

Ultra-Low Latency Mobile Networking

Kwang-Cheng Chen, Tao Zhang, Richard D. Gitlin, and Gerhard Fettweis

ABSTRACT

Mobile networking to achieve the ultra-low latency goal of 1 msec enables massive operation of autonomous vehicles and other intelligent mobile machines, and emerges as one of the most critical technologies beyond 5G mobile communications and state-of-the-art vehicular networks. Introducing fog computing and proactive network association, realizing virtual cell by integrating open-loop radio transmission and error control, and innovating anticipatory mobility management through machine learning, opens a new avenue toward ultra-low latency mobile networking.

INTRODUCTION

Mobile communications have been one of the most important technological innovations in recent decades. Currently, commercial deployment of fifth generation (5G) mobile communications is within our sight, starting as early as 2018 or 2019, after 3GPP R14 and the coming R15. There are three pillar technologies for 5G: enhanced mobile broadband (eMBB), massive machine type communication (MTC), and ultra-reliable and low-latency communication (uRLLC). eMBB has been well supported by technology innovations such as cloud radio access networks (C-RANs), massive MIMO, small cell, and so on. However, effective massive access and low latency networking remains open to further improve beyond 5G [1].

In the meantime, vehicular communications as one critical application scenario of Internet of Things (IoT) is approaching the stage of deployment. One example is dedicated short-range communication (DSRC) for connected vehicles. More exciting situations are autonomous driving vehicles, while level-3 (conditional automation defined by SAE) autonomous driving has been trialed and soon commercially deployed. Nevertheless, safety and reliable deployment of autonomous vehicles (AVs) on a massive scale must rely on ultra-low latency mobile networking to keep end-to-end networking latency down to the 1 msec range, which is even more challenging than the tactile Internet [2] due to the need to support high mobility of AVs. With state-of-the-art mobile networks supporting networking latency in the range of 100 msec or higher, we have to think out of the box to develop new technological solutions.

In this article, we focus on ultra-low latency mobile networking for AVs, which also provides many benefits for other intelligent mobile machines (IMMs) such as service robots. Starting

from computing scenarios, a new networking architecture has been identified, and then re-innovations of open-loop wireless communications beyond state-of-the-art low latency techniques have been introduced. The proposed approach integrates the idea of virtual cell for each AV, network virtualization, and proactive network association to reduce end-to-end networking latency toward 1 msec. Subsequent challenges in asynchronous multiple access can be resolved by multiuser detection (MUD). Finally, machine learning enabling anticipatory mobility management to serve proactive network association and open-loop wireless communications is shown to be effective. Such a holistic mobile network architecture accomplishing ultra-low latency fundamentally innovates the technology of wireless networks and communications.

FOG COMPUTING AND HETEROGENEOUS NETWORKING

A state-of-the-art AV operates based on computer vision with the assistance of LIDARs on the vehicles to achieve autonomous driving as a reliable and never tired human driver. By referring to the view from the National Highway Traffic Safety Administration (NHTSA) on vehicle-to-vehicle communications, we use Fig. 1 to illustrate the insufficiency without networking when multiple AVs are interacting with each other and the environment. As a large-size animal (e.g., a deer) on the roadside jumps into the road, which is not rare in North America, the silver vehicle has to turn left into the path of the red car as the silver vehicle's vision devices may not be able to observe such a situation due to not being in line-of-sight. One solution is vehicle-to-vehicle (V2V) communication and multi-hop ad hoc networking, which suffers from scalability limitations when the number of vehicles becomes large and consequently prohibits reaching low networking latency [3]. A new approach is to leverage fog computing that combines computing at the edge and predictive data analytics of sensor and traffic data. As illustrated in Fig. 1, an anchor node (AN) governing a number of wireless access points (APs) to roadside information infrastructure works together with fog computing. One or more sensors may sense the animal that poses a potential danger to the vehicles and send messages to fog computing through AP(s). Fog computing analyzes the risks to inform the silver vehicle and red vehicle, and to warn other vehicles to take proper safety actions. Without going through the cloud, which requires much longer latency, appropriate actions of ultra-low latency can be possibly achieved in a more

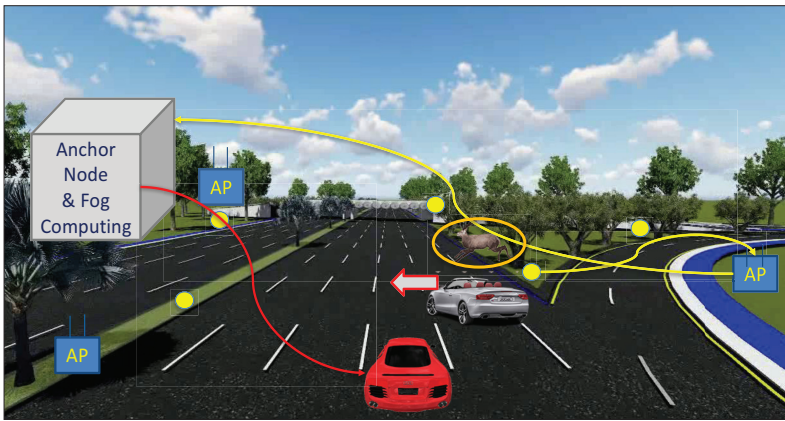


FIGURE 1. Fog/edge computing and networking enhances safety and reliability. A deer is on the roadside to jump into the road. Yellow dots represent various types of sensors with wireless communication capabilities including cameras. The anchor node to govern APs shall be co-located with fog computing facilities.

reliable way than V2V communications (note many vehicles can be on the road) via fog/edge computing and networking [4]. We adopt fog computing in this article to integrate data analytics for predictive safety and real-time actions, traffic analysis, real-time control and management of AV operations, and anticipatory mobility management to achieve ultra-low latency mobile networking. Such a mechanism can also resolve the dilemma for recent fatal accidents by testing AV that just relies on visionary and line-of-sight sensors.

For the purpose of safety and reliability, fog/edge computing and networking is therefore a promising architecture to achieve ultra-low latency. Since the deployment of APs might not always guarantee successful networking due to radio resource constraints, interference, fading, packet coverage of fog networking, and so on, falling back to cellular coverage suggests forming a heterogeneous networking architecture [5] (Fig. 2), to ensure reliability of the networking and computing. The overall multi-scale computing environment can be summarized in three ways:

- Cloud computing as a baseline architecture for reliable and secure management/control of various AVs, through rich computing power, database, machine learning and artificial intelligence, at the price of end-to-end latency on the order of seconds or even longer.
- Edge/Fog computing to assist AVs via *low-latency* command and control that incorporates information from the physical surroundings and from local agents (such as other AVs and humans) [6]. In this article we will demonstrate the effectiveness by introducing edge/fog computing into traditional cloud-on-board computing scenario.
- Onboard computing in an AV to decide and act based on local sensors and actuators, such as vision and speech recognition, external sensors from roadside units and surveillance cameras, *ad hoc information* from other agents and so on, and management/command/control messages from the cloud and fog for given missions (e.g., autonomous driving to a designated destination).

The consequent networking architecture to support above multi-scale computing, consistent

with the vision in [7], shall integrate sensors, cellular network consisting of high power nodes (HPNs), radio access network (RAN), and core network, and fog network including access points (APs) governed by the anchor node (AN). AN is connected to the core network under its management such that fog network can supply ultra-low latency mobile networking while cellular network can fall back in case of no coverage of APs.

With the above advantages of using fog computing and networking to achieve ultra-low latency mobile networking, there still exist the following technological challenges:

- Control signaling storm by wireless closed-loop physical layer (PHY) communication [8, 9], involving thousands of control messages of power control or channel estimation messages in a single PHY link.
- Protocol stacks and optimization of routing and scheduling in networking consuming significant time on computation and information gathering.
- DSRC using carrier sense multiple access as the medium access control protocol, which is inherently unstable [10] to result in random-access and networking delay much longer than several milliseconds.

RE-INNOVATION OF TECHNOLOGIES

Current state-of-the-art designs toward low latency wireless communication primarily considers techniques for fast re-transmission, uplink grant free transmission, and short scheduling [11]. In this section, we will present methods to further push technologies in these directions to the extreme.

OPEN-LOOP WIRELESS COMMUNICATION

First, let us re-consider how to speed up re-transmission. The latency to transmit a packet or a frame of bits is not just the transmission time in the air. The communication latency mainly comes from signal processing and coding/decoding of bits and error control of the link. If a packet is correctly received, a positive acknowledgment shall be sent. If a packet is incorrectly received by passing CRC check, the receiver has to send a negative Acknowledgment to the transmitter and the transmitter will re-transmit the packet again, while hybrid ARQ can further improve. A closed-loop is actually formed between transmitter and receiver. If multiple access or medium access control is considered with PHY, it is known that four-way handshaking such as the distributed coordination function of IEEE 802.11 wireless local area networks is the way to retain reliable communication of packets [10]. With minimal protocol execution, the consequent latency is at least four times the single transmission/reception period that includes preamble processing/detection, signal processing, modulation/demodulation, coding/decoding and protocol execution time by processors at both the transmitter and receiver ends.

The fundamental solution toward the lowest possible latency in radio transmission is to adopt *open-loop communication* by avoiding re-transmission, in which the transmitter simply selects radio resource units (RRUs) and transmits the packet to the designated receiver without any further processing. Optimal operation of open-loop communication has been initially investigated in [8, 9], including the adaptive operation with convention-

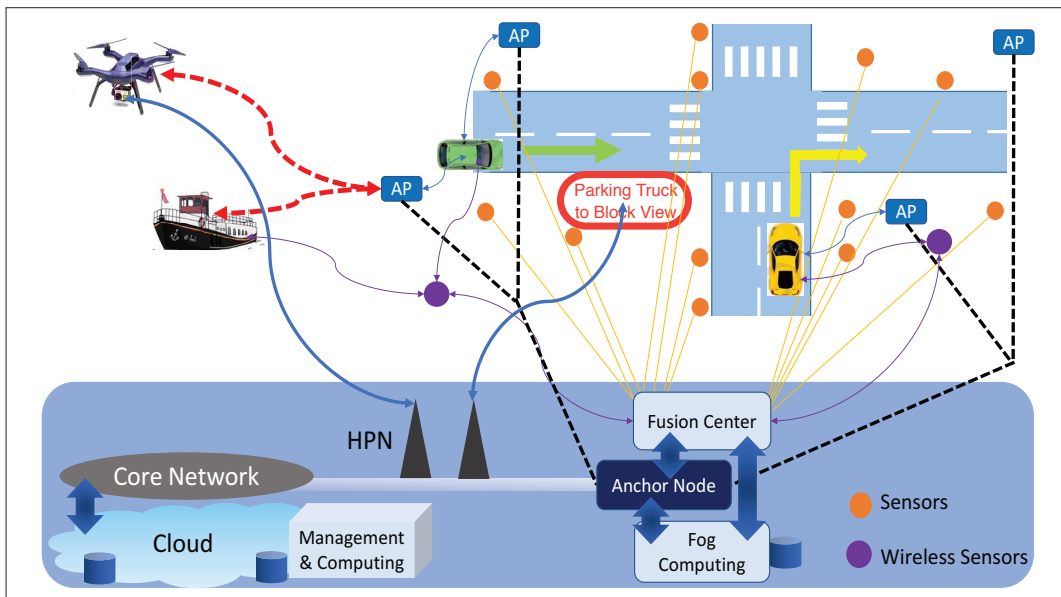


FIGURE 2. Heterogeneous network architecture and multi-scale computing environment for AVs, IMMs, and Smart City Applications (e.g., street security and safety, elderly cares, harbor and airport security, smart street light, parking, and so on). The yellow vehicle turns red is another scenario requiring vehicular networking by NHTSA.

al closed-loop means. Open-loop communication has another advantage to significantly reduce control signaling overhead compared to the closed loop for power control in 3G cellular and channel estimation in 4G.

ERROR-CONTROL OF OPEN-LOOP MULTI-PATH TRANSMISSIONS

In spite of the minimal latency by employing open loop wireless communication, error control is still required, because reliability is critical to ultra-low latency mobile networking applications. Let us design error control from the networking perspective, rather than traditional physical point-to-point link. The wireless networking unit (WNU) in a vehicle exchanges information with the fog computing through AP(s). For the uplink, WNU can send the packet via multiple paths, that is, multiple APs under networking architecture, to enhance reliability as the common technique in ad hoc networking. Uplink grant-free transmission is therefore facilitated. Similarly, fog computing can send a packet to the WNU of a vehicle via multiple APs in the downlink, while mobility management will be described later in this article. Open-loop multi-path transmissions could be reliable provided that at least one of the multi-path transmissions can correctly reach the destination.

The real situation might be more challenging since open-loop wireless communication may have to randomly select RRU(s) or a radio slice(s) to transmit packet(s) without centralized optimization of radio resource allocation, so as to realize the shortest possible scheduling. For example, as shown in Fig. 3, the orange vehicle (i.e., AV) is connected to three APs using the radio slices in orange color, to communication with the anchor node (AN). The abstract communication scenario is depicted in the right side of the figure by three networking paths of two-hops in each networking path. The red dot block can be analogous to a physical communication channel. Then, AV-AN multi-hop networking may correspond to a multi-input-multi-output (MIMO) channel.

To minimize the latency, the transmitter of open-loop wireless communication may randomly select RRUs, which favors large available spectrum bandwidth such as mmWave frequency bands. The optimal selection mechanism in general networking scenarios is subject to further study.

To minimize the latency, the transmitter of open-loop wireless communication may randomly select RRUs, which favors large available spectrum bandwidth such as mmWave frequency bands. The optimal selection mechanism in general networking scenarios is subject to further study. However, a side challenge in multiple access communication is the hit, a terminology borrowed from frequency hopping multiple access, if one or more RRUs are used by simultaneous transmissions within radio range. Therefore, the success of an open-loop transmission is stochastic, or more precisely, opportunistic. Due to potential hits, interference, fading, and high mobility (i.e., leaving earlier AP(s) in the downlink and arriving at new APs in the uplink), the transmitted packet in any path may be lost. We can model the mechanism of packet loss in each PHY transmission in an opportunistic way. The end-to-end error control between source and destination (i.e., vehicle and AN as Fig. 3) is therefore similar to MIMO communication error control over the network session. The links between AV and APs are actually opportunistic, and thus each networking path is opportunistic. This is fundamentally the same as error control of multi-hop cognitive radio networking. In addition to hybrid ARQ, the concept of space-time codes for the PHY transmission can be generalized to path-times codes (PTC) over source-destination of multi-path multi-hop networking [12]. Note that opportunistic selection of RRU(s) to transmit by open-loop communication may result in the hits, simultaneous utilization of the same frequency-time radio resource (or portion of the RRUs) from multiple physical trans-

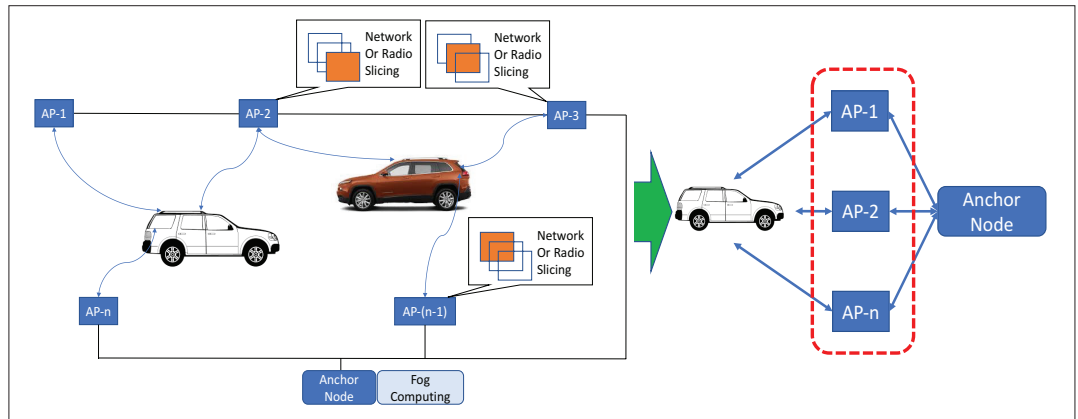


FIGURE 3. Multi-Path Networking between AV (white vehicle) and AN, Analog to MIMO PHY Transmission, and AV as the Center of Virtual Cell Cooperatively Served by APs. A single AV (orange vehicle) forms a virtual cell by network function virtualization (NFV) to use network slices (i.e., AP-2, AP-3, and AP-(n - 1)) and corresponding radio slices in orange color.

To realize ultra-low latency in each open-loop radio transmission, proactive network association is suggested, while an AV just selects appropriate APs in range for network association and thus no handover mechanism is needed anymore to significantly reduce networking latency among small cells.

missions from/to different AVs. Hits cause packet loss, particularly in the uplink since the WNU of an AV selects RRU(s) to proceed open-loop communication without full information of centralized scheduling or cognitive radio resource allocation [13]. Together with imperfect synchronization and such possible hits, the packet over such opportunistic link and thus path might be lost, while AN might realize the loss prior to decoding; sphere decoding or MAP decoding is therefore adopted in the receiver. As indicated in [13] and later radio slicing reception, successful reception of a packet can be reliably achieved by cooperative multi-path open-loop communication.

FAST RESOURCE ALLOCATION

Cognitive radio resource allocation leverages local sensing in scheduling [13], and such an allocation can directly instruct PHY to simplify the protocol stack toward ultra-low latency [14], but require rapid processing architecture for PHY and network operation [15]. By using the virtual cell concept in the next section, rapid scheduling of multiple objectives suggests technological opportunities in the processing architecture and programming.

VEHICLE-CENTRIC NETWORKING

In traditional mobile networks, the networking functions between an AV's WNU and network infrastructure is generally centrally controlled by the network, while complicated protocol stacks typically cost significant networking latency. Using open-loop PHY, a breakthrough idea is to treat each AV as a virtual cell [16]. That is, there is only one mobile station in each virtual cell and multiple APs cooperatively serve this mobile node using the coordinated multi-point transmission/reception [17]. With network virtualization, each AP is designated as a network and radio slice to this virtual cell, and simultaneously serves multiple virtual cells via other radio slices. Consequently, APs and subsequently the AN shall be facilitated by network virtualization in software defined networking

(SDN), while we temporally keep exact realization as future research in multi-core computing and parallel programming.

PROACTIVE NETWORK ASSOCIATION

To realize ultra-low latency in each open-loop radio transmission, proactive network association is therefore suggested, while an AV just selects appropriate APs in range for network association and thus no handover mechanism is needed anymore to significantly reduce networking latency among small cells. Figure 3 also depicts such scenario and the orange vehicle is the center of the virtual cell and communicating with three APs. To minimize networking latency, proactive network association allows the virtual cell (i.e., the orange vehicle) to select proper APs (i.e., network slice) to access and to proceed uplink transmission via selected radio slice. Both uplink and downlink transmissions are implemented like cooperative multi-point transmission/reception (CoMP) [18]. Each radio transmission from one vehicle to any AP is realized as open-loop communication, that is, no Acknowledgment for PHY transmission, to minimize communication latency and to forward packets to higher layers as fast as possible. In the downlink, cooperative communication similar to CoMP proceeds, while each AP allocates an appropriate radio slice to the virtual cell. AN is under the instruction of edge/fog computing and sends packets to those APs associated with the virtual cell and finally to the virtual cell, again by open-loop communication without feedback Acknowledgment.

For the sake of ultra-low networking latency, proactive network association is preferred among APs, which is known as horizontal association. However, a virtual cell might not be able to get the appropriate network slices (APs), due to radio resource constraints (for example, no radio slice physically available in AP(s)), no AP being deployed, or handover between two ANs. Though such a chance might be small, to ensure high reliability, the virtual cell must fall back to the high-power node (HPN) in the cellular network of heterogeneous networking architecture for such scenarios by *vertical association*, which can be viewed as a cellular V2X (C-V2X) network. This is ultimately related to the choice between cloud and edge networking, which determines the latency,

link reliability, and satisfaction of real-time computing needs for AVs and IMMs. Considering interference analysis by stochastic geometry, queuing analysis, and Lyapunov optimization theory, Fig. 4 shows that end-to-end networking latency toward 1 msec is realizable [19], where the simulations are taken from random deployed APs as the worst-case scenario and the horizontal axis represents the average number of APs in use over simulations.

NETWORK FUNCTION VIRTUALIZATION AND SOFTWARE DEFINED NETWORKS

In a practical implementation, a physical mobile heterogeneous network might be configured into several virtual networks, and one is reserved for ultra-low latency networking with dedicated PHY radio resources. Some innovative dedicated multi-core processors to serve such traffic of ultra-low latency would be expected as hardware of network function virtualization and software defined networks.

RADIO SLICING RECEPTION USING MULTIUSER DETECTION

Up to this point, most primary network functions of heterogeneous mobile network architecture have been successfully innovated, and could be compatible with any multiple access technology in the air-interface, though the above discussions imply OFDMA. However, there exists a major challenge to utilize open-loop wireless communication in multiple access of mobile communications, losing perfect synchronization or alignment in timing, frequency, and phase. Furthermore, once losing perfect synchronization with network, interference alignment is not possible and multiple-access interference (MAI) emerges, which means the receiver experiencing non-orthogonal multiple access even orthogonal multiple access is actually in use. Figure 5 illustrates such a technological challenge in the downlink.

Multiuser detection (MUD) is well known to overcome the challenge of MAI, including applications to asynchronous NOMA [20], and can be applied to multiuser synchronization [21] and receiver design of multi-path open-loop communication. The well known concern of applying MUD is its complexity exponentially growing proportional to the number of simultaneous users. In the downlink, this is not a major issue as a very limited number of APs are sending signals to the WNU of an AV due to the nature of virtual cell. For the uplink, AP is expected to equip short range radio, and the subsequent number of simultaneous users much less than that of a base station (HPN) in cellular networks. To illustrate by the case of two APs to vehicle, let A_1 , A_2 be mean signal powers at the vehicle, H_1 , H_2 be corresponding fading coefficients, X_1 , X_2 be transmitted symbols, respectively, and Y be the net received signal. The maximum-likelihood MUD demodulates the desired signals from MAI as

$$\left(\hat{X}_1, \hat{X}_2\right) = \arg \min_{(X_1, X_2) \in \mathcal{M}_1 \times \mathcal{M}_2} \left| Y - \sqrt{A_1} H_1 X_1 - \sqrt{A_2} H_2 X_2 \right|^2$$

where \mathcal{M}_1 , \mathcal{M}_2 denote signal constellations used at two APs, and can be different due to different SINRs and rate requests. By further leveraging multiple-antenna receive diversity, detection performance can be near-far resistant and approach the theoretical limit even for high-order modulated signal in 64-QAM. Similarly, each AP can still apply MUD and finally combine the packets at the anchor node, with more possible techniques such as HARQ, Raptor

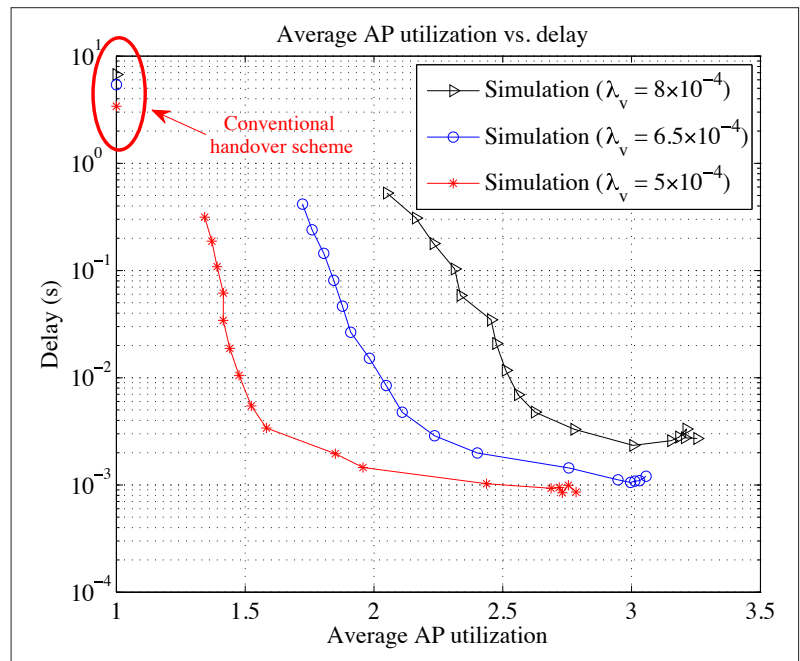


FIGURE 4. Networking latency vs. average number of APs in use (IV: density of APs per m^2) [19].

codes, and so on, to be explored. Multiple access can be successfully accomplished by applying MUD and earlier error control, while losing orthogonality in each open-loop radio transmission.

ANTICIPATORY MOBILITY MANAGEMENT VIA MACHINE LEARNING

Virtual cell using proactive network association and open-loop communication forms the core technology to accomplish ultra-low latency mobile networking. However, a significant technological challenge has hardly been considered in the literature. Conventional mobility management in cellular networks is handled by network infrastructure in a rather centralized manner. By proactive network association with a virtual cell frequently accessing new APs due to high mobility and small cell structure [22], the network infrastructure, including RAN, core network, and AN/APs, does not know the APs that the virtual cell is going to associate with in the next time instant. Uplink communication is proactive, but the downlink communication typically with critical control information will encounter this technological challenge. We shall develop the methodology that AN of fog computing determines potential APs in the next time instant to serve a virtual cell in the downlink. Even more challenging, the entire function must be completed in almost real-time. The only way is to execute vehicular data analytics inside fog such that the mobility management directs downlink packets via those APs anticipating connections from the virtual cell.

MACHINE LEARNING ON BIG VEHICULAR DATA

More precisely, the success of ultra-low latency networking must rely on AN's anticipatory mobile management (AMM) that predicts likely APs for each mobile node to proactively connect in the next time instant, given ultra-low latency. This is a fundamentally new challenge different from mobility pattern

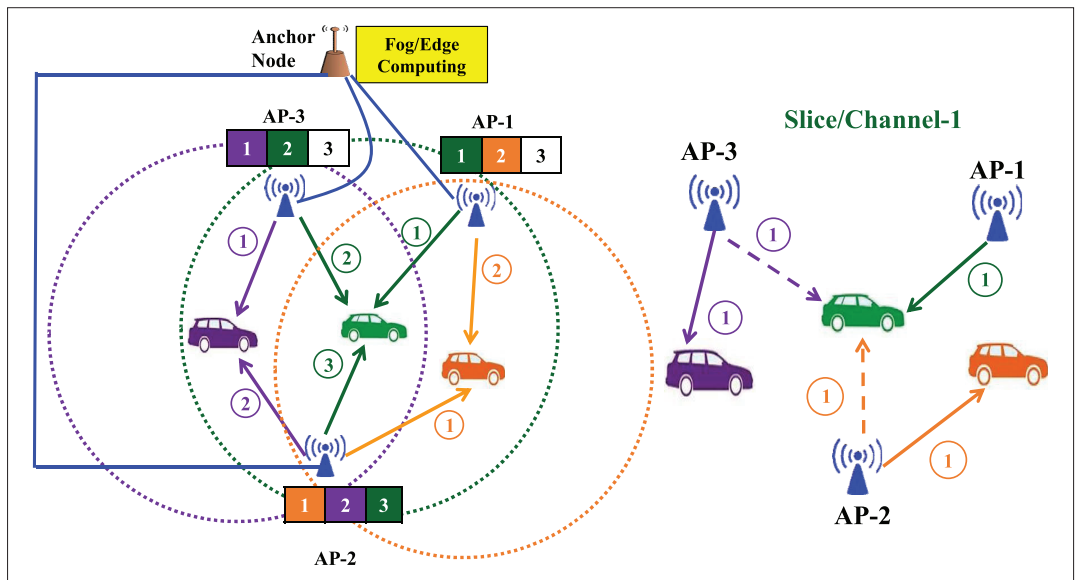


FIGURE 5. (Left) Each AP supporting the virtual cell allocates a radio slice to an individual AV. (Right) Physical co-channel multiple access interference in slice/channel-1, owing to lacking coordination among open-loop transmissions.

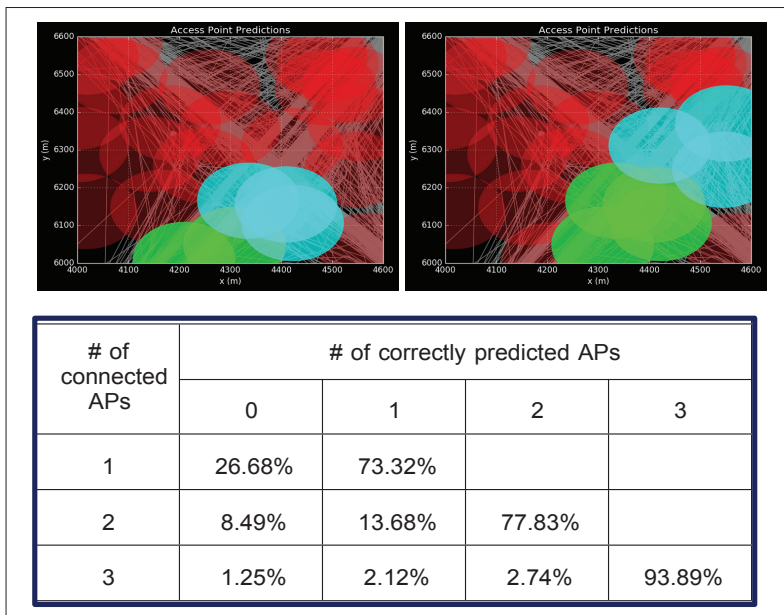


FIGURE 6. (Top) An illustration of anticipatory mobility management, where green circles represent the coverage of APs in use and blue circles represent the coverage of predicted APs to be used in the next time instant. Upper left to upper right illustrates a time trace. (Bottom) Prediction accuracy aided by machine learning and big data analytics, conditioned on the number of connected APs.

prediction, and machine learning on big vehicular data appears an attractive approach. Though applying machine learning to enhance the performance of wireless networks has attracted recent research interest [23], AMM might be the very first effort to develop wireless networking *relying on machine learning* and big data analytics. The unique aspect of machine learning to realize AMM is that AMM or AN only knows past traces of vehicles (i.e., the record that different APs are used by past vehicles) and all other information such as GPS data must be carefully utilized due to privacy and security concerns, as illustrated in the top part of Fig. 6. AMM

typically handled by fog computing and AN should learn from past data of earlier passing vehicles to predict anticipated APs in the next time instant based on APs being used for this AV. Since such predictive analytics must be done in almost real-time, popular deep learning neural networks might not directly fit the purpose. Instead, the Naïve Bayesian prediction and by recursive Bayesian estimation of low complexity with the aid of location-based information (i.e., map) can be adopted through the hidden Markov model [24]. A taxi dataset consisting of 12,000 vehicles running over one month has been considered for evaluation of methods. The APs are randomly deployed over the region, which represents the worst-case scenario with possible service holes existing. Indeed, recurrent neural networks cannot deliver satisfactory performance. Location-assisted Naive Bayesian estimation with velocity estimation gives satisfactory results as shown in the table in Fig. 6, where a higher failure rate for only one AP can be easily improved by well-planned deployment of APs by referencing the street map.

AMM in the uplink communication is straightforward since an AV just connects to APs in range via proactive network association. However, for ultra-low latency packets in the downlink, AMM must predict APs to be connected by the vehicle in the next time instant. Both uplink and downlink require multi-path error control as [12] together with further error control techniques. The AMM consists of prediction by machine learning by the anchor node in the fog, and falling back mechanism to HPN(s) in the heterogeneous network if no AP is successfully predicted. The extremely simple AMM operation except challenging prediction warrants ultra-low latency mobile networking.

Security Consideration: Security concerns may arise due to distributed control of network functioning and open-loop wireless communications, since traditional network security is accomplished by sophisticated cryptography and secure protocols. However, a different view of network security has been recently developed by treating the dynamics of the network to enhance security.

Proactive network association with randomized selection of RRUs for the uplink and fog predictive cooperative communication via multiple APs that only fog knows in the downlink suggest further dynamic randomness for attackers.

CONCLUSION REMARKS

To ultimately achieve ultra-low latency mobile networking, a collection of technologies that quite differ from conventional design of communication systems and networks have been presented across network layers in this article. By thinking out of the box and integrating views of computing and networking, it suggests further frontiers of wireless networking engineering knowledge.

REFERENCES

- [1] R. Ford et al., "Achieving Ultra-Low Latency in 5G Millimeter Wave Cellular Networks," *IEEE Commun. Mag.*, vol. 55, no. 3, Mar. 2017, pp. 196–203.
- [2] G. P. Fettweis, "The Tactile Internet: Applications and Challenges," *IEEE Vehicular Technology Mag.*, vol. 9, no. 1, Mar. 2014, pp. 64–70.
- [3] P. Gupta and P. Kumar, "The Capacity of Wireless Networks," *IEEE Trans. Information Theory*, vol. 46, no. 2, Mar. 2000, pp. 388–404.
- [4] D.-J. Deng et al., "Latency Control in Software-Defined Mobile-Edge Vehicular Networking," *IEEE Commun. Mag.*, vol. 55, no. 8, Aug. 2017, pp. 87–93.
- [5] K. Zhang et al., "Reliable and Efficient Autonomous Driving: The Need for Heterogeneous Vehicular Networks," *IEEE Commun. Mag.*, Dec. 2015, pp. 72–79.
- [6] M. Satyanarayanan et al., "An Open Ecosystem for Mobile-Cloud Convergence," *IEEE Commun. Mag.*, vol. 53, no. 3, Mar. 2015.
- [7] M. Chiang and T. Zhang, "Fog and IoT," *IEEE Internet of Things J.*, vol. 3, no. 6, Dec. 2016, pp. 854–64.
- [8] S.-C. Hung, S.-Y. Lien, and K.-C. Chen, "Low Latency Communication for Internet of Things," *IEEE Int'l. Conf. Communications in China*, 2015.
- [9] S.-Y. Lien, S.-C. Hung, and K.-C. Chen, "Optimal Radio Access for Fully Packet Switching 5G Networks," *IEEE Int'l. Conf. Commun.*, 2015.
- [10] K.-C. Chen, "Medium Access Control of Wireless Local Area Networks for Mobile Computing," *IEEE Networks*, vol. 8, no. 5, Sept. 1994, pp. 50–63.
- [11] C. She, C. Yang, and T. Q. S. Quek, "Radio Resource Management for Ultra-Reliable and Low-Latency Communications," *IEEE Commun. Mag.*, vol. 55, no. 6, June 2017, pp. 72–78.
- [12] I.W. Lai et al., "End-to-End Virtual MIMO Transmission in Ad Hoc Cognitive Radio Networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 1, Jan. 2014, pp. 330–41.
- [13] S.-Y. Lien et al., "Cognitive Radio Resource Management for Future Cellular Networks," *IEEE Wireless Commun.*, vol. 21, no. 1, Feb. 2014, pp. 70–79.
- [14] S.-Y. Lien et al., "Ultra-Low Latency Ubiquitous Connections in Heterogeneous Cloud Radio Access Networks," *IEEE Wireless Commun.*, vol. 22, no. 3, June 2015, pp. 22–31.
- [15] S. Haas et al., "Heterogeneous SDR MPSoC in 28nm CMOS for Low Latency Wireless Applications," *IEEE Design Automation Conf.*, 2017.
- [16] A. Behnad and X. Wang, "Virtual Small Cell Formation in 5G Networks," *IEEE Commun. Lett.*, vol. 21, no. 3, Mar. 2017, pp. 616–19.
- [17] J. Lee et al., "Coordinated Multipoint Transmission and Reception in LTE-advanced Systems," *IEEE Commun. Mag.*, vol. 50, no. 11, Nov. 2012, pp. 44–50.
- [18] P. Marsch, S. Khattak, and G. Fettweis, "A Framework for Determining Realistic Capacity Bounds for Distributed Antenna Systems," *Proc. IEEE Information Theory Workshop*, 2006.
- [19] S.-C. Hung et al., "Delay Guaranteed Network Association for Mobile Machines in Heterogeneous Cloud Radio Access Networks," *IEEE Trans. Mobile Computing*, vol. 17, no. 12, pp. 2744–60, Dec. 2018.
- [20] J. Cui et al., "Asynchronous NOMA for Downlink Transmission," *IEEE Commun. Lett.*, vol. 21, no. 2, pp. 402–5, Feb. 2017.

The AMM consists of prediction by machine learning by the anchor node in the fog, and falling back mechanism to HPN(s) in the heterogeneous network if no AP is successfully predicted. The extremely simple AMM operation except challenging prediction warrants ultra-low latency mobile networking.

- [21] C.-M. Chang and K.-C. Chen, "Joint Linear Timing and Carrier Phase Estimation of DS-CDMA Multiuser Communications," *IEEE JSAC*, vol. 18, no. 11, Jan. 2000, pp. 87–98.
- [22] A. J. Fehske et al., "Small-Cell Self-Organizing Wireless Networks," *Proc. IEEE*, vol. 102, no. 3, Mar. 2014, pp. 334–50.
- [23] C. Jiang et al., "Machine Learning Paradigms for Next-Generation Wireless Networks," *IEEE Wireless Commun.*, vol. 24, no. 2, Apr. 2017, pp. 98–105.
- [24] C.-Y. Lin et al., "Anticipatory Mobility Management by Big Data Analytics for Ultra-Low Latency Mobile Networking," *Proc. IEEE Int'l. Conf. Commun.*, 2018.

BIOGRAPHY

KWANG-CHENG CHEN (email: kwangcheng@usf.edu) is a professor with the Department of Electrical Engineering, University of South Florida, after an academic career in Taiwan and an industrial career in the U.S. He has contributed essential technology to various IEEE 802, Bluetooth, and LTE and LTE-A wireless standards. In addition to service with IEEE journals and conferences, he founded and then chaired the TC on Social Networks of the IEEE Communications Society. He is an IEEE Fellow and has received a number of awards, such as the 2011 IEEE ComSoc WTC Recognition Award, the 2014 IEEE Jack Neubauer Memorial Award, and the 2014 IEEE ComSoc AP Outstanding Paper Award. His recent research interests include wireless networks, artificial intelligence and machine learning, IoT/CPS, social networks and data analytics, and cybersecurity.

TAO ZHANG has been in various technical and executive positions leading strategies, research, and product development, over 30 years. He helped build Cisco's first Smart Connected Vehicles Business Unit and served as its CTO/Chief Scientist. He was the Chief Scientist for vehicular networking at Telcordia Technologies (formerly Bell Communications Research). Tao is Chair Professor at the National Chiao Tung University. He co-founded the OpenFog Consortium that has been leading fog computing. He co-founded and was a founding Board Director for the Connected Vehicle Trade Association (CVTA) in the US. Tao is an IEEE Fellow. He served as the CIO of the IEEE Communications Society in 2016 and 2017. He holds over 50 US patents and has co-authored two books.

RICHARD D. GITLIN is a State of Florida 21st Century World Class Scholar, a Distinguished Professor and the Agere endowed chair at the University of South Florida. He received a Doctor of Science from Columbia University and was at Bell Labs/Lucent Technologies for 32 years, including serving as Senior Vice President for Communications and Networking. He holds 65 US patents, is responsible for the co-invention of the digital subscriber line (DSL), and is a member of the National Academy of Engineering (NAE), a Fellow of the IEEE, a Bell Laboratories Fellow, a Charter Fellow of the National Academy of Inventors (NAI) and a 2017 inductee in the Florida Inventors Hall of Fame. He co-authored a seminal textbook in electrical engineering, and he has published more than 150 papers.

GERHARD P. FETTWEIS has been Vodafone Chair Professor at TU Dresden since 1994, and heads the Barkhausen Institute since 2018, respectively. He earned his Ph.D. under H. Meyr's supervision from RWTH Aachen in 1990. After one year at IBM Research in San Jose, CA, he moved to TCSI Inc., Berkeley, CA. He coordinates the 5G Lab Germany, and two German Science Foundation (DFG) centers at TU Dresden, namely cfaed and HAEC. Gerhard is an IEEE Fellow, a member of the German Academy of Sciences (Leopoldina), the German Academy of Engineering (acatech), and he has received multiple IEEE recognitions as well as the VDE ring of honor. He co-chairs the IEEE 5G Initiative, and has helped organize IEEE conferences, most notably as TPC Chair of ICC 2009 and of TTM 2012, and as General Chair of VTC Spring 2013 and DATE 2014. His research focusses on wireless transmission and chip design for wireless/IoT platforms.