task Force, Jun. 1988, updated by RFCs 1388, 1723, available: <http://www.ietf.org/rfc/rfc1058.txt>.
- [8] M. Zink *et al.*, "Characteristics of YouTube Network Traffic at a Campus Network — Measurements, Models, and Implications," *Comp. Net.*, vol. 53, no. 4, Mar. 2009, pp. 501–14.
- [9] GEANT project home, available: <http://www.geant.net/>.
- [10] S. Uhlig *et al.*, "Providing Public Intradomain Traffic Matrices to the Research Community," *SIGCOMM Comp. Commun. Rev.*, vol. 36, no. 1, Jan. 2006, pp. 83–86.
- [11] X. Cheng, J. Liu, and C. Dale, "Understanding the Characteristics of Internet Short Video Sharing: A YouTube-based Measurement Study," *IEEE Trans. Multimedia*, vol. 15, no. 5, Aug. 2013, pp. 1184–94.
- [12] B. Fortz and M. Thorup, "Optimizing OSPF/IS-IS Weights in a Changing World," *IEEE JSAC*, vol. 20, no. 4, Sept. 2006, pp. 756–67.

BIOGRAPHIES

GEORGE KAMEL (g.kamel@surrey.ac.uk) is a research fellow at the Institute for Communication Systems (ICS), University of Surrey, UK. He received his M.Eng. (honours) in telecommunications engineering and his Ph.D. in mobile communications from King's College London, UK in 2005 and 2010, respectively. He is currently involved in research efforts at the 5G Innovation Centre (5GIC), University of Surrey, UK, and has been involved in a number of EU, EPSRC, and industry-funded projects. His research interests include information-centric networking, context-aware network management, and cellular traffic offloading.

NING WANG (n.wang@surrey.ac.uk) is a reader at the Institute for Communication Systems (ICS), University of Surrey, UK. He received his B.Eng. (honours) from Changchun University of Science and Technology, P.R. China in 1996, his M.Eng. from Nanyang Technological University, Singapore in 2000, and his Ph.D. from the University of Surrey, UK in 2004. He currently leads the work area on content, user and network context at the 5G Innovation Centre (5GIC), University of Surrey, UK. His research interests include information-centric networking, context-aware network management, and mobile content delivery.

VASSILIOS VASSILAKIS (vv274@cl.cam.ac.uk) is a research associate at the Computer Laboratory, University of Cambridge, UK. He received his Ph.D. in electrical and computer engineering from the University of Patras, Greece in 2011. From 2011 to 2013 he was with the Network Convergence Laboratory, University of Essex, UK where he conducted research on information-centric networking. From 2013 to 2015 he was with the Institute for Communication Systems (ICS), University of Surrey, UK, where he conducted research on 5G wireless networks. He has been involved in a number of EU and industry-funded R&D projects. His research interests are in the areas of future Internet technologies and next-generation wireless networks.

ZHLI SUN (z.sun@surrey.ac.uk) has been with the University of Surrey, UK, since 1993, and is currently a professor at the university's Institute for Communication Systems (ICS). He received his B.Sc. in mathematics from Nanjing University, China in 1982, and his Ph.D. in computer science from Lancaster University, UK. He worked as a postdoctoral research fellow at Queen Mary, University of London, UK. He has been principle investigator in many projects funded by the EU, ESA, EPSRC, and industry, and has published more than 195 papers and three books. His research interests include wireless sensor networks, mobile ad hoc networks, satellite networks, mobile operating systems, future Internet protocols, and security.

PIRABAKARAN NAVARATNAM (piraba.navaratnam@rspb.co.uk) is an experienced research engineer with a background in mobile communication technologies and Internet of Things. He received his B.Sc.Eng. in electronic and telecommunications engineering from the University of Moratuwa, Sri Lanka, and his Ph.D. in mobile communications from the University of Surrey, UK. He worked as a research fellow at the Institute for Communication Systems (ICS), University of Surrey, UK. He is currently working for RSSB, UK on European research and innovation programs. His research interests include machine-to-machine communication, resource management, information management systems, and future Internet architecture design.

CHONGGANG WANG (chonggang.wang@interdigital.com) is currently a member of technical staff at InterDigital Communications, USA, where his focus is on Internet of Things (IoT) R&D activities including technology development and standardization. He received his Ph.D. from Beijing University of Posts and Telecommunications (BUPT), China in 2002. He is the founding Editor-in-Chief of *IEEE Internet of Things Journal*, and is on the editorial board for several journals including *IEEE Transactions on Big Data and IEEE Access*. He is an IEEE ComSoc Distinguished Lecturer (2015–2016). His research interests include IoT, mobile communication and computing, and big data management and analytics.

LIJUN DONG (lijun.dong@interdigital.com) is currently a staff engineer at InterDigital Communications, USA. She received her Ph.D. in electrical and computer engineering from Rutgers University, USA in 2011. Her research interests include machine-to-machine communication, Internet of Things, and information-centric networking.

RAHIM TAFAZOLLI (r.tafazolli@surrey.ac.uk) is a professor and the Director of the Institute for Communication Systems (ICS) and the 5G Innovation Centre (5GIC), University of Surrey, UK. He has published more than 500 research papers in refereed journals, international conferences, and as an invited speaker. He is the editor of two book volumes published by Wiley in 2004 and 2006, entitled *Technologies for the Wireless Future*. In April 2011 he was appointed a Fellow of the Wireless World Research Forum (WWRF), in recognition of his personal contribution to the wireless world.

Through simulation, we have shown that CAINE can achieve optimized network utilization and effective load-balancing between servers, particularly when their utilizations are at critical levels, albeit at the cost of increased network path cost.

On the Use of Radio Environment Maps for Interference Management in Heterogeneous Networks

Jordi Perez-Romero, Andreas Zalonis, Lila Boukhatem, Adrian Kliks, Katerina Koutlia, Nikos Dimitriou, and Reben Kurda

ABSTRACT

This article addresses the use of REMs to support interference management optimization in heterogeneous networks composed of cells of different sizes and including both cellular and non-cellular (e.g. WiFi) technologies. After presenting a general architecture for including REM databases in different network entities, the article analyzes the achievable benefits in relation to specific interference management techniques, including a discussion on practical considerations such as information exchange requirements, REM ownership, and security aspects. Finally, several research directions derived from the proposed framework are identified.

INTRODUCTION

In recent years there has been an exponential growth in the demand for mobile broadband services associated with the massive penetration of mobile devices and the proliferation of bandwidth-intensive applications. This trend is expected to continue to increase in the future with applications involving high definition video, virtual reality, and so on. For the provisioning of such demanding services, the classical cellular network concept is being shifted toward heterogeneous networks (HetNets) [1] that combine large macrocells with smaller cells of different sizes. HetNets are expected to provide high capacities in densely populated areas and enhance coverage at specific locations (e.g. in indoor environments). They may also involve other access technologies, such as WiFi, that can offload traffic from the cellular network. Although HetNets are already a reality in current fourth generation (4G) systems with the use of microcells, picocells, and femtocells, they are expected to also be fundamental in future fifth generation (5G) systems that envisage extreme densification of small cells in certain scenarios.

The widespread introduction of small cells requires enhanced inter-cell interference coordi-

nation (eICIC) methods for controlling and mitigating interference whenever the same frequency is shared by different cells. Various techniques have been developed and can be classified as frequency-domain techniques, time-domain techniques, and power control techniques [2]. Frequency-domain techniques usually assign different frequency resources to users in the cells that can potentially cause interference [3]. In time-domain techniques, users suffering from interference are assigned resources in specific time periods where the interference is suppressed. Power control techniques adjust the transmitted power to reduce the interference generated to the victim users [4]. In addition, a proper user-to-cell association is essential to ensure that users are connected to the most convenient cell, and correspondingly, less interference is generated/received to/from the other cells. Moreover, when considering cellular and WiFi networks (or other non-3GPP networks), it is also possible to reduce the interference in the cellular network by offloading traffic to the WiFi network. This usually relies on the application of access network discovery and selection function (ANDSF), and of the solutions proposed for the so called hotspot 2.0 [5].

Given the randomness associated with propagation effects, user mobility, and traffic generation, the development of optimized eICIC techniques requires proper knowledge about the environment wherein HetNets are deployed. In this direction, the term radio environment map (REM) refers to a database that dynamically stores information about the environment wherein a cognitive radio system operates [6]. It includes information about propagation conditions, locations of active transmitters in the area, traffic density, and so on. This information can be exploited to optimize wireless networks, as in [7], where different applicability areas of the REM concept were identified.

Along these lines, this article focuses on the use of REMs to support the optimization of eICIC for HetNets. In contrast to previous works, which have either identified different

Jordi Perez-Romero and Katerina Koutlia are with Universitat Politècnica de Catalunya.

Andreas Zalonis is with National Kapodistrian University of Athens.

Lila Boukhatem is with University of Paris-Sud 11.

Adrian Kliks is with Poznan University of Technology.

Nikos Dimitriou is with National Center of Scientific Research "Demokritos."

Reben Kurda is with University of Koya.

The coordination between the global REMs at MME and HeNB GW can be achieved through the S1 interface. Similarly, the possibility of having only the REM at the MME controlling both the eNBs and the HeNBs can also be considered, although this may not be efficient in the case of a very large deployment of HeNBs.

architecture of Fig. 1. One option would be a totally distributed architecture with only local REMs located at each eNB/HeNB and no global REMs. Coordination between REM instances to exchange information about neighboring cells could be achieved through the X2 interface. This solution may be efficient from the REM storage and management perspective because each REM only has to account for its local area. However, signaling associated with coordinating the different REM instances has to be considered as a function of the required information update rate, the amount of exchanged information, the computation and processing cost at each cell, and so on.

Another option arising from Fig. 1 would be a fully centralized approach with only global REMs at the MME and HeNB GW. Each REM will contain information about all of the cells in the controlled area. This will facilitate coordination but may involve large complexities and storage requirements if the number of cells is high. Moreover, REM-based decisions may be executed at a lower rate than in the distributed or layered case because of the latencies associated with accessing the REM. The coordination between the global REMs at MME and HeNB GW can be achieved through the S1 interface. Similarly, the possibility of having only the REM at the MME controlling both the eNBs and the HeNBs can also be considered, although this may not be efficient in the case of a very large deployment of HeNBs.

REM-BASED INTERFERENCE COORDINATION TECHNIQUES IN HETNETS

This section analyses the use of REMs in relation to specific eICIC techniques, illustrates the benefits, and discusses the practical and architectural implications.

REM-BASED eICIC TECHNIQUES

The use of REM as a support tool in interference management in HetNets is applicable to a variety of different categories of eICIC techniques.

Power control techniques: These techniques adjust the transmit power of certain base stations of the network to reduce the generated interference. In a HetNet topology, with eNBs and HeNBs that serve a closed subscriber group of users, a critical challenge is the interference from an HeNB to nearby co-channel macrocell user equipments (MUEs). In this case, the transmit power of the HeNBs should be adjusted to avoid or reduce interference to the victim MUE. In [11], a baseline approach was presented, wherein each HeNB autonomously adjusts its transmit power based on its own received power measurements from the eNB. The introduction of a local REM in HeNB can enhance the effectiveness of this approach using the radio propagation characteristics of the surrounding area and the location of the neighboring MUEs, HeNBs, and eNBs. This is demonstrated in the REM-based autonomous HeNB power control

(RAHPC) technique in which the HeNB uses the local REM to detect and locate the victim MUE and then it adjusts its transmission power to maintain a predefined signal-to-interference-and-noise ratio (SINR) target for the MUE [12]. Another proposal is REM-based macrocell-assisted power control (RMAPC), where the eNBs support the power adjustment of HeNBs by considering the contribution of each HeNB to the total interference and the impact on the outage of both MUEs and HeNB user equipments (HUEs) [13]. In this case, the HeNBs' local parameters are stored in the REM, and the eNB can use them to achieve the globally coordinated power adjustment.

Frequency domain techniques: In this category, the REM can support the optimal selection of sub-bands to be used in the macrocells and small cells, therein targeting the minimization of inter-cell interference. One approach is the Gibbs sampler-based technique originally proposed in [14] in the context of macrocell scenarios. Here, it is extended to a HetNet scenario by also considering its implementation based on the REM concept. In this technique, which is denoted in the following as REM-based frequency optimization (RFO), small cells use a single sub-band, while eNBs use two different sub-bands, one for the users located in the inner part of the cell and the other for the users located in the outer part of the cell. Then, at random instants defined by an exponential timer, each cell modifies the used sub-bands following a Gibbs-Boltzmann distribution that selects with higher probability those sub-bands where the cell receives and generates less interference to their neighbor cells. This is performed iteratively so that the system progressively reduces the total inter-cell interference. The estimation of the received and generated interference is performed based on the propagation losses between users and neighbor cells that are obtained from the REM.

Time-domain techniques: REM-related information can also be helpful in developing an optimal configuration of the muting periods of the macrocells, that is, almost blank subframes (ABS) [2], to enable interference-free small cell transmission. REM will help in the identification of small-cell users that are more sensitive to macrocell interference and in deciding how many ABSs are needed.

WiFi offloading: Efficient data offloading from the cellular network via the WiFi network (or, in general, via other non-3GPP networks) allows a decrease in the HeNB/eNB load, and consequently, it can simplify interference management in the cellular network [15]. REM can be considered as a technical enabler for this offloading because it can contain information about the detailed locations of available WiFi networks and their characteristics.

ARCHITECTURAL CONSIDERATIONS

The selection of an eICIC technique impacts the choices derived from the architecture of Fig. 1 concerning the use of the local/global REMs or the type of stored information. The architectural considerations of the techniques described in the previous sub-section are presented in the follow-

Category of parameters		REM-based Autonomous HeNB Power Control (RAHPC)	REM-based Macrocell-Assisted Power Control (RMAPC)	REM-based Frequency Optimization (RFO)	WiFi offloading
Local REM	Locations of nodes	HeNB, UE locations	UE locations		AP location, UE location
	Radio-related information	Propagation losses between network elements	Macro and HeNB UE quality indicators, signal strength of victim MUEs, HUEs in outage	Propagation losses between users and cells	Received power level from APs and cells
	Transmit power		Transmit power per Resource Block in each HeNB	Transmit power in each sub-band	Transmit power of different nodes
	List of neighbor cells		Set of interfering HeNBs for each macro UE victim	Set of interfering cells for each cell	Set of neighbor APs
Global REM	QoS metrics	SINR target	SINR target		
	Locations of nodes	eNB locations			eNB locations
	WiFi-related information				Occupancy of the WiFi channels, utilization level of IP connections, type of AP (fee-based access or free access, public or private)

Table 1. Categorization of the parameters in the Local/global REM databases and how they are used in the different strategies.

ing, and Table 1 lists the specific REM information used in each case.

Local REM information: The more dynamic parameters of the radio environment and the information that only affects a reduced number of nodes will preferably be stored in the local REMs because this would facilitate the REM updates. As shown in the examples in Table 1, this mainly includes radio-propagation-related information, such as propagation losses, signal strengths, and the locations of certain nodes such as mobile terminals or HeNBs.

Global REM information: Usually, global REMs will store the less dynamic parameters or the parameters that may affect a high number of network nodes. As observed in Table 1, the stored information includes certain quality of service (QoS) metrics, the positions of eNBs, and information related to available WiFi access networks such as the ownership (e.g. private or public, fee-based access or free access, and with or without authorization), the quality of the IP addressing options, and the parameters of the available backhauling options for each WiFi.

Mapping of REM entities: In all of the strategies considered here, the local REM of each cell includes the REM manager and REM SA, and MCDs will be the mobile terminals and the cells, whose measurements will be used to build the REM data.

BENEFITS

The use of a REM in the above mentioned techniques improves the HetNet performance in terms of various metrics.

Capacity and throughput improvement: The information stored in the REMs helps increase the cell capacity and/or user throughput. To illustrate this, two examples are discussed. First,

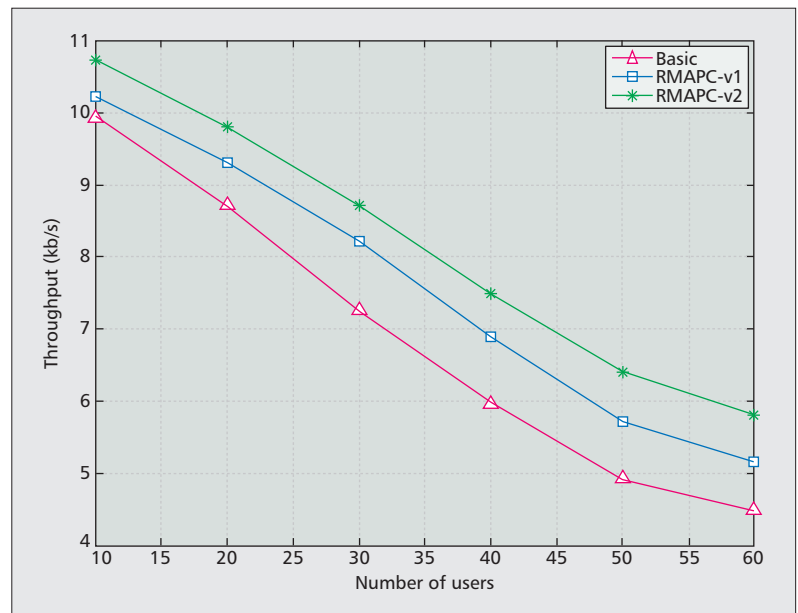


Figure 2. Throughput gain for macrocell users for constant bit rate flows with the RMAPC technique.

Fig. 2 depicts the throughput gains that are obtained for MUEs with the RMAPC strategy in relation to a baseline solution where no REM-based power control is applied. Two versions of the RMAPC mechanism are considered, which differ in their selection strategy of the set of HeNBs used to execute the power adjustment and the amount of power to reduce. The first version is more aware of HUEs' performance degradation, and the second version prioritizes MUEs. Gains of 18 percent for MUEs are observed at the cost of a slight degradation of

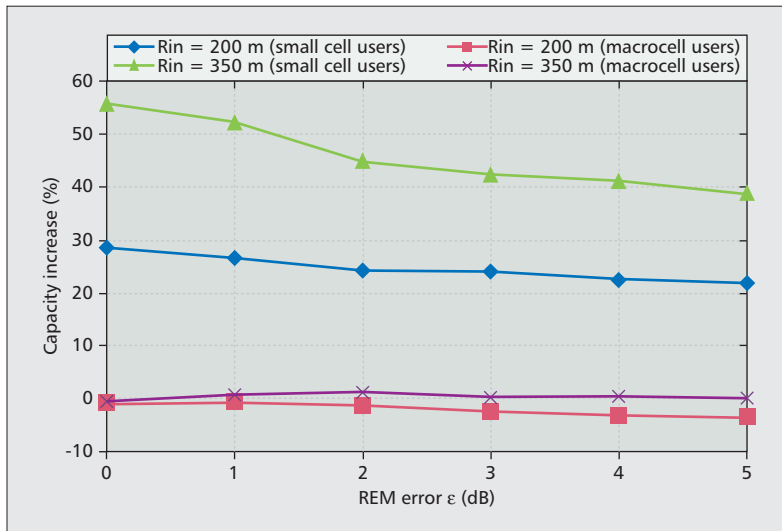


Figure 3. Increase in the downlink average capacity per user achieved by the RFO technique for different values of the REM information error ϵ . The results are presented for a scenario with 12 eNBs and 8 small cells with radius 500 m and 100 m, respectively. Two values of the eNB inner cell radius R_{in} are considered, that is, 200 and 350 m, corresponding to transmit powers of the inner part of 28 and 37 dBm, respectively. The transmit power for the outer users is 43 dBm, and the transmit power of the small cells is 20 dBm. The propagation models are from [14].

the HUEs' throughput, which remains under their minimum requirements. Concerning the RFO technique, Fig. 3 presents the capacity increase with respect to a reference scheme that assumes that eNBs follow a classical fractional frequency reuse and that the sub-band allocated to a small cell is randomly selected among those not used by the closest eNB. Figure 3 shows that, for the ideal case without errors where the REM information matches the real propagation losses (i.e. $\epsilon = 0$ dB in the figure), very significant capacity gains of between 30 and 55 percent for small-cell users, depending on the inner cell radius of the eNBs, were observed with negligible impact on the capacity of MUEs.

MUE outage reduction: In the RAHPC technique, the use of REM can significantly reduce the MUE outage (i.e. the probability of being below the SINR target) with respect to the baseline scheme of [11]. This is illustrated in Fig. 4 in a specific scenario with a transmitting co-channel HeNB in close proximity to an MUE [12]. The biggest improvement is observed for the case when the victim MUE is located close to the eNB so that it receives a strong signal from this eNB. In this case, the outage is reduced by approximately 15 percent in the presented scenario. In this respect, it is envisaged that this type of technique can be useful in future scenarios with extreme densification of cells in certain areas, in which the situations with low signal strength will be reduced and the main challenge will be in the interference control.

Traffic offloading: The use of REM in WiFi/cellular scenarios can lead to efficient traffic offloading from the cellular network to the WiFi network, which in turn simplifies interference management in the cellular network. In particular, it is shown in [15] that up to 30 per-

cent of the total traffic can be shifted to the WiFi network without violating the QoS, assuming that the REM has perfect knowledge about current WiFi channel utilization, WiFi APs, and the locations of small-cell base-stations.

PRACTICAL ASPECTS

Interference management techniques compel the REM manager to obtain fast and reliable access to information from various sources (e.g. from cellular or other non-3GPP network elements or from dedicated sensor networks). Depending on the scenario and the applied algorithms, both dynamic and static information can be considered. The singularities of each case will dictate the best practical approach, thereby attempting to balance the accuracy of the utilized information, latency issues, processing complexity, and related management and security aspects. This section analyzes some of these practical considerations.

Information exchange: Because interference optimization techniques usually require short-time-scale interactions, direct interfaces, such as X2, enable fast data exchange between local REM entities. Signaling requirements will depend on the utilized technique. For example, in the RFO approach, each time the algorithm is executed, the local REM of a cell needs to receive information from its neighbor cells, including the propagation losses to the users of these cells and the transmit power in each sub-band. Let us assume that the propagation loss for a user to a cell is encoded with L_P bits, the transmit power per sub-band is encoded with L_T bits, and the header of the REM message has L_H bits. Then, for a total of N_{bands} sub-bands, N_U users per cell, N_{neigh} neighbor cells and a decision procedure executed on average every t_a s, the local REM signaling requirements per cell will be $(L_H + N_U \cdot L_P + N_{bands} \cdot L_T)N_{neigh}t_a$ (b/s). In the scenario considered in Fig. 3, assuming $L_H = 16$, $L_P = 7$, $L_T = 4$, $N_{neigh} = 20$, $N_{bands} = 4$, $N_U = 10$, and $t_a = 30$ s, the REM signaling requirement is 68 bits/s per cell, which can be considered a quite acceptable value. In the case of more centralized solutions involving global REMs, solutions should be developed under the premise that only local databases are updated frequently, whereas information about global REMs is updated at a lower rate. This will reduce the signaling traffic in the backhaul links.

Building REM information: Various methods can be used to collect REM-related information. In feedback-based mechanisms, the network elements (eNBs, HeNBs, and terminals) collect/measure and report information related to channel gains, location, sub-band use, and so on. However, additional mechanisms may be needed for certain techniques, such as the RAHPC approach, where the local REM in HeNB should estimate the location of victim MUEs that do not communicate their location to the HeNB. In this case, the incorporation of sensing capabilities in the HeNBs or the use of a dedicated sensor network in the HeNB vicinity should be considered.

Robustness against errors: The above mentioned REM building process will impact the reliability of the stored information, which in

turn will influence the performance achieved by a REM-based interference management technique. To illustrate this point, Fig. 3 presents the capacity gain achieved by the FCO technique as a function of the error in the propagation losses stored in the REM. For each value, the error is modeled as a uniformly distributed random variable in the range $[-\epsilon, \epsilon]$ dB. It is observed that, as the REM error increases, the capacity improvements are progressively reduced, although continuing to maintain significant values, revealing the robustness against errors in the considered approach.

REM ownership and management: In the context of interference management, a natural approach is that the REM is owned and managed by the cellular network operator, who will have complete control over the REM functionalities. In the case of multi-operated HetNets, such as when a WiFi network belongs to a different provider than the cellular operator, or when different cellular operators cooperate for interference coordination purposes (e.g. when shared or unlicensed spectrum is used), three possible solutions can be identified for REM management and ownership. The first option is that each operator possesses its own databases, and a dedicated and secured protocol is used for data exchange among 3GPP and non-3GPP networks.¹ In this case, a REM user will have access to the REMs of its operator but will also indirectly benefit from the local and global REMs of cooperating operators. A second option is a hybrid solution whereby some operators decide to merge their REM databases or apply techniques for transparent data sharing. Finally, another option is the establishment of a third-party dedicated provider responsible for REM construction and management. This solution does not exclude the existence of local and global REMs by each operator.

Security and privacy: From the user perspective, because the REM databases may store sensitive information for interference coordination purposes, security and privacy constitute significant challenges. Privacy threats against personal information, such as fine-grained user locations, should be addressed in order not to disclose this information against the users' will. When the REM owner is the operator, the REM should be accessible only from entities residing within the network operator itself, which will ensure that the REM contents will have a similar level of security than other elements of the operator network. Then, users' data integrity and confidentiality can be guaranteed at the same level as the private data of all the mobile users. The same situation occurs for data that the operator does not want to disclose to others. In the hybrid solution whereby some operators merge their databases, information exchange with non-3GPP networks can be realized by dedicated secured protocols (e.g. IPSec), but again, the ownership and security assurances remain under the auspices of the network operators. In that sense, an extension of existing intra-network security solutions could be envisioned for the secure access of the REM. In the case where the interconnecting databases are under the management of a third-party entity, additional security mecha-

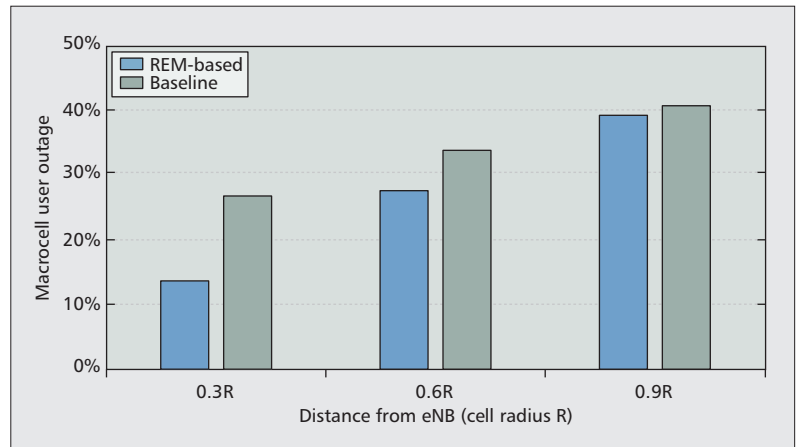


Figure 4. Average MUE outage for various distances between the eNB and the MUE (eNB radius $R = 1$ km). The MUE is placed in a random position close to a co-channel HeNB (inside or outside the house) with an SINR target of 3 dB in a suburban scenario. Details of simulation parameters can be found in [12].

nisms should be designed. Therefore, the relevant regulatory bodies and the operators must determine the best approach to protect the information.

RESEARCH DIRECTIONS

The analysis of previous sections paves the way for further research directions related to REM applicability for interference management in HetNets. In this respect, the following main directions are identified.

- More sophisticated eICIC techniques can be developed by combining the dimensions of power, frequency, time, and user-to-cell association into a multi-parameter optimization framework supported by REMs. REM information will allow the identification of users that are more sensitive to interference in each situation and adjustment of different parameters accordingly. For example, the adequate number of ABSs should be decided to properly protect these users while simultaneously not excessively degrading the performance of the remaining users of the macrocells. Similarly, depending on the characteristics of the scenario, it may be more convenient to allow transmission at a reduced power in certain subframes instead of completely muting the transmission during ABSs. In addition, some subframes can be configured by reserving certain resources in the frequency domain for users located in certain areas of the cell (e.g. the small-cell users located at the edge). These combined techniques require setting multiple parameters (e.g. number of ABSs, transmit powers, and fraction of reserved resources); thus, their optimization requires accurate knowledge of the conditions that each user/cell is experiencing to properly assess the impact of the variations in each parameter on the achieved performance. Specific challenges in this context include the placement of the optimization algorithms, the representation of the REM parameters, the update rate, and the accuracy of the stored information.

- Interference management will be a key com-

¹ For example, the Next Generation Hotspot certified with WiFi Certified Passpoint™, from which rich information on the WiFi network can be obtained.

REMs may support advanced optimization algorithms for interference coordination, efficient sharing of unlicensed bands, energy-efficient user-to-cell associations, as well as backhauling aspects, e.g. to identify/predict the existence of a line-of-sight link between access points.

ponent in scenarios with extreme densification of cells in certain areas (e.g. malls and stadiums). In these scenarios, a terminal may receive high signal levels from a very large number of cells, thus requiring more efficient solutions to interference coordination and user-to-cell association. In such complex environments, the multi-layered REM will address the optimization of interference coordination in a simpler manner and facilitate scalability, e.g. by clustering the local REMs of cells with stronger interactions under the control of a global REM.

• Novel architectural paradigms, such as network function virtualization (NFV) and software defined networking (SDN), are envisaged to facilitate the introduction of the REM concept. NFV refers to the software implementation of network functions running on general purpose computing/storage resources. This can be applied to radio access in the form of a cloud-radio access network (C-RAN). SDN refers to decoupling the network control plane and the data plane, thereby enabling the implementation of control functions as software applications running on top of an SDN controller that provides a programmatic interface to the network. Through SDN/NFV, local/global REMs could be implemented as virtual databases supporting the interference control functions implemented as software packages. In this manner, the system becomes substantially more adaptive and flexible and can be used to dynamically optimize the information split between local/global REMs and quickly introduce new stored parameters as needed by the interference control functions.

• Further challenges arise in the evolution toward heterogeneous 5G systems involving multiple technologies, a wider range of spectrum bands (e.g. millimeter waves and unlicensed bands), additional types of wireless links (e.g. backhaul and device-to-device), and new interworking challenges. REMs may support advanced optimization algorithms for interference coordination, efficient sharing of unlicensed bands, energy-efficient user-to-cell associations, as well as backhauling aspects, e.g. to identify/predict the existence of a line-of-sight link between access points. The vision of a converged network, or network of networks with close cooperation between them, can also be facilitated through REMs that provide the contextual information to make decisions affecting such cooperation.

CONCLUSIONS

This article has proposed the use of REMs to support interference management in HetNets. A general layered architecture including global and local REM databases in the context of a network with both cellular and WiFi technologies has been suggested. The benefits and architectural implications of REMs have been illustrated for various specific techniques, and we have discussed the achievable gains in terms of capacity or outage reduction, including considerations of ownership and security. The article has concluded with a list of research challenges derived from the presented framework.

ACKNOWLEDGMENTS

This work has been supported by the European Commission in the framework of the FP7 NEW-COM# project (contract no. 318306) and by the Spanish Research Council and FEDER funds under RAMSES (ref. TEC2013-41698-R) grant.

REFERENCES

- [1] A. Ghosh *et al.*, "Heterogeneous Cellular Networks: From Theory to Practice," *IEEE Commun. Mag.*, vol. 50, no. 6, June, 2012, pp. 54–64.
- [2] D. Lopez-Perez *et al.*, "Enhanced Intercell Interference Coordination Challenges in Heterogeneous Networks," *IEEE Wireless Commun.*, vol. 18, no. 3, June 2011, pp. 22–30.
- [3] T. Novlan *et al.*, "Analytical Evaluation of Fractional Frequency Reuse for OFDMA Cellular Networks," *IEEE Trans. Wireless Commun.*, vol. 10, no. 12, Dec., 2011, pp. 4294–4305.
- [4] 3GPP R1-103823, "HeNB Power Setting Performance under Different Access Constraints," July, 2010.
- [5] B. Orlandi and F. Scahill, "WiFi-Roaming — Building on ANDSF and Hotspot 2.0," White paper of Alcatel-Lucent and BT, 2012, available at <http://www.tmcnet.com/tmc/whitepapers/documents/whitepapers/2013/6686-wi-fi-roaming-building-andsfand-hotspot20.pdf>, accessed 29th April 2015.
- [6] Y. Zhao *et al.*, "Applying Radio Environment Maps to Cognitive Wireless Regional Area Networks," *Proc. 2nd IEEE Int'l. Symposium on New Frontiers in Dynamic Spectrum Access Networks (DySPAN 2007)*, Dublin, Ireland, April, 2007.
- [7] J. van de Beek *et al.*, "How a Layered REM Architecture Brings Cognition to Today's Mobile Networks," *IEEE Wireless Commun.*, vol. 19, no. 4, Aug., 2012, pp. 17–24.
- [8] B. Sayrac (editor), "D2.4: Final System Architecture," Deliverable of FARAMIR project, December, 2011, available at http://www.ict-faramir.eu/fileadmin/user_upload/deliverables/FARAMIR-D2.4-Final.pdf, accessed 29th April 2015.
- [9] 3GPP TS 36.300 v11.7.0 "Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description; Stage 2 (Release 11)," Sept., 2013.
- [10] SmallCellForum, "Integrated Femto-WiFi (IFW) Networks," document number 033.05.01, Dec., 2013.
- [11] Femtoforum, "Interference Management in OFDMA Femtocells," March, 2010.
- [12] A. Zalonis *et al.*, "Femtocell Downlink Power Control Based on Radio Environment Maps," *IEEE Wireless Commun. Netw. Conf. (WCNC)*, Paris, France, April 2012.
- [13] R. Kurda *et al.*, "Power Adjustment Mechanism Using Context Information for Interference Mitigation in Two-Tier Heterogeneous Networks," *Proc. IEEE 19th Symposium on Computers and Communication (ISCC)*, Madeira, Portugal, June 2014.
- [14] K. Koutlia *et al.*, "On the Use of Gibbs Sampling for Inter-Cell Interference Mitigation under Partial Frequency Reuse Schemes," *Proc. The Third International Conference on Mobile Services, Resources, and Users (MOBILITY)*, Lisbon, Portugal, Nov., 2013.
- [15] A. Kliks *et al.*, "WiFi Traffic Offloading for Energy Saving," *Proc. 20th Int'l. Conf. Telecommunications (ICT)*, Casablanca, Morocco, 6–8 May, 2013.

BIOGRAPHIES

JORDI PEREZ-ROMERO (jorperez@tsc.upc.edu) is an associate professor in the Department of Signal Theory and Communications at the Universitat Politècnica de Catalunya (UPC) in Barcelona, Spain, where he received his telecommunications engineering and Ph.D. degrees in 1997 and 2001, respectively. His research interests are in the field of mobile communication systems, especially radio resource and QoS management, heterogeneous wireless networks, and cognitive networks. He has been involved in different European projects and in projects for private companies.

ANDREAS ZALONIS has a B.Sc. degree in physics (2000) and an M.Sc. degree in telecommunications and electronics (2002), both from National Kapodistrian University of

Athens (NKUA), Greece. Since 2003 he has been working as a research associate at the Institute of Accelerating Systems and Applications (IASA) at NKUA. He has participated in several research projects. His main research interests are radio resource management and flexible/adaptive transceiver design for wireless communication systems.

LILA BOUKHATEM (Lila@Iri.fr) is an associate professor at LRI Laboratory — University of Paris-Sud 11. She received her M.Sc degree in 1998 from the University of Versailles (UVSQ) and her Ph.D. degree in 2001 from the University of Paris 6. She joined the LRI Laboratory in 2002 where she is developing various research projects on mobile and wireless networks, including cross-layer design, modeling and performance evaluation, resource allocation, interference mitigation, mobility management, and routing and energy in ad hoc, vehicular, and sensor networks. She has been involved in several national and international projects and programs.

ADRIAN KLIKS (akliks@et.put.poznan.pl) received with honors his M.Sc. and Ph.D. degrees in telecommunications from Poznan University of Technology (PUT), Poland, in 2005 and 2011, respectively. Currently he is employed as an assistant professor at the Chair of Wireless Communications, PUT. He has been involved in various industrial and international research projects. His scientific interests include advanced multicarrier communications, waveform design, small cells in heterogeneous networks, WiFi offloading, and cognitive radio.

KATERINA KOUTLIA received her B.A. in electronics engineering (2009) from the Technological Institution of Thessaloniki, Greece, and her M.Sc. with distinction (2011) in wireless communication systems from Brunel University, Uxbridge, United Kingdom. Currently she is a Ph.D. student in the Mobile Communication Research Group (GRCM) at UPC. Her current research interests include inter-cell interference coordination (ICIC) techniques, resource management, and heterogeneous networks.

NIKOS DIMITRIOU (SM.IEEE, 2011) holds a diploma in electrical and computer engineering from NTUM (1996), an M.Sc. with distinction in mobile and satellite communications (1997) from the University of Surrey, United Kingdom, and a Ph.D. in mobile communications from the same university (2001). He is currently a research fellow at the Institute of Informatics and Telecommunications at the National Center of Scientific Research "Demokritos." His research interests include radio resource management for dense HetNets and robust routing for mobile ad hoc networks.

REBEN KURDA earned his Bachelor's degree (B.A.) in computer science from the University of Sulaymaniyah, Kurdistan-Iraq (2003), and his M.Sc. degree in informatics from the University of Koya, Kurdistan-Iraq (2007). He recently obtained his Ph.D. degree in wireless communications from the University of Paris-Sud 11, France. He is currently employed as a lecturer at the University of Koya. He has been involved in various industrial and international research projects. His scientific interests include LTE-Advanced multicarrier communications, interference mitigation in small cell and heterogeneous networks, and WiFi offloading.

Vehicle-to-Vehicle Communication in C-ACC/Platooning Scenarios

Alexey Vinel, Lin Lan, and Nikita Lyamin

ABSTRACT

Cooperative adaptive cruise control (C-ACC) and platooning are two emerging automotive intelligent transportation systems (ITS) applications. In this tutorial article we explain their principles, describe related ongoing standardization activities, and conduct performance evaluation of the underlying communication technology.

INTRODUCTION

Sensor-based cruise control (CC) systems are currently deployed worldwide as common driver assistance systems. CC allows a predefined speed to be maintained and thus reduces a driver's workload in free flowing traffic. Conventional adaptive cruise control (ACC), which is also on the market, is an enhancement of CC. ACC enables a preset distance from the preceding vehicle to be maintained. The measurements of the distance are handled by automotive radar mounted on the front of the vehicle (Fig. 1a). A line of vehicles connected by the ACC system is subject to the adverse effect of shockwaves because information on the acceleration/breaking of the first vehicle propagates along the caravan with significant radar measurement-induced delays [1].

Recent advances in vehicular networking [2] make it possible to further enhance ACC in order to avoid shockwaves propagating along a caravan of vehicles. This is achieved by direct vehicle to vehicle (V2V) wireless connectivity and information exchange with one or more of the preceding vehicles so as to maintain the predefined inter-vehicle distances (Fig. 1b). Such a system is referred to as cooperative adaptive cruise control (C-ACC). The information that is transmitted over the wireless connection includes the vehicle's position, velocity, and acceleration. ACC and C-ACC with automatic longitudinal control only, assume that the driver controls the car using the steering wheel. Thanks to inter-vehicle wireless connectivity, information about the maneuvering of the lead vehicle is available almost instantly to the caravan members.

C-ACC can be further enhanced if automatic lateral control of the vehicle is also provided

(Fig. 1c). In such a case, a professional, specially-trained driver manually controls the first vehicle in the caravan, while the others follow it automatically. Such a highly automated system means that drivers revert to manual control in certain situations, although most of the time they are not involved in any driving tasks. Further intelligence, e.g. protocols for joining/leaving the caravan or assisting other vehicles during on-ramp highway merging, can be added to such a system into what results as a platooning application. The differences between the C-ACC and platooning are discussed further in the next section.

Another motivation for C-ACC/platooning is to further reduce the inter-vehicle distances in the caravan, thereby decreasing air drag, which leads to lower fuel consumption [3]. Typically, an interval of 0.5 seconds (12.5 m in 90 km/h) for platoons is considered, while in a typical ACC (no wireless communication involved and based on radar measurements only) the minimum interval is set at 1.6 seconds. The recommended safety interval in Sweden, for example, with no support is set at 3 seconds. Indeed, under Swedish law "the police impose a fine when the safe distance is less than 1 second. If the safe distance is less than 0.5 seconds, the driver's driving license can be revoked." Therefore, a re-evaluation to or an amendment of the legal framework is key to the future development and deployment of automated driving systems.

Several projects on vehicle C-ACC/platooning have recently been carried out. These include Connect&Drive [4], Grand Cooperative Driving Challenge (GCDC) [5], and Safe Road Trains for the Environment (SARTRE). The Connect&Drive and the GCDC projects have C-ACC employing longitudinal control, while the SARTRE project has platoons consisting of heavy-duty vehicles and ordinary passenger cars with both automated longitudinal and lateral control. Demonstration of smart platooning functionalities, such as the merging of two platoons, is planned for 2016 within the framework of the GCDC II (i-Game) project.

The rest of this article is organized as follows. In the following section we provide an overview of the relevant standardization activities. The third section briefly discusses V2V communica-

Alexey Vinel and Nikita Lyamin are with Halmstad University

Lin Lan is with Hitachi Europe SAS.

tion patterns enabling C-ACC/platooning. A system model, performance metrics, and simulation results for the platooning scenario are presented next. We then conclude the article with a discussion of plans for future work.

STANDARDIZATION ACTIVITIES

V2V and vehicle to infrastructure (V2I) are also referred to as cooperative ITS (C-ITS). Key stakeholders in North America and the EU have been driving research and development of C-ITS for more than a decade. Standardization is one of the key building blocks of the C-ITS deployment roadmap. In 2014 the European Telecommunication Standard Institute (ETSI) and the European Committee for Standardization (CEN) jointly delivered the first release of C-ITS standards, enabling deployment of a set of day-one applications. The main target applications supported by the release one standard can be summarized as the cooperative awareness application and the road hazard signaling applications. These applications do not require any intervention to the vehicle electronic systems, but focus instead on providing information or a warning to the driver of a hazardous road situation (decentralized environmental notification message (DENM) [6]) as well as the kinematic state of other vehicles (cooperative awareness message (CAM) [7]). Release one standards also enable transmission of infrastructure information to vehicles via a set of infrastructure to vehicle (I2V) messages, such as signal phase and timing information (SPAT), road topology information (MAP), and road signage information (in-vehicle information). The communication of V2V and V2I messages requires the establishment of direct vehicle to vehicle and vehicle to infrastructure wireless ad hoc network and low latency media access. Therefore, release one standards also include specifications on a specific networking communication stack (geoNetworking, networking functionalities with addressing scheme based on the geographical position of nodes), and access technologies (e.g. EU profile of IEEE 802.11p operating in the 5.9 GHz spectrum band allocated for ITS applications).

In addition, special attention has been paid during the standard specification phase to optimize the network resource usage, given the expected network density level and the amount of data being exchanged between nodes to satisfy the application requirements. For example, ETSI TC ITS operates a decentralized congestion control mechanism to dynamically measure the network load in real time and also to implement functionalities to keep the load below a threshold level. It should be noted that even though ITS-specific technologies standards are made available, the C-ITS does not preclude the use of other technologies, particularly when the penetration rate of the ITS-equipped nodes is low and when the application requirements may be met by other applications (e.g. for nonsafety applications). In fact, legacy communication stacks (e.g. TCP/IPv6 stack) and communication technologies (e.g. cellular network) are also included in the overall ITS communication architecture. Nevertheless, in order to ensure com-

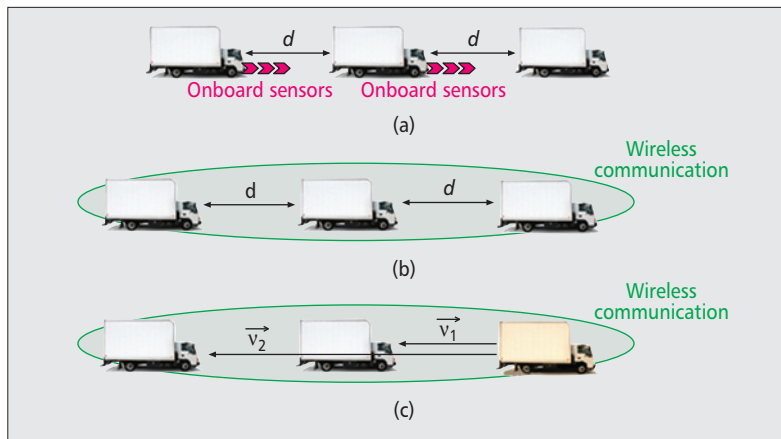


Figure 1. Illustration of ACC/C-ACC/platooning concepts: a) ACC; b) C-ACC; c) platooning.

munication interoperability between vehicles from different vendors from the beginning of the deployment, a common agreement among stakeholders is required. The European Car 2 Car Communication Consortium (C2C-CC) is currently developing recommendations based on release one standards, with the aim of specifying a minimum set of standardized features and minimum sets of system performance to be implemented by all major car manufacturers and system providers in the EU and worldwide.

Among the various messages mentioned, CAM is one of the key basic features required for day one deployment. This is a high-frequency (1-10 Hz) periodic heart-beat message, broadcast by every equipped vehicle to its immediate communication neighbors, providing the vehicle is in the traffic flow and the C-ITS system is in operation. CAM includes the following content:

- Highly dynamic vehicle kinematic data such as position, time, heading, speed, acceleration, and status of acceleration control systems.
- Vehicle attributes such as vehicle width, length, vehicle type, vehicle role.
- Vehicle movement data, including vehicle historical path and path prediction data, e.g. yaw rate and curvatures.
- Additional information for special vehicle types, e.g. emergency vehicles, buses, road maintenance vehicles, and so on.

In the published standard [7], the CAM generation rate is dynamically adjusted between 1 Hz and 10 Hz according to vehicle speed, movement heading, and changes in acceleration. The generation rate is increased whenever there is an increase in the vehicle movement dynamics, to ensure the movement dynamic is correctly reflected in the message content update rate. During its development phase, CAM and the corresponding protocol have been tested, validated in several ETSI conformance and the interoperability test event ETSI Plugtest, as well as in multiple EU R&D and Field Operational Test Projects (FOT). It was published as a European Norm in late 2014.

European ITS standard organizations are currently preparing for release two of ITS standards. Among many potential fields of stakeholder interest, one is the development of

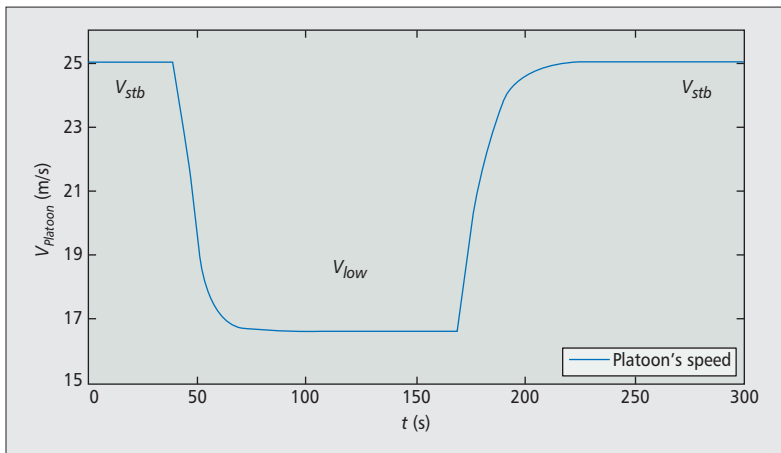


Figure 2. Illustration of the disturbance scenario speed pattern for the platooning application.

C-ITS standards for connected automated driving applications and C-ITS-based advanced driver assistance applications. For example, since April 2014, ETSI TC ITS has established three new work units: C-ACC (TR 103 299), Vulnerable Road User safety (TR 103 300), and Platooning (TR 103 301). The focus of these projects is to conduct a pre-standardization study of these three applications. Instead of developing brand new standards from the very beginning, the pre-standardization study provides an overview of the applications, including their functional and operational requirements (e.g. performance requirements, data exchange requirements, communication requirements, and communication security requirements). The requirements analysis is essential for estimating the applicability of existing standards for these applications, as well as any new standard features that are needed (message sets specifications, communication protocol specifications, communication security features, congestion control requirements, etc). The expected outcomes of these projects are the recommended specifications for future standards required for C-ACC, platooning, and vulnerable road users safety applications.

The initial technical work of the C-ACC and platooning applications in ETSI TC ITS focuses on the development of a high-level definition of C-ACC and platooning applications. This high-level definition is similar to those assumed by us, and can be summarized as follows.

- C-ACC is an embedded in-vehicle system that extends the ACC function so as to further reduce the time gap between the preceding vehicle or preceding traffic. The operation of C-ACC is based on kinematic data directly transmitted from the preceding and/or following vehicles via a V2V communication link. Multiple C-ACC-equipped vehicles may be aligned together to form a convoy (or caravan). Each vehicle is, however, responsible for its own maneuvering. In summary, C-ACC is a distributed automated driving or ADAS system.

- In platooning, a group of vehicles sharing a similar itinerary over a period of time form a vehicle fleet train, coordinated by a platoon

leader. With increased levels of automation, the platoon leader may coordinate with platoon members for group maneuvering (platoon joining/leaving/group speed), or even make decisions for members in certain situations. The platoon leader is also in charge of monitoring the driving environment not only for him/herself, but also for the platoon members. Members of the platoon may be responsible for following the vehicle ahead, so in this respect, C-ACC may be considered as one technology for platoon operations.

- Both longitudinal and lateral control functions may be used in two applications, to further increase operation stability.

- Different levels of automation should be considered in C-ACC and platooning applications.

Several R&D projects have demonstrated that minor extensions to CAM and DENM may be sufficient to support C-ACC and platooning applications. For a platooning application, new messages/protocols may also be needed to enable platoon group operations such as negotiation for joining/leaving the platoon or merging different platoons. In addition, for new features such as cooperative sensing (exchange of vehicle environment perception with other vehicles), cooperative maneuvering would be helpful in realizing automated driving applications. Such projects would bring important technical inputs for standard development work. For example, in January 2015, a new work item on cooperative sensing (TS 103 324: Cooperative Observation Service) was established by ETSI TC ITS.

COMMUNICATIONS FOR PLATOONING/C-ACC

In a platoon situation, the platoon leader should be aware of the kinematic state of the platoon members in real time for monitoring purposes. In addition, the platoon leader may transmit “a platoon control message” to the platoon members for cooperative maneuvering, e.g. platoon group target speed, configured time gap between platoon members, and so on. The present study focuses on the leader receiving messages from all the other vehicles. The “platoon control message” is not, therefore, considered here.

In a C-ACC situation, a C-ACC vehicle follows the preceding vehicle/s and maintains a target time gap with the preceding vehicle. For this purpose, the C-ACC vehicle receives kinematic data on the preceding vehicle/s. In the present study, we assume that kinematic status information is transmitted between vehicles by CAM messages. The simulation work is done for CAM messages, which represent the most stringent cases.

We assume that the CAM messages are broadcast by all the nodes, but we are mainly interested in ensuring that the data age deadline of the leader is met by all the caravan members.

We also assume that the transmitting vehicle should be able to provide kinematic data at an update rate equal to or higher than the maximum CAM transmission rate, to ensure that the transmitted CAM always contains actual vehicle’s kinematic data. It is assumed that both transmitting and receiving vehicles are equipped

with HW/SW solutions that meet certain performance requirements for the processing of CAM messages, including processing at protocol stacks (networking, MAC etc.) and at security. For example, according to TS 101 559 – 1 (RHS) [8], the end to end latency of CAM should be $u_{max} = 300$ ms for a road hazard signaling application. For platooning and C-ACC applications, this end to end time latency requirement may be further reduced.

PERFORMANCE EVALUATION

SYSTEM MODEL

In the model it is assumed that the platoon has a leading vehicle that is steered by a human and $N - 1$ following automated vehicles moving together along a highway. To enable functioning of the platoon control systems, each vehicle executes the following steps:

- Generate CAMs in accordance with ETSI EN 302 637-2 specification [7] (the generation moment is denoted as t_0).
- Generate random transmission delay \sim uniform(0, 50 ms) (processing delay).
- Transmit CAMs on a dedicated channel in accordance with IEEE 802.11p Medium Access Control (MAC) specification [9].

On the receiver side, a random message verification delay \sim uniform(50, 100 ms) is introduced (the moment of time that the verification ends is denoted as t_1).

Following our previous work [10, 11], the following assumptions are made in the present study:

- All the vehicles in the platoon are within each other's communication range. This is a valid assumption for the realistic set-up of a platoon with 20–25 vehicles, when the IEEE 802.11p communication range is in the order of 400–500 m, inter-vehicle distance is 7 m, and truck length is 13 m.
- The kinematic parameters of the leading vehicle are modelled via the constant-acceleration heuristic (CAH) state-of-the-art car-following mobility model [12].
- Random deviations in the velocities of the following vehicles in the caravan are modelled by applying the following approach: we add a random delay $\delta \sim$ uniform[0, $k\sigma$] to a CAM generation moment in order to reflect the non-perfect synchronization between their velocities, where $k = 500$ is the maximum delay expressed in $\sigma = aTimeSlot$ ($aTimeSlot$ is defined in the standard [9]).
- We add independent packet losses to our MAC layer for each pair of nodes (for this work we only need PER values for each ordinary vehicle transmitting to the leader). The Nakagami-m ($m = 1$) propagation model is used.
- The decentralized congestion control (DCC) functionality is disabled.
- Each vehicle is able to update the CAM content for each generated CAM.

PERFORMANCE METRICS

Data Age: The data age u_n is a random variable defined as the time elapsed since the last successfully received packet of vehicle $2 \leq n \leq N$

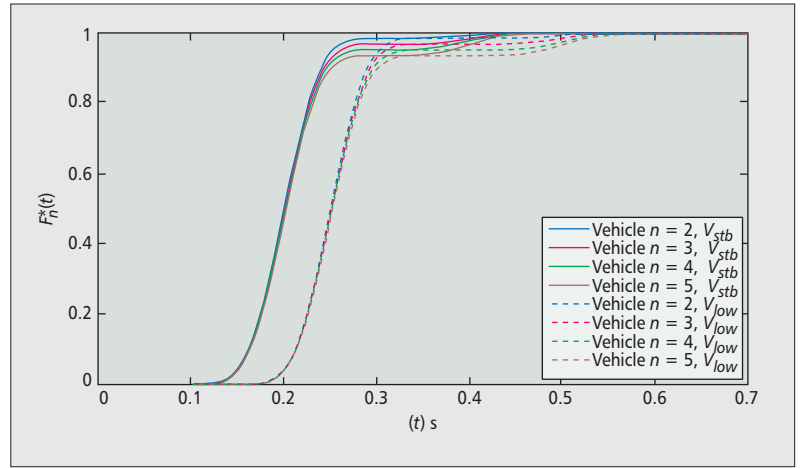


Figure 3. Empirical cumulative distribution functions of data age for a platoon of $N = 5$ vehicles.

by the leading vehicle. Data age is the difference between $t - t_1$, where t is the current moment of time and t_1 is the moment when the last successfully received packet of vehicle n was received by the leading vehicle.

Note: data age relates to the leader and is computed for an ordinary vehicle. We assume that the platoon leader determines the inter-distance for all platoon members. Platoon members use automated driving to maintain the distance specified by the leader.

Cumulative Distribution Function (CDF) of the Data Age: For a particular vehicle n :

$$F_n(t) = \Pr\{u_n \leq t\}$$

We denote the respective *empirical CDF (ECDF)* of the data age for a particular vehicle n as $F_n^*(t)$.

Data Age Deadline: The data age deadline u_{max} is the maximum acceptable data age of a vehicle from the leader's perspective.

Probability to Meet the Deadline: The probability that data age value will not exceed deadline:

$$U_n^* = \Pr\{u_n \leq u_{max}\}$$

CURRENT STANDARDIZATION

Let us evaluate if the current CAM generation rules are sufficient to meet the platoon/C-ACC needs.

We chose the parameters sampling period (Δ) and disturbance parameter (δ) so that the CAM generation moments synchronization effect discussed in [11] is eliminated. We fix the mobility pattern to the disturbance scenario presented in [7]. In the scenario in Fig. 2, the leading vehicle decelerates from the desired steady speed ($V_{stb} \sim 90$ km/h) to a lower speed ($V_{low} \sim 60$ km/h), maintains this speed for some time, and then accelerates back to the initial speed. The disturbance scenario could be regarded as a pattern to describe the appearance of a slow moving vehicle in front of the platoon or a road speed limit. This corresponds to a CAM generation rate change from $1/[4/V_{stb}] = 6.25$ Hz (generation interval of 160 ms) to $1/[4/V_{low}] = 4.25$ Hz (generation interval of 240ms). Addi-

tionally we provide results for the scenario widely used in the literature when CAMs are generated with a fixed frequency of 10 Hz and compare performance of both approaches.

In Fig. 3 ECDF of data age (hereafter, data age ECDF) of each ordinary vehicle in the platoon composed of $N = 5$ vehicles is shown. Solid lines show data age ECDF when the platoon maintains (V_{stb}) speed while dashed lines indicate (V_{low}) speed. Since the message-triggering process according to ETSI EN 302 637-2 relies on the current values of kinematic parameters, the data age of each vehicle will proportionally decrease/increase as the respective speed increases/decreases.

Obviously, vehicles located farther from the leader experience higher packet loss due to fading and as a consequence have higher data age. Later in this article we will focus on the data age of the last vehicle $n = N$ in the platoon.

Figure 4 shows the frequency distribution of data age for the last vehicle in a platoon of $N = 25$ vehicles when the platoon maintains V_{stb} . In our setup the data age for the most distant vehicle may exceed 1 second. The reason for such high values and the form of the distri-

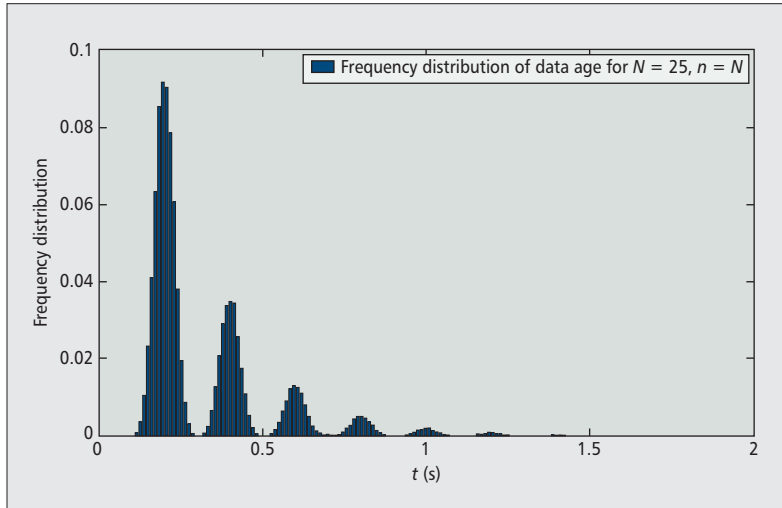


Figure 4. Frequency distribution of data age for the last vehicle in platoon of $N = 25$ vehicles when the platoon maintains V_{stb} .

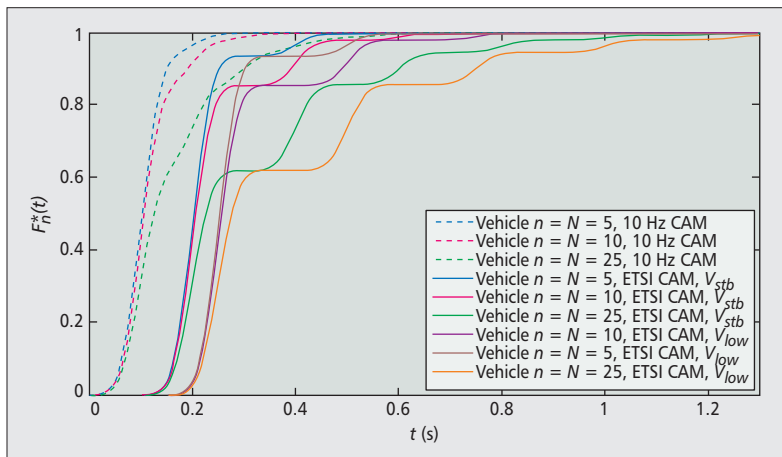


Figure 5. Empirical cumulative distribution functions of data age for a platoon of $N = 5, 10,$ and 25 vehicles with fixed and variable CAM rates.

bution is that data age may increase proportionally to the number of subsequently lost CAMs.

Figure 5 illustrates data age ECDF's for the last vehicle in the platoons of length $N = 5, 10, 25$ vehicles. With an increase in platoon size, data age could increase significantly, which could make operation of the control system difficult.

Table 1 shows the probabilities U_n^* of meeting the deadline $u_{max} = 300$ ms for the last vehicle in a platoon of length $N = 5, 10, 25$. Since the CAM generation rate for ETSI EN 302 637-2 is kinematic-dependent, the probability of missing the deadline when the speed is V_{low} becomes higher. The situation becomes even more acute if the platoon decelerates to lower speed values. In contrast, U_n^* for a fixed 10 Hz mechanism is predictable and depends only on the size of the platoon.

Figure 5 shows a comparative data age distribution between when CAMs are triggered in accordance with ETSI EN 302 637-2 (solid lines) and when employing a constant frequency of 10 Hz (dashed lines). It should be noted that when the platoon moves at V_{stb} , the corresponding generation frequency is about 6.25 Hz ($1/[4/V_{stb}]$). Since we propose a dedicated communication channel for platoon coordination, even for $N = 25$ members, 10 Hz will always outperform the ETSI EN 302 637-2 approach (they may perform equally when the platoon's speed exceeds $1/[4/V_{stb}] = 10, V_{stb} = 40\text{m/s} = 144$ km/h, which is an unrealistic speed pattern for a platooning application). The main conclusion to be drawn is that a 10 Hz CAM rate would be preferable to the current triggering condition, particularly when platoon speed is high. Another conclusion is that although the platoon leader receives CAMs from the platoon members, the current standard CAM rates tend to be insufficient for the leader to maintain the desired 0.5 second distance for safe operation. Therefore, the CAM rate should be further increased.

RECOMMENDATIONS FOR IMPROVEMENT

- Enable constant CAM generation rates exceeding 10 Hz in a platoon, especially at higher speeds.
- Further reducing the processing delay at the receiving vehicle may be beneficial. In particular, the security-related processing delay has an important impact on data age.

FUTURE PLANS

- In our future work we will:
- Take DCC into account in future simulation studies.
 - Improve CAM message content so it can distinguish between platoon and non-platoon members, e.g. group identification.
 - Introduce messages and protocols for platoon control in the overall traffic flow, e.g. space reservation for platoon lane change.

ACKNOWLEDGMENTS

This study is supported by NFITS — the National ITS Postgraduate School (Sweden), and is a part of the “ACDC: Autonomous Cooperative Driving: Communications Issues” project (2014-2016) funded by the Knowledge Foundation

(Sweden) in cooperation with Volvo GTT, Volvo Cars, Scania, Kapsch TrafficCom, and Qamcom Research & Technology.

The authors also express their gratitude to Denis Kleyko from Lulea University of Technology for his valuable comments, which helped to improve the quality of the manuscript.

REFERENCES

- [1] L. Xiao and G. Feng "Practical String Stability of Platoon of Adaptive Cruise Control Vehicles," *IEEE Trans. Intell. Transp. Syst.*, vol. 12, no. 4, 2011.
- [2] G. Karagiannis et al., "Vehicular Networking: A Survey and Tutorial on Requirements, Architectures, Challenges, Standards and Solutions," *IEEE Commun. Surveys & Tutorials*, vol. 4, no. 13, 2011, pp. 584–616.
- [3] A. Al Alam et al., "An Experimental Study on the Fuel Reduction Potential of Heavy Duty Vehicle Platooning," *Proc. 13th Int'l. IEEE Conf. Intelligent Transportation Systems (ITSC)*, 2010.
- [4] J. Ploeg, A. F. A. Serrarens, and G. J. Heijenk, "Connect & Drive: Design and Evaluation of Cooperative Adaptive Cruise Control for Congestion Reduction," *J. Modern Transportation*, vol. 19, no. 3, 2011, pp. 207–13.
- [5] J. Ploeg et al., Introduction to the Special Issue on the 2011 Grand Cooperative Driving Challenge," *IEEE Trans. Intell. Transp. Syst.*, vol. 13, no. 3, 2012, pp. 989–93.
- [6] ETSI EN 302 637-3 V1.2.2 (2014-11) Intelligent Transport Systems (ITS); Vehicular Communications; Basic Set of Applications; Part 2: Specification of Decentralized Environmental Notification Basic Service.
- [7] ETSI EN 302 637-2 V1.3.2 (2014-11) Intelligent Transport Systems (ITS); Vehicular Communications; Basic Set of Applications; Part 2: Specification of Cooperative Awareness Basic Service.
- [8] ETSI TS 101 539-1 V1.1.1 (2013-08) Intelligent Transport Systems (ITS); V2X Applications; Part 1: Road Hazard Signalling (RHS) application requirements specification.
- [9] IEEE Std. 802.11-2012, Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) specifications, Mar. 2012.
- [10] N. Lyamin, A. Vinel, and M. Jonsson, "Does ETSI Beaconing Frequency Control Provide Cooperative Awareness?" *Proc. IEEE ICC 2015 — Workshop on Dependable Vehicular Communications (DVC)*.
- [11] N. Lyamin et al., "On the Performance of ETSI EN 302 637-2 CAM Generation Frequency Management," *Proc. 2014 IEEE Vehicular Networking Conference (VNC)*.
- [12] A. Kesting et al., "Enhanced Intelligent Driver Model to Access the Impact of Driving Strategies on Traffic Capacity," *Philosophical Trans. of the Royal Society A*, 368, 2010, pp. 4585–4605.

ADDITIONAL READING

- [1] J. Dongyao, K. Lu, and J. Wang, "A Disturbance-Adaptive Design for VANET-Enabled Vehicle Platoon," *IEEE Trans. Veh. Technol.*, 63.2, 2014.

$U_n^*, n = N$	ETSI V_{stb}	ETSI V_{low}	10 Hz
$N = 5$	0.9337	0.8950	0.9974
$N = 10$	0.8530	0.8178	0.9869
$N = 25$	0.6182	0.5943	0.8995

Table 1. Probability $u_n \leq u_{max}$ to meet the deadline for vehicle n .

BIOGRAPHIES

ALEXEY VINEL (M'07, SM'12) received the bachelor's (Hons.) and masters' (Hons.) degrees in information systems from Saint-Petersburg State University of Aerospace Instrumentation, Saint Petersburg, Russia, in 2003 and 2005, respectively, and the Ph.D. degrees in technology from the Institute for Information Transmission Problems, Moscow, Russia, in 2007, and the Tampere University of Technology, Tampere, Finland, in 2013. He is currently a professor of data communications at the School of Information Technology, Halmstad University, Halmstad, Sweden. He has been involved in research projects on vehicular networking standards, advanced driver assistance systems, and autonomous driving. He has been an associate editor for the *IEEE Communications Letters* since 2012.

LAN LIN received her bachelor's degree at Tongji University in China and a master's degree at the Civil Engineering School (Ecole Nationale des Ponts et Chaussees) in France at 2005. She joined the Hitachi Europe Big Data Lab in Sophia Antipolis, France shortly after. She is senior researcher at Hitachi Europe, actively involved in R&D activities on intelligent transport system, smart mobility, big data and digital manufacturing. She has been involved in R&D projects in the field of ITS, vehicular communications, smart mobility, and smart cities such as iTETRIS, COVEL, SCOREF, PREDIVE C2X, DRIVE C2X, ECO-FEV, AutoNet2030, and so on. She also participates in standardization activities in different standards organizations such as ETSI, SAE, IEEE, and IEC. She is currently chair of ETSI TC ITS WG1 and rapporteur of several Work Items of the WG 1 (application and user requirements), as well as vice-chair of European Car to Car Communication Consortium (C2C-CC) application working group (WG).

NIKITA LYAMIN received a B.S. degree (Hons.) and M.S. degree (Hons.) in 2011 and 2013, respectively, in telecommunications from Siberian State University of Telecommunications and Information Sciences, Novosibirsk, Russia. He is now a Ph.D. student at the School of Information Technology, Halmstad University, Halmstad, Sweden. His current research interest is in the areas of vehicular ad-hoc networks, platooning, and multiple-access protocols.

Greater Reliability in Disrupted Metropolitan Area Networks: Use Cases, Standards, and Practices

Ming-Tuo Zhou, Masayuki Oodo, Vinh Dien Hoang, Liru Lu, Xin Zhang, and Hiroshi Harada

ABSTRACT

Unlike most general wireless users enjoying broadband access to the Internet and so on, a number of mission-critical applications such as PPDR are using narrowband systems that are only capable of transmitting voice and low-rate data. Broadband wireless networks with greater reliability are increasingly demanded by these applications including the emerging Smart Grid. With this vision, the IEEE 802.16 Working Group completed two standard amendments recently, IEEE Std 802.16n-2013 and IEEE Std 802.16.1a-2013. 802.16 network reliability is significantly increased by technical developments dealing with failures of network infrastructure stations, radio path and backhaul connectivity, and so on. This article serves to analyze the envisaged typical use cases of greater reliable broadband wireless networks, outline the main technical developments of the two standards, and introduce a recently developed practical system for PPDR applications in Japan.

INTRODUCTION

With the popularity of advanced wireless technologies, particularly international mobile telecommunications (IMT), general wireless consumers in many countries have gained broadband access to databases, media, and the Internet; however, in contrast, most existing wireless systems for particular mission-critical applications, such as public protection and disaster relief (PPDR), avionics, airport surface communications, maritime safety, and surveillance, are narrowband and only capable of transmitting voice and low-rate data. With considerable growing demand for video and high-speed data, and the emerging trend toward voice over Internet Protocol (VoIP), the need for broadband wireless systems is becoming ever more important to the stakeholders of these applications. On the other hand, because of the essential importance of critical missions, the serving broadband wireless systems need to be much more reliable than their general counterparts. In the case of any network disruptions

such as failure of the network infrastructures and/or the radio path connections, the networks are required to be able to recover from degradations quickly so that the services can be uninterrupted.

IEEE 802.16, one IMT technologies, has great potential to serve for the above applications since it is capable of providing a broadband IP-based network supporting voice, data, and video services. It has also been discussed for serving as backhaul and last-mile access networks for the emerging smart grid, which has similar need for broadband communications for meter readings, facility monitoring, real-time demand response, and so on. However, by 2009, none of the completed 802.16-family standards had considered the requirement for greater reliability for these applications. Realizing this need and after about one year of full study, the IEEE 802.16 Working Group launched a project with an objective to provide Greater Reliability in a Disrupted Metropolitan Area Network, namely GRIDMAN, in July 2010. The Working Group completed two GRIDMAN standards in 2013, that is, IEEE Std 802.16n-2013 [1] and IEEE Std 802.16.1a-2013 [2], as amendments to the two base standards, IEEE Std 802.16-2012 and IEEE Std 802.16.1-2012, respectively.

The rest of this article presents the developed GRIDMAN mechanisms improving 802.16 network reliability based on illustrating the two most typical use cases, PPDR and smart grid, and introduces a practical wireless system developed in Japan based on the GRIDMAN protocols recently.

TYPICAL USE CASES AND SYSTEM REQUIREMENTS

A number of use cases were identified by the GRIDMAN Task Group to demonstrate the need for more reliable 802.16 networks and to facilitate discussion of the system requirements [3]. Based on the study, the functional mechanisms needed to achieve higher reliability were summarized, and later the standard protocols with technical details were developed.

The authors are with the National Institute of Information and Communications Technology.

DISRUPTED METROPOLITAN AREA NETWORK

A typical 802.16 system consists of one or several networks co-located at a geographic location. Each network comprises a base station (BS), tens or hundreds of mobile stations (MS), and possibly a number of relay stations (RS). Generally, a BS is connected to a backhaul network, by which the 802.16 network is managed and connected to other networks. A disrupted network refers to failure of one or more infrastructure nodes (BS or RS), loss of BS connection to backhaul, failure of MS connection to BS, and so on. To be robust and reliable, a more reliable network needs to automatically recover from degradation quickly.

TYPICAL USE CASES

The two most typical use cases recognized by the GRIDMAN Task Group were PPDR [4] and smart grid (SG) [5]. The two use cases are both of critical missions and need broadband wireless communications with coverage range of several to tens kilometers.

Public Protection and Disaster Relief: PPDR refers to radiocommunications used by responsible agencies and organizations dealing with maintenance of law and order, protection of life and property, disaster relief, and emergency response. PPDR communications are predominantly mission-critical because these aid in the protection and safety of life or property on a day-to-day basis as well as in response, rescue, and recovery efforts before, during, and after emergencies and disasters. A PPDR network should be available whenever users need the services.

A PPDR wireless system is usually a land mobile radio network. In the past, several robust, secure, and highly reliable digital standards have been established for PPDR applications. However, these systems are mostly narrowband (less than 25 kHz), and can only support voice and low-data-rate applications. Although there will continue to be narrowband requirements, there is increasing need for broadband (data rates on the order of 1–100 Mb/s) for future PPDR applications (e.g., video surveillance). IEEE 802.16 may play an important role in broadband PPDR applications since it supports data rates of several to tens of megabits per second, and a communication range of several to tens of kilometers.

As illustrated in Fig. 1a, an 802.16 network with greater reliability may enable far more reliable PPDR services. For example, when a BS is not functional due to an earthquake, an MS may act as a serving BS to maintain normal network operation; when a BS loses backhaul connection because of flooding, it can relay traffic to a neighbor BS that has backhaul connection; an alternative path can be maintained and activated quickly when an MS loses direct connection to its serving BS; an MS that is out of coverage of a serving BS may connect to the network via relay or forwarding of an in-coverage MS; and a standalone network can be set up quickly for communications of firefighters and commanders when a fire disaster happens, and so on.

Smart Grid: Smart grid enables intelligent management and control of the generation, transmission, distribution, and consumption of electricity.

Initiatives on smart grid have arisen in the United States, Europe, and many other countries in recent years. Communication infrastructure is one of the critical technologies enabling smart grid. With high-speed long-distance coverage, 802.16 is a very promising technology to serve as both the backhaul and last-mile access network for smart grid. As shown in Fig. 1b, both homogeneous and heterogeneous networks can be employed. In the former, each home appliance and power meter can be equipped with an 802.16 MS end, as well as the facilities of a bulk/distributed power generator (for generation), a substation (for transmission), and a microgrid (for distribution). In the latter, a home-area energy management network and another equipment/asset monitoring network can employ other technologies such as 802.15.4, and then be connected to an 802.16 network via a gateway.

A smart grid requires very high reliability or availability that is usually represented by percentage targets, such as backhaul at 99.999 percent (about 5 min/yr downtime), substation at 99.99 percent (a little less than 1 h/yr downtime), or feeder devices at 99.9 percent (about 9 h/yr downtime). With greater reliability, a smart grid communication infrastructure may recover quickly from degradation so that disruptions of the grid equipment can be quickly reported and attended to, and real-time services such as demand response can be uninterrupted. As shown in Fig. 1b, similar to PPDR, better reliability technologies can be used to handle failures of 802.16 infrastructure stations, loss of BS connection to backhaul, failure of MS connections to BS, and so on, in smart grid.

SYSTEM REQUIREMENTS

General and functional requirements of GRIDMAN have been discussed based on the study of use cases [6]. Generally, a more reliable 802.16 network is required to be backward compatible and to support services that require a higher degree of assurance of maintaining sufficient connectivity than legacy 802.16 networks. Considering the wide range of radio frequency bands allocated for the envisaged applications, the operation frequencies range from 200 MHz up to 6 GHz. Table 1 lists the main functions that are required to handle the envisaged degradations. Detailed descriptions of the functional requirements and the corresponding realized mechanisms are presented later.

STANDARD PROTOCOLS

Protocols with technical details were developed by the Task Group after identifying the system requirements. Although 802.16n and 802.16.1a have different medium access control (MAC) and physical layer (PHY) protocols, the ideas and main procedures of the new technical developments in the two standards are basically the same. The following focuses on the main basic technical principles, rather than going into technical details. Table 2 lists the mechanisms, operations, effects, and illustrations of the main standard protocols developed, and the following presents the corresponding technical descriptions.

MULTI-MODE OPERATIONS

Relay Function of an HR-BS: The relay function of an HR-BS allows an HR-BS to lose backhaul connection to relay traffic to a neighbor HR-BS

With greater reliability, a smart grid communication infrastructure may recover quickly from degradation so that disruptions of the grid equipment can be quickly reported and attended to, and real-time services such as demand response can be uninterrupted.

While operating as an HR-RS, the station may maintain certain HR-MS functionalities. A mode switch to HR-RS shall be commanded by its superordinate HR-BS. If an HR-MS is released from its role from the relay mode as a relay, it may perform handover to any infrastructure station.

that has backhaul connection. When the connection to a backhaul is lost, an affected HR-BS scans its environment and establishes a relay link with a selected superordinate HR-BS. The relay link is configured and maintained by the superordinate HR-BS. Finally, when the affected HR-BS recovers backhaul connection, the relay link will be released and hand over back to the recovered HR-BS.

The affected HR-BS may operate in either time-division transmit and receive (TTR) mode

or simultaneous transmit and receive (STR) mode. In TTR mode, the affected HR-BS maintains connectivity with its subordinate HR-MSs by performing dual-role BS/RS operation; while in STR mode, the affected HR-BS maintains BS functionality and therefore supports its subordinate HR-MSs continually. In both cases, an affected HR-BS performs the following activities:

- To establish a relay link with a serving HR-BS

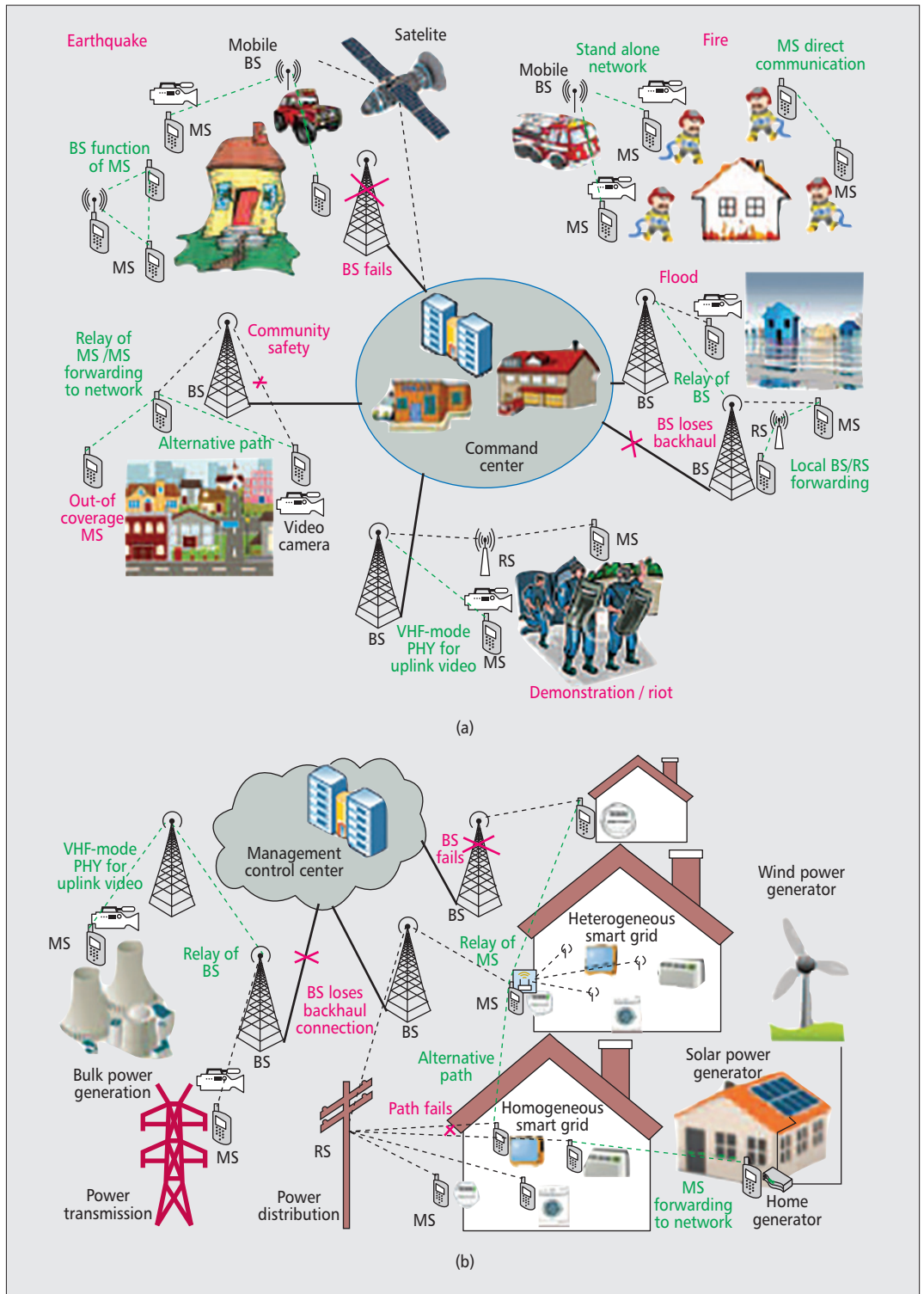


Figure 1. Use cases of GRIDMAN: a) PPDR; b) smart grid.

- If necessary, to inform some subordinate stations to perform handover
- If necessary, to reconfigure the physical frame
- To commence operation in relay mode

Relay Function of HR-MS: An HR-MS may operate as relay to provide connectivity for HR-MSs out of network coverage. During basic capability negotiation at network entry, an HR-MS that is capable of role changing to HR-RS shall report such capability to the superordinate HR-BS/HR-RS. While operating as HR-RS, the station may maintain certain HR-MS functionalities. A mode switch to HR-RS shall be commanded by its superordinate HR-BS. If an HR-MS is released from its role from the relay mode as a relay, it may perform handover to any infrastructure station.

To support relay function for HR-MS, an HR-MS capable of relay function may establish a relay link with its serving HR-BS. An HR-MS acting as HR-RS operates in either TTR mode or STR mode, and its relay mode is determined by the serving HR-BS.

BS Function for HR-MS: An HR-MS may operate as an HR-BS to provide connectivity for itself and other HR-MSs. During basic capability negotiation at network entry, an HR-MS that is capable of role changing to HR-BS shall report such capability to the superordinate HR-BS/HR-RS. While operating as an HR-BS, the station may maintain certain HR-MS functionalities.

An HR-MS may change to an HR-BS in a proactive or reactive approach. In proactive operation, the mode switch is directed by its superordinate HR-BS, while in a reactive approach, the mode switch is initiated by the HR-MS itself. When a superordinate HR-BS fails, several HR-MSs may contend to operate as an HR-BS, and carrier sense multiple access with collision avoidance (CSMA/CA) protocol with backoff algorithm is used to resolve collisions.

ALTERNATIVE PATH DISCOVERY AND MANAGEMENT

An HR-MS may discover and maintain an alternative path to a neighbor network or to its serving HR-BS/HR-RS. In the former case, if its current serving HR-BS/HR-RS fails, the HR-MS enters the neighbor network quickly. In the latter case, if its current connection to its serving HR-BS/HR-RS fails, it reenters the network quickly by going through a forwarding neighbor HR-MS that has connection to the serving HR-BS/HR-RS.

Alternative Path to a Neighbor HR-BS/HR-RS: In order to discover an alternative path to a neighbor HR-BS/HR-RS for an HR-MS, a serving HR-BS/HR-RS scans for neighbor HR-BSs/HR-RSs and evaluates the scan report. It determines an alternative path if it exists and sends this information to the HR-MS using a handover request message. The message contains a recommended ranging code and slot from the ranging region of the target HR-BS/HR-RS, thus to enable fast network reentry if needed since no contention will occur when multiple

Degradation		Functions required to support
Infrastructure failure	An HR-BS loses connection to backhaul	Multi-mode operation: relay of HR-BS
		Forwarding between infrastructure stations
	Standalone network	
	An HR-BS is not functional	Multi-mode operation: BS function of HR-MS
Failure/lack of radio connectivity/network redundancy	An HR-MS fails in connecting to its serving HR-BS	Alternative path discovery and management
	Two HR-MSs are within each other's coverage range	HR-MS direct communication
	An HR-MS is out coverage of an HR-BS	Multi-mode operation: relay function of HR-MS
		HR-MS forwarding to network
	Two HR-MSs are within coverage of an HR-BS/HR-RS	Local HR-BS/HR-RS forwarding
Too large DL delay spread and/or heavy uplink load		VHF-mode PHY

Table 1. Main GRIDMAN functional requirements for handling network degradations.

HR-MSs perform fast network recovery. The alternative path may be updated or replaced by a new one by the serving HR-BS/HR-RS based on the most recent scanning results.

Alternative Path Going through a Forwarding HR-MS: An HR-MS may establish an alternative path to its serving HR-BS/HR-RS through a neighbor HR-MS, which is called a forwarding HR-MS. To do so, it requests the serving HR-BS/HR-RS to discover neighbor HR-MSs. Based on neighbor information received from the serving HR-BS/HR-RS, the HR-MS selects one neighbor as the forwarding HR-MS and informs the serving HR-BS/HR-RS. If the serving HR-BS/HR-RS detects loss of direct connectivity to the HR-MS, it initiates a coverage extension process, that is, instructing the forwarding HR-MS to facilitate network entry of HR-MSs outside coverage of the serving HR-BS/HR-RS by broadcasting preambles and control messages. On the other hand, when it detects loss of direct connectivity to the serving HR-BS/HR-RS, the HR-MS starts to scan preambles and control messages from the target forwarding HR-MS for fast network reentry.

In the above, to detect connectivity with an HR-MS, a serving HR-BS/HR-RS schedules the HR-MS to transmit ranging code periodically, and it measures the quality of the received ranging signal. Loss of connectivity with the HR-MS can be detected if the received ranging signal is below a certain strength or no ranging code from the HR-MS detected at all.

Mechanism		Degradation	Operation	Effect	Illustration
Multi-mode operation	Relay function of HR-BS	A HR-BS loses connection to backhaul	The affected HR-BS sets up relay connection to a neighbor HR-BS that has connection to backhaul	Network connectivity of the subordinate stations of the affected HR-BS could be maintained	
	Relay function of HR-MS	An HR-MS is out of coverage of an HR-BS	An HR-MS relays traffic between its serving HR-BS and the HR-MS out of the HR-BS's coverage	The out-of-coverage HR-MS has network connection	
	BS function of HR-MS	The serving HR-BS fails operation	An HR-MS acts as BS	Connections of the subordinate HR-MSs/HR-RSs could be maintained	
Alternative path discovery and management	Alternative path to a neighbor HR-BS/HR-RS	An HR-MS's connection to its serving HR-BS/HR-RS fails	Switch to the alternative path to a neighbor network's HR-BR/HR-RS	The connection of a degraded HR-MS to the backhaul maintains	
	Alternative path going through a forwarding HR-MS		Switch to the alternative path through the forwarding HR-MS to the serving HR-BS/HR-RS		
Standalone network		A serving HR-BS has no ways connecting to a backhaul	The degraded HR-BS provides local connectivity to its subordinate stations within its coverage	All service flows established between the HR-MSs of the standalone network can be maintained	
Forwarding between infrastructure station		A serving HR-BS loses connection to a backhaul	The degraded HR-BS connects to a neighbor HR-BS with backhaul connection through forwarding of an HR-MS in between	The connectivity of the degraded HR-BS network can be maintained	
HR-MS direct communication		Two HR-MSs are within each other's coverage	The two HR-MSs set up direct communication link	The two HR-MSs can communicate directly	
HR-MS forwarding to network		An HR-MS is out of coverage of an HR-BS or it has low data rate to an HR-BS	The forwarding HR-MS forwards user data and control information between the forwarded HR-MS and the serving HR-BS	The forwarded HR-MS has connection to the serving HR-BS or the data rate is increased	
Local HR-BS/RS forwarding		Two HR-MSs within coverage of an HR-BS/HR-RS have no connection	The serving HR-BS/HR-RS sets up a forwarding path between the two forwarded HR-MSs	The two HR-MSs can communicate without going through backhaul	
VHF-mode PHY		An HR network operates on a VHF channel and experiences too large delay spread due to wide coverage	The HR network adopts VHF-mode PHY, where DL tile has more pilot symbols along frequency dimension and UL tile has more data symbols along time dimension	Downlink performance is improved as the severe DL delay spread is combated by the increased DL pilot symbols. Heavy traffic uplink can be supported due to more UL data symbols.	

Table 2. The main developed mechanisms of GRIDMAN.

To maintain an alternative path through a forwarding HR-MS, a serving HR-BS/HR-RS schedules periodic ranging code transmission and reception by an HR-MS and its forwarding HR-MS. The forwarding HR-MS periodically reports measurement results, and if connectivity between the HR-MS and the forwarding HR-MS is lost, the serving HR-BS/HR-RS helps the HR-MS to discover another alternative path.

STANDALONE NETWORK

When an HR-BS, including an HR-MS acting as HR-BS with multimode operation, loses backhaul connection and has no other way to connect to any neighbor HR-BS, it forms a standalone network together with its subordinate HR-MSs and HR-RSs. In a standalone network, local connectivity can be provided for all subordinate stations within coverage of the degraded HR-BS so that all service flows established between the HR-MSs of the standalone network are maintained. In 802.16n, local connectivity is maintained by maintaining a network topology table at the degraded HR-BS. The network topology table is updated by periodically broadcasting (by the degraded HR-BS) and acknowledging (by the HR-MSs and HR-RSs) network topology messages. In 802.16.1a, local connectivity is established by a dynamic service addition procedure with possible help of blind paging. If the degraded HR-BS has no connectivity information on a targeted HR-MS when it establishes local connectivity between two HR-MSs, it tries to confirm connectivity with the HR-MS by blind paging.

In a standalone network, an unassociated HR-MS is still allowed entry to the network and to establish connection with the degraded HR-BS.

FORWARDING BETWEEN INFRASTRUCTURE STATIONS

To provide higher reliability and robustness against a single point of failure such as an HR-BS that loses backhaul connection to a core network, an HR network (i.e., a GRIDMAN network) supports forwarding between infrastructure stations (FBIS). With this function, an HR-BS losing backhaul connection may connect to another HR-BS that has backhaul connection through a subordinate HR-MS.

When a degraded HR-BS loses backhaul connection and is outside coverage of any neighbor HR-BS, it may instruct one or more of its subordinate HR-MSs to discover a neighbor HR-BS that has backhaul connection. Upon receiving such an order, an HR-MS scans its environment for neighbor HR-BSs and reports the result to the degraded HR-BS. With information collected, the degraded HR-BS selects a target HR-BS and a corresponding HR-MS (called FBIS HR-MS). It then instructs the selected FBIS HR-MS to carry out handover to the target HR-BS. After completion of handover, the FBIS HR-MS maintains two connections, one to the degraded HR-BS, the other to the target HR-BS. It switches access between the two HR-BSs.

The degraded HR-BS uses the FBIS connection to serve its subordinate HR-MSs. Establishment of FBIS connection is by dynamic service

addition (DSA), which is carried out first between the degraded HR-BS and the FBIS HR-MS, and second between the FBIS HR-MS and the targeted HR-BS. When the degraded HR-BS recovers a connection to the backhaul network, an FBIS connection can be terminated by a procedure of dynamic service deletion (DSD) and network handover, which enables the FBIS HR-MS to connect back to the previous degraded HR-BS normally.

HR-MS DIRECT COMMUNICATIONS

In HR-MS direct communication, data packets are exchanged between two HR-MSs directly. It is applicable when both HR-MSs, only one of them, or neither of them are within coverage of a serving HR-BS/HR-RS. BS-controlled HR-MS direct communication is mainly introduced in this article. In case two (or more) HR-MSs are not within coverage of any HR-BS/HR-RS, one of them can change mode as HR-BS and then take on the role of management.

Neighbor Discovery: By neighbor discovery, an HR-MS knows the possibility of communicating with other HR-MSs directly. The serving HR-BS/HR-RS schedules one or more associated HR-MSs to broadcast a ranging sequence, and arranges some other associated HR-MSs to receive. Upon receiving the ranging sequence, and if it meets reporting criteria, an HR-MS reports its measurement to the serving HR-BS/HR-RS.

Connection Setup and Management: A direct communication link shall be set up to support direct communication between two HR-MSs. By checking a neighbor table, a serving HR-BS/HR-RS knows the possibility of setting up direct communication link between two HR-MSs, and then schedules channel measurement between the two HR-MSs if necessary. Based on direction from the HR-BS/HR-RS, the HR-MSs carry out channel measurement and report results to the serving HR-BS/HR-RS. The serving HR-BS/HR-RS decides to set up a direct communication link between the two HR-MSs and sends link creation messages. If necessary, it also helps the two HR-MSs to establish a security association over the direct link. Once the security association is set up, the direct communication link is considered to be set up. The HR-BS/HR-RS may also request the two HR-MSs to report status of a direct communication link and may delete the link when necessary.

After a direct communication link has been set up between two HR-MSs, service flows can be set up between them. The process is initialized by a source HR-MS upon receiving a request for a new flow from its application layer. It then sends a request for service flow addition to the serving HR-BS/HR-RS. The HR-BS/HR-RS negotiates with the destination HR-MS. If the destination HR-MS accepts the request, the HR-BS/HR-RS confirms the addition of the service flow over the direct link with the source HR-MS. Otherwise, the flow shall be set up on uplink in a normal way. The service flows over a direct communication link may be modified or removed upon agreement of the two HR-MSs and the serving HR-BS/HR-RS.

When the degraded HR-BS recovers a connection to the backhaul network, an FBIS connection can be terminated by a procedure of dynamic service deletion and network handover, which enables the FBIS HR-MS to connect back to the previous degraded HR-BS normally.

An inside coverage forwarded HR-MS requests bandwidth as a normal HR-MS. An outside coverage forwarded HR-MS transmits a bandwidth request to its forwarding HR-MS using resource specified by the serving HR-BS/HR-RS.

Synchronization: Synchronization with each other is a must for two HR-MSs involved in direct communication. There are two levels of synchronization:

- Frame-level synchronization — The two HR-MSs share same understanding of frame timing and configuration.
- Symbol-level synchronization — Data transmissions between the two HR-MSs can be received within appropriate timing threshold.

Two HR-MSs within coverage of a common serving HR-BS/HR-RS achieve frame-level synchronization simply by receiving preambles and control messages from the serving HR-BS/HR-RS. To achieve symbol-level synchronization, they transmit ranging sequences to each other under the schedule of the serving HR-BS/HR-RS. During direct data transmission, a transmitting HR-MS times its transmission as if it is a normal uplink transmission, while the receiving HR-MS estimates and adjusts its time offset with the transmitting HR-MS.

If one HR-MS is inside and the other is outside the coverage of a serving HR-BS/HR-RS, the inside coverage HR-MS first achieves frame-level synchronization by receiving preambles and control messages from the serving HR-BS/HR-RS. Subsequently, the inside coverage HR-MS broadcasts preambles and network configuration information for the outside coverage HR-MS to achieve frame-level synchronization. The two HR-MSs achieve symbol-level synchronization by transmitting a ranging signal to each other under the schedule of the serving HR-BS/HR-RS.

HR-MS FORWARDING TO NETWORK

HR-MS forwarding to network allows an HR-MS (forwarding HR-MS) to forward user data and control signaling between an HR-BS/HR-RS and an HR-MS (forwarded HR-MS). The user data and control signaling do not go through higher layer at a forwarding HR-MS. It is applicable when both the forwarding HR-MS and forwarded HR-MS are within coverage of an HR-BS/HR-RS, as well as when only the forwarding HR-MS is in coverage of an HR-BS/HR-RS while the forwarded HR-MS is outside coverage of the HR-BS/HR-RS. With this function, an HR-BS/HR-RS may extend network coverage or provide higher throughput to an inside-of-coverage forwarded HR-MS.

Network Entry of Out-of-Coverage Forwarded HR-MS: An inside-coverage forwarded HR-MS follows normal procedure to enter a network, while an outside-coverage forwarded HR-MS enters network with the help of a forwarding HR-MS. The procedure starts when a serving HR-BS/HR-RS instructs a forwarding HR-MS to extend network coverage by broadcasting preambles and control/configuration information.

When a new outside-coverage HR-MS receives preambles and control/configuration information, it acquires synchronization and uplink transmission parameters. After that it starts initial ranging through the forwarding HR-MS by sending out ranging code. Upon receiving ranging code from a new HR-MS, a forwarding

HR-MS responds accordingly. After receiving a response of SUCCESS, the new HR-MS sends a ranging request to the serving HR-BS/HR-RS via the forwarding HR-MS. The initial ranging process is over when the new HR-MS receives a response to its ranging request from the serving HR-BS/HR-RS via the forwarding HR-MS.

After successful initial ranging, a new HR-MS exchanges control messages with the serving HR-BS/HR-RS via a forwarding HR-MS to complete basic capability negotiation, authorization and key exchange, and registration.

Bandwidth Request of a Forwarded HR-MS:

An inside-coverage forwarded HR-MS requests bandwidth as a normal HR-MS. An outside-coverage forwarded HR-MS transmits a bandwidth request to its forwarding HR-MS using resource specified by the serving HR-BS/HR-RS. A forwarding HR-MS forwards a bandwidth request to its serving HR-BS/HR-RS, and also forwards a bandwidth grant from the serving HR-BS/HR-RS to the forwarded HR-MS.

Synchronization of a Forwarded HR-MS:

In HR-MS forwarding to network, synchronization between the forwarded and forwarding HR-MSs uses the same methods as HR-MS direct communication.

Connection Management Involving a Forwarding HR-MS:

Management and transport between a forwarded HR-MS and a serving HR-BS/HR-RS are defined in the same way as is specified in the basis standards (802.16-2012 and 802.16.1-2012). Each forwarding HR-MS keeps track of all connection identifiers (CID, in 802.16n) or station identifiers (STID, in 802.16.1a) of all forwarded HR-MSs associated with it.

A serving HR-BS/HR-RS assigns resources for data and control signaling transmissions between a forwarded HR-MS and the serving HR-BS/HR-RS. As the corresponding forwarding HR-MS monitors all connection identifiers (in 802.16n) or station identifiers (in 802.16.1a), it is aware of the resource allocation related to any of its forwarded HR-MSs, and then receives the traffic in corresponding resource and forwards the traffic to the serving HR-BS/HR-RS (uplink) or the forwarded HR-MS (downlink). The resource used for traffic forwarding is also allocated by the serving HR-BS/HR-RS.

LOCAL HR-BS/HR-RS FORWARDING

With local forwarding, an HR-MS may communicate with one or more HR-MSs via an HR-BS/HR-RS without going through the backhaul. Opportunity for local forwarding can be found by a serving HR-BS/HR-RS or an upper layer. When a serving HR-BS/HR-RS detects a local forwarding opportunity, it may request permission of the network entities to perform local forwarding.

When local forwarding at a serving HR-BS/HR-RS is allowed, a local forwarding path may be set up during or after data service flow establishment. In the former case, when a serving HR-BS/HR-RS receives a request to set up data service flow from an HR-MS, and if local forwarding is permitted, it establishes uplink from the source HR-MS and downlink to the destina-

It is well known that Japan is a country with frequent disasters such as earthquakes, typhoons, floods, tsunamis, and so on. Therefore, the public wireless communication system plays a very critical and important role in disaster relief in Japan, especially in cases where the general communication systems are not functional.

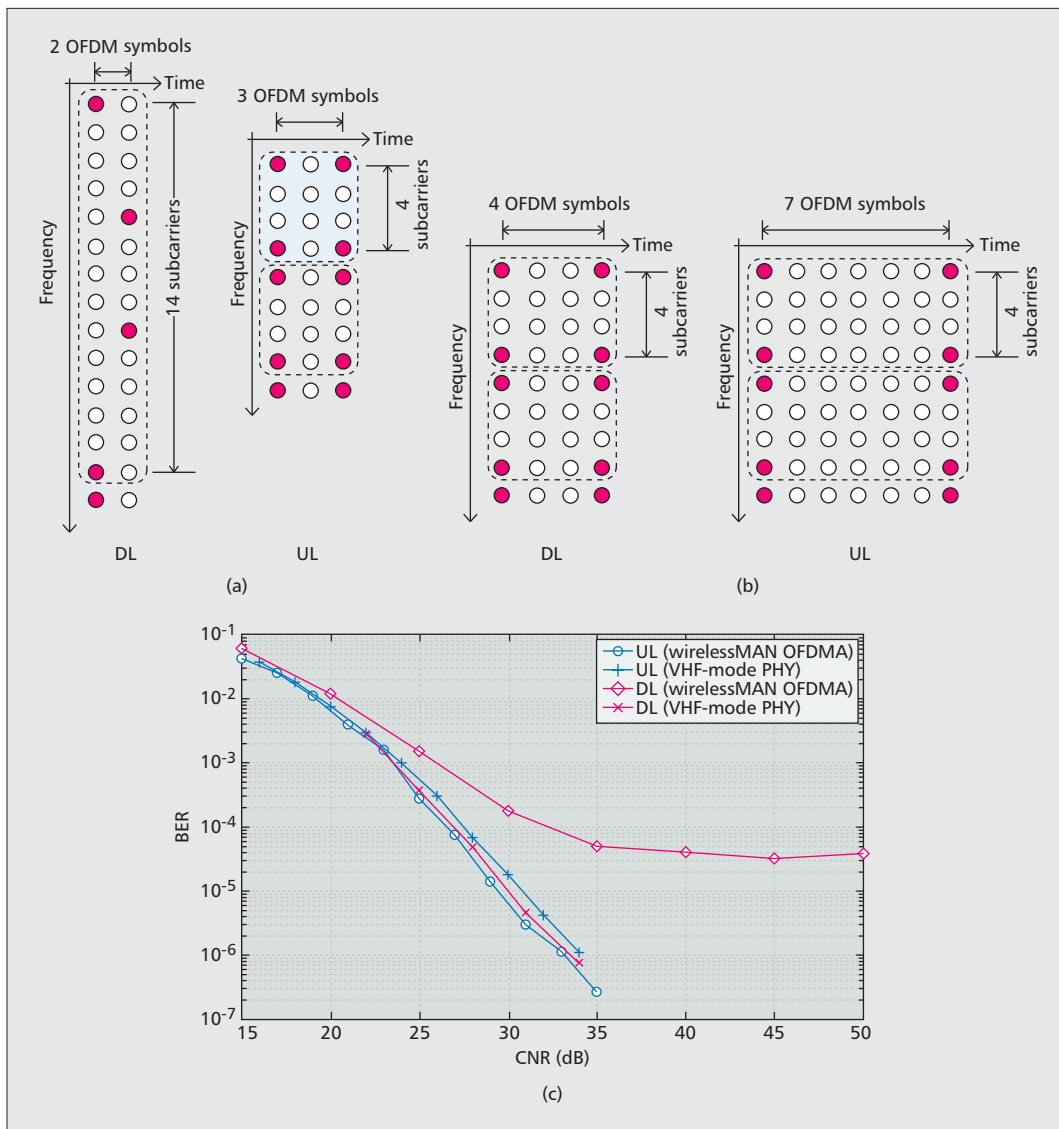


Figure 2. VHF-mode PHY extension of DL cluster and UL tile structure in 802.16n. [7]; a) wirelessMAN-OFDMA; b) VHF-mode PHY; c) comparison of link performance (simulation results).

tion HR-MS(s) by using the dynamic service addition (DSA) procedure individually. In the latter case, a serving HR-BS/HR-RS performs local forwarding directly and may request the network entities to remove the data service flow of the corresponding uplink and downlink. Once a local forwarding path is set up, a serving HR-BS/HR-RS forwards traffic from the source HR-MS to destination HR-MS(s). A local forwarding path is terminated by dynamic service deletion (DSD).

PHY EXTENSION: VHF-MODE PHY

The operation band of 802.16 is extended to 200 MHz by GRIDMAN to cover the wider area that may be needed for PPDR and smart grid [5]. Compared to normal 802.16 operation frequencies above 1 GHz, a VHF-frequency signal may transmit further, therefore experiencing larger delay spread. Meanwhile, PPDR and smart grid may need remote video monitoring, so a GRIDMAN uplink needs to transmit video signal with acceptable quality to a command/control center. As examples, a standard digital video (DV) with

resolution of 720×480 requires a data rate of 2.47 Mb/s; a high-quality DV (HDV) with resolution of 1280×720 and 1440×1080 requires data rates of 6.6 Mb/s and 11.2 Mb/s, respectively.

IEEE 802.16n developed a PHY extension of the WirelessMAN-OFDMA (i.e., orthogonal frequency-division multiple access) air interface, that is, VHF-mode PHY, to support more robust downlink and higher throughput uplink. VHF-mode PHY modifies DL cluster and UL tile structure and the DL:UL symbol ratio of WirelessMAN-OFDMA air interface. As shown in Fig. 2a, a legacy WirelessMAN-OFDMA air interface downlink cluster is of a structure of [14 subcarriers] \times [2 symbols], and each symbol consists of two pilot subcarriers and 12 data subcarriers; an uplink tile is a structure of [4 subcarriers] \times [3 symbols] consisting of four pilot subcarriers, and each corner is a pilot subcarrier. A recent simulation result (Fig. 2c) shows that with the legacy WirelessMAN-OFDMA downlink tile, the downlink bit error rate (BER) cannot achieve 10^{-6} even at a high carrier-to-noise ratio (CNR)

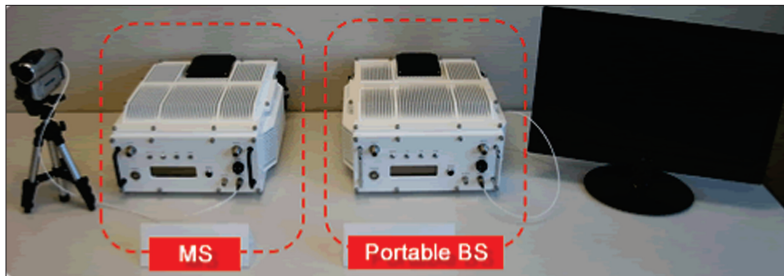


Figure 3. The developed PBB portable BS and MS.

due to limited capability of channel estimation [7]. Moreover, with a bandwidth of 5 MHz and the legacy 802.16 DL:UL symbol ratio of 26:21 specified by the WiMAX Forum, the uplink data rate is 2.9 Mb/s when 16-quadrature amplitude modulation (QAM)-CC3/4 is used; this is only capable of transmitting one standard DV.

With VHF-mode PHY, more pilot subcarriers are inserted in the frequency domain for DL, and fewer pilot symbols are included in the time domain for both DL and UL since VHF signal has smaller Doppler shift. The DL and UL tile structure of the VHF-mode PHY is shown in Fig. 2b, where a downlink cluster has a structure of [4 subcarriers] × [4 symbols], and an uplink tile has a structure of [4 subcarriers] × [7 symbols]. Simulation results (Fig. 2c) shows that with the inserted pilot subcarriers, a VHF-mode PHY DL can achieve a BER of 10^{-6} . Meanwhile, a DL:UL symbol ratio of 9:38 was proposed with VHF-mode PHY, leading to an uplink data rate of 7.4 Mb/s when employing 16QAM-CC3/4, with which three standard DV streams or one HDV stream can be supported. The bandwidth of VHF-mode PHY is 5 MHz, and it employs fast Fourier transform (FFT) size of 1024, which is double of the legacy WirelessMAN-OFDMA air interface. VHF-mode PHY has smaller (half) subcarrier spacing and longer (double) symbol duration, so larger channel spread delay at VHF frequencies can be combated. In the above simulation, the channel model used in simulation is built based on recent measurement at 200 MHz in Japan [8], and the mobility scenario has a speed of 80 km/h.

PRACTICES: JAPANESE PUBLIC BROADBAND SYSTEM

It is well known that Japan is a country with frequent disasters such as earthquakes, typhoons, floods, tsunamis, and so on. Therefore, the public wireless communication system plays a very critical and important role in disaster relief in Japan, especially in cases where the general communication systems are not functional. Hence, applications of a public wireless system include public safety, fire rescue, and so on. Requirements of such a public wireless system include high robustness and reliability, easy and fast deployment, long-range coverage, and support to uplink video transmission so that command headquarters can view the field and direct rescue effectively. In Japan, the legacy public wireless system is narrowband, and only capable

of transmitting voice and low-rate data; thus, it cannot meet the high data rate needs of applications such as video surveillance. In 2011, with the spectrum dividend brought by the digital television transition, a 35 MHz band between 170 and 205 MHz was allocated for public broadband (PBB) applications by the Ministry of Internal Affairs and Communications (MIC) of Japan. In 2013, a PBB wireless communication system was developed by the Smart Wireless Laboratory of NICT and partners, and the system has been in practical use since then [9].

The Japanese PBB system is based on the IEEE 802.16 standard, since this standard supports broadband transmission over distance of several kilometers, and, most importantly, the recently developed GRIDMAN supports higher reliability, heavy uplink data transmission, and 200 MHz frequency band. Selected GRIDMAN protocols have been implemented in the Japanese PBB system. Both BSs and MSs are portable and can be powered by battery so that a standalone network can be set up very quickly. Due to the excellent propagation properties of the 200 MHz band signal, the coverage of a PBB network could be about 10 km. With the equipped GRIDMAN VHF-mode PHY, the uplink data rate could be as high as 7.6 Mb/s when the DL:UL symbol ratio is set to 9:38, and the relatively large delay spread due to long coverage is overcome by the increased pilot subcarriers design. The data rate of 7.6 Mb/s supports the high-quality uplink video surveillance and supervision very well. The developed portable BS and MS are shown in Fig. 3, and the technical specifications of the Japanese PBB system are listed in Table 3.

Performance measurements of the Japanese PBB system were conducted in Numazu City, Japan recently [10]. The measured uplink data rate was 2 Mb/s at a BS-MS distance of 9.3 km. The BS antenna and MS antenna were 57 m and 21 m above sea level, respectively. Transmission power of the BS was 5 W, and for the MS it was automatically controlled by the BS. Both the BS and MS employed a 5-element Yagi-Uda antenna with a maximum gain of 9 dBi. The measurement center frequency was 190 MHz (band of 187.5–192.5 MHz), and QPSK-CC1/2 was employed for uplink transmission.

CONCLUSION

Mechanisms of the main technical developments of the two GRIDMAN standards, IEEE Std 802.16n-2013 and IEEE Std 802.16.1a-2013, are introduced. The new developments greatly improve the reliability and performance of an 802.16 network, and meet the needs for high-speed and high reliability of a number of mission-critical applications such as PPDR and SG. A recently developed PBB system in Japan confirmed the merits of the new developments and showed that GRIDMAN actually meets the requirements set with the use cases in practice.

REFERENCES

- [1] "IEEE Standard for Air Interface for Broadband Wireless Access Systems — Amendment 2: Higher Reliability Networks," 2013.

- [2] "IEEE Standard for WirelessMAN-Advanced Air Interface for Broadband Access Systems — Amendment 2: Higher Reliability Networks, Std., 2013.
- [3] M. Sherman, "Nrr Draft Report," Tech. Rep. IEEE 802.16gman-10/0019r1, 2010.
- [4] L. Lu *et al.*, "Japanese Public Safety Use Case Model for 16n," tech. rep. IEEE C802.16n-10/0067r1, 2011.
- [5] C. Rodine, "Smart Grid Requirements on Man Infrastructure," tech. rep. IEEE C802.16gman-10/0007, 2010.
- [6] "802.16n system Requirements Document Including Sarm Annex," tech. rep. IEEE 802.16n-10/0048r3, 2011.
- [7] H. Harada *et al.*, "A Public Broadband Wireless Communication System on VHF TV Band," *Proc. 6th Int'l ICST Conf. Cognitive Radio Oriented Wireless Networks and Communications*, 2011.
- [8] M. Oodo *et al.*, "Channel Model for Broadband Wireless Communication in the VHF-Band" IEICE tech. rep., 2010.
- [9] "Police-Info Communications," Info-Commun.Bureau, Nat'l. Police Agency, tech. rep., 2013; <http://www.npa.go.jp/joutuu/index.htm>
- [10] M. Oodo and H. Harada, "Current Status of 200 MHz-Band Public Broadband Wireless Communication System in Japan," *Proc. 17th Int'l Symp. Wireless Personal Multimedia Commun.*, 2014.

BIOGRAPHIES

MING-TUO ZHOU [S'01, M'04, SM'11] (zhou.mingtuo@ieee.org) received his B.E. degree from Hunan University in 1997, his M.E. degree from Chongqing University of Posts and Telecommunications in 2000, and his Ph.D. degree from the Asian Institute of Technology in 2003. During January to June 2004, he was a research specialist with the Finland Government Program, Asian Institute of Technology. In July 2004 he joined the National Institute of Information and Communications Technology (NICT) and is now a senior research scientist at the Smart Wireless Laboratory of NICT Singapore Representative Office. He served as co-editor of two books; Co-Chair, committee member, and Finance Chair of a number of international conferences. He was Technical Co-Editor of IEEE 802.16n and IEEE 802.16.1a, a voting member and technical contributor of the IEEE 802.11, 802.15, and 802.16 Working Groups, and is currently a member of the Test and Certification Working Group and Chair of South-East Asia Marketing Subcommittee of the Wi-SUN Alliance. He is Treasurer of the IEEE Vehicular Technology Society (VTS) Chapter, Singapore Section. His research interests include industrial wireless networks, wireless smart utility/grid networks, broadband wireless access, and radio over fiber communications.

MASAYUKI OODO [M] (moodo@nict.go.jp) received B.E., M.E., and D.E. degrees in electrical and electronic engineering from the Tokyo Institute of Technology, Japan, in 1992, 1994, and 1997, respectively. In 1997, he joined the Communications Research Laboratory (CRL), now NICT, Tokyo. During July 2004 and July 2005, he was with the University of York as a visiting research fellow. His research interests are radio propagation, antennas, and physical layer design for broadband wireless communications. He is a member of IEICE of Japan.

VINH DIEN HOANG (hvdien@nict.com.sg) graduated from Hanoi National University, Vietnam, with a B.Sc. degree in mathematics and informatics in 1996, and an M.Sc. degree in mathematics in 1999. In 2002, he received an M.Sc. in communication software and networks from the School of Electrical and Electronics Engineering, Nanyang Technological University, Singapore. He is currently working as a research scientist at the Smart Wireless Laboratory, NICT Singapore Representative Office. He is a Technical Editor of IEEE 1900.7 standards and involved in other IEEE standards such as 802.16. His research interests include cognitive radio and white space communications.

LIRU LU (lirul@ieee.org) received her B.E. (Hons.) degree in electrical and communications engineering from Harbin Institute of Technology, China, in 2000 and her Ph.D. degree in communications engineering from Nanyang Technological University in 2006. She is currently a patent examiner at the Intellectual Property Office of Singapore (IPOS). She was a senior algorithm engineer at Huawei Technologies before joining IPOS. From 2010 to 2014, she was with the NIT Smart Wireless Laboratory situated in Singapore. She had been actively involved in IEEE standardization activities in 802.11, 802.15, 802.16 Working Groups and Task Groups. She served as the Assistant Secretary of the IEEE 802.15.4g/4m Task Group and Ballot Resolution

Item	Specification		
Center frequencies	175, 180, 185, 190, 195, 200 MHz		
Channel bandwidth	5 MHz		
Multiple access/duplex	OFDMA/TDD		
Maximum transmission power	5 W (portable BS, MS); 20 W (fixed BS)		
FFT size	1024		
Subcarrier spacing	5.47 kHz		
Cyclic prefix length	22.9 μ s		
Frame length	10 ms		
OFDM symbol ratio (DL:UL)	37:10	23:24	9:38
Maximum UL rate @16QAM-CC3/4	1.5 Mb/s	4.5 Mb/s	7.6 Mb/s
FPGA	XilinxXC5VLXC330, XC5VLX220		
Battery lifetime	1.5 hours (continuous transmission mode)		

Table 3. Technical specifications of the Japanese PBB system.

Committee. Her previous research areas include the general topics of wireless communications, including wireless PAN, wireless LAN, wireless MAN, and wireless RAN networks with emphasis on modulation and coding.

XIN ZHANG (xin.zhang@ifn.et.tu-dresden.de) received her B.S and Ph.D. degree from the School of Electrical and Electronic Engineering, Nanyang Technological University in 2005 and 2011, respectively. One year before her Ph.D. graduation, she joined the NICT Singapore office. She has been active in IEEE 802 wireless standards for four years. She participated extensively in 802.11ad, 802.22b, and 802.16n standard development. Currently, she is working as a research staff member with the Vodafone Chair at Technische Universität Dresden, Germany. Her research interests include PHY layer system design, beamforming, coding and modulation, and signal processing.

HIROSHI HARADA (harada@nict.go.jp) is a professor at Kyoto University and was director of the Smart Wireless Laboratory of NICT. He joined the Communications Research Laboratory, Ministry of Posts and Communications, in 1995 (currently NICT). Since 1995, he has researched software defined radio, cognitive radio, dynamic spectrum access networks, smart utility networks, and broadband wireless access systems on the VHF, TV white space, microwave, and millimeter-wave bands. He has also joined many standardization committees and forums in the United States as well as in Japan, and has filled important roles for them. He currently serves on the Boards of Directors of the Wireless Innovation Forum (formerly SDR Forum), White Space Alliance, and Wi-SUN Alliance, and has also been Chair of the IEEE DySPAN Standards Committee (formerly IEEE SCC41 and IEEE 1900) since 2009 and Vice Chair of IEEE P1900.4, IEEE P802.15.4g, TIA TR-51, and IEEE P802.15.4m since 2008, 2009, 2011, and 2011, respectively. He was also Chair of the IEICE Technical Committee on Software Radio (TCSR) in 2005–2007 and has been Chair of the Public Broadband Mobile Communication Development Committee, ARIB since in 2010. He is also involved in many other activities related to telecommunications. He has been a visiting professor at the University of Electro-Communications, Tokyo, since 2005 and is the author of *Simulation and Software Radio for Mobile Communications* (Artech House, 2002). He received the Achievement Award from and became a Fellow of IEICE in 2006 and 2009, respectively, and the Achievement Award of ARIB and Funai Prize for Science in 2009 and 2010, respectively, on the topic of cognitive radio research and development.

ADVERTISERS' INDEX

COMPANY	PAGE
CTIA.....	39
ICC 2016 CFP.....	Cover 3
Keysight.....	Cover 2, 1
Keysight Tutorial.....	89
National Instruments.....	3
Siemens Industries.....	Cover 4
Stanford Research.....	5

ADVERTISING SALES OFFICES

Closing date for space reservation: 15th of the month prior to date of issue

<p>NATIONAL SALES OFFICE James A. Vick Sr. Director Advertising Business, IEEE Media EMAIL: jv.ieeemediamedia@ieee.org</p> <p>Marion Delaney Sales Director, IEEE Media EMAIL: md.ieeemediamedia@ieee.org</p> <p>Mark David Sr. Manager Advertising & Business Development EMAIL: m.david@ieee.org</p> <p>Mindy Belfer Advertising Sales Coordinator EMAIL: m.belfer@ieee.org</p> <p>NORTHERN CALIFORNIA George Roman TEL: (702) 515-7247 FAX: (702) 515-7248 EMAIL: George@George.RomanMedia.com</p> <p>SOUTHERN CALIFORNIA Marshall Rubin TEL: (818) 888 2407</p>	<p>FAX:(818) 888-4907 EMAIL: mr.ieeemediamedia@ieee.org</p> <p>MID-ATLANTIC Dawn Becker TEL: (732) 772-0160 FAX: (732) 772-0164 EMAIL: db.ieeemediamedia@ieee.org</p> <p>NORTHEAST Merrie Lynch TEL: (617) 357-8190 FAX: (617) 357-8194 EMAIL: Merrie.Lynch@celassociates2.com</p> <p>Jody Estabrook TEL: (77) 283-4528 FAX: (774) 283-4527 EMAIL: je.ieeemediamedia@ieee.org</p> <p>SOUTHEAST Scott Rickles TEL: (770) 664-4567 FAX: (770) 740-1399 EMAIL: srickles@aol.com</p>	<p>MIDWEST/CENTRAL CANADA Dave Jones TEL: (708) 442-5633 FAX: (708) 442-7620 EMAIL: dj.ieeemediamedia@ieee.org</p> <p>MIDWEST/ONTARIO, CANADA Will Hamilton TEL: (269) 381-2156 FAX: (269) 381-2556 EMAIL: wh.ieeemediamedia@ieee.org</p> <p>TEXAS Ben Skidmore TEL: (972) 587-9064 FAX: (972) 692-8138 EMAIL: ben@partnerspr.com</p> <p>EUROPE Christian Hoelscher TEL: +49 (0) 89 95002778 FAX: +49 (0) 89 95002779 EMAIL: Christian.Hoelscher@husonmedia.com</p>
---	--	---

CURRENTLY SCHEDULED TOPICS

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

www.comsoc.org/commag/call-for-papers



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

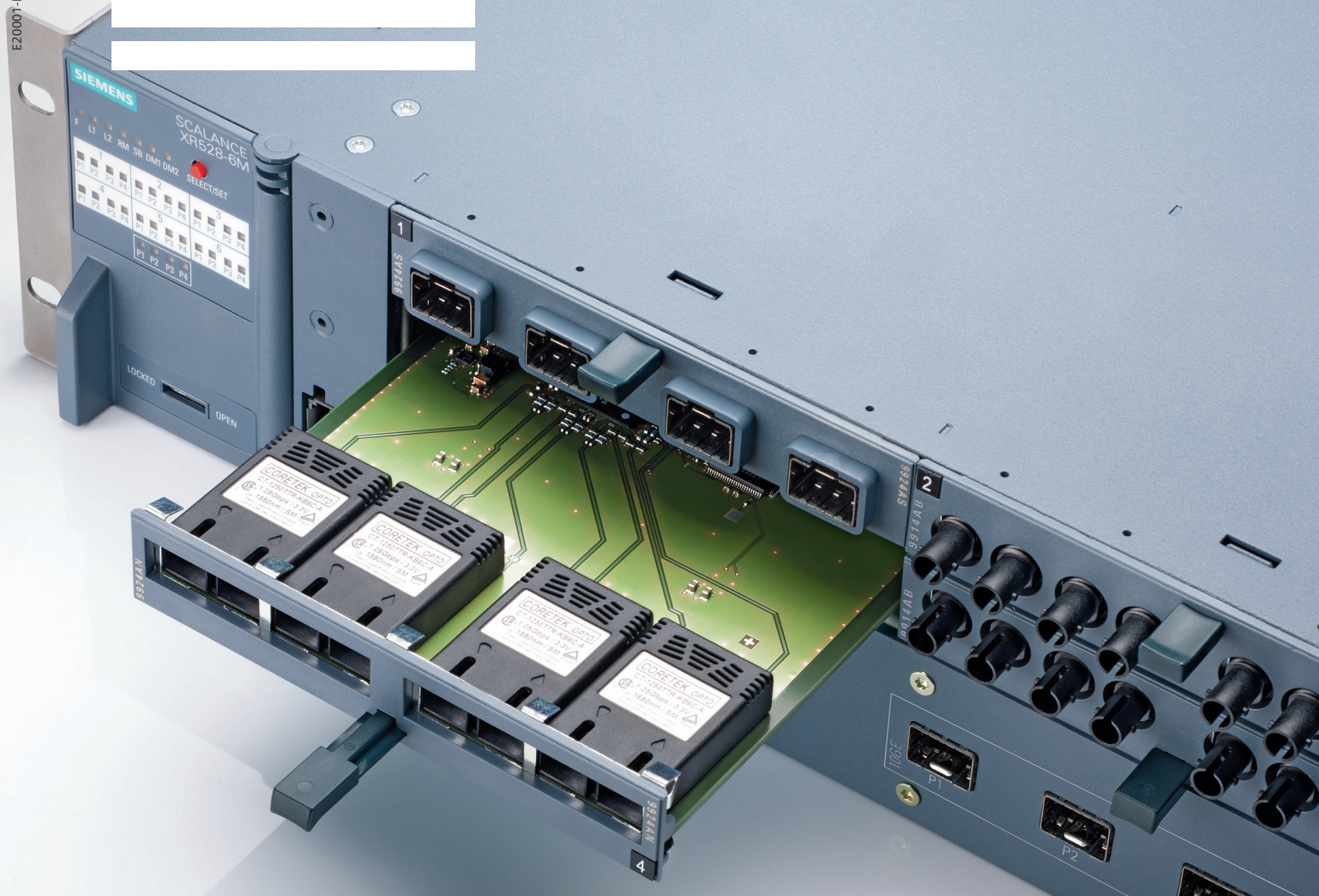
Paper Acceptance Notification:
29 January 2016

Camera-Ready Papers:
29 February 2016

Accepted and presented technical and workshop papers will be published in the IEEE ICC 2016 Conference Proceedings and submitted to IEEE Xplore®. See the website for author requirements of accepted authors. Full details of submission procedures are available at www.ieee-icc.org.

SIEMENS

E20001-F600-P820-X-7600



Automation meets IT

SCALANCE X-500: Two worlds – one switch

SCALANCE X-500 Industrial Ethernet switches ensure seamless transition between the automation and IT networks. They also support industry-specific functions as well as typical IT communication standards.

Whether in typical industrial plants or in industry-related applications: SCALANCE X-500 Industrial Ethernet switches guarantee maximum network availability.

Read the QR code with the QR code reader in your mobile.



[siemens.com/switches](https://www.siemens.com/switches)