



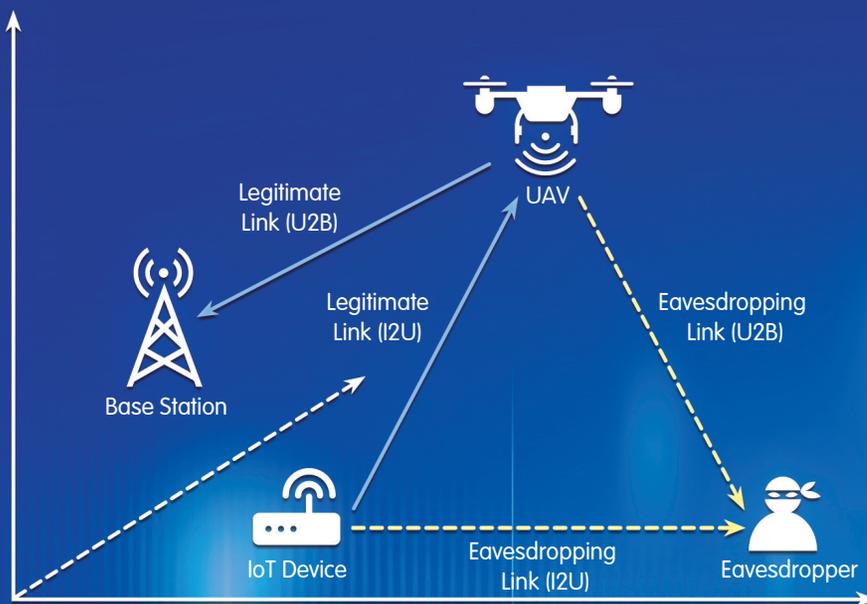
ZTE COMMUNICATIONS

中兴通讯技术(英文版)

<http://tech-en.zte.com.cn>

March 2021, Vol. 19 No. 1

Special Topic: Energy Consumption Challenges and Prospects on B5G Communication Systems



0 3>



9 771673 518215

The 8th Editorial Board of ZTE Communications

Chairman GAO Wen, Peking University (China)
Vice Chairmen XU Ziyang, ZTE Corporation (China) | XU Chengzhong, University of Macau (China)

Members (Surname in Alphabetical Order)

AI Bo	Beijing Jiaotong University (China)
CAO Jiannong	Hong Kong Polytechnic University (China)
CHEN Chang Wen	The State University of New York at Buffalo (USA)
CHEN Yan	Northwestern University (USA)
CHI Nan	Fudan University (China)
CUI Shuguang	UC Davis (USA) and The Chinese University of Hong Kong, Shenzhen (China)
GAO Wen	Peking University (China)
GAO Yang	Nanjing University (China)
GE Xiaohu	Huazhong University of Science and Technology (China)
HWANG Jenq-Neng	University of Washington (USA)
Victor C. M. LEUNG	The University of British Columbia (Canada)
LI Guifang	University of Central Florida (USA)
LI Xiangyang	University of Science and Technology of China (China)
LI Zixue	ZTE Corporation (China)
LIN Xiaodong	ZTE Corporation (China)
LIU Chi	Beijing Institute of Technology (China)
LIU Jian	ZTE Corporation (China)
LIU Ming	Institute of Microelectronics of the Chinese Academy of Sciences (China)
MA Jianhua	Hosei University (Japan)
MA Zheng	Southwest Jiaotong University (China)
NIU Zhisheng	Tsinghua University (China)
PAN Yi	Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences (China)
REN Fuji	Tokushima University (Japan)
REN Kui	Zhejiang University (China)
SHENG Min	Xidian University (China)
SONG Wenzhan	University of Georgia (USA)
SUN Huifang	Mitsubishi Electric Research Laboratories (USA)
SUN Zhili	University of Surrey (UK)
TAO Meixia	Shanghai Jiao Tong University (China)
WANG Haiming	Southeast University (China)
WANG Xiang	ZTE Corporation (China)
WANG Xiaodong	Columbia University (USA)
WANG Xiyu	ZTE Corporation (China)
WANG Yongjin	Nanjing University of Posts and Telecommunications (China)
WANG Zhengdao	Iowa State University (USA)
XU Chengzhong	University of Macau (China)
XU Ziyang	ZTE Corporation (China)
YANG Kun	University of Essex (UK)
YUAN Jinhong	University of New South Wales (Australia)
ZENG Wenjun	Microsoft Research Asia (China)
ZHANG Chengqi	University of Technology Sydney (Australia)
ZHANG Honggang	Zhejiang University (China)
ZHANG Jianhua	Beijing University of Posts and Telecommunications (China)
ZHANG Yueping	Nanyang Technological University (Singapore)
ZHOU Wanlei	City University of Macau (China)
ZHUANG Weihua	University of Waterloo (Canada)

CONTENTS

ZTE COMMUNICATIONS March 2021 Vol. 19 No. 1 (Issue 73)

Special Topic

Energy Consumption Challenges and Prospects on B5G Communication Systems

Editorial 01

GE Xiaohu, YANG Yang

Saving Energy for Wireless Transmission: An Important Revelation from Shannon Formula 02

To develop the energy-saving technologies for future wireless transmissions and networks, this paper presents two basic study points: The multiple events are merged into a single event; The high-order mode is changed to the low-order mode. For this reason, the authors seek that multiple events in wireless transmission links are fused into a single event from Shannon formulas. The authors also analyze the relationship between the information modulation and the error correction, and give a fusion structure of error-corrected modulation. The results of numerical analysis demonstrate the wireless energy saving methods for wireless systems based on Shannon formulas are the achievable efficient schemes.

ZHU Jinkang, ZHAO Ming

Efficient Network Slicing with Dynamic Resource Allocation 11

With the rapid development of wireless network technologies and the growing demand for a high QoS, the effective management of network resources has attracted a lot of attention. Specifically the authors reduce the time consumed for routing by slicing, but the routing success rate after slicing is reduced compared with the unsliced case. In this context, the authors propose a two-stage dynamic network resource allocation framework that first makes decisions on the slices to which flows are assigned, and coordinates resources among slices to ensure comparable routing success rate as in the unsliced case, while taking advantage of the time efficiency gains from slicing.

JI Hong, ZHANG Tianxiang, ZHANG Kai, WANG Wanyuan, WU Weiwei

20 Enabling Energy Efficiency in 5G Network

This paper introduces NR cell switching on/off schemes in 3GPP to achieve energy efficiency in 5G RAN, including intra-system ES scheme and inter-system ES scheme. Additionally, NR architectural features including CU/DU split and dual connectivity are also considered in NR energy saving. How to apply artificial intelligence application in 5G networks is a new topic in 3GPP, and the authors also propose a machine learning based scheme to save energy by switching off the cell selected relying on the load prediction. According to the results of experiments in the real wireless environment, the ML based ES scheme can reduce more power consumption than the conventional ES scheme without load prediction.

LIU Zhuang, GAO Yin, LI Dapeng, CHEN Jiajun, HAN Jiren

30 Cluster Head Selection Algorithm for UAV Assisted Clustered IoT Network Utilizing Blockchain

An unmanned aerial vehicle (UAV) network assisted clustered IoT system is proposed, and a corresponding UAV cluster head (CH) selection algorithm is designed. In this scheme, UAVs are selected as CHs to serve IoT clusters. The proposed CH selection algorithm considers the maximal transmit power, residual energy and distance information of UAVs, which can greatly extend the working life of IoT clusters. Through Monte Carlo simulation, the key performance indexes of the system, including energy consumption, average secrecy rate and the maximal number of data packets received by the base station (BS), are evaluated. The simulation results show that the proposed algorithm has great advantages compared with the existing CH selection algorithms.

LIN Xinhua, ZHANG Jing, LI Qiang

Submission of a manuscript implies that the submitted work has not been published before (except as part of a thesis or lecture note or report or in the form of an abstract); that it is not under consideration for publication elsewhere; that its publication has been approved by all co-authors as well as by the authorities at the institute where the work has been carried out; that, if and when the manuscript is accepted for publication, the authors hand over the transferable copyrights of the accepted manuscript to *ZTE Communications*; and that the manuscript or parts thereof will not be published elsewhere in any language without the consent of the copyright holder. Copyrights include, without spatial or timely limitation, the mechanical, electronic and visual reproduction and distribution; electronic storage and retrieval; and all other forms of electronic publication or any other types of publication including all subsidiary rights.

Responsibility for content rests on authors of signed articles and not on the editorial board of *ZTE Communications* or its sponsors.

All rights reserved.

CONTENTS

ZTE COMMUNICATIONS March 2021 Vol. 19 No. 1 (Issue 73)

Green Air-Ground Integrated Heterogeneous Network in 6G Era 39

The integration of aerial network and terrestrial network has been an inevitable paradigm in the 6G era. However, energy-efficient communications and networking among aerial network and terrestrial network face great challenges. This paper is dedicated to discussing green communications of the air-ground integrated heterogeneous network. The authors first provide a brief introduction to the characteristics of AGIHN in 6G networks, and then analyze the challenges of green AGIHN from the aspects of green terrestrial networks and green aerial networks. Finally, several solutions and key technologies to the green AGIHN are discussed.

WU Huici, LI Hanjie, TAO Xiaofeng

Kinetic Energy Harvesting Toward Battery-Free IoT: Fundamentals, Co-Design Necessity and Prospects 48

This paper gives a brief introduction to the configurations and basic principles of practical Kinetic energy harvesting IoT systems, including their mechanical, electrical, and computing parts. Although there are already a few commercial products in some specific application markets, the understanding and practice in the co-design and optimization of a single KEH-IoT device are far from mature, let alone the conceived multiagent energy-autonomous intelligent systems. Future research and development of the KEH-IoT system beckons for more exchange and collaboration among mechanical, electrical, and computer engineers toward general design guidelines to cope with these interdisciplinary engineering problems.

LIANG Junrui, LI Xin, YANG Hailiang

Review

Next Generation Semantic and Spatial Joint Perception—Neural Metric-Semantic Understanding 61

The author attempts to summarize the recent trends and applications of neural metric-semantic understanding. Starting with an overview of the underlying computer vision and machine learning concepts, he discusses critical aspects of such perception approaches. Specifically, the empha-

sis is on fully leveraging the joint semantic and 3D information. Later on, many important applications of such perception capability such as novel view synthesis and semantic AR contents manipulation are also presented. Finally, the author concludes with a discussion of the technical implications of such technology under a 5G edge computing scenario.

ZHU Fang

Research Paper

72 Integrating Coarse Granularity Part-Level Features with Supervised Global-Level Features for Person Re-Identification

A robust coarse granularity part-level network for person Re-ID, which extracts robust regional features and integrates supervised global features for pedestrian images is proposed. CGPN gains two-fold benefit toward higher accuracy for person Re-ID. On one hand, CGPN learns to extract effective regional features for pedestrian images. On the other hand, CGPN learns to extract more accurate global features with a supervision strategy. The single model trained on three Re-ID datasets achieves state-of-the-art performances.

CAO Jiahao, MAO Xiaofei, LI Dongfang, ZHENG Qingfang, JIA Xia

82 Adaptability Analysis of Fluctuating Traffic for IP Switching and Optical Switching

This paper establishes a multi-layer network architecture through Clos network model and discusses the impacts of maximum allowable blocking rate and service bandwidth standard deviation on CAPEX of IP network and OTN network to find CAPEX demarcation point in different situations. As simulation results show, when the bandwidth deviation mean rate is 0.3 and the maximum allowable blocking rate is 0.01, the hardware cost of OTN switching will exceed IP switching as the average bandwidth is greater than 6 100 Mbit/s. When the service bandwidth fluctuation is severe, the hardware cost of OTN switching will increase and exceed IP switching as the single port rate is allowed in optical switching. The increasing of maximum allowable blocking rate can decrease the hardware cost of OTN switching. Finally, it is found that Flex Ethernet (FlexE) can be used to decrease CAPEX of OTN switching greatly at this time.

LIAN Meng, GU Rentao, JI Yuefeng, WANG Dajiang, LI Hongbiao

Serial parameters:CN 34-1294/TN*2003*q*16*90*en*P*¥ 20.00*2200*10*2021-03

Statement

This magazine is a free publication for you. If you do not want to receive it in the future, you can send the "TD unsubscribe" mail to magazine@zte.com.cn. We will not send you this magazine again after receiving your email. Thank you for your support.



Editorial: Special Topic on Energy Consumption Challenges and Prospects on B5G Communication Systems



Guest Editor

GE Xiaohu received the Ph.D. degree in communication and information engineering from the Huazhong University of Science and Technology (HUST), China in 2003. He has been working with HUST since November 2005. Prior to that, he was a researcher with Ajou University, South Korea, and the Politecnico di Torino, Italy, from January 2004 to October 2005. He is currently a full professor with the School of Electronic Information and Communications, HUST. He is also an adjunct professor with the Faculty of Engineering and Information Technology, University of Technology Sydney (UTS), Australia. He has authored more than 200 papers in refereed journals and conference proceedings and has been granted 35 patents in China. He is leading several projects funded by the NSFC, China MOST, and industries in China. He is taking part in several international joint projects, such as WINDOW and CROWN sponsored by the EU FP7-PEOPLE-IRSES. His research interests include mobile communications, traffic modeling in wireless networks, green communications, and interference modeling in wireless communications. He was a recipient of the best paper awards from the IEEE GLOBECOM 2010. He serves as an associate editor of *IEEE Transactions on Vehicular Technology* and *IEEE Wireless Communications*. He is a senior member of the IEEE.

The objective of this special issue was to attract high quality research articles on energy consumption challenges and prospects for beyond fifth generation (B5G) communication systems. We have received approximately 10 papers in different areas. The submitted papers were rigorously reviewed, and six papers were finally accepted.

The first paper entitled “Saving Energy for Wireless Transmission: An Important Revelation from Shannon Formula” by ZHU et al. presents two basic study points for wireless saving energy and provides the error-corrected modulation method based on extending the Shannon formulas. The numerical analysis shows that the given error-corrected modulation method greatly improves the energy-saving effect of the traditional method in theory. The second paper entitled “Efficient Network Slicing with Dynamic Resource Allocation” by JI et al. proposes a two-stage dynamic resource allocation framework that first makes decisions on the slices to which flows are assigned, and then coordinates resources adjustment among slices to overcome the resource imbalance. The proposed algorithm is evaluated in simulation environments for hierarchical ring 5G networks. The third paper entitled “Enabling Energy



Guest Editor

YANG Yang is currently a full professor at ShanghaiTech University, China, serving as the master of Kedao College and the director of Shanghai Institute of Fog Computing Technology (SHIFT). He is also an adjunct professor with the Research Center for Network Communication at Peng Cheng Laboratory, China. Before joining ShanghaiTech University, he has held faculty positions at the Chinese University of Hong Kong, Brunel University (UK), University College London (UCL, UK), and SIMIT, CAS (China). His current research interests include fog computing networks, service-oriented collaborative intelligence, wireless sensor networks, IoT applications, and advanced testbeds and experiments. He has published more than 200 papers and filed more than 80 technical patents in these research areas. He has been the chair of the Steering Committee of Asia-Pacific Conference on Communications (APCC) since January 2019. In addition, he is a general co-chair of the IEEE DSP 2018 conference and a TPC vice-chair of the IEEE ICC 2019 conference. He is a fellow of the IEEE.

Efficiency in 5G network” by LIU et al. focuses on the energy efficiency of radio access networks and introduces NR cell switching on/off schemes in 3GPP to achieve energy efficiency in 5G RAN. The proposed scheme is experimented in the real wireless environment, whose power consumption can be reduced significantly. The fourth paper entitled “Cluster Head Selection Algorithm for UAV Assisted Clustered IoT Network Utilizing Blockchain” by LIN et al. proposes a designed unmanned aerial vehicle (UAV) cluster head selection algorithm for UAV networks assisted clustered IoT system. The simulation results show that the proposed algorithm has great advantages compared with the existing cluster head selection algorithms. The fifth paper entitled “Green Air-Ground Integrated Heterogeneous Network in 6G Era” by WU et al. is dedicated to discussing green communications of air-ground integrated heterogeneous network (AGIHN). From the aspects of green terrestrial network and green aerial network, challenges of green AGIHN are analyzed and several promising green techniques which can be employed in AGIHN are discussed. The final paper entitled “Kinetic Energy Harvesting Toward Battery-Free IoT: Fundamentals, Co-Design Necessity and Prospects” by LIANG et al. gives a brief introduction to the configurations and basic principles of practical KEH-IoT systems, including their mechanical, electrical, and computing parts.

We would like to thank all the authors for their valuable contributions. We hope that our readers will enjoy reading the articles and find this special issue helpful to their own research work.

DOI: 10.12142/ZTECOM.202101001

Citation (IEEE Format): X. H. Ge and Y. Yang, “Editorial: special topic on energy consumption challenges and prospects on B5G communication systems,” *ZTE Communications*, vol. 19, no. 1, pp. 1 – 1, Mar. 2020. doi: 10.12142/ZTECOM.202101001.

Saving Energy for Wireless Transmission: An Important Revelation from Shannon Formula



ZHU Jinkang^{1,2}, ZHAO Ming^{1,2}

(1. Key Laboratory of Wireless–Optical Communications, Chinese Academy of Sciences, Hefei 230027, China;
2. School of Information Science and Technology, University of Science and Technology of China, Hefei 230027, China)

Abstract: The reduction of power consumption is important for wireless communications and networks. To develop the energy-saving technologies for future wireless transmissions and networks, this paper presents two basic study points: 1) The multiple events are merged into a single event; 2) the high-order mode is changed to the low-order mode. For this reason, we seek that multiple events in wireless transmission links are fused into a single event from Shannon formulas. We also analyze the relationship between the information modulation and the error correction, and give a fusion structure of error-corrected modulation. The energy-saving performance of the error-corrected modulation method is further analyzed through comparison with the traditional methods of modulation plus error correction. The results of numerical analysis demonstrate the wireless energy saving methods for wireless systems based on Shannon formulas are the achievable efficient schemes.

Keywords: wireless saving energy; extension of Shannon formula; error-corrected modulation; energy-saving performance

DOI: 10.12142/ZTECOM.202101002

<https://kns.cnki.net/kcms/detail/34.1294.TN.20210220.1358.002.html>, published online February 20, 2021

Manuscript received: 2020–12–10

Citation (IEEE Format): J. K. Zhu and M. Zhao, “Saving energy for wireless transmission: an important revelation from Shannon formula,” *ZTE Communications*, vol. 19, no. 1, pp. 02 – 10, Mar. 2021. doi: 10.12142/ZTECOM.202101002.

1 Introduction

1.1 Motivation

The energy saving, also said as the power efficiency of wireless communications, has always been an important aim pursued for wireless communications and networks.

From 3G, 4G to 5G, the energy consumption per information bit has dropped significantly.

However, on the other hand, 5G networks are pursuing extremely high peak rates and cloud network uniform manage-

ment, which requires high power consumption. The failure to basically seek a solution to this problem will seriously affect the operation of 5G and the future B5G/6G development. Therefore, various energy efficient methods have been researched and developed, to optimize and reduce the energy consumption in various links of wireless communications and networks. However, some proposals and methods seem too scattered or specific, and show no systematic support from the basic theory.

In the face of increasing demands for higher transmitting rates and energy, it is necessary to find solutions by seeking enlightenment from expansion and augmentation of the most basic Shannon formula. This is an important issue facing B5G/6G in the future, which is worth studying carefully.

This work was supported by National Natural Science Foundation of China (NSFC) under Grant No. 61631018.

1.2 Related Work

The Green Wireless Conference held in Huangshan, China in 2009 is still vivid in our memory. As an outcome of this conference, Ref. [1] summarized the preliminary research on antenna design, service transmission, network design and energy-saving function design, reflecting the results and thinking on green wireless communication technology at that time. A related R&D research project supported by National Key Basic Research Program of China ("973" Program) was subsequently launched in 2010 and has made outstanding contributions to promoting the development of energy-saving technology for wireless communications in China. Moreover, among the achievements are several representative important papers such as "Cell Zooming for Cost-Efficient Green Cellular Networks"^[2] and "Traffic-Aware Network Planning and Green Operation (TANGO)"^[3]. The international communities have also been studying green wireless communications actively.

Many research papers focusing on energy efficiency and energy saving have been published around the world. These publications concern three main aspects: fundament research, cellular networks and sensor networks.

The basic concepts of energy-efficient communications can be found in Ref. [4], which also summarized some fundamental works and advanced techniques for energy efficiency, including information-theoretic analysis and multiple transmission technologies.

Based on energy conservation and the Shannon capacity theorem, the capacity-power consumption formula was proposed in Ref. [5]. The network spectral efficiency and energy efficiency functions of the cellular network were researched, and the relationship between the power consumption and the spectrum efficiency in the cellular networks was also revealed^[6]. A consumption factor theory to analyze and compare energy efficient design choices for wireless communication networks was presented in Ref. [7]. These approaches provide new methods for analyzing and comparing the power efficiency of communication systems.

For 5G development, the optimization solutions to energy and cost efficiency were investigated for wireless communication systems with a large number of antennas and radio frequency (RF) chains^[8]. The overall power transfer efficiency (PTE) and the energy efficiency (EE) of a wirelessly powered massive multiple input multiple output (MIMO) system were investigated in Ref. [9]. Moreover, a novel quadrature space-frequency index modulation (QSF-IM) scheme was proposed as a promising energy-efficient radio-access technology for 5G wireless systems^[10]. Using dual antenna constellation, the proposed scheme can enhance data rates with no extra cost of energy consumption.

Recently, the energy-saving research on sensor networks has made further progress. A novel inter-cluster routing was proposed in Ref. [11], which simultaneously takes the energy efficiency in both intra-cluster and inter-cluster phases into

account. Moreover, a novel concept of energy efficiency welfare was introduced^[11]. The nonlinear fractional programming for the optimal solution to energy efficiency maximization was presented, based on which a particle-swarm optimization-based solution algorithm was proposed in Ref. [12]. An analytical framework for studying the energy efficiency trade-off of cooperation in sensor networks was presented in Ref. [13]; this trade-off is shown to depend on several parameters such as the received power, processing power and the power amplifier loss. The analytical and numerical results reveal that for small distance separation between the source and destination, direct transmission is more energy efficient than relaying.

The joint research of spectrum efficiency and energy efficiency in wireless communications is also one of the most important topics in the next-generation wireless networking area, which is attracting more and more attention from industry, research, and academia^[14]. In Ref. [15], the energy efficiency and spectrum efficiency in underlay device-to-device (D2D) communications enabled cellular networks were investigated.

In summary, since 2009, the research on improving energy efficiency and energy saving has achieved many results. However, compared with the prediction made by Green Touch's research that the net energy consumption in communications networks would be reduced by up to 98% by 2020 relative to 2010^[16] or that the energy efficiency would be increased by a factor of 1 000 compared to the 2010 level^[17], it is far from being achieved. Therefore, from the enlightenment of expending Shannon formulas, this paper will study the foundation and new methods of energy saving in wireless transmission and network coverage to meet the needs of future B5G/6G development.

1.3 Contributions

Energy saving has always been an important aim pursued from 3G, 4G to 5G. The energy consumption per information bit has dropped significantly. However, it is still far away from the future B5G/6G development requirement.

To develop the energy-saving technologies for future wireless transmissions and networks, this paper presents two basic study points: 1) The multiple events are merged into a single event, or the opposite; 2) The high-order mode is changed to the low-order mode, or the opposite.

Making the joint study of the two points above, we seek that the multiple events are merged into a single event in wireless transmission links, from Shannon channel capacity formula, to obtain a new relationship between the information modulation and the error correction, and give a new method of fusing constellation structures of error correction and modulation. Further, the energy-saving performance of the given fusion structure is analyzed, and compared with traditional method of modulation plus error correction.

The research results indicate the given method of wireless

saving energy with the revelation from the Shannon formula has high energy efficiency.

The remainder of this paper is organized as follows. Section 2 is the problem formulation. Section 3 gives a fusion method of error correction and modulation with revelation from the Shannon channel capacity formula. Section 4 analyzes the energy-saving performance of the given method. Finally, in Section 5, we conclude this paper.

2 Problem Formulation

Facing the future communications, high transmission speed, low energy consumption and short time delay are important requirements that must be met. The historical experience tells us that the solution to major problems must begin from the analysis and demonstration of basic theories.

From the perspective of theoretical analysis of saving energy, the topology structures of two basic study points presented by this paper can be written into two expressions.

The first topology structure is to transform the processing with two or more sub-events into a simple event (or vice versa). It can be expressed as

$$A = \sum_{j=1}^g A_g \Leftrightarrow B, \tag{1}$$

where the g sub-events of event A are turned into event B ; the event B is turned into g sub-events of event A .

The second expression of topology structure is that the high-order event is transformed into the multiple low-order sub-events to improve the energy efficiency (or vice versa), which can be expressed as

$$A^e \Leftrightarrow q(A^{e-1}), \tag{2}$$

where the exponential order $e - 1$ of event A^{e-1} is lower than the exponential order e of A^e , and not as complicated as A^e . q is the coefficient of the parallel lower-order.

Therefore, this paper discusses the mathematic expressions of energy-saving ability of the two topology structures, including the performance evaluation of energy saving.

2.1 Evaluation Function of Power Consumption

When wireless communication event A is considered, such as modulation/demodulation (Mod/Dem), the required power consumption P_a , for achieving transmission capability S_a , can be expressed as

$$P_a = f_a(S_a; Q_a), \tag{3}$$

where Q_a is other resource consumption items required for achieving expected capability S_a . This formula represents the energy consumption to realize the transmission capacity S_a . In general, the unit of power consumption is mW and the unit of

transmission capability is bit.

Given the other resource consumption items Q_a , such as the frequency bandwidth and the time delay, the relationship between the fluctuation in achievable performance and the increase or the decrease in power consumption can be derived by the partial differentiation of the power consumption in Eq. (3) as

$$\frac{\partial P_a}{\partial S_a} = \frac{\partial f_a(S_a; Q_a)}{\partial S_a} \Big|_{Q_a=Q}, \tag{4}$$

where Q is a given value of other resource consumption. This formula represents the energy consumption for one-bit increase of the transmitted information, which is the incremental relationship between energy consumption and information bits.

Therefore, we define the energy-saving evaluation function of event A as

$$\eta_a = \frac{1}{\frac{\partial P_a}{\partial S_a} \Big|_{Q_a=Q}} = \frac{\partial S_a}{\partial P_a} \Big|_{Q_a=Q}, \tag{5}$$

where η_a is the amount of information that can be obtained per added unit of power, and it must be greater than zero. As long as $\eta_a > 1$, the performance improvement will be greater than the increased energy consumption, and it is possible for improving the energy-saving effect. The larger η_a , the greater the energy efficiency, or vice versa.

Obviously, Eqs. (3), (4) and (5) are also suitable for event B .

As in Eq. (1), wireless event A consists of g sub-events and the energy consumption is $P_a = P_{a1} + \dots + P_{ag}$. Then the incremental relationship between energy consumption and information bits is

$$\frac{\partial P_a}{\partial S_a} = \sum_{j=1}^g \frac{\partial f_{a_j}(S_{a_j}; Q_{a_j})}{\partial S_{a_j}} \Big|_{Q_{a_j}=Q}, \tag{6}$$

and the energy-saving evaluation function of event A is rewritten as

$$\eta_a = \sum_j \frac{\partial S_{a_j}}{\partial P_a} \Big|_{Q_{a_j}=Q}. \tag{7}$$

Therefore, in wireless communications, how to seek an achievable technical method to obtain high energy-saving efficiency is an important problem.

2.2 Energy Saving of Combining Multiple Events

Now, we consider to develop a new event (event B), which is synthesized by the g sub-events of event A . Also we will complete the design for selecting event B or original event A depending on the consumed energy P_b . Based on the principle of minimum energy expenditure, Eq. (8) can be used to

choose the best design of a new event according to the principle of consuming less energy.

$$\begin{aligned} &\text{if } P_a > P_b, \text{ choose } B, \\ &\text{if } P_a < P_b, \text{ choose } A. \end{aligned} \tag{8}$$

In fact, the design above is not that simple. For example, are the changes of energy consumption of event A and event B for the change of transmission ability the same? When the transmission capability S of event A and that of event B are the same, the answer to this problem is

$$P_a > P_b \text{ and } \frac{\partial \sum_{i=1}^g P_{a_i}}{\partial S_a} > \frac{\partial P_b}{\partial S_b}, \tag{9}$$

and then we must choose event B , and vice versa.

Therefore, in-depth research is needed to find a better method and effective design for achieving the given wireless event, which facilitates minimizing the energy expenditure.

2.3 Energy Saving of High-Order Event

There is a wireless event with e -order, denoted as A^e , of which the power consumption is P_{a^e} . For the order reduction processing, we transform event A^e into event A^{e-1} , reducing the event from e -order to $(e-1)$ -order. Generally, the energy consumption of event A^{e-1} will be less than that of event A^e , and its performance will also be less than the performance of event A^e .

Thence, we need to confirm how many events A^{e-1} have the same performance with the single event A^e , and carefully study if their power consumption is less than that of the single event A^e . The comparison of the achievable performance and the energy consumption between the high-order event and q low-order events, when the transmission capability S of event A and that of event B are the same, can be expressed as

$$P_{a^e} > P_{a^{e-1}} \text{ and } \frac{\partial P_{a^e}}{\partial S_{a^e}} > q \frac{\partial P_{a^{e-1}}}{\partial S_{a^{e-1}}}, \tag{10}$$

and then, we must choose event B , and vice versa.

Eq. (10) represents that the higher the energy efficiency, the lower the energy consumption required for performance improvement and the better the design.

Thence, the energy-saving issue of wireless communications and networks is divided into two research topics:

1) When an event having multiple sub-events compares with another single event, which one has smaller energy consumption?

2) Comparing a high-order event and multiple low-order sub-events, which one has smaller energy consumption?

Here, we have presented the mathematical expressions of

two types of energy-saving problems. The next sections will show the revelation from the Shannon formula and accordingly provide energy-saving solutions to the problems.

3 Revelation from Shannon Channel Capacity Formula

As is well known, the channel capacity formula of Shannon theory is a very important theoretical foundation of wireless communications. It is also very important for our research on energy saving for wireless transmission links, wireless area coverage, wireless networking, etc.

Here, we discuss energy-saving issues of the wireless transmission link that includes two parts: the error correction coding/decoding (codec) and the modulation/demodulation (Mod/Dem), as shown in **Fig. 1**. This link is a stable transmission flow for a given channel. In this regard, some researchers have made considerable efforts, trying to combine the error correction and the modulation into one event. However, they have not yet obtained good usable results. From the perspective of saving energy, it is worth deep studying.

Therefore, we suggest seeking the inspiration and methods by extending the Shannon formula, study the fusion of the error correction codec and the modulation/demodulation, and analyze the relationship between the information rate and power consumption in the fading channel.

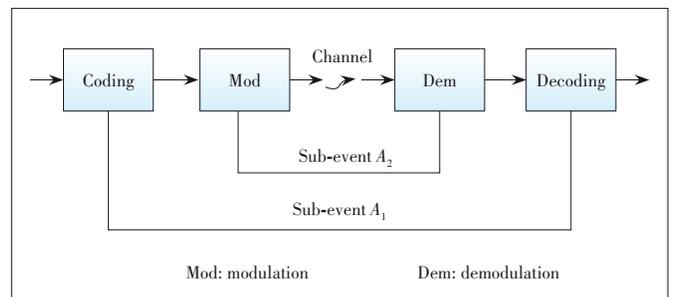
3.1 Fusion of Error Correction and Modulation

If the input signal is $x(t)$ and the output signal is $y(t)$ through the Gaussian fading channel, the characteristic of the Gaussian channel is h , and the channel noise is N_0 , the relationship between input and output is

$$y(t) = hx(t) + N_0. \tag{11}$$

According to the Shannon formula of channel capacity^[18], the wireless transmission capacity $C_{(x,y)}$ can be written as

$$C_{(x,y)} = \text{Max}_{P(x)} H(y) - H(N_0), \tag{12}$$



▲ **Figure 1.** Current wireless transmission link

where $C_{(x,y)}$ is the entropy of the output signal y when the input is x , i.e., the channel capacity; $H(N_0)$ is the lost entropy due to channel noise; $P(x)$ is the statistics function determined by the transmitted signal source $x(t)$. Generally, $\text{Max} H(y)$ is the entropy of the input signal x , i.e. $\text{Max} H(y) = H(x)$. $H(x)$ is an integral from negative infinity to positive infinity, which is unavailable in practical applications.

For this reason, we can define the confidence of the cumulative probability distribution as a reference variable, which is denoted as ω , and get the accurate entropy under the given confidences.

We assume ω is the achievable confidence of the signal $x(t)$, the reliable channel capacity C_ω under the given confidence is the difference of the entropy of input signal $H(x)$ and the entropy of noise $H_{N_0,(1-\omega)}$, which contains the noise entropy and the out-of-confidence discarding entropy. Therefore, the achievable transmission capacity C_ω under the confidence ω is expressed as

$$C_\omega = H_x - H_{N_0,(1-\omega)}. \tag{13}$$

If $H_{N_0,(1-\omega)} > 0$, $C_\omega < H_x$ and $\omega < 1$. It was only when $\omega = 100\%$ and $N_0 = 0$ that we may achieve the lossless capacity, $C_{x,\omega} = H_x$.

When the input signal $x(t)$ has n symbols, (x_1, \dots, x_n) and the Gaussian channel is a normal distributed channel, the mean probability of errors appearing at the Gaussian channel is $\bar{p}(N_{0,i}) = 1/C_i^n$, where C_i^n is the number of combinations of i in n . Then the entropy of N_0 , which causes i error symbols in the output, is expressed as

$$H_{N_0,i} = \log \frac{n!}{i!(n-i)!}. \tag{14}$$

In this way, the reliable transmission capacity C_ω based on the confidence ω can be expressed as

$$C_\omega = n - \log \sum_{i=0}^k C_i^n, C_i^n = \frac{n!}{i!(n-i)!}, \tag{15}$$

where k is the maximum number of the error symbols that can be corrected at the same time.

If there is only one error or error-free in n output symbols, the $n + 1$ symbol combination states in the output will be only received, and the reliable transmission capacity C_ω can be simplified to

$$C_\omega = n - \log_2(1 + n), \text{ for } k = 1. \tag{16}$$

The channel capacity C_ω is the amount of receivable information (denoted as m) transmitted by n symbols. For example, if the signal x has three symbols ($n = 3$), x_1, x_2 and x_3

will be treated as a block, including one information symbol ($m=1$), and input into the Gaussian fading channel. The possible states received are one error-free state, and three states with a single error. The total is four combination output states (Fig. 2).

Here we express the probability of the right symbol and the wrong symbol as p_i and q_i , respectively. If the appearing probabilities of the four output states are all the same, i.e., $p_i = q_i = 1/4$ for $i = 1, 2, 3$, the confidence of this block is $\omega = 3/4 = 75\%$ and the reliable channel capacity of correcting one error is $C_\omega = 3 - \log_2(1+3) = 1$ bit.

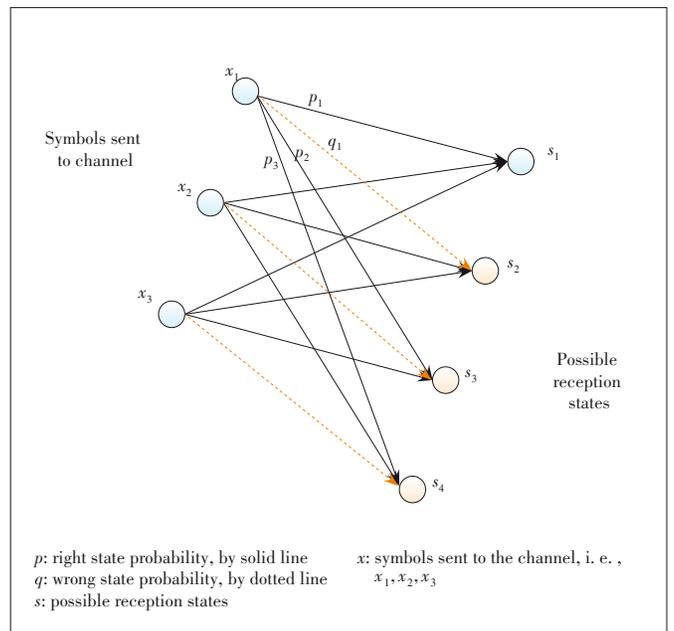
Therefore, based on Eqs. (15) and (16), we can give an error-corrected modulation method. For example, in the 1/3 code block shown in Fig. 2, x_2 is the information bit while the others are check symbols. In this way, one symbol error can be corrected and the transmission efficiency is 1/3. Similarly, we can build error-corrected modulation of 2/5 code, 3/7 code, 4/9 code, 5/11 code, 6/13 code, 7/15 code, etc.

3.2 Constellation Diagram of Error-Corrected Modulation

Based on the above modulation method, we can combine the error-correction function with modulation structure.

The error-corrected modulation method is a constellation modulation structure with the error correction capability.

To construct the constellation diagram of the error-corrected modulation, the processing steps are divided into three parts: 1) planning the constellation point with the information bit plus check bit as a code; 2) choosing a location of the constellation point suitable for transmitting information bits; 3) dividing the constellation area of the correctable error, where the erroneous information bit can be directly detected by the receiver.



▲ Figure 2. Modulation block with one error correction

Fig. 3 shows the (3,1) modulation code with 3 symbols as an example, where one information bit and two check symbols are included and the single error can be corrected. Therefore, the information symbol of the (3,1) code is 0 or 1 and the added error check symbols can be 00, 11, or 01, 10. The modulation coding has one of two structures with no error: (0,0,0) (1,1,1) or (0, 0, 1) (1, 1, 0). This modulation code with the constellation points can correct one error.

4 Analysis of Energy-Saving Efficiency

4.1 Energy Consumption of Two Modulation Methods

Based on the above processing, this section analyzes the power consumption of two structures of the error correction plus modulation and the error-corrected modulation for wireless transmission links, to find which method saves more energy.

First, let us consider the traditional transmission link, in which the error correction codec is event A_1 and the Mod/Dem is event A_2 , to analyze the energy-saving efficiency.

The power consumption of event A_1 can be expressed as

$$P_{a_1} = \left(f_T(S_{a_1}; Q_{a_1}) + f_R(S_{a_1}; Q_{a_1}) \right) \Big|_{Q_{a_1}=1} \approx \left(f_T(S_{a_1}) + f_R(S_{a_1}) \right) P_0. \quad (17)$$

In general, the coding process is addition operation depending on the coding length n , and the power consumption of the coding process can be expressed as the function of i_T -order coding length n . For simplicity, the power consumption of the decoding process can be expressed as the function of i_R -order coding length n . Then, the power consumption of the coding and decoding processing of the (n, m) code can be respectively simplified to

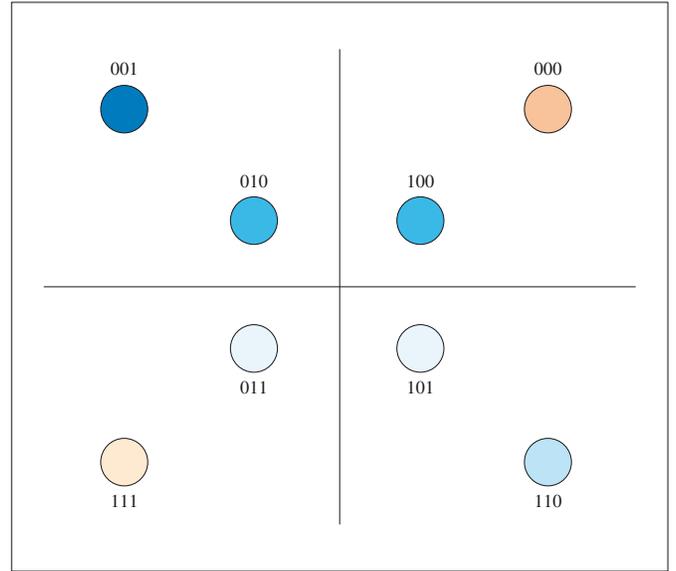
$$f_T(S_{a_1}) \approx \alpha_T(n^{i_T}), \quad f_R(S_{a_1}) \approx \alpha_R(n^{i_R}). \quad (18)$$

Therefore, the total power consumption of the coding and decoding processing of the (n, m) code for the corrected error $k = 1$, i.e., the total of event A_1 , is

$$P_{a_1} \approx \left(\alpha_T(n)^{i_T} + \alpha_R(n)^{i_R} \right) P_0, \quad (19)$$

where the subscript T means the transmitting process, R means the receiving process, and P_0 is the power consumption of a single addition operation ($i = 1$) of one symbol. Moreover, $i = 1$ means the addition operation, and $i = 2$ and $i = 3$ are respectively the multiplication operation and the convolution or iteration operation.

Second, event A_2 is the n -order quadrature amplitude modulation (QAM) modulation. By the Shannon theory, the transmitted signal symbol rate in unit bandwidth and unit time is



▲ **Figure 3. Error-corrected modulation constellation of (3, 1) coding**

$$n = \log \left(1 + \frac{P_{a_2}}{N_0} \right) \approx \log \frac{P_{a_2}}{N_0}, \quad P_{a_2} \approx 2^n N_0, \quad (20)$$

where n is the number of transmitted signal symbols in a block. The receiving demodulation of event A_2 is similar to code demodulation, with only multiplication and comparison; the power consumption can be expressed as $(\beta_R n^{i_R}) P_0$.

Therefore, the power consumption of Mod/Dem with n symbols is

$$P_{a_2} = f_T(2^n(N_0)) + f_R(n^r) \approx (\beta_T 2^n) N_0 + (\beta_R n^r) P_0, \quad (21)$$

where N_0 is the power of channel noise; $f_R(n^r)$ is the power consumption of the receiver, which is proportional to the code length n .

Therefore, the total power consumption of coding/decoding plus Mod/Dem can be expressed as

$$P_a = P_{a_1} + P_{a_2} \approx \left(\alpha_T n^{i_T} + \alpha_R n^{i_R} + \beta_R n^r \right) P_0 + (\beta_T 2^n) N_0. \quad (22)$$

When $i_T = i_R = r = 2$, all power operations in P_0 item of Eq. (22) are multiplication operation. The power consumption of the operation of one symbol is denoted as P_0 , and is equal to the unit noise power N_0 . If $P_0 = N_0$, Eq. (22) can be simplified as

$$P_a \approx N_0 \left((\alpha_T + \alpha_R + \beta_R) n^2 + \beta_T 2^n \right). \quad (23)$$

Fig. 4 shows the power consumptions of event A under different sending/receiving parameters with $k = 1$ and $N_0 = P_0 = 1 \mu\text{W}$. Obviously, the power consumption increases exponentially as n increases. With the increase of the sending/

receiving parameters, the power consumption also increases significantly.

Fig. 4 demonstrates that for every additional bit of information in error correction coding, from m to $m+1$, the code length must be increased by two symbols at least, from n to $n+2$, that is $m=(n-1)/2$. Therefore, (3, 1) code, (5, 2) code, (7, 3) code, (9, 4) code, (11, 5) code, (13, 6) code, (15, 7) code, etc. are all such coding.

Based on Eq. (5) in Section 2, the evaluation function of power consumption of event A can be expressed as

$$\eta_a = \frac{\partial S_a}{\partial P_a} \approx \frac{1/N_0}{(\alpha_T + \alpha_R + \beta_R)((n+2)^2 - n^2) + \beta_T(2^{n+2} - 2^n)} \cdot (24)$$

If $\alpha_T = \beta_T = 1$ and $\alpha_R = \beta_R = 3$, the evaluation function of power consumption of event A can be expressed as

$$\eta_a \approx \frac{1/N_0}{7((n+2)^2 - n^2) + (2^{n+2} - 2^n)} \cdot (25)$$

Based on different lengths of symbol blocks and information bits (n, m), the evaluation function of the power consumption for error correction and modulation separation in the traditional transmission link is shown in Fig. 5.

Now, let us consider the fusion structure of error-corrected modulation of event B . Only Mod/Dem processing is taken for C_j^n combination states shown in Eq. (15), where $j = 0, 1, \dots, k$ (k is the number of error correction symbols of a block). Then, the power consumption of event B can be simplify as

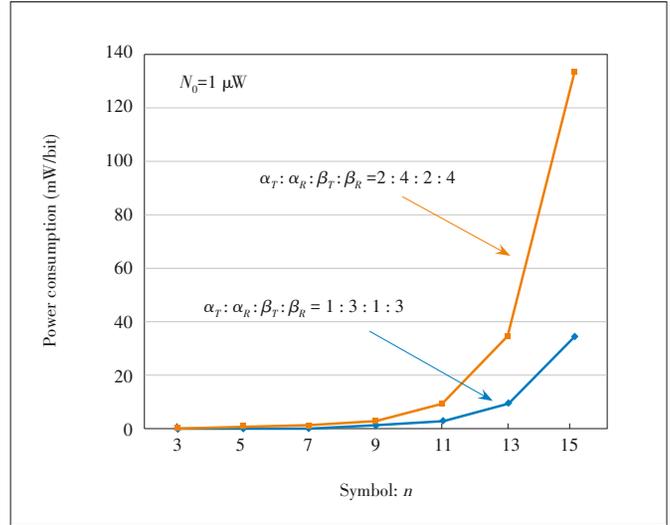
$$P_b \approx \beta_T(N_0) \sum_{j=0}^{k-1} C_j^n 2^m + \beta_R(n)^i P_0 \approx N_0(\beta_R n^2 + \beta_T(1+n)2^m), (26)$$

where $i = 2$ and $P_0 = N_0$.

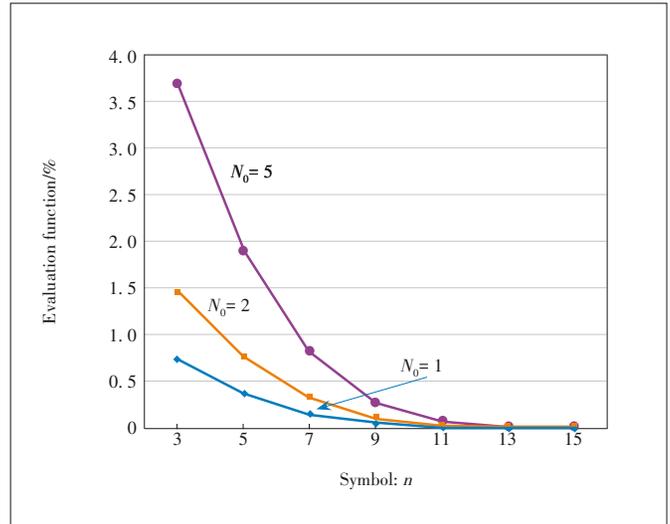
The power consumptions of event B under different sending/receiving parameters are shown in Fig. 6, where the parameters are $\alpha_T, \beta_T, \alpha_R,$ and $\beta_R, k = 1,$ and $N_0 = P_0 = 1 \mu\text{W}$. Obviously, the power consumption increases exponentially with m . Along with the increase of the sending/receiving parameters, the power consumption also increases. However, compared with event A shown in Fig. 4, the power consumption is significantly reduced.

Similar to Eq. (25), the evaluation function of the power consumption of event B can be expressed as

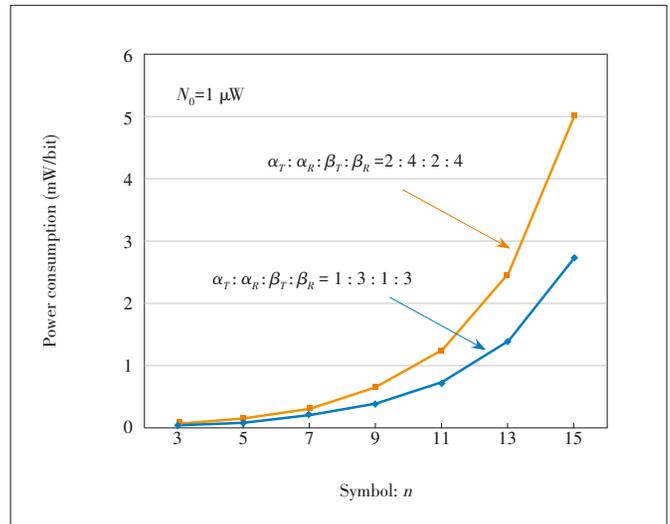
$$\eta_b \approx \frac{1/N_0}{(\beta_R)((n+2)^2 - n^2) + \beta_T((3+n)2^{m+1} - (1+n)2^m)} \cdot (27)$$



▲ Figure 4. Power consumption of event A under different parameters



▲ Figure 5. Energy-saving evaluation function of the traditional link



▲ Figure 6. Power consumption of event A under different parameters

When $\alpha_T = \beta_T = 1$ and $\alpha_R = \beta_R = 3$, the evaluation function of power consumption of event B can be simplified as

$$\eta_b \approx \frac{1/N_0}{3((n+2)^2 - n^2) + ((3+n)2^{m+1} - (1+n)2^m)} \quad (28)$$

Based on different length of symbol blocks and information bits (n, m), the evaluation function of power consumption of the given modulation method for event B is shown in **Fig. 7**.

4.2 Energy-Saving Comparison of Two Modulation Methods

According to the above analysis, the modulation link of (n, m) code structure is divided into two modes, event A and event B . The amount of information transmitted in a block with a coding length n is m , that is, the transmission rate is m/n . Then, the power consumption of event A is $P_a \approx N_0((\alpha_T + \alpha_R + \beta_R)n^2 + \beta_T 2^n)$ and that of event B is $P_b \approx N_0(\beta_R n^2 + \beta_T(1+n)2^m)$.

Therefore, the improved energy-saving degree from event A to event B at the same information rate m/n is defined as

$$\mu_{B/A} = \frac{\eta_b - \eta_a}{\eta_a} = \frac{(\alpha_T + \alpha_R + \beta_R)((n+2)^2 - n^2) + \beta_T(2^{n+2} - 2^n)}{(\beta_R)((n+2)^2 - n^2) + \beta_T((3+n)2^{m+1} - (1+n)2^m)} - 1. \quad (29)$$

When $\alpha_T = \beta_T = 1$ and $\alpha_R = \beta_R = 3$, the improved energy-saving degree of conversion of event A into event B is

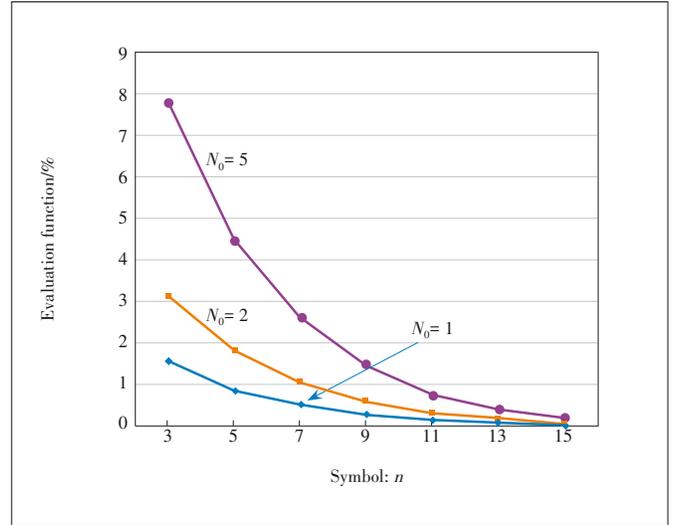
$$\mu_{B/A} = \frac{7((n+2)^2 - n^2) + (2^{n+2} - 2^n)}{3((n+2)^2 - n^2) + ((3+n)2^{m+1} - (1+n)2^m)} - 1. \quad (30)$$

The improved degree of energy saving is shown in **Fig. 8**, which demonstrates that the longer the code length, the higher the improved degree in energy saving. If the length of coding is 15, the improved energy-saving degree reaches up to 35% in theory.

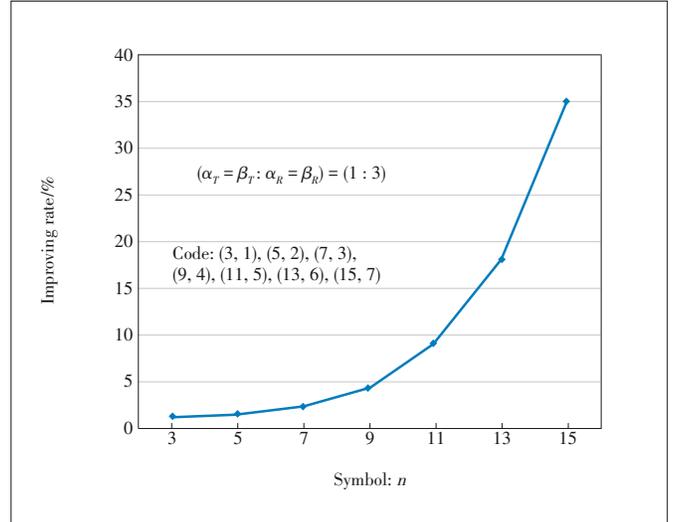
5 Conclusions

With the widespread deployment and application of 5G networks, the requirements for wireless energy saving are getting higher and higher. This paper introduces two basic study points for wireless energy saving and gives the error-corrected modulation method and its fusing constellation structure based on extending the Shannon formulas.

This paper also analyzes and compares the energy-saving performance of two wireless transmission chains, the traditional and the proposed. The numerical analysis shows that the pro-



▲ **Figure 7.** Energy-saving evaluation function of the proposed modulation mode



▲ **Figure 8.** Error-corrected modulation constellation of (3,1) coding

posed error-corrected modulation method improves the energy-saving effect of the traditional method by 35% in theory.

References

- [1] YOU X H, WANG J, ZHANG P, et al. Study and ideas for green wireless mobile communications [J]. Journal of university of science and technology of China, 2009, 39(10): 1009 - 1015
- [2] NIU Z S, WU Y Q, GONG J, et al. Cell zooming for cost-efficient green cellular networks [J]. IEEE communications magazine, 2010, 48(11): 74 - 79. DOI: 10.1109/mcom.2010.5621970
- [3] NIU Z S. TANGO: traffic-aware network planning and green operation [J]. IEEE wireless communications, 2011, 18(5): 25 - 29. DOI: 10.1109/mwc.2011.6056689
- [4] LI G, XU Z K, XIONG C, et al. Energy-efficient wireless communications: tutorial, survey, and open issues [J]. IEEE wireless communications, 2011, 18

- (6): 28 – 35. DOI: 10.1109/mwc.2011.6108331
- [5] ZHU J K. Capacity-power consumption and energy-efficiency evaluation of green wireless networks [J]. China communications, 2012, 9(2): 13 – 21
- [6] ZHU J K, XU L. Spectrum-efficiency and energy-efficiency functions of green cellular networks [J]. Journal of communications, 2013, 34(1): 1 – 7
- [7] MURDOCK J N, RAPPAPORT T S. Consumption factor and power-efficiency factor: a theory for evaluating the energy efficiency of cascaded communication systems [J]. IEEE journal on selected areas in communications, 2014, 32(2): 221 – 236. DOI: 10.1109/jsac.2014.141204
- [8] ZI R, GE X H, THOMPSON J, et al. Energy efficiency optimization of 5G radio frequency chain systems [J]. IEEE journal on selected areas in communications, 2016, 34(4): 758 – 771. DOI: 10.1109/jsac.2016.2544579
- [9] KHAN T A, YAZDAN A, HEATH R W. Optimization of power transfer efficiency and energy efficiency for wireless-powered systems with massive MIMO [J]. IEEE transactions on wireless communications, 2018, 17(11): 7159 – 7172. DOI: 10.1109/twc.2018.2865727
- [10] PATCHARAMANEEPAKORN P, WANG C X, FU Y, et al. Quadrature space-frequency index modulation for energy-efficient 5G wireless communication systems [J]. IEEE transactions on communications, 2018, 66(7): 3050 – 3064. DOI: 10.1109/tcomm.2017.2776956
- [11] LIN D Y, MIN W D, XU J F. An energy-saving routing integrated economic theory with compressive sensing to extend the lifespan of WSNs [J]. IEEE internet of things journal, 2020, 7(8): 7636 – 7647. DOI: 10.1109/ijot.2020.2987354
- [12] MIN S, MENG Z. Energy efficiency optimization for wireless powered sensor networks with nonorthogonal multiple access [J]. IEEE sensors letters, 2018, 2(1): 1 – 4. DOI: 10.1109/lsens.2018.2792454
- [13] SADEK A K, YU W, LIU K J R. On the energy efficiency of cooperative communications in wireless sensor networks [J]. ACM transactions on sensor networks, 2009, 6(1): 1 – 21. DOI: 10.1145/1653760.1653765
- [14] YI Q. Spectrum efficiency and energy efficiency in wireless communication networks [J]. IEEE wireless communications, 2020, 27(5): 2 – 3. DOI: 10.1109/mwc.2020.9241874
- [15] CAI Y, NI Y, ZHANG J, et al. Energy efficiency and spectrum efficiency in underlay device-to-device communications enabled cellular networks [J]. China communications, 2019, 16(4): 16 – 34
- [16] GreenTouch. Reducing the net energy consumption in communications networks by up to 98% by 2020 [R]. Murray Hill, USA: GreenTouch Consortium, 2015
- [17] ELMIRGHANI J M H, KLEIN T, HINTON K, et al. GreenTouch GreenMeter core network energy-efficiency improvement measures and optimization [J]. Journal of optical communications and networking, 2018, 10(2): A250 – A269
- [18] SHANNON C E. A mathematical theory of communication [J]. Bell system technical journal, 1948, 27(3): 379 – 423. DOI: 10.1002/j.1538-7305.1948.tb01338.x

Biographies

ZHU Jinkang (jkzhu@ustc.edu.cn) received the B.E. degree from Sichuan University, China in 1966. He has been a professor of University of Science and Technology of China (USTC) since 1992 and committed to research on wireless mobile communications and networks. He was the leader of the Personal Communication Group of the Communication Subject Expert Group of the National High-Tech Development Program, China. His current research focus is the green wireless communications, wireless big data and wireless AI, emerging theory and technology of wireless communications.

ZHAO Ming received the B.E., M.E. and Ph.D. degrees in electronic engineering and information science from the University of Science and Technology of China (USTC) in 1999, 2002 and 2018, respectively. Now he is an associate professor with the Department of Electronic Engineering and Information Science, USTC. His research interests include non-orthogonal multiple access, heterogeneous networks and green communications.

Efficient Network Slicing with Dynamic Resource Allocation



JI Hong¹, ZHANG Tianxiang², ZHANG Kai¹, WANG Wanyuan¹, WU Weiwei¹

(1. School of Computer Science and Engineering, Southeast University, Nanjing 211189, China;
2. ZTE corporation, Shenzhen 518057, China)

Abstract: With the rapid development of wireless network technologies and the growing demand for a high quality of service (QoS), the effective management of network resources has attracted a lot of attention. For example, in a practical scenario, when a network shock occurs, a batch of affected flows needs to be rerouted to respond to the network shock to bring the entire network deployment back to the optimal state, and in the process of rerouting a batch of flows, the entire response time needs to be as short as possible. Specifically, we reduce the time consumed for routing by slicing, but the routing success rate after slicing is reduced compared with the unsliced case. In this context, we propose a two-stage dynamic network resource allocation framework that first makes decisions on the slices to which flows are assigned, and coordinates resources among slices to ensure a comparable routing success rate as in the unsliced case, while taking advantage of the time efficiency gains from slicing.

Keywords: network slicing; dynamic resource allocation; reinforcement learning

DOI: 10.12142/ZTECOM.202101003

<https://kns.cnki.net/kcms/detail/34.1294.TN.20210208.1133.002.html>, published online February 8, 2021

Manuscript received: 2020-12-09

Citation (IEEE Format): H. Ji, T. X. Zhang, K. Zhang, et al. "Efficient network slicing with dynamic resource allocation," *ZTE Communications*, vol. 19, no. 1, pp. 11 - 19, Mar. 2021. doi: 10.12142/ZTECOM. 202101003.

1 Introduction

With the fast growth of network technologies, such as 5G and data center networks, and ever-increasing demand for the high quality of service (QoS)^[1], efficient employment of existing network infrastructure becomes a challenging task. Network slicing provides an effective method that can introduce flexibility and faster resource deployment in network resource management^[2]. A slice is a horizontal subset of the entire network which is set to satisfy resource requests, for example, bandwidth for flows. A flow is a certain amount of bandwidth requirement on the passing links^[3].

Many studies have been devoted to making full use of QoS and network resource utilization for traffic scheduling, for

which queue management and scheduling are widely used. Generally, the flows in the queue may have different priorities, for example, the preceding flows may have a higher priority than other flows^[4]. In floodlight-based software defined networks (SDN)^[5], queue-based scheduling technology is widely used to implement QoS support. In the above research, the traffic to be transmitted is divided into QoS flows and optimal flows, and assigned to the queue according to their priorities^[6].

However, the queuing nature of the above resource allocation or scheduling strategy has shortcomings. It is necessary to calculate the feasible route of the flow/request in real-time by checking the remaining available bandwidth on the network link. However, it is time-consuming to calculate the feasible route for the flows in the queue which need to be processed in

sequence. Inspired by network slicing, which has revealed an acceleration effect on routing^[7], we introduce a network slicing model to parallel and speed up routing calculations. In the proposed network slicing model, different slices (e.g., horizontal slices) share the same topology, but have different bandwidth resources on the link. Flows are allocated to different slices, and routed in an independent manner in each slice. Finally, the routing process can be calculated in parallel, thus speeding up the routing process.

The main challenge of network slicing is to determine to which slice the request should be deployed. The objective of resource allocation is to ensure a high deployment success rate, for which some flows may not be deployed due to scarce resources. Moreover, using the checking mechanism to test whether the route is feasible or not will degrade the efficiency of parallel routing. Against this background, the first contribution of this paper is that we propose a two-stage resource allocation algorithm, which consists allocation to slice and the resource adjustment among slices. In the first stage of the flow allocation, given the current arrival flows, a reinforcement learning (RL) based mechanism is proposed to deploy these flows. In the second stage of resource adjustment, a dynamic and parallel network resource reallocation among slices is proposed. Another contribution is that we conduct an experimental evaluation in a hierarchical ring 5G network similar to a real large-scale network. Simulation results show that this method can still reduce the routing calculation time and maintain the deployment success rate when dealing with large-scale networks.

The rest of this paper is organized as follows. Section 2 reviews the related work on resource allocation of network technologies. Section 3 describes the problem and forms the model. We propose the resource allocation algorithm in Section 4. Section 5 presents the results of the experiments. Finally, we conclude this paper in Section 6.

2 Related Works

In recent years, network slicing and effective use of network resources have received a lot of attention. SDN is the candidate technology for implementing network slicing on common network infrastructure for deploying a number of services with different requirements. Deployed slices guarantee the isolation in terms of connectivity and performance^[8].

Network resource virtualization (slicing) has been regarded as one of the main development trends of 5G cellular networks and data center networks, which can improve QoS, quality of experience (QoE), and network resource utilization^[9]. Through network slicing, the whole network can easily accommodate the different needs of divergent service types, applications, and services in support of vertical industries^[10]. A virtualized infrastructure is provided in Ref. [11], which is a network infrastructure in which some or all of the elements are virtual-

ized by the data center network, such as servers, links, switches and topology. The cloud computing platform consists of single or multiple virtualized infrastructure, which relies on virtualization technology to divide available resources and share them among users.

For the limitation of network resources, a resource allocation plan can be implemented to improve communication reliability and network utilization. However, slicing may cause severe network performance degradation. ZHANG et al.^[12] use a supply-demand model to quantify slicing interference. In order to maximize the total throughput of accepted requests, an adaptive interference awareness (AIA) heuristic method is proposed to automatically place slices in network slices customized for 5G services. CHEN et al.^[13] also develop a dynamic network slice resource scheduling and management method based on SDN to meet the services' requirements with time-varying characteristics. A resource combination and allocation mechanism is proposed to maximize the total rate of the virtualized network based on its channel state information^[14]. An algorithm based on iterative slice provision has been proposed, which adjusts the minimum slice requirement based on channel state information, but does not consider the global resource utilization rate of the network and slice priority. A centralized joint power and resource allocation scheme for priority multi-layer cellular networks has been proposed^[15], which works by allowing users with higher priority to maximize the number of service users. Priority is only considered at the user level, and different priorities between slices are ignored.

In this paper, we formulate the network slicing model, which enables flow requests to be dispatched among different slices, thereby speeding up routing computations. Moreover, the proposed algorithm is validated to have the advantage of maximizing the success rate of flow deployment.

3 Problem Formulation

3.1 Network Formulation

We model the wide area network (WAN) as an undirected graph $G = \{V, E, A\}$, where $V = \{v_1, v_2, \dots, v_n\}$ is a set of nodes which can represent switches or routers in the WAN, and $E = \{(v_s, v_d) | i = 1, 2, \dots, m\}$ is a set of edges in the graph which can represent link connections between WAN nodes. Let $A = \{a_1, a_2, \dots, a_m\}$ represent the available bandwidth per link. We use the abbreviation e_{sd} or abbreviation of edge number e_i , $i = 1, 2, \dots, m$, to denote an edge between v_s and v_d . We use a_{sd} or a_i to denote the available bandwidth of a specific link. The notations to be used in this section can be found in **Table 1**.

Define a flow f_j as users' requests, indicating that at each time period, a user requests certain channels of a pre-defined transmission rate (bandwidth demand). Let $F = \{(v_s, v_d) | j = 1, 2, \dots, h\}$ denote a flow set, and d_j denote the pre-defined

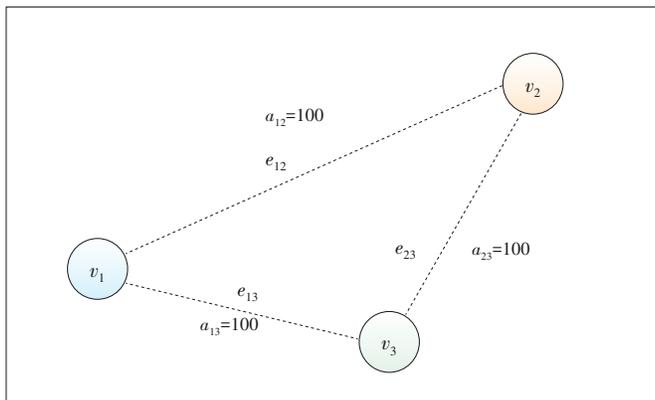
transmission rate of a flow. Each flow can be represented by a path that consists of a series of links, which is generated by the shortest path (SP) routing scheme. Note that a feasible SP route scheme depends on the current network state especially each link's available bandwidth. Let the path calculated by a certain routing scheme for f_j as P_j . After determining a routing scheme, we can get all links in P_j , so P_j also can be represented as a set of links. These links should also satisfy the constraints described above: $d_j \leq a_i, \forall e_i \in P_j$.

3.2 Network Slicing

Assuming we have a real network and several flows, after allocating some bandwidth resources to a flow, the network state will be changed, so the next flow's routing scheme must depend on the changed network state. If they are planned together, some conflicts may happen, for the reason that, to deploy a bunch of flows, the shortest path should be calculated sequentially to avoid deployment failures of following flows. For example, there are two flows, namely f_1 and f_2 , established on a network in **Fig. 1**, where $f_1 = \{v_1, v_3, 80\}$ and $f_2 = \{v_1, v_3, 60\}$ if their shortest paths are calculated at the same time. Both of them will have the same path $v_1 \rightarrow v_2 \rightarrow v_3$, but the link be-

▼ **Table 1. Notation overview**

Notation	Description
$V = \{v_1, \dots, v_n\}$	the set of nodes
$E = \{(v_i, v_j) i = 1, 2, \dots, m\}$	the set of edges
e_i or e_{sd}	a certain link
$A = \{a_1, \dots, a_m\}$	available bandwidth per link
$G = \{V, E, A\}$	the set of a network
$F = \{(v_i, v_j) j = 1, \dots, h\}$	the set of flows
f_j or f_{sd}	a certain flow
d_j	demand of a flow
P_j	a path the flow j takes
$S = \{s_1, \dots, s_k, \dots, s_K\}$	set of slices
$f_j^k \in F^k$	set of flows to be deployed in slice k



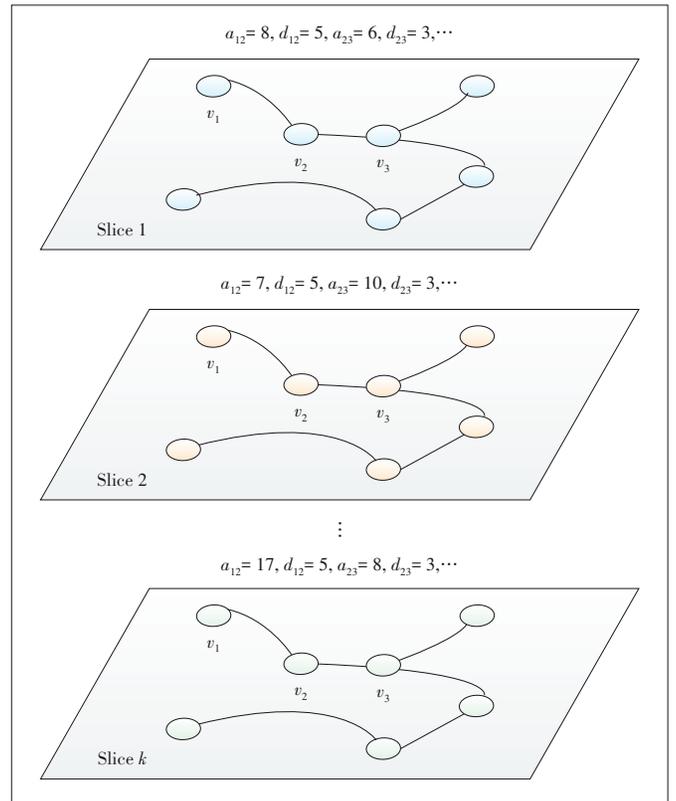
▲ **Figure 1. A network with 3 nodes and 3 edges**

tween v_1 and v_2 cannot hold these two flows, so we need to calculate paths sequentially.

Due to this, the total time for planning all flows is the sum of time planning for each flow. In the real production environment, we always expect a minimum deployment time. To achieve this goal, we need to establish some resource isolation, so that paralleling computing can be applied to the separate resource part. This is a method called slicing which is what we are going to introduce in the following. First, we introduce the slice set $S = \{s_1, \dots, s_k, \dots, s_K\}$, which represents a set of slices. Each slice has the same topology as G , and each links' capacity is part of the original network. We call the process of dividing bandwidth into parts slicing. Assuming that the original network is divided into K slices, in this scenario, the total time of calculating route paths is the maximum time consumed among all slices. If we have K slices, the time consumed can be approximately reduced to $1/K$ of original time-consuming.

All slices can be seen as virtual networks based on the same physical network. In detail, we can also formulate a slice as an undirected graph $s_k = \{V, E, A_k\}$, which has the same topology as the original network. Additionally, the delay of an edge is the same across slices.

Different slices have the same network topology and delay, while their bandwidths and available bandwidths can be different, which depends on the slicing method. **Fig. 2** illustrates



▲ **Figure 2. Network slices in our slicing model**

the slicing method in our slicing model.

3.3 Resources Allocation

Given the network slicing, by randomly allocating users' flow requests to slices, these flows in different slices can be routed by SP in parallel, which in turn can reduce total routing time. However, compared with routing on the original network, this network slicing-based routing mechanism will lead to a higher failure rate, since the network resources might be sliced in unavailable fragments. Therefore, this method of directly slicing reduces the flexibility of routing compared with routing on the original network without slicing. In other words, this method improves the calculation time performance at the expense of the success rate of routing.

To improve the flow deployment success rate, we design an appropriate method to determine which slice the flow should be deployed in. We use f_j^k to denote that f_j will be deployed in s_k , $F^k = \{f_1^k, \dots, f_h^k\}$ represent the set of flows deployed in s_k , and we also use F^k to represent the set of flows F_{succ}^k successfully deployed in s_k .

Our objective is to design a flow allocation algorithm to maximize the success rate of flow deployment on slices, which can be formulated as $\max \sum_{i=1}^K |F_{succ}^i|$. Particularly, we hope that the success rate after slicing can be close to the success rate of deployment on the original topology.

4 Algorithm

In this section, we propose a two-stage network resource allocation algorithm. In the first stage, once the flow requests are received, we propose a RL-based mechanism to allocate these flows to slices in sequence. In the second stage, aware of the imbalance of resources among slices, we propose a real-time resource adjustment mechanism to balance the resources dynamically.

4.1 RL Based Flow Allocation

We use reinforcement learning to train a general policy for specific topology and flow, which is used to determine on which slice each flow is deployed to maximize the success rate.

RL is a field of machine learning that emphasizes how to act based on the environment to maximize the expected benefits^[16]. The basic reinforcement learning model is defined by the tuple $\langle S, A, R, P \rangle$, where S is the set of states, A is the set of actions, R is the reward function and P is the state transition probability. We formulate the flow deployment with network slicing problem as a Markov decision process (MDP), shown as follows:

- Observation state: The observation obtained by the agent from the simulation environment is composed of two parts. One is the current state of all slices, and the other is the infor-

mation of the flow to be deployed. We use obs and obf to represent these two parts. So, the state $s \in S$ can be denoted by $s = \{ob_s, ob_f\}$. More specifically, ob_s is represented by the adjacency matrix of the available bandwidth of each slice. ob_f includes the source node and the target node, and the size of the current flow f_j .

- Action: The action space is a discrete space of slice index. Specifically, our RL agent selects the slice number to deploy for each flow, so we use $A = \{1, 2, \dots, K\}$ to represent the action space, where K is the number of slices.

- Reward: Reward is numerical feedback obtained by the agent after taking an action according to the current state and applying it in the environment. The magnitude of the value can reflect the quality of the action, and the reinforcement learning design is used to maximize the reward value in the long-term range. In our situation, the specific form of the reward is shown in Eq. (1) below.

$$R(f_j, a) = s_j^a \times \frac{K}{\sum_{k=1}^K \mathbb{I}(s_j^a = s_j^k)}, \quad (1)$$

in which $s_j^k \in \{-1, 1\}$ is a binary variable, when a flow f_j to be deployed is successfully deployed on the slice s_k . The value of s_j^k is 1, otherwise the value is -1. The $\mathbb{I}(\cdot)$ is an indicator function, and the condition is 1, otherwise, it is 0.

The reward function designed in this way can properly reflect the correctness of the decision made by the RL agent. When most of the slices can correctly deploy the flow, it is easy to make a decision, so a small reward will be generated. Moreover, if most of the slices cannot deploy the flow correctly, and the RL agent makes a correct decision, it will produce a large positive reward. Conversely, the wrong decision will always produce a negative reward. The absolute value of the reward is related to the difficulty of making a wrong decision.

- Transition probability: The MDP transition probability $p(s_i + 1 | s_i, a_i)$ is deterministic, which depends on the way that the flow routing method and the resource adaptation between slices after the flow is deployed. It will be introduced in the next section.

We want to train an agent to decide which slice to assign based on the state of each slice. The process that the RL policy interacts with the environment and the specific elements of the learning environment can be seen in **Fig. 3**.

1) RL training mechanism: We use an asynchronous proximal policy optimization (APPO) algorithm^[17] based on the importance weighted actor-learner architecture (IMPALA)^[18] to train our agent. Compared with synchronous proximal policy optimization (PPO), APPO is more efficient in wall-clock time due to its use of asynchronous sampling. Using a clipped loss also allows for multiple stochastic gradient descent (SGD) passes, and therefore the potential for better sample efficiency compared with IMPALA. Meantime, V-trace can also be en-

abled to correct for off-policy samples. The PPO-based flow allocation algorithm is described in Ref. [17]. It repeatedly uses the data obtained by sampling to update the strategy parameters with gradient ascent until the strategy is no longer updated. The specific proof of convergence can be seen in the original paper^[17].

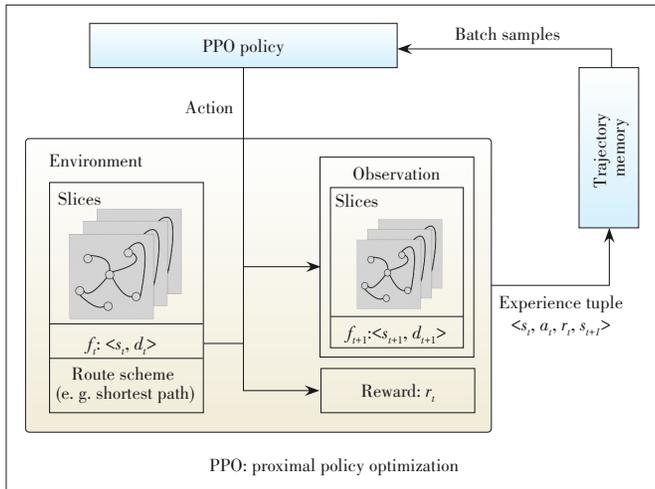
The PPO is more stable than the Q-learning algorithm learning process. This training framework can use distributed learning to solve actual large-scale network topology scenarios.

2) RL Training Architecture: We use a framework called ray^[20-22] for our reinforcement learning training, which can provide simple primitives for building and running distributed applications. With the help of ray, we can build a distributed framework to accelerate the training process of reinforcement learning. The specific framework is shown in Fig. 4. Each rollout worker contains a simulation environment and the latest policy, where each policy exchanges weights with the central SGD learning process at regular intervals to obtain the latest learned policy.

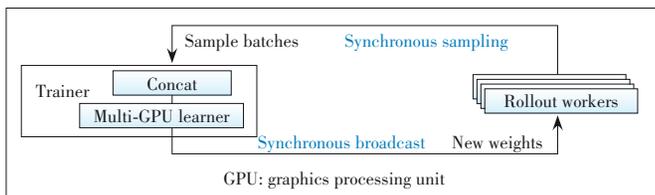
In this distributed reinforcement learning architecture, the central learner runs SGD in a tight loop, while asynchronously extracting sample batches from many participant processes, and also supports multiple GPU learners and experience replay.

4.2 Dynamic Resource Adjustment Among Slices

We first introduce the traditional resource allocation in the slice scenario. In a network that already has some flows, each link in the network has different remaining bandwidths. At this time, we have to start the slicing operation. The common



▲ Figure 3. Interaction between policy and environment



▲ Figure 4. Reinforcement learning (RL) training architecture

method is to evenly allocate the remaining bandwidth of each link to each slice. This relatively fair initialization strategy is intuitive; however, the resources of slices are unchanged during deployment. Therefore, the flexibility of flow routing on each slice is greatly reduced compared with when it is not sliced.

We now propose the dynamic resource adjustment strategy to increase the flexibility of resource allocation method. Compared with the traditional (static) method, our proposed method dynamically adjusts resources among slices in a real-time manner, aiming at balancing the type of flow deployed in the slice.

The specific algorithm description can be seen in Algorithm 1, and the deployment in each slice can be processed in separate processes. In the beginning, in each slice process, we expand the link bandwidth of the slice in step 3, and this step is to apply resources transferred from other slices to the current slice before the flow deployment. Then in step 4, we calculate the path by the specified routing strategy. In the algorithm, we use the simple strategy of the SP as an example. In fact, it can be any routing strategy. After this, in step 5 we check the bandwidth request from another slice. If the available bandwidth is greater than twice the maximum flow size, we then transfer the required bandwidth and reduce it in the link. At last, in step 6, we send requests to other slices to increase the bandwidth resources of each link in P_n .

Algorithm 1 Dynamic resource adjustment

Input: F -flows, S -slices.

Output: Resource adjustment

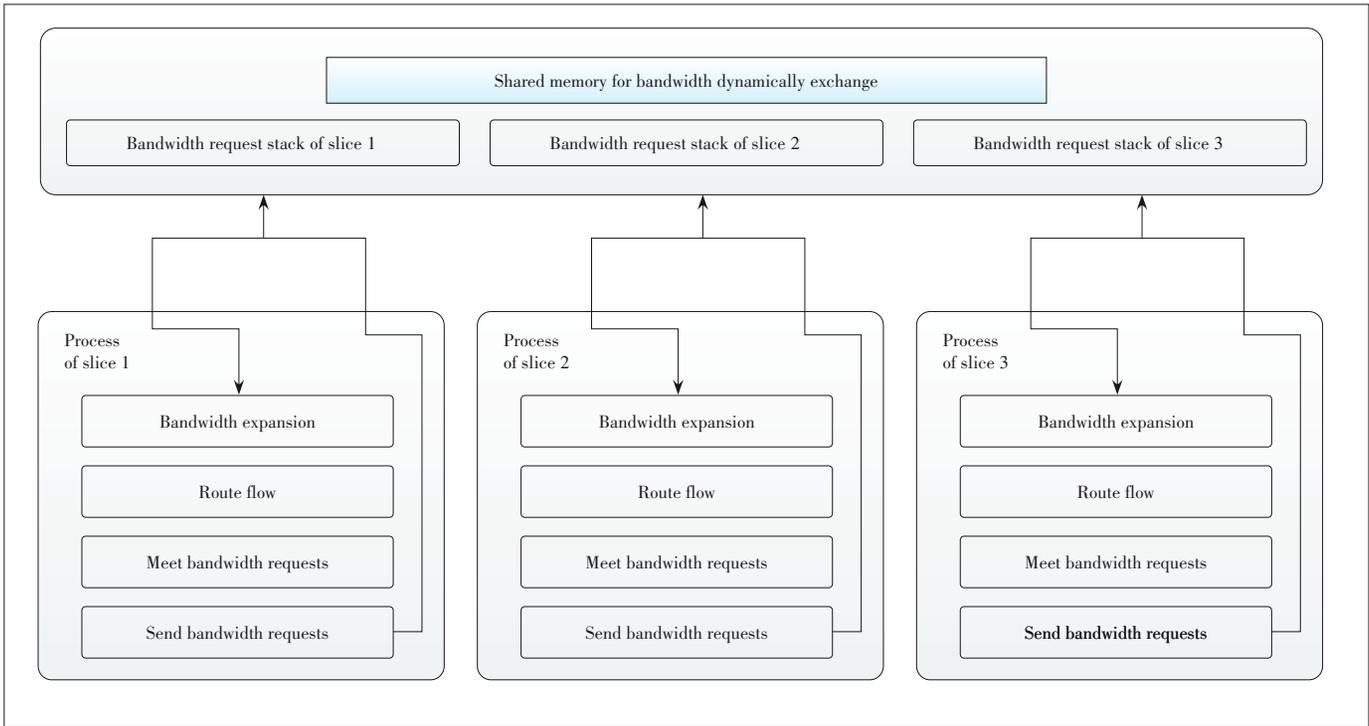
```

1: for all  $s_k$  in  $S$  do
2:   for  $f_n$  in  $F^k$  do
3:     Expand link bandwidth in  $s_k$ 
4:      $P_n \leftarrow SPf_n$ 
5:     Meet the bandwidth transfer requirements of other slices
6:     Request bandwidth resources of  $P_n$  from other slices
7:   end for
8: end for

```

The specific structure of Algorithm 1 can be seen in Fig. 5, where we use a shared memory stack to asynchronously exchange bandwidth resources between different slice processes. This is an example of three slices, which can be adjusted to any number of slices. Each slice first obtains the bandwidth resources transferred by other slices, and then deploys the flow. After that, the slice determines whether the current bandwidth resources can meet the requests of other slices, and if so, transfers the bandwidth to the other slices. The slice finally sends bandwidth transfer requests to other slices according to the path taken by the current deployment flow.

In the flow classification process of RL training, the environment that RL agent learns from also contains a resource adaptation strategy between slices. With the cooperation of the two methods, RL can learn a policy to allocate flows that re-



▲ Figure 5. Dynamic resource adjustment

peatedly pass through certain links to the same slice. One slice focuses on transmitting the same type of flows, and other infrequently used links' bandwidth resources are automatically adjusted to other more needed slices.

5 Experiment

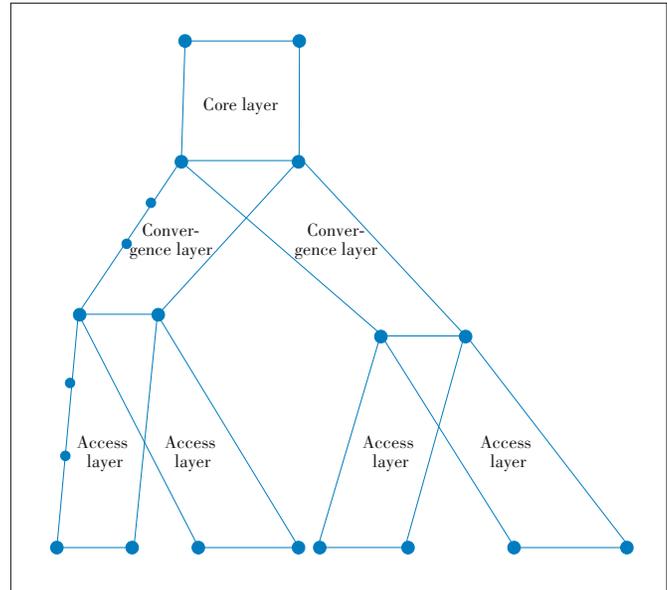
In this section, we conduct some experiments to explain the effectiveness of our algorithm.

5.1 Experiment Setup

We first introduce the information of network topology and flow.

Network topology: We use the ring network similar to a real 5G network as the experimental environment. This ring network consists of three kinds of rings, namely, the core layer, the convergence layer, and the access layer. The bandwidth and delay of links in three kinds of layers are different. The core layer, a ring network composed of the core equipment, has the largest bandwidth of the network, and it is the destination of most services in the network. The convergence layer, the bandwidth of which is smaller than that of the core layer, connects the core layer and the access layer, and aggregates each access layer network. The access layer has the smallest bandwidth of the network, which is composed of users and terminals. An example of the network topology of the ring network is shown in Fig. 6. In this experiment, we use a ring network with 6 245 nodes and 8 135 links.

Flows: The source and end nodes of the flow are randomly gen-

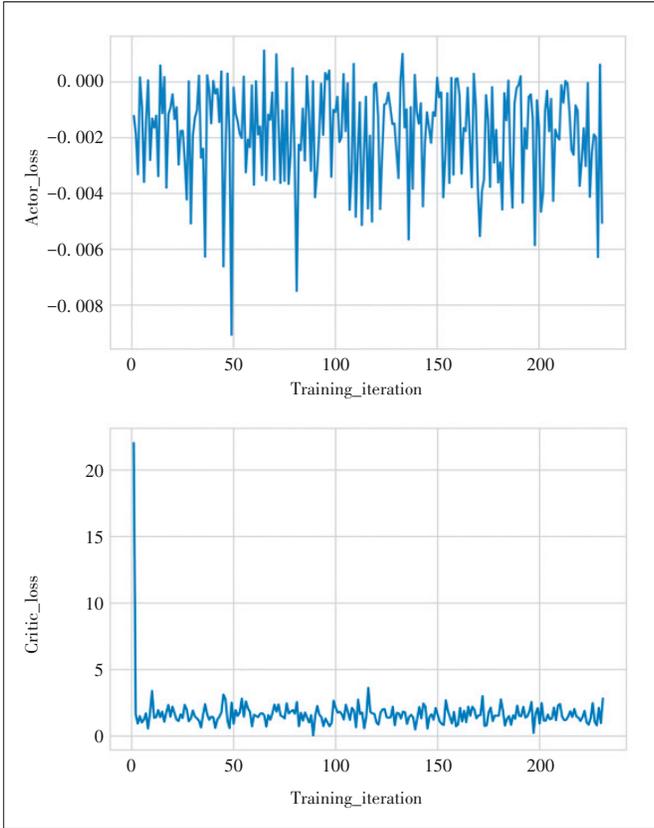


▲ Figure 6. An example of network topology of the ring network

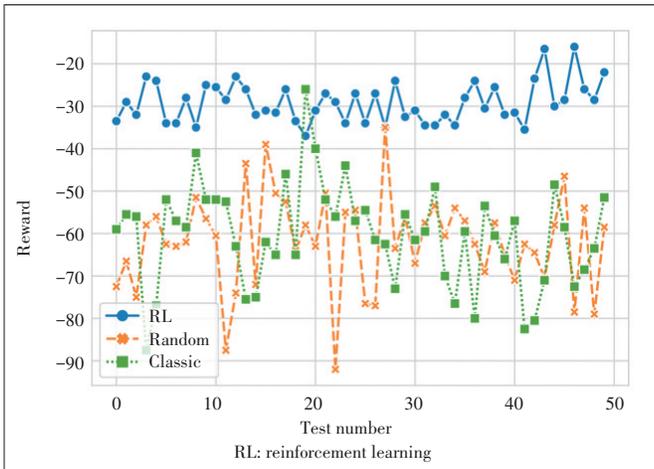
erated, and the size of the flow is $U(0.1,0) * \max \text{bandwidth}_{\text{access}}$, where $U(0.1, 1)$ is a uniform distribution.

5.2 Result of RL Agent

In the training process, thanks to the ray framework and careful hyperparameter adjustment, our algorithm can tend to converge within five iterations. The specific training loss is shown in Fig. 7. The critic loss can reflect the accuracy of RL



▲ Figure 7. Loss during training is based on a 2080 Ti graphics processing unit (GPU) and 120 sampling processes (2 Xeon CPUs)



▲ Figure 8. Total reward of different methods

agent's estimation of observation, and the policy loss is close to zero, which indicates that the strategy is close to the optimal strategy.

In RL test, we use accumulated rewards to show whether the classification effect is good or not, and we also convert other classification into the same measurement indicators. From Fig. 8, we can see that the cumulative reward of the converged RL scheme fluctuates around -30 (the larger the

better). The cumulative reward of the classic pooling method and random choice method used for comparison is around -70 , so we can see that RL can effectively choose the correct slice for the flow.

The above is a unified measurement in the reinforcement learning environment, and what we actually care about is whether the success rate of the deployment can be improved with the help of RL, which relates to whether our reward design is reasonable.

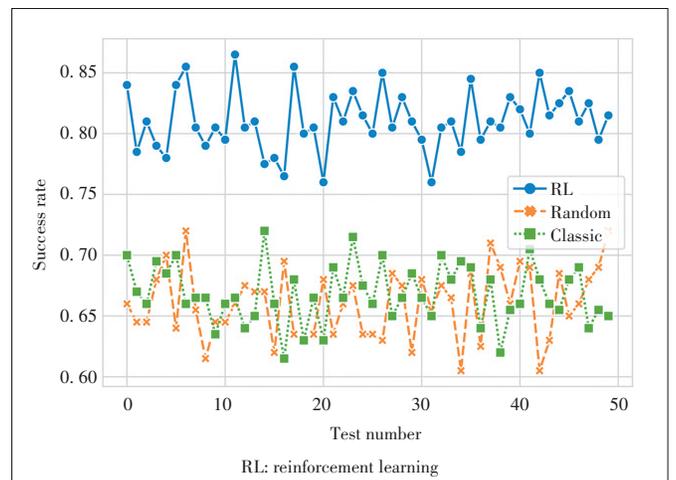
So we test the success rate of using the RL agent for flow classification, as well as the success rate of the classic method and random method as a comparison, and the results are shown in Fig. 9. We can see that the strategy learned using RL is significantly better than the classic method while maintaining the same trend as in Fig. 8, which shows that the reward we designed for RL is reasonable.

5.3 Result of Dynamic Resource Adjustment Strategy

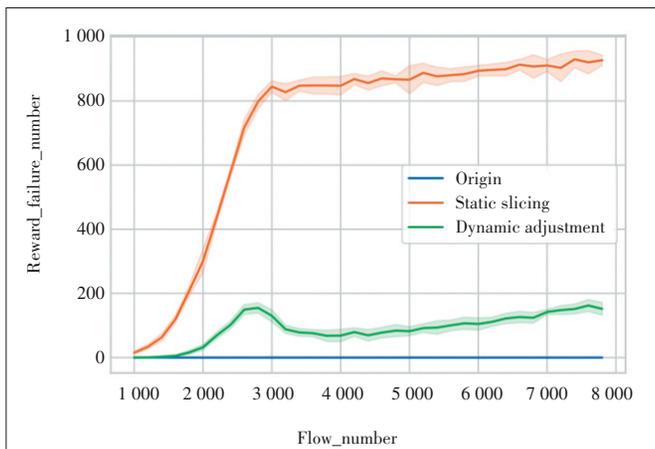
In the topology set above, we continuously increase the number of flows to compare the effects of different strategies by comparing the number of failed deployments. The first is the non-slicing case. After slicing, due to the decrease of routing flexibility, the failure rate will inevitably be greater than the case without slicing, so we regard the non-slicing case as a benchmark. Then we test the static slicing strategy and compare the dynamically adapted slicing strategy we proposed. The result can be seen in Fig. 10.

From Fig. 10, we can see that the dynamic adjustment strategy can provide a failure rate that is almost close to that of an unsliced case, and the effect is much better than the static solution.

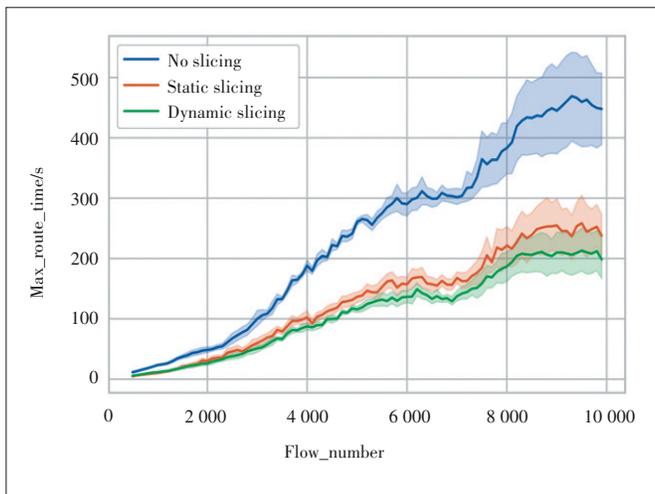
The purpose of slicing is to save the calculation time for parallel calculation, but the dynamic adjustment strategy will obviously increase some time consumption, so we test the average time of routing each flow, and the results are shown in Fig. 11. We can get that dynamic action almost bring no additional time consumption.



▲ Figure 9. Success Rate of different methods



▲ Figure 10. Relative failure number of dynamic resource adaptation vs. static slicing



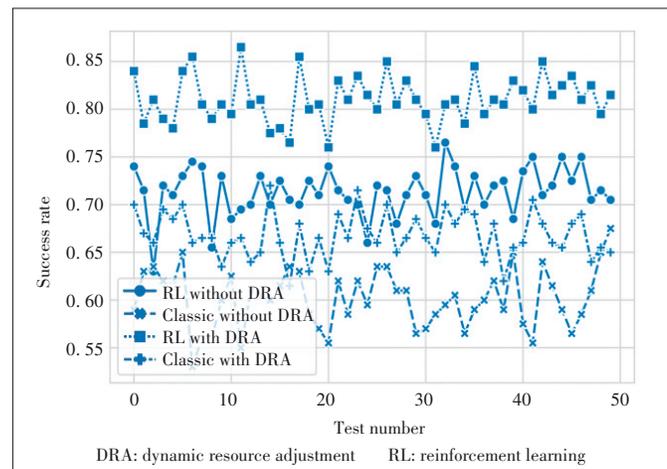
▲ Figure 11. Comparison of calculation time of different methods

5.4 Comparison of Effectiveness of Combined Methods

Although we have confirmed the effectiveness of the dynamic resource adjustment strategy, we still have uncertainty about the combined effect of the two algorithms. Therefore, we conduct the following experiments to compare the performance of the RL agent flow classification and the classic method with and without the dynamic resource adjustment strategy. The result is shown in **Fig. 12**. It can be seen that no matter which method is used, the dynamic adjustment strategy can indeed further improve the success rate of flow deployment.

6 Conclusions

To solve the time-consuming problem of path calculation, a slicing operation is used to separate resources so that parallel path calculation can be performed to shorten the time. At the same time, after slicing, although the calculation time is shortened, the more serious problem occurs, which is the decline in



▲ Figure 12. Comparison of the effectiveness of combined methods

the success rate of flow deployment. In response to this problem, we propose an efficient network slicing with dynamic resource allocation algorithm to let the success rate of flow deployment close to the level when not sliced. It combines the flow classification of reinforcement learning and resource adaptation between slices. The dynamic resource adaptive strategy between slices enables slice resources to be exchanged during flow deployment and gradually adapt to the characteristics of the flows to be deployed in each slice. In this way, the parallel calculation path can be shortened without sacrificing the success rate. Besides, we have also tested our method on a large-scale actual network structure, and the effect exceeds the previous classic methods.

References

- [1] JEFFREY G A, BUZZI S, CHOI W, et al. What will 5G be? [J]. IEEE Journal on selected areas in communications, 2014, 32(6): 1065. DOI: 10.1109/JSAC.2014.2328098
- [2] EINSIEDLER H J, GAVRAS A, SELLSTEDT P, et al. System design for 5G converged networks [C]//2015 European Conference on Networks and Communications. Paris, France: IEEE, 2015: 391 - 396. DOI:10.1109/EuCNC.2015.7194105
- [3] HAWILO H, SHAMI A, MIRAHMADI M, et al. NFV: state of the art, challenges, and implementation in next generation mobile networks [J]. IEEE network, 2014, 28(6): 18 - 26. DOI:10.1109/MNET.2014.6963800
- [4] KARAKUS M, DURRESI A. Quality of service (QoS) in software defined networking (SDN): a survey [J]. Journal of network and computer applications, 2017, 80: 200 - 218. DOI: 10.1016/j.jnca.2016.12.019
- [5] WALLNER R, CANNISTRA R. An SDN approach: quality of service using big switch's floodlight open-source controller [J]. Proceedings of the Asia-Pacific advanced network, 2013, 35: 14. DOI:10.7125/APAN.35.2
- [6] XU C, CHEN B, QIAN H. Quality of service guaranteed resource management dynamically in software defined network [J]. Journal of communications, 2015: 843 - 850. DOI:10.12720/jcm.10.11.843-850
- [7] XU C, GAMAGE S, LU H. vTurbo: accelerating virtual machine I/O processing

- using designated turbo-sliced core [C]//2013 USENIX Annual Technical Conference. San Jose, USA: USENIX, 2013:243 – 254
- [8] SCANO D, VALCARENGHI L, KONDEPU K, et al. Network slicing in SDN networks [C]//2020 22nd International Conference on Transparent Optical Networks (ICTON). Bari, Italy: IEEE, 2020: 1 – 4. DOI: 10.1109/ICTON51198.2020.9203184
- [9] ZHU K, HOSSAIN E. Virtualization of 5G cellular networks as a hierarchical combinatorial auction [J]. IEEE transactions on mobile computing, 2016, 15(10): 2640 – 2654. DOI:10.1109/TMC.2015.2506578
- [10] KATSALIS K, NIKAEIN N, SCHILLER E, et al. Network slices toward 5G communications: Slicing the LTE network [J]. IEEE communications magazine, 2017, 55(8): 146 – 154. DOI:10.1109/MCOM.2017.1600936
- [11] BARI M F, BOUTABA R, ESTEVES R, et al. Data center network virtualization: A survey [J]. IEEE communications surveys & tutorials, 2013, 15(2): 909 – 928. DOI:10.1109/SURV.2012.090512.00043
- [12] ZHANG Q X, LIU F M, ZENG C B. Adaptive interference-aware VNF placement for service - customized 5G network slices [C]//IEEE INFOCOM 2019 IEEE Conference on Computer Communications. Paris, France: IEEE, 2019: 2449 – 2457. DOI:10.1109/INFOCOM.2019.8737660
- [13] CHEN J J, TSAI M H, ZHAO L Q, et al. Realizing dynamic network slice resource management based on SDN networks [C]//2019 International Conference on Intelligent Computing and its Emerging Applications (ICEA). Tainan, Taiwan, China: IEEE, 2019: 120 – 125. DOI:10.1109/ICEA.2019.8858288
- [14] PARSAEFFARD S, JUMBA V, DERAKHSHANI M, et al. Joint resource provisioning and admission control in wireless virtualized networks [C]//2015 IEEE Wireless Communications and Networking Conference. New Orleans, USA: IEEE, 2015: 2020 – 2025. DOI:10.1109/WCNC.2015.7127778
- [15] MONEMI M, RASTI M, HOSSAIN E. Low-complexity SINR feasibility checking and joint power and admission control in prioritized multitier cellular networks [J]. IEEE transactions on wireless communications, 2016, 15(3): 2421 – 2434. DOI:10.1109/TWC.2015.2504084
- [16] SUTTON R S, BARTO A G. Reinforcement learning: an introduction [M]. Cambridge, UK: MIT Press, 2018
- [17] SCHULMAN J, WOLSKI F, DHARIWAL P, et al. Proximal policy optimization algorithms [EB/OL]. (2017-08-28) [2021-12-09]. <https://arxiv.org/abs/1707.06347>
- [18] ESPEHOLT L, SOYER H, MUNOS R, et al. IMPALA: scalable distributed deep-RL with importance weighted actor-learner architectures [EB/OL]. (2018-06-28) [2021-12-09]. <https://arxiv.org/abs/1802.01561?context=cs>
- [19] STOOKE A, ABBEEL P. Accelerated methods for deep reinforcement learning [EB/OL]. (2019-12-09) [2021-01-18]. <https://arxiv.org/abs/1803.02811>
- [20] MORITZ P, NISHIHARA R, WANG S, et al. Ray: A distributed framework for emerging AI applications [C]//13th USENIX Symposium on Operating Systems Design and Implementation. Carlsbad, USA: OSDI, 2018: 561 – 577
- [21] LIANG E, LIAW R, NISHIHARA R, et al. RLLIB: abstractions for distributed reinforcement learning [EB/OL]. (2018-07-29) [2021-12-09]. <https://arxiv.org/abs/1712.09381>
- [22] LIAW R, LIANG E, NISHIHARA R, et al. Tune: a research platform for distributed model selection and training [EB/OL]. (2018-07-13) [2021-12-09]. <https://arxiv.org/abs/1807.05118>

Biographies

JI Hong (hong_ji@seu.edu.cn) received the B.S. degree from School of Economics and Management from Xidian University, China in 2018. He is currently a master student at School of Cyber Science and Engineering, Southeast University, China. His research interests lie in network optimization and intelligent decision making.

ZHANG Tianxiang received the B.S. degree from Nanjing University of Aeronautics and Astronautics, China. He is currently an engineer with ZTE Corporation. His research interests include traffic scheduling and graph neural networks.

ZHANG Kai is currently a master student at School of Cyber Science and Engineering, Southeast University, China. His research interests include graph neural networks and resource allocation and reinforcement learning.

WANG Wanyuan is an assistant professor with the School of Computer Science and Engineering, Southeast University, China. He received his Ph. D. degree in computer science from Southeast University in 2016. He has published several articles in refereed journals and conference proceedings, such as the *IEEE Transactions on Mobile Computing*, *IEEE Journal on Selected Areas in Communications*, *IEEE Transactions on Cybernetics*, *AAAI*, and *AAMAS*. He won the best student paper award from ICTAI14. His main research interests include artificial intelligence, multiagent systems, and game theory.

WU Weiwei is a professor in the School of Computer Science and Engineering, Southeast University, China. He received his B.Sc. degree from South China University of Technology, China and the Ph.D. degree from City University of Hong Kong (CityU), China and University of Science and Technology of China (USTC) in 2011, and went to Nanyang Technological University, Singapore for post-doctorial research in 2012. He has published over 50 peer-reviewed papers in international conferences/journals, and serves as TPCs and reviewers for several top international journals and conferences. His research interests include optimizations and algorithm analysis, wireless communications, crowdsourcing, cloud computing, reinforcement learning, game theory and network economics.

Enabling Energy Efficiency in 5G Network



LIU Zhuang^{1,2}, GAO Yin^{1,2}, LI Dapeng¹, CHEN Jiajun¹, HAN Jiren¹

(1. R&D Center of ZTE Corporation, Shanghai 201203, China;

2. State Key Laboratory of Mobile Network and Mobile Multimedia, Shenzhen 518057, China)

Abstract: The mobile Internet and Internet of Things are considered the main driving forces of 5G, as they require an ultra-dense deployment of small base stations to meet the increasing traffic demands. 5G new radio (NR) access is designed to enable denser network deployments, while leading to a significant concern about the network energy consumption. Energy consumption is a main part of network operational expense (OPEX), and base stations work as the main energy consumption equipment in the radio access network (RAN). In order to achieve RAN energy efficiency (EE), switching off cells is a strategy to reduce the energy consumption of networks during off-peak conditions. This paper introduces NR cell switching on/off schemes in 3GPP to achieve energy efficiency in 5G RAN, including intra-system energy saving (ES) scheme and inter-system ES scheme. Additionally, NR architectural features including central unit/distributed unit (CU/DU) split and dual connectivity (DC) are also considered in NR energy saving. How to apply artificial intelligence (AI) into 5G networks is a new topic in 3GPP, and we also propose a machine learning (ML) based scheme to save energy by switching off the cell selected relying on the load prediction. According to the experiment results in the real wireless environment, the ML based ES scheme can reduce more power consumption than the conventional ES scheme without load prediction.

Keywords: cell switch off; energy efficiency; energy saving; 5G; machine learning

DOI: 10.12142/ZTECOM.202101004

<https://kns.cnki.net/kcms/detail/34.1294.TN.20210220.1440.004.html>, published online February 20, 2021

Manuscript received: 2020-12-10

Citation (IEEE Format): Z. Liu, Y. Gao, D. P. Li, et al., "Enabling energy efficiency in 5G network," *ZTE Communications*, vol. 19, no. 1, pp. 20 - 29, Mar. 2021. doi: 10.12142/ZTECOM.202101004.

1 Introduction

To cope with expected drastic data traffic growth, 5G new radio (NR) is designed to enable denser network deployments, while the densification of networks has implied higher energy expenditure. In a typical radio access network (RAN), most energy is consumed by base stations. However, with the foreseen NR deployment of more base stations with massive multiple-input multiple-output (MIMO), energy efficiency (EE) in NR becomes even more urgent

and challenging.

Energy consumption (EC) is a main part of operational expense (OPEX). The telecommunication operators are seeking for a better way to expand market shares while energy consumption in networks can be decreased to lower their OPEX. Energy efficiency in NR networks is also a significant research topic in 3GPP. Switching off cells is a widely used strategy to reduce the energy consumption of networks during off-peak conditions. Thus network elements with low power consumption become more and more important and the shut-

down of unused capacity cells is also valuable. The important aspect of RAN energy efficiency is, during the network running, how to ensure cell switching-off without affecting the customer satisfaction e.g., calls dropped; quality of service (QoS) degraded. A typical energy saving (ES) scenario is that capacity booster cells are deployed under the umbrella of cells providing basic coverage and that the capacity booster cells can be switched off to enter into the dormant mode when its capacity is no longer needed and to be reactivated on a need basis. This paper introduces the 3GPP schemes for switching on/off NR cells to achieve energy efficiency in 5G networks, including the 5G intra-system energy saving scheme and 4G/5G inter-system energy saving scheme involving different core networks (CN), e.g., evolved packet core (EPC) and 5G core (5GC) networks. We also propose a machine learning (ML) based scheme to save energy by switching off the cell selected relying on load prediction. According to the experiment results in the real wireless environment, the ML based ES scheme can reduce more power consumption than the conventional ES scheme without load prediction.

2 Cell Switch on/off for Energy Saving

An NR cell, which acts as a capacity booster, may be switched off and enter into the ES dormant state if there is radio coverage by another cell. **Fig. 1** shows an example of the next-generation Node B (gNB) capacity booster cell fully overlaid by a coverage providing cell. The gNB is a node providing the NR user plane and control plane with protocol terminations towards user equipment (UE), and connected to the 5GC via the next-generation (NG) interface. In the figure, Cell A is deployed to provide continuous coverage of the area, while Cell B provides more capacity only for special sub-areas, such as hot spots. The ES activation procedure of Cell B may be triggered in case that light traffic in Cell B is detected. Then the cell will be switched off and enter into the ES dormant state; if there are some users in service in Cell B, the cell will be switched off only after the handover actions to offload its traffic to Cell A is completed. The ES activation of Cell B may be triggered, that is, the cell is

switched on again, when the traffic of the ES area (measured by Cell A) resumes to a high level^[1].

In real network deployment, ES can be divided into centralized ES and distributed ES. For the distributed ES, the NR capacity booster cell may decide to switch off when it detects that its traffic load is below a certain threshold, and its coverage can be provided by the coverage providing cell. The coverage providing cell decides to reactivate the NR capacity booster cell when it detects additional capacity is needed. For the centralized ES, a centralized entity, such as the operation and maintenance (O&M) entity, collects the traffic load performance measurements from the NR capacity booster cell and coverage providing cells, and may request a NR capacity booster cell to switch off when its traffic is below certain threshold.

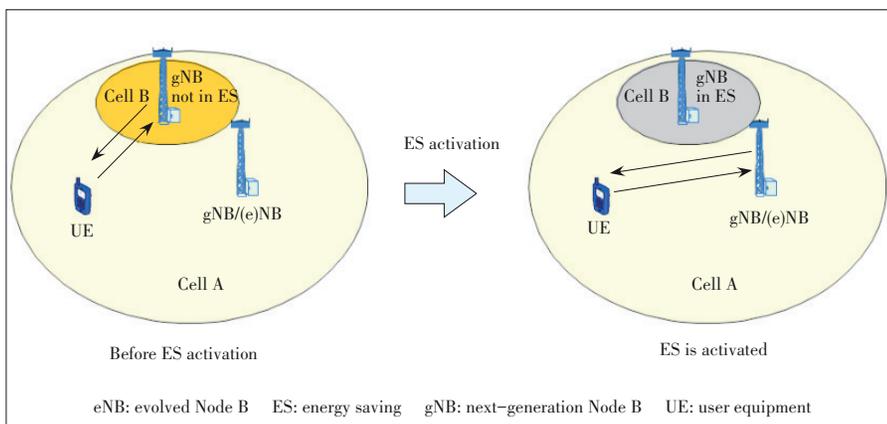
In general, NR energy saving solutions include the 5G intra-system energy saving scheme and 4G/5G inter-system energy saving scheme, involving different core networks (CN), such as EPC and 5GC. Additionally, NR architectural features including central unit/distributed unit (CU/DU) split and dual connectivity (DC) are also considered in NR energy saving. We will discuss these NR ES scenarios in the following sections.

3 5G Intra-System Energy Saving

3.1 Scenarios

In a 5G network, a next-generation RAN (NG-RAN) node is either a gNB or a next-generation evolved Node B (ng-eNB), providing Long Term Evolution (LTE) services or Evolved Terrestrial Radio Access Network (E-UTRAN) services towards the UE. The gNB provides services of the NR user plane and control plane; the ng-eNB provides the E-UTRA user plane and control plane with protocol terminations towards the UE, and connected to the 5GC via the NG interface. The gNBs and ng-eNBs are inter-connected with each other by means of the Xn interface, while they are connected to the 5GC by means of the NG interfaces. The scenarios for NR intra-system energy saving are summarized in **Table 1**.

In 5G intra-radio access technology (Intra-RAT) ES cases (Scenario 1; Scenario 2), some gNB (or ng-eNB) cells are deployed to provide basic coverage, while the other gNB (or ng-eNB) cells boost the capacity (**Fig. 2**). Therefore, the coverage provider and the capacity provider are using the same RAT, e.g., 5G NR or LTE, to provide NR services or LTE/E-UTRAN services to UE. The NG-RAN cell providing the capacity booster can decide to switch off autonomously; the switch-off decision may also be taken by the O&M entity that will inform the



▲ **Figure 1.** Capacity booster cell overlaid by coverage providing cell

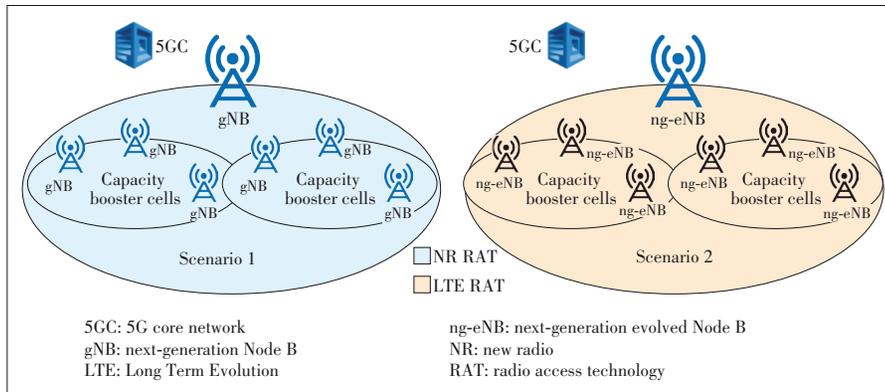
neighbor NG-RAN cell about its deactivation action over the Xn Application Protocol (XnAP). On the other hand, the NG-RAN node providing the basic coverage can request to reactivate the switched-off booster cell over XnAP.

In 5G inter-radio access technology (Inter-RAT) ES cases (Scenario 3; Scenario 4), some gNB (or ng-eNB) cells are deployed to provide basic coverage, while the other ng-eNB (or gNB) cells boost the capacity (Fig. 3). Obviously, the booster cells and coverage cells are in different RAT networks. That is, the booster cells are in NR RAT while the coverage cells in LTE RAT, or the booster cells are in LTE RAT while the coverage cells in NR RAT. The Xn signaling support for inter-RAT ES is the same as that for intra-RAT ES.

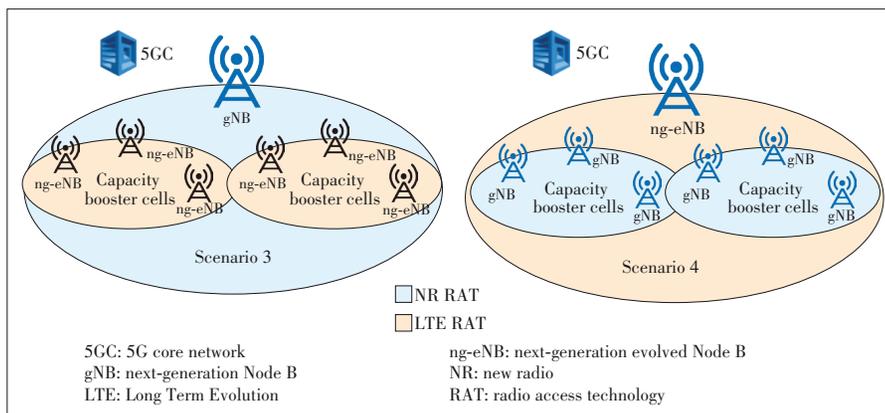
▼Table 1. 5G intra-system energy saving scenarios (only connected with 5GC)

5G Intra-System ES Scenario	Coverage Provider	Capacity Booster Provider	Description
1	gNB connected with 5GC	gNB connected with 5GC	intra-RAT ES
2	ng-eNB connected with 5GC	ng-eNB connected with 5GC	
3	gNB connected with 5GC	ng-eNB connected with 5GC	inter-RAT ES
4	ng-eNB connected with 5GC	gNB connected with 5GC	

5GC: 5G core network
 eNB: evolved Node B
 ES: energy saving
 gNB: next-generation Node B
 ng-eNB: next-generation evolved Node B
 NR: new radio
 RAT: radio access technology



▲Figure 2. 5G Intra-RAT energy saving



▲Figure 3. 5G Inter-RAT energy saving

3.2 Signaling Support

The Xn signaling support for both the intra-RAT ES and inter-RAT ES scenarios is same. As shown in Fig. 4, NG-RAN Node 1 that owns a capacity booster cell and can autonomously decide whether to switch off this cell based on cell load information, while the switch-off decision may also be taken by the O&M entity. All neighbor NG-RAN nodes are informed by the NG-RAN Node 1 owning the concerned cell about the switch-off actions over the Xn interface, by means of the NG-RAN node configuration update message.

The purpose of the cell activation procedure is to enable an NG-RAN node to request the neighboring NG-RAN node to switch on one or more cells that are previously reported as inactive due to energy saving reasons.

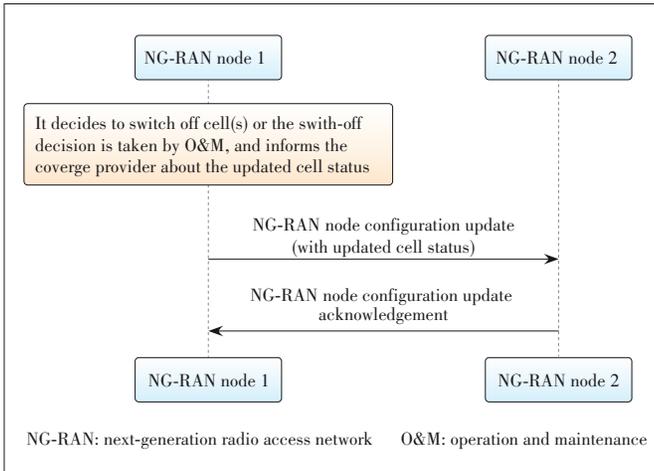
As shown in Fig. 5, if the basic coverage is ensured by NG-RAN node cells, the NG-RAN node owning non-capacity boosting cells may request a reactivation over the Xn interface via the cell activation procedure if needed. Upon receipt of a cell activation request message, the booster NG-RAN node activates the cells indicated in the message and these cells are also indicated in the cell activation response message when the request is fulfilled.

3.3 Intra-System Energy Saving for Multi-Radio Dual Connectivity

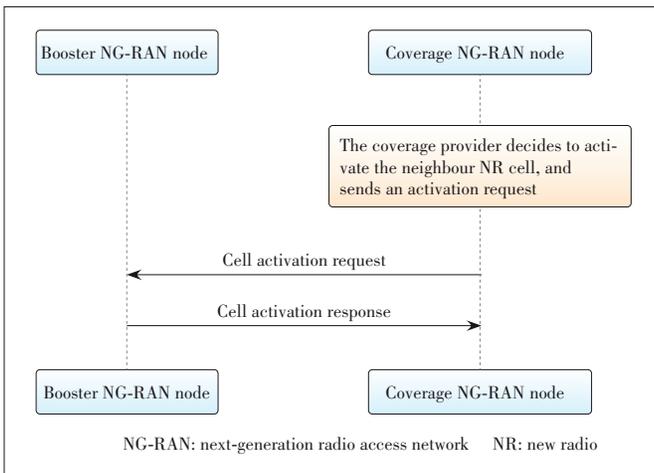
In order to implement multi-radio (MR) dual connectivity (DC), UE may be configured to utilize the resources provided by two different nodes, one providing NR access and the other one providing either E-UTRA or NR access. One node acts as the master node (MN) and the other as the secondary node (SN). The MN and SN are connected via an Xn or X2 interface and at least the MN is connected to a core network. In 3GPP Release 15, the MR DC energy saving is already supported, which is intra-system ES.

When the MN is connected to the 5GC, the DC cases include NG-RAN E-UTRA-NR dual connectivity (NGEN-DC), NE-NR-E-UTRA dual connectivity (DC) and NR-NR dual connectivity (NR-DC):

- NGEN-DC: one ng-eNB is connected with the 5GC and acts as an MN, while one gNB acts as an SN;
- NE-DC: one gNB is connected with the 5GC and acts as an MN, and one ng-eNB



▲ Figure 4. NG-RAN node informs the neighbor NG-RAN node about cell status over Xn



▲ Figure 5. Coverage NG-RAN node requests to activate booster cells over Xn

acts as an SN;

- NR-DC: one gNB is connected with the 5GC and acts as an MN, and another gNB acts as an SN.

The energy saving scheme for the above cases is similar to 5G intra-system ES (Section 3.1), where the SN can act as a capacity booster provider and MN provides continuous coverage of the area. The SN can autonomously decide to switch off cell(s) based on cell load information or the switch-off decision is taken by O&M; it then informs the MN about the cell deactivation action over XnAP. The MN can request to reactivate the switched-off booster cell at the SN over XnAP. The Xn ES Signaling support for MR DC with 5GC is the same with 5G Intra-system ES Signaling described in Section 3.2.

In the case that the MN is connected to the EPC, that is E-UTRA-NR dual connectivity (EN-DC), one eNB acts as an MN and one en-gNB acts as an SN. The eNB is connected to the EPC via the S1 interface and to the en-gNB via the X2 interface. The EN-DC scenario is shown in Fig. 6.

The EN-DC configuration update procedure and EN-DC cell

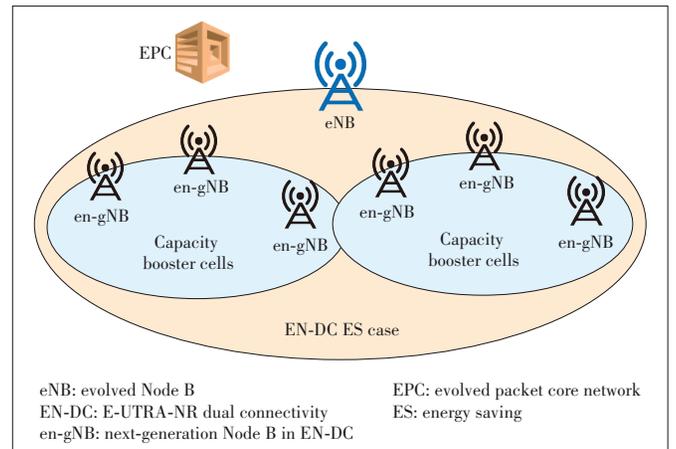
activation procedure are used to support EN DC for intra-system energy saving over the X2 interface. The EN-DC configuration update procedure can be used to exchange updated cell statuses of eNB and en-gNB over the X2 interface. The EN-DC cell activation procedure enables an eNB to request the neighboring en-gNB to switch on one or more cells that are previously reported as inactive due to energy saving reasons. Upon receipt of this message, the en-gNB should activate the cell/s indicated both in the cell activation request message and in the EN-DC cell activation response message sent after the activation request is fulfilled. Fig. 7 shows the detailed signaling flows.

4 Inter-System Energy Saving of 4G and 5G Systems

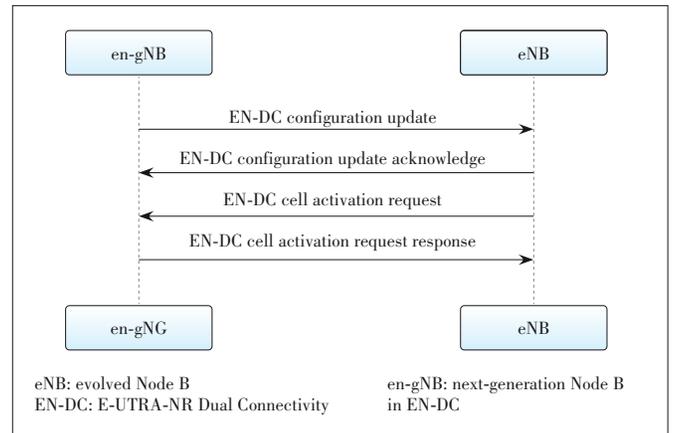
4.1 Scenarios

As shown in Fig. 8, 3G users continued to move to 4G from 2009 to Q3, 2018. Although 4G users kept growing steadily, the number of 2G and 3G users could not be ignored for a long time yet^[2].

Similar to what happened in the 4G era shown in Fig. 8, it



▲ Figure 6. EN-DC energy saving



▲ Figure 7. EN-DC energy saving signaling over X2

can be predicted that 4G users of operators will gradually decrease when 5G networks are deployed, but 4G networks will coexist with 5G networks for a long time. Therefore, inter-system energy saving solutions to 4G and 5G coexisting scenarios (Table 2) should be considered.

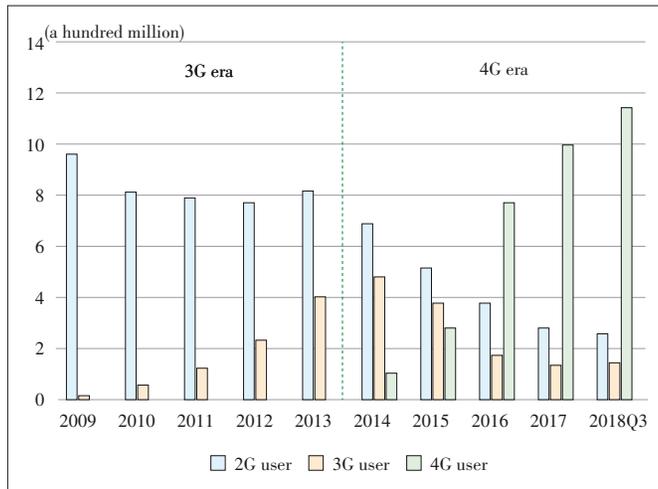
For the inter-system ES cases (Scenario 1 and Scenario 2 in Table 2), the NG-RAN node (a gNB, providing NR services; or an ng-eNB, providing E-UTRAN services) owns a capacity booster cell and can autonomously switch off this cell to the dormant state. The switch-off decision is typically based on cell load information, but may also be taken by the O&M entity. The NG-RAN node indicates the switch-off action to the eNB over the NG and S1 interfaces. The NG-RAN node could also indicate the switch-on action to the eNB over the NG and S1 interfaces. The eNB providing basic coverage may request an NG-RAN node's cell reactivation based on its own cell load information or neighbor cell load information, and the switch-on decision may also be taken by O&M. The eNB requests an NG-RAN node's cell reactivation and receives the NG-RAN node's cell reactivation reply from the NG-RAN node over the S1 and NG interfaces. The scenarios in Table 2 are shown in Fig. 9, where the E-UTRAN cell associated eNB and the NR-RAN cell associated gNB are connected to the EPC and the 5GC.

4.2 Signaling Support

3GPP Release15 (R15) defines signaling for cell activation/deactivation over X2 and Xn interfaces for intra-system ener-

gy saving. However, the signaling defined by R15 fails to support inter-system energy saving scenarios (Fig. 10^[3]) without a direct interface between the eNB and gNB/ng-eNB. The coverage eNB cannot directly send a request to reactivate the switched-off NR booster cell.

Therefore, NG and S1 interfaces are enhanced in 3GPP Release 16 (R16) to support the inter-system scenarios. Specifically, when the NR capacity booster cell is switched off, the LTE eNB for basic coverage should be informed by the gNB via the NG/S1 message; when the LTE eNB is going to acti-



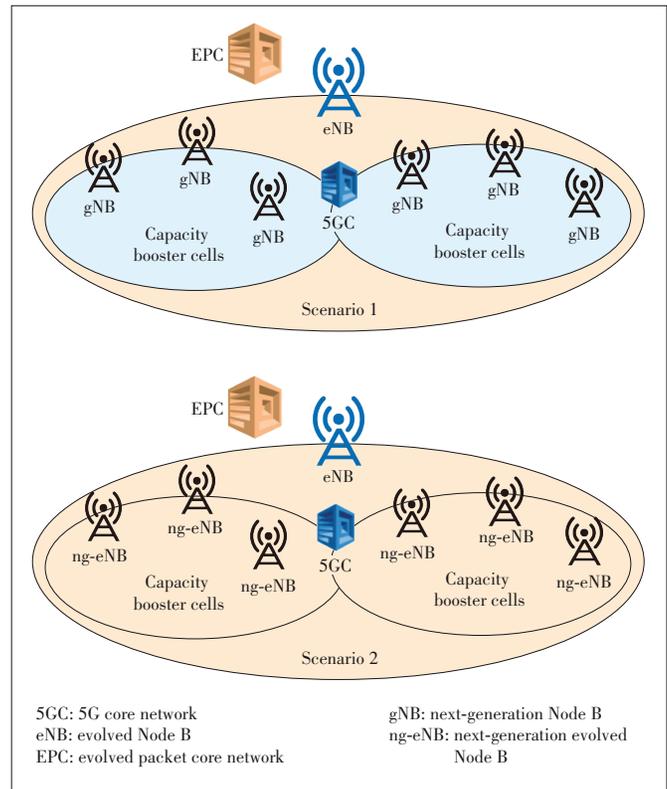
▲ Figure 8. Users in China continued to move from 2G and 3G networks to 4G networks

▼ Table 2. The inter-system ES scenarios of 4G and 5G systems (involving EPC and 5GC)

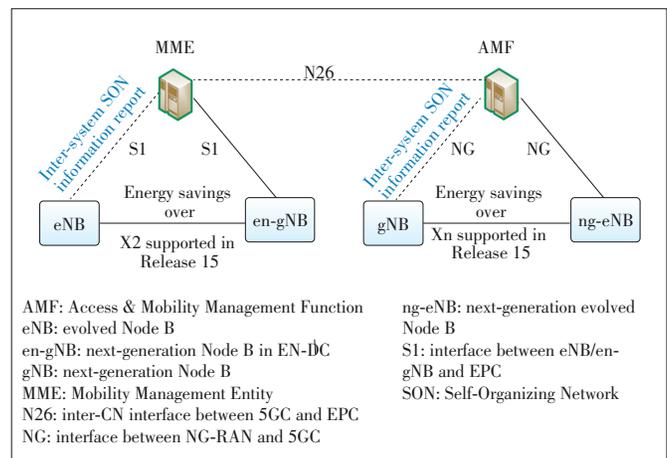
Scenario	Coverage Provider	Capacity Booster Provider
1	eNB connected with EPC	gNB connected with 5GC
2	eNB connected with EPC	ng-eNB connected with 5GC

5GC: 5G core network
eNB: evolved Node B
EPC: evolved packet core network

ES: energy saving
gNB: next-generation Node B
ng-eNB: next-generation evolved Node B



▲ Figure 9. Inter-system energy saving of 4G and 5G systems



▲ Figure 10. Inter-system energy saving is not supported by signaling defined by 3GPP Release 15

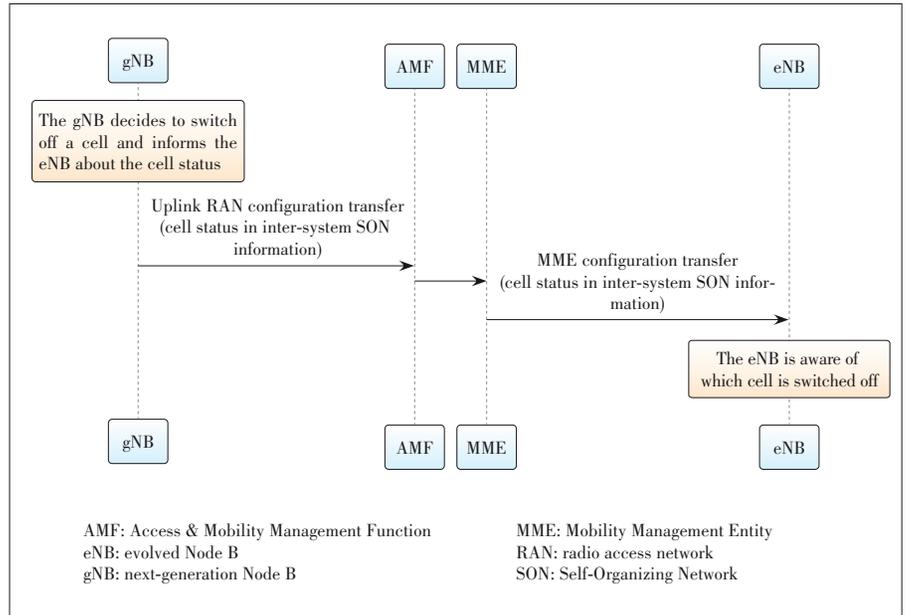
vate the NR cell, the gNB should be informed by the LTE eNB via the NG/S1 message. In R16, the inter-system Self-Organizing Network (SON) configuration transfer Information Element (IE) is introduced over both the S1 and NG interfaces in the following messages:

- eNB Configuration Transfer (TS36.413^[4]);
- Mobility Management Entity (MME) Configuration Transfer (TS36.413^[4]);
- Uplink RAN Configuration Transfer (TS38.413^[5]);
- Downlink RAN Configuration Transfer (TS38.413^[5]).

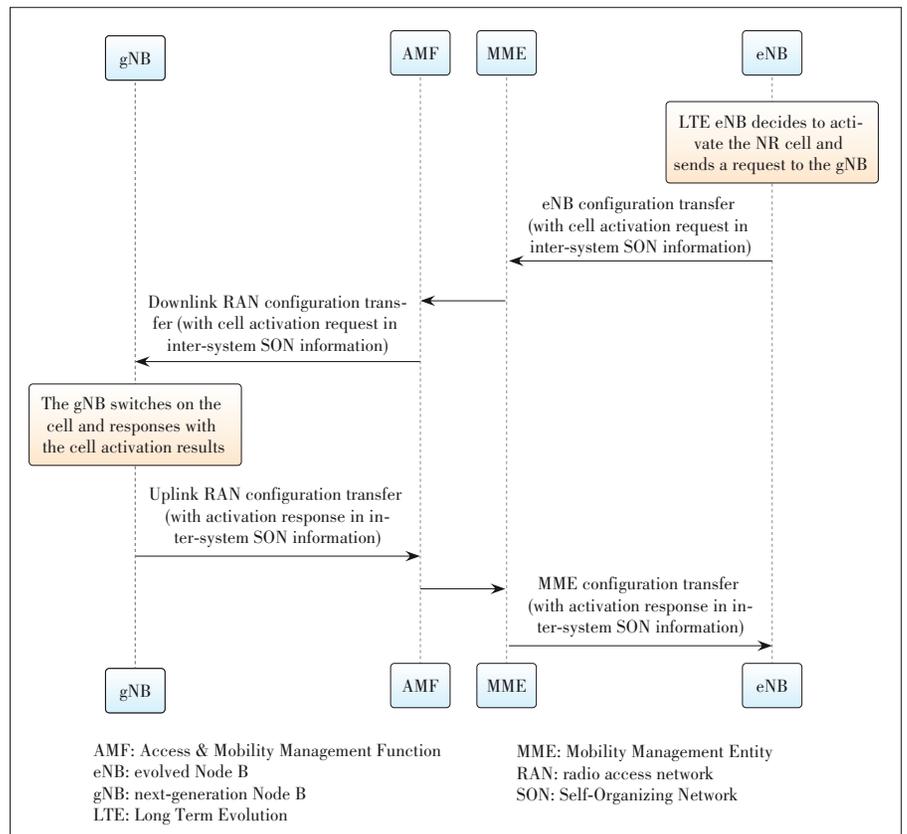
The messages of SON configuration transfer procedures, e.g., the eNB/MME configuration transfer message over S1 and uplink/downlink RAN configuration transfer message over NG, can be re-used for R16 inter-system energy saving, as the inter-system SON Information IE in such messages is extended to support inter-system cell status transfer and cell activation request and response between the eNB and NG-RAN node. The detailed signaling flows are shown in **Figs. 11 and 12** respectively.

In Fig. 11, the NG-RAN node owns a capacity booster cell and can autonomously switch off this cell to the dormant state. The switch-off decision is typically based on cell load information and consistent with the configured information. This decision may also be taken by the O&M entity. The NG-RAN node indicates either the switch-off or switch-on actions to the eNB over the NG and S1 interfaces.

In Fig. 12, the eNB providing basic coverage may request an NG-RAN node's cell reactivation based on its own cell load information or neighbor cell load information, and the switch-on decision may also be taken by O&M. The eNB requests an NG-RAN node's cell reactivation and receives the NG-RAN node's cell reactivation reply from the NG-RAN node over the S1 and NG interfaces.



▲ **Figure 11. Next-generation radio access network (NG-RAN) node connected with 5GC informs cell status to Long Term Evolution (LTE) eNB connected with evolved packet core (EPC) network**



▲ **Figure 12. LTE eNB connected with EPC requests to activate an NR CELL connected with 5GC**

5 Energy Saving in CU/DU Split Architecture

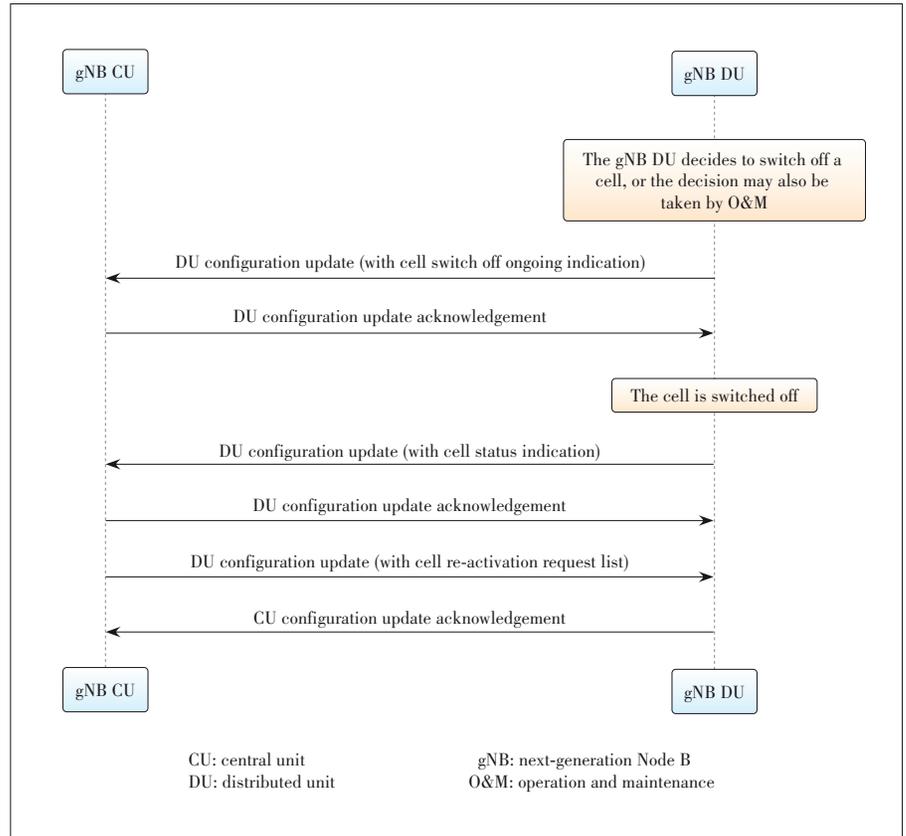
For the intra-system ES described in Section 3 and inter-system ES described in Section 4, if a gNB is deployed with

CU/DU split architecture, the F1 interface shall be enhanced to support the cell reactivation procedure and cell status exchange (**Fig. 13**).

When the booster gNB with CU/DU split decides to switch

off cell(s) to the dormant state, the decision is typically made by the gNB-DU based on cell load information or by the O&M entity. Before the cell in the gNB-DU enters into the dormant mode, the gNB-DU will send the gNB-DU configuration update message to the gNB-CU to indicate that the gNB-DU will switch off the cell after some time. During the switch-off period, the gNB-CU shall offload the UE to a neighboring cell and simultaneously not accept any incoming UE towards this switch-off ongoing cell. After the cell at gNB-DU enters into the dormant mode, the gNB-DU sends a new gNB-DU configuration update message to inform the “inactive” status of this cell to the gNB-CU. The gNB-CU needs to inform the updated cell status to the coverage provider node.

When the gNB-CU receives the cell activation request/EN-DC cell activation request from a coverage provider node over the Xn or X2 interface, or the gNB-CU decides to activate the dormant cell by itself, it will trigger the gNB-CU configuration update message to the gNB-DU with a list of the cells to be activated.



▲ Figure 13. CU/DU energy saving signaling support over F1 interface

6 Energy Efficiency KPI

EE Key Performance Indicator (KPI) shows data energy efficiency in NG-RAN. The EE KPI is defined as the data volume (in kbits) divided by energy consumption (in kWh) of the considered network elements. The unit of this KPI is bit/J^[6].

$$EE = \frac{\sum_{Samples} (DRB.PdcpSduVolumnUL + DRB.PdcpSduVolumnDL)}{\sum_{Samples} PEE.Energy}$$

(for non-split gNBs). (1)

$$EE = \frac{[(F_1uPdcpSduVolumeUL + XnuPdcpSduVolumeUL + \sum_{Samples} X_2uPdcpSduVolumeUL) + (F_1uPdcpSduVolumeDL + XnuPdcpSduVolumeDL + \sum_{Samples} X_2uPdcpSduVolumeDL)]}{\sum_{Samples} PEE.Energy}$$

(for split gNBs). (2)

For non-split gNBs (the gNBs without CU/DU split), the defined $DRB.PdcpSduVolumnUL$ in Eq. (1) is the measured data volume of Packet Data Convergence Protocol (PDCP) Service Data Unit (SDU) of a DRB in the uplink, delivered from the PDCP layer to Service Data Adaptation Protocol (SDAP) layer; $DRB.PdcpSduVolumnDL$ in the equation is the

measured data volume of PDCP SDU of a DRB in the downlink, delivered to the PDCP layer. The total data volume (in kbit) is obtained by measuring the amount of uplink and downlink PDCP SDU bits of all DRBs of the non-split gNBs over the measurement period.

For gNBs with CU/DU split, the defined $F_1uPdcpSduVolumeUL$ in Eq. (2) is the measured data volume of PDCP SDU in the uplink, delivered to gNB-CU-UP (gNB-CU-User Plane entity) from gNB-DU via F1-U (F1 User plane interface), $XnuPdcpSduVolumeUL$ is that in the uplink delivered from external gNB-CU-UP via Xn-U (Xn User plane interface), and $X_2uPdcpSduVolumeUL$ is that in the uplink delivered from external eNB via X2-U; the defined $F_1uPdcpSduVolumeDL$ in the equation is the measured data volume of PDCP SDU in the downlink, delivered from GNB-CU-UP to GNB-DU via F1-U, $XnuPdcpSduVolumeDL$ is that in the downlink delivered to external gNB-CU-UP via Xn-U, and $X_2uPdcpSduVolumeDL$ is that in the downlink delivered to external eNB via X2-U. The total data volume (in kbit) is obtained by measuring the amount of uplink and downlink PDCP SDU bits of all interfaces (F1-U, Xn-U and X2-U) of the split gNBs over the measurement period.

The energy consumption (in kWh) is obtained by measuring the Power, Energy and Environmental (PEE) of the considered network elements over the same period of time.

7 Machine Learning Based Energy Saving

In this paper, we introduce 3GPP energy saving schemes by cell switching on/off in ultra-dense networks. However, the traditional cell switching on/off relying on real-time load information is not accurate enough, the inappropriate switching-off of the cells may seriously deteriorate network performance because other active cells have to serve some extra traffic. In order to solve the potential issues of existing ES methods and to achieve intelligent ES, we propose a machine learning (ML) based scheme that relies on the load prediction to save energy by switching off the selected cell^[7]. As the data used for load prediction have various features, such as the current and history load and the neighbor cells' load, different techniques are used for the different features of load prediction. The auto-regression integrated moving average (ARIMA), Prophet, Random Forest (RF), Long Short-Term Memory (LSTM), ensemble learning model, and linear regression are used as the input models for load prediction in this paper.

1) ARIMA: It is a time series analysis model, which is fitted to time series data either for better understanding the data or for predicting future points in the time series. When the trend change (T), cyclic change (C), seasonal change (S) and irregular change (I) are used to characterize the time features, the time sequences can be described as

$$Y_t = f(T, C, S, I) = T_t + C_t + S_t + I_t. \quad (3)$$

As the ARIMA model has certain requirements for data stability, if considerable changes happen in the load distribution, the model may cause the forecast deviation. Hence, the data could be filtered based on the data stability, so as to select the cells with better response to ARIMA, thereby improving the accuracy of prediction. ARIMA is generally denoted as

$$ARIMA(p, d, q), p, q \in \{0, 1, 2, 3\}, d \in \{0, 1\}, \quad (4)$$

where p is the order of the autoregressive model, q is the order of the moving-average model, p and q are determined by the lowest Bayesian information criterion (BIC), and d is the degree of differentiation to make the data stationary.

2) Prophet: The Prophet model, which is similar to ARIMA mode, is expressed as

$$Y_t = f(g_t, s_t, h_t, e_t), \quad (5)$$

where g_t denotes non-periodic changes, such as linear growth or logical growth; s_t is cyclic changes, like seasonality; h_t is irregular changes caused by users; e_t is the error used to describe the abnormal changes in the model.

3) RF: The preliminary of the random forest prediction model is the decision tree learning that segments the features based on their characteristics. Combining the random subspace method with the decision tree, the RF model selects the features to enhance the prediction, increasing the correlation among the se-

lected features. The model exploits the historical loads to predict future loads; in this way, loads in the past and neighbor loads are taken into consideration when constructing the model.

4) LSTM: It is an artificial recurrent neural network (RNN) architecture used in the field of deep learning. Different from the standard convolution neural network, LSTM could process the single data points like images, as well as sequences of data such as video or speech. The prediction system in this paper is composed of three layers, two LSTM layers and one fully connection layer.

5) Ensemble learning: This mode combines multiple learning algorithms to achieve better performance for a particular intelligence problem. In other words, ensemble learning can combine several weak models that get poor prediction to produce a strong learning model. While some simple models only learn part of the data, the ensemble method can strategically divide the data set into small data sets, train them separately, and then combine them with certain strategy.

In order to compare the algorithms mentioned above, the mean absolute error (MAE) is used to measure the difference between the forecast and the real load. It can be described as

$$MAE = \frac{\sum_{i=1}^N |P_i - R_i|}{N} = \frac{\sum_{i=1}^N |e_i|}{N}, \quad (6)$$

where N is the number of points, P is the predicted load output by an algorithm, and R is the real load. The intuitive meaning of the function MAE is quite clear: the greater the distance between the predicted value P and the true value R , the larger the loss, and vice versa.

6) Linear regression: It is a linear algorithm to map the relationship between a scalar response and one or more explanatory variables. In the linear regression, unknown model parameters are also estimated from the data. If the goal is to predict or forecast the state, the linear regression is able to fit a predictive model to an observed data. Given a data set $\{y_i, x_{i1}, \dots, x_{ip}\}_{i=1}^n$, the linear regression model can be expressed as:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, i = 1, \dots, n, \quad (7)$$

where T denotes the transpose, $\mathbf{x}_i^T \boldsymbol{\beta}$ is the inner product between vectors \mathbf{x}_i and $\boldsymbol{\beta}$.

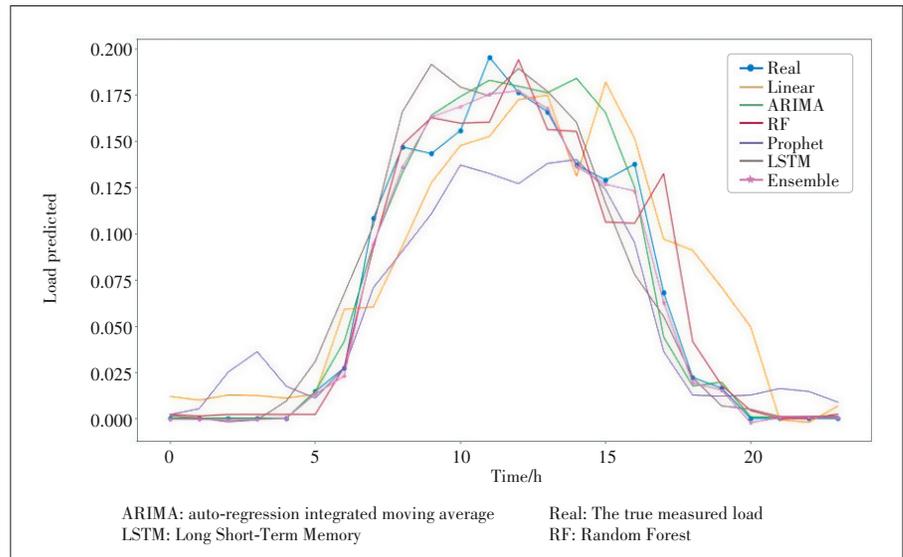
Fig. 14 compares the load prediction in one cell with different prediction models. The simulation results of load prediction are based on the physical resource block (PRB) utilization in 50 cells. It can be seen that the ensemble learning model has further improved the prediction accuracy compared to each independent sub-model. The average MAE of the ensemble learning method is reduced by an average of 0.008.

The comparison and analysis results of the machine learning models mentioned above is listed in **Table 3**. These different load prediction models are suitable for dealing with differ-

ent radio access networks.

We evaluate the application of machine learning techniques with real scenarios, and time efficiency is taken into account, ARIMA is implemented in our ML based ES scheme. There are 1 089 cells and 329 base stations (BSs) in the test area. Different switch-off strategies including the symbol switch-off, channel switch-off and carrier switch-off are applied for different groups of the measured BSs. A cell is considered as the switch-off as the cell carrier(s) is/ are all switched-off, and the BS shall indicate the switch-off action to the BS providing basic coverage. In the real deployment, different BS types may cause different power consumption, so power saving would be averaged to all cells of the measured BS groups. As expected, the artificial intelligence (AI) energy saving scheme predicts load accurately and switch off the cell in time to achieve better performance on energy saving. In addition, the actual energy saving of each cell per day is also significantly increased. **Fig. 15** shows that the AI based power saving could reach up to 1.24 kWh each cell per day, and no matter what switch-off strategy is used, AI based ES is a better solution to power saving.

Table 4 shows the power consumption and electricity charge saving with different kinds of ES methods and without ES. We can see that the power consumption is totally 25 988 kWh every week if any ES method is not used, while the power consumption with the AI ES methods is 22 304 kWh. Electricity charge saving with the AI ES methods increases more than that with the

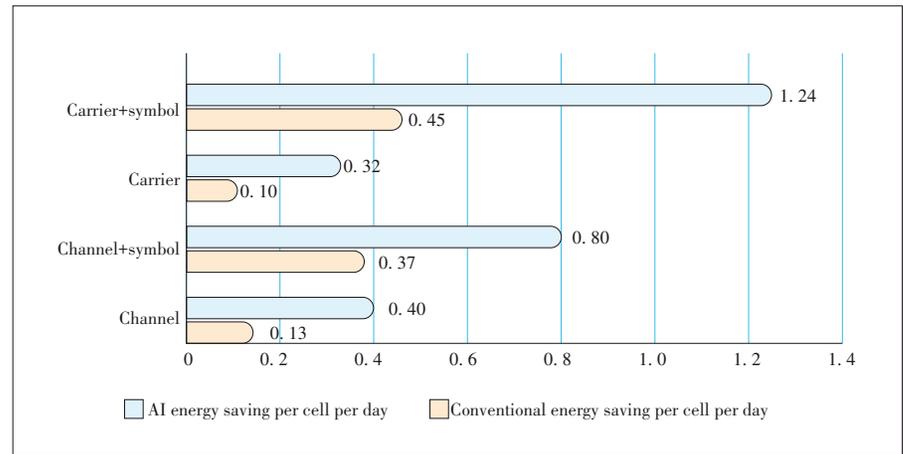


▲ **Figure 14.** Load prediction by using different models

▼ **Table 3.** Comparison and analysis of the machine learning models

Model	Accuracy	Speed	Complexity
ARIMA	Medium	Fast	Low
Prophet	Medium	Fast	Low
LSTM	High	Slow	High
RF	High	Slow	High
Ensemble	High	Extremely slow	High

ARIMA: auto-regression integrated moving average LSTM: Long Short-Term Memory RF: Random Forest



▲ **Figure 15.** Statistics of the power saving (kWh)

▼ **Table 4.** Comparison of power consumption and electricity charge saving with/without ES methods

Switch-off Strategy	Number of Measured Cells	Power Consumption of Measured Cells (kWh/Week)			Electricity Charge Saving of Measured Cells (CNY/Week)		
		No ES	Conventional ES	AI ES	Conventional ES	AI ES	Increase
Carrier	8	382	377	364	5	18	13
Carrier+symbol	7	366	344	305	22	61	39
Channel	633	16 853	16 265	15 872	588	981	393
Channel+symbol	327	8 387	7 541	6 555	846	1 832	986
Total	975	25 988	24 527	22 304	1 461	3 684	2 223

AI: artificial intelligence ES: energy saving

conventional ES, and the saving is totally increased by 2 223 CNY every week.

8 Conclusions and Future work

In this paper, we introduce the 3GPP energy saving schemes by switching on/off cells in ultra-dense networks. In order to achieve intelligent ES, we also propose a machine learning based ES scheme by switching off the cell selected based on load prediction. How to apply AI into the 5G network is a new topic for 3GPP and future works might focus on potential solutions to smart energy saving based on AI and the corresponding 3GPP standard impacts on data collection and interface between NG-RAN nodes.

References

- [1] 3GPP. Technical specification—energy efficiency of 5G release 16: TS28.310 V16.2.0 [S]. 2020
- [2] ZTE. Consideration on inter-RAT energy saving: 3GPP R3-191453 [R]. 2019
- [3] Qualcomm. Inter-system inter-RAT energy saving: 3GPP R3-204802 [R]. 2020
- [4] 3GPP. Technical specification—S1 application protocol release 16: TS36.413 V16.2.0 [S]. 2020
- [5] 3GPP. Technical specification—NG application protocol release 16: TS38.413 V16.3.0 [S]. 2020
- [6] 3GPP. Technical specification—5G end to end key performance indicators release 16: TS 28.554 16.5.0 [S]. 2020
- [7] GAO Y, CHEN J, LIU Z, et al. Machine learning based energy saving scheme in wireless access networks [C]//16th International Wireless Communications and Mobile Computing Conference (IWCMC). Limassol, Cyprus: IEEE, 2020: 1573 – 1578

Biographies

LIU Zhuang (liu.zhuang2@zte.com.cn) received the master's degree in computer science from Xidian University, China in 2003. He is currently a 5G senior research engineer at the R&D center of ZTE Corporation and the State Key Laboratory of Mobile Network and Mobile Multimedia, China. His research interests include 5G wireless communications and signal processing. He has filed more than 100 patents.

GAO Yin received the master's degree in circuit and system from Xidian University, China in 2005. She has been engaged in the study of 4G/5G technology since 2005 and is currently a wireless expert and project manager at the R&D center of ZTE Corporation and the State Key Laboratory of Mobile Network and Mobile Multimedia, China. She has authored or co-authored about hundreds of proposals for 3GPP meetings and journal papers in wireless communications and has filed more than 200 patents. In August 2017, she was elected as 3GPP RAN3 Vice Chairman.

LI Dapeng received the master's degree in computer science from University of Electronic Science and Technology of China in 2003. He is currently a senior researcher at the R&D center of ZTE Corporation and mainly focuses on research and implementation of wireless access network systems.

CHEN Jiajun received the master's degree in electronics and communications engineering from Shanghai University, China in 2019. He has been a technology pre-research engineer at the R&D center of ZTE Corporation. His research interests include next-generation radio access network and deep learning.

HAN Jiren received the master's degree in wireless communication systems from University of Sheffield, UK in 2016. He is currently a technology pre-research engineer at the R&D center of ZTE Corporation. His research focuses on next-generation radio access networks.

Cluster Head Selection Algorithm for UAV Assisted Clustered IoT Network Utilizing Blockchain



LIN Xinhua, ZHANG Jing, LI Qiang

(Huazhong University of Science and Technology, Wuhan 430074, China)

Abstract: To guarantee the security of Internet of Things (IoT) devices, the blockchain technology is often applied to clustered IoT networks. However, cluster heads (CHs) need to undertake additional control tasks. For battery-powered IoT devices, the conventional CH selection algorithm is limited. Based on the above problem, an unmanned aerial vehicle (UAV) network assisted clustered IoT system is proposed, and a corresponding UAV CH selection algorithm is designed. In this scheme, UAVs are selected as CHs to serve IoT clusters. The proposed CH selection algorithm considers the maximal transmit power, residual energy and distance information of UAVs, which can greatly extend the working life of IoT clusters. Through Monte Carlo simulation, the key performance indexes of the system, including energy consumption, average secrecy rate and the maximal number of data packets received by the base station (BS), are evaluated. The simulation results show that the proposed algorithm has great advantages compared with the existing CH selection algorithms.

Keywords: cluster head selection; unmanned aerial vehicle; blockchain; IoT; average secrecy rate

DOI: 10.12142/ZTECOM.202101005

<https://kns.cnki.net/kcms/detail/34.1294.TN.20210223.1206.002.html>, published online February 20, 2021

Manuscript received: 2020-12-10

Citation (IEEE Format): X. H. Lin, J. Zhang, and Q. Li, "Cluster head selection algorithm for UAV assisted clustered IoT network utilizing blockchain," *ZTE Communications*, vol. 19, no. 1, pp. 30 - 38, Mar. 2021. doi: 10.12142/ZTECOM.202101005.

1 Introduction

With the development of fifth-generation (5G) networks, which provide extended coverage, higher throughput, lower latency and higher connection density with a massive bandwidth, 5G based Internet of Things (5G-IoT) devices have emerged as small-size, low-cost, typically battery-powered and densely distributed devices to support large-scale information exchange. Therefore, the 5G-IoT is a core component of the future network. On the one hand, the evolution of 5G networking not only has paved the way for the connection of massive IoT nodes to the Internet to facilitate the advancement of various IoT applications from theory to reality, but also has led to the proposal of various potential technologies, such as millimeter-wave, massive multi-

ple-input multiple-output and device-to-device. On the other hand, over 75 billion devices will be connected to the IoT by 2025, which is expected to have a dramatic impact on our lives in the near future^[1]. This will be beneficial for supporting networks in generating enormous amounts of information traffic, enabling humans to obtain messages about anything and anyone at any time and any place (4A)^[2]. Despite the fruitful developments in 5G-IoT communications, several issues that hamper effective IoT communication in 5G networks remain unsolved, including redundancy in data, dynamic size of the network, less reliable medium, heterogeneous network, and multiple base stations (BSs) or sink nodes. To process data in a distributed way, remove redundant data and improve the energy efficiency, the IoT system needs to adopt clustering

technology^[3]. Clustering builds a hierarchy of clusters or groups of sensing nodes that collects and transfers the data to its respective cluster heads (CHs). The CH then groups and sends the data to the sink node or BS. The CHs act as middleware between the end user and the network, so the selection of CHs is particularly important^[4].

Due to the limited computing capacity and energy of IoT devices in the process of data transmission, it is difficult to adopt highly complex algorithms and frameworks to ensure the data security. Therefore, IoT devices face many security issues that include the authenticity and confidentiality of data^[1]. Blockchain, which can guarantee the integrity, transparency and security of data in industrial data processing, has attracted great attention in the application of IoT^[5]. Integrating blockchain and IoT has many advantages. Firstly, it can improve resilience and adaptability of the IoT system. Blockchain can store redundant replicas of data in the form of transactions over blockchain nodes, which helps to maintain data integrity and provide resilience to the IoT system. Secondly, since blockchain is a distributed ledger, using blockchain as the data management mechanism for the IoT can adapt to varying environments and use cases to meet the growing needs and demands of IoT devices, which improves the adaptability of the system. Finally, integrating blockchain and IoT can enhance the fault tolerance and security of the whole system. However, due to the verification of blockchain, IoT devices will perform additional computing tasks, which will greatly increase energy consumption and reduce the service life of IoT systems.

Considering the energy limitation of IoT devices, there are many clustering technologies and CH selection algorithms to reduce energy consumption of the IoT system. HEINZELMAN et al. proposed a low-energy adaptive clustering hierarchy (LEACH) protocol^[6], in which CHs were randomly selected in each round. Since the selection of CHs is random, the nodes with low energy are at the same priority as those with high energy. If the nodes with low energy are selected as CHs, they will fail quickly, thus shortening the network life. Based on the LEACH, TRUPTI et al. proposed a CH selection algorithm based on residual energy, which is to choose the devices with more residual energy as the CH^[4]. YOUNIS et al. adopted the hybrid energy efficient distributed clustering (HEED) algorithm, which could select devices with high battery power as the CH through the proposed iterative CH selection algorithm^[7]. In the above works, wireless sensors or IoT devices are selected as CHs. Although the energy limits of devices are considered in these algorithms when selecting CHs, the energy limitations of IoT devices will lead to frequent failure, resulting in more system energy consumption. AADIL et al. proposed energy aware link-based clustering (EALC), which adds two other parameters (energy level and distance) to the neighborhood to select the optimal CH. EALC extends cluster life and reduces energy consumption^[8].

Due to the large difference of devices and limited resources

of the blockchain-based IoT system, which needs to perform additional blockchain computing tasks, the choice of IoT devices as the CHs will have great limitations. The emergence of unmanned aerial vehicles (UAVs) provides new opportunities for the blockchain-based IoT system. When UAVs are used as flying BSs, they can support the connectivity of existing ground wireless networks to help land systems achieve good coverage and effectively reduce the data traffic of other BSs. Moreover, as devices with flexible deployment, UAVs are equipped with high-performance calculators with high computing capacity, which can quickly respond to the communication and computing needs of IoT devices, thus improving the quality of service^[9]. In addition, solar-powered UAVs can convert solar energy into electric energy, thus increasing its service time^[10]. Moreover, a large number of UAVs can cooperate with each other through relay nodes to build a self-organizing intelligent UAVs network to complete complex tasks^[11]. Therefore, it has great advantages to choose UAVs as CHs.

To solve the problems of limited resources and security faced by IoT clusters, the main contributions of this paper are as follows. Firstly, the UAV network served IoT cluster system is built. To ensure the security of data, the IoT devices in the system use blockchain technology to store data. Secondly, we propose a UAV CH selection algorithm. The algorithm jointly considers the distance between UAVs and IoT devices, the distance between UAVs and BSs, residual energy, and the maximal transmit power of UAVs. The IoT devices calculate the corresponding weighted value of the UAV through the proposed algorithm, and choose the UAV with the smallest weighted value to vote. The UAV with the most votes serves the IoT cluster as the CH. Finally, based on the proposed algorithm, this paper evaluates several performance indicators such as the energy consumption of the IoT cluster, the average secrecy rate and the maximal number of packets received by the BS, and compares the performance with several existing CH selection algorithms, which demonstrates the superiority of the proposed algorithm.

This paper is structured as follows. Section 2 presents our system model and the basic procedure of the practical Byzantine fault tolerance (PBFT) consensus algorithm. Moreover, we use received signal strength (RSS) technology to estimate the distance between IoT devices and UAVs in this section. In Section 3, we propose a UAV selection algorithm based on a private blockchain and introduce performance evaluation indicators. The simulation results are analyzed in Section 4. Finally, Section 5 concludes this paper.

2 System Model

2.1 Network Topology

To enable secure energy-efficient communication, a blockchain-based CH selection algorithm is proposed in this paper.

As shown in **Fig. 1**, a blockchain-based clustered IoT network, where a UAV swarm composed of M UAVs, denoted by set $\{U_1, U_2, U_3, \dots, U_M\}$, is deployed to serve S IoT clusters, denoted by set $\{C_1, C_2, C_3, \dots, C_S\}$. Each IoT cluster contains K IoT devices, denoted by set $\{D_1, D_2, D_3, \dots, D_K\}$. To reduce the power consumption of the IoT devices during data transmission, UAVs hovering in the sky collect data from the IoT devices before transmitting the collected data to the BS. Meanwhile, the IoT clusters adopt private blockchain technology to protect their collected data and to facilitate secure communication. Eavesdroppers coexisting with the IoT clusters may intercept the transmitted data. We assume that the IoT devices and the UAVs can establish Line-of-Sight (LoS) communication links for data transmission. In contrast, an eavesdropper may experience a Rayleigh fading channel while eavesdrop on the IoT devices. The UAVs first broadcast a message containing the pilot signal and the UAV information to the IoT devices. The IoT devices in each cluster use blockchain technology to verify the information received from UAVs and estimate the distance to each UAV. Using the proposed CH selection algorithm, the IoT devices in each cluster then vote through the PBFT consensus algorithm. According to the voting results, the UAV that receives the most votes from the cluster is selected as the CH. When different clusters choose the same UAV as the CH, it is assumed that when the energy of the UAV selected by a cluster is exhausted, that cluster will select a different UAV. The IoT devices in each cluster communicate using orthogonal frequency division multiple access (OFDMA) technology, with a system bandwidth of B Hz.

2.2 Message Broadcasted by UAVs

As the first step of the blockchain-based CH selection process, the M UAVs first broadcast message I_m to all IoT devices

in the area. The message content includes the serial number of the UAVs, U_m , which ranges from 1 to M for the considered UAV swarm; the maximal transmit power of the UAV, P_m ; the remaining energy in the battery of the UAV E_m ; the distance between UAV m and BS d_{BU}^m , which is estimated at the BS by measuring the pilot signal of the UAV. The IoT devices can obtain the message of the UAV from the mark bit named Mark Signa.

2.3 Distance Estimation Based on RSS

In practice, the channel state information (CSI) between IoT devices and UAVs is unknown. To evaluate the CSI, the IoT devices usually adopt distance estimation. On the one hand, a UAV flying in the sky can provide a LoS propagation environment, which is beneficial for evaluating the distances between the IoT devices and the UAV. On the other hand, the IoT devices are powered by batteries and have a simple hardware structure, and it is difficult for them to perform complex signal processing to obtain the CSI. Therefore, distance estimation based on RSS of an IoT device is suitable for UAV-assisted IoT communication^[12].

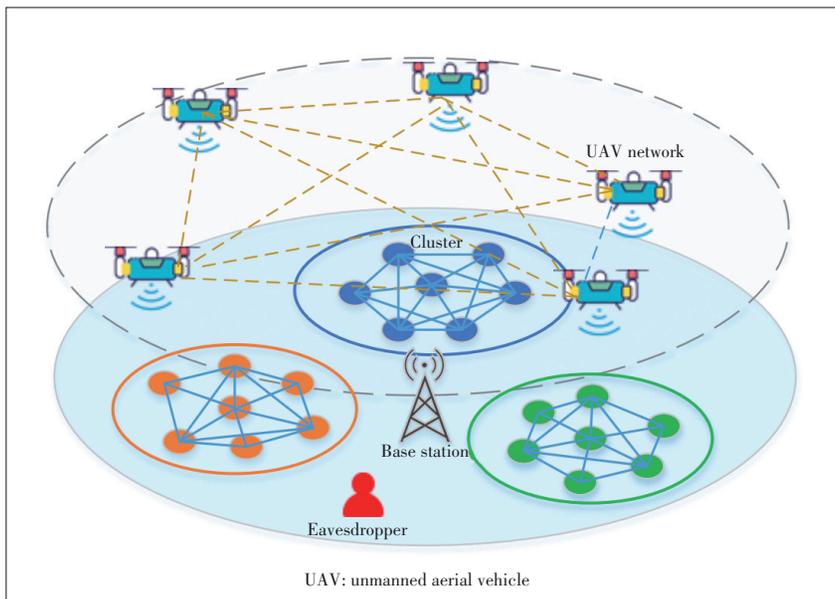
In particular, after the messages broadcast by the UAVs are received by the IoT devices, each IoT device will estimate the distance to each UAV based on the RSS. By utilizing the RSS at the IoT device and the information of the UAV's maximal transmit power provided in the broadcast message, the distance between an IoT device and a UAV can be formulated through maximum likelihood estimation as follows:

$$\hat{d}_{mk} = \left(\frac{P_{mk}}{P_m} \right)^{\frac{-1}{n_p}}, \tag{1}$$

where P_{mk} is the signal strength received by IoT device k from UAV m , i.e., the RSS; P_m is the maximal transmit power of UAV m , which is specified in the broadcast message; n_p is the path loss factor.

2.4 Private Blockchain Constructed for IoT Clusters

To ensure data security, the IoT devices adopt a private blockchain to verify their collected data. Specifically, the IoT devices transmit their received messages broadcast from the UAVs to other IoT devices in the same cluster as transactions and then apply the PBFT algorithm to reach agreement. The consensus data will be stored in blocks in the form of transaction, and each block contains the hash code of the previous block, thus forming a blockchain. To ensure the privacy and security of the data in the PBFT process, a hash algorithm and an asymmetric encryption



▲ Figure 1. UAV assisted clustered IoT network utilizing blockchain

algorithm are introduced. We use the elliptic curve encryption (ECC) algorithm and Secure Hash Algorithm-256 (SHA-256) to detect whether transactions have been tampered with during data transmission^[13]. The encryption process is shown in **Fig. 2**. By using SHA-256, we can generate a Merkle root, which can be used to effectively compress the amount of data to link each block. When using encrypted data, the IoT devices can perform the same hash calculation and compare the hash codes to verify the data. The ECC algorithm is used to generate public and private keys to encrypt the data. The data in the database will be encrypted into ciphertext by using the public key. Each user should provide his or her own private key to decrypt the encrypted message used for the custom service. These two algorithms can ensure the privacy of the data and prevent illegal operations. The consensus process of PBFT is shown in Fig. 2 which contains the following phases:

- Request phase: We refer to the IoT device that needs to publish transactions as a client. In our model, IoT devices not on-

ly act as the publisher of transactions, but also as the verifier of transactions. Before IoT devices transmit data to other nodes for verification, they need to encrypt the data.

- Pre-prepare phase: After receiving the message from the client, the primary node will assign an integer sequence number to the request, and then generate the pre-prepare message. The primary node then broadcasts the pre-prepared message to replica nodes.
- Prepare phase: The replica nodes verify the message that has not been tampered with, and then send a prepare message to other nodes.
- Commit phase: After verifying that all the prepared messages have not been tampered with, all nodes will broadcast the confirm message to other nodes.
- Reply phase: After verifying the message, all nodes will return the result to the client.

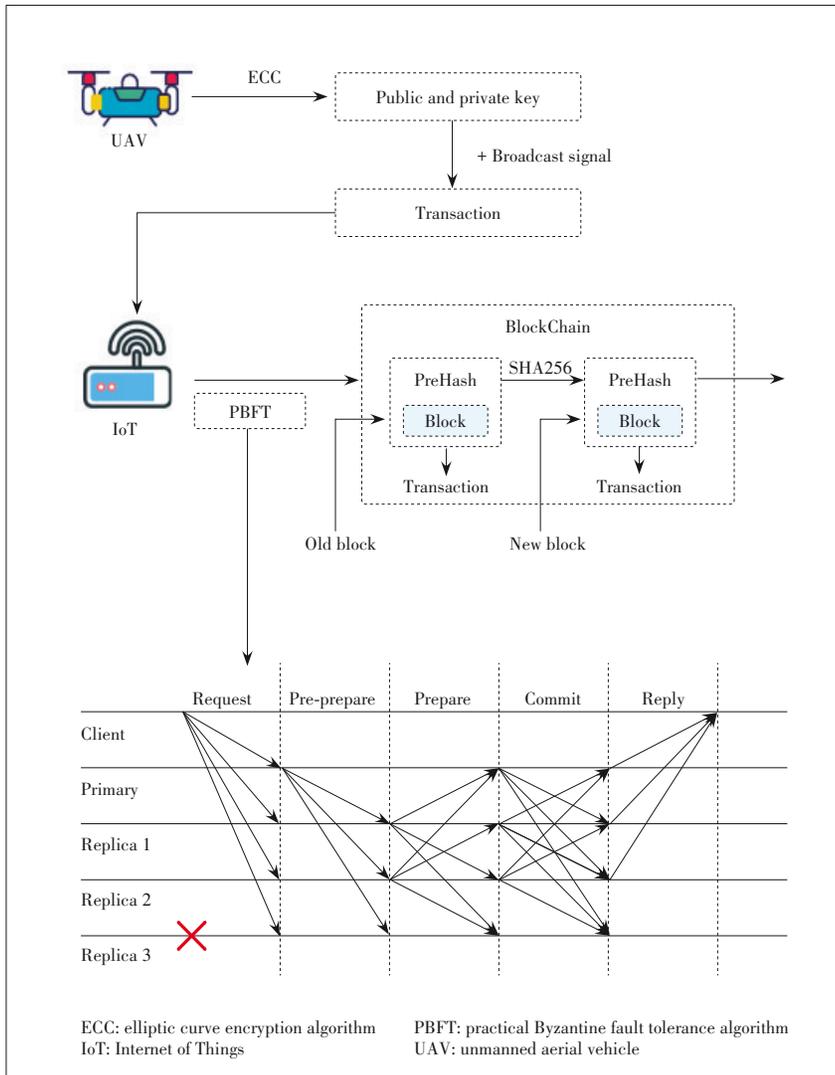
2.5 System Performance Metrics

To analyze the performance of the CH selection algorithm proposed in this paper, several important system parameters will be used. The CH selection process consists of four steps: An IoT cluster receives the broadcast messages from the UAVs, each IoT device chooses a UAV to serve as the CH based on the received messages, the IoT devices in the IoT cluster achieve consensus through the PBFT algorithm, and the IoT cluster sends a response message to the UAV swarm indicating the chosen CH. We utilize the energy consumption for CH selection to represent the system resource consumption of the IoT cluster. Meanwhile, system performance metrics, i. e., the average secret rate and the maximal number of received packets, will be used to evaluate the gain of our proposed algorithm.

- CH selection energy consumption: CH selection delay refers to the time taken by a cluster from receiving UAV signal to reaching consensus and finally sending the selection result to the selected UAV, which reflects the effectiveness of CH selection algorithm. Large time delay will lead to large energy consumption, which will affect the life cycle of the device.

- Average secrecy rate: The secrecy rate is a key design metric for IoT networks that is widely adopted for evaluating physical layer security is the secrecy rate. High secrecy rate will reduce the probability of data eavesdropping.

- Maximal number of packets received by the BS: The number of packets received by the BS reflects the throughput of the system, which is a very important measure of the system.



▲ Figure 2. Private blockchain constructed for IoT clusters and PBFT process

3 Working Model

3.1 UAV Selection Algorithm

Based on the messages broadcast by the UAVs, all IoT devices in a cluster adopt the PBFT algorithm to verify their received messages to achieve consensus. Then, the IoT cluster chooses a UAV from the UAV swarm as its CH by following steps.

- Step 1: According to Eq. (3), each of the K IoT devices in the IoT cluster estimates the distance d_{mk} from each UAV based on the RSS.

- Step 2: The distance between the BS and the m -th UAV, denoted by d_{BU}^m , can be determined from the broadcast information sent by the UAV.

- Step 3: From the broadcast messages sent by the UAVs, the IoT devices are informed of the remaining energy of each UAV, E_m . For each UAV, the energy ratio of the total remaining energy of all UAVs to the remaining energy of that UAV is

$$E_p = \frac{\sum_{m=1}^M E_m}{E_m}. \quad (2)$$

- Step 4: From the messages broadcast by the UAVs, the IoT devices are also informed of the maximal transmit power of each UAV, P_m . For each UAV, the ratio of the total maximal output power of all UAVs to the transmit output power of that UAV is

$$P_p = \frac{\sum_{m=1}^M P_m}{P_m}. \quad (3)$$

- Step 5: The weighted value of each UAV is computed as follows:

$$F_k = \alpha \hat{d}_{mk} + \beta d_{BU}^m + \varsigma E_p + \theta P_p, \quad (4)$$

where α, β, ς and θ are weighting factors that satisfy $\alpha + \beta + \varsigma + \theta = 1$.

- Step 6: Each IoT device calculates its corresponding weighted value F_k for each UAV following the above method. Then, the k -th IoT device votes for the UAV with the smallest F_k to serve as the CH. All IoT devices in the same IoT cluster use the PBFT algorithm to vote for consensus. Finally, the UAV with the most votes is chosen to serve the entire cluster. The proposed CH selection process is presented in Algorithm 1.

Algorithm 1. Proposed UAV CH selection algorithm

Input: UAVs $U_m (m \in 1, 2, \dots, M)$ blockchain-based IoT clusters $C_s (s \in 1, 2, \dots, S)$ and IoT devices D_k in each cluster, $D_k (k \in 1, 2, \dots, K)$.

Output: UAVs chosen as the CH, $CH_s (s \in 1, 2, \dots, S)$.

/* Initialization Phase */

Assign each UAV m transmit power P_m and residual energy E_m .

Assign the distance between UAV m and BS d_{BU}^m .

Assign the weighting factors $\alpha, \beta, \varsigma, \theta$.

/* Computation Phase */

While $(k++ < K+1)$ **do**

for each IoT device $D_k (k \in 1, 2, \dots, K)$ **do**

 Estimate the distance between IoT device k and each UAV m , d_{mk} , using (1)

end for

for each IoT device $D_k (k \in 1, 2, \dots, K)$ **do**

 Measure the energy ratios of the UAVs, E_p using (2).

end for

for each IoT device $D_k (k \in 1, 2, \dots, K)$ **do**

 Measure the maximal transmit power ratios of the UAVs, P_p using (3).

end for

 Calculate F_k using (4).

 Vote for the optimal UAV with smallest F_k .

end while

return the UAV with the most votes

3.2 Performance Metrics

3.2.1 Energy Consumption for CH Selection

In our system, the energy consumption of an IoT device mainly includes three components: the energy consumption for data transmission, E_{total}^{tx} ; the energy consumption for data reception, E_{total}^{rx} ; and the energy consumption for computing using the PBFT algorithm, E_{total}^c . For a UAV swarm composed of M UAVs and an IoT cluster with K IoT devices, each IoT device in the cluster will transmit $3K$ transactions, receive $(M + 2K - 1)$ transactions, and perform $(2K - 1)$ computing operations during the PBFT process.

The energy consumption of each IoT device during the process of transmitting transactions is calculated as^[14]

$$E_{total}^{tx} = \begin{cases} 3KL \cdot (E_{elec} + \varepsilon_{fs} d_{mk}^2), & d_{mk} < d \\ 3KL \cdot (E_{elec} + \varepsilon_{fs} d_{mk}^4), & d_{mk} \geq d \end{cases}, \quad (5)$$

where E_{elec} is the energy dissipated per bit to run the transmitter or receiver circuit, $\varepsilon_{fs} d_{mk}^2$ and $\varepsilon_{fs} d_{mk}^4$ are the energy cost of a single amplifier under the two communication models depending on the distance between the transmitter and receiver, and d is the threshold value.

The energy consumption of each IoT device during the process of receiving transactions is calculated as

$$E_{total}^{rx} = L \cdot E_{elec} (M + 2K - 1). \quad (6)$$

The energy consumption of each IoT device during the process of verifying transactions is calculated as

$$E_k^c = k_m s_m f_k^2 LM (2K - 1), \quad (7)$$

where s_m is the number of rotations required to calculate 1 bit of data, k_m is the calculation efficiency and f_k is the computing capacity of the k -th IoT device.

Considering that there are K IoT devices in the IoT cluster, the total energy consumed by all devices in the cluster for PBFT processing is calculated as

$$E_{total}^c = \sum_{k=1}^K E_k^c = \sum_{k=1}^K k_m s_m f_k^2 LM (2K - 1). \quad (8)$$

In the end, the total energy consumption of the IoT cluster for selecting the m -th UAV in the UAV swarm as the CH is

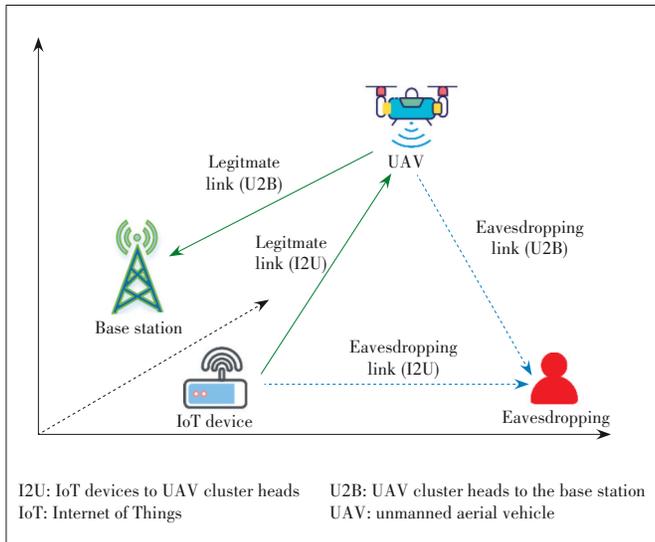
$$E_k^s = E_{total}^c + E_{total}^{tx} + E_{total}^{rx}. \quad (9)$$

For all S IoT clusters, the total energy consumption is obtained as follows:

$$E_t = \sum_{s=1}^S E_k^s. \quad (10)$$

3.2.2 Average Secrecy Rate of IoT Clusters

A diagram of the system secrecy rate is shown in **Fig. 3**. When the strength of legitimate links is greater than that of eavesdropping links, a nonzero secrecy rate will be achieved^[15]. When the m -th UAV in the UAV swarm is chosen as the CH for an IoT cluster, the information transmitted from the IoT devices to the UAV CH and from the UAV CH to the BS can be eavesdropped on. Therefore, the secrecy rates of both types of links are analyzed in the following.



▲ **Figure 3. Diagram of system secrecy rate**

For the information transmission from IoT device k to the UAV CH, i.e., UAV m , the achievable rate is obtained as follows:

$$R_{km} = \log \left(1 + \frac{P_k \beta_0}{\sigma^2 d_{mk}^2} \right), \quad (11)$$

where P_k is the transmission power of IoT device k , β_0 is the signal gain at a distance $d_0 = 1$ m, σ^2 is the noise power and d_{mk} is the distance between IoT device k and the chosen UAV m .

For the information transmission from UAV CH to the BS, the achievable rate is obtained as follows:

$$R_{mb} = \log \left(1 + \frac{P_m \beta_0}{\sigma^2 (d_{BU}^m)^2} \right). \quad (12)$$

When there is an eavesdropper, the transmission rate from an IoT device to the eavesdropper is calculated as

$$R_{ke} = \log \left(1 + \frac{P_k \beta_0}{\sigma^2 d_{ke}^\gamma} \right), \quad (13)$$

where d_{ke} is the distance between IoT device k and the eavesdropper, and γ is the path loss exponent.

For the eavesdropping link from a UAV to the eavesdropper, the transmission rate is calculated as

$$R_{ue} = \log \left(1 + \frac{P_m \beta_0}{\sigma^2 d_{me}^\gamma} \right), \quad (14)$$

where d_{me} is the distance between UAV m and the eavesdropper.

The average secrecy rate of an IoT cluster is

$$R_{sec}^{av} = \frac{\sum_{k=1}^K [R_{km} + R_{mb} - (R_{ke} + R_{ue})]^+}{K}, \quad (15)$$

where $[x]^+ = \max(x, 0)$.

For all the S IoT clusters, the total average secrecy rate can be calculated as follows:

$$R_{sec}^{total} = \sum_{s=1}^S R_{sec}^{av}. \quad (16)$$

3.2.3 Maximal Number of Packets Received by BS

The number of packets received by the BS is an important indicator of the total throughput for an IoT cluster. When an IoT device transmits data packets comprising L_k bits to its UAV CH in each time slot and the UAV CH retransmits these data packets to the BS, the time consumed for the IoT device to transmit data to the UAV CH is

$$t_k = \frac{L_k}{(B/K) \log \left(1 + \frac{P_k \beta_0}{(B/K) \sigma^2 d_{mk}^2} \right)}. \quad (17)$$

The IoT devices within an IoT cluster utilize the OFDMA scheme to transmit their data packets to the UAV CH, and hence, the total time consumed by the IoT cluster to transmit data to the UAV CH is

$$t'_k = \max(t_k). \quad (18)$$

When the UAV CH retransmits these data packets to the BS, the UAV CH can utilize the whole usable bandwidth, and hence, the consumed time is

$$t_u = \frac{KL_k}{B \log \left(1 + \frac{P_m \beta_0}{B \sigma^2 (d_{BU}^m)^2} \right)}. \quad (19)$$

Meanwhile, the UAV also consumes propulsion power to support it as it flies in the sky. It should be noted that if the UAV uses up its energy between communication and propulsion, the UAV will be unable to retransmit the data packets, and the CH will break down. According to Ref. [16], the power consumed for propulsion is calculated as

$$P_v = \frac{\delta_d}{8} \rho s A \Omega^3 R^3, \quad (20)$$

where δ_d is the profile drag coefficient, ρ is the air density, s is the robustness of the rotor, A is the area of the rotor, and R is the radius of the rotor.

Each time a data packet is sent, the energy consumption of the UAV is

$$E_u = P_v \cdot (t_k + t_u) + KL_k \cdot E_{Rx} + P_m t_u. \quad (21)$$

The maximal number of packets that UAV m can transmit is

$$n_m^s = \frac{E_m}{E_u}. \quad (22)$$

Thus, the maximal number of packets received by the BS is

$$n_m^{total} = \sum_{s=1}^S n_m^s. \quad (23)$$

4 Performance Analysis and Simulation Results

In this section, numerical results are pre-

sented for evaluating the performance of the proposed CH selection algorithm. We compare our proposed CH selection scheme with other existing CH selection schemes, such as Low-Energy Adaptive Clustering Hierarchy (LEACH)^[6], Hybrid Energy-Efficient Distributed Clustering (HEED)^[7] and Energy Aware Link-Based Clustering (EALC)^[8]. To illustrate the advantages of our proposed algorithm for blockchain-based IoT clusters, we use the existing CH selection algorithms as baselines for selecting UAV CHs and compare the system performance in each case with that achieved using the algorithm proposed in this paper. The simulation parameters are shown in **Table 1**.

4.1 Analysis of Security

1) Data trustworthiness: Data trustworthiness greatly affects the security of the collected data. Malicious nodes may insert fake data into the network and interfere with normal nodes, which may cause node failure. Our system uses the PBFT consensus algorithm, which has an error tolerance rate of $(N - 1)/3$. As long as the number of failed nodes does not exceed this tolerance value, the system's data can be transmitted once the correct consensus has been reached, which can effectively guarantee the credibility of the data.

2) Privacy: Privacy is extremely important to the system. If a user's private information is leaked, this may result in enormous losses. Our system uses private blockchain technology; thus, devices will be authenticated by blockchain, and data will be stored in the blocks in the form of transactions. A block cannot be tampered with or deleted, thereby guaranteeing the undeniability and confidentiality of the data. Our proposed blockchain-based CH selection algorithm does not require the intervention of a trusted third party, thereby ensuring the robustness and privacy of the system.

▼Table 1. Simulation parameters

Parameter	Value
Network area	100 m×100 m
Number of UAVs	5 - 30
Total number of IoT devices	10 - 150
The number of IoT in a cluster	4 - 20
UAV transmit power P_m	2 - 4 W
UAV remaining energy E_m	400 - 900 kJ
IoT transmit power P_k	0.5 - 1.5 W
Computational capability of an IoT device f_k	0.1 GHz CPU cycles/bit
Computational energy efficiency coefficient of the processors chip in an IoT device k_m	10^{-26}
Computation workload/intensity s_m	18 000 CPU cycles/bit
Sizes of transaction L	256 bit
Size of a packet transmitted by an IoT device L_k	4 000 bit
Noise power, σ^2	-100 dBm
Weighting factors, $\alpha, \beta, \varsigma, \theta$	0.3, 0.2, 0.3, 0.2

IoT: Internet of Things UAV: unmanned aerial vehicle

4.2 Analysis of Energy Consumption of IoT Devices

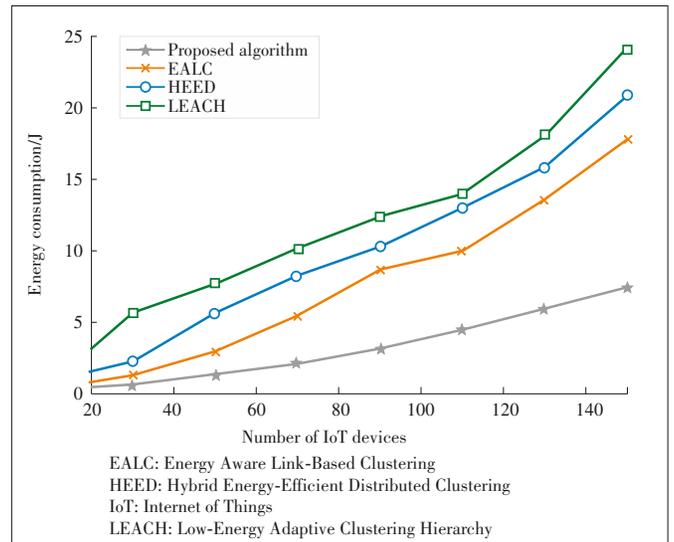
Fig. 4 plots the energy consumption versus the number of IoT nodes for the various CH selection algorithms. The number of UAVs is fixed at 6. From Fig. 4, it can be seen that our proposed strategy achieves the minimal energy consumption compared with the other existing CH algorithms. Meanwhile, as the number of IoT devices increases, the gaps between our proposed algorithm and the other existing algorithms become larger. This gap enlargement occurs because the average distance between the IoT devices and the UAVs is small and the uplink transmission between the IoT devices and the UAV CHs can take advantage of the LoS channel environment in our proposed algorithm. Therefore, our proposed algorithm incurs significantly less energy consumption for communication than the other schemes do. Moreover, compared with the typical random selection algorithms LEACH and HEED, which select CHs by means of multiple votes and therefore cause the IoT devices to consume more energy, our proposed algorithm needs each device to vote only once; hence, the energy consumed to reach consensus within an IoT cluster is reduced.

4.3 Analysis of Average Secrecy Rate

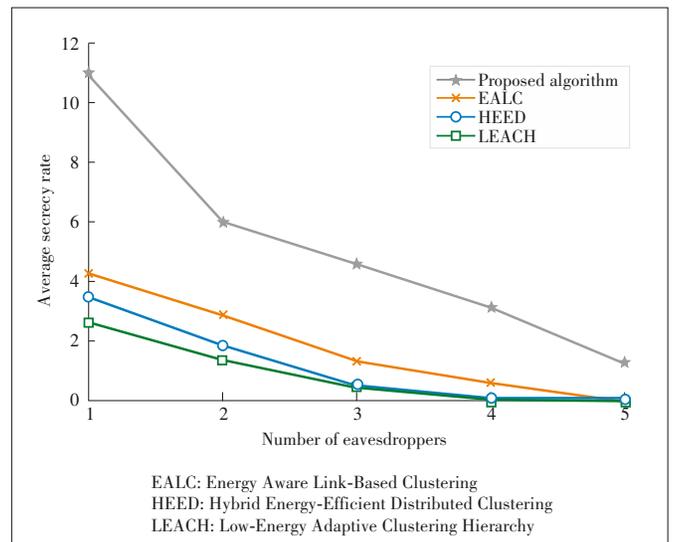
Fig. 5 shows the average secrecy rates achieved with the different CH selection algorithms versus the number of eavesdroppers. The number of UAVs is fixed at 6, and the number of IoT devices is 50. The figure shows that as the number of eavesdroppers increases, the average secrecy rate decreases. The presence of more eavesdroppers will cause the eavesdropping rate for an IoT cluster to increase. Hence, the secrecy rate of the IoT cluster will inevitably decrease. As shown in Fig. 5, the average secrecy rate of our proposed algorithm is significantly better than those of the other existing CH algorithms. This is because the distance and transmit power of each UAV are considered in our proposed algorithm. In this way, both the legitimate transmission rate and the secrecy rate of the IoT clusters can be increased.

4.4 Analysis of the Maximal Number of Packets Received by BS

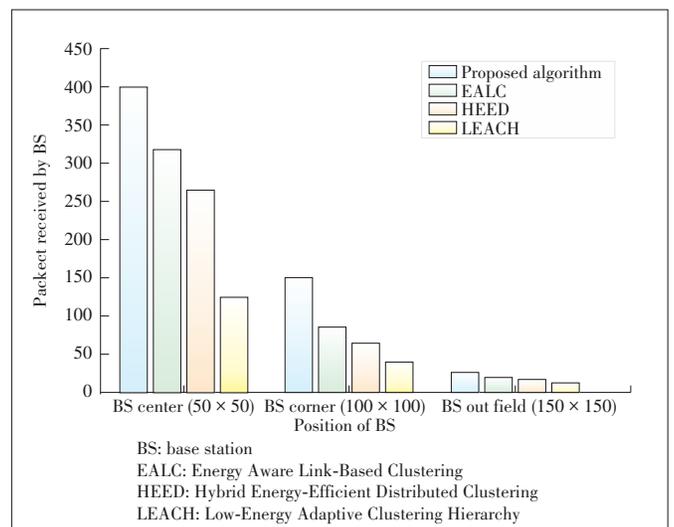
For 150 IoT devices and 6 UAVs, **Fig. 6** analyzes the number of data packets sent by the IoT devices and received by the BS. We compare the maximal numbers of data packets that the BS can receive at three typical locations: the center of the IoT network, the corner of the IoT network, and the outside of the IoT network. Our proposed algorithm performs significantly better than the other algorithms in terms of this metric. This is because the maximal number of data packets received at the BS strongly depends on the energy of the CHs and IoT devices. Less energy consumption of the IoT clusters will lead to a longer lifetime of the IoT network, and hence, more data packets can be transmitted in the system. Our proposed algorithm can reduce the energy consumed by the IoT devices for communications, including data transmission, data reception,



▲ **Figure 4.** Energy consumption versus the number of IoT devices



▲ **Figure 5.** Average secrecy rate versus the number of eavesdroppers



▲ **Figure 6.** Packets received by BS versus different positions of BS

and data processing. Moreover, the remaining energy of the UAVs is also considered in our proposed CH selection algorithm. Consequently, the proposed algorithm can prolong the lifetime of the system by reducing the probability of UAV CH breakdown. Thus, the number of data packets received by the BS increases.

5 Conclusions

In this paper, we propose a novel UAV CH selection algorithm for IoT clusters based on blockchain technology. Our proposed algorithm considers the combined effect of the distances between the IoT devices and the UAVs, the distances between the UAVs and the BS, the maximal transmission power of the UAVs, and the remaining energy of the UAVs; it has the flexibility to assign different weights to these different contributing factors. Each IoT device votes for its optimal UAV through our proposed CH selection algorithm, and then, all IoT devices in a cluster use blockchain technology to achieve consensus to ensure the correctness and security of the vote data. The UAV with the most votes among the devices in an IoT cluster will act as the CH to serve the IoT cluster. Simulation results illustrate the system performance that are compared with corresponding results of the existing algorithms, such as LEACH, HEED and EALC. The simulation results show that our proposed algorithm outperforms the existing algorithms in terms of the energy consumption of the IoT clusters, the average secrecy rate of the IoT clusters and the maximal number of data packets received by the BS.

References

- [1] BUTUN I, ÖSTERBERG P, SONG H B. Security of the Internet of Things: vulnerabilities, attacks, and countermeasures [J]. *IEEE communications surveys & tutorials*, 2020, 22(1): 616 – 644. DOI: 10.1109/COMST.2019.2953364
- [2] AGIWAL M, ROY A, SAXENA N. Next generation 5G wireless networks: a comprehensive survey [J]. *IEEE communications surveys & tutorials*, 2016, 18(3): 1617 – 1655. DOI: 10.1109/COMST.2016.2532458
- [3] XU L N, COLLIER R, O’HARE G M P. A survey of clustering techniques in WSNs and consideration of the challenges of applying such to 5G IoT scenarios [J]. *IEEE Internet of Things journal*, 2017, 4(5): 1229 – 1249. DOI: 10.1109/JIOT.2017.2726014
- [4] BEHERA T M, MOHAPATRA S K, SAMAL U C, et al. Residual energy-based cluster-head selection in WSNs for IoT application [J]. *IEEE Internet of Things journal*, 2019, 6(3): 5132 – 5139. DOI: 10.1109/JIOT.2019.2897119
- [5] ALI M S, VECCHIO M, PINCHEIRA M, et al. Applications of blockchains in the Internet of Things: a comprehensive survey [J]. *IEEE communications surveys & tutorials*, 2019, 21(2): 1676 – 1717. DOI: 10.1109/COMST.2018.2886932
- [6] HEINZELMAN W B, CHANDRAKASAN A P, BALAKRISHNAN H. An application-specific protocol architecture for wireless microsensor networks [J]. *IEEE transactions on wireless communications*, 2002, 1(4): 660 – 670. DOI: 10.1109/TWC.2002.804190
- [7] YOUNIS O, FAHMY S. HEED: a hybrid, energy-efficient, distributed clustering approach for ad hoc sensor networks [J]. *IEEE transactions on mobile computing*, 2004, 3(4): 366 – 379. DOI: 10.1109/TMC.2004.41
- [8] AADIL F, KHAN M F, MAQSOOD M, et al. Energy aware cluster-based routing in flying ad-hoc networks [J]. *Sensors*, 2018, 18(5): 1413. DOI: 10.3390/s18051413
- [9] MOZAFFARI M, SAAD W, BENNIS M, et al. A tutorial on UAVs for wireless networks: applications, challenges, and open problems [J]. *IEEE communications surveys & tutorials*, 2019, 21(3): 2334 – 2360. DOI: 10.1109/COMST.2019.2902862
- [10] SUN Y, DONGFANG X, NG D W K, et al. Optimal 3D-trajectory design and resource allocation for solar-powered UAV communication systems [J]. *IEEE transactions on communications*, 2019, 67(6): 4281 – 4298. DOI: 10.1109/TCOMM.2019.2900630
- [11] GUPTA L, JAIN R, VASZKUN G. Survey of important issues in UAV communication networks [EB/OL]. (2016-03-28)[2020-10-16]. <https://arxiv.org/abs/1603.08462>
- [12] MAO G, FIDAN B, ANDERSON B D O. Wireless sensor network localization techniques [J]. *Computer networks*, 2007, 51(10): 2529 – 2553. DOI: 10.1016/j.comnet.2006.11.018
- [13] FERNÁNDEZ-CARAMÉS T M, FRAGA-LAMAS P. A review on the use of blockchain for the Internet of Things [J]. *IEEE access*, 2018, 6: 32979 – 33001. DOI: 10.1109/ACCESS.2018.2842685
- [14] ARAFAT M Y, MOH S. Localization and clustering based on swarm intelligence in UAV networks for emergency communications [J]. *IEEE Internet of Things journal*, 2019, 6(5): 8958 – 8976. DOI: 10.1109/JIOT.2019.2925567
- [15] WU Q Q, MEI W D, ZHANG R. Safeguarding wireless network with UAVs: a physical layer security perspective [EB/OL]. (2019-07-24)[2020-10-16]. <https://arxiv.org/abs/1902.02472>
- [16] ZENG Y, XU J, ZHANG R. Energy minimization for wireless communication with rotary-wing UAV [EB/OL]. (2018-04-06)[2020-10-16]. <https://arxiv.org/abs/1804.02238>

Biographies

LIN Xinhua is a graduate student of Huazhong University of Science and Technology, China. His main research interests include UAV communications, blockchain technology and IoT networks.

ZHANG Jing (zhangjing@hust.edu.cn) received the M.S. and Ph.D. degrees in electronics and information engineering from Huazhong University of Science and Technology (HUST), China in 2002 and 2010, respectively. He is currently an associate professor with HUST. He has conducted research in the areas of multiple-input multiple-output, CoMP, beamforming, and next-generation mobile communications. His current research interests include HetNet in 5G, green communications, energy harvesting, IoT network, optimization and performance analysis in networks.

LI Qiang received the Ph.D. degree in electrical and electronic engineering from Nanyang Technological University (NTU), Singapore in 2011. He is currently an associate professor with Huazhong University of Science and Technology (HUST), China. His current research interests include next generation mobile communications, fog computing, edge caching, cognitive radios/spectrum sharing, wireless cooperative communications, full-duplex techniques, simultaneous wireless information and power transfer.

Green Air-Ground Integrated Heterogeneous Network in 6G Era



WU Huici, LI Hanjie, TAO Xiaofeng

(Beijing University of Posts and Telecommunications, Beijing 100876, China)

Abstract: The research of three-dimensional integrated communication technology plays a key role in achieving the ubiquitous connectivity, ultra-high data rates, and emergency communications in the sixth generation (6G) networks. Aerial networking provides a promising solution to flexible, scalable, low-cost and reliable coverage for wireless devices. The integration of aerial network and terrestrial network has been an inevitable paradigm in the 6G era. However, energy-efficient communications and networking among aerial network and terrestrial network face great challenges. This paper is dedicated to discussing green communications of the air-ground integrated heterogeneous network (AGIHN). We first provide a brief introduction to the characteristics of AGIHN in 6G networks. Further, we analyze the challenges of green AGIHN from the aspects of green terrestrial networks and green aerial networks. Finally, several solutions to and key technologies of the green AGIHN are discussed.

Keywords: air-ground integrated heterogeneous network; 6G; green communications

DOI: 10.12142/ZTECOM.202101006

<https://kns.cnki.net/kcms/detail/34.1294.TN.20210309.0907.004.html>, published online March 09, 2021

Manuscript received: 2021-01-29

Citation (IEEE Format): H. C. Wu, H. J. Li, and X. F. Tao, "Green air-ground integrated heterogeneous network in 6G era," *ZTE Communications*, vol. 19, no. 1, pp. 39 - 47, Mar. 2021. doi: 10.12142/ZTECOM.202101006.

1 Introduction

The improvement of network capacity, coverage, delay, security, etc. has always been a key and core task in the development of ground mobile communication networks. Three-dimensional integrated communication is one of the key research directions in achieving the ubiquitous connectivity, ultra-high data rates, and emergency communications in the six generation (6G) networks^[1]. Aerial and space networks facilitate adaptive, flexible, scalable, efficient, and re-

liable three-dimensional wireless coverage for wireless terminals, which have attracted much attention from industry and academia societies. The air-ground integrated heterogeneous network (AGIHN) integrating various aerial communication platforms and terrestrial infrastructures is a cost-efficient paradigm to facilitate extended wireless coverage, ultra-high data rates, post-disaster communication assistance and recovery etc.

A typical AGIHN architecture is shown in **Fig. 1**, where aerial communication platforms such as airships, balloons, and unmanned aerial vehicles (UAVs) acting as the carrier of information collection, transmission and processing can provide broadband wireless communications and supplement terrestrial networks. Terrestrial networks mainly consisting of

This work was supported by National Natural Science Foundation of China under Grant Nos. 61901051 and 61932005.

heterogeneous cellular networks, wireless local area networks (WLAN), and mobile ad hoc networks (MANET) support various applications and services in the areas where infrastructures are easy and low-cost to be deployed.

Nowadays, industry and academia societies have started research and implementation of AGIHN. For example, in 2016, Nokia Bell Labs demonstrated the world’s first flying cell (F-Cell) based on UAV, which was powered by solar energy and could wirelessly transmit high-definition video^[2]. In 2017, EE, the British Telecom Operator, broadcast the mountain bike race live on the mini mobile site “Air Mast” connected to the helium balloon^[3]. As of December 2019, the FirstNet communications platform jointly built by AT&T and First Responder Network Authority had reached more than 1 million connections^[4]. Flying Cells on Wings (COWs) in the platform is an

ideal choice for wildfire and mountain rescue missions. During Hurricane Michael, a COW provided services to first responders on the battered Mexican beach in Florida to support disaster recovery. One Aerostat, which was launched later, can provide more than twice the coverage area compared with COWs, helping responders keep connected in the event of a large-scale catastrophic event.

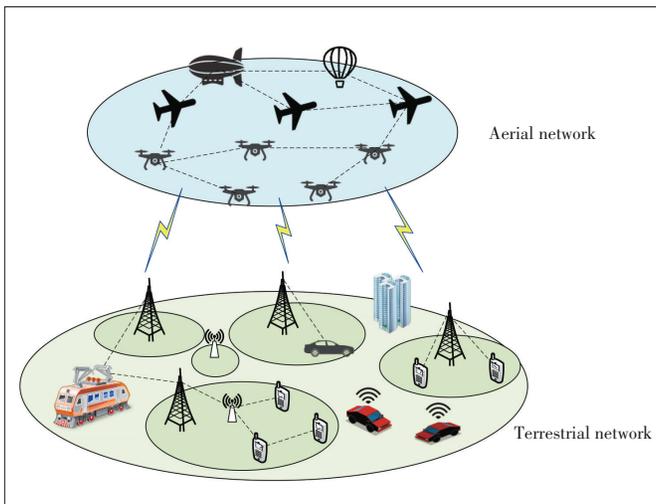
1.1 Characteristics of AGIHN

1.1.1 Heterogeneity

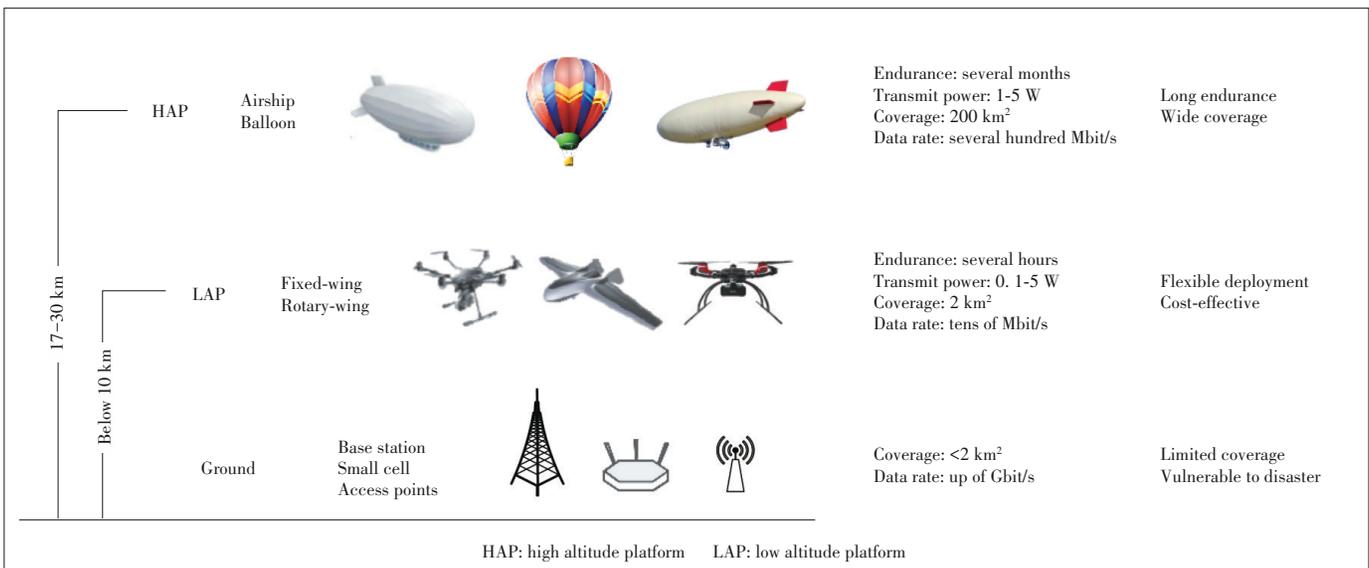
In addition to terrestrial heterogeneous cellular networks, high altitude platforms (HAPs) such as airship and balloon and low altitude platforms (LAPs) such as UAVs are employed in the aerial networks to achieve seamless wireless coverage and to meet differentiated data rate requirements. This heterogeneous integrated network enables diverse systems to cooperate, coordinate, and share information for serving mobile terminals with individualization service anytime and anywhere.

As shown in **Fig. 2**, AGIHN is a large-scale and multi-layer 3D heterogeneous network. HAPs such as airships and balloons are distributed in remote rural areas with imperfect terrestrial infrastructure or disaster areas, with an altitude of 17 – 30 km^[5]. UAVs are used for high-speed services in hot spots or wireless connection in disaster areas, with an altitude below 10 km^[6]. Terrestrial base stations (BSs) and access points (APs) are typically deployed in the area with an altitude below 1 km.

Terrestrial heterogeneous cellular networks realize coverage optimization and capacity improvement by deploying dense small cells with lower transmission power, such as microcell, picocell and femtocell. Due to economic cost and terrain constraints, these terrestrial communication infrastructures are established according to the human habitation and living habits,



▲ Figure 1. Architecture of air-ground integrated heterogeneous network (AGIHN)



▲ Figure 2. Heterogeneity of air-ground integrated heterogeneous network (AGIHN)

which makes wireless traffic a stable and periodical spatial-temporal distribution.

Aerial networks, as a supplement to terrestrial networks, have to deal with the complex and diverse application scenarios with uneven spatial-temporal distributed traffic, diverse service demand, and sudden surge of wireless traffic. Flexible movement is a key feature of AGIHN. HAPs are quasi-static (relative to the ground) platforms. The moving speed of UAVs is 0 – 460 km/h^[7]. Moreover, UAVs can move freely in the 3D space with random trajectory. The ground BSs are typically fixed and deployed in buildings or on high mountains. Ground terminals such as vehicles and terminals on high-speed rails typically have speeds of 0 – 350 km/h and move with relatively fixed trajectories^[8]. Power supply of network nodes in AGIHN is also diverse from each other. There is no continuous power supply source for aerial nodes. Battery, wind, solar, and other combined power supply are the main energy sources for balloons. The endurance of such platforms can reach 150 – 200 days^[9]. UAVs generally use battery power supply and the endurance is only about half an hour to 24 hours^[10]. Ground BSs and APs are driven by grid power system for continuous operation.

For the frequency bands and radio propagations of ground 4G and 5G networks, the frequency resources occupied by 4G system include 1 880 – 1 900 MHz, 2 320 – 2 370 MHz and 2 575 – 2 635 MHz, while the frequency resources occupied by 5G system include 3.3 – 4.2 GHz, 4.4 – 5.0 GHz, the millimeter wave band, 26 GHz, 28 GHz and 39 GHz^[5]. The aerial nodes such as UAVs, balloons and airships work at the Long Term Evolution (LTE) or Wi-Fi communication bands^[11]. They can also work in the unlicensed Industrial, Scientific and Medical (ISM) band defined by the ITU Radio-communication Sector (ITU-R)^[12]. The electromagnetic propagation of different frequency bands also differs from each other. Radio attenuation on high frequency bands is more serious. Compared with the electromagnetic fading at 2 GHz, an additional 22.9 dB of fading exists at 28 GHz^[13]. Good line-of-sight (LoS) transmission links exist in air-to-air, air-to-ground, and ground-to-air channels while the radio propagation in ground transmissions faces more serious signal fading due to rich reflection, refraction, scattering, etc.

1.1.2 High-Dynamically Changed Network Topology

Flexibility is one of the most key features of aerial networks. Payload, height, speed and endurance are the four key factors influencing communication performance of aerial platforms. Payload represents the maximum carrying weight that the platform can hold. HAPs and LAPs carry different communication equipment with different weight. Height refers to the maximum altitude that the aerial platform can be reached, which is closely related to the coverage of the aerial platform. Endurance refers to the maximum flight duration without charging and refueling. As mentioned above, the height, speed and endurance of difference components in AGIHN differ

from each other, which results in the high-dynamic change of the network topology.

Different height and moving speed of diverse platforms make the network topology more stereoscopic. In order to provide flexible services for wireless terminals, communication platforms change their positions and height adaptively. As a result, the network topology changes dramatically with the moving of platforms. The endurance is another key factor having great impact on the network topology. The network topology of ground networks changes slightly since ground BSs and APs are generally fixed located and are powered by grid system while that of aerial networks changes frequently and rapidly due to the energy depletion and battery charging. Besides, AGIHN are more vulnerable to malicious attacks such as wiretapping, hijacking, masking, and jamming, which cause disconnection and interruption of aerial links and re-connection of surviving nodes. The disconnection and re-connection of networks nodes in AGIHN also contribute to the changes of network topology.

1.1.3 Random Perturbation of Aerial Platforms

Due to the lack of fixed infrastructure, the aircrafts are susceptible to airflow and body vibration, leading to random perturbation of aerial communication platforms. According to the tests and measurements, the variation of roll angle (i.e., the elevation angle in this paper) is ± 0.02 rad. The variation of pitch angle (i.e., the azimuth angle in this paper) is ± 0.1 rad^[14]. The random perturbations of aerial platforms may cause error to the estimation of angle of departure (AOD) and angle of arrival (AOA) between transceivers, further leading to the error of channel state information estimation and distortion of coverage area. Consequently, the perturbation of aerial platforms will cause non-robust transmission links, inefficient energy consumption, and serious information leakage, etc.

The perturbation angle of UAV was assumed as high as 10 degrees in Ref. [14]. It shows that the jitter of UAV causes inaccurate estimation of deviation angle between the UAV and ground users and increases the error of AOD estimation. The influence of wind on UAVs was then simulated by using on-board sensors in Ref. [15]. The maximum amplitude of sideslip angle and trajectory angle jitter was approximately 10 degrees, which verifies the previous hypothesis. Considering the impact of UAV jitter, energy-saving secure communications in a downlink A2G wiretap system was investigated in Ref. [16].

1.2 Integrating AGIHN in 6G

With the commercialization of 5G networks, various groups from worldwide countries and regions have initialized plans and programs on potential key technologies for 6G networks. Space-air-ground integrated networking (SAGIN) is acknowledged as a key direction in achieving global connectivity. AGIHN, as an important part of SAGIN, is seen as a cost-effective approach to meeting the requirement of ultra-high data rate

and ubiquitous coverage. To realize the expected goal, the integration of ground networks and aerial networks has to solve the problem of new network architecture design and the challenge of disruptive technology innovation.

1.2.1 Directions of Network Architecture Design

Network architecture design is the first step to realize the integration of ground networks and aerial networks. Efficient coordination of resources and fully exploitation of cooperation between ground networks and aerial networks are the main goals in network architecture design. In order to solve the complex interoperation in the management of heterogeneous networks, the software defined network (SDN) and network function virtualization (NFV) are applied in 5G networks. SDN and NFV are still seen as efficient solution to the network management in 6G networks. In the AGIHN with high-dynamically changed topology, SDN and NFV based core network management architecture can provide distributed and on-demand resource allocation, service guaranteed network slicing, flexible programming of network functions, and security management^[17-20]. SDN and NFV are also seen as promising technologies for providing flexible and reconfigurable green satellite services in space-air-ground integrated networks^[21].

Efficient energy utilization and low energy consumption are always key concepts in network architecture design. Green AGIHN architecture design can be carried out from the aspects of green communications and green computing. For green communications in AGIHN, the aerial platforms can provide cost-effective and energy-saving transmissions for the wireless terminals with appropriate cooperation, trajectory design, user scheduling, power allocation, and combination with improved wireless technologies^[22-25]. Coordination and cooperation architecture of AGIHN and resource management of heterogeneous network nodes are the keys to green communications of AGIHN. For the green computing, aerial communication platforms with mobile edge computing (MEC) can greatly improve the data rate and latency performance in AGIHN^[26]. Moreover, distributed cloud architecture can achieve seamless handover and effective task offloading among UAVs and ground terminals^[27]. Green computing in AGIHN can be realized with energy-efficient MEC and green cloud architecture.

Intelligence is a core idea in the 6G era. To realize ubiquitous intelligent mobile society in the 6G era, artificial intelligence (AI) is expected to fully penetrate the network evolution. With AI applied in AGIHN architecture design, efficient resource management and network optimization can be realized by exploiting the potential information in wireless big data and with less or even no human intervention. Moreover, enhanced privacy preserving can be achieved by leveraging AI in aerial networks^[28].

1.2.2 Directions of Key Wireless Technologies

1) Terahertz communications

AGIHN is facing the contradiction between limited spectrum resources and the rapid growth of high-speed traffic demand. Terahertz communications is an important direction to break through the resource limitations in 6G networks^[29]. The terahertz frequency resources is from 0.1 THz to 10 THz. Terahertz communications has the advantages of ultra-low delay, excellent directivity, anti-interference, wide bandwidth, and strong penetration. Moreover, since the terahertz wavelength is greatly reduced, the antenna size can be greatly reduced, which is beneficial to antenna integration. Leveraging terahertz technology to aerial networks can further improve the data rate with highly concentrated beams, strong LoS path and wide bandwidth resources. However, the short terahertz wave and the weak diffraction introduce quite high path loss of radio propagation. Thus, denser BSs and APs are required to achieve seamless coverage, which means more energy consumption will be introduced.

2) Intelligent reflecting surface

Reconfigurable intelligent reflecting surface (IRS) is attracting attention for wireless networks since it can significantly improve the wireless channel quality by adaptively reconfiguring wireless propagations with massive low-cost passive reflecting elements integrated into a planar surface^[30]. Combining IRS with non-orthogonal multiple access (NOMA) and MEC can greatly improve network throughput and reduce latency^[31-32]. Leveraging IRS in AGIHN, the received signals at the UAV from cellular BSs can be greatly improved by configuring IRS deployed on building walls. The secrecy rate can be enhanced by jointly optimizing phase shifters of IRS, UAV trajectory, and UAV power^[33]. Moreover, leveraging UAVs with IRS can provide energy-efficient communications, which provides new sights in green communications in 6G networks^[34].

3) Spectrum sharing

Spectrum resources are the treasure for wireless communications. In addition to terahertz communications and visible light communications, spectrum sharing is another approach in extending spectrum resources and improving spectrum efficiency in 6G networks with flexible and intelligent frequency allocation and reuse^[35]. Employing blockchain in spectrum sharing can further prevent jamming from malicious users^[36]. Spectrum reuse between dense ground BSs and flexible-mobility UAVs makes interference management more challenging. With appropriate spectrum sharing between aerial networks and ground networks, the area spectrum efficiency and network throughput can be significantly improved.

4) Energy harvesting technologies

With the proliferation of mobile devices and the denser deployment of network BSs and APs, prolonged battery life and improved energy harvesting efficiency are the keys to realize green AGIHN. Simultaneous wireless information and power transfer (SWIPT) is one of the popular energy harvesting technologies studied in recent years^[37]. SWIPT can charge wireless devices while supporting communications and is a promising

energy charging technology for sensors nodes. Leveraging SWIPT in AGIHN, the aerial platforms can effectively charge ground terminals and improve energy efficiency by exploiting the benefit of flexible mobility. In addition, SWIPT combined with IRS in aircrafts can further improve the energy harvesting efficiency^[38].

2 Challenges of Green AGIHN

According to the forecasts of International Energy Agency (IEA) and Global Electronic Sustainability Initiative (GESI), the information and communication technology (ICT) industry currently consumes 2% – 4% of the global energy, which is equivalent to the amount of energy consumption of the aviation industry^[39]. The huge energy consumption not only increases operating costs, but also brings series of resource and environment problems. It is requested by the ITU that the global ICT industry reduce greenhouse gas (GHG) emissions by 45% from 2020 to 2030^[40]. How to improve energy efficiency of the communication industry is an urgent problem to be solved in the 6G era.

In AGIHN, aircrafts acting as information carriers can reduce more energy consumption and provide better energy-efficient services compared with terrestrial networks due to the reduced energy consumption of ancillary facilities such as the air conditioner at BSs. Nevertheless, the green communication in AGIHN still faces many challenges.

2.1 Challenges of Green Terrestrial Networking in AGIHN

2.1.1 Optimized Deployment of BS and AP

In AGIHN of the 6G era, more ground BSs and APs will be deployed to accomplish the requirement of ultra-high data rate services. More infrastructures will be established to support the deployment and operation of BSs and APs. Accordingly, more energy will be consumed. According to energy efficiency requirements for telecommunications proposed by Verizon^[41], BSs consume nearly 80% of the energy consumed in cellular networks while the power amplifiers and air conditioners consume almost 70% of the total energy at BSs^[42]. It is reported by Huawei that the maximum energy consumption of a 5G BS is about 11.5 kW, which is 10 times of that of a 4G BS^[43]. However, according to Daiwa's prediction, the number of 5G BSs will be four times that of 4G BSs in their respective eras^[44]. With the employment of advanced wireless technologies, the deployment of BSs will be much denser in 6G networks. It is foreseen that there will be up to 40 000 sub-networks per km² in 6G networks^[45]. Optimizing the deployment of BSs and APs is one of the key approaches in reducing the energy consumption in AGIHN.

In addition, the ground wireless traffic is non-uniformly distributed in both time and space. It is predicted that the data traffic in the downtown of Milan is 4 times of that in Bocconi University located in suburb^[46]. An analysis report of data traf-

fic in Shanghai reveals that the data traffic in residential areas is 1.5 times of that in office areas^[47]. Moreover, the ratio of day-time traffic amount to night-time traffic is around 0.8 in residential areas while it is up to 1.4 in office areas. Although aerial networks can provide flexible services for ground terminals in such areas with non-uniform traffic, fixed terrestrial communication infrastructure is still the main approach for providing robust and cost-effective services. To satisfy the requirement of temporal and spatial non-uniform traffic and to simultaneously save energy, flexible BS sleeping and awake schemes play important roles in saving energy at BSs^[42].

2.1.2 Increased Mobile Devices and Wireless Access

According to Cisco, there were 8.8 billion mobile devices and wireless connections in 2018, including 4.9 billion smart phones and 1.1 billion IoT devices^[48]. 42% of the devices enjoyed wireless services through 4G cellular networks. Due to the continuous prosperity of sensors, intelligent furniture, Internet of vehicles, smart city and medical applications, and the continuous penetration of vertical industry with 5G networks, there will be 28.5 billion wireless devices by 2022^[49], among which 51% (14.6 billion) of the devices are battery-powered machine-to-machine (M2M) devices. It is estimated that by 2025, the global network standby energy consumption of IoT edge devices will approach 46 TWh^[50], which is about equivalent to the annual electricity consumption of Portugal in 2019^[51].

The massive wireless connectivity brings a great challenge to the wireless access networks due to massive connection requests, ultra-heavy traffic load, limited battery of wireless devices, and diverse levels of subscribers. Moreover, burst access attempts may happen due to some unexpected events such as power failure, which will lead to a sharp increase of control signaling, network congestion, and further increased energy consumption. In addition, handover of massive devices among heterogeneous networks introduces complex interoperation and resource managements, which also causes the increase of energy consumption. Low energy consumption and improved energy efficiency are crucially important for wireless communications in future networks.

2.2 Challenges of Green Aerial Networking

Although aerial platform-based communications can save more energy than ground communications, the explosive growth of aircrafts in wireless communications and the high-dynamically changed topology bring new challenges to the green communications. The UAV is the most widely applied aircraft for wireless communications due to its flexible mobility, cost-efficient and rapid deployment, and LoS communication support. The Federal Aviation Administration (FAA) pointed out that, by 2024, the world will see the emergence of 1.48 million recreational Unmanned Aircraft System (UAS) fleets^[52].

A small UAV usually needs 20 to 200 W/kg to fly^[53]. It is usually powered by on-board batteries. Different types of UAVs carry different battery capacities. For example, the Skywalker fixed-wing UAV carries four 8 500 mAh batteries in series, while the AKS Raven X8 multi-rotor UAV carries two 10 000 mAh batteries^[54]. The battery power consumption of 30 kg to 35 kg UAVs to complete a flight mission (i.e., lifting up, hover, flight, and landing) is about 12.53% and 13.82% of the full amount^[55]. The power consumption of Wi-Fi and GPS communications of a small UAV with the weight of 865 g is about 8.3 W, and that of horizontal flight is 310 W^[56]. Therefore, in the case of limited battery, reducing the power consumption of UAV can provide longer endurance and communication services.

2.2.1 Green Communication Module

Limited battery limits the performance of aerial transmissions. Improving the energy efficiency of communication modules of aircrafts and increasing the battery capacity are the two main approaches for green communications of aerial networking. Aircraft placement, trajectory optimization, power control, flight duration optimization, resource and interference management, and handover in high-dynamic networks are the main factors influencing the energy consumption and energy saving in AGIHN.

Radio propagation and load balance among BSs or APs are closely related with the placement and trajectory design of aerial communication platforms. For HAPs, the placement optimization directly influences the load balance between HAPs and ground BSs. Appropriate placement of HAPs can extend wireless coverage, improve network throughput, and further reduce energy consumption of the communication modules. For LAPs, optimized trajectory design can improve network throughput, reduce energy consumption, and extend flight duration. The trajectory design of multi-UAV networking can further provide seamless connectivity for wireless terminals and extend the coverage area of aerial networks. It should be noted that collision avoidance among multiple UAVs should be considered in the trajectory design. Frequency reuse, spectrum sharing, and power control are key factors in resource and interference management between aerial networks and terrestrial cellular networks. Flexible and adaptive frequency reuse and power control can decrease interference among aerial nodes and ground nodes, which can further improve network throughput and improve energy efficiency.

2.2.2 Green Flight Module

A UAV consumes more propulsion power in flight than in hover^[56]. LAPs have the advantage of mobility compared with HAPs while they consume more power due to their high-dynamic mobility. The flight power consumption of a fixed-wing UAV can be expressed as functions of its velocity V and acceleration $a(t)$ ^[57]:

$$P(V, a(t)) = c_1 V^3 + \frac{c_2}{V} + \frac{c_2}{Vg^2} a^2(t), \quad (1)$$

where $c_1 = \rho C_{D_0} S/2$ and $c_2 = 2W^2/(\pi e_0 A_R \rho S)$ are two parameters (where ρ and C_{D_0} are the air density and zero-lift drag coefficient, respectively, S and W are the wing area and aircraft weight, respectively, and e_0 and A_R denote the wing span efficiency and aspect ratio of the wing, respectively); g is the gravitational acceleration. According to Eq. (1), we can see that the energy consumption increases with the increase of velocity and the absolute value of acceleration, which means that more energy will be consumed with higher flight speed and faster change of the speed. Therefore, in order to achieve green flight of UAV and further to achieve green aerial communications, the UAV should move as smoothly as possible while improving the transmission performance. The flight power consumption of a rotary-wing UAV is a function of the UAV velocity, which is given by Ref. [58]:

$$P(V) = P_0 \left(1 + \frac{3V^2}{U_{\text{tip}}^2} \right) + P_i \left(\sqrt{1 + \frac{V^4}{4v_0^4}} - \frac{V^2}{2v_0^2} \right)^{1/2} + \frac{1}{2} d_0 \rho s A V^3, \quad (2)$$

where $P_0 = \delta \rho s A \Omega^3 R/8$ and $P_i = (1+k)W^{3/2}/\sqrt{2\rho A}$ are two parameters representing the blade profile power and induced power of the UAV in hover, respectively (where δ is the profile drag coefficient, s and A are the rotor stiffness and rotor disc area, respectively, Ω and R denote the blade angular velocity and rotor radius, respectively, and k denotes the incremental correction factor to induced power); U_{tip} is the rotor tip velocity; v_0 is the average rotor induced velocity in hover; d_0 denotes the fuselage drag ratio. According to Eq.(2), the first and the third terms are the power required to overcome the profile drag of the blades and the fuselage drag, respectively, which increases with the square and cubic of the velocity. The second term is the power required to overcome the induced drag of the blades, which decreases with the velocity. It is verified in Ref. [58] that the total power consumption first decreases and then increases with the increase of UAV velocity, which means that the energy consumption can be minimized with the optimal UAV velocity.

As a result, the improvement of transmission quality and energy efficiency in AGIHN should take flight consumption into consideration. There are three main research directions for energy-saving UAV communications: the optimization of flight radius and velocity with a fixed trajectory^[59]; the joint optimization of UAV trajectory, acceleration, and velocity^[60]; the trade-off between performance improvement and energy consumption^[61].

3 Key Technologies in Green AGIHN

How to realize energy-efficient network collaboration and integration of heterogeneous networks in AGIHN is one of the keys to SAGIN. In this section, several promising technologies for harvesting energy and reducing energy consumption in AGIHN are analyzed.

3.1 Energy Harvesting Technology

Energy harvesting technologies are promising approaches to prolonging battery life and providing extra power in green communications by collecting external energy resources such as light, heat, electromagnetic, and mechanical. Wireless power transfer (WPT) and SWIPT are two energy harvesting technologies for transporting energy with electromagnetic energy. WPT is a basic energy harvesting technology while SWIPT is an energy harvesting technology that integrates WPT and the communication function, i.e., SWIPT can simultaneously transmit information signals while transporting electromagnetic energy.

Leveraging WPT and SWIPT in AGIHN can realize remote charging and mutual charging among aircrafts. However, the mutual energy transport definitely causes increased energy consumption. Traditional energy collection technologies such as solar and wind systems can be combined as a RF energy source to reduce consumption of non-green energy, which is a promising research direction for AGIHN. For example, charging LAPs with solar-powered satellites and HAPs can prolong the flight duration of LAPs and enhance the stability of aerial networking. Moreover, the LAPs can further charge each other by exploiting the benefits of LoS links and charge ground terminals with two-hop wireless energy transmission. In addition, leveraging IRS in SWIPT-assisted LAP system can simultaneously improve the network performance and enhance the energy transmission efficiency^[38].

3.2 Cooperative Communications

Cooperative transmission can improve network performance by coordinating multiple diverse nodes to exploit the multiplexing gain and diversity gain. In AGIHN, coordinating multiple aerial nodes as information signal sources can provide significantly improved capacity and coverage performance with LoS links, flexible-changed topology, and adaptive resource coordination. The transmission power of aerial nodes can be greatly reduced due to the decreased path loss fading. Further, coordinating multiple aerial nodes as electromagnetic energy sources can realize rapid power charging and network restoration and reconstruction. Coordinating aerial nodes and ground BSs can provide more energy-efficient transmissions than coordinating only ground BSs due to the reduced transmission power of aerial nodes and the exemption of energy consumption of BS infrastructures. Especially, the energy efficiency of cell edge users can be significantly improved by deploying periodic UAVs at the edge of ground BSs and periodically coordinating UAVs-enabled BSs or relays with the ground BSs^[62].

Coordination and cooperation of ground nodes and aerial nodes can realize improved energy efficiency in AGIHN. However, the high-dynamic network topology, non-uniform data traffic, and random perturbation of aerial platforms bring great challenges to the cooperative communications in AGIHN. The high-dynamic network topology introduces frequent disconnection and reconnections, which brings huge signaling overhead. Blockchain based access registration provides a way in reducing the signaling overhead^[63]. Non-uniform data traffic brings challenges to the user scheduling and traffic load coordination among cooperative nodes, which further causes tradeoff between balanced traffic offloading and efficient energy consumption. The random perturbation of aerial platforms causes non-robust and inefficient transmission links between a pair of transceivers. Accurate estimation of platform perturbation, appropriate compensation for perturbation, and adaptive power allocation and beam adjusting are required for energy-efficient communications in AGIHN.

3.3 Integrating Intelligence into Green AGIHN

AI technologies, especially machine learning (ML) and big data analysis, are promising and widely acknowledged solutions to smart network control and network management. AI has already been applied in mobile networks from physical layer design to network layer control. Leveraging AI in 6G networks is an inevitable trend. In green AGIHN, AI can be applied in many aspects including intelligent architecture design, real-time network data analysis, flexible aerial platform control, secure aerial platform tracing, smart platform position, high-efficient energy harvesting, intelligent routing, intelligent caching and computing, intelligent sleeping and wake-up mechanisms, adaptive and efficient resource allocation and scheduling, etc.

Leveraging intelligence into AGIHN can realize enhanced energy efficiency by globally optimizing the network control and management. However, mechanisms adopted to realize intellectualization will bring additional energy consumption, especially for the intelligent network management through global control or massive data analysis. Therefore, energy-efficient intelligent network control and management mechanisms are required in green AGIHN.

4 Conclusions

Facing the high complexity and diversity of future networks and the proliferation of wireless devices, AGIHN is considered as a cost-effective approach to ubiquitous coverage and ultra-high throughput. With the increasing scarcity of natural resources and complexity of the radio environment, it is urgent to solve the problems of green AGIHN to reduce energy consumption and to improve energy efficiency. In this paper, we first introduced the integration of AGIHN and 6G networks. Then, the challenges of green AGIHN were analyzed from the

aspects of green terrestrial network and green aerial network. Following the analysis, several promising green technologies that can be employed in AGIHN have been discussed.

References

- [1] ZHANG Z Q, XIAO Y, MA Z, et al. 6G wireless networks: vision, requirements, architecture, and key technologies [J]. IEEE vehicular technology magazine, 2019, 14(3): 28 – 41. DOI: 10.1109/MVT.2019.2921208
- [2] F-Cell technology from Nokia Bell Labs revolutionizes small cell deployment by cutting wires, costs and time [EB/OL]. (2016-10-03)[2020-12-09]. <https://www.nokia.com/about-us/news/releases/2016/10/03/f-cell-technology-from-nokia-bell-labs-revolutionizes-small-cell-deployment-by-cutting-wires-costs-and-time>
- [3] Launches world's first 4G "AIR MAST" to connect red bull foxhunt mountain bike even in rural wales [EB/OL]. (2017-10-11)[2020-12-09]. <https://newsroom.ee.co.uk/ee-launches-worlds-first-4g-air-mast-to-connect-red-bull-foxhunt-mountain-bike-event-in-rural-wales>
- [4] FirstNet reaches over 1 million connections [EB/OL]. (2019-12-03)[2020-12-09]. https://about.att.com/story/2019/fn_hits_one_million.html
- [5] ZHOU D, GAO S, LIU R Q, et al. Overview of development and regulatory aspects of high altitude platform system [J]. Intelligent and converged networks, 2020, 1(1): 58 – 78. DOI: 10.23919/ICN.2020.0004
- [6] AL-HOURANI A, KANDEEPAN S, LARDNER S. Optimal LAP altitude for maximum coverage [J]. IEEE wireless communications letters, 2014, 3(6): 569 – 572. DOI: 10.1109/LWC.2014.2342736
- [7] WIREN R. 5G and UAV use cases [EB/OL]. (2017-09-18)[2020-12-09]. [http://www.5gsummit.org/helsinki/docs/RWiren-5G and UAV use cases-2017-09-18.pdf](http://www.5gsummit.org/helsinki/docs/RWiren-5G%20and%20UAV%20use%20cases-2017-09-18.pdf)
- [8] SHARMA V, BENNIS M, KUMAR R. UAV-assisted heterogeneous networks for capacity enhancement [J]. IEEE communications letters, 2016, 20(6): 1207 – 1210. DOI: 10.1109/LCOMM.2016.2553103
- [9] QIU J F, GRACE D, DING G R, et al. Air-ground heterogeneous networks for 5G and beyond via integrating high and low altitude platforms [J]. IEEE wireless communications, 2019, 26(6): 140 – 148. DOI: 10.1109/MWC.0001.1800575
- [10] HWANG M H, CHA H R, JUNG S Y. Practical endurance estimation for minimizing energy consumption of multirotor unmanned aerial vehicles [J]. Energies, 2018, 11(9): 2221. DOI: 10.3390/en11092221
- [11] 3GPP. Technical Specification Group Radio Access Network; Study on Enhanced LTE Support for Aerial Vehicles: TR36.777 [R]. Sophia - Antipolis, France: 3GPP, 2017.
- [12] CHRIKI A, TOUATI H, SNOUSSI H, et al. Centralized cognitive radio based frequency allocation for UAVs communication [C]//15th International Wireless Communications & Mobile Computing Conference (IWCMC). Tangier, Morocco: IEEE, 2019: 1674 – 1679. DOI: 10.1109/IWCMC.2019.8766481
- [13] WANG G Y, LIU Y J, QI X Y. Study on the propagation characteristics of 28GHz radio wave in outdoor microcellular [C]//Asia-Pacific Microwave Conference (APMC). Nanjing, China: IEEE, 2015: 1 – 3. DOI: 10.1109/APMC.2015.7413403
- [14] AHMED B, POTA H R, GARRATT M. Flight control of a rotary wing UAV using backstepping [J]. International journal of robust and nonlinear control, 2010, 20(6): 639 – 658. DOI: 10.1002/rnc.1458
- [15] CHOI H S, LEE S, RYU H, et al. Dynamics and simulation of the effects of wind on UAVs and airborne wind measurement [J]. Transactions of the Japan society for aeronautical and space sciences, 2015, 58(4): 187 – 192. DOI: 10.2322/tjsass.58.187
- [16] WU H C, WEN Y, ZHANG J Z, et al. Energy-efficient and secure air-to-ground communication with jittering UAV [J]. IEEE transactions on vehicular technology, 2020, 69(4): 3954 – 3967. DOI: 10.1109/TVT.2020.2971520
- [17] MUNOZ R, VILALTA R, CASELLAS R, et al. Integrated SDN/NFV management and orchestration architecture for dynamic deployment of virtual SDN control instances for virtual tenant networks [J]. IEEE/OSA journal of optical communications and networking, 2015, 7(11): B62 – B70. DOI: 10.1364/JOCN.7.000B62
- [18] XILOURIS G K, BATISTATOS M C, ATHANASIADOU G E, et al. UAV-assisted 5G network architecture with slicing and virtualization [C]//IEEE Globecom Workshops (GC Wkshps). Abu Dhabi, United Arab Emirates: IEEE, 2018: 1 – 7. DOI: 10.1109/GLOCOMW.2018.8644408
- [19] WHITE K J S, DENNEY E, KNUDSON M D, et al. A programmable SDN+NFV-based architecture for UAV telemetry monitoring [C]//14th IEEE Annual Consumer Communications & Networking Conference (CCNC). Las Vegas, USA: IEEE, 2017: 522 – 527. DOI: 10.1109/CCNC.2017.7983162
- [20] HERMOSILLA A, ZARCA A M, BERNABE J B, et al. Security orchestration and enforcement in NFV/SDN-aware UAV deployments [J]. IEEE access, 2020, 8: 131779 – 131795. DOI: 10.1109/ACCESS.2020.3010209
- [21] HERMOSILLA A, ZARCA A M, BERNABE J B, et al. Security orchestration and enforcement in NFV/SDN-aware UAV deployments [J]. IEEE access, 2020, 8: 131779 – 131795. DOI: 10.1109/ACCESS.2020.3010209
- [22] NIE Y W, ZHAO J H, LIU J, et al. Energy-efficient UAV trajectory design for backscatter communication: a deep reinforcement learning approach [J]. China communications, 2020, 17(10): 129 – 141. DOI: 10.23919/JCC.2020.10.009
- [23] LIU T, CUI M, ZHANG G C, et al. 3D trajectory and transmit power optimization for UAV-enabled multi-link relaying systems [J]. IEEE transactions on green communications and networking, 8135, PP(99): 1. DOI: 10.1109/TGCN.2020.3048135
- [24] YIN S X, ZHAO Y F, LI L H, et al. UAV-assisted cooperative communications with power-splitting information and power transfer [J]. IEEE transactions on green communications and networking, 2019, 3(4): 1044 – 1057. DOI: 10.1109/TGCN.2019.2926131
- [25] WANG Y S, HONG Y W P, CHEN W T. Trajectory learning, clustering, and user association for dynamically connectable UAV base stations [J]. IEEE transactions on green communications and networking, 2020, 4(4): 1091 – 1105. DOI: 10.1109/TGCN.2020.3005290
- [26] LV Z, HAO J J, GUO Y J. Energy minimization for MEC-enabled cellular-connected UAV: trajectory optimization and resource scheduling [C]//IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS). Toronto, Canada: IEEE, 2020: 478 – 483. DOI: 10.1109/INFOCOM-WKSHPS50562.2020.9162853
- [27] SUN X, ANSARI N. Green cloudlet network: a distributed green mobile cloud network [J]. IEEE network, 2017, 31(1): 64 – 70. DOI: 10.1109/MNET.2017.1500293NM
- [28] WANG Y T, SU Z, ZHANG N, et al. Learning in the air: secure federated learning for UAV-assisted crowdsensing [J]. IEEE transactions on network science and engineering, 4385, PP(99): 1. DOI: 10.1109/TNSE.2020.3014385
- [29] CHEN Z, MA X Y, ZHANG B, et al. A survey on terahertz communications [J]. China communications, 2019, 16(2): 1 – 35. DOI: 10.12676/jcc.2019.02.001
- [30] WU Q Q, ZHANG R. Towards smart and reconfigurable environment: intelligent reflecting surface aided wireless network [J]. IEEE communications magazine, 2020, 58(1): 106 – 112. DOI: 10.1109/MCOM.001.1900107
- [31] SONG D, SHIN W, LEE J. A maximum throughput design for wireless powered communications networks with IRS-NOMA [J]. IEEE wireless communications letters, 6722, PP(99): 1. DOI: 10.1109/LWC.2020.3046722
- [32] ZHOU F S, YOU C S, ZHANG R. Delay-optimal scheduling for IRS-aided mobile edge computing [J]. IEEE wireless communications letters, 2189, PP(99): 1. DOI: 10.1109/LWC.2020.3042189
- [33] FANG S S, CHEN G J, LI Y H. Joint optimization for secure intelligent reflecting surface assisted UAV networks [J]. IEEE wireless communications letters, 2021, 10(2): 276 – 280. DOI: 10.1109/LWC.2020.3027969
- [34] MOHAMED Z, AISSA S. Leveraging UAVs with intelligent reflecting surfaces for energy-efficient communications with cell-edge users [C]//IEEE International Conference on Communications Workshops (ICC Workshops). Dublin, Ireland: IEEE, 2020: 1 – 6. DOI: 10.1109/iccworkshops49005.2020.9145273
- [35] ZHANG J Z, CHEN Y, LIU Y X, et al. Spectrum knowledge and real-time observing enabled smart spectrum management [J]. IEEE access, 2020, 8: 44153 – 44162. DOI: 10.1109/ACCESS.2020.2978005
- [36] ZHANG Y Y, FANG Z J. Dynamic double threshold spectrum sensing algorithm based on block chain [C]//3rd International Conference on Electronic Information Technology and Computer Engineering (EITCE). Xiamen, China: IEEE, 2019: 1090 – 1095. DOI: 10.1109/EITCE47263.2019.9094864
- [37] SUN X L, YANG W W, CAI Y M, et al. Physical layer security in millimeter

- wave SWIPT UAV-based relay networks [J]. *IEEE access*, 2019, 7: 35851 – 35862. DOI: 10.1109/ACCESS.2019.2904856
- [38] LIU J X, XIONG K, LU Y, et al. Energy efficiency in secure IRS-aided SWIPT [J]. *IEEE wireless communications letters*, 2020, 9(11): 1884 – 1888. DOI: 10.1109/LWC.2020.3006837
- [39] WEBB M. *Smart 2020: enabling the low carbon economy in the information age* [R]. London, UK: The Climate Group, 2008.
- [40] ITU-T. Greenhouse gas emissions trajectories for the information and communication technology sector compatible with the UNFCCC paris agreement: L.1470 [R]. 2020
- [41] Verizon NEBSTM Compliance: Energy Efficiency Requirements for Telecommunications [EB/OL]. (2009-08-07)[2020-12-09]. <http://www.tiab-online.co.uk/pdf/L5428.pdf>
- [42] GANDOTRA P, JHA R K, JAIN S. Green communication in next generation cellular networks: a survey [J]. *IEEE access*, 2017, 5: 11727 – 11758. DOI: 10.1109/ACCESS.2017.2711784
- [43] Huawei Technology Co., Ltd. 5G power white paper [R]. Shenzhen, China: Huawei, 2019
- [44] Daiwa Capital Markets. *Asia mobile communication: let's talk about 5G* [R]. Hong Kong, China: Daiwa Capital Markets Hong Kong Limited, 2018
- [45] BERARDINELLI G, MOGENSEN P, ADEOGUN R O. 6G subnetworks for life-critical communication [C]//2nd 6G Wireless Summit (6G SUMMIT). Levi, Finland: IEEE, 2020: 1 – 5. DOI:10.1109/6GSUMMIT49458.2020.9083877
- [46] ZENG Q T, SUN Q, CHEN G, et al. Traffic prediction of wireless cellular networks based on deep transfer learning and cross-domain data [J]. *IEEE access*, 2020, 8: 172387 – 172397. DOI: 10.1109/ACCESS.2020.3025210
- [47] ZHANG M Y, FU H H, LI Y, et al. Understanding urban dynamics from massive mobile traffic data [J]. *IEEE transactions on big data*, 2019, 5(2): 266 – 278. DOI: 10.1109/TBDATA.2017.2778721
- [48] Cisco. *Cisco annual internet report* [R]. San Jose, USA: Cisco Systems, Inc., 2020
- [49] Cisco. *Cisco visual networking index (VNI) complete forecast update* [R]. San Jose, USA: Cisco Systems, Inc., 2018
- [50] FRIEDLI M, KAUFMANN L, PAGANINI F, et al. Energy efficiency of the internet of things [R]. Luzern, Switzerland: Lucerne University of Applied Sciences, 2016
- [51] Enerdata. *Portugal energy report* [R]. Grenoble, France: Enerdata, 2020
- [52] FAA. *FAA aerospace forecast fiscal years 2020-2040* [R]. Washington DC, USA: Federal Aviation Administration, 2020
- [53] BERTRAN E, SÀNCHEZ-CERDÀ A. On the tradeoff between electrical power consumption and flight performance in fixed-wing UAV autopilots [J]. *IEEE transactions on vehicular technology*, 2016, 65(11): 8832 – 8840. DOI:10.1109/TVT.2016.2601927
- [54] BOON M A, DRIJFHOUT A P, TESFAMICHAEL S. Comparison of a fixed-wing and multi-rotor UAV for environmental mapping applications: A case study [J]. *ISPRS-international archives of the photogrammetry, remote sensing and spatial information sciences*, 2017, XLII-2/W6: 47 – 54. DOI: 10.5194/isprs-archives-xlii-2-w6-47-2017
- [55] CHAN C W, KAM T Y. A procedure for power consumption estimation of multi-rotor unmanned aerial vehicle [C]//*Journal of Physics: Conference Series*, Volume 1509, 10th Asian-Pacific Conference on Aerospace Technology and Science & 4th Asian Joint Symposium on Aerospace Engineering (APCATS'2019 / AJSAAE'2019). Bristol, UK: IOP Publishing, 2020, 1509(1): 012015
- [56] ABEYWICKRAMA H V, JAYAWICKRAMA B A, HE Y, et al. Comprehensive energy consumption model for unmanned aerial vehicles, based on empirical studies of battery performance [J]. *IEEE access*, 2018, 6: 58383 – 58394. DOI: 10.1109/ACCESS.2018.2875040
- [57] ZENG Y, ZHANG R. Energy-efficient UAV communication with trajectory optimization [J]. *IEEE transactions on wireless communications*, 2017, 16(6): 3747 – 3760. DOI: 10.1109/TWC.2017.2688328
- [58] ZENG Y, XU J, ZHANG R. Energy minimization for wireless communication with rotary-wing UAV [J]. *IEEE transactions on wireless communications*, 2019, 18(4): 2329 – 2345. DOI: 10.1109/TWC.2019.2902559
- [59] HU Y L, YUAN X P, XU J, et al. Optimal 1D trajectory design for UAV-enabled multiuser wireless power transfer [J]. *IEEE transactions on communications*, 2019, 67(8): 5674 – 5688. DOI: 10.1109/TCOMM.2019.2911294
- [60] AHMED S, CHOWDHURY M Z, JANG Y M. Energy-efficient UAV relaying communications to serve ground nodes [J]. *IEEE communications letters*, 2020, 24(4): 849 – 852. DOI: 10.1109/LCOMM.2020.2965120
- [61] YANG D C, WU Q Q, ZENG Y, et al. Energy tradeoff in ground-to-UAV communication via trajectory design [J]. *IEEE transactions on vehicular technology*, 2018, 67(7): 6721 – 6726. DOI: 10.1109/TVT.2018.2816244
- [62] HUA M, WANG Y, LI C G, et al. Energy-efficient optimization for UAV-aided cellular offloading [J]. *IEEE wireless communications letters*, 2019, 8(3): 769 – 772. DOI: 10.1109/LWC.2019.2891727
- [63] LI Z Y, HAO J L, LIU J, et al. An IoT-applicable access control model under double-layer blockchain [J]. *IEEE transactions on circuits and systems II: Express briefs*, 5031, PP(99): 1. DOI: 10.1109/TCSII.2020.3045031

Biographies

WU Huici (dailywu@bupt.edu.cn) received the Ph.D. degree from Beijing University of Posts and Telecommunications (BUPT), China in 2018. From 2016 to 2017, she visited the Broadband Communications Research (BBRC) Group, University of Waterloo, Canada. She is now an associate professor at BUPT. She served as the publication co-chair of APCC 2018 and TPC member of IEEE ICC 2019/2020 and IEEE/CIC ICC 2019/2020. Her research interests are in the area of wireless communications and networks, with current emphasis on collaborative air-to-ground communications and wireless access security.

LI Hanjie received her B.E. degree from North China University of Water Resources and Electric Power, China in 2017. She is currently pursuing M.S. degree in Beijing University of Posts and Telecommunications, China. Her research interests are in the area of physical layer security and UAV power saving communications.

TAO Xiaofeng received the B.S. degree in electrical engineering from Xi'an Jiaotong University, China in 1993 and the M.S.E.E. and Ph.D. degrees in telecommunication engineering from Beijing University of Posts and Telecommunications (BUPT), China in 1999 and 2002, respectively. He was a visiting professor with Stanford University, USA from 2010 to 2011; the chief architect of the Chinese National FuTURE Fourth-Generation (4G) TDD working group from 2003 to 2006; and established the 4G TDD CoMP trial network in 2006. He is currently a professor at BUPT and the Fellow of the Institution of Engineering and Technology (IET). He is the inventor or co-inventor of 50 patents and the author or co-author of 120 papers in 4G and beyond 4G.

Kinetic Energy Harvesting Toward Battery-Free IoT: Fundamentals, Co-Design Necessity and Prospects



LIANG Junrui, LI Xin, YANG Hailiang

(School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China)

Abstract: Energy harvesting (EH) technology is developed with the purpose of harnessing ambient energy in different physical forms. Although the available ambient energy is usually tiny, not comparable to the centralized power generation, it brings out the convenience of on-site power generation by drawing energy from local sources, which meets the emerging power demand of long-lasting, extensively-deployed, and maintenance-free Internet of Things (IoT). Kinetic energy harvesting (KEH) is one of the most promising EH solutions toward the realization of battery-free IoT. The KEH-based battery-free IoT can be extensively deployed in the smart home, smart building, and smart city scenarios, enabling perceptivity, intelligence, and connectivity in many infrastructures. This paper gives a brief introduction to the configurations and basic principles of practical KEH-IoT systems, including their mechanical, electrical, and computing parts. Although there are already a few commercial products in some specific application markets, the understanding and practice in the co-design and optimization of a single KEH-IoT device are far from mature, let alone the conceived multi-agent energy-autonomous intelligent systems. Future research and development of the KEH-IoT system beckons for more exchange and collaboration among mechanical, electrical, and computer engineers toward general design guidelines to cope with these interdisciplinary engineering problems.

Keywords: kinetic energy harvesting; battery-free solution; Internet of Things; co-design

DOI: 10.12142/ZTECOM.202101007

<https://kns.cnki.net/kcms/detail/34.1294.TN.20210309.0902.002.html>, published online March 9, 2021

Manuscript received: 2021-01-22

Citation (IEEE Format): J. R. Liang, X. Li, and H. L. Yang, "Kinetic energy harvesting toward battery-free IoT: fundamentals, co-design necessity and prospects," *ZTE Communications*, vol. 19, no. 1, pp. 48 - 60, Mar. 2021. doi: 10.12142/ZTECOM.202101007.

1 Background

The on-going development trend of the Internet of Things (IoT) technology has attracted a lot of interests from the communication and computer research communities, as well as numerous investments in the related industries. The topics of wireless communication and computational intelligence have caught the most attention. For wireless communication, existing technologies can be categorized according to their spatial scopes. Rooted from the radio-

frequency identification (RFID) technology, the near field communication (NFC) is already very mature for connecting things in a very short distance, that is to say, less than several centimeters. For a longer distance around a human user, ranging from a few centimeters to a few meters, the wireless personal area network (WPAN) is also very mature. Bluetooth and Zigbee are two of the most prevailing WPAN technologies for the IoT. Unlike general Wi-Fi communication, the low-power feature is usually emphasized in the WPAN for the IoT. To

connect things in an even longer distance, we need the low power wide area network (LPWAN). Narrow-band Internet of Things (NB-IoT) and Long Range (LoRa) technologies are the two most representative technologies of LPWAN. In particular, NB-IoT is compatible with the existing cellular networks and operates in licensed frequency bands; it might be constructed as essential infrastructure for communication in the visible future.

Wireless communication is the main battlefield of IoT technology. As mentioned above, different wireless communication technologies for different targeted ranges have already been proposed and received extensive investigations. On the other hand, we should also keep in mind that all things with sensing, computing, or communication capabilities need electric power to run. Compared with wireless communication, the wireless power transfer or power generation technologies in different distance ranges are far from mature. The near-field case is the most mature, given the rapid development of wireless power transfer (WPT) technology in the last two decades^[1]. For the WPAN range, electric power can hardly be transferred efficiently to passive devices, which do not carry long-term energy storage, on edge. The WPAN-level backscatter communication has recently attracted a lot of research interest^[2-3]. These backscatter devices only scavenge tiny energy from the transmitter to alter their radio frequency (RF) characteristics and deliver low throughput data, rather than actively send back signals^[4]. For the things in the wider WPAN or LPWAN scopes, wireless power transfer is invalid, as the power attenuation of the RF signal is proportional to the cube of the distance between the transmitter and receiver. To keep things connected in such a wide range, they must bring their own energy storage or be able to harvest energy by themselves in their surroundings^[5-6]. In the long run, the energy harvesting (EH) technology is the only option to provide continuous power supply to many scattered things for supporting their everlasting operation^[7-8]. Once all things become energy-autonomous, many benefits, such as maintenance-free operation and ubiquitous deployment, can be achieved. The EH technology is naturally linked with the battery-free IoT applications; otherwise, the tiny ambient energy means almost nothing compared with the centralized large-scale power generation or even the storable energy in chemical batteries¹. Therefore, at the current stage and from the practical point of view, EH should be first regarded as a backbone technology toward the realization of massively deployed and everlasting low-power IoT, rather than taken as a substitution of general power generation technology. The ambient energy sources are usually characterized as volatile sources, like most renewable sources such as solar and wind power. The EH devices should be designed not only to handle these volatile sources but also to

closely meet the demand of low-power IoT in time. The scarce and volatile energy supply in practice and the reliable and timely information demand in expectation together impose the biggest challenge in the co-design and optimization of EH-IoT systems.

2 Kinetic Energy Harvesting

Besides harvesting energy from the conventional renewable ambient energy sources such as the solar source², the other most popular and extensively investigated EH technologies include the RF EH, thermoelectric EH (TEH), and kinetic EH (KEH). Among those energy sources, the kinetic one is unique. On the one hand, electric machine technology has maturely been used for many years for large-scale electromechanical power generation. Kinetic energy is a must step in most centralized power generation processes. For example, in coal-fired, nuclear, and hydro-power plants, energy is converted through turbines and electromagnetic generators into electricity. On the other hand, the small-scale kinetic energy harvesting technology has caught people's attention for just one or two decades. Different from the centralized generation and power grid facilities, in which the power flows among the power generation, transmission, distribution, and utilization stages are intensively and precisely controlled, a small-scale KEH-IoT is a self-contained and energy-autonomous system. It has to be sustainable against environmental uncertainties (intermittent or random vibrations), and execute the sensing, computing, and communication functions correctly and timely according to the information demand.

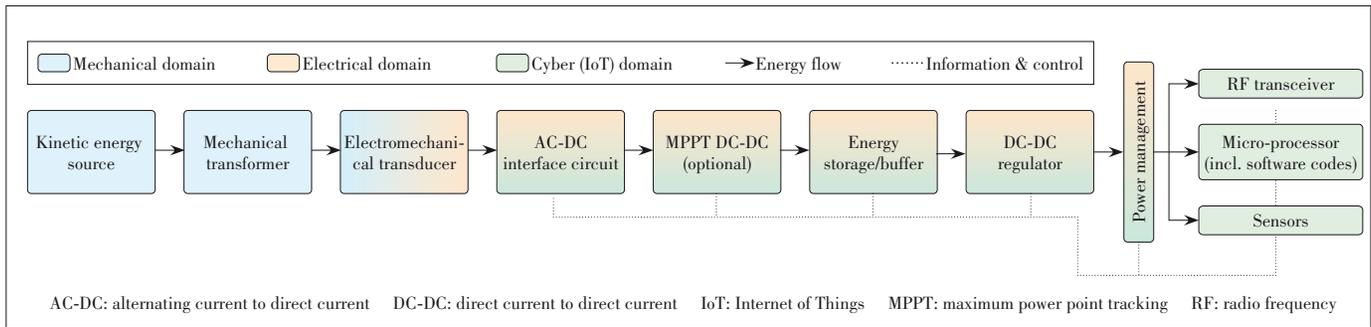
In a word, toward practical application, it is inappropriate to merely emphasize the KEH design as either a mechanical dynamics problem, or an electrical power conversion problem, or a computational problem. Perfect coordination among mechanical dynamics, electrical power conversion and information processing is what we are looking for toward the successful development and application of a KEH-IoT system.

2.1 System Overview

A KEH-IoT system connects three major physical domains: the mechanical, electrical, and cyber domains. **Fig. 1** shows the general block diagram of a KEH-IoT system, in which the three domains are distinguished with different colors. The mechanical kinetic energy excites the vibration or movement of the mechanical structure. A mechanical transformer transfers the excitation to the transducer input with a specific force or displacement ratio to match the kinetic source and transducer mechanical characteristics. An electromechanical transducer converts the mechanical energy into an electrical one, which is in an AC (alternating current) form. Following the transduc-

1. The available power from the ambient energy sources is very tiny, usually from micro-watt to watt scales.

2. Different from the large-scale solar panel, the solar energy harvesting is usually referred to as the on-site power generation using small photovoltaic panels.

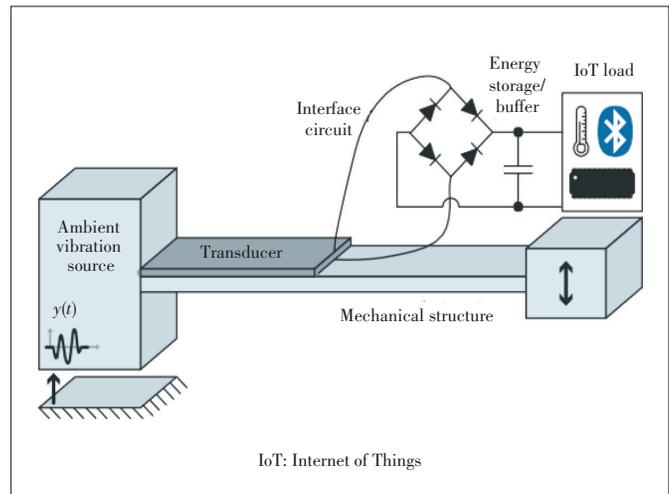


▲ Figure 1. Block diagram of a kinetic energy harvesting (KEH)-IoT system

er electrical output, there are several stages of power electronic modules for adapting the extracted power into a stable logic voltage level. The interface circuit carries out the AC-DC (alternating current to direct current) conversion, which is usually called rectification. At the same time, it can collaborate with a following and optional DC-DC maximum power point tracking (MPPT) stage for extracting more power from the vibrating mechanical structure. The extracted energy is stored in a storage device, such as a chemical battery, a super-capacitor, or sometimes just a small buffer capacitor for immediate utilization. Since the storage level might be floating, which is likely to give a fluctuating output voltage, a DC-DC regulator is necessary for providing a stable and reliable voltage level for powering digital electronics.

In the previous studies, people mostly studied a KEH system by considering the theoretical model and practical design of its mechanical structure and/or electrical circuit^[9-11]. However, recent research on energy harvesting based IoT systems has shown that, when doing computation with scarce and volatile energy supply, there might be more stories than most people, who are working on mechanical or/and electrical designs, have thought^[12]. Therefore, the effect of an IoT load should not be simply taken as an equivalent constant resistive load. In a KEH-IoT system, the dynamics produced by its computing demand is as important as the electromechanically coupled dynamics. In Fig. 1, the cyber part is usually composed of a microprocessor, sensors, the RF transceiver, the digital interface and control of the power converters, and more importantly and necessarily, the software codes. The power management hardware design and the embedded software act crucial roles in moderating the KEH harvested power (as an income) and IoT information demand (as an outcome). The energy and information flows within a KEH-IoT system are also illustrated in Fig. 1.

To give a more intuitive idea of a KEH-IoT system, Fig. 2 shows a typical piezoelectric KEH system. The base vibration corresponds to the kinetic energy source. The cantilever beam and proof mass together act as a mechanical transformer. The electrical part includes a bridge rectifier for AC-DC rectification and a filter capacitor as the energy storage/buffer. The IoT load is not specified here. It might carry out some sensing, computing, or communication functions in general. Given the



▲ Figure 2. A typical piezoelectric KEH system

scarce energy that can be harvested with a piezoelectric energy harvester, the hardware and software of the IoT load must be carefully designed by referring to the energy-aware computing technology.

2.2 Mechanical Design

The modeling and design of KEH systems have attracted a lot of research interest from mechanical engineers during the last two decades^[10, 13-15]. It is intuitively understandable that removing energy from a vibration system must result in an increase of its damping coefficient. The damping effect was quantitatively investigated in the previous studies^[16-18]. In some designs, harvesting energy might also change the system dynamics from the electrical side as well^[19-21]. A large branch of studies focuses on the theoretical modeling of mechanical dynamics when the system is subjected to different loading conditions, i.e., to understand the behavior of the KEH system in terms of vibration dynamics and harvested power^[22-23]. Another biggest category of studies investigated how to modify or redesign the structural configurations toward some objectives, such as

- enhancing the harvesting capability by matching the stroke or force between the mechanical energy source and the transducer mechanical-side input^[24];

- broadening the harvesting bandwidth^[25];
- adapting to vibrations from different spatial directions^[26];
- adapting to different types of vibration patterns, such as harmonic, intermittent, or random vibrations^[27];
- better exploring different kinetic sources, such as human motion^[24], fluid^[28], or gust^[29];
- focusing elastic wave energy for better feeding the energy harvester^[30-31].

As illustrated in Fig. 1, two of the most important mechanical aspects of a KEH system is the kinetic energy source and mechanical transformer. Different from their solar and RF counterparts, and even the thermoelectric energy sources, kinetic sources might have a lot of different possible patterns and affluent dynamic characteristics. Environmental vibrations are usually characterized as random excitation^[32-33], which has a very broad power spectrum. Some machine vibrations are harmonic^[22-23] or periodic^[27,34], i.e., only component vibrations at some narrow frequency bands. Besides, some movements, in particular the fluid movements, are with constant current^[28,35]; while some others are shock- or impact-based^[36-37]. Therefore, there seems no such a universal structure that can optimally harness energy from all kinds of vibration conditions. Customization is a must for any given vibration/motion scenario. Understanding the characteristics of different mechanical vibration sources is very important for designing their specific KEH systems.

Considering the mechanical mismatch between the source and transducer, in terms of strain or stress range, resonance frequency, etc., a mechanical transformer (a kind of rigid-body or deformable mechanism) must be used to accommodate the vibration source and transducer input. Many scholars have made in-depth discussions on strain range matching^[38-39] and frequency matching^[14-15]. In particular, numerous mechanical designs were proposed to broaden the energy harvesting bandwidth. These designs include the resonant tuning^[40], multi-mode energy harvester^[41], frequency-up design^[27] and nonlinear structure^[42].

The electromechanical transducer, which enables the power transformation from mechanical to electrical form, is one of the essential components in a KEH system. A transducer has both mechanical and electrical characteristics, which are mutually influenced or coupled in short. The most investigated electromechanical transducers used for KEH include the electromagnetic^[43], piezoelectric^[10,13,43-44], and electrostatic^[45-46] ones. **Table 1** summarizes the pictures and features of these three types of transducers. Due to their different coupling features, people tend to use them in some specific or preferred scenarios. For instance, when the kinetic energy is associated with rapid movements, it is more appropriate to use electromagnetic harvesters; when it is associated with large force or structural deformation, it is more suitable to use piezoelectric harvesters. A mechanical transformer, such as a lever or gear-box, can help match the vibration source and transducer in

terms of their force or displacement ranges. Electrostatic harvesters are more suitable in micro-electromechanical system (MEMS) scale designs. Some flexible transducers, such as triboelectric and soft piezoelectric materials, are very popular in the research communities nowadays^[47-49]. Their electrical principle can be referred to as either of these three types of transducers. The studies of high power-density and flexible transducers have also attracted a lot of research interests from material scientists^[50-53].

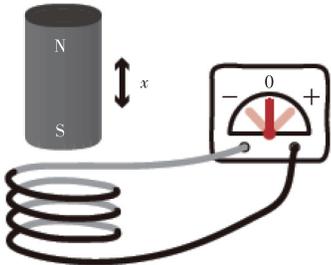
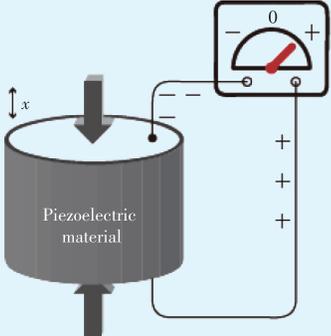
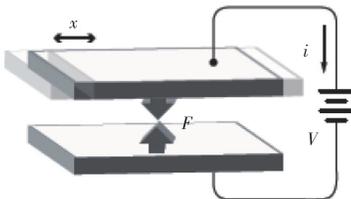
2.3 Power Conditioning Circuit

The electric circuit design was another research hot spot during the last two decades. As shown in Fig. 1, the electrical parts in green connect the mechanical parts in khaki and the cyber parts in blue. The modern power electronics technology provides an effective solution to better harness the incoming fluctuating power flow and adequately satisfying the energy demand for timely information processing. The power conditioning circuits built from fundamental power electronics play an essential role in the KEH-IoT system^[54]. They should be designed toward some objectives, such as

- extracting more power by making proper intervention to the source dynamics;
- achieving high power conversion efficiency under load variation for maintaining efficient use of the extracted energy;
- better mediating the volatile energy supply and the asynchronous energy demand in edge computing.

Generally speaking, all of the aforementioned micro-energy generators produce AC voltages at the open-circuit condition. The power conditioning circuit should be first designed for adapting to the source characteristics. As shown in Fig. 1, the immediate circuit block after the transducer is called the AC-DC interface circuit. The simplest way to convert an AC voltage into DC is through a diode bridge rectifier^[55-58]. Regarding the internal impedance of these sources, they can be classified into two big categories, inductive and capacitive sources. The electromagnetic harvesters are inductive with small internal impedance. Their generated voltage is relatively small (usually from several mV to several V). The piezoelectric harvesters are capacitive with large internal impedance. They give a relatively high voltage output (several V to hundreds of V). The electrostatic ones are also capacitive with an even larger internal impedance. They give the highest voltage output (from tens of V to several kV). Since a diode is a semiconductor device with an almost constant sub-1-V forward voltage drop, the voltage range covered by the piezoelectric case is the easiest to handle, neither too small as the electromagnetic case does, nor too large as the electrostatics case does. The electromagnetic cases usually need a voltage multiplier to passively boost the voltage level for easy utilization. The electrostatic cases usually need better protection to avoid the breakdown under a high voltage beyond the device rating. Moreover, both electromagnetic and piezoelectric sources are self-

▼Table 1. Three types of major electromechanical transducers for kinetic energy harvesting (KEH) and their features

Transducer	Picture	Mechanical Feature	Electrical Feature
Electromagnetic		<ul style="list-style-type: none"> • Large velocity preferred • Complex assembly • Small- to large-scale systems • Need no contact • Bidirectional force 	<ul style="list-style-type: none"> • Small voltage (mV - V) • Large current (mA - A) <ul style="list-style-type: none"> • Inductive source • Small output impedance • Self-generation
Piezoelectric		<ul style="list-style-type: none"> • Large force (hard materials); small force (soft materials) <ul style="list-style-type: none"> • Simple structure • Small- to middle-scale systems <ul style="list-style-type: none"> • Need contact • Bidirectional force 	<ul style="list-style-type: none"> • Large voltage (V - kV) • Small current (nA - μA) <ul style="list-style-type: none"> • Capacitive source • Large output impedance • Self-generation
Electrostatic including triboelectric		<ul style="list-style-type: none"> • Small displacement (out-of-phase); large displacement (in-phase) <ul style="list-style-type: none"> • Simple structure • Small-scale system • Need no contact • Unidirectional force 	<ul style="list-style-type: none"> • Very high voltage (kV) • Very small current (nA) <ul style="list-style-type: none"> • Capacitive source • Very large output impedance • Need a bias-voltage to run (self-generation for triboelectric generator)

generating, while the electrostatic case needs an external source to provide an initial static charge. The external source is used to realize charge separation in a general electrostatic generator^[46]. Therefore, a general electrostatic harvester requires a more complicated interface circuit compared with electromagnetic and piezoelectric cases. However, in some special electrostatic-like generators, such as the triboelectric^[49] and electret^[25] ones, where the charge separation is an inherent capability due to the material property. For these systems, the interface circuits can be simpler without using an external source.

Like most practical electric power sources, a KEH source usually has a maximum power point when the load impedance is at a value between the extremely short-circuit and open-circuit conditions. Besides the detailed technical solutions to the very low voltage in the electromagnetic case or very high voltage in the electrostatic case, more studies focused on the maximum harvested power issue. The most straight-forward investigations considered the resistive matching since harvesting energy from a source more or less brings in some damping ef-

fect^[16-17]. The resistive load was emulated using an energy-dissipative resistor in some early studies^[59-60], and later using a DC-DC buck-boost converter in some later studies toward practical applications^[61]. On the other hand, it was proven that the dynamic effect of a bridge rectifier is different from a resistive load^[16]. It is slightly capacitive in the piezoelectric case^[17] and slightly inductive in the electromagnetic case^[58]. No matter for the pure resistive load or those with the slight reactive component, the power optimization can be realized by implementing an additional maximum power point tracking (MPPT) module to tune the duty cycle of the DC-DC converter, which usually corresponds to a specific one-dimensional impedance trajectory, toward the optimal value.^[23, 55-56, 61]

The term impedance matching is usually referred to in KEH technology to emphasize the importance of load tuning. Different from the MPPT concept, which chases the optimal point on a specific one-dimensional trajectory on the two-dimensional impedance plane, the impedance matching technology seeks the global maximum in the entire two-dimensional complex impedance plane, which is formed by the resistive and reac-

tive components^[62]. The maximum power transfer theorem is a fundamental concept in circuit analysis; however, there were still some studies that only qualitatively quoted such a concept without quantitative analysis. More studies discussing conjugate impedance matching were conducted for piezoelectric KEH systems^[23, 63–64]. For the electromagnetic case, the air coil usually gives small inductance^[65]. The inductive reactance is much lower than the resistance; therefore, the electromagnetic case usually needs resistive matching. For the electrostatic case, the capacitance is usually extremely small^[46]. It produces a large capacitive reactance, which is hard to realize a conjugate matching.

Theoretical studies discussed the impedance matching using linear reactive components, whose values are too large to be put in practical use^[66–67]. The non-dissipative or less-dissipative impedance matching has been realized using the so-called synchronized switch harvesting technology, which is implemented by using modern switched-mode power electronics. Given the capacitive input impedance of piezoelectric transducer, the synchronous electric charge extraction (SECE)^[68] and synchronized switch harvesting on inductor (SSHI)^[69] were proposed. The synchronized switch technology can produce a matched or nearly-matched inductive load in a semi-passive and energy-efficient way; therefore, they can enhance the harvesting capability by several folds under the same excitation condition. These switched-mode circuit solutions have stimulated a large batch of following studies^[11, 41]. Since the interface design in piezoelectric KEH has such a significant effect, the synchronized switch technology has also been utilized in electromagnetic^[70–71], and triboelectric (electrostatic) cases^[72–73] for boosting the harvesting capability.

The KEH technology is designed to replace the large battery storage in some applications. However, this does not mean that the KEH systems have no energy storage devices. Because the input power may be unstable in the KEH-powered system, a relatively small storage device is needed as an energy buffer or filter to improve the reliability in operation. In the previous studies, it was found that the super-capacitors have higher power density and longer lifetime than chemical batteries. Therefore, the super-capacitors are more suitable to provide energy buffers in a KEH system^[74].

Before connecting to a digital circuit, the DC storage voltage needs to be regulated to a specific logical voltage level. We must choose a suitable regulator according to the characteristics of every specific KEH source. For example, piezoelectric transducers provide a relatively high voltage output; therefore, the buck converter is usually used for voltage regulation of piezoelectric systems, while the boost converter is mostly used in electromagnetic systems. LTC3588 (Linear Technology Co.) is the most widely used DC-DC regulator integrated circuit (IC) in piezoelectric KEH^[75]. Different from the conventional pulse-width modulation (PWM) driven DC-DC converter, LTC-3588 works in Burst Mode (a registered trademark of

Linear Technology Corporation), which provides high efficiency at light loads. In the LTC3588, there is an internal under-voltage locking (UVLO) threshold, which is fixed at 5 V (LTC3588-1) or 16 V (LTC3588-2). Only by carefully selecting the appropriate storage capacitor can a good trade-off be made between larger energy storage capacity and rapid charging response. Because the LTC3588 lacks a storage-level indicating signal, to better monitor and manage the stored energy, we often need to add some additional power management circuit design to ensure the energy build-up process^[76–77]. Such a storage-level indicating signal had better be programmable according to the software scheduling. In some other commercial ICs, e.g., BQ25505 (Texas Instrument Co.)^[78], there is a user-defined battery-good indicator, whose threshold can be adjusted offline with a resistor network.

3 KEH-IOT

Given the fluctuating feature of most kinetic energy sources as well as the scarce power output capability, the available energy from a kinetic energy harvester might be unstable and sometimes unpredictable. When building a KEH-IoT system based on these KEH sources, power failures are very likely to happen from time to time. The sensing, computing, and communication tasks carried out by a KEH-IoT system must take sufficient consideration of such limitations and fluctuations. In other words, these battery-free devices are fundamentally different from most conventional computing systems because they violate one of the most basic computing assumptions—a stable power supply^[79]. Therefore, building an organic KEH-IoT system can neither easily replace the battery with an energy harvester nor take an IoT node simply as an equivalent constant resistive load. Exploring the synergy among the mechanical, electrical, and cyber parts is necessary for a comprehensive co-design. To build the cyber part of a KEH-IoT system, one has to know several aspects about the emerging battery-free computing technology, i.e., energy-aware operations, computing modes, and useful development platforms.

3.1 Energy-Neutral and Power-Neutral Operations

In the last decade, tremendous efforts have been made to cope with KEH dynamics and achieve a demand-supply balance. The most commonly adopted model for energy utilization is the energy-neutral operation. In such an operation, we use relatively large energy storage for smoothing out the temporally fluctuating dynamics from both energy supply and load consumption. Another alternative model is power-neutral operation. In this operation model, the computational demand is instantaneously adapted to the harvested power, by which the need for large energy storage is eliminated. Only a small energy buffer is needed for a power-neutral system.

1) Energy-neutral operation: The idea of the energy-neutral operation is to decouple the high-frequency dynamics of the

ambient energy source and power consumption of the IoT load, such that the energy supply and demand can be balanced over a long period. Neglecting the power dissipation during conversion, the energy relation within a period T can be formulated as follows:

$$E(0) + \int_0^T P_h(t)dt = E(T) + \int_0^T P_c(t)dt, \quad (1)$$

where $P_h(t)$ and $P_c(t)$ are the instantaneous harvested power and consumed power at time t ; $E(0)$ and $E(T)$ are the stored energy in the zero and T instants, respectively. Since E_0 and E_T are finite, when T approaches infinite or a sufficiently large number, the energy-neutrality is achieved^[80], i.e.,

$$\int_0^\infty P_h(t)dt = \int_0^\infty P_c(t)dt. \quad (2)$$

This concept of energy-neutral operation has been extensively accepted by many KEH-IoT designers^[81]. For example, Trinity^[82] is a self-sustaining and self-contained indoor sensing system, which is powered by airflow-induced vibration. By adding an energy storage device with suitable capacity, it can smooth the short-term energy variations between supply and demand. Without experiencing a power outage over a reasonable period T , the system might “look like” a battery-powered system. It can perform continuous wind speed and temperature monitoring and wireless communications. In addition, several hardware and software synergistic approaches have been proposed to improve the average power generation \bar{P}_h or reduce the average power consumption \bar{P}_c , such as the dynamic energy burst scaling technology proposed in Refs. [7] and [83]. This technology uses a simple hardware interface consisting of only a few digital inputs. It allows the system to dynamically adjust the supply voltage according to the needs in operation to minimize the energy consumption. Thus, the system using this technology can operate efficiently with smaller energy storage, although in some intervals, the harvested power is much lower than the minimum power requirement of the load.

2) Power-neutral operation: Different from the energy-neutral operation, the power-neutral operation adaptively modulates the instantaneous power consumption to follow the dynamic variation of harvested power; thereby, the required energy storage capacity can be largely reduced. A small energy buffer, such as a μF -level capacitor, is sufficient to support the “once come, immediately serve” operation. Such a concept can mathematically express by taking T as an infinitesimally small number. When $T \rightarrow 0$, the power relation in Eq. (1) can be simplified by differentiating the both sides, such that we can obtain^[80]

$$P_h(t) = P_c(t). \quad (3)$$

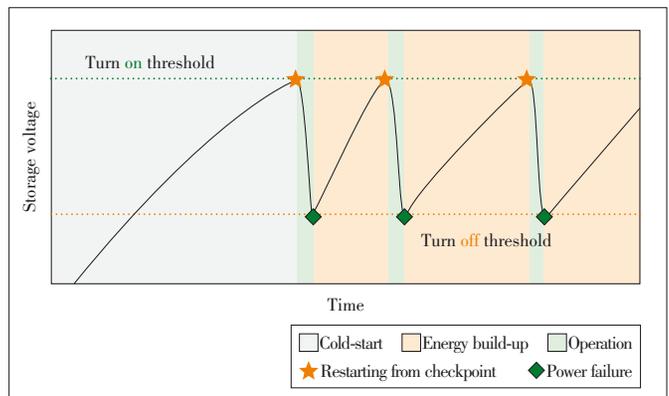
This equation implies that an instantaneous and timely sup-

ply demand balance should be implemented in a power-neutral system. Otherwise, if the application requirement varies independently of the harvested power, it might cause a waste of harvested power in some good-harvest periods or degradation of system performance, or even crash in some other bad-crop periods.

In order to satisfy Eq. (3), the power-neutral system must have a more flexible hardware and/or software strategy according to the stringent energy requirement. It can be carried out by moderating its computing tasks using various controls provided by the digital hardware^[81]. For example, Hibernus^[84-85] achieves a power-neutral behavior on an advanced computing platform. By controlling the central processing unit (CPU) frequency and the number of active cores, Hibernus can modulate the power consumption in real time to match the harvesting dynamics. Compared with the energy-neutral KEH-IoT systems, the power-neutral ones remove the large energy storage unit and its associating charging, monitoring, and conversion circuitry, and accordingly reduce the hardware complexity and improve the energy efficiency.

3.2 Intermittent Computing

Energy-neutral and power-neutral operations only describe the energy service condition of a KEH-IoT system. Each target can be realized by utilizing different software or hardware strategies. Intermittent computing is one of the most extensively studied solutions that ensure reliable and sustainable computing under unpredictable power failures. The working principle of intermittent computing is illustrated in Fig. 3. During the intermittent computing execution, the CPU is activated whenever the storage voltage exceeds the turn-on threshold; it is shut down when the storage voltage is below the turn-off threshold. Therefore, the energy build-up periods and computing periods appear alternatively. By separating the computational progress, an intermittent computing system allows long-running applications to progress incrementally under fluctuat-



▲ Figure 3. Energy picture during intermittent computing. The energy availability depends on environmental conditions and sometimes also the loading effect. Intermittent execution is a side-effect strategy to cope with this battery-free model, where the blackout periods separating the bursts of execution are unknown

ing or scarce energy conditions. In recent years, intermittent computing has attracted much attention in several related areas, such as programming compiler, hardware-software co-design, and non-volatile technique.

From the typical operation picture of intermittent computing shown in Fig. 3, the turning points marked by red pentagrams and green diamonds represent the energy-aware state checkpointing. When a power outage is imminent, a snapshot of the system states, including the program counter, processor registers, static random access memory (SRAM) contents, etc., is immediately saved in the non-volatile memory (e.g., flash memory, ferroelectric random access memory (RAM), magnetoresistive RAM). During the next power burst, the system reboots and restores the states from the stored checkpoint such that program execution resumes. As a result, long-running programs execute gradually in small increments once the accumulated energy is sufficient^[86]. Checkpointing technology has been extensively employed in recent intermittent computing systems. For example, Mementos^[87] supports instrument programs with energy checks in the compiling stage. It provides automatic state checkpointing and recovery in runtime. Ref. [88] implemented a battery-free Game Boy with a just-in-time differential checkpointing scheme. It can efficiently preserve game progress despite power failures.

3.3 Development Platforms

Designing a KEH-IoT system from scratch is time-consuming, costly, and leaves little room for error^[89-90]. In particular, one should be good at all aspects from mechanical structure designs, to power conditioning circuits, to embedded hardware and software. In real-world deployments, it has to not only adapt to various ambient kinetic sources, such as intermittent vibration of bridges, transient human motion actions, or continuous wind excitation. It must also meet the needs of actual IoT end-users for different customized functions, including sensing, computing, and communication. The design of a practical KEH system, even the simplest one, covers the domain knowledge of mechanical engineering, electrical engineering, and computer engineering. It is desirable to divide the design tasks into their individual domain by taking a “modular way of thinking”. However, decoupling by taking simple abstraction usually results in oversimplification and misinterpretation of some practical issues. Most existing studies only demonstrated their performance in one, or at most two, domains. Few have offered a holistic perspective regarding all mechanical, electrical, and cyber considerations, let alone some studies that have also involved materials scientists, civil engineers, biomedical engineers, etc. Building a development platform or test-bed of the KEH-IoT studies is beneficial and important for uniting people from different disciplines. Such infrastructure might

lead to the growth and thriving of all related research and development communities^[89].

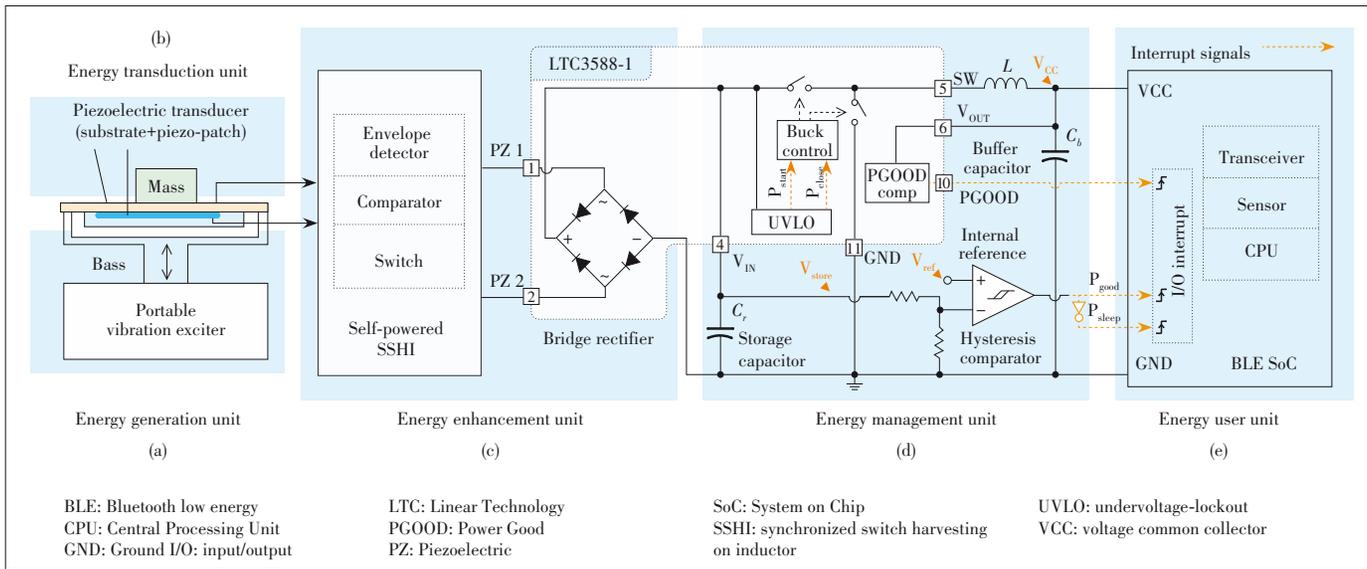
The wireless identification and sensing platform (WISP)^[91] has been developed and become a popular platform for the research of battery-free IoT systems; however, it only supports RF energy harvesting devices. Flicker^[92] is a more versatile battery-free IoT platform. It supports solar, RF, and kinetic energy sources, as well as a variety of communication peripherals. For KEH, Flicker has fully adopted a typical commercialized KEH specific regulator LTC3588^[75, 93], which was designed by Linear Technology Co. ten years ago. The development tool eZ430-RF2500^[94], which was released in 2007 by Texas Instrument Co., is one of the earliest electronic platforms designed for energy harvesting applications.

A vibration-power sensing node (ViPSN)^[77] is the first open-source battery-free IoT platform specified for KEH. It offers all mechanical, electrical, and cyber parts for an efficient demonstration. The major modules are replaceable and extensible toward rapid prototyping and customization of KEH-IoT systems under different excitation conditions and different application scenarios. More importantly, owing to the enhanced energy management unit, which is developed based on the commercialized LTC3588, an energy build-up phase is inserted between the cold-start and normal operation phases; therefore, a reliable and robust computation is reinforced. ViPSN can efficiently and robustly operate under various kinetic excitations, such as harmonic, intermittent, or transient vibrations. It also supports many kinetic energy transducers, such as piezoelectric cantilever^[77], bistable electromagnetic switch^[95-97], and triboelectric generator^[98]. The mechanical computer-aided design (CAD) model, circuit schematics and printed circuit board (PCB) layout, and fundamental firmware codes are all open-source by the Mechatronics and Energy Transformation Laboratory (METAL) of ShanghaiTech University³. **Fig. 4** shows the hardware architecture of ViPSN. It provides all necessary modules from a small vibration generator to Bluetooth low energy (BLE) user unit. It can be regarded as an embodiment of the block diagram shown in Fig. 1. Any peripheral module within the specific energy constraint, such as a low-power sensor, can be added to the KEH-IoT system under this framework for rapid prototyping and efficient customization.

4 Applications

A lot of fundamental and engineering issues as well as possible application designs of KEH^[9, 44] and battery-free IoT^[12] have been discussed in the related research communities during the last decade. On the other hand, KEH-IoT systems have also attracted extensive attention from the industry. There are several scenarios where KEH technologies can significantly

3. <https://github.com/METAL-ShanghaiTech/ViPSN>



▲ Figure 4. ViPSN hardware architecture^[77]

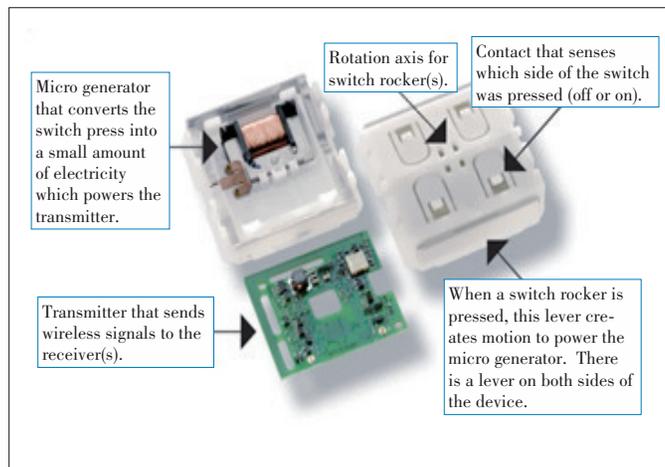
enhance the system’s performances in terms of less or free maintenance, better information acquisition, lower cost, and environmental friendliness.

First, KEH-IoT systems are the most necessary in applications where the IoT devices need to be deployed in the hard-to-reach or too-far-to-reach areas. For example, Perpetuum Ltd.^[99] produced a KEH-based train monitoring system. It can be installed near the wheels for monitoring the vibration conditions of the carriage body to estimate the health condition of the bearings. Since there are many wheels on a train, it is not convenient for frequent maintenance. ReVibe Energy Ltd.^[100] and Xidas Ltd.^[101] have proposed similar products by harnessing vibration energy from machines, railways, and heavy trucks/vehicles. Enervibe Ltd.^[102] manufactures energy harvesters for smart tires and smart shoes. Most of these commercialized harvesters are based on linear electromagnetic harvesters. Each of them can only work in a narrow frequency range.

Second, KEH-IoT systems can be used in applications that require massive deployed IoT nodes (hundreds or even thousands of them). Frequent and massive battery replacement or recharge for all devices are unaffordable. For example, in structural health monitoring, there might be lots of battery-free IoT systems embedded in large structures, such as bridges and highways. The IC solutions are a must toward a low-cost and integrated KEH-IoT system. Besides the big IC companies, such as Analog Devices and Texas Instruments, nowadays, we can also find more and more startups dedicated to the energy harvesting IC technology, ranging from power management IC, such as Nowi Ltd.^[103], to the entire system on chip (SoC) solution, such as Atmosic Ltd.^[104]. Recently, Renesas Co. has just released its RE Family of 32-bit microcontrollers (MCUs) based on the Silicon on Thin Buried Oxide (SOTB) process technology^[105]. These MCUs are equipped with an en-

ergy harvesting control circuit inside, which enables further integration toward many possible applications.

Third, KEH-IoT systems are widely adopted in many human-factored applications, where kinetic energy is extracted from human motions, such as walking, running, jumping, and even finger tapping, to power the motion sensing and transmitting tasks, i.e., toward the self-powered and self-sensing wireless motion sensors. A typical case is the motion-powered light switch first developed by EnOcean GmbH^[95], whose modules are shown in **Fig. 5**. In such a design, an electromagnetic instantaneous-flipping structure is used for generating energy under finger press excitation. Only one finger press can provide sufficient energy for transmitting several wireless packets. There are several similar products on the market, such as those manufactured by ZF GmbH^[106], Alps Alpine Co.^[107], Linptech Co.^[96], and Chlorop Co.^[97]. These battery-free wire-



▲ Figure 5. One of the best commercialized KEH-IoT designs, the battery-free wireless light switch developed by EnOcean GmbH^[95]

less switches are mostly designed for use in smart home and smart manufacturing scenarios. In addition, Pavegen Ltd.^[108] designed a human-powered floor tile, which can harvest footsteps energy to power some environmental monitoring and lighting devices. BinoicPower Ltd.^[109] released an energy harvesting knee brace to generate electricity from natural walking. The company claimed that, over the course of an hour, walking at a comfortable pace, users wearing a harvester on each leg could generate enough power to charge up to four smartphones over an hour^[110]. It is also worth noting that Binoic Power Ltd. and HoverGlide Ltd. are two startups producing suspending backpacks and had two pieces of fundamental research work about energy harvesting from human motions, which were published in the prestigious journals *Science and Nature*^[111-113]. Therefore, the spring of technical transfer of KEH technology is very likely to arrive in the coming decade.

5 Concluding Remarks

KEH-powered IoT is an interesting and emerging topic toward an inspiring vision of ubiquitous battery-free IoT systems for better facilitating our living, manufacturing, etc. It can energize “every sand” in our globe and equip them with perceptivity, intelligence, and connectivity. Tremendous research efforts have been made on either KEH technology or battery-free IoT during the last two decades. The investigations range from pioneering scientific exploration, rigorous engineering modeling and design, to practical application-level development. If the available ambient power is sufficiently large and steady, the easiest way is to take the “modular way of thinking” and consider the mechanical, electrical, and cyber parts as black boxes, as the designs of many modern engineering products did. However, given the practical constraints on size, low power level, volatile source, etc., those modules are usually mutually influenced. People who design the mechanical structure should also have some ideas about the dynamics of power conditioning circuits. Those who design the power conditioning circuits should also understand the possible behaviors of both the mechanical source and the computing behavior. Therefore, the cross-domain knowledge and cyber-electromechanical synergy are beneficial for marching toward a comprehensive and practical KEH-IoT design.

As the topic of KEH-IoT has attracted people across at least five research fields, it is very difficult to enclose all the related work in this short paper. As this paper has cited many comprehensive review articles, readers who are interested in any aspects of KEH-IoT can refer to those detailed papers. It is human nature to start work on the solution of an interdisciplinary problem from his/her own domain knowledge. Some people might think we need more advanced materials for power generation; some might think we need a precise mechanical model and excellent design for better understanding and harnessing kinetic energy; some might think power conditioning is the

most significant difference between these systems and the battery-powered ones; some might consider the computing difficulties without getting sufficient energy, etc. However, instead of thinking in a one-sided style, we should open our horizons, be humble, and respect other shareholders’ opinions, such that we can arrive at inclusive knowledge and holistic solutions toward valuable co-designs and successful products.

Last but not least, energy harvesting is applied research. Although some large-scale systems, such as ocean-wave power generators, are also called energy harvesting, the majority of studies under this umbrella are about small-scale devices. Given its immediate contribution toward practical applications and productivity enhancement, the information attribute of KEH technology should receive prior attention. It directly enables the devices to function and lets the market healthily run. Power conversion and management unit is necessary, but those issues such as conversion efficiency and harvesting capability are less urgent compared to its functionality. Therefore, the energy attribute comes second. Mechanical designs and modeling are necessary too. However, since only a mechanical structure cannot independently generate useful and transmittable information, mechanical dynamics comes the third. Material scientists have made many pioneering inventions and published many high-impact-factor papers; nevertheless, those new findings usually should go through strict engineering approaches (such as modeling, calibration, standardization and fault test), cost-performance evaluation, etc., before made into a product. Therefore, the material advancement also shares the third urgency, in terms of an immediate contribution to the possible industry. In general, the research and investment strategies on the KEH-IoT technology should closely follow the level of productivity and, at the same time, keep a good balance between application development and fundamental research.

References

- [1] BHATTI N A, ALIZAI M H, SYED A A, et al. Energy harvesting and wireless transfer in sensor network applications [J]. *ACM transactions on sensor networks*, 2016, 12(3): 1 - 40. DOI: 10.1145/2915918
- [2] LIU V, PARKS A, TALLA V, et al. Ambient backscatter [J]. *ACM SIGCOMM computer communication review*, 2013, 43(4): 39 - 50. DOI: 10.1145/2534169.2486015
- [3] LI X, MA X Y, ZHANG P L, et al. Escape or exploit? a noise-modulation-based communication under harsh interference [C]//Proc. 7th International Workshop on Real-World Embedded Wireless Systems and Networks. New York, USA: ACM, 2018: 31 - 36. DOI: 10.1145/3277883.3277890
- [4] WANG A, IYER V, TALLA V, et al. FM backscatter: enabling connected cities and smart fabrics [C]//14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17). Boston, USA: USENIX, 2017: 243 - 258
- [5] TALLA V, KELLOGG B, GOLLA KOTA S, et al. Battery-free cellphone [J]. *Proceedings of the ACM on interactive, mobile, wearable and ubiquitous technologies*, 2017, 1(2): 1 - 20. DOI: 10.1145/3090090
- [6] SAFFARI A, HESSAR M, NADERIPARIZI S, et al. Battery-free wireless video

- streaming camera system[C]//IEEE International Conference on RFID (RFID). Phoenix, USA: IEEE, 2019: 1 – 8. DOI: 10.1109/RFID.2019.8719264
- [7] GOMEZ A, SIGRIST L, SCHALCH T, et al. Efficient, long-term logging of rich data sensors using transient sensor nodes [J]. *ACM transactions on embedded computing systems*, 2018, 17(1): 1 – 23. DOI: 10.1145/3047499
- [8] AFANASOV M, BHATTI N A, CAMPAGNA D, et al. Battery-less zero-maintenance embedded sensing at the mithraeum of circus maximus [C]//Proc. 18th Conference on Embedded Networked Sensor Systems. New York, USA: ACM, 2020: 368 – 381. DOI: 10.1145/3384419.3430722
- [9] LIU H C, FU H L, SUN L N, et al. Hybrid energy harvesting technology: from materials, structural design, system integration to applications [J]. *Renewable and sustainable energy reviews*, 2021, 137: 110473. DOI: 10.1016/j.rser.2020.110473
- [10] SAFAEI M, SODANO H A, ANTON S R. A review of energy harvesting using piezoelectric materials: State-of-the-art a decade later (2008 – 2018) [J]. *Smart materials and structures*, 2019, 28(11): 113001. DOI: 10.1088/1361-665x/ab36e4
- [11] BRENES A, MOREL A, JUILLARD J, et al. Maximum power point of piezoelectric energy harvesters: a review of optimality condition for electrical tuning [J]. *Smart materials and structures*, 2020, 29(3): 033001. DOI: 10.1088/1361-665x/ab6484
- [12] MA D, LAN G H, HASSAN M, et al. Sensing, computing, and communications for energy harvesting IoTs: a survey [J]. *IEEE communications surveys & tutorials*, 2020, 22(2): 1222 – 1250. DOI: 10.1109/COMST.2019.2962526
- [13] ANTON S R, SODANO H A. A review of power harvesting using piezoelectric materials (2003 – 2006) [J]. *Smart materials and structures*, 2007, 16(3): R1. DOI: 10.1088/0964-1726/16/3/r01
- [14] TANG L H, YANG Y W, SOH C K. Toward broadband vibration-based energy harvesting [J]. *Journal of intelligent material systems and structures*, 2010, 21(18): 1867 – 1897. DOI: 10.1177/1045389x10390249
- [15] HARNE R L, WANG K W. A review of the recent research on vibration energy harvesting via bistable systems [J]. *Smart materials and structures*, 2013, 22(2): 023001. DOI: 10.1088/0964-1726/22/2/023001
- [16] LESIEUTRE G A, OTTMAN G K, HOFMANN H F. Damping as a result of piezoelectric energy harvesting [J]. *Journal of sound and vibration*, 2004, 269(3/4/5): 991 – 1001. DOI: 10.1016/S0022-460X(03)00210-4
- [17] LIANG J R, LIAO W H. Piezoelectric energy harvesting and dissipation on structural damping [J]. *Journal of intelligent material systems and structures*, 2009, 20(5): 515 – 527. DOI: 10.1177/1045389x08098194
- [18] LIANG J R, LIAO W H. Energy flow in piezoelectric energy harvesting systems [J]. *Smart materials and structures*, 2011, 20(1): 015005. DOI: 10.1088/0964-1726/20/1/015005
- [19] LEFEUVRE E, BADEL A, BRENES A, et al. Analysis of piezoelectric energy harvesting system with tunable SECE interface [J]. *Smart materials and structures*, 2017, 26(3): 035065. DOI: 10.1088/1361-665x/aa5e92
- [20] MOREL A, PILLONNET G, GASNIER P, et al. Frequency tuning of piezoelectric energy harvesters thanks to a short-circuit synchronous electric charge extraction [J]. *Smart materials and structures*, 2019, 28(2): 025009. DOI: 10.1088/1361-665x/aaf0ea
- [21] ZHAO B, WANG J H, LIANG J R, et al. A dual-effect solution for broadband piezoelectric energy harvesting [J]. *Applied physics letters*, 2020, 116(6): 063901. DOI: 10.1063/1.5139480
- [22] SHU Y C, LIEN I C. Analysis of power output for piezoelectric energy harvesting systems [J]. *Smart materials and structures*, 2006, 15(6): 1499 – 1512. DOI: 10.1088/0964-1726/15/6/001
- [23] LIANG J R, LIAO W H. Impedance modeling and analysis for piezoelectric energy harvesting systems [J]. *IEEE/ASME transactions on mechatronics*, 2012, 17(6): 1145 – 1157. DOI: 10.1109/TMECH.2011.2160275
- [24] FAN K Q, CAI M L, WANG F, et al. A string-suspended and driven rotor for efficient ultra-low frequency mechanical energy harvesting [J]. *Energy conversion and management*, 2019, 198: 111820. DOI: 10.1016/j.enconman.2019.111820
- [25] TAO K, TANG L H, WU J, et al. Investigation of multimodal electret-based MEMS energy harvester with impact-induced nonlinearity [J]. *Journal of microelectromechanical systems*, 2018, 27(2): 276 – 288. DOI: 10.1109/JMEMS.2018.2792686
- [26] HU G B, LIANG J R, LAN C B, et al. A twist piezoelectric beam for multi-directional energy harvesting [J]. *Smart materials and structures*, 2020, 29(11): 11LT01. DOI: 10.1088/1361-665x/abb648
- [27] FANG S T, FU X L, DU X N, et al. A music-box-like extended rotational plucking energy harvester with multiple piezoelectric cantilevers [J]. *Applied physics letters*, 2019, 114(23): 233902. DOI: 10.1063/1.5098439
- [28] ZHAO L Y, TANG L H, LIANG J R, et al. Synergy of wind energy harvesting and synchronized switch harvesting interface circuit [J]. *IEEE/ASME transactions on mechatronics*, 2017, 22(2): 1093 – 1103. DOI: 10.1109/TMECH.2016.2630732
- [29] WANG Y, INMAN D J. Experimental validation for a multifunctional wing spar with sensing, harvesting, and gust alleviation capabilities [J]. *IEEE/ASME transactions on mechatronics*, 2013, 18(4): 1289 – 1299. DOI: 10.1109/TMECH.2013.2255063
- [30] AHMED R, MIR F, BANERJEE S. A review on energy harvesting approaches for renewable energies from ambient vibrations and acoustic waves using piezoelectricity [J]. *Smart materials & structures*, 2017, 26(8): 085031
- [31] CHEN Z S, GUO B, YANG Y M, et al. Metamaterials-based enhanced energy harvesting: a review [J]. *Physica B: condensed matter*, 2014, 438: 1 – 8. DOI: 10.1016/j.physb.2013.12.040
- [32] YU J. Review of nonlinear vibration energy harvesting: duffing, bistability, parametric, stochastic and others [J]. *Journal of intelligent material systems and structures*, 2020, 31(7): 921 – 944. DOI: 10.1177/1045389x20905989
- [33] SCRUGGS J T. On the causal power generation limit for a vibratory energy harvester in broadband stochastic response [J]. *Journal of intelligent material systems and structures*, 2010, 21(13): 1249 – 1262. DOI: 10.1177/1045389x10361794
- [34] GUAN M J, LIAO W H. Design and analysis of a piezoelectric energy harvester for rotational motion system [J]. *Energy conversion and management*, 2016, 111: 239 – 244. DOI: 10.1016/j.enconman.2015.12.061
- [35] TAN T, YAN Z. Analytical solution and optimal design for galloping-based piezoelectric energy harvesters [J]. *Applied physics letters*, 2016, 109(25): 253902. DOI: 10.1063/1.4972556
- [36] FU X L, LIAO W H. Nondimensional model and parametric studies of impact piezoelectric energy harvesting with dissipation [J]. *Journal of sound and vibration*, 2018, 429: 78 – 95. DOI: 10.1016/j.jsv.2018.05.013
- [37] GU L, LIVERMORE C. Impact-driven, frequency up-converting coupled vibration energy harvesting device for low frequency operation [J]. *Smart materials and structures*, 2011, 20(4): 045004. DOI: 10.1088/0964-1726/20/4/045004
- [38] EVANS M, TANG L H, TAO K, et al. Design and optimisation of an underfloor energy harvesting system [J]. *Sensors and actuators A: physical*, 2019, 285: 613 – 622. DOI: 10.1016/j.sna.2018.12.002
- [39] LIU H L, HUA R, LU Y, et al. Boosting the efficiency of a footstep piezoelectric-stack energy harvester using the synchronized switch technology [J]. *Journal of intelligent material systems and structures*, 2019, 30(6): 813–822. DOI: 10.1177/1045389x19828512
- [40] LELAND E S, WRIGHT P K. Resonance tuning of piezoelectric vibration energy scavenging generators using compressive axial preload [J]. *Smart materials and structures*, 2006, 15(5): 1413 – 1420. DOI: 10.1088/0964-1726/15/5/030
- [41] SHAHRUZ S M. Design of mechanical band-pass filters for energy scavenging: multi-degree-of-freedom models [J]. *Journal of vibration and control*, 2008, 14(5): 753 – 768. DOI: 10.1177/1077546307083274
- [42] DAQAQ M F, MASANA R, ERTURK A, et al. On the role of nonlinearities in vibratory energy harvesting: a critical review and discussion [J]. *Applied mechanics reviews*, 2014, 66(4): 040801
- [43] BEEBY S P, TORAH R N, TUDOR M J, et al. A micro electromagnetic generator for vibration energy harvesting [J]. *Journal of micromechanics and microengineering*, 2007, 17(7): 1257 – 1265. DOI: 10.1088/0960-1317/17/7/007
- [44] YANG Z B, ZHOU S X, ZU J, et al. High-performance piezoelectric energy harvesters and their applications [J]. *Joule*, 2018, 2(4): 642 – 697. DOI: 10.1016/j.joule.2018.03.011
- [45] TORRES E O, RINCON-MORA G A. Electrostatic energy-harvesting and battery-charging CMOS system prototype [J]. *IEEE transactions on circuits and systems I: regular papers*, 2009, 56(9): 1938 – 1948. DOI: 10.1109/TC-SI.2008.2011578
- [46] BASSET P, BLOKHINA E, GALAYKO D. Electrostatic kinetic energy harvesting [M]. Hoboken, USA: John Wiley & Sons, Inc., 2016. DOI: 10.1002/

- 978119007487
- [47] WANG Z L, SONG J H. Piezoelectric nanogenerators based on zinc oxide nanowire arrays [J]. *Science*, 2006, 312(5771): 242 – 246. DOI: 10.1126/science.1124005
- [48] QIN Y, WANG X D, WANG Z L. Microfibre – nanowire hybrid structure for energy scavenging [J]. *Nature*, 2008, 451(7180): 809 – 813. DOI: 10.1038/nature06601
- [49] WANG Z L, LIN L, CHEN J, et al. *Triboelectric nanogenerators* [M]. Heidelberg, Germany: Springer, 2016
- [50] NARITA F, FOX M. A review on piezoelectric, magnetostrictive, and magneto-electric materials and device technologies for energy harvesting applications [J]. *Advanced engineering materials*, 2018, 20(5): 1700743. DOI: 10.1002/adem.201700743
- [51] PRIYA S, SONG H C, ZHOU Y, et al. A review on piezoelectric energy harvesting: materials, methods, and circuits [J]. *Energy harvesting and systems*, 2019, 4(1). DOI: 10.1515/EHS-2016-0028
- [52] WANG Z L, JIANG T, XU L. Toward the blue energy dream by triboelectric nanogenerator networks [J]. *Nano energy*, 2017, 39: 9 – 23. DOI: 10.1016/j.nanoen.2017.06.035
- [53] XU W H, ZHENG H X, LIU Y, et al. A droplet-based electricity generator with high instantaneous power density [J]. *Nature*, 2020, 578(7795): 392 – 396. DOI: 10.1038/s41586-020-1985-6
- [54] SZARKA G D, STARK B H, BURROW S G. Review of power conditioning for kinetic energy harvesting systems [J]. *IEEE transactions on power electronics*, 2012, 27(2): 803 – 815. DOI: 10.1109/TPEL.2011.2161675
- [55] OTTMAN G K, HOFMANN H F, LESIEUTRE G A. Optimized piezoelectric energy harvesting circuit using step-down converter in discontinuous conduction mode [J]. *IEEE transactions on power electronics*, 2003, 18(2): 696 – 703. DOI: 10.1109/TPEL.2003.809379
- [56] OTTMAN G K, HOFMANN H F, BHATT A C, et al. Adaptive piezoelectric energy harvesting circuit for wireless remote power supply [J]. *IEEE transactions on power electronics*, 2002, 17(5): 669 – 676. DOI: 10.1109/TPEL.2002.802194
- [57] LOONG C N, CHANG C C, DIMITRAKOPOULOS E G. Circuit nonlinearity effect on the performance of an electromagnetic energy harvester-structure system [J]. *Engineering structures*, 2018, 173: 449 – 459. DOI: 10.1016/j.engstruct.2018.06.090
- [58] LIANG J R, GE C, SHU Y C. Impedance modeling of electromagnetic energy harvesting system using full-wave bridge rectifier [C]//*Smart Structures and Materials + Nondestructive Evaluation and Health Monitoring Conference*. Proc. SPIE 10164, Active and Passive Smart Structures and Integrated Systems, Portland, USA: SPIE, 2017: 101642N. DOI: 10.1117/12.2259870
- [59] ROUNDY S, LELAND E S, BAKER J, et al. Improving power output for vibration-based energy scavengers [J]. *IEEE pervasive computing*, 2005, 4(1): 28 – 36. DOI: 10.1109/MPRV.2005.14
- [60] MITCHESON P D, YEATMAN E M, RAO G K, et al. Energy harvesting from human and machine motion for wireless electronic devices [J]. *Proceedings of the IEEE*, 2008, 96(9): 1457 – 1486. DOI: 10.1109/JPROC.2008.927494
- [61] KONG N, HA D S, ERTURK A, et al. Resistive impedance matching circuit for piezoelectric energy harvesting [J]. *Journal of intelligent material systems and structures*, 2010, 21(13): 1293 – 1302. DOI: 10.1177/1045389x09357971
- [62] ZHAO B, LIANG J. On the circuit solutions towards broadband and high-capability piezoelectric energy harvesting systems [M]//*Active and Passive Smart Structures and Integrated Systems XII*, vol. 10595. Bellingham, USA: SPIE, 2018: 105950E
- [63] BRUFAU-PENELLA J, PUIG-VIDAL M. Piezoelectric energy harvesting improvement with complex conjugate impedance matching [J]. *Journal of intelligent material systems and structures*, 2009, 20(5): 597 – 608
- [64] KIM H, PRIYA S, STEPHANOU H, et al. Consideration of impedance matching techniques for efficient piezoelectric energy harvesting [J]. *IEEE transactions on ultrasonics, ferroelectrics, and frequency control*, 2007, 54(9): 1851 – 1859. DOI: 10.1109/TUFFC.2007.469
- [65] CHENG S, WANG N G, ARNOLD D P. Modeling of magnetic vibrational energy harvesters using equivalent circuit representations [J]. *Journal of micromechanics and microengineering*, 2007, 17(11): 2328 – 2335. DOI: 10.1088/0960-1317/17/11/021
- [66] FLEMING A J, BEHRENS S, MOHEIMANI S O R. Synthetic impedance for implementation of piezoelectric shunt-damping circuits [J]. *Electronics letters*, 2000, 36(18): 1525. DOI: 10.1049/el:20001083
- [67] PARK C H, INMAN D J. Enhanced piezoelectric shunt design [J]. *Shock and vibration*, 2003, 10(2): 127 – 133. DOI: 10.1155/2003/863252
- [68] LEFEUVRE E, BADEL A, RICHARD C, et al. Piezoelectric energy harvesting device optimization by synchronous electric charge extraction [J]. *Journal of intelligent material systems and structures*, 2005, 16(10): 865 – 876. DOI: 10.1177/1045389x05056859
- [69] GUYOMAR D, BADEL A, LEFEUVRE E, et al. Toward energy harvesting using active materials and conversion improvement by nonlinear processing [J]. *IEEE transactions on ultrasonics, ferroelectrics, and frequency control*, 2005, 52(4): 584 – 595. DOI: 10.1109/TUFFC.2005.1428041
- [70] LOMBARDI G, LALLART M. Synchronous electric charge and induced current extraction (SECICE): a unified nonlinear technique combining piezoelectric and electromagnetic harvesting [J]. *Smart materials and structures*, 2021, 30(2): 025029. DOI: 10.1088/1361-665x/abd346
- [71] LALLART M, LOMBARDI G. Synchronized switch harvesting on electromagnetic system: a nonlinear technique for hybrid energy harvesting based on active inductance [J]. *Energy conversion and management*, 2020, 203: 112135. DOI: 10.1016/j.enconman.2019.112135
- [72] LI X, SUN Y. An SSHI rectifier for triboelectric energy harvesting [J]. *IEEE transactions on power electronics*, 2020, 35(4): 3663 – 3678. DOI: 10.1109/TPEL.2019.2934676
- [73] XU S X, DING W B, GUO H Y, et al. Boost the performance of triboelectric nanogenerators through circuit oscillation [J]. *Advanced energy materials*, 2019, 9(30): 1900772. DOI: 10.1002/aenm.201900772
- [74] GUAN M J, LIAO W H. Characteristics of energy storage devices in piezoelectric energy harvesting systems [J]. *Journal of intelligent material systems and structures*, 2008, 19(6): 671 – 680. DOI: 10.1177/1045389x07078969
- [75] LTC3588 - 1 nanopower energy harvesting power supply [R]. Milpitas, USA: Linear Technologie, 2010
- [76] HUANG Q Y, MEI Y, WANG W, et al. Toward battery-free wearable devices: the synergy between two feet [J]. *ACM transactions on cyber-physical systems*, 2018, 2(3): 20. DOI: 10.1145/3185503
- [77] LI X, TENG L, TANG H, et al. ViPSN: A vibration-powered IoT platform [J]. *IEEE Internet of Things journal*, 2021, 8(3): 1728 – 1739. DOI: 10.1109/JIOT.2020.3016993
- [78] BQ25505: ultra low-power boost charger with battery management and autonomous power multiplexer for primary battery in energy harvester applications [R]. Dallas, USA: Texas Instruments, 2013
- [79] HESTER J, SORBER J. Batteries not included [J]. *XRDS: crossroads, the ACM magazine for students*, 2019, 26(1): 23 – 27. DOI: 10.1145/3351474
- [80] SLIPER S T, CETINKAYA O, WEDDELL A S, et al. Energy-driven computing [J]. *Philosophical transactions of the royal society A: mathematical, physical and engineering sciences*, 2020, 378(2164): 20190158. DOI: 10.1098/rsta.2019.0158
- [81] MERRETT G V, AL-HASHIMI B M. Energy-driven computing: rethinking the design of energy harvesting systems [C]//*Design, Automation & Test in Europe Conference & Exhibition (DATE)*. Lausanne, Switzerland: IEEE, 2017: 960 – 965. DOI: 10.23919/DATE.2017.7927130
- [82] XIANG T, CHI Z C, LI F, et al. Powering indoor sensing with airflows: a trinity of energy harvesting, synchronous duty-cycling, and sensing [C]//*Proc. 11th ACM Conference on Embedded Networked Sensor Systems*. New York, USA: ACM, 2013: 16. DOI: 10.1145/2517351.2517365
- [83] GOMEZ A, SIGRIST L, MAGNO M, et al. Dynamic energy burst scaling for transiently powered systems [C]//*Design, Automation & Test in Europe Conference & Exhibition (DATE)*. Dresden, Germany: IEEE, 2016: 349 – 354
- [84] BALSAMO D, WEDDELL A S, MERRETT G V, et al. Hibernus: sustaining computation during intermittent supply for energy-harvesting systems [J]. *IEEE embedded systems letters*, 2015, 7(1): 15 – 18. DOI: 10.1109/LES.2014.2371494
- [85] BALSAMO D, WEDDELL A S, DAS A, et al. Hibernus: a self-calibrating and adaptive system for transiently-powered embedded devices [J]. *IEEE transactions on computer-aided design of integrated circuits and systems*, 2016, 35(12): 1968 – 1980. DOI: 10.1109/TCAD.2016.2547919
- [86] JAYAKUMAR H, RAHA A, RAGHUNATHAN V. QUICKRECALL: A low overhead HW/SW approach for enabling computations across power cycles in transiently powered computers [C]//*27th International Conference on VLSI Design and 13th International Conference on Embedded Systems*. Mumbai, India:

- IEEE, 2014: 330 – 335. DOI: 10.1109/VLSID.2014.63
- [87] RANSFORD B, SORBER J, FU K. Mementos: system support for long-running computation on RFID-scale devices [C]//Proc. sixteenth international conference on architectural support for programming languages and operating systems. New York, USA: ACM, 2011: 159 – 170. DOI: 10.1145/1950365.1950386
- [88] DE WINKEL J, KORTBEEK V, HESTER J, et al. Battery-free game boy [J]. Proceedings of the ACM on interactive, mobile, wearable and ubiquitous technologies, 2020, 4(3): 1 – 34. DOI: 10.1145/3411839
- [89] HESTER J, SORBER J. The future of sensing is batteryless, intermittent, and awesome[C]//Proc. 15th ACM Conference on Embedded Network Sensor Systems. New York, USA: ACM, 2017: 1 – 6. DOI: 10.1145/3131672.3131699
- [90] LUCIA B, BALAJI V, COLIN A, et al. Intermittent computing: challenges and opportunities [C]//2nd Summit on Advances in Programming Languages (SNAPL 2017). Asilomar, USA: PL Community, 2017. DOI: 10.4230/LIPIcs.SNAPL.2017.8
- [91] SAMPLE A P, YEAGER D J, POWLEDGE P S, et al. Design of an RFID-based battery-free programmable sensing platform [J]. IEEE transactions on instrumentation and measurement, 2008, 57(11): 2608 – 2615. DOI: 10.1109/TIM.2008.925019
- [92] HESTER J, SORBER J. Flicker: rapid prototyping for the batteryless internet-of-things [C]//Proc. 15th ACM Conference on Embedded Network Sensor Systems. New York, USA: ACM, 2017: 19. DOI: 10.1145/3131672.3131674
- [93] WHITAKER M. Energy harvester produces power from local environment, eliminating batteries in wireless sensors [J]. Journal of analog innovation, 2010, 20(1): 1 – 8
- [94] Texas Instruments Co. eZ430-RF2500 development tool user's guide [EB/OL]. [2020-12-12]. <https://www.ti.com/lit/ug/slau227f/slau227f.pdf>
- [95] Technology of EnOcean GmbH [EB/OL]. [2020-12-13]. <https://www.enocean.com/en/products>
- [96] Technology of Linptech Ltd. [EB/OL]. [2020-12-13]. <http://www.linptech.com>
- [97] Technology of Chlorop Ltd. [EB/OL]. [2020-12-13]. <http://www.chlorop.com>
- [98] ZHANG J X, GONG S B, LI X, et al. A wind-driven poly (tetrafluoroethylene) electret and polylactide polymer-based hybrid nanogenerator for self-powered temperature detection system [J]. Advanced sustainable systems, 2021, 5(1): 2000192. DOI: 10.1002/adsu.202000192
- [99] Technology of Perpetuum Ltd. [EB/OL]. [2020-12-15]. <https://perpetuum.com/technology>
- [100] Technology of ReVibe Energy Ltd. [EB/OL]. [2020-12-15]. <https://revibeenergy.com>
- [101] Technology of Xidas Ltd. [EB/OL]. [2020-12-15]. <https://xidasiot.com>
- [102] Technology of Enervibe Ltd. [EB/OL]. [2020-12-15]. <https://enervibe.co>
- [103] Technology of NOWI Ltd. [EB/OL]. [2020-12-16]. <https://www.nowi-energy.com>
- [104] Technology of Atmosic Ltd. [EB/OL]. [2020-12-16]. <https://www.atmosic.com>
- [105] Renesas Electronics Co. RE Cortex-M0+Ultra-low Power SOTB MCUs. [EB/OL]. [2020-12-16]. <https://www.renesas.com>
- [106] Technology of ZF GmbH [EB/OL]. [2020-12-20]. <https://www.zf.com>
- [107] Alps Alpine Co. Spga series energy harvester [EB/OL]. [2020-12-20]. <https://tech.alpsalpine.com/prod/e/html/harvester>
- [108] Technology of Pavegen Ltd. [EB/OL]. [2020-12-20]. <https://pavegen.com>
- [109] Technology of Bionic Power Ltd. [EB/OL]. [2020-12-20]. <https://www.bionic-power.com>
- [110] MaizeKennedy. Power from the People? A Long Way to Go. [EB/OL]. [2020-12-20]. <https://www.powermag.com>
- [111] ROME L C, FLYNN L, YOO T D. Biomechanics: rubber bands reduce the cost of carrying loads [J]. Nature, 2006, 444(7122): 1023 – 1024. DOI: 10.1038/4441023a
- [112] ROME L C, FLYNN L, GOLDMAN E M, et al. Generating electricity while walking with loads [J]. Science, 2005, 309(5741): 1725 – 1728. DOI: 10.1126/science.1111063
- [113] DONELAN J M, LI Q, NAING V, et al. Biomechanical energy harvesting: generating electricity during walking with minimal user effort [J]. Science, 2008, 319(5864): 807 – 810. DOI: 10.1126/science.1149860

Biographies

LIANG Junrui (liangjr@shanghaitech.edu.cn) received the B.E. and M.E. degrees in instrumentation engineering from Shanghai Jiao Tong University, China in 2004 and 2007, respectively, and the Ph.D. degree in mechanical and automation engineering from the Chinese University Hong Kong, China in 2010. He is currently an assistant professor with the School of Information Science and Technology, ShanghaiTech University, China. His research interests include energy conversion and power conditioning circuits, kinetic energy harvesting and vibration suppression, IoT devices, and mechatronics. Dr. LIANG has published 84 technical papers in the leading international academic journals and conferences. He has received two Best Paper Awards in the IEEE International Conference on Information and Automation in 2009 and 2010, respectively. He is an associate editor of *IET Circuits, Devices and Systems* and the general chair of the Second International Conference on Vibration and Energy Harvesting Applications in 2019.

LI Xin received the B.E. and B.Ec. degrees from the North University of China, 2016. He is currently pursuing the Ph.D. degree with the School of Information Science and Technology, ShanghaiTech University, China. His research interests include vibration energy harvesting, ubiquitous computing, and Internet of Things. He was a recipient of the First Place of the International Conference on Embedded Wireless Systems and Networks Dependability Competition in 2019, the First Runner Up of the IEEE Industrial Electronics Society Inter-Chapter Paper Competition in 2019, and Best Student Hardware Award Finalist in ASME Smart Materials, Adaptive Structures, and Intelligent Structures Conference (SMASIS) 2020.

YANG Hailiang received the B.E. degree in electronic information science and technology from Wuhan University of Technology, China in 2020. He is now working towards his master's degree at ShanghaiTech University, China. His research interests include the Internet of Things and energy harvesting.



Next Generation Semantic and Spatial Joint Perception — Neural Metric–Semantic Understanding

ZHU Fang^{1,2}

(1. State Key Laboratory of Mobile Network and Mobile Multimedia Technology, Shenzhen 518057, China;
2. ZTE Corporation, Shenzhen 518057, China)

DOI: 10.12142/ZTECOM.202101008

<http://kns.cnki.net/kcms/detail/34.1294.TN.20210218.1753.002.html>, published online February 19, 2021

Manuscript received: 2020–12–25

Abstract: Efficient perception of the real world is a long-standing effort of computer vision. Modern visual computing techniques have succeeded in attaching semantic labels to thousands of daily objects and reconstructing dense depth maps of complex scenes. However, simultaneous semantic and spatial joint perception, so-called dense 3D semantic mapping, estimating the 3D geometry of a scene and attaching semantic labels to the geometry, remains a challenging problem that, if solved, would make structured vision understanding and editing more widely accessible. Concurrently, progress in computer vision and machine learning has motivated us to pursue the capability of understanding and digitally reconstructing the surrounding world. Neural metric-semantic understanding is a new and rapidly emerging field that combines differentiable machine learning techniques with physical knowledge from computer vision, e.g., the integration of visual-inertial simultaneous localization and mapping (SLAM), mesh reconstruction, and semantic understanding. In this paper, we attempt to summarize the recent trends and applications of neural metric-semantic understanding. Starting with an overview of the underlying computer vision and machine learning concepts, we discuss critical aspects of such perception approaches. Specifically, our emphasis is on fully leveraging the joint semantic and 3D information. Later on, many important applications of the perception capability such as novel view synthesis and semantic augmented reality (AR) contents manipulation are also presented. Finally, we conclude with a discussion of the technical implications of the technology under a 5G edge computing scenario.

Keywords: visual computing; semantic and spatial joint perception; dense 3D semantic mapping; neural metric-semantic understanding

Citation (IEEE Format): F. Zhu, “Next generation semantic and spatial joint perception—neural metric-semantic understanding,” *ZTE Communications*, vol. 19, no. 1, pp. 61 – 71, Mar. 2021. doi: 10.12142/ZTECOM. 202101008.

1 Introduction

The perception of the real world in a meaningful reconstructive way has been one of the primary driving forces for the development of sophisticated computer vision techniques. The semantic and spatial joint perception of a variety of scenes is shown in **Fig. 1**. Computer vision approaches span a range from real-time mapping, which enables the latest generation of robots, to sophisticated semantic identification for the meaningfully structured information in various big data applications. In both cases, one of the main bottlenecks is the exact and consistent context understanding in terms of occlusion, view-angle, and illumination conditions, i.e., despite of the noticeable progress in fine-grained semantic scene understanding tasks like detection and instance seg-

mentation, computers still perform unsatisfactorily on visually understanding humans in crowded scenes. Concurrently, powerful consistent context understanding models have emerged in the computer vision and machine learning communities. The seminal works related to semantic and spatial joint perception, the so-called dense 3D semantic mapping framework by HERMANS et al.^[1], have evolved in recent years into joint volumetric 3D reconstruction and semantic segmentation formulas for both the unmanned system and the human-involved virtual/augmented reality (VR/AR) immersive experience. Here, the synthesis of more plausible depth in parts of the scene or more reliable semantic image classification can be achieved by jointly optimizing geometry and semantics in 3D. Very recently, such an area has been explored as “metric-semantic understanding”. One of the first publications that used

the term metric-semantic understanding is Kimera^[2]. It enables machines to learn to perceive their surroundings by combining the-state-of-the-art geometric and semantic understanding into a modern perception way. Furthermore, the authors also argue that the semantic information based on the geometric information provides the ideal level of abstraction to provide humans with models of the environment that are easy to understand. Instead of implicitly combining the geometry and semantic segmentation of 3D, a variety of other methods more explicitly follow this notion of collaboration to exploit components of the perception pipeline.

While classical computer vision starts from the affine imaging of the physical world to addressing the geometrical consistency by modeling, for example, the camera's viewpoint, odometry, and depth map properties, machine learning comes from an end-to-end trainable (differentiable) and statistical perspective. It is a well-known fact that the differentiable machine learning technique can capture more complex dependencies and achieve a high level of expressiveness, while, if used only, cannot be metric or explicitly follow the strict consistency behind the physical world. To this end, the quality of mainly traditional computer vision-based dense 3D semantic mapping relies on the physical correctness of the employed models. Direct joint estimation of geometry and semantics in a multi-view 3D reconstruction setting, which implicitly combines the geometry and semantic information in the scenes, is hard and error-prone and leads to artifacts in the reconstructed map. Thus, the classic computer vision-based geometry reconstructions suffer from not only classical issues, such as poorly textured areas, repetitive patterns, and occlusions, but also several additional challenges, such as higher noise level, and, often, the presence of shake and motion blur. To this end, traditional metric-semantic understanding methods try to overcome these issues by using heuristic regularization, like convex anisotropic regularizers, to combine captured imagery. But in the complex scenery, these methods require thousands of iterations for convergence or are unable to fully capture the complex semantic and geometric dependencies behind them. Neural metric-semantic understanding brings the promise of addressing both geometry reconstruction and fusion of geometry and semantic information by using deep networks to learn complex mappings from captured images to 3D semantic maps. The underlying principle is to combine the differentiable machine learning techniques with physical knowledge from computer vision to yield new and powerful algorithms for semantic and spatial collaborative perception.

Neural metric-semantic understanding does not yet have a clear definition in the literature. Here, we define neural metric-semantic understanding as: deep image or video semantic & spatial collaborative perception approaches and also sub-modules that enable the explicit or implicit fusion of semantic and geometric context properties of the scene, such as deep convolutional neural networks in volumetric space for 3D se-



▲Figure 1. Semantic and spatial joint perception of a variety of scenes^[2-3]

semantic segmentation, incorporation of conventional multi-view stereo concepts within a deep learning framework, fine-tuning of the deep network by using the extracted geometric constraints, and a representation of semantics as an invariant scene for medium-term continuous tracking of large scale 3D scanning.

This paper defines the components of the semantic and spatial collaborative perception pipeline and exploits the different directions of neural metric-semantic understanding formulations, embedded in corresponding components. One central scheme around which we structure this paper is the combination of computer vision imaging principles and learning-based primitives to yield new and powerful algorithms for visual content's consistent understanding, since consistency in the real-world understanding is essential for many media editing and structural data indexing applications. We start by discussing previous explorations' fundamental concepts and components of metric-semantic understanding, which are prerequisites for the semantic and spatial collaborative perception pipeline. Afterwards, we discuss critical aspects of emerging neural-based metric-semantic understanding approaches, fusions of learning-based primitives and affine imaging principles, such as type of fusion, how the fusion is provided, which components of the metric-semantic understanding pipeline are learned, and explicit v.s. implicit fusion. Following, we discuss the panorama of applications that is enabled by semantic and spatial collaborative perception. The applications range from relighting, novel view synthesis, to the manipulation of semantic contents for augmented reality (AR). The semantic manipulation of AR contents, achieving natural interaction between the virtual and real world and finally facilitating natural interaction between "digital twins" and the real world, has many technical implications on the evolving storage-computing network, especially when instant response computing and privacy preserving strategies can be carried out with the help of edge computing based on 5G. We then conclude with these implications.

2 Related Surveys

Metric-semantic understanding, sometimes called “dense 3D semantic mapping”, has been continuously studied in the literature, such as Ref. [2] and Refs. [4 – 8]. It includes robot perception and mixed reality. The perceptual understanding using classic computer vision or with some convolutional neural networks (CNNs) as classification assistance has been studied extensively. The thorough analysis survey^[9] of such classical computer vision methods, for the implicit combination of the geometry and semantic segmentation of 3D, focuses on specific heuristic regularization, such as surface normal directions^[10] and special treatment for highly reflective objects^[11]. Recent explorations regarding explicitly semantic and spatial collaborative perception through the components of the perception pipeline, with the emerging machine learning techniques, have also been discussed in Refs. [12 – 15]. Recent reports, like Refs. [16 – 19], discuss various applications with the help of metric-semantic understanding techniques, such as novel view synthesis, relighting, and semantic AR contents manipulation. However, none of the above reports or literature provides a structured or comprehensive look into the new and rapidly emerging field, neural metric-semantic understanding, which combines differentiable machine learning techniques with physical knowledge. Such a comprehensive approach, especially linking clues from classic computer vision to the “new” neural assistance, is critical, since the “next generation” semantic and spatial collaborative perception can reach new heights in the performance of these tasks, which motivates us to pursue the modern computer vision capability of understanding and digitally reconstructing the surrounding world.

3 Theoretical Fundamentals

In this section, we discuss the theoretical fundamentals of working in the semantic and spatial collaborative perception space. First, we discuss dense depth map formation models in computer vision, followed by the classic methods of high-quality 3D scanning of large-scale scenes. Next, we discuss approaches to semantic generative models in deep learning. In the end, we discuss the core principles of volumetric semantic 3D reconstruction.

3.1 Dense Depth Map Formation

Classical computer vision methods approximate the reverse prediction process of image formation in the real world. Light sources emit photons that interact with the objects in the scene, as a function of their geometry and material properties, before being recorded by multiple cameras with overlapping views. This process is known as dense depth estimation. Early passive stereo methods, referred to as an in-depth analysis in Ref. [20], relied on at least two recorded frames based on the known camera geometry to extract stereo correspondence, the so-called dense disparity map. Among them, some multi-view

stereo methods use multi-valued, voxel-based, or layer-based presentations, while most stereo correspondence methods compute a univalued disparity function $d(x, y)$ with respect to a reference image. The central element to methods that produce a univalued disparity map $d(x, y)$ is the concept of a disparity space (x, y, d) . The term disparity describes the difference in the location of corresponding features seen by the left and right eyes. The correspondence between a pixel (x, y) in reference image r and a pixel (x', y') in matching image m is then given by Eq. (1). And the common steps in the stereo algorithms generally include matching cost computation, support aggregation, disparity computation, and disparity optimization. The actual sequence of steps taken depends on the specific algorithm.

$$x' = x + sd(x, y), \quad y' = y. \quad (1)$$

Passive stereo matching algorithms work well on textured scenes but require demanding computation. Later on, active stereo methods (e.g. Kinect), which triangulate correspondences between a structured active illumination and a camera, have raised a lot of interest. While unstructured surfaces are no longer a problem, the lateral resolution of the active stereo-only methods is limited by the resolution of the projection system under the constraint of size or power. Currently, accurate real-time dense depth estimation is mostly fulfilled with the fusion of sensors, which ultimately improves speed, robustness and quality. A thorough re-inspection regarding the classical paradigm and the fusion between the time of flight (ToF) and stereo, can refer to Ref. [21]. To exploit the complementary strengths, accurate but sparse active range measurements and the ambiguous but dense passive stereo information must be fused under the principle described in Eq. (2) below.

$$E(d) = w_{\text{stereo}} E_{\text{Stereo}}(d) + w_{\text{ToF}} E_{\text{ToF}}(d | d_{\text{ToF}}) + R_{\text{smooth}} + R_{\text{temp}}, \quad (2)$$

where w_{stereo} and w_{ToF} represent confidence/weights, E represents the objective energy to be minimized, and R represents the regularizer.

Different optimization strategies can refer to variably concrete formulas corresponding to the principle described in Eq. (2), such as the local method in Eq. (3) and the variational framework in Eq. (4).

$$E(z_i) = w E_{\text{ToF}}(z_i | z_i^{\text{ToF}}) + (1 - w) E_{\text{stereo}}(z_i), \quad (3)$$

$$E_{\text{data}}(u) := \int_{\Omega} \chi_{\text{ToF}}(x) \rho_{\text{ToF}}(u(x)) + \chi_{\text{Stereo}}(x) \rho_{\text{Stereo}}(u(x)) dx. \quad (4)$$

In Eq. (4), ρ represents the local term for penalizing the deviation from the ToF or stereo data, and X represents spatial

indicator functions for valid/trusted ToF/stereo.

3.2 3D Scanning of Large-Scale Scenes

Given the accurate dense depth map of the observed view, high-quality consistent 3D scanning of large-scale scenes is the next key step to the geometric and photometric registration between the virtual and real world. The most important tasks under the objective are estimating globally optimized poses, robust tracking with recovery from gross tracking failures, and re-estimating the 3D model to ensure global consistency, as mentioned by DAI et al.^[22]. The core of the above tasks is a robust pose estimation strategy, which globally optimizes the camera trajectory per frame, considering the complete history of the single view depth and image input in an efficient local-to-global hierarchical optimization framework, as described in Refs. [22 – 24]. While each has trade-offs, global optimization methods based on implicit bundle adjustment (BA) are the de facto methods for the highest quality reconstructions. Finally, the optimization for both dense photometric and geometric alignment is based on the energy illustrated in Eq. (5):

$$\begin{aligned}
 E_{icp} &= \sum_k ((v^k - \exp(\hat{\xi})Tv^k) \cdot n^k)^2, \\
 E_{rgb} &= \sum_{u \in \Omega} (I(u, C_t^i) - I(\pi(K \exp(\hat{\xi})Tp(u, D_t^i)), \hat{C}_{t-1}^i))^2, \\
 E_{track} &= E_{icp} + w_{rgb}E_{rgb},
 \end{aligned} \quad (5)$$

where v^k represents the back-projection of the k -th vertex and n^k is the corresponding normal; D represents the live depth map and C represents the live color image; ξ is the motion parameter and $\exp(\xi)$ is the matrix exponential that maps a member of the Lie Algebra $se3$ to a member of the corresponding Lie group $SE3$; T is the current estimate of the transformation from the previous camera pose to the current one; E represents the cost function that needs to be minimized and w represents manually defined weights.

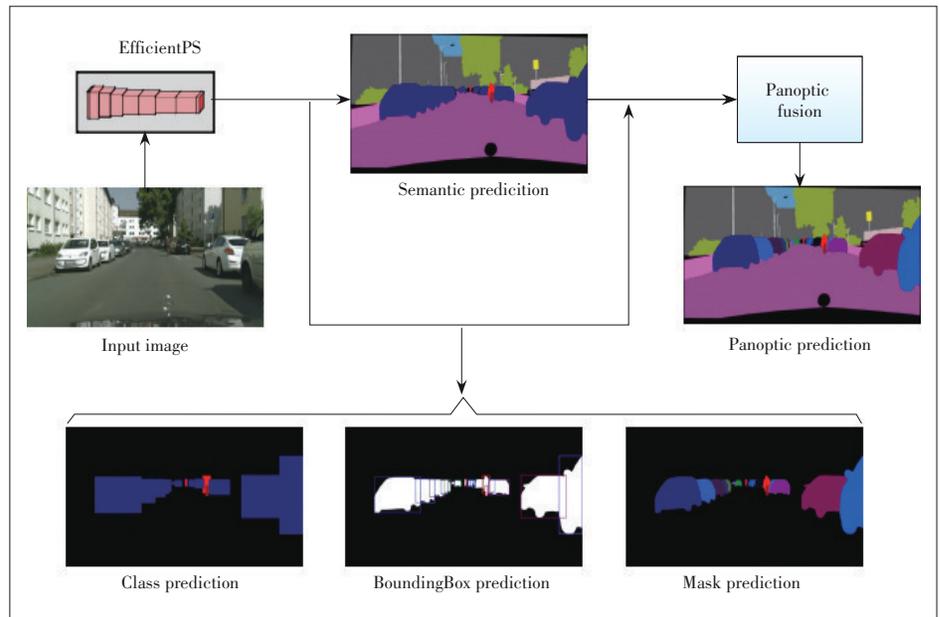
3.3 Semantic Understanding

Besides the geometric and photometric registration following the above methods, semantic generative models assist in semantic content registration of the corresponding large-scale scenes. Such scene comprehension, which necessitates recognizing instances of scene participants along with general scene semantics, can be addressed by the panoptic segmentation task with corresponding semantic generative models such as those in

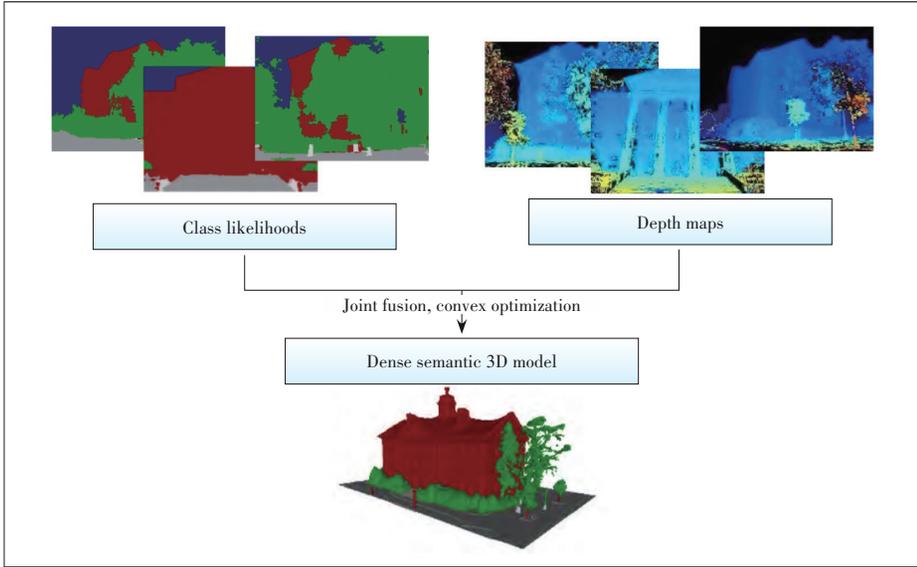
Refs. [25 – 26]. Such semantic generative models generally need a deep neural network (e.g, Feature Pyramid Network) as a backbone to efficiently encode and fuse semantically rich multi-scale features, which is followed by a panoptic head network to extract coherently understandable visual scenes at both the fundamental pixel level and distinctive object instance level, as shown in **Fig. 2**. The model predicts four outputs: semantics prediction from the semantic head, class, bounding box, and mask prediction from the instance head. All the aforementioned predictions are then fused in the panoptic fusion module to yield the final panoptic segmentation output. Moreover, advances in the state-of-the-art deep learning methods continually boost the performance of these tasks to new heights.

3.4 Volumetric Semantic 3D Reconstruction

With the above programs, depth maps and pixel-wise semantic classification scores are achieved as inputs to the final objective, the semantic understanding of 3D environments. The core processing will be carried by the volumetric semantic reconstruction, which is cast as a volumetric fusion of depth maps and pixel-wise semantic classification scores. In practical applications, 3D reconstruction systems or semantic segmentation algorithms are not robust enough and often lead to challenging results given surfaces observed under very certain viewing angles. Many of these limitations under such fusion processes can be overcome by casting dense 3D reconstruction and semantic segmentation as a joint optimization formulation, shown in **Fig. 3**. The general idea of the formulation is that each of the voxels gets assigned one out of $L + 1$ labels where label $i = 0$ denotes the free space label and the L



▲ **Figure 2. Overview of the overall architecture for the classical panoptic segmentation (pictures taken from Ref. [26])**



▲ Figure 3. Dense semantic 3D reconstruction^[9]

labels with $i > 0$ indicate the occupied space, which is segmented into several semantic classes. Such formulation, so-called objective function of the volumetric multi-label approaches, can be resolved with the objective function of the convex multi-label energy extended from the volumetric 3D reconstruction energy, as described in Eq. (6). The energy $E(x)$ consists of two parts, in which the former data term is a function of a given label, and is parameterized by the internal probability distribution of the voxel/surfel. The subsequent pairwise smoothness term is a function of the labeling of two connected voxels/surfels in the graph, and is parameterized by the geometry of the map.

$$E(x) = \sum_{s \in \Omega} \left(\sum_i \rho_s^i x_s^i + \sum_{i,j: i < j} \Phi_s^{ij} (x_s^i - x_s^j) \right), \quad (6)$$

where E represents the objective function of the convex multi-label energy, X_s^i represents the label assigned to voxel s , ρ_s^i represents a cost for assigning label i to voxel s , and Φ_s^i represents transition-specific, direction-and-location-dependent penalizer of the surface area formed as an interface between labels i and j .

This type of formulation describes a convex relaxation procedure, which is closely related to linear programming (LP) relaxations for approximate maximum a posteriori (MAP) estimation inference in Markov random fields (MRFs). The classical solutions to addressing this procedure include the Bayesian, conditional random field (CRF), MRF and variation framework. The work of HÄNE et al.^[9] can be referred to for thorough exploration regarding such formations and approaches. HAN et al.^[15] also address some latest emerging technique problems, inspired by the continually boosted deep learning achievement.

4 Neural Metric-Semantic Understanding

Following the above overview of the underlying computer vision and machine learning concepts, we will discuss the new explorations regarding fully leveraging the joint semantic & 3D information, neural metric-semantic understanding. Given the high-quality geometric and semantic scene understanding specification, classic semantic and spatial collaborative perception methods can reconstruct global 3D semantic dense maps for a variety of real-world scenes. Moreover, such dense 3D semantic mapping techniques give us explicit editing control over all the elements of the perception pipeline, and strictly

follow physical knowledge from computer vision—camera viewpoint, lighting, geometry and materials. However, building high-quality semantic & 3D reconstruction, especially directly from poorly textured areas, under a higher noise level, in dynamic surrounding environments, requires significant manual effort, and automated high consistent context understanding from images is an open research problem. On the other hand, the emerging learning-based techniques are now starting to produce a plausible dense depth map or even 3D scanning of scenes, which is either from random noise or conditioned on certain user specifications. However, they do not yet allow geometrical consistency and cannot always handle the true depth by a single scale factor. In contrast, neural metric-semantic understanding brings the promise of combining these approaches to enable high quality co-consistency under both semantic and geometric scenarios. Neural metric-semantic understanding techniques are diverse, differing in the fusion that they provide over the perception pipeline, the type of fusion and the network structures they utilize. A typical neural metric-semantic understanding approach takes red-green-blue depth (RGBD) sequences corresponding to certain scenes as input, builds a dense 3D reconstruction from them, and adopts the volumetric 3D convolution for point cloud segmentation to extract the final semantic 3D understanding. The dense 3D reconstruction is not restricted by directly using classical computer vision methods to geometric modeling of the environment and can be optimized with the combination of differentiable machine learning techniques for high quality consistent understanding. At the same time, neural metric-semantic understanding approaches incorporate ideas from classical computer vision in the form of orthogonal approaches to reduce drift, traditionally-obtained geometric constraints, and network architectures—to make the learning task easier and the output more consistent.

We propose a taxonomy of neural metric-semantic understanding approaches along the axes that we consider the most important:

- Joint volumetric multi-label formulation
- Semantically geometric and photometric registration
- Semantical depth map regulation

In the following, we will discuss current state-of-the-art methods under these axes.

4.1 Neural Joint Volumetric Multi-Label Formulation

According to the general pipeline of metric-semantic understanding, depth maps and pixel-wise semantic classification scores are achieved as inputs of the final objective, the “semantic understanding of 3D environments”. Various approaches are proposed to tackle joint optimization formulation. Authors in Refs. [15, 27 – 28] directly use 3D convolutional neural networks approach on voxels (representation of 3D scenes), like 2D convolution on pixels, while the methods in Ref. [13], such as variational methods for convex relaxation, incorporate the physical knowledge to an emerging differentiable learning network.

3D convolutional neural network methods rely on generic 3D convolutional neural network architectures, and take the three-dimensional representation of 3D scenes as input. The curse of dimensionality applies, in particular, to data that lives on grids, which have three or more dimensions. The number of points on the grid grows exponentially with its dimensionality. In such scenarios, as the counterpart of 2D convolutional processing for two-dimensional pictures, it becomes increasingly important to reduce the computational resources needed for 3D data convolutional processing, such as exploiting sparsity and reduces the number of global memory accesses. Prior work in Ref. [28] implements sparse convolutions (SCs) and introduces a novel convolution operator termed submanifold sparse convolution (SSC) that restricts computation and storage to “active” sites. The utilization of the sparsity nature of points in the 3D volumetric space forms the basis for a new mainstream solution, submanifold sparse convolutional networks (SSCNs), which are optimized for efficient semantic segmentation of 3D representation of scenes. A later trial in Ref. [15] extends the SSCN with explorations in addressing the efficiency bottleneck of sparse 3D CNN, which lies in the unorganized memory access of the sparse convolution steps, for the demand of online computations.

Directly applying 3D convolutional neural networks to voxels like 2D convolution on pixels will introduce some limitations, such as the insufficient capacity of deep learning techniques to delineate visual objects. This, for instance, can result in non-sharp boundaries and blob-like shapes in semantic segmentation tasks. While in the classical perception pipeline, probabilistic graphical models have been developed as effective methods to enhance the accuracy of the above task, as illustrated in Section 3.4. To this end, com-

pared with the classic convex relaxation procedure which always requires regularizers with hand-designed priors, a new differentiable learning network method^[13] combines the advantages of classical variational approaches with recent advances in deep learning, and improves the inference/optimization formulation from hand-tuned and not-easy convergence to a simple, generic, and substantially more scalable way. A reason for the improvement is that previously employed priors are not rich enough to capture the complex relationships of our 3D world, while learning-based differentiable networks break through automatically in an end-to-end trainable model. Furthermore, such an explicitly reused concept of variational energy minimization has led to great advances when dealing with noise and missing information.

On a separate track to the progress of joint optimization with neural deep learning techniques, some novel frameworks in Ref. [29] aggregate inputs from the initial stage of the previous pipeline and the information of multiple 2D observations from different view angles, and straightly reconstruct the final 3D semantic results with full deep learning framework. Rather than using the above methods, projecting color data into a volumetric grid and operating solely in 3D, with end-to-end network architecture, directly extracting feature maps from associated RGB images and then mapping into the volumetric feature grid of a 3D network using a differentiable back projection layer can result in more sufficient details.

4.2 Neural Semantically Geometric and Photometric Registration

Despite of the full exploration of the joint optimization formulation with geometry and semantic map as the input, emerging neural network techniques have also tried to leverage the combination of differentiable machine learning techniques with physical knowledge from computer vision in the submodules of the perception pipeline, to enable the classic metric-semantic understanding performance in complex scenes. The seminal methods in Refs.[14] and [31] aim to address the underlying key challenges of such scenarios, namely globally consistent geometric and photometric registration, with some revolutionary thinking, such as fine-tuning the deep network by using the extracted geometric constraints and representing semantics as an invariant scene for medium-term continuous tracking of large scale 3D scanning.

Robust data association is a core problem of visual odometry and the cornerstone of large-scale geometry reconstructions. Currently, the state-of-the-art classic metric-semantic understanding methods use short-term tracking to obtain continuous frame-to-frame constraints, while long-term constraints are established using loop closures, as illustrated in Ref. [14]. Although these two approaches are orthogonal and greatly reduce drift by collaboration, invariant representation of scenes to viewpoint and illumination changes cannot always be guaranteed, because of the gap between action in-

terval spans. The author originally proposes using semantics for medium-term continuous tracking of points to improve the first drift correction strategy. The underlying intuition is that changes in viewpoint, scale, illumination, etc., only affect the low-level appearance of objects but not their semantic meaning. By readily integrating semantic reprojection errors into existing video odometry (VO) approaches and combining differentiable machine learning techniques with physical knowledge from computer vision, translational drift in fast or complex scenes has reduced significantly, as reported in the literature.

The reverse thinking of the above method, emerging as another optimizing direction of deep learning in computer vision, is reflected in the method proposed by LUO et al.^[31]. The method leverages a convolutional neural network trained for single-image depth estimation along with conventional structure-from-motion reconstruction to establish geometric constraints on pixels in the image sequence. The authors firstly train a single-image depth estimation network to synthesize plausible depth for general color images, and then fine-tune the network by using the extracted geometric constraints via traditional reconstruction methods at the test time. This novel formula, which combines the strengths of traditional techniques and learning-based techniques, addresses the geometrical consistency of the reconstruction over time even under a gentle amount of dynamic scene motion.

4.3 Neural Semantically Depth Map Regulation

As the basic input of the semantic understanding of 3D environments, input geometry and semantic maps, recorded by the overlapped views or “active” sensing, always suffer from inaccuracy and incompatible resolutions because of the different sensing schemes. Plenty of progress as shown in Refs. [30, 32 – 33] has been made to reduce the noise and boost geometric details, especially after consumer depth sensors coming into our daily lives, marked by the recent integration in the latest iPhone. In many classic metric-semantic understanding approaches, volumetric depth map “fusion” has become a standard method, which shows geometric details boosting with sparse depth and dense RGB information, based on truncated signed distance functions. Due to the disadvantages and the real-time requirement of related classic methods, neural-based novel depth map regulation approaches emerge in multiple ways for new heights of performance: 1) semantic information which enriches the scene representation and is incorporated into the fusion process; 2) leveraging the multi-frame fused geometry and the accompanying high-quality color image through a joint training strategy; 3) depth upsampling method which is tolerant to outlier factors (such as mismeasured depth points, flipping points, and disocclusion) and to spontaneously adapt to each scene by a self-learning framework in an online update manner.

Instead of explicitly combining the geometry and semantic

segmentation of 3D in the former, others follow that by including this notion of collaboration more implicitly. However, efficiently encoding and fusing “semantically” rich multi-scale features from an end-to-end trainable (differentiable) way is abnormally obvious. Furthermore, recently there has also been immense progress on learning-based methods that operate on single images. These methods result in the pleasing ability to synthesize plausible depth, in particular, in dynamic scenes as well as limitations of the sensing range. In order to construct fine-grained depth sensing, one of the seminal works by TULSIANI et al.^[3] specializes those object’s representation in scenes to some particular instances, signaling that both top-down and bottom-up cues influence the perception, and perfectly deform into shapes even slightly different from those in the training. Fig. 1 illustrates the pleasing semantic object reconstruction result, which reflects the impressive influence introduced by neural semantic depth map regulation.

5 Applications of Semantic and Spatial Collaborative Perception

Semantic and spatial collaborative perception has many important use cases including, but not limited to, relighting, novel view synthesis, as well as semantic AR contents manipulation. The following is a detailed discussion of various applications.

5.1 Relighting

Relighting is known as a procedure for the photo-realistically rendering of a scene under a novel illumination. It is a fundamental component for a number of media editing applications including AR and visual effects. The previously challenging settings like large-scale outdoor scene relighting can be addressed with the help of multi-view-based semantic and spatial collaborative perception. Relighting in the wild^[18] casts the problem as a multi-modal image synthesis problem, which takes a rendered deep buffer as input, containing depth and color channels, together with a semantic label (also known as an “appearance code”), and outputs realistic views of tourist landmarks under various lighting conditions, as shown in **Fig. 4**. Fig. 4a shows that the model is rendered into a deep buffer of depth, color and semantic labels, and Fig. 4b shows that a relighting method translates these buffers into realistic renderings under multiple appearances. The input views including depth and color channels are used to reconstruct the 3D geometry of the scene; the semantic labels are also taken as the input to indicate the location of transient objects like pedestrians. Using the above corresponding rendered deep buffers and pairs of real photos, a multi-modal image synthesis pipeline learns an implicit model of appearance, which represents the time of the day, weather conditions and other properties not presented in the 3D model. A similar principle is also adopted by the multi-view relighting method^[34]. Furthermore, the author considers

that such geometry is coarse and erroneous, and directly relighting it would produce poor results. Instead, the geometry is used to construct intermediate buffers—normals, reflection features, and RGB shadow maps—as auxiliary inputs to guide a neural network-based relighting method. The above methods all generalize real scenes, producing high-quality results for applications like the creation of time-lapse effects from multi-

ple images.

5.2 Novel View Synthesis

Rendering of a scene under novel camera perspectives of the scene with a fixed set of images given—a procedure known as “novel view synthesis” or “free viewpoint videos”—is a critical component of the emerging media entertainment applications, 360 VR. The topic has gained a lot of interest in the research community and reached compelling quality results with the work of COLLET et al.^[35] and its real-time counterpart by DOU et al.^[36–37]. Key challenges of such applications are inferring the scene’s 3D structure through given sparse observations, for example, the painting of unseen parts of the scene. Recently, reconstructing a learned representation of the scene from the observations, and learning of priors on geometry, appearance and other scene properties in learned feature space with a differentiable renderer, has become a hot topic and made significant progress in previously open challenges such as learning from extremely sparse observations, as shown in **Fig. 5**. Such semantic and spatial collaborative perception-based approaches range from explicit 3D disentanglement of multi-plane images^[38] to proposing 3D-structured representations such as voxel grids of features in Refs. [16] and [17]. Among them, HoloGAN^[16] implements an explicit affine transformation layer that directly applies view manipulations to learn 3D features to build an unconditional generative model that allows explicit viewpoint changes. Scene representation networks (SRNs)^[17] encode both scene geometry and appearance in a single fully connected neural network, to parameterize surface geometry via an implicit function. Although such approaches show better results compared with previous ones, they still have limitations, i.e., they are restricted to a specific use case and limited by the training data.

5.3 Semantic AR Contents Manipulation

Semantic AR contents manipulation is, but not only, the key procedure of the emerging AR experience paradigm, the so-called “retargetable AR”^[19]. As the authors illustrate, re-



▲ Figure 4. Relighting in the wild^[18] reconstructs a proxy 3D model from a large-scale Internet photo collection



▲ Figure 5. Scene representation networks^[17] allow full 3D reconstruction from a single image (bottom row, surface normals and color render) by learning strong priors via a continuous, 3D-structure-aware neural scene representation

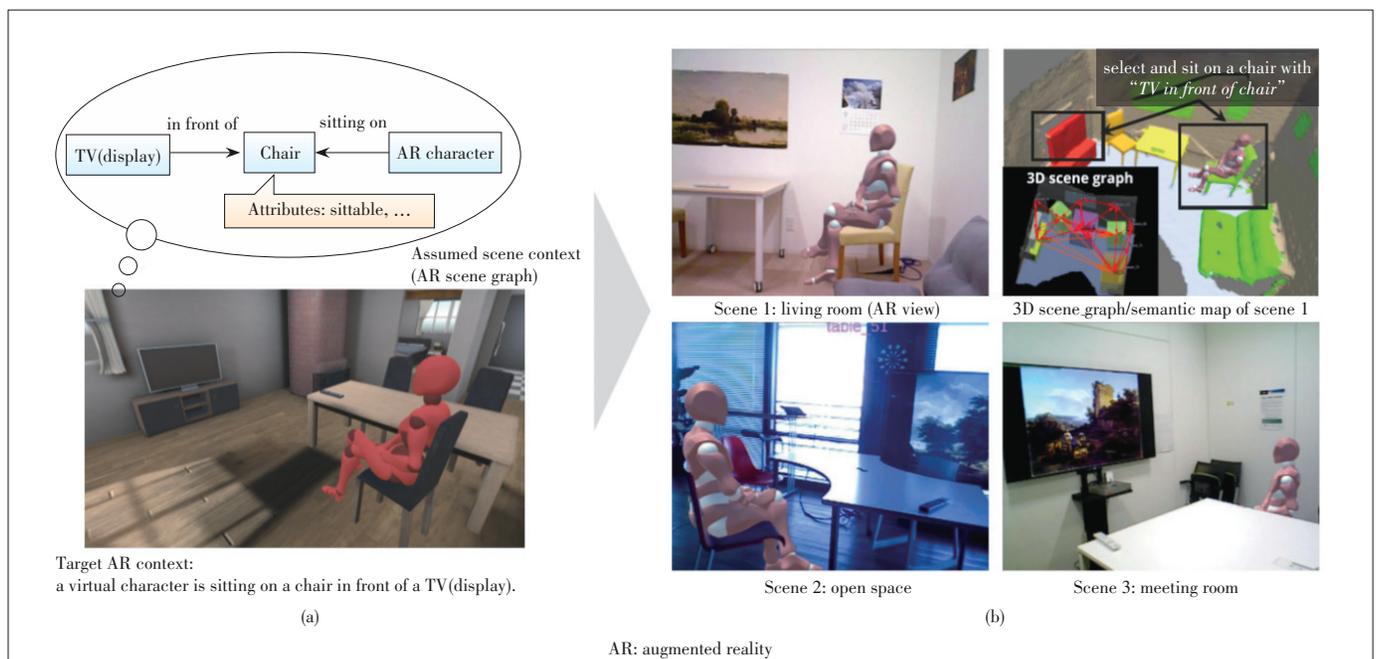
targetable AR is a novel AR framework that yields an AR experience that is aware of scene contexts set in various real environments, achieving natural interaction between the virtual and real world, as shown in Fig. 6, in which images are taken from Ref. [19]. It is expressed as an abstract AR scene graph based on the relationships among objects. Such a retargetable correspondence, which is between the realistic scene and the constructed graph, provides a semantically registered content arrangement, and finally facilitates natural interaction between “digital twins” and the real world. The key procedure, semantic AR contents manipulation, is an extension of the original solution, only a geometric and photometric registration between the virtual and real world^[39], to the integration of virtual objects into real environments accurately and naturally. It is achieved by the integration of the advanced abstraction (3D scene graph), and the accurately underlying semantic and spatial collaborative perception, which is the fusion of geometric and semantic information densely reconstructed and labeled in the scene. A similar idea is also proposed by ROSINOL et al.^[21], stating that the ideal level of abstraction will be more practical and crucial for the later augmented reality/mixed reality (AR/MR) systems. Even more, linked by such mechanism, the massive knowledge map combined with natural language expressions, and also the above deep understanding of physical environments can be collaboratively learned and managed.

6 Technical Implications

In the above sections, we present a multitude of applica-

tions with various target domains by semantic and spatial collaborative perception. While some applications are mostly insensitive to the processing time and response time, others, with legitimate and extremely useful use cases, should be used in an instant reaction manner (e.g., semantic AR contents manipulation). Methods for image and video manipulation are as old as the media themselves, and understanding-based structured visual editing is currently common, for example, in the Internet industry. Neural metric-semantic understanding approaches have the potential to lower the barrier for entry, making manipulation technology accessible to non-experts with limited resources. Although we believe that all the methods discussed in this paper have the potential to positively influence the world via better content creation and storytelling, we must not be complacent. It is important to proactively discuss and devise a plan to systematically arrange the submodules of the above methods under the 5G edge computing scenario for instant reaction and also privacy protection purpose. We believe it is critical that understanding-based synthesizing images and videos are extremely resource- and power-consuming. We also believe that it is essential to raise significant privacy concerns before directly uploading visual raw data to cloud-based semantic and spatial collaborative perception systems, like, for the localization purpose, even if only derived image features are uploaded.

Such related topics regarding “to cloud or not to cloud” were first explored by NAQVI et al.^[40], and then extended to edge computing architectures, even with 5G, by BARESI et al.^[41-43]. Given the evaluation regarding the added value of cloud computing as a key enabler for AR applications on mo-



▲ Figure 6. Illustration of semantic AR contents manipulation: (a) retargetable AR; (b) framework that retargets the AR scene to various real scenes by comparing the AR scene graph with 3D scene graphs constructed in each of the scenes^[19]

mobile devices^[40], the authors disclose an important principle that the latency due to connectivity type and the amount of data to be communicated is a major trade-off, and the dynamic deployment and reconfiguration of the framework components between mobile and cloud ends are really important. Furthermore, with respect to the final quality of experience requirements, context-awareness based resource allocation at the wireless network edge^[40, 42] and the adoption of serverless edge computing architecture^[41 - 43] become the consensus. With the deployment of services to the cloud, the initially widely ignored privacy concerns become an emerging key challenge. The possibility was strikingly demonstrated in Ref. [44], even when only the extracted features are uploaded.

The importance of developing corresponding safe disclosure technologies and building corresponding communities has risen to an urgent position. Such safeguarding measures would reduce the potential for misuse while allowing creative uses of semantic and spatial collaborative perception technologies. In one recent example in the field of image-based localization^[45], the authors adopted a cloud-based “obfuscate upload” approach, refraining from uploading the full 3D points of structure-from-motion maps immediately, instead of uploading random line features, lifted from 2D/3D feature points.

Learning from this example, we believe researchers and related business operators must make privacy preserving strategies a key part of all the edge-based semantic and spatial collaborative perception systems with a potential for misuse, but not an afterthought. Also, it is important that we, as a community, continue to develop responsible neural metric-semantic understanding techniques to enable cloud-based semantic and spatial collaborative perception solutions without sacrificing the privacy of users by hiding the privacy concerning contents of the uploading media information.

7 Conclusions

Neural metric-semantic understanding and also the newly neural extension have raised a lot of interest in the past few years. This paper investigates the linkage between the classical and concurrent explorations and a variety of directions related to the topic, which reflects the immense increase of research in this field. The target application is not bound to a specific one but spans a variety of use cases that range from novel view synthesis, relighting, to the manipulation of semantic contents for AR. We believe that metric-semantic understanding will have a profound impact on making complex structured vision understanding and editing tasks accessible to a much broader audience. We hope that this article, which especially focuses on neural metric-semantic understanding, can introduce such modern perception capability to a large research community, which in turn will help to develop the next generation of computer vision applications under the direction.

References

- [1] HERMANS A, FLOROS G, LEIBE B. Dense 3D semantic mapping of indoor scenes from RGB-D images [C]//2014 IEEE International Conference on Robotics and Automation (ICRA). Hong Kong, China: IEEE, 2014: 2631 - 2638. DOI: 10.1109/ICRA.2014.6907236
- [2] ROSINOL A, ABATE M, CHANG Y, et al. Kimera: an open-source library for real-time metric-semantic localization and mapping [C]//2020 IEEE International Conference on Robotics and Automation (ICRA). Paris, France: IEEE, 2020: 1689 - 1696. DOI: 10.1109/ICRA40945.2020.9196885
- [3] TULSIANI S, KAR A, CARREIRA J, et al. Learning category-specific deformable 3D models for object reconstruction [J]. IEEE transactions on pattern analysis and machine intelligence, 2017, 39(4): 719 - 731. DOI: 10.1109/TPAMI.2016.2574713
- [4] TATENO K, TOMBARI F, NAVAB N. Real-time and scalable incremental segmentation on dense SLAM [C]//2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). Hamburg, Germany: IEEE, 2015: 4465 - 4472. DOI: 10.1109/IROS.2015.7354011
- [5] MCCORMAC J, HANDA A, DAVISON A, et al. SemanticFusion: dense 3D semantic mapping with convolutional neural networks [C]//2017 IEEE International Conference on Robotics and Automation (ICRA). Singapore, Singapore: IEEE, 2017: 4628 - 4635. DOI: 10.1109/ICRA.2017.7989538
- [6] NAKAJIMA Y, TATENO K, TOMBARI F, et al. Fast and accurate semantic mapping through geometric-based incremental segmentation [C]//2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). Madrid, Spain: IEEE, 2018: 385 - 392. DOI: 10.1109/IROS.2018.8593993
- [7] NARITA G, SENO T, ISHIKAWA T, et al. PanopticFusion: online volumetric semantic mapping at the level of stuff and things [C]//2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). Macao, China: IEEE, 2019: 4205 - 4212. DOI: 10.1109/IROS40897.2019.8967890
- [8] PHAM Q H, HUA B S, NGUYEN T, et al. Real-time progressive 3D semantic segmentation for indoor scenes [C]//2019 IEEE Winter Conference on Applications of Computer Vision (WACV). Waikoloa Village, USA: IEEE, 2019: 1089 - 1098. DOI: 10.1109/WACV.2019.00121
- [9] HÄNE C, POLLEFEYS M. An overview of recent progress in volumetric semantic 3D reconstruction [C]//2016 23rd International Conference on Pattern Recognition (ICPR). Cancun, Mexico: IEEE, 2016: 3294 - 3307. DOI: 10.1109/ICPR.2016.7900143
- [10] LADICKÝ L, ZEISL B, POLLEFEYS M. Discriminatively trained dense surface normal estimation [C]//European Conference on Computer vision. Zurich, Switzerland: ECCV, 2014: 0906 - 0912. DOI: 10.1007/978-3-319-10602-1_31
- [11] GÜNEY F, GEIGER A. Displets: Resolving stereo ambiguities using object knowledge [C]//2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Boston, USA: IEEE, 2015: 4165 - 4175. DOI: 10.1109/CVPR.2015.7299044
- [12] LI R H, GU D B, LIU Q, et al. Semantic scene mapping with spatio-temporal deep neural network for robotic applications [J]. Cognitive computation, 2018, 10(2): 260 - 271. DOI: 10.1007/s12559-017-9526-9
- [13] CHERABIER I, SCHÖNBERGER J L, OSWALD M R, et al. Learning priors for semantic 3D reconstruction [M]//European Conference on Computer vision. Munich, Germany: ECCV, 2018: 325 - 341. DOI: 10.1007/978-3-030-01258-8_20
- [14] LIANOS K N, SCHÖNBERGER J L, POLLEFEYS M, et al. VSO: visual semantic odometry [M]//Computer Vision - ECCV 2018. Cham, Switzerland: Springer International Publishing, 2018: 246 - 263. DOI: 10.1007/978-3-030-01225-0_15
- [15] HAN L, ZHENG T, ZHU Y H, et al. Live semantic 3D perception for immersive augmented reality [J]. IEEE transactions on visualization and computer graphics, 2020, 26(5): 2012 - 2022. DOI: 10.1109/TVCG.2020.2973477
- [16] NGUYEN-PHUOC T, LI C, THEIS L, et al. HoloGAN: unsupervised learning of 3D representations from natural images [C]//2019 IEEE/CVF International Conference on Computer Vision (ICCV). Seoul, South Korea: IEEE, 2019: 7587 - 7596. DOI: 10.1109/ICCV.2019.00768
- [17] SITZMANN V, ZOLLHÖFER M, WETZSTEIN G. Scene representation networks: continuous 3D-structure-aware neural scene representations [EB/OL]. [2021-01-05]. <https://arxiv.org/abs/1906.01618>
- [18] MESHRY M, GOLDMAN D B, KHAMIS S, et al. Neural rerendering in the wild [C]//2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Long Beach, CA: IEEE, 2019: 1000 - 1009. DOI: 10.1109/CVPR42620.2019.00100

- tion (CVPR). Long Beach, CA, USA: IEEE, 2019: 6871 – 6880. DOI: 10.1109/CVPR.2019.00704
- [19] TAHARA T, SENO T, NARITA G, et al. Retargetable AR: context-aware augmented reality in indoor scenes based on 3D scene graph [EB/OL]. (2020-08-18) [2021-01-05]. <https://arxiv.org/abs/2008.07817>
- [20] SCHARSTEIN D, SZELISKI R, ZABIH R. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms [C]/Proceedings IEEE Workshop on Stereo and Multi-Baseline Vision (SMBV 2001). Kauai, HI, USA: IEEE, 2001: 131 – 140. DOI: 10.1109/SMBV.2001.988771
- [21] NAIR R, RUHL K, LENZEN F, et al. A Survey on time-of-flight stereo fusion [J]. Time-of-flight and depth imaging. sensors, algorithms, and applications, 2013, 8200:105 – 127. DOI: 10.1007/978-3-642-44964-2_6
- [22] DAI A, NIEBNER M, ZOLLHÖFER M, et al. BundleFusion [J]. ACM transactions on graphics, 2017, 36(4): 1. DOI: 10.1145/3072959.3126814
- [23] WHELAN T, SALAS-MORENO R F, GLOCKER B, et al. ElasticFusion: Real-time dense SLAM and light source estimation [J]. The international journal of robotics research, 2016, 35(14): 1697 – 1716. DOI: 10.1177/0278364916669237
- [24] HAN L, FANG L. FlashFusion: real-time globally consistent dense 3D reconstruction using CPU computing [C]/Robotics: Science and Systems XIV. Robotics: Science and Systems Foundation, 2018. DOI: 10.15607/rss.2018.xiv.006
- [25] DE GEUS D, MELETIS P, DUBBELMAN G. Fast panoptic segmentation network [J]. IEEE robotics and automation letters, 2020, 5(2): 1742 – 1749. DOI: 10.1109/LRA.2020.2969919
- [26] MOHAN R, VALADA A. EfficientPS: efficient panoptic segmentation [EB/OL]. (2020-05-19) [2021-01-05] <https://arxiv.org/abs/2004.02307>
- [27] ARMENI I, SENER O, ZAMIR A R, et al. 3D semantic parsing of large-scale indoor spaces [C]/2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Las Vegas, USA: IEEE, 2016: 1534 – 1543. DOI: 10.1109/CVPR.2016.170
- [28] GRAHAM B, ENGELCKE M, MAATEN L V D. 3D semantic segmentation with submanifold sparse convolutional networks [C]/2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. Salt Lake City, UT, USA: IEEE, 2018: 9224 – 9232. DOI: 10.1109/CVPR.2018.00961
- [29] DAI A, NIENER M. 3DMV: joint 3D-multi-view prediction for 3D semantic scene segmentation [C]/Computer vision. Munich, Germany: ECCV, 2018: 0908 – 0914. DOI: 10.1007/978-3-030-01249-6_28
- [30] ROZUMNYI D, CHERABIER I, POLLEFEYS M, et al. Learned semantic multi-sensor depth map fusion [C]/2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW). Seoul, South Korea: IEEE, 2019: 2089 – 2099. DOI: 10.1109/ICCVW.2019.00264
- [31] LUO X, HUANG J B, SZELISKI R, et al. Consistent video depth estimation [J]. ACM transactions on graphics, 2020, 39(4): 1 – 13. DOI: 10.1145/3386569.3392377
- [32] SHIM I, OH T H, KWEON I. High-fidelity depth upsampling using the self-learning framework [J]. Sensors, 2018, 19(1): 81. DOI: 10.3390/s19010081
- [33] YAN S, WU C L, WANG L Z, et al. DDRNet: depth map denoising and refinement for consumer depth cameras using cascaded CNNs [C]/European Conference on Computer vision. Munich, Germany: ECCV, 2018. DOI: 10.1007/978-3-030-01249-6_10
- [34] PHILIP J, GHARBI M, ZHOU T H, et al. Multi-view relighting using a geometry-aware network [J]. ACM transactions on graphics, 2019, 38(4): 1 – 14. DOI: 10.1145/3306346.3323013
- [35] COLLET A, CHUANG M, SWEENEY P, et al. High-quality streamable free-viewpoint video [J]. ACM transactions on graphics, 2015, 34(4): 1 – 13. DOI: 10.1145/2766945
- [36] DOU M, KHAMIS S, DEGTAREV Y, et al. Fusion4D: Real-time performance capture of challenging scenes [J]. ACM transactions on graphics, 2016, 35(4): 1 – 13. DOI: 10.1145/2897824.2925969
- [37] DOU M S, DAVIDSON P, FANELLO S R, et al. Motion2Fusion [J]. ACM transactions on graphics, 2017, 36(6): 1 – 16. DOI: 10.1145/3130800.3130801
- [38] XU Z X, BI S, SUNKAVALLI K, et al. Deep view synthesis from sparse photometric images [J]. ACM transactions on graphics, 2019, 38(4): 1 – 13. DOI: 10.1145/3306346.3323007
- [39] KIM K, BILLINGHURST M, BRUDER G, et al. Revisiting trends in augmented reality research: a review of the 2nd decade of ISMAR (2008 – 2017) [J]. IEEE transactions on visualization and computer graphics, 2018, 24(11): 2947 – 2962. DOI: 10.1109/TVCG.2018.2868591
- [40] NAQVI N Z, MOENS K, RAMAKRISHNAN A, et al. To cloud or not to cloud: a context-aware deployment perspective of augmented reality mobile applications [C]/Proceedings of the 30th Annual ACM Symposium on Applied Computing. Salamanca Spain. New York, USA: ACM, 2015: 0413 – 0417. DOI: 10.1145/2695664.2695880
- [41] BARESI L, FILGUEIRA MENDONÇA D, GARRIGA M. Empowering low-latency applications through a serverless edge computing architecture [C]/Service-oriented and cloud computing. Oslo, Norway: ESOC, 2017: 0927 – 0929. DOI: 10.1007/978-3-319-67262-5_15
- [42] CHATZIELEFTHERIOU L E, IOSIFIDIS G, KOUTSOPOULOS I, et al. Towards resource-efficient wireless edge analytics for mobile augmented reality applications [C]/2018 15th International Symposium on Wireless Communication Systems (ISWCS). Lisbon, Portugal: IEEE, 2018: 1 – 5. DOI: 10.1109/ISWCS.2018.8491206
- [43] BARESI L, FILGUEIRA MENDONÇA D. Towards a serverless platform for edge computing [C]/2019 IEEE International Conference on Fog Computing (ICFC). Prague, Czech Republic: IEEE, 2019: 1 – 10. DOI: 10.1109/ICFC.2019.00008
- [44] PITTALUGA F, KOPPAL S J, KANG S B, et al. Revealing scenes by inverting structure from motion reconstructions [C]/2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Long Beach, CA, USA: IEEE, 2019: 145 – 154. DOI: 10.1109/CVPR.2019.00023
- [45] GEPPERT M, LARSSON V, SPECIALE P, et al. Privacy preserving structure-from-motion [C]/16th European Conference Computer Vision. Glasgow, United Kingdom: EVVC, 2020:0823 – 0828. DOI: 10.1007/978-3-030-58452-8_20

Biography

ZHU Fang (zhu.fang@zte.com.cn) received the B.Eng. degree in electrical engineering and the M.Sc. degree in information and system from Xi'an Jiaotong University, China and the Ph.D. degree in electronic engineering from Southeast University, China. He is currently the director of the technical committee in digital video and vision of ZTE Corporation, and also the deputy director in multimedia of State Key Laboratory of Mobile Network and Mobile Multimedia Technology. He is a senior member of IEEE, focusing on circuits and systems for video technology and smart vision. His research interests include core technology, cloud architecture and acceleration chipset for specific application of XR & Smart Vision based on mobile computing.



Integrating Coarse Granularity Part-Level Features with Supervised Global-Level Features for Person Re-Identification

Abstract: Person re-identification (Re-ID) has achieved great progress in recent years. However, person Re-ID methods are still suffering from body part missing and occlusion problems, which makes the learned representations less reliable. In this paper, we propose a robust coarse granularity part-level network (CGPN) for person Re-ID, which extracts robust regional features and integrates supervised global features for pedestrian images. CGPN gains two-fold benefit toward higher accuracy for person Re-ID. On one hand, CGPN learns to extract effective regional features for pedestrian images. On the other hand, compared with extracting global features directly by backbone network, CGPN learns to extract more accurate global features with a supervision strategy. The single model trained on three Re-ID datasets achieves state-of-the-art performances. Especially on CUHK03, the most challenging Re-ID dataset, we obtain a top result of Rank-1/mean average precision (mAP)=87.1%/83.6% without re-ranking.

Keywords: person Re-ID; supervision; coarse granularity

CAO Jiahao^{1,2}, MAO Xiaofei^{1,2},
LI Dongfang^{1,2}, ZHENG Qingfang^{1,2},
JIA Xia^{1,2}

(1. State Key Laboratory of Mobile Network and Mobile Multimedia Technology, Shenzhen 518057, China;
2. ZTE Corporation, Shenzhen 518057, China)

DOI: 10.12142/ZTECOM.202101009

<http://kns.cnki.net/kcms/detail/34.1294.TN.20210225.1146.002.html>, published online February 25, 2021

Manuscript received: 2020-12-05

Citation (IEEE Format): J.H. Cao, X. F. Mao, D. F. Li, et al., "Integrating coarse granularity part-level features with supervised global-level features for person re-identification," *ZTE Communications*, vol. 19, no. 1, pp. 72 - 81, Mar. 2021. doi: 10.12142/ZTECOM.202101009.

1 Introduction

Person re-identification (Re-ID) aims to retrieve a given person among all the gallery pedestrian images captured by different cameras. It is a challenging task to learn robust pedestrian feature representations as realistic scenarios are highly complicated with regards to the illumination, the background, and occlusion problems. In recent years, person Re-ID has achieved great progress^[1-7]. However, person Re-ID methods are still suffering from occluded or body part missing pedestrian images, where they fail to extract discriminative deep features for person Re-ID. Intuitively, the complexity of realistic scenarios increases the difficulty in making correct retrieval for person Re-ID^[8-10]. Therefore, existing person Re-ID methods usually decline a lot in performance when dealing with realistic person Re-ID dataset like CUHK03, which contains a lot of occluded or body part missing pedestrian images, as illustrated in **Fig. 1**.



▲ Figure 1. Pedestrian images in CUHK03-labeled dataset

As it is well known, part-based methods^[5-7] like multiple granularity network (MGN)^[5] are widely used in person Re-ID and have achieved promising performance. Generally, part-based methods learn to combine global features and discriminative regional features for person Re-ID. The global features in part-based methods are usually extracted directly from the whole person image by the backbone network, while the regional features are generated by directly partitioning feature maps of the whole body into a fixed number of parts. Nevertheless, the overall performance of such part-based methods seriously depends on that all person images are well-bounded holistic person images with few occlusion or body part missing. As real-world scenarios are complicated, the bounding boxes detected by the detection algorithm may not be accurate enough, which usually leads to occluded or body part missing pedestrian images as Fig.1 shows. In Fig.1, ID-A means the ID of the pedestrian is A. We can see that for the same person in the realistic scenario, occluded and body-part missing pedestrian images are both captured as people are moving around the cameras. When dealing with such occluded or body part missing pedestrian images, the global features extracted from the whole image directly by the backbone network become less accurate; moreover, the regional features generated by directly partitioning feature maps of the whole body may focus on occluded parts and become ineffective, which impair the person Re-ID accuracy evidently.

To address the above problems, in this paper, we propose the coarse granularity part-level network (CGPN) for person Re-ID model that learns discriminative and diverse feature representations without using any third models. Our CGPN model can be trained end-to-end and performs well on three person Re-ID datasets. Especially on CUHK03, which contains a lot of occluded or body part missing pedestrian images, our method achieves state-of-the-art performances and outperforms the current best method by a large margin. CGPN has three branches, and each branch consists of a global part and a local part. The global part is supervised to learn more accurate global features by part-level body regions. With the supervision strategy, the global part can learn more proper global features for occluded or body part missing pedestrian images. For the local part, as pedestrian images detected in realistic scenarios are often occluded or body-part missing, too many fine grained local features generated by partitioning the whole body feature maps may decrease model performance. Therefore we propose a coarse grained part-level feature strategy that can extract effective regional features and perform better on the three person Re-ID datasets.

CGPN gains two-fold benefit toward higher accuracy for person Re-ID. Firstly, compared with extracting global features directly by backbone network, CGPN learns to extract more accurate global features with the supervision strategy. Secondly, with the coarse grained part-level feature strategy, CGPN is

capable of extracting effective body part features as regional features for person Re-ID. Besides, our method is completely an end-to-end learning process, which is easy for learning and implementation. Experimental results confirm that our method achieves state-of-the-art performances on several mainstream Re-ID datasets, especially on CUHK03, the most challenging dataset for person Re-ID, in single query mode, and we obtain a top result of Rank-1/mean average precision (mAP)=87.1%/83.6% without re-ranking.

The main contributions of our work are summarized as follows:

- We propose a novel framework named CGPN, which effectively integrates coarse grained part-level features and supervised global-level features and is more robust for person Re-ID.
- We develop the coarse grained part-level feature strategy for person Re-ID.
- We prove that the integration model of coarse grained part-level features and supervised global-level features achieves state-of-the-art results on three Re-ID datasets, especially on the CUHK03 dataset, in which our model outperforms the current best method by a large margin.

2 Related Works

2.1 Part-Based Re-ID Model

As deep learning is widely used in person Re-ID nowadays, most existing methods^[11-12] choose to extract feature maps by directly applying a deep convolution network such as ResNet^[13]. However, the single global feature extracted from the whole person image by a deep convolution network does not perform as well as expected. The reason is that person images captured by cameras usually contain random background information and are often occluded or body part missing, which impairs the performance a lot. Then part-based methods are proposed to get additional useful local information from person images for person Re-ID. As an effective way to extract local features, part-based methods^[5-7,14] usually benefit from person structure and together with global features, push the performance of person Re-ID to a new level. The common solution of part-based methods is to split the feature maps horizontally into several parts according to human body structure and concatenate the feature maps of each part. However, when dealing with occluded or body part missing pedestrian images, we find that part-based methods like MGN^[5], which has received state-of-the-art results on person Re-ID datasets, face the problem of performance decrease. Obviously, part-based methods are common solutions to holistic pedestrian images as they can get correct body parts by uniform partitioning, however, these methods are less effective to occluded or body part missing pedestrian images.

2.2 Attention-Based Re-ID Model

Recently, some attention-based methods try to address the

occlusion or body-part missing problems with the help of attention mechanisms. Attention module is developed to help extract more accurate features by locating the significant body parts and learning the discriminative features from these informative regions. LI et al.^[15] propose a part-aligning convolutional neural network (CNN) network for locating latent regions (hard attention) and then extract these regional features for Re-ID. ZHAO et al.^[16] employ the spatial transformer network^[17] as the hard attention model to find discriminative image parts. LI et al.^[18] use multiple spatial attention modules (by softmax function) to extract features at different spatial locations. XU et al.^[19] propose to mask the convolutional maps via a pose-guided attention module. LI et al.^[14] jointly learn multi-granularity attention selection and feature representation for optimizing person Re-ID in deep learning. However, most of the attention-based methods are often more prone to higher feature correlations, as these methods tend to have features focusing on a more compact subspace, which makes the extracted features attentive but less diverse, and therefore leads to sub-optimal matching performance.

2.3 Pose-Driven Re-ID Model

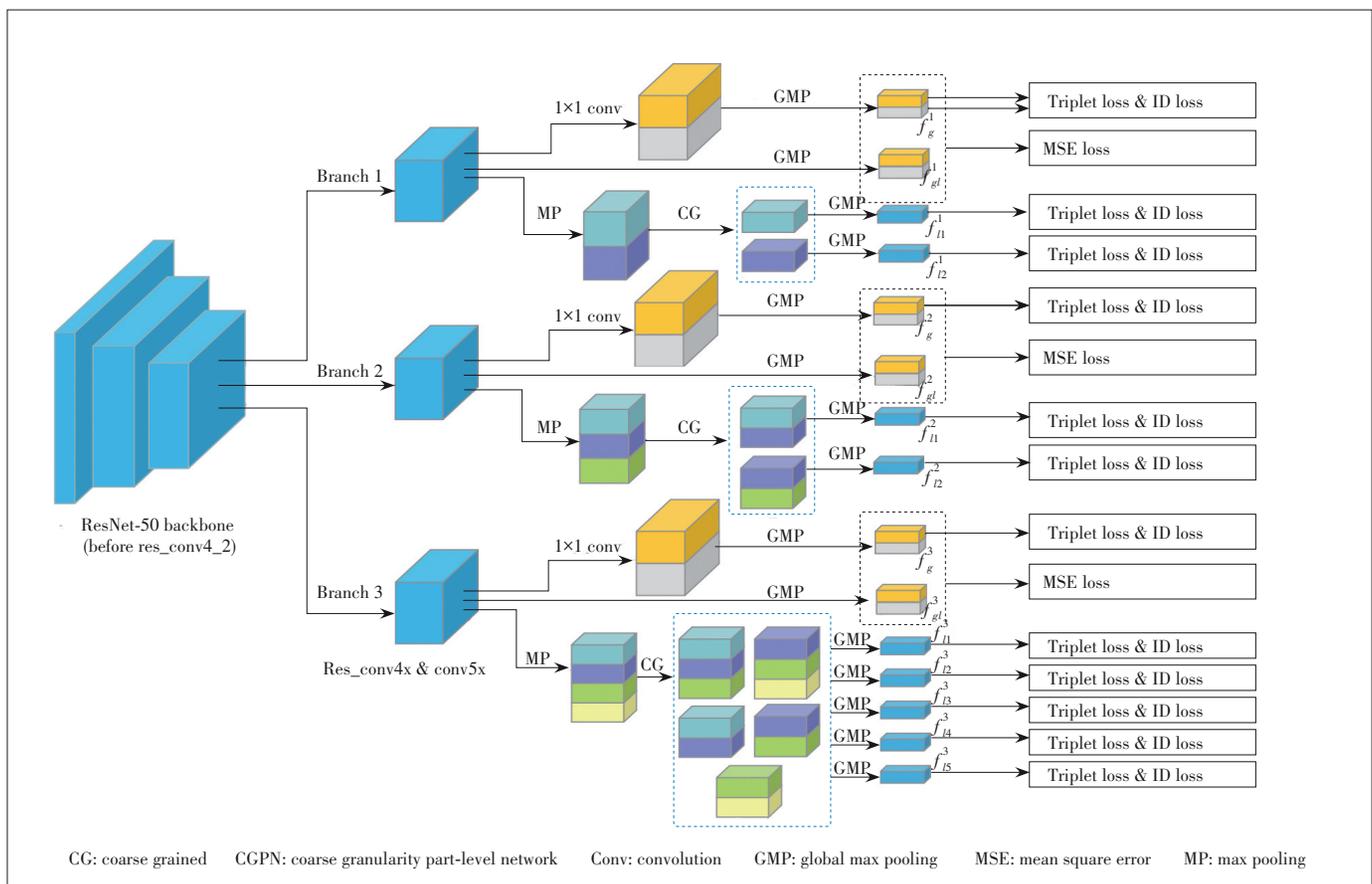
Some pose-driven methods utilize pose information to tackle the occlusion or body-part missing problems. In these meth-

ods, pose landmarks are introduced to help to align body parts as pose landmarks indicate the body position of persons. ZHENG et al.^[10] propose to use a CNN-based external pose estimator to normalize person images based on their pose, and the original and normalized images are then used to train a single deep Re-ID embedding. SARFRAZ et al.^[20] directly concatenate fourteen landmarks confidence maps with the image as network input, letting the model automatically learn alignment way. HUANG et al.^[21] propose a part aligned pooling (PAP) that utilizes seventeen human landmarks to enhance alignment. MIAO et al.^[22] learn to exploit pose landmarks to disentangle the useful information from the occlusion noise. However, the landmarks of persons are obtained usually by a third pose estimation model trained on an extra dataset, which increases the complicity of the whole Re-ID network. What's more, standard pose estimation datasets may not cover the drastic viewpoint variations in surveillance scenarios, and besides, surveillance images may not have sufficient resolution for stable landmarks prediction.

3 Proposed Method

3.1 Structure of CGPN

In this part, we present our CGPN structure in **Fig. 2**, in



▲ Figure 2. Structure of CGPN

which the ResNet-50 backbone is split into three branches after res_conv4_1 block. Each branch consists of a global part and a local part. In the global part, we apply two 1×1 convolutional layers and global max pooling (GMP) to generate global features. While in the local part, we apply a max pooling (MP) with different kernel sizes, split the feature maps into different spatial horizontal stripes, and then apply a coarse grained (CG) strategy and GMP to generate local features. The backbone of our network is a CNN structure, such as ResNet^[13], which achieves competitive results in many deep learning tasks. Like MGN^[5], we divide the output of res_conv4_1 into three different branches. Through the backbone network, CGPN transfers the input image into a 3D tensor T with size of $c \times h \times w$ (c is the channel number, h is the height, and w is the width). Each of the three branches contains a global part and a local part. The global part in three branches shares the same structure, while in every local part the output feature maps are uniformly partitioned into different stripes.

In the global part, two 1×1 convolution layers are applied to the output feature maps to extract regional features. Each of the 1×1 convolution layers will output c -channel features and be supervised by the corresponding part features. In further detail, for i -th branch's global part, to supervise the global features, the output feature maps are uniformly divided into two parts in the vertical direction, and a global pooling is applied to each of them to extract two part features $\{f_{g1}^i, f_{g2}^i\}$. The two part features $\{f_{g1}^i, f_{g2}^i\}$ are utilized in the training stage to supervise global features $\{f_{g1}^i, f_{g2}^i\}$ generated by the two 1×1 convolution layers. After the training stage finishes, the two part features are no longer needed. The first c -channel global features f_{g1}^i should be closer to the upper part features f_{g1}^i , and in the same way, the second c -channel global features f_{g2}^i should be closer to the bottom part features f_{g2}^i . In the inference stage, the first c -channel global features f_{g1}^i and the second c -channel global features f_{g2}^i are concatenated to form 2 c -channel features as final global features f_g^i . As the global parts of the three branches all share the same structure, we can get three global features $\{f_g^1, f_g^2, f_g^3\}$ in total. With the supervision of the part features, in the final 2 c -channel global features, the first c -channel global features are forced to focus on the upper part of the human body, while the second c -channel global features focus on the bottom part of the human body, which makes final global features more robust to person image occlusion or body-part missing.

For the local part in three branches, the output feature maps are divided into N stripes in the vertical direction with each stripe having the size of $c \times (h/N) \times w$, from which we prepare to extract local features. However, for person images that are occluded or body part missing, it might be harmful and decrease the performance of person Re-ID, if the granularity of local features is too fine. To alleviate the drawbacks of fine-grained local features, we choose

to extract local features in a bigger receptive field that contains enough body structure information to well represent the corresponding body region. In this paper, we propose a coarse-grained part-level feature strategy in which the combined part stripes must be adjacent and the minimum height proportion of combined local features should be no less than half of the output feature maps. The detail of the coarse grained strategy is illustrated in **Fig. 3**. In the local part of the first branch, the output feature maps are divided into two stripes in vertical direction as shown in Fig. 3a, and then pooling operations are performed to get local feature representations $\{f_{l1}^1, f_{l2}^1\}$ corresponding to the size of $c \times (h/2) \times w$. In the local part of the second branch, the output feature maps are divided into three stripes but we combine two adjacent stripes to get two $2/3$ proportion local features $\{f_{l1}^2, f_{l2}^2\}$ corresponding to the size of $c \times (2h/3) \times w$. For the local part of the third branch, the output feature maps are divided into four stripes, and then we combine two and three adjacent stripes to get three $1/2$ proportion and two $2/3$ proportion local features respectively, with $\{f_{l1}^3, f_{l2}^3, f_{l3}^3, f_{l4}^3, f_{l5}^3\}$ corresponding to the size of $c \times (3h/4) \times w, c \times (3h/4) \times w, c \times (3h/4) \times w, c \times (h/2) \times w, c \times (h/2) \times w$ respectively.

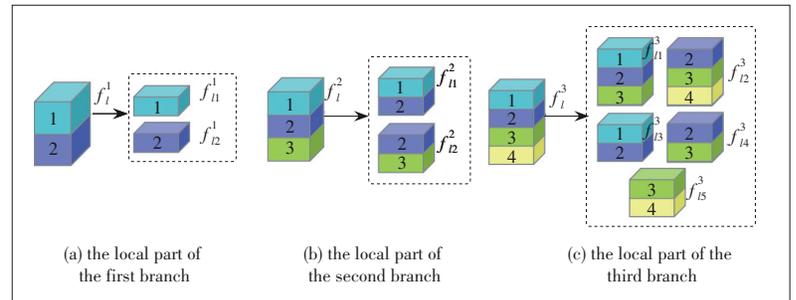
During the test, all features $\{f_g^1, f_g^2, f_g^3, f_{l1}^1, f_{l2}^1, f_{l1}^2, f_{l2}^2, f_{l1}^3, f_{l2}^3, f_{l3}^3, f_{l4}^3, f_{l5}^3\}$ generated by the global part and the local part in each branch are reduced to 256-dimension and are concatenated together as the final features, as different branches in CGPN actually learn representing information with different granularities which can cooperatively supplement discriminating information after the concatenation operation.

3.2 Loss Functions

Like various deep Re-ID methods, we employ softmax loss for classification, and triplet loss^[23] for metric learning. For the supervision of the global part in each branch, we use mean square error (MSE) loss in the training stage.

To be precise, in each branch, the local part is trained with the combination of softmax loss and triplet loss while the global part is trained with MSE loss, softmax loss and triplet loss as illustrated in Fig. 2.

For i -th learned features f_i , W_k is a weight vector for class k with the total class number C . N is the number of training examples in a mini-batch, and the softmax loss is formulated as:



▲ **Figure 3. Coarse-grained part-level feature strategy**

$$L_{\text{softmax}} = - \sum_{i=1}^N \log \frac{e^{W_{y_i}^T f_i}}{\sum_{k=1}^C e^{W_k^T f_i}}. \quad (1)$$

We employ the softmax loss to all global features $\{f_g^1, f_g^2, f_g^3\}$, and all coarse grained local features $\{f_{li}^1, f_{li}^2, f_{li}^3, f_{li}^5\}$.

Besides, all the global features $\{f_g^1, f_g^2, f_g^3\}$ are also trained with triplet loss. In the training stage, an improved batch hard triplet loss is applied with formula as follows:

$$L_{\text{triplet}} = \sum_{i=1}^P \sum_{a=1}^K \left[\alpha + \max_{p=1, \dots, K} \|f_a^{(i)} - f_p^{(i)}\|_2 - \min_{\substack{n=1, \dots, K \\ j=1, \dots, P \\ j \neq i}} \|f_a^{(i)} - f_n^{(j)}\|_2 \right]. \quad (2)$$

In the above formula, P is the number of selected identities and K is the number of images from each identity in a mini-batch. f_a^i is the anchor sample, f_p^i is the positive sample, f_n^i is the negative sample and α is the margin parameter to control the differences of intra- and inter-distances, which is set to 1.2 in our implementation.

To supervise the global features, we employ the MSE loss to all global features $\{f_g^1, f_g^2, f_g^3\}$ and the supervision local features $\{f_{gl}^1, f_{gl}^2, f_{gl}^3\}$ with the formula as follows:

$$L_{\text{mse}} = \sum_{i=1}^B \sum_{p=1}^M \|f_{gp}^i - f_{glp}^i\|_2^2, \quad (3)$$

where f_{gp}^i is the p -th c -channel features of global features in i -th branch, f_{glp}^i is the supervision p -th c -channel part features in the same branch. As the global part consists of two c -channel features and there are three branches in our network, M is set to 2 and B is set to 3 in our implementation.

The overall training loss is the sum of above three losses, which is formulated by:

$$L = L_{\text{softmax}} + L_{\text{triplet}} + L_{\text{mse}}. \quad (4)$$

4 Experiment

4.1 Datasets and Protocols

We train and test our model respectively on 4 mainstream Re-ID datasets: Market-1501^[24], DukeMTMC-reID^[25], CUHK03^[26] and Occluded-DukeMTMC^[22]. Especially the CUHK03 dataset, which is the most challenging realistic scenario Re-ID dataset as it consists of a lot of occluded or body part missing pedestrian images as illustrated in Fig. 1.

Market-1501 is captured by six cameras in front of a campus supermarket, which contains 1 501 person identities, 12 936 training images from 751 identities and 19 732 test-

ing images from 750 identities. All provided pedestrian bounding boxes are detected by deformable part models (DPM)^[27].

DukeMTMC-reID contains 1 812 person identities captured by 8 high-resolution cameras. There are 1 404 identities in more than two cameras and the other 408 identities are regarded as distractors. The training set consists of 16 552 images from 702 identities and the testing set contains 17 661 images from the rest 702 identities.

CUHK03 contains 1 467 person identities captured by six cameras on campus of The Chinese University of Hong Kong (CUHK). Both manually labeled pedestrian bounding boxes and automatically detected bounding boxes are provided. In this paper, we use the manually labeled version and follow the new training/testing protocol proposed in Ref. [28], with 7 368 images from 767 identities for training and 5 328 images from 700 identities for testing.

Occluded-DukeMTMC is re-segmented from the original DukeMTMC-reID dataset. The training set contains 15 618 images, and the gallery set and query set contain 17 661 and 2 210 images, respectively, in which all query images and some gallery images are occluded images, and these occluded images retain their occluded regions without being manually cropped.

In our experiment, we report the average cumulative match characteristic (CMC) at Rank-1, Rank-5, Rank-10 and mean average precision (mAP) on all the candidate datasets to evaluate our method.

4.2 Implementation Details

All images are re-sized into 384×128 px and the backbone network is ResNet-50^[13], pre-trained on ImageNet with the original fully connected layer discarded. In the training stage, the mini-batch size is set to 64, in which we randomly select 8 identities and 8 images for each identity ($P=8, K=8$). Besides, we deploy a randomly horizontal flipping strategy to images for data augmentation. Different branches in the network are all initialized with the same pre-trained weights of corresponding layers after res_conv4_1 block. Our model is implemented on Pytorch platform. We use stochastic gradient descent (SGD) as the optimizer with the default hyper-parameters (momentum=0.9, weight decay factor=0.0005) to minimize the network loss. The initial learning rate is set to 1e-2 and we decay it at epoch 60 and 80 to 1e-3 and 1e-4 respectively. The total training takes 240 epochs. During the evaluation, we use the average of original image features and horizontally flipped image features as the final features. All of our experiments on different datasets follow the settings above.

4.3 Comparison with State-of-the-Art Methods

In this section, we compare our proposed approach with current state-of-the-art methods on the three main-stream Re-ID datasets.

The statistical comparison between our PGCN network and

the state-of-the-art methods on Market-1501, DukeMTMC-reID and CUHK03 datasets is shown in **Table 1**.

On Market-1501 dataset, semantics aligning network (SAN) achieves the best published result without re-ranking, but our CGPN achieves 89.9% on the metric mAP, exceeding SAN by +1.9%. On the metric Rank-1, our CGPN achieves 96.1%, on a par with SAN, while our model is trained in an easier and end-to-end way. Compared with multiple granularity network (MGN) which is also a multiple branches method, our model surpasses MGN by +0.4% on the metric Rank-1 and by +3.0% on the metric mAP.

Among the performance comparisons on DukeMTMC-reID dataset, Pyramid achieved the best published result on metrics Rank-1 and mAP respectively. Our CGPN achieves the state-of-the-art result of Rank-1/mAP = 90.4%/80.9%, outperforming Pyramid by +1.4% on the metric Rank-1 and +1.9% on the metric mAP.

From Table 1, our CGPN model achieves Rank-1/mAP = 87.1%/83.6% on the most challenging CUHK03 labeled dataset under the new protocol. On the metric Rank-1, our CGPN outperforms the best published result of SAN by +7.0% and outperforms the best published result of Pyramid by +6.7% on mAP.

▼ **Table 1. Performance comparisons with the state-of-the-art results on Market-1501, DukeMTMC-reID and CUHK03 datasets in single query mode without re-ranking**

Method	Market-1501		DukeMTMC-reID		CUHK03	
	Rank-1/%	mAP/%	Rank-1/%	mAP/%	Rank-1/%	mAP/%
IDE ^[29]	-	-	-	-	22.2	21.0
PAN ^[30]	-	-	-	-	36.9	35.0
SVDNet ^[31]	-	-	-	-	40.9	37.8
IDE(R)+DM ^[32]	73.4	51.8	-	-	-	-
MGCAM ^[33]	83.8	74.3	-	-	50.1	50.2
DHA-Net + ISO(Aggr) ^[34]	88.2	70.1	74.2	54.5	-	-
HA-CNN ^[14]	91.2	75.7	80.5	63.8	44.4	41.0
VPM ^[35]	93.0	80.8	83.6	72.6	-	-
SCP ^[36]	94.1	-	84.8	-	-	-
PCB+RPP ^[6]	93.8	81.6	83.3	69.2	-	-
SphereReID ^[37]	94.4	83.6	83.9	68.5	-	-
MGN ^[5]	95.7	86.9	88.7	78.4	68.0	67.4
DSA ^[38]	95.7	87.6	86.2	74.3	78.9	75.2
Pyramid ^[7]	95.7	88.2	89.0	79.0	78.9	76.9
SAN ^[39]	96.1	88.0	87.9	75.5	80.1	76.4
CGPN	96.1	89.9	90.4	80.9	87.1	83.6

CGPN: coarse granularity part-level network
DHA: deep hidden attribute
DM: discrepancy matrix and matrix metric
DSA: densely Semantically Aligned
HA-CNN: harmonious attention convolutional neural network
IDE: ID-discriminative embedding
ISO: Identity-preserving, Sparsity constraints and the Orthogonal generation module
mAP: mean average precision

MGCAM: mask-guided contrastive attention model
MGN: multiple granularity network
PAN: pedestrian alignment network
PCB: part-based convolutional baseline
RPP: refined part pooling
SAN: semantics aligning network
SCP: spatial-channel parallelism
SVDNet: singular vector decomposition network
VPM: visibility-aware part model

In summary, our proposed CGPN can always outperform all other existing methods and shows strong robustness over different Re-ID datasets. According to the comparative experiments on the three datasets, especially on CUHK03 dataset, our approach can consistently outperform all other models by a large margin. Therefore, we can conclude that our method can effectively extract robust deep features from occluded or body part missing pedestrian images in person Re-ID.

We also conduct an experiment and compare the performances with the existing methods on Occluded-DukeMTMC. The results are listed in **Table 2**. As can be seen that, CGPN gets the top performance among the compared approaches, and obtains 58.5%/50.9% in rank-1/mAP. CGPN surpasses pose-guided feature alignment (PGFA) by +7.1% rank-1 accuracy and +13.6% mAP, which is a large margin. Therefore, we can conclude that our proposed CGPN integrated with supervised global-level features can effectively address the occlusion problem in person Re-ID.

4.4 Importance of Coarse-Grained Part-Level Features

To verify the effectiveness of coarse-grained part-level feature strategy in the CGPN model, we train two mal-functioned CGPN for comparison:

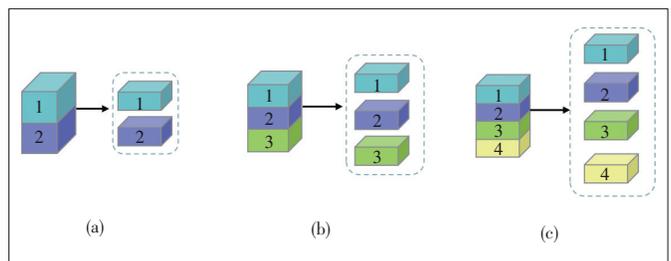
- CGPN-1 abandons the local parts in three branches and keeps only global parts.
- CGPN-2 replaces coarse-grained part-level features with fine-grained part-level features. It abandons coarse-grained strategy in local parts of three branches, compared with the normal CGPN model. Its local parts in three branches directly divide the output feature maps into two, three and four parts as shown in **Fig. 4**.

From the comparison of CGPN-1 with CGPN, we can see

▼ **Table 2. Performance comparisons with the state-of-the-art results on Occluded-DukeMTMC dataset in single query mode without re-ranking.**

Method	Occluded-DukeMTMC			
	Rank-1/%	Rank-5/%	Rank-10/%	mAP/%
HA-CNN ^[14]	34.4	51.9	59.4	26.0
PCB+RPP ^[6]	42.6	57.1	62.9	33.7
PGFA ^[22]	51.4	68.6	74.9	37.3
CGPN	58.5	73.4	78.4	50.9

CGPN: coarse granularity part-level person Re-ID network
HA-CNN: harmonious attention convolutional neural network
mAP: mean average precision
PCB: part-based convolutional baseline
PGFA: pose-guided feature alignment
RPP: refined part pooling



▲ **Figure 4. Fine grained local part structure in CGPN-2**

a significant performance decrease on Rank-1/mAP by -1.2% / -2.0% , -1.1% / -2.5% and -4.7% / -3.7% on Market-1501, DukeMTMC-reID and CUHK03 datasets respectively. Especially on CUHK03, we can observe a sharp decrease by -4.7% / -3.7% on the metric Rank-1/mAP. As CGPN-1 is trained in exactly the same procedure with the CGPN model and the CUHK03 dataset typically consists of many occluded or body part missing person images, we can infer that the coarse-grained local part is critical for CGPN model, especially on the dataset which contains a lot of occluded or body part missing person images.

Comparing CGPN-2 with CGPN, we can still observe a performance decrease by -0.8% / -0.5% , -0.1% / -0.7% and -1.8% / -1.1% on the metric Rank-1/mAP on Market-1501, DukeMTMC-reID and CUHK03 datasets respectively. Compared with fine-grained part-level features, coarse-grained part-level features contain enough body structure information to better represent the corresponding body regions, which makes CGPN learn more robust local features. Besides, on CUHK03, we can also see a sharper performance decrease compared with the other two datasets. The reason is that Market-1501 and DukeMTMC-reID consist of mainly holistic person images with little occlusions or body part missing, and these images keep complete body structure and make fine-grained part-level features achieve comparable performance with coarse-grained part-level features. While on CUHK03, as it consists of a lot of occluded or body part missing person images, coarse-grained part-level features outperform fine-grained part-level features evidently. Our experiments clearly prove that our coarse-grained part-level feature strategy can improve model performance significantly and is critical for model robustness, especially for occluded or body part missing person images.

4.5 Importance of Supervised Global Part

To further verify the effectiveness of the supervised global part in CGPN model, we train another two mal-functioned CGPN for comparison:

- CGPN-3 abandons global parts in all three branches and keeps only local parts that are trained with triplet loss and softmax loss.

▼ **Table 3. Ablation study of CGPN coarse grained part-level feature strategy and supervised global part, with comparison results on Market-1501, DukeMTMC-reID and CUHK03-labeled at evaluation metrics of Rank-1 and mAP in single query mode without re-ranking**

Method	Market-1501		DuckMTMC-reID		CUHK03	
	Rank-1/%	mAP/%	Rank-1/%	mAP/%	Rank-1/%	mAP/%
CGPN-1	94.9	87.9	89.3	78.4	82.4	79.9
CGPN-2	95.3	89.4	90.3	80.2	85.3	82.5
CGPN-3	94.2	86.2	88.6	76.9	84.3	81.0
CGPN-4	95.2	89.3	90.0	79.9	83.4	80.7
CGPN	96.1	89.9	90.4	80.9	87.1	83.6

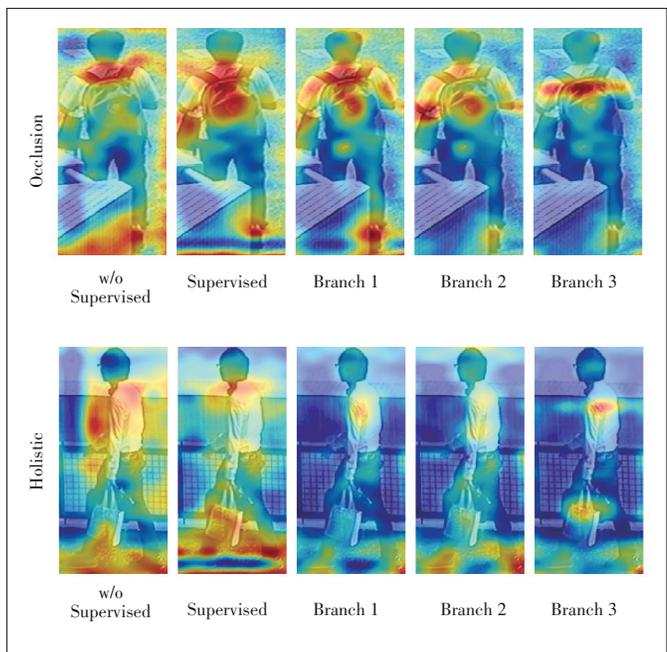
CGPN: coarse granularity part-level network mAP: mean average precision

- CGPN-4 keeps the global parts but abandons the supervision learning of all global parts in three branches, and these global parts are trained only with triplet loss and softmax loss.

Comparing CGPN-3 with CGPN, we observe a dramatic performance decrease on all three datasets. The performance on the metric Rank-1/mAP decreases by -1.9% / -3.7% , -1.8% / -4.0% and -2.8% / -2.6% on Market-1501, DukeMTMC-reID and CUHK03 respectively. As the three models are trained in exactly the same procedure, we conclude that the global part is critical to CGPN.

Comparing CGPN-4 with CGPN, after abandoning global supervision, we observe a performance decrease on Rank-1/mAP of -0.9% / -0.6% and -0.4% / -1.0% on Market-1501 and DukeMTMC-reID. While on CUHK03 we observe a dramatic performance decrease by -3.7% / -2.9% . The reason of such a different performance decrease is that Market-1501 and DukeMTMC-reID mainly consist of holistic person images from which the global part can get enough good global features directly even without supervision, while CUHK03 contains a lot of occluded or body part missing person images and the supervised global part is much more important for extracting accurate global features.

Comparing CGPN-4 with CGPN-3, after adding unsupervised global parts, we see a large performance improvement on Rank-1/mAP of $+1.0\%$ / $+3.1\%$ and $+1.4\%$ / $+3.0\%$ on Market-1501, DukeMTMC-reID. But on CUHK03 we observe a significant performance decrease by -0.9% / -0.3% unexpectedly. As analyzed above, the image type is quite different in the three datasets, especially CUHK03 which contains a lot of occluded or body part missing person images.



▲ **Figure 5. Feature visualization of different branches in the two cases of occluded pedestrian images and holistic pedestrian images**

The unexpected performance decrease on CUHK03 further proves that unsupervised global features can be harmful and certainly impair model performance. We conclude that the supervision of global features is critical for high performance of person Re-ID and that unsupervised global features will result in inaccurate global features which impair model performance evidently. As shown in **Fig. 5**, we can find that the unsupervised global features may receive interference from the background or occlusion. But, our proposed supervised global features are more robust to person image occlusion or body-part missing.

4.6 Branch Settings Ablation Study

In this section, we conduct a large number of comparative experiments on CUHK03 dataset to verify the effectiveness of the numbers of 1×1 convolution layers in the global part and multi-branch architecture settings.

In the global part, the number of 1×1 convolution layers is a hyper-parameter and influences the receptive field of its corresponding supervision part features. To evaluate the effect of various numbers of 1×1 convolution layers in the global part, as three branches' global part all share the same structure, we only keep branch 1 of CGPN and abandon the other two branches of CGPN. We also abandon the local part of branch 1 and only keep the global part of branch 1 denoted as Branch1-Global.

Table 4 shows the results of Branch1-Global with different numbers of 1×1 convolution layers, i.e. 2, 3, 4, 8. From these results, we can find that Branch1-Global reaches the best performance with two 1×1 convolution layers, which achieves Rank-1/mAP = 77.5%/74.7% on the CUHK03 dataset. The experiment results further illustrate that coarse grained features can make full use of local information and preserve more semantic information, and thus help to extract more accurate global features. Therefore, We finally adopt two 1×1 convolution layers in our CGPN architecture.

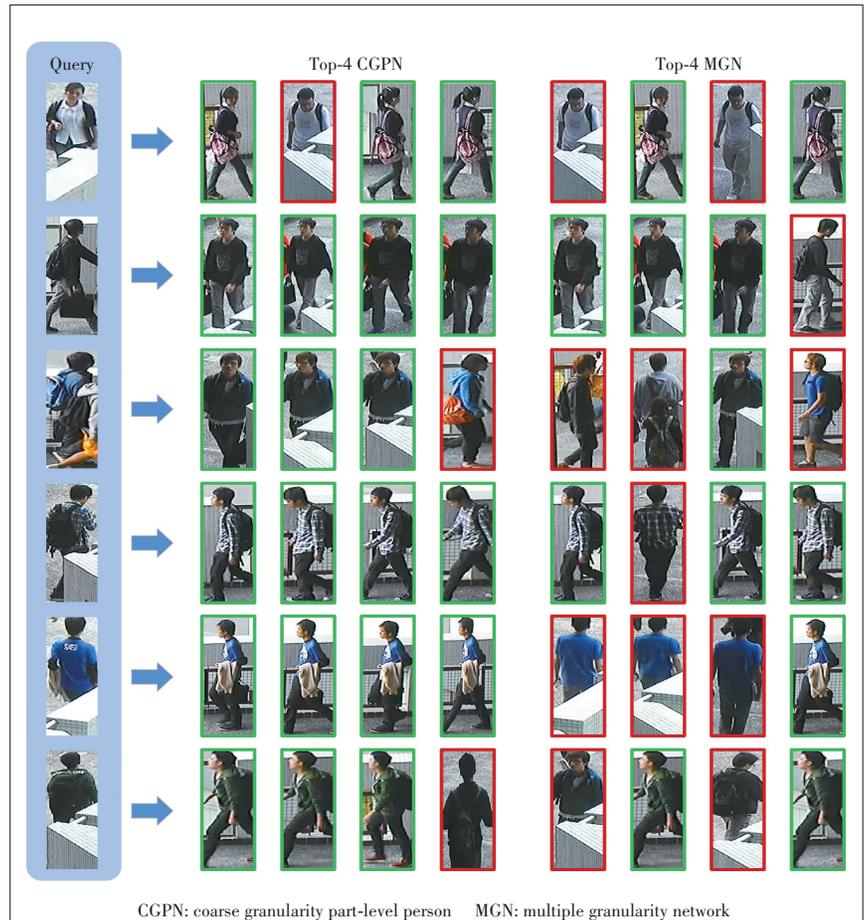
We further perform experiments to verify the importance of various branch settings in CGPN. Here, Branch x means we only keep CGPN's backbone and branch x after res_conv4_1. For example, Branch 1 means just preserving CGPN's backbone and branch 1 of CGPN and removing branches 2 and branch 3. With the increasing number of branches, Rank-1/mAP is significantly improved from 82.9%/79.4% to 85.4%/82.2% even to 87.1%/83.6%, as illustrated in Table 4. But, when we try more branches, such as CGPN + Branch 4 and CGPN + Branch 4 +

Branch 5, we observe a significant performance decrease on Rank-1/mAP of $-1.7\%/-1.2\%$ and $-1.3\%/-1.0\%$ unexpectedly. Therefore, we can conclude that the carefully-designed net-

▼ **Table 4. Comparison results of different number of 1×1 convolution layers in the global part and multi-branch settings on CUHK03 dataset at evaluation metrics of Rank-1 and mAP in single query mode without re-ranking**

Model	Rank-1/%	mAP/%
Branch1-Global w/2-part supervised	77.5	74.7
Branch1-Global w/3-part supervised	78.2	74.5
Branch1-Global w/4-part supervised	76.6	73.4
Branch1-Global w/8-part supervised	76.1	73.5
Branch 1	82.9	79.4
Branch 2	81.8	79.2
Branch 3	82.6	78.7
Branch 2 & Branch 3	84.5	82.1
Branch 1 & Branch 3	83.6	81.1
Branch 1 & Branch 2	85.4	82.2
CGPN + Branch 4	85.4	82.4
CGPN + Branch 4 + Branch 5	85.8	82.6
CGPN	87.1	83.6

CGPN: coarse granularity part-level network mAP: mean average precision



▲ **Figure 6. Top-4 ranking list for some query images on CUHK03-labeled dataset by CGPN and MGN**

work architecture is also the main contributor to performance improvement, and three branches can effectively and efficiently capture enough complement information.

Besides, from Fig. 5, we can find that the three branches of CGPN focus on different parts of the pedestrian, and the extracted features are complementary to each other.

4.7 Visualization of Re-ID Results

We visualize the retrieval results by CGPN and MGN for some given query pedestrian images of CUHK03-labeled dataset in Fig. 6, in which the retrieved images are all from the gallery set, but not from the same camera shot. The images with green borders belong to the same identity as the given query, and those with red borders do not. These retrieval results show the great robustness of our CGPN model, regardless of the occlusions or body part missing of detected pedestrian images. CGPN can robustly extract discriminative features for different identities.

5 Conclusions

In this paper, we propose a coarse-grained part-level features learning network integrated with supervised global-level features for person Re-ID. With the coarse-grained part-level strategy, the local parts in three branches learn more discriminative local features. With the supervision learning of global parts in three branches, the global parts learn to extract more accurate and suitable global features for pedestrian images. Experiments have confirmed that our model not only achieves state-of-the-art results on all three mainstream person Re-ID datasets, but pushes the performance to an exceptional level.

References

- [1] CHANG X B, HOSPEDALES T M, XIANG T. Multi-level factorisation net for person re-identification [EB/OL]. (2018-04-17) [2020-12-05]. <https://arxiv.org/abs/1803.09132>
- [2] LIU H, FENG J, QI M, et al. End-to-end comparative attention networks for person re-identification [J]. *IEEE transactions on image processing*, 2017, 26(7): 3492 – 3506. DOI: 10.1109/tip.2017.2700762
- [3] SARFRAZ M S, SCHUMANN A, EBERLE A, et al. A pose-sensitive embedding for person re-identification with expanded cross neighborhood re-ranking [C]//*IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, USA: IEEE, 2018: 420 – 429. DOI: 10.1109/cvpr.2018.00051
- [4] SHEN Y T, LI H S, XIAO T, et al. Deep group-shuffling random walk for person re-identification [C]//*IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, USA: IEEE, 2018: 2265 – 2274. DOI: 10.1109/CVPR.2018.00241
- [5] WANG G S, YUAN Y F, CHEN X, et al. Learning discriminative features with multiple granularities for person re-identification [C]//*Proceedings of the 26th ACM International Conference On Multimedia*. Seoul, Korea: ACM, 2018: 274 – 282. DOI: 10.1145/3240508.3240552
- [6] SUN Y F, ZHENG L, YANG Y, et al. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline) [C]//*Proceedings of the European Conference on Computer Vision*. Munich, German: ECCV, 2018: 480 – 496. DOI: 10.1007/978-3-030-01225-0_30
- [7] ZHENG F, DENG C, SUN X, et al. Pyramidal person re-identification via multi-loss dynamic training [C]//*Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Long Beach, USA: CVPR, 2019: 8514 – 8522. DOI: 10.1109/cvpr.2019.00871
- [8] SHEN Y, LIN W, YAN J, et al. Person re-identification with correspondence structure learning [C]//*Proceedings of the IEEE international conference on computer vision*. Santiago, Chile: IEEE, 2015: 3200 – 3208. DOI: 10.1109/iccv.2015.366
- [9] VARIOR R R, SHUAI B, LU J W, et al. A siamese long short-term memory architecture for human re-identification [C]//*European Conference on Computer Vision*. Amsterdam, Netherlands, ECCV, 2016: 135 – 153. DOI: 10.1007/978-3-319-46478-7_9
- [10] ZHENG L, HUANG Y J, LU H C, et al. Pose-invariant embedding for deep person re-identification [J]. *IEEE transactions on image processing*, 2019, 28(9): 4500 – 4509. DOI: 10.1109/tip.2019.2910414
- [11] LI W, ZHAO R, XIAO T, et al. DeepReID: deep filter pairing neural network for person re-identification [C]//*2014 IEEE Conference on Computer Vision and Pattern Recognition*. Columbus, USA: IEEE, 2014: 152 – 159. DOI: 10.1109/CVPR.2014.27
- [12] YI D, LEI Z, LIAO S C, et al. Deep metric learning for person re-identification [C]//*2014 22nd International Conference on Pattern Recognition*. Stockholm, Sweden: IEEE, 2014: 34 – 39. DOI: 10.1109/icpr.2014.16
- [13] HE K M, ZHANG X Y, REN S Q, et al. Deep residual learning for image recognition [C]//*2016 IEEE Conference on Computer Vision and Pattern Recognition*. Las Vegas, USA: CVPR, 2016: 770 – 778. DOI: 10.1109/cvpr.2016.90
- [14] LI W, ZHU X T, GONG S G. Harmonious attention network for person re-identification [C]//*2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, USA: CVPR, 2018: 2285 – 2294. DOI: 10.1109/cvpr.2018.00243
- [15] LI D W, CHEN X T, ZHANG Z, et al. Learning deep context-aware features over body and latent parts for person re-identification [C]//*2017 IEEE Conference on Computer Vision and Pattern Recognition*. Honolulu, USA: CVPR, 2017: 7398 – 7407. DOI: 10.1109/CVPR.2017.782
- [16] ZHAO L M, LI X, ZHUANG Y T, et al. Deeply-learned part-aligned representations for person re-identification [C]//*2017 IEEE International Conference on Computer Vision (ICCV)*. Venice, Italy: IEEE, 2017: 3219 – 3228. DOI: 10.1109/iccv.2017.349
- [17] JADERBERG M, SIMONYAN K, ZISSERMAN A. Spatial transformer networks [C]//*Proceedings of the 28th International Conference on Neural Information Processing Systems*. Montreal, Canada, 2015: 2017 – 2025
- [18] LI S, BAK S, CARR P, et al. Diversity regularized spatiotemporal attention for video-based person re-identification [C]//*2018 IEEE Conference on Computer Vision and Pattern Recognition*. Salt Lake City, USA: IEEE, 2018: 369 – 378. DOI: 10.1109/CVPR.2018.00046
- [19] XU J, ZHAO R, ZHU F, et al. Attention-aware compositional network for person re-identification [C]//*Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Salt Lake City, USA: CVPR, 2018: 2119 – 2128. DOI: 10.1109/cvpr.2018.00226
- [20] SARFRAZ M S, SCHUMANN A, EBERLE A, et al. A pose-sensitive embedding for person re-identification with expanded cross neighborhood re-ranking [C]//*2018 IEEE Conference on Computer Vision and Pattern Recognition*. Salt Lake City, USA: CVPR, 2018: 420 – 429. DOI: 10.1109/CVPR.2018.00051
- [21] HUANG H J, YANG W J, CHEN X T, et al. EANet: enhancing alignment for cross-domain person re-identification [EB/OL]. (2018-12-19) [2020-12-05]. <https://arxiv.org/abs/1812.11369>
- [22] MIAO J X, WU Y, LIU P, et al. Pose-guided feature alignment for occluded person re-identification [C]//*2019 IEEE International Conference on Computer Vision*. Seoul, South Korea: IEEE, 2019: 542 – 551. DOI: 10.1109/ICCV.2019.00063
- [23] HERMANS A, BEYER L, LEIBE B. In defense of the triplet loss for person re-identification [EB/OL]. (2017-03-22) [2020-12-12]. <https://arxiv.org/abs/1703.07737v4>
- [24] ZHENG L, SHEN L Y, TIAN L, et al. Scalable person re-identification: A benchmark [C]//*2015 IEEE International Conference on Computer Vision*. Santiago, Chile: IEEE, 2015: 1116 – 1124. DOI: 10.1109/ICCV.2015.133
- [25] RISTANI E, SOLERA F, ZOU R, et al. Performance measures and a data set

- for multi-target, multi-camera tracking [C]//European Conference on Computer Vision. Amsterdam, Netherlands: ECCV, 2016: 17 - 35. DOI: 10.1007/978-3-319-48881-3_2
- [26] LI W, ZHAO R, XIAO T, et al. DeepReID: deep filter pairing neural network for person re-identification [C]//2014 IEEE Conference on Computer Vision and Pattern Recognition. Columbus, USA: IEEE, 2014: 152 - 159. DOI: 10.1109/CVPR.2014.27
- [27] FELZENSZWALB P, MCALLESTER D, RAMANAN D. discriminatively trainedA, multiscale, deformable part model [C]//2008 IEEE Conference on Computer Vision and Pattern Recognition. Salt Lake City, USA: IEEE, 2008: 1 - 8. DOI: 10.1109/CVPR.2008.4587597
- [28] ZHONG Z, ZHENG L, CAO D L, et al. Re-ranking person re-identification with k-reciprocal encoding [C]//2017 IEEE Conference on Computer Vision and Pattern Recognition. Honolulu, USA: IEEE, 2017: 3652 - 3661. DOI: 10.1109/CVPR.2017.389
- [29] ZHENG L, YANG Y, HAUPTMANN A G. Person re-identification: past, present and future [EB/OL]. [2020-12-05]. <https://www.arxiv-vanity.com/papers/1610.02984/>
- [30] ZHENG Z D, ZHENG L, YANG Y. Pedestrian alignment network for large-scale person re-identification [J]. IEEE transactions on circuits and systems for video technology, 2019, 29(10): 3037 - 3045. DOI: 10.1109/TCSVT.2018.2873599
- [31] SUN Y F, ZHENG L, DENG W J, et al. SVDNet for pedestrian retrieval [C]//2017 IEEE International Conference on Computer Vision. Venice, Italy: IEEE, 2017: 3820 - 3828. DOI: 10.1109/ICCV.2017.410
- [32] WANG Z, HU R M, CHEN C, et al. Person re-identification via discrepancy matrix and matrix metric [J]. IEEE transactions on cybernetics, 2018, 48(10): 3006-3020. DOI: 10.1109/TCYB.2017.2755044
- [33] SONG C F, HUANG Y, OUYANG W L, et al. Mask-guided contrastive attention model for person re-identification [C]//2018 IEEE Conference on Computer Vision and Pattern Recognition. Salt Lake City, USA: IEEE, 2018: 1179 - 1188. DOI: 10.1109/CVPR.2018.00129
- [34] WANG Z, JIANG J J, WU Y, et al. Learning sparse and identity-preserved hidden attributes for person re-identification [J]. IEEE transactions on image processing, 2019, 29: 2013-2025. DOI: 10.1109/TIP.2019.2946975
- [35] SUN Y F, XU Q, LI Y L, et al. Perceive where to focus: Learning visibility-aware part-level features for partial person re-identification [C]//2019 IEEE Conference on Computer Vision and Pattern Recognition. Long Beach, USA: IEEE, 2019: 393 - 402. DOI: 10.1109/CVPR.2019.00048
- [36] FAN X, LUO H, ZHANG X, et al. SCPNet: spatial-channel parallelism network for joint holistic and partial person re-identification [C]//Asian Conference on Computer Vision. Perth, Australia: ACCV, 2018: 19 - 34. DOI: 10.1007/978-3-030-20890-5_2
- [37] FAN X, JIANG W, LUO H, et al. SphereReID: Deep hypersphere manifold embedding for person re-identification [J]. Journal of visual communication and image representation, 2019, 60: 51 - 58. DOI: 10.1016/j.jvcir.2019.01.010
- [38] ZHANG Z Z, LAN C L, ZENG W J, et al. Densely semantically aligned person re-identification [C]//2019 IEEE Conference on Computer Vision and Pattern Recognition. Long Beach, USA: IEEE, 2019: 667 - 676. DOI: 10.1109/CVPR.2019.00076
- [39] JIN X, LAN C L, ZENG W J, et al. Semantics-aligned representation learning for person re-identification [EB/OL]. (2019-05-30) [2020-12-05]. <https://arxiv.org/abs/1905.13143>

Biographies

CAO Jiahao (cao.jiahao@zte.com.cn) received the M.S. degree from Northeastern University, China in 2019 and joined ZTE corporation after he graduated. His current research interests include image processing and deep learning technologies.

MAO Xiaofei received the M.S. degree from TELECOM ParisTech, France in 2017. His current research interests include person re-identification, image processing and deep learning technologies.

LI Dongfang received the M.S. degree in electronics and communications engineering from Harbin Engineering University, China in 2017. He has been engaged in deep learning technologies in ZTE Corporation since his graduation.

ZHENG Qingfang received the B.S. degree in civil engineering and computer applications from Shanghai Jiaotong University, China in 2002, and the Ph.D. degree in computer sciences from Chinese Academy of Sciences, China in 2008. He is currently the chief scientist of video technology in ZTE Corporation. His research interests include computer vision, video codec, video streaming and multimedia content analysis and retrieval. He has published around 10 papers in various journals and conferences.

JIA Xia received her B.S. degree and M.S. degree in control theory and control engineering from Taiyuan University of Technology and Dalian University of Technology, China in 1995 and 2001, respectively. She joined ZTE Corporation in 2001 and worked in the State Key Laboratory of Mobile Network and Mobile Multimedia Technology. Her main research interests include deep learning techniques, face detection and recognition, Re-ID, and activity detection and recognition.



Adaptability Analysis of Fluctuating Traffic for IP Switching and Optical Switching

Abstract: The technological development of smart devices and Internet of Things (IoT) has brought ever-larger bandwidth and fluctuating traffic to existing networks. The analysis of network capital expenditure (CAPEX) is extremely important and plays a fundamental role in further network optimizing. In this paper, an adaptability analysis is raised for IP switching and optical transport network (OTN) switching in CAPEX when the service bandwidth is fluctuating violently. This paper establishes a multi-layer network architecture through Clos network model and discusses impacts of maximum allowable blocking rate and service bandwidth standard deviation on CAPEX of IP network and OTN network to find CAPEX demarcation point in different situations. As simulation results show, when the bandwidth deviation mean rate is 0.3 and the maximum allowable blocking rate is 0.01, the hardware cost of OTN switching will exceed IP switching as the average bandwidth is greater than 6 100 Mbit/s. When the service bandwidth fluctuation is severe, the hardware cost of OTN switching will increase and exceed IP switching as the single port rate is allowed in optical switching. The increasing of maximum allowable blocking rate can decrease hardware cost of OTN switching. Finally, it is found that Flex Ethernet (FlexE) can be used to decrease CAPEX of OTN switching greatly at this time.

Keywords: bandwidth fluctuating; CAPEX; OTN switching; IP switching; FlexE

LIAN Meng¹, GU Rentao¹, JI Yuefeng¹,
WANG Dajiang², LI Hongbiao²

(1. Beijing Laboratory of Advanced Information Network, Beijing University of Posts and Telecommunications, Beijing 100876, China;
2. Wireline Product Planning Department, ZTE Corporation, Shenzhen 518057, China)

DOI: 10.12142/ZTECOM.202101010

<http://kns.cnki.net/kcms/detail/34.1294.TN.20210204.1350.002.html>, published online February 4, 2021

Manuscript received: 2020-04-08

Citation (IEEE Format): M. Lian, R. T. Gu, Y. F. Ji, et al., "Adaptability analysis of fluctuating traffic for IP switching and optical switching," *ZTE Communications*, vol. 19, no. 1, pp. 82 - 90, Mar. 2021. doi: 10.12142/ZTECOM.202101010.

1 Introduction

Network data are growing dramatically due to the development of smart devices and multimedia application technologies. By 2021, global mobile traffic will reach 6.7 times that in 2016, exceeding 48.3 Exabytes per month^[1], as a result, existing networks will not be able to afford such huge business growth. To address this issue, 5G cellular networks are developed^[2].

With the advent of the 5G era, the development and application of various technologies are becoming possible. The 4G long term evolution (LTE) cellular system is difficult to support the transmission of such a large number of services brought by those applications, such as Internet of Things (IoT) and e-health, and the next generation datacenter needs to solve the problem of switching and transmission of massive services^[3-4]. In order to correspond to the exponential rise in

user and traffic capacity, datacenter capital expenditure (CAPEX) will increase significantly. The current datacenter inter-network is required to evolve into a high capacity, low-cost and low-latency network platform.

Packet services based on IP switching are emerging, which raises new requirements for transport networks. Optical networks can provide high-capacity end-to-end communication pipes for upper-layer services like IP services. Therefore, service transmission needs multiple electro-optic (E/O) and optical-electronic (O/E) conversions. Among the total services required processed by IP nodes, 55% - 85% of them only need to transit through these nodes.

Researchers have made significant research and proposed several methods to reduce CAPEX. In general, the problem of network cost is considered as integer linear programming (ILP) optimization problems^[5-7]. Heuristic approaches are usually applied to solve cost problems, including increasing restructure capacity^[8], adopting optical bypass to avoid traffic in the optical domain from contacting IP routers^[9-10], using software define network (SDN)^[11] for business grouping and multi-

This work was supported by National Key Research and Development Program (2018YFB1800504) and the ZTE Research Fund, China.

layer joint optimization to reduce costs^[12-13]. But increasing the number and capacity of infrastructure will lead to the rising cost of network construction as well^[8,13]. The method of adding optical bypass is not significant, and usually only saves 10% - 15% of the CAPEX cost^[13]. In the joint optimization process, multiple links are affected in the optimization process of a certain link, so multiple optimizations should be considered^[13-15].

Edge micro datacenters (mDC) are often used to handle small bandwidth and low latency service requests^[16-17]. On the other hand, service requests that interact between multiple datacenters (DC) typically have high bandwidth characteristics. The optical layer switching is considered to have large capacity, high bandwidth, and low latency. Therefore, optical transmission technologies^[18-19], like optical transport networks (OTN)^[20-21], are used extensively in datacenter networks to save cost and reduce latency while enabling high-speed data transmission.

However, OTN switching is very inflexible due to the large switching granularity. This inflexibility is more prominent when the bandwidth of the service fluctuates greatly. In order to ensure the traffic transmission blocking rate, the OTN switching must meet the largest bandwidth service requirement in the network, which will result in a large amount of cost waste, and therefore OTN switching is not necessarily the optimal choice. As a result, the cost analysis of large volatility bandwidth services needs further research.

Flex Ethernet (FlexE) is a good solution to this problem. FlexE supports the binding of multiple links and the link aggregation. For example, a 50 Gbit/s service can be transmitted by one 40 Gbit/s Ethernet port and one 10 Gbit/s Ethernet port instead of one 100 Gbit/s Ethernet port. Link binding using FlexE can greatly reduce resource waste of OTN switching.

In this paper, by calculating CAPEX, we perform adaptive analysis on IP switching and optical switching. We hope to obtain the bandwidth adaptation range of IP switching and optical switching when the average bandwidth and the standard deviation of the service are both large. Next, by abstracting the switching behavior of the IP router and the OTN device, a network connection switching model is established to perform CAPEX. Then, we introduce the adaptive assessment process to determine the bandwidth cross-point. After that, modifying the average bandwidth and the service standard of the service, we analyze hardware costs to accommodate the range of two switching methods. We analyze the impact of network service quality and service bandwidth fluctuations on network hardware costs by changing the maximum allowable blocking rate and the bandwidth standard deviation. Finally, we use FlexE technology in the model to explore the impact of port bonding on optical switching CAPEX.

The remaining parts of this paper are organized as follows. Section 2 introduces the structure and problems of the datacenter network. In Section 3, a mathematical model is de-

signed for calculating the cost of the switching network. Section 4 performs a performance evaluation based on simulation results. Finally, in Section 5, we summarize this paper.

2 Distribution of Datacenter Network Architecture and Cost

With the continuous development of cloud computing and 5G technologies, the granularity of service traffic increases. Under this trend, optical layer switching may have advantages in cost and energy consumption compared with IP switching.

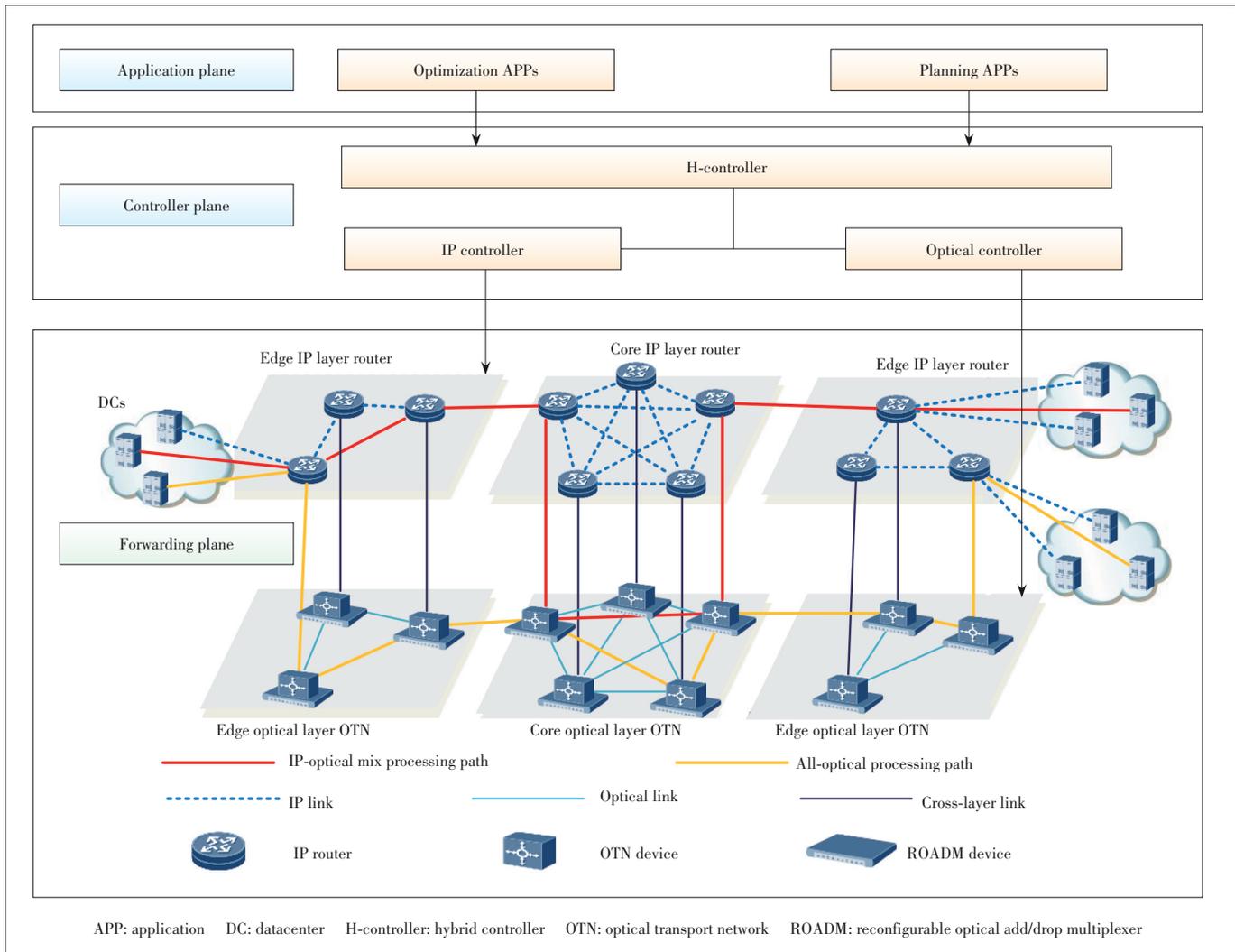
In the existing network, most of the transport layers select OTN and wavelength division multiplexing (WDM)^[22] system as the infrastructure. In this paper, we assume that the OTN is used as the network infrastructure^[20] technology.

2.1 Datacenter Network Architecture

The datacenter network architecture is mainly implemented through the coordination of the three major planes of application management, network control and service forwarding. As shown in **Fig. 1**, used for forwarding plane, DCs are interconnected with hosting services through a multi-domain network request based on the strategy provided by the controller plane. The SDN control technology is adopted in the controller plane, while centralized management and cross-layer joint network planning and optimization techniques are adopted in the application plane. In the forwarding plane, differentiated delivery of service with different quality of service (QoS) requirements is mainly achieved through the coordination of IP network and OTN network controllers, ensuring the efficiency of the entire network resources and management. In the controller plane, SDN technology is mainly used to realize the Openflow protocol. The smooth transition of the network will accommodate the flexible scheduling of multi-layer multi-domain distributed networks in the future. The main feature of the management plane is the unified and centralized management mode, which greatly simplifies network management processes, reduces network maintenance costs and improves network management and maintenance efficiency, thereby saving the total cost of the network.

The IP services in the datacenter network are diversified, and optical channels use wavelength switching. With large capacity, there is a fundamental mismatch between service traffic and optical flow capacity. Usually, the solution is using the traffic aggregation function of the core router in the IP layer to provide services to convergence processing, and then through the optical layer to achieve large-capacity long-distance transmission. This solution offers a high-speed information transmission channel for the IP layer.

However, service transmission requires times of E/O and O/E conversions, which greatly increases network CAPEX. Appropriately reducing IP routers processed can significantly reduce the cost of the network between datacenters. The optical



▲ Figure 1. Architecture of the datacenter network

layer switching has larger capacity, higher bandwidth, and lower latency compared with the IP layer switching. However, if the granularity of service requests is small, the large granularity of optical switching service will reduce resource utilization and further increase network CAPEX. Therefore, this architecture proposes the following strategy: When the service bandwidth is in the interval where the IP switching is more advantageous, the service transmission route selects the traditional transmission solution and transmits through both the IP equipment and the OTN equipment. Conversely, when the service bandwidth is in the interval where the optical switching is more advantageous, the service transmission route selection is directly transmitted through the OTN equipment. The equipment provides optical channel data unit (ODU) k -class hard pipes directly, and then transmits high-speed services through ROADM devices in the optical layer, thereby reducing equipment investment, saving energy consumption and promoting sustainable development of the transmission network.

This paper considers the direct transmission of services

through OTN equipment, avoiding multiple E/O and O/E conversions. OTN equipment is divided into edge-based aggregate OTN equipment and core switching OTN equipment. The aggregate OTN equipment uses the cross-point switching matrix to perform optical mixing processing of the service, and the service performs E/O and O/E conversion at the edge. At the core layer, the OTN equipment performs all-optical processing of the service through the ROADM to ensure reliable and high-quality service (yellow path in Fig. 1). Correspondingly, the IP network is also divided into two layers: the edge and the core. The edge aggregation router supports aggregation services, and the core layer uses the core router to provide a higher level of switching services (red path in Fig. 1). We want to get the range of adaptability of the two switching methods and evaluate the capital expenditures.

2.2 Composition of Network Equipment Cost

In order to facilitate comparison and calculation, this paper calculates the equipment cost after normalization. The cost pa-

parameters of the equipment are shown in **Tables 1** and **2**, and the cost is expressed as a standardized value. OTN equipment is divided into common units and service units. Accordingly, IP equipment includes base systems, processor cards, line cards, and optical modules. The common unit and the basic system are independent of the business service, and they are used to build the basic platform. The service unit is related to the specific service. Line cards and processor cards need to be used together. A mother card can carry two daughter cards. For example, if we need an IP router with two 100 G ports, we choose a base system as the platform, an A-type network processor card, two 100 G line cards and two 100 G optical modules. Therefore, the total cost is calculated as: $10 + 21.14 + 2 \times 58.57 + 2 \times 53.97 = 256.22$.

3 Problem Formalization

In this section, a mathematical model is proposed to obtain the network CAPEX for adaptive analysis.

3.1 Mathematical Model

Due to the large switching granularity of optical switching, the key is how to ensure connectivity in the network. In this paper, a strict non-blocking Clos network structure is used to ensure full network connectivity.

The Clos network was proposed by Charles CLOS in 1952^[23]. The Clos network refers to a switching network that tries to reduce the number of intermediate cross points as much as possible in order to reduce the cost of a multi-level switching network. The basic idea is to use multiple smaller-scale switching units according to a certain connection method. The advantage of Clos network is that its crossbar architecture and the Clos network can provide a non-blocking network.

Under the condition of non-blocking, the three-level Clos network is split into switching matrixes and constitutes a multi-level Clos network, each switching matrix representing a piece of equipment. As shown in **Figs. 2a** and **2b**, on the premise of ensuring strict non-blocking, datacenter network can be mapped to Clos network, and a three-layer symmetric switching

network is established according to the Clos model. Then, as shown in **Fig. 2c**, the switching network is divided into multi-

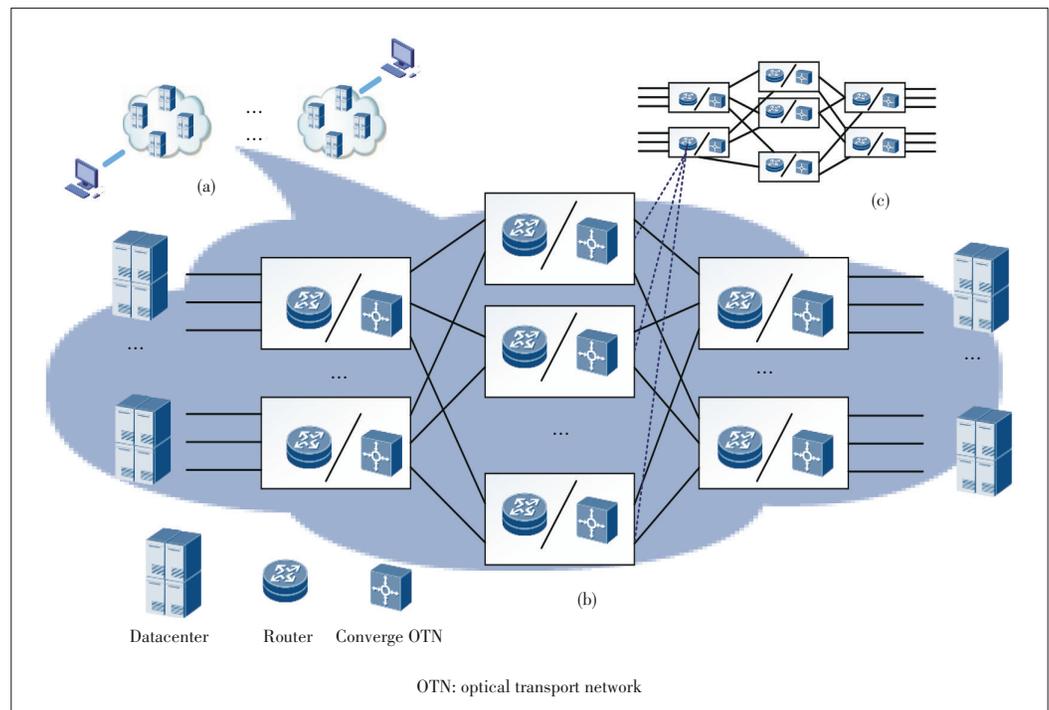
▼ **Table 1. Cost parameters for normalized OTN equipment**

OTN Equipment	Composition		Normalized Cost
Common unit	Common unit		3.83
Service unit	Port	Capacity/G	
	16	2.5	1.25
	12	10	2.22
	2	40	2.29
	2	100	1.59
	1	100	1.83

OTN: optical transport network

▼ **Table 2. Cost parameters for normalized IP router**

IP Router	Composition		Cost
Base system	Base system		10
Processor card	A-type network processor card		21.14
	B-type network processor card		84.13
Line card	Port	Capacity/G	
	1	100	58.57
	3	40	52.06
	10	10	24.57
	8	2.5	97.62
Optical module		100	53.97
		40	8.41
		10	0.22
		2.5	0.16



▲ **Figure 2. Multi-level Clos network: (a) datacenter network; (b) datacenter network mapped to Clos network model; (c) middle layer of Clos network for splitting**

ple layers based on the service parameters and network parameters. Then the minimum number of IP switching or OTN switching in the network is determined, and the number of links required by the network and by the number of devices and the network topology is calculated. The OTN equipment or IP equipment is mapped to the switching matrix (SM). The connections between them correspond to the optical links in the real network.

We assume that the service bandwidth obeys the normal distribution $N(\mu_{sd}, \sigma_{sd}^2)$, where the mean value is μ_{sd} , the standard deviation is σ_{sd}^2 , and the probability density function is $f(x)$.

The blocking rate $\alpha = \int_{b_{sd}^{\alpha}}^{+\infty} f(x)dx$ is allowed, and the value b_{sd}^{α} , the maximum bandwidth of the service request, can be obtained according to the blocking rate α . In order to make the model realistic, we assume that the bandwidth mean and the standard deviation of the service are proportional, which means that larger service has a larger standard deviation.

For the IP switching, with the dynamic bandwidth allocation function, bandwidth should be reserved for dynamic service. For OTN switching, it is necessary to reserve bandwidth resources for the service according to the maximum bandwidth to prevent random service.

For the OTN switching, the service request needs to be mapped to electrical cross-connection and ODU k before being multiplexed. Assuming that in ODU k , $k = 0, 1, 2, 3, 4$, the multiplexing relationship between the ODU k s can be expressed as: $2U_0 = U_1, 4U_1 = U_2, 4U_2 = U_3, 2U_3 = U_4$.

In **Table 3**, We define the required symbols, costs and variables. In this paper, a piece of core equipment can carry eight mother cards.

After ensuring the connectivity through the Clos network structure, the objective function can obtain the total hardware cost of

the network convergence layer and the core layer. And then the minimum cost for different speed card configurations is found.

$$\text{Min } C_{\text{total_ip}} = 2 \cdot N_{\text{edge_router}} \cdot C_{\text{ip}} + N_{\text{core_router}} \cdot C_{\text{ip}}, \quad (1)$$

$$\text{Min } C_{\text{total_otn}} = 2 \cdot N_{\text{edge_otn}} \cdot C_{\text{otn}} + N_{\text{core_otn}} \cdot C_{\text{otn}}. \quad (2)$$

The adaptation constraints of IP layer and OTN layer are evaluated as:

$$B_{sd} = U_k, U_{k-1} < b_{sd}^{\alpha} \leq U_k, \forall sd \in r, U_{-1} = 0. \quad (3)$$

When the FlexE technology is applied, OTN switching allows multiple port bindings to be used. At this time, the U_k is equal to $U_{k-1} + U_1$.

Strict non-blocking constraints are shown below. Under the premise of strict non-blocking, the network can establish a service request at any time, if both the source server and the destination server are free.

$$N_{\text{port_ip}}^j = \lceil (2 \cdot N_{\text{ip_card}}^j + 1) / 2 \rceil, \forall j \in P, \quad (4)$$

$$N_{\text{port_otn}}^j = \lceil (2 \cdot N_{\text{otn_card}}^j + 1) / 2 \rceil, \forall j \in P. \quad (5)$$

According to the requirements of the equipment, the first three types of line cards need to be matched with the type A network processor card, and the fourth type of the line card needs to be equipped with the type B network processor card. The cost of a single IP equipment C_{ip} and that of a single OTN equipment C_{otn} are calculated as follows:

$$C_{\text{ip}} = C_{\text{ip-base}} + \sum_{j=1}^3 N_{\text{ip_card}}^j \times (C_{\text{ip-a}} + 2C_{\text{ip-interface}}^j + 2C_{\text{ip-port}}^j \times C_{\text{ip-module}}^j) + N_{\text{ip_card}}^4 \times (C_{\text{ip-b}} + 2C_{\text{ip-interface}}^4 + 2C_{\text{ip-port}}^4 \times C_{\text{ip-module}}^4), \quad (6)$$

$$C_{\text{otn}} = C_{\text{otn-base}} + 2 \sum_{j=1}^4 N_{\text{otn_card}}^j \times (C_{\text{otn-business}}^j + \lfloor C_j \times C_{\text{otn-port}}^j / C_4 \rfloor \times C_{\text{otn-module}}^j). \quad (7)$$

After the switching network is split, the number of equipment in the convergence layer can be calculated:

$$N_{\text{edge_router}} = \left\lceil \sum_{\forall sd \in r} (N_r \times \mu_{sd} / \sum_{\forall j \in P_i} (N_{\text{port_ip}}^j \times C_j \times C_{\text{ip-port}}^j)) \right\rceil, \quad (8)$$

$$N_{\text{edge_otn}} = \left\lceil \sum_{\forall sd \in r} N_r \times B_{sd} / \sum_{\forall j \in P_i} (N_{\text{port_otn}}^j \times C_j \times C_{\text{otn-port}}^j) \right\rceil. \quad (9)$$

The middle tier n_{iter} should be split until $r(i) \leq m$. For IP switching or OTN switching, the number of iterations after

▼ **Table 3. Notations, costs and variables**

Notation	Description
N_r	Total number of service requests
$P = \{1, 2, 3, 4\}$	Card type set for IP/OTN equipment
C_j	Port rate set for different types of cards
$N_{\text{ip_card}}^j / N_{\text{otn_card}}^j$	Number of mother card j on one piece of IP/OTN equipment
$N_{\text{port_ip}}^j / N_{\text{port_otn}}^j$	Number of port for IP/OTN equipment on one card j
$C_{\text{ip-base}}$	Cost for base system of IP equipment
$C_{\text{ip-a}} / C_{\text{ip-b}}$	Cost for A-type/B-type card of IP equipment
$C_{\text{ip-interface}}^j / C_{\text{ip-module}}^j$	Cost for line card/optical module of IP equipment on one card j
$C_{\text{otn-base}}$	Cost for common unit of OTN equipment
$C_{\text{otn-business}}^j / C_{\text{otn-module}}^j$	Cost for service unit/optical module of OTN equipment
$C_{\text{total_ip}} / C_{\text{total_otn}}$	Total hardware cost for IP/OTN equipment
$N_{\text{edge_router}} / N_{\text{core_router}}$	Number of convergence/core IP equipment
$N_{\text{edge_otn}} / N_{\text{core_otn}}$	Number of convergence/core OTN equipment
B_{sd}	Bandwidth after the OTN multiplexing
n_{iter}	Iteration time for splitting the middle layer SM

OTN: optical transport network SM: switching matrix

splitting the middle tier n_{iter} is expressed as follows:

$$n_c(n_{iter}) = 2 \cdot r(n_{iter}) + m \cdot n_c(n_{iter} - 1), n_c(0) = 1, \quad (10)$$

$$n = \left\lceil \frac{m+1}{2} \right\rceil, \quad (11)$$

$$r(i) = \lfloor r(i-1)/n \rfloor, i = 2, 3, \dots, n_{iter}. \quad (12)$$

For IP switching:

$$N_{core_router} = 2 \cdot \sum_{\forall j \in P_i} N_{ip_card}^j \cdot n_c(n_{iter}), \quad (13)$$

$$m = 2 \cdot \sum_{\forall j \in P_i} N_{ip_card}^j, \quad (14)$$

$$r(1) = \lfloor N_{edge_router}/n \rfloor. \quad (15)$$

For OTN switching:

$$N_{core_otn} = 2 \cdot \sum_{\forall j \in P_i} N_{otn_card}^j \cdot n_c(n_{iter}), \quad (16)$$

$$m = 2 \cdot \sum_{\forall j \in P_i} N_{otn_card}^j, \quad (17)$$

$$r(1) = \lfloor N_{edge_otn}/n \rfloor, \quad (18)$$

where m and n are Clos network parameters.

Algorithm 1 Network Adaptability Assessment Process

Given: μ_{sd} (Mbit/s), α , σ_{sd} , service request number set X_t

Output: hardware cost C_{total_ip} and C_{total_otn} :

- 1: **For** each t in X_t , do
- 2: $\mu_{sd} = 1$;
- 3: **While** $\mu_{sd} \geq 1$ do
- 4: Calculate the B_{sd} according to μ_{sd} , σ_{sd} and α ;
- 5: Establish a three-level symmetric Clos network under strict non-blocking principle;
- 6: Map the first and third layer SM to the convergence layer equipment, and then calculate the N_{edge_router} and N_{edge_otn} by Eq. (8) and Eq. (9);
- 7: Split and map the middle layer SM based on the port rate of core devices;
- 8: Calculate the C_{ip} and C_{otn} by Eq. (6) and Eq. (7);
- 9: Calculate the N_{core_router} and N_{core_otn} by Eq. (13) and Eq. (16);
- 10: Calculate the $\text{Min } C_{total_ip}$ and $\text{Min } C_{total_otn}$ by Eq. (1) and

Eq. (2);

11: $\mu_{sd} = \mu_{sd} + 1$;

12: **end While**

13: **end For**

14: **Return** $\text{Min } C_{total_ip}$ and $\text{Min } C_{total_otn}$

3.2 Adaptive Assessment Procedure

For IP switching, since the service bandwidth is normally distributed and according to the symmetry of the normal distribution, the large service can always be combined with a small service for transmission. So the maximum capacity of the IP equipment in the network is the average bandwidth. The hardware cost of the IP switching is only related to the average bandwidth of the service in the network and is independent of the standard deviation of the service bandwidth in the network. The OTN switching focuses on the maximum service bandwidth in the network. In order to transmit the largest bandwidth services, the network must use OTN equipment with higher capacity, which will result in a waste of transmission resource. The hardware cost of OTN switching is related to not only the average bandwidth of network services, but also the standard deviation of service bandwidth. When the services with large bandwidth and large bandwidth standard deviation are transmitted, the hardware cost of OTN switching may be greater than the IP switching.

In order to describe the results of cost-adaptive evaluation more accurately, we introduce the concept of bandwidth cross-point, which is the cost demarcation point for service bandwidth between IP switching and OTN switching. We divide the cross-points into first bandwidth cross-point and second bandwidth cross-point. When the IP switching cost exceeds the OTN switching, this demarcation point is the first bandwidth cross-point. When the OTN switching cost exceeds the IP switching, this demarcation point is the second bandwidth cross-point. Since the first bandwidth cross-point is usually small and we are more concerned with the large service bandwidth, the second bandwidth cross-points are mainly discussed and analyzed.

In the adaptive evaluation process, the minimum cost value of the total hardware of the network is calculated under the premise of ensuring the network blocking rate. For IP switching, we allocate port resources based on the average bandwidth value for the service request. After that, we set the service standard deviation to be proportional to the business mean value. When the blocking rate is guaranteed to α , the service granularity is changed, and the total hardware cost can be calculated.

4 Simulation Results and Performance Analysis

In this section, IP equipment and OTN equipment support

mixed-port rate transmission (2.5/10/40/100 Gbit/s) in the simulation, and get the lowest cost by calculating the lowest cost connection topology in multiple hybrid modes. In this section, we first explore the relationship between the number of businesses and the cost of network hardware. Then we study the impact of the maximum allowable blocking rate and the ratio of the standard deviation to the service mean of bandwidth (deviation-mean rate) to the hardware cost of IP switching and OTN switching. Finally, we apply the FlexE technology in the network to explore whether the port binding can reduce the OTN switching hardware cost.

As the average bandwidth is 6 300 Mbit/s, the deviation-mean rate is 0.3, and the maximum allowable blocking rate of the network is 0.01. The hardware cost of IP switching and OTN switching is shown in Fig. 3a. The deviation-mean rate is 0.5, and the maximum allowable blocking rate of the network is 0.1. The hardware cost of IP switching and OTN

switching is shown in Fig. 3b. When the number of services is greater than 1 250, the OTN switching hardware cost is higher than the IP switching. Therefore, when the services with large bandwidth and large bandwidth standard deviation are transmitted, the hardware cost of OTN switching will exceed the IP switching.

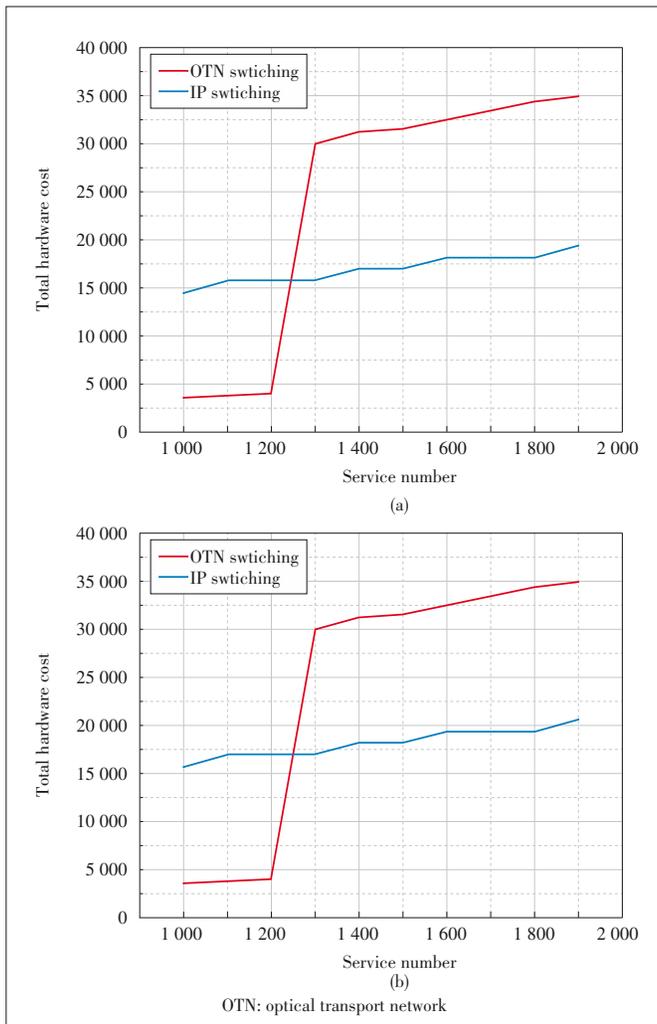
After that, we change the average bandwidth of the service, the deviation-mean rate, and the maximum blocking rate of the network when when the number of services circulated in the network is 1 800. As shown in Fig. 4a, when the deviation-mean rate is 0.3 and the maximum allowable blocking rate of the network is 0.01, as the average bandwidth of the service increases, if the service bandwidth is less than 6 100 Mbit/s, the hardware cost for OTN switching will be less than IP switching. However, as the average bandwidth is greater than 6 100 Mbit/s, the hardware cost of OTN switching will exceed the IP switching.

As shown in Fig. 4b, the maximum allowable blocking rate of the network is 0.01 and the deviation-mean rate is increased to 0.5. As the average bandwidth of the service increases, and the average bandwidth is greater than 4 800 Mbit/s, the hardware cost of OTN switching will exceed IP switching. The second bandwidth cross-point is smaller than the second bandwidth cross-point when the deviation-mean rate is 0.3.

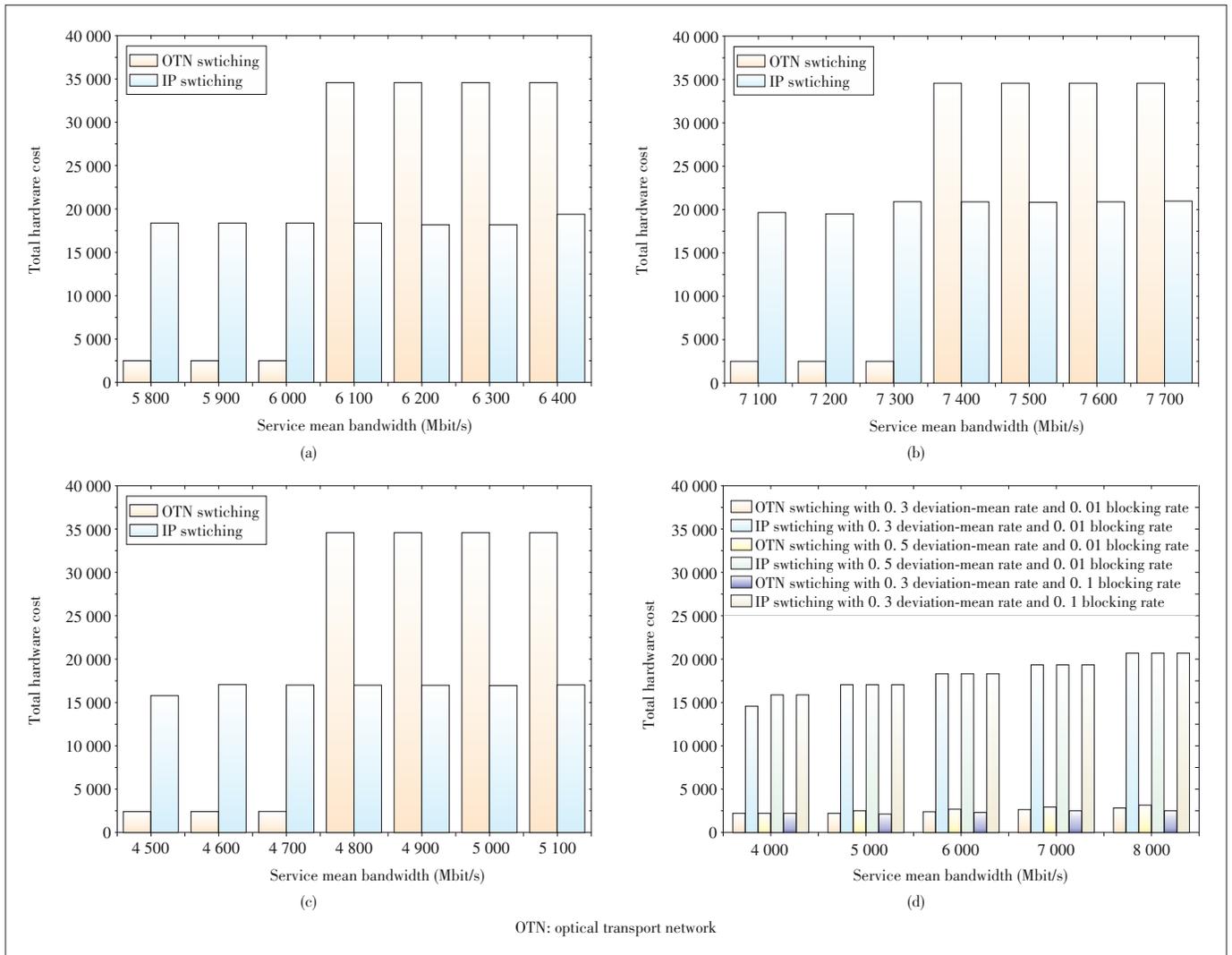
As shown in Fig. 4c, the deviation-mean rate is 0.3 and the maximum allowable blocking rate of the network increases to 0.1. As the average bandwidth of the service increases, and the service bandwidth is greater than 7 400 Mbit/s, the hardware cost of OTN switching will exceed IP switching. The second bandwidth cross-point is larger than the second bandwidth cross-point at the maximum blocking rate of 0.01.

An important breakthrough in FlexE technology is physical channel bonding. For example, it supports 200 Gbit/s media access control (MAC) based on two bonded 100 Gbit/s physical layer channels. In this model, FlexE technology can be implemented by the combination binding of different ODUks and a bonded ODUk can be used as a new port with a different granularity. Using the FlexE technology to bind OTN equipment ports can reduce the hardware cost of OTN switching effectively. Because the multiplexing relationship of the ODUks is $2U_0 = U_1, 4U_1 = U_2, 4U_2 = U_3, 2U_3 = U_4$, where the three ports bindings can meet the vast majority of granularity. As shown in Fig. 4d, when we allow three-port bindings, the hardware cost of OTN switching is significantly reduced and the second bandwidth cross-point is no longer present.

Since the maximum blocking rate of the network should be guaranteed, in the model we assume that the services with larger bandwidth are blocked, which is also practical. Therefore, the smaller the maximum allowable blocking rate of the network is, the larger the maximum bandwidth of the service which can actually pass through the network is. Similarly, increasing the standard deviation of service bandwidth can also increase the maximum bandwidth of services in the network.



▲ Figure 3. Hardware cost of IP switching and OTN switching: (a) with 6 300 Mbit/s mean bandwidth, 0.3 deviation-mean rate and 0.01 blocking rate; (b) with 6 300 Mbit/s mean bandwidth, 0.5 deviation-mean rate and 0.1 blocking rate



▲ Figure 4. Total hardware cost when service number is 1 800 with: (a) 0.3 deviation-mean-rate and 0.01 blocking rate; (b) 0.5 deviation-mean-rate and 0.01 blocking rate; (c) 0.3 deviation-mean-rate and 0.1 blocking rate; (d) cost of IP switching and OTN switching after using FlexE

At this time, the OTN switching needs to transmit the largest bandwidth service in the network, so the hardware cost of OTN switching will increase significantly, which leads to resources waste in the increased hardware cost.

FlexE technology can be used to bind OTN equipment ports. When OTN equipment needs to transmit large bandwidth services, other untagged port resources can be also used. In this way, the resource waste can be reduced and the OTN switching hardware cost can be decreased.

5 Conclusions

This paper analyzes the hardware cost of datacenter networks. In most solutions, optical switching is considered for large-bandwidth service transmission. This paper not only considers the service bandwidth, but also the fluctuation of the service bandwidth. The hardware cost of IP switching and

OTN switching when transmitting large bandwidth and large fluctuation services are discussed. The simulation results show that although OTN switching is more suitable for transmitting large bandwidth services in principle, the hardware cost of OTN switching will exceed IP switching when large bandwidth services have large bandwidth fluctuations. After that, this paper analyzes the impact of maximum allowable blocking rate of the network and the standard deviation of service bandwidth on OTN switching hardware cost.

When the deviation-mean rate is 0.3 and the maximum allowable blocking rate of the network is 0.01, as the average bandwidth is greater than 6 100 Mbit/s, the hardware cost of OTN switching will exceed the IP switching. As the deviation-mean rate increases, the OTN switching hardware cost increases; as the maximum allowable blocking rate of the network increases, the OTN switching hardware cost decreases. Therefore, the fluctuation of service bandwidth will increase the

hardware cost of OTN switching and eventually exceed IP switching. The increase of maximum allowable blocking rate of the network can decrease the OTN switching cost. Finally, the simulation results show using FlexE can effectively reduce the OTN hardware cost.

References

- [1] Cisco visual networking index: forecast and methodology, 2016 – 2021 [EB-OL]. [2020-02-01]. <https://www.cisco.com/c/en/us/solutions/collateral/service-provider/visualnetworking-index-vni/complete-white-paper-c11481360.html>
- [2] LUONG N C, WANG P, NIYATO D, et al. Applications of economic and pricing models for resource management in 5G wireless networks: A survey [J]. IEEE communications surveys & tutorials, 2019, 21(4): 3298 – 3339. DOI: 10.1109/COMST.2018.2870996
- [3] HOLMA H, TOSKALA A, REUNANEN J. LTE small cell optimization [M]. Chichester, UK: John Wiley & Sons Ltd., 2015. DOI: 10.1002/9781118912560
- [4] eLTEHuawei2.2 DBS3900 LTE configuration principles [EB-OL]. [2020-02-01]. <http://e.huawei.com/au/marketingmaterial/onLineView?MaterialID=%7B7C9C0FCC-2359-4709-83C2-7E4F7C16A495%7D>
- [5] GERSTEL O, FILSFILS C, TELKAMP T, et al. Multi-layer capacity planning for IP-optical networks [J]. IEEE communications magazine, 2014, 52(1): 44 – 51. DOI: 10.1109/MCOM.2014.6710063
- [6] GKAMAS V, CHRISTODOULOPOULOS K, VARVARIGOS E. A joint multi-layer planning algorithm for IP over flexible optical networks [J]. Journal of lightwave technology, 2015, 33(14): 2965 – 2977. DOI: 10.1109/JLT.2015.2424920
- [7] KATIB I, MEDHI D. IP/MPLS-over-OTN-over-DWDM multilayer networks: an integrated three-layer capacity optimization model, a heuristic, and a study [J]. IEEE transactions on network and service management, 2012, 9(3): 240 – 253. DOI: 10.1109/TNSM.2012.12.110124
- [8] ZHAO X X, VUSIRIKALA V, KOLEY B, et al. The prospect of inter-data-center optical networks [J]. IEEE communications magazine, 2013, 51(9): 32 – 38. DOI: 10.1109/MCOM.2013.6588647
- [9] GHONAIM F A, DARCIÉ T E, GANTI S. Impact of SDN on optical router bypass [J]. IEEE/OSA journal of optical communications and networking, 2018, 10(4): 332 – 343. DOI: 10.1364/JOCN.10.000332
- [10] HUANG Y C, YOSHIDA Y, IBRAHIM S, et al. Bypassing route strategy for optical circuits in OPS-based datacenter networks [C]//Photonics in Switching. Florence, Italy: IEEE, 2015. DOI: 10.1109/PS.2015.7328965
- [11] Software-defined networking: the new norm for networks [EB-OL]. [2020-02-01]. <https://www.semanticscholar.org/paper/Software-Defined-Networking-The-New-Norm-for-Tank-Dixit/6457799bfda12f18c6f3f6cdaad1848bcc4e3aa2>
- [12] LU P, ZHU Z Q. Data-oriented task scheduling in fixed-and flexible-grid multi-layer inter-DC optical networks: a comparison study [J]. Journal of lightwave technology, 2017, 35(24): 5335 – 5346. DOI: 10.1109/JLT.2017.2777605
- [13] DAS S, PARULKAR G, MCKEOWN N. Rethinking IP core networks [J]. Journal of optical communications and networking, 2013, 5(12): 1431. DOI: 10.1364/jocn.5.001431
- [14] NIKOLAIDIS A. Leveraging FlexGrid and advanced modulations in a multilayer, inter-datacenter network [C]//Optical Fiber Communications Conference and Exhibition (OFC). Los Angeles, USA: IEEE, 2017: 1 – 3
- [15] PAPANIKOLAOU P, CHRISTODOULOPOULOS K, VARVARIGOS E. Joint multi-layer survivability techniques for IP-over-elastic-optical-networks [J]. Journal of optical communications and networking, 2016, 9(1): 85. DOI: 10.1364/jocn.9.000a85
- [16] LIU Z, ZHANG J W, BAI L, et al. Joint jobs scheduling and routing for metro-scaled micro datacenters over elastic optical networks [C]//2018 Optical Fiber Communications Conference and Exposition (OFC). San Diego, USA: IEEE, 2018: 1 – 3
- [17] CAO X Y, CHEN G, POPESCU I, et al. Adaptive DC interconnection provisioning in distributed all-optical micro-datacenters using holistically SDN orchestration for dynamic access [J]. Photonic network communications, 2018, 35(2): 129 – 140. DOI: 10.1007/s11107-017-0735-7
- [18] ASENSIO A, RUIZ M, VELASCO L. Requirements to support cloud, video and 5G services on the telecom cloud [C]//2016 21st European Conference on Networks and Optical Communications (NOC). Lisbon, Portugal: IEEE, 2016: 64 – 69. DOI: 10.1109/NOC.2016.7506987
- [19] TAN Y X, GU R T, LIAN M, et al. Cost comparison and adaptability analysis for OTN switching and IP switching [C]//Fiber Optic Sensing and Optical Communication. Kunming, China: 2018: 1084 DOI: 10.1117/12.2502269
- [20] International Telecommunication Union (ITU - T). Interfaces for the Optical Transport Network (OTN): Rec. ITU-T G.709/Y.1331 [S], 2012
- [21] LÓPEZ VIZCAÍNO J, YE Y, LÓPEZ V, et al. OTN switching for improved energy and spectral efficiency in WDM MLR networks [C]//Optical Fiber Communication Conference. Anaheim, USA: OSA, 2016. DOI: 10.1364/ofc.2016.m3k.3
- [22] YANG H, ZHANG J, ZHAO Y, et al. SUDO: software defined networking for ubiquitous datacenter optical interconnection [J]. IEEE Communications. 54(2): 86 – 95
- [23] CLOS C. A study of non-blocking switching networks [J]. Bell labs technical journal, 1953, 32(2): 406 – 424

Biographies

LIAN Meng received the B.S. degree from Beijing University of Posts and Telecommunications, China in 2018. He is currently pursuing the master degree with the State Key Laboratory of Information Photonics and Optical Communications, Beijing University of Posts and Telecommunications. His research interests include intelligent optical network and IP-optical network survivability.

GU Rentao (rentaogu@bupt.edu.cn) received the B.E. and Ph.D. degrees from Beijing University of Posts and Telecommunications (BUPT), China in 2005 and 2010, respectively. From 2008 to 2009, he was a visiting scholar with the Georgia Institute of Technology, USA. He is currently a professor and Vice Dean in the School of Information and Telecommunication Engineering, BUPT. His current research interests include optical networking and intelligent information processing. He is a senior member of IEEE, China Institute of Communications, and Chinese Institute of Electronics.

JI Yuefeng received the Ph.D. degree from Beijing University of Posts and Telecommunications, China. He is currently a professor and the deputy director of the State Key Lab of Information Photonics and Optical Communications there. His research interests are primarily in the area of broadband communication networks and optical communications, with emphasis on key theory, realization of technology, and applications. He is a Fellow of the China Institute of Communications, Chinese Institute of Electronics, and IET.

WANG Dajiang is currently the product planning manager of ZTE's Bearer Network. His main research interests include intelligent network analysis, intent based optical network, and the applications of digital twin in the field of optical network. He participated in 3 national science and technology projects related to intelligent optical network management and control, published more than 10 technical documents, and obtained more than 20 Chinese and International granted patents.

LI Hongbiao is currently the chief engineer of ZTE Bearer Network Product Line Planning. He has been engaged in the research and planning on SDN/NFV, IP + optical, and cloud datacenter products for many years. The related products and solutions have won many awards in SDN/NFV Global Conference, China Institute of Communication, etc.

ZTE Communications Guidelines for Authors

Remit of Journal

ZTE Communications publishes original theoretical papers, research findings, and surveys on a broad range of communications topics, including communications and information system design, optical fiber and electro-optical engineering, microwave technology, radio wave propagation, antenna engineering, electromagnetics, signal and image processing, and power engineering. The journal is designed to be an integrated forum for university academics and industry researchers from around the world.

Manuscript Preparation

Manuscripts must be typed in English and submitted electronically in MS Word (or compatible) format. The word length is approximately 3 000 to 8 000, and no more than 8 figures or tables should be included. Authors are requested to submit mathematical material and graphics in an editable format.

Abstract and Keywords

Each manuscript must include an abstract of approximately 150 words written as a single paragraph. The abstract should not include mathematics or references and should not be repeated verbatim in the introduction. The abstract should be a self-contained overview of the aims, methods, experimental results, and significance of research outlined in the paper. Five carefully chosen keywords must be provided with the abstract.

References

Manuscripts must be referenced at a level that conforms to international academic standards. All references must be numbered sequentially in-text and listed in corresponding order at the end of the paper. References that are not cited in-text should not be included in the reference list. References must be complete and formatted according to *ZTE Communications Editorial Style*. A minimum of 10 references should be provided. Footnotes should be avoided or kept to a minimum.

Copyright and Declaration

Authors are responsible for obtaining permission to reproduce any material for which they do not hold copyright. Permission to reproduce any part of this publication for commercial use must be obtained in advance from the editorial office of *ZTE Communications*. Authors agree that a) the manuscript is a product of research conducted by themselves and the stated co-authors; b) the manuscript has not been published elsewhere in its submitted form; c) the manuscript is not currently being considered for publication elsewhere. If the paper is an adaptation of a speech or presentation, acknowledgement of this is required within the paper. The number of co-authors should not exceed five.

Content and Structure

ZTE Communications seeks to publish original content that may build on existing literature in any field of communications. Authors should not dedicate a disproportionate amount of a paper to fundamental background, historical overviews, or chronologies that may be sufficiently dealt with by references. Authors are also requested to avoid the overuse of bullet points when structuring papers. The conclusion should include a commentary on the significance/future implications of the research as well as an overview of the material presented.

Peer Review and Editing

All manuscripts will be subject to a two-stage anonymous peer review as well as copyediting, and formatting. Authors may be asked to revise parts of a manuscript prior to publication.

Biographical Information

All authors are requested to provide a brief biography (approx. 100 words) that includes email address, educational background, career experience, research interests, awards, and publications.

Acknowledgements and Funding

A manuscript based on funded research must clearly state the program name, funding body, and grant number. Individuals who contributed to the manuscript should be acknowledged in a brief statement.

Address for Submission

<http://mc03.manuscriptcentral.com/ztecom>

ZTE COMMUNICATIONS

中兴通讯技术(英文版)

ZTE Communications has been indexed in the following databases:

- Abstract Journal
- Cambridge Scientific Abstracts (CSA)
- China Science and Technology Journal Database
- Chinese Journal Fulltext Databases
- Index of Copernicus
- Inspec
- Ulrich's Periodicals Directory
- Wanfang Data

ZTE COMMUNICATIONS

Vol. 19 No. 1 (Issue 73)

Quarterly

First English Issue Published in 2003

Supervised by:

Anhui Publishing Group

Sponsored by:

Time Publishing and Media Co., Ltd.

Shenzhen Guangyu Aerospace Industry Co., Ltd.

Published by:

Anhui Science & Technology Publishing House

Edited and Circulated (Home and Abroad) by:

Magazine House of ZTE Communications

Staff Members:

General Editor: WANG Xiyu

Editor-in-Chief: JIANG Xianjun

Executive Editor-in-Chief: HUANG Xinming

Editor-in-Charge: ZHU Li

Editors: REN Xixi, LU Dan, XU Ye, YANG Guangxi

Producer: XU Ying

Circulation Executive: WANG Pingping

Liaison Executive: LU Dan

Assistant: WANG Kun

Editorial Correspondence:

Add: 12F Kaixuan Building, 329 Jinzhai Road,
Hefei 230061, P. R. China

Tel: +86-551-65533356

Email: magazine@zte.com.cn

Website: <https://tech-en.zte.com.cn>

Annual Subscription: RMB 80

Printed by:

Hefei Tiancai Color Printing Company

Publication Date: March 25, 2021

China Standard Serial Number: ISSN 1673-5188
CN 34-1294/TN