

Joint Uplink and Downlink Resource Configuration for Ultra-Reliable and Low-Latency Communications

Changyang She^{ID}, *Member, IEEE*, Chenyang Yang^{ID}, *Senior Member, IEEE*,
and Tony Q. S. Quek^{ID}, *Fellow, IEEE*

Abstract—Supporting ultra-reliable and low-latency communications (URLLC) is one of the major goals for the fifth-generation cellular networks. Since spectrum usage efficiency is always a concern, and large bandwidth is required for ensuring stringent quality-of-service (QoS), we minimize the total bandwidth under the QoS constraints of URLLC. We first propose a packet delivery mechanism for URLLC. To reduce the required bandwidth for ensuring queueing delay, we consider a statistical multiplexing queueing mode, where the packets to be sent to different devices are waiting in one queue at the base station, and broadcast mode is adopted in downlink transmission. In this way, downlink bandwidth is shared among packets of multiple devices. In uplink transmission, orthogonal subchannels are allocated to different devices to avoid strong interference. Then, we jointly optimize uplink and downlink bandwidth configuration and delay components to minimize the total bandwidth required to guarantee the overall packet loss and end-to-end delay, which includes uplink and downlink transmission delays, queueing delay, and backhaul delay. We propose a two-step method to find the optimal solution. Simulation and numerical results validate our analysis and show remarkable performance gain by jointly optimizing uplink and downlink configuration.

Index Terms—Ultra-reliable and low-latency communications, resource configuration, packet delivery mechanism.

Manuscript received August 2, 2017; revised November 14, 2017 and January 3, 2018; accepted January 3, 2018. The work of C. She was supported in part by the National Natural Science Foundation of China under Grant 61731002, in part by the SUTD-ZJU Research Collaboration under Grant SUTDZJU/RES/01/2016, and in part by the MOE ARF Tier 2 under Grant MOE2015-T2-2-104. The work of C. Yang was supported by the National Natural Science Foundation of China under Grant 61731002. The work of T. Q. S. Quek was supported in part by the SUTD-ZJU Research Collaboration under Grant SUTDZJU/RES/01/2016, and in part by the MOE ARF Tier 2 under Grant MOE2015-T2-2-104. This paper was presented in part at the workshop on ultra-reliable and low-latency communications in wireless networks with IEEE Global Communications Conference 2016 [1]. The associate editor coordinating the review of this paper and approving it for publication was V. Aggarwal. (*Corresponding author: Changyang She.*)

C. She was with the School of Electronics and Information Engineering, Beihang University, Beijing 100191, China. He is now with the Information Systems Technology and Design Pillar, Singapore University of Technology and Design, Singapore 487372 (e-mail: shechangyang@gmail.com).

C. Yang is with the School of Electronics and Information Engineering, Beihang University, Beijing 100191, China (e-mail: cyyang@buaa.edu.cn).

T. Q. S. Quek is with the Information Systems Technology and Design Pillar, Singapore University of Technology and Design, Singapore 487372 (e-mail: tonyquek@sutd.edu.sg).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCOMM.2018.2791598

I. INTRODUCTION

ULTRA-RELIABLE and low-latency communications (URLLC) are required in the emerging application scenarios of the fifth-generation (5G) cellular networks [2]. Different from existing communication systems that are designed for human-to-human (H2H) communications, URLLC target to ultra-reliable machine-type-communications and human-to-machine communications that require haptic interactions, such as autonomous vehicles, factory automation, and remote control [3], [4]. All those applications have strict requirements on end-to-end (E2E) or round trip delay (say around 1 ms) and reliability (say around 10^{-6} packet loss probability), which can not be satisfied in Long Term Evolution systems.

The E2E delay consists of various delay components that depend on communication scenarios. In long distance communication scenarios, the E2E delay consists of transmission delay and queueing delay in radio access network, routing delay in backhaul and core networks, and also propagation delay that is hard to control. For example, when the communication distance is longer than 300 km, the propagation delay is longer than 1 ms since light travels 300 km per millisecond in vacuum [5]. Nevertheless, it is worth noting that ensuring the ultra-short E2E delay is not easy even in local communication scenarios, where communication is only required by devices in adjacent cells with short backhaul delay and negligible propagation delay.

A core difference between radio resource allocation for URLLC and that for traditional real-time service comes from both transmission delay and packet size [2]. In H2H communications, transmission delay is relative long (say 10 ms) and the packet size is large (say 1500 bytes) [6]. As a result, Shannon's Capacity is widely applied in existing literatures to characterize achievable rate of traditional services with long packets (e.g., [7] and references therein). In URLLC, to satisfy short transmission delay, short packets are transmitted (say 20 bytes [2]). As a result, the blocklength of channel coding is short. In addition, to ensure ultra-high reliability, decoding error with short blocklength channel codes cannot be ignored. Therefore, decoding error probability with short blocklength channel codes should be applied [8], which is with very complicated expression. Fortunately, approximate decoding error probability in finite blocklength regime has been obtained with

simple expressions in [9] and [10], which are shown accurate for quasi-static fading channels. Yet these approximations are neither convex nor concave in transmit power or bandwidth. As a result, the resource allocation optimization for URLLC is much more challenging than H2H communications.

Similar to H2H communications, ensuring short queueing delay is also necessary for URLLC, where queueing delay requirement should be characterized by the queueing delay bound and its violation probability. Considering that packets are randomly generated and the service rate of a wireless link could be random, queueing delay has been considered in single-user scenarios [11]–[13] and multi-user scenarios [14], where achievable rate in finite blocklength regime was applied in their analyses. A packet scheduling policy was proposed under strict delay bound constraint on queueing delay in [11], which cannot be satisfied with probability one due to channel fading. To show when the delay bound can be satisfied, a feasible condition was obtained, but delay bound violation probability can not be derived under the framework in [11]. To analyze queueing delay bound violation probability, network calculus was applied to obtain an upper bound of the delay violation probability in [12]. Simulation results in [12] validated that if Shannon's Capacity is applied, then the delay violation probability will be underestimated, and hence the quality-of-service (QoS) cannot be guaranteed. The performance of relay systems was analyzed in [13], where effective capacity was applied to characterize the queueing delay. More recently, the transmission policy in both single-user and multi-user scenarios was optimized to minimize the maximal transmit power required to satisfy QoS requirement in [14]. To study how to serve multiple users with multiple BSs, a signal-to-interference-plus-noise ratio (SINR) model was applied to simplify the reliability requirement in [15], where multi-connectivity was exploited to improve reliability.

The study in [14] focuses on downlink (DL) transmission design, and implicitly assume that the uplink (UL) transmission can be finished in a short time with guaranteed reliability. However, ensuring ultra-reliable and low-latency for UL transmission is not easy as well. To ensure the QoS for URLLC, UL and DL resource allocation should be jointly optimized [16]. To guarantee queueing delay violation probability, effective bandwidth and effective capacity were applied in [16], where the Shannon's Capacity was used as the service rate (and hence decoding error probability was not considered), and the global optimal solutions of the problem was not found. If the achievable rate in finite blocklength regime is applied in the joint UL and DL resource allocation, it could be more challenging to find the optimal solution [17].

On the other hand, spectrum is scarce resource for wireless communications. Due to stringent QoS requirement, the resource allocation for URLLC is inevitably conservative, and hence the bandwidth usage efficiency for URLLC will be very low without judicious control. In some application scenarios, the overall bandwidth can become unaffordable. For example, the packet arrival rate in tactile internet is very high [18], and in smart factory the number of devices can be very large [19]. In order to improve the spectrum usage efficiency, the packet delivery mechanism for URLLC,

including queueing and transmission modes, should be reconsidered. To guarantee the QoS of each user, the packets to different users are waiting in different queues (i.e., individual queueing mode) in [14] and [16]. However, to achieve the same average queueing delay, the required service rate of the individual queueing mode is much higher than that of a statistical multiplexing queueing mode, where the packets to different users are waiting in one common queue, and hence the DL bandwidth can be shared among multiple users [20]. Besides, most of existing studies assumed that channel state information at transmitter (CSIT) is perfectly known, which incurs training/feedback overhead that linearly increases with the number of receivers [11], [12], [14], [16]. Since the packet size in URLLC is usually very small, large signaling overhead leads to low spectrum efficiency. With a transmission mechanism without CSIT, the overhead can be reduced, but whether the QoS requirement can be guaranteed becomes a question.

In this paper, we study how to jointly optimize UL and DL resource configuration and delay components to minimize the total bandwidth to support URLLC in local communication scenarios. We focus on orthogonal multiple access systems to avoid strong interference. The major contributions are summarized as follows:

- We propose a packet delivery mechanism for URLLC. To save bandwidth in UL transmission, bandwidth is not reserved for all the UL devices since they may stay dumb for a long time. Only the devices that request sending packets will be assigned bandwidth. To reduce the bandwidth for ensuring queueing delay requirement, statistical multiplex queueing mode is adopted, where the packets to different devices are waiting in one queue at the buffer of the base station (BS). By taking Poisson arrival process as an example, we prove that under the same queueing delay bound and queueing delay violation probability, the required service rate of a statistical multiplex queue is less than the sum of the required service rates of all individual queues, where the packets to different users are waiting in different queues. To reduce overhead, broadcast is applied for DL transmission.
- We jointly optimize the bandwidth assignment for UL and DL transmissions and delay components to minimize the total bandwidth required by the packet delivery mechanism to ensure the E2E delay and overall packet loss probability, where the achievable rate in finite blocklength regime is applied, and routing and propagation delays are characterized by a deterministic backhaul delay. The E2E delay includes UL and DL transmission delays, queueing delay, and backhaul delay. The overall packet loss includes packet loss in UL and DL transmissions and queueing delay violation. A two-step method is proposed, where the bandwidth assignment is first optimized with given delay components and then the uplink and downlink bandwidth are optimized given the E2E delay. Numerical results show that the joint configuration requires half of the total bandwidth of the non-joint optimization.

The remainder of this paper is organized as follows. System model is described in Section II. A packet delivery mechanism

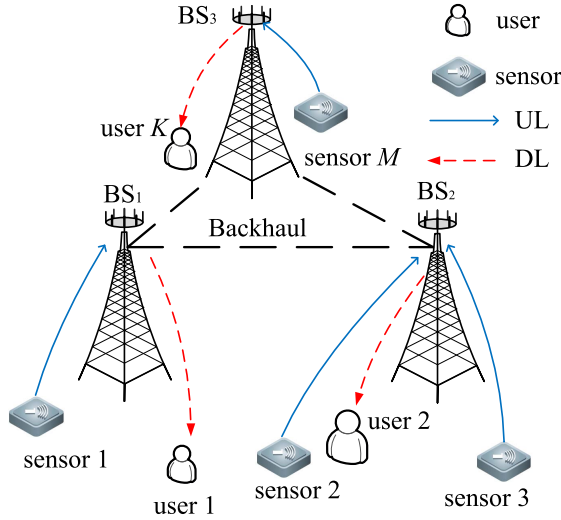


Fig. 1. Illustration of a local communication scenario.

is proposed in Section III. Section IV formulates the optimization problem. Section V shows how to obtain the optimal solution. Simulation and numerical results are provided in Section VI. Section VII concludes this work.

II. SYSTEM MODEL

A. Local Communication Scenarios

We consider local communication scenarios, where BSs are connected with one-hop backhaul and the communication distance is less than a few kilometers. In such scenarios, propagation delay is negligible. By deploying high-capacity backhaul links such as fiber, the backhaul delay is around 0.1 ms [21].

As illustrated in Fig. 1, we consider a cellular system with $K + M$ single-antenna devices and three BSs. Each device is served by one of the BSs, which are equipped with N_t antennas. To avoid strong interference among adjacent BSs, the frequency-reuse factor of the network is assumed to be 1/3. Frequency-division multiple access is applied to avoid interference among different devices. There are two kinds of devices. The first kind of devices are K users, which need to download packets from the BSs. The second kind of devices are M sensors, which generate and upload packets to the BSs. After receiving the packets successfully, the BSs send the packets to users.¹ Depending on application scenarios, each user may require packets from one or more sensors in the M sensors, and the packets of one sensor may be required to one or more target users in the K users. If a sensor and a target user are connected to two BSs, then the required packets need to be forwarded from the BS connected with the sensor to the BS connected with the user via backhaul. Frequency division duplex (FDD) systems is considered, because they are widely deployed. Our studies can be easily extended to time division duplex (TDD) systems. In TDD systems, one can adjust the

¹Device-to-device communications are possible for URLLC, which will not be addressed in this work.

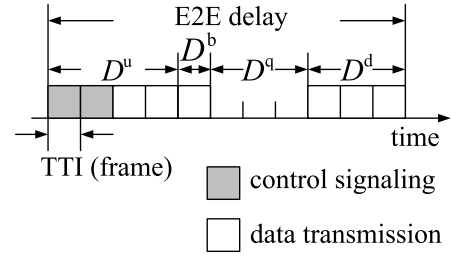


Fig. 2. Frame duration and delay components of event-driven packets.

ratio of UL and DL transmission durations, which is equivalent to adjusting UL and DL bandwidth in FDD systems.

B. QoS Requirement

The QoS requirement of URLLC is characterized by an ultra-short E2E delay D_{\max} and an overall packet loss probability ε_{\max} that are imposed on each packet. For local communication scenario, the E2E delay includes UL and DL transmission delays, queueing delay at the BSs and backhaul delay by assuming negligible propagation delay and processing delay. The overall packet loss comes from decoding errors and queueing delay violation. In Long Term Evolution systems, transmission time interval (TTI) is 1 ms, and hence the E2E delay of URLLC cannot be supported. To reduce transmission delay, we consider the short frame structure as illustrated in Fig. 2, where TTI equals to the frame duration T_f , which is the minimal time granularity of the system (i.e., subframe in [22]). Therefore, the transmission delays and queueing delay should be divisible by frame duration. The E2E delay requirement is given by

$$D^u + D^d + D^q + D^b \leq D_{\max}, \quad (1)$$

where D^u , D^d , D^q and D^b are the UL and DL transmission delays, queueing delay and backhaul delay, respectively.

In application scenarios of URLLC such as smart factory and vehicle networks, there are two kinds of packets: event-driven packets and periodic packets [23], [24]. For the event-driven packets, the transmission delay includes those caused by control signaling and data transmission as illustrated in Fig. 2. The UL transmission procedure includes the following steps: (i) generation of a packet by a sensor; (ii) UL scheduling request from the sensor; (iii) bandwidth assignment and transmission grant by the BS; (iv) UL data transmission. The scheduling requests are triggered by some urgent events, and the BSs need to grant the transmission immediately when a request is received. To ensure that the UL scheduling request can be successfully received by the BS (i.e., to avoid scheduling collision), control channel should be reserved for each sensor with event-driven packets [22], and hence only two frames are occupied by control signaling in steps (ii) and (iii). For periodic packets, the inter-arrival time between packets are known at the BSs. By reserving data transmission resource to each sensor, there is no control signaling.²

²Since random access could cause long latency for machine-type communications [25], it will not be used for both kinds of packets in URLLC [2].

Denote the requirements on the packet loss probabilities in UL transmission, DL transmission and queueing as ε^u , ε^d and ε^q , respectively. The ultra-high reliability can be ensured if

$$(1 - \varepsilon^u)(1 - \varepsilon^d)(1 - \varepsilon^q) \approx 1 - \varepsilon^u - \varepsilon^d - \varepsilon^q \leq \varepsilon_{\max}, \quad (2)$$

where the approximation is accurate since ε^u , ε^d , and ε^q are extremely small.

III. PACKET DELIVERY MECHANISM

In this section, we propose a packet delivery mechanism for event-driven packets in URLLC, which can be easily extended to periodic packets since their arrival processes are deterministic. To reduce the total bandwidth required to ensure QoS requirement, we consider statistical multiplexing mode for queueing at BSs, broadcasting mode for downlink transmission, and bandwidth assignment for UL and DL transmissions.

A. Queueing Mode

As shown in [26] and [27], the packet arrival process in vehicle networks and some M2M communications can be modeled as a Poisson process, which is an aggregation of packets generated by multiple sensors. Denote the average arrival rate of the Poisson process as λ packets/frame. The arrival process of each sensor is modeled as Bernoulli process. Denote $a_m(n)$ as the number of packets arrived at a BS from the m th sensor in the n th frame, $m = 1, \dots, M$. With probability p_m , $a_m(n) = 1$, and with probability $1 - p_m$, $a_m(n) = 0$. Then, the average total arrival rate of the M sensors is $\lambda = \sum_{m=1}^M p_m$ packets/frame. According to the result in [14], the effective bandwidth of the arrival process can be expressed as follows,

$$E_B = \frac{T_f \ln(1/\varepsilon^q)}{D^q \ln \left[\frac{T_f \ln(1/\varepsilon^q)}{\lambda D^q} + 1 \right]} \text{ packets/frame}, \quad (3)$$

which is the minimal constant packet service rate required to ensure queueing delay D^q and queueing delay violation probability ε^q . It is widely believed that effective bandwidth is applicable when the queue length or queueing delay is long. However, the results in [28] imply that for Poisson process and the arrival processes that are more bursty than Poisson process, a short delay requirement (D^q, ε^q) can be satisfied with a constant packet service rate that is equal to or higher than E_B . This implication is validated in [14] for typical arrival processes in URLLC, such as Poisson process, interrupted Poisson processes that is more bursty than Poisson process, and Switched Poisson process that is autocorrelated.

In existing papers that study how to ensure queueing delay requirement of each users, the packets target to different users wait in different queues before DL transmission as in Fig. 3(a) [14], [16]. Such a queueing mode is referred to as *individual queueing mode*. It has been shown that if $\tilde{\lambda} = \lambda/L$ and $\tilde{E}_B = E_B/L$, then the average queueing delay with the individual queueing mode will be L times of that with the *statistical multiplexing mode* in Fig. 3(b) [20]. The following proposition indicates that a similar conclusion can

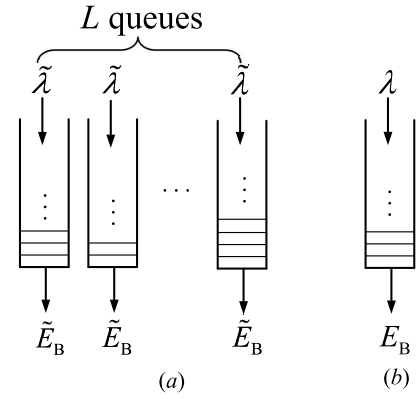


Fig. 3. Queueing modes. (a) Individual queueing mode. (b) Statistical multiplexing queueing mode.

be obtained when the delay requirement is characterized by (D^q, ε^q), which is imposed on each packet.

Proposition 1: Given the requirement on D^q and ε^q , if $\tilde{\lambda} = \lambda/L$, then $L\tilde{E}_B > E_B$.

Proof: See proof in Appendix A. \square

If the queueing delay requirement (D^q, ε^q) can be satisfied with the statistical multiplexing queue, then for any packet that comes from any of the M sensors, the probability that the queueing delay of the packet exceeds D^q is smaller than ε^q . In other words, no matter which sensor a packet came from, the delay requirement of it can be satisfied with the statistical multiplexing queue. Proposition 1 indicates that to guarantee the queueing delay requirement imposed on each packet, the required effective bandwidth of statistical multiplexing queue is less than the total effective bandwidth of individual queues. Since effective bandwidth is the minimal constant service rate that can ensure queueing delay requirement, and the required bandwidth decreases with service rate, we consider statistical multiplexing mode for saving bandwidth.

B. Transmission Mode and Bandwidth Assignment

In typical scenarios of URLLC, the channel coherence time is longer than the typical E2E delay (i.e., 1 ms) [14]. Then, the channel gain changes little before the deadline of conveying each packet, and hence simply retransmitting a packet in multiple consecutive frames can hardly improve the reliability. When channel is in deep fading, frequency diversity can be applied. To illustrate how the overall reliability can be ensured with diversity, we consider a simple method that transmits each packet multiple times over multiple separated subchannels [29].³

1) *Subchannel Assignment for UL Transmission:* Considering that a sensor may stay dumb for a long duration between the transmissions of short packets [26], the bandwidth is only

³For saving bandwidth, a simple idea is that all sensors transmit packet under the same spectrum. In the spectrum sharing network, interference is random and strong, which cannot be treated as additive noise or simply ignored. Unfortunately, the achievable rate with finite blocklength codes in the interference environment is unavailable in existing literatures. Therefore, whether or not the QoS requirement of URLLC can be satisfied in spectrum sharing systems is unknown and deserves further study.

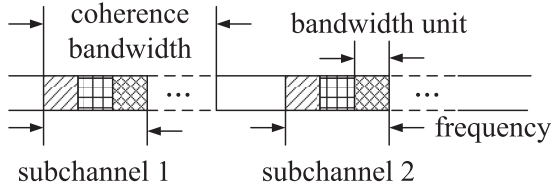


Fig. 4. Bandwidth of each subchannel.

assigned to the active sensors immediately after receiving the scheduling requests.

To exploit frequency diversity, the BS assigns N_m^u separated subchannels to the m th sensor if it has a packet to transmit. The packet is repeatedly transmitted over the N_m^u subchannels. If one of the N_m^u transmissions is successful, then the packet is successfully received at the BS. Since the interference among sensors causes severe deterioration in QoS, we assume that different subchannels are assigned to the sensors that request transmissions concurrently. To maximize frequency diversity gain, the instantaneous channel gains on the N_m^u subchannels assigned to transmit one packet should be independent, i.e., the separation of the N_m^u subchannels should exceed the channel coherence bandwidth W_c , as shown in Fig. 4. In real-world systems, frequency is discretized into basic bandwidth units, e.g., subcarriers in orthogonal frequency division multiple access systems, and then each subchannel consists of multiple bandwidth units with bandwidth B_0 . Denote B_m^u as the bandwidth of each subchannel allocated to the m th sensor. Then, B_m^u is divisible by B_0 . By adjusting the number of bandwidth units in one subchannel, the bandwidth of each subchannel can be controlled. The bandwidth assigned to the sensor for transmitting a packet is $N_m^u B_m^u$. We assume that $B_m^u < W_c$, such that each subchannel is frequency-flat fading.

2) *Subchannel Assignment for DL Transmission:* We consider broadcast for DL transmission.⁴ Without access control, acknowledgment feedback, and CSIT, the control and training/feedback overhead is negligible. To guarantee reliability, each packet is repeatedly transmitted over N^d subchannels each with bandwidth $B^d < W_c$, and different packets are transmitted over different subchannels. Each user can receive the signals on all the DL subchannels. We assume that the channel coding on each subchannel is independent of the others. With independent channel coding, decoding errors on different subchannels are uncorrelated. If one packet is lost due to decoding error, other packets can still be decoded successfully.

C. Illustration of the Packets Delivery Mechanism

The proposed packets delivery mechanism is illustrated in Fig. 5, where sensors 1, 2, and M are served by BSs 1, 2 and 3, respectively. In the considered time slots, sensor M has no packet to transmit and stays dumb, and the other

⁴In some applications like factory automation and tactile internet each user only requests packets from M' sensors, where $M' < M$. For these applications, multi-cast is an option for DL transmission, where users and sensors are grouped into multiple clusters. Our method can be extended into multi-cast systems by applying it in each cluster.

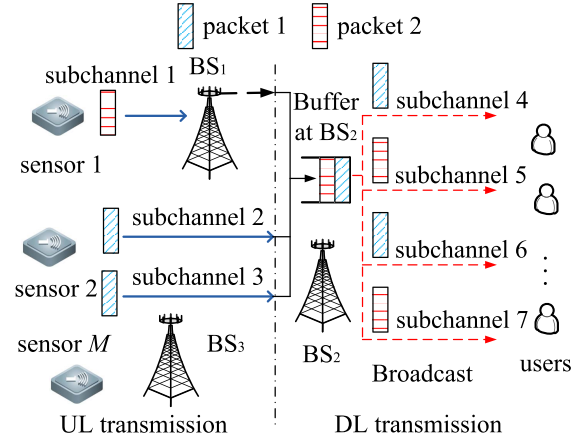


Fig. 5. Illustration of packets delivery mechanism.

two sensors are active. In the UL transmission phase, sensor 1 sends a packet to BS₁ over subchannel 1, and sensor 2 sends two copies of its packet to BS₂ over subchannels 2 and 3. Then, BS₁ forwards the packet from sensor 1 to BS₂ via backhaul. After arriving at the buffer of BS₂, both packets wait in the buffer before DL transmission. In the DL transmission phase, packets in the buffer are broadcast to all the users associated with BS₂ over multiple subchannels at rate E_B^+ packets/frame, where E_B^+ is the minimal integer that is equal to or higher than E_B .

IV. PROBLEM FORMULATION

In this section, we formulate the optimization problem to minimize the total bandwidth by jointly optimizing UL and DL transmission delays, queueing delay and subchannel assignment.

A. Ensuring UL and DL Packet Loss Requirements

When analyzing the reliability of URLLC, Shannon capacity was applied in many existing studies such as [15] and [30], which cannot characterize the decoding error probability.

1) *Constraint on UL Transmission Packet Loss Probability:* Denote the large-scale channel gain of the m th sensor as α_m^u , and the instantaneous channel gain on the i th subchannel allocated to the m th sensor as $g_{m,i}^u = (\mathbf{h}_{m,i}^u)^H \mathbf{h}_{m,i}^u$, where $[\cdot]^H$ denotes the conjugate transpose and $\mathbf{h}_{m,i}^u \in \mathbb{C}^{N_i \times 1}$ is the channel vector whose elements are independent and identically complex Gaussian distributed with zero mean and unit variance. To avoid feedback overhead, CSIT is not assumed available at the sensors, and then the maximal transmit power at each sensor is equally allocated among N_m^u subchannels. For single-input-multiple-output system, the achievable rate from the m th sensor to the BS over the i th subchannel can be accurately approximated by [10]

$$R_{m,i}^u \approx \frac{T_f B_m^u}{\ln 2} \left\{ \ln \left(1 + \frac{\alpha_m^u P_{\max}^u g_{m,i}^u}{\phi N_0 B_m^u N_m^u} \right) - \sqrt{\frac{V_{m,i}^u}{(D^u - 2T_f) B_m^u}} f_Q^{-1}(e_{m,i}^u) \right\} \text{ bits/frames,} \quad (4)$$

where P_{\max}^u is the maximal transmit power of each sensor, $\phi > 1$ reflects the signal-to-noise ratio (SNR) loss due to the errors of channel estimation at receiver,⁵ N_0 is the single-sided noise spectral density, $e_{m,i}^u$ is the decoding error probability (i.e., the block error probability) on the i th subchannel of the m th sensor, $f_Q^{-1}(x)$ is the inverse of the Q-function, and $V_{m,i}^u = 1 - \left[1 + \frac{\alpha_m^u P_{\max}^u g_{m,i}^u}{\phi N_0 B_m^u N_m^u}\right]^{-2}$ [10]. The blocklength of channel coding is determined by the UL bandwidth of each subchannel and transmission duration according to $(D^u - 2T_f)B_m^u$. When the blocklength is large, (4) approaches to the Shannon's capacity.

As shown in [9] and [10], the approximation in (4) is very accurate in quasi-static channel when $e_{m,i}^u \in [10^{-3}, 10^{-6}]$. In typical scenarios, the required transmission delay in URLLC is shorter than the channel coherence time [14], i.e., the channel is quasi-static.⁶ This suggests that the approximation in (4) is applicable for URLLC.

When transmitting a packet that contains b bits over one subchannel, the decoding error probability can be obtained by substituting (4) into $R_{m,i}^u(D^u - 2T_f)/T_f = b$ as

$$e_{m,i}^u = f_Q \left\{ \sqrt{(D^u - 2T_f)B_m^u} \left[\ln \left(1 + \frac{\alpha_m^u P_{\max}^u g_{m,i}^u}{\phi N_0 B_m^u N_m^u} \right) - \frac{b \ln 2}{B_m^u (D^u - 2T_f)} \right] \right\}, \quad (5)$$

where $V_m^u \approx 1$ is applied. $e_{m,i}^u$ in (5) depends on channel, and hence is a random variable. To show the relation between the decoding error probability and packet loss probability, we use indicator functions to represent whether a packet is successfully transmitted over multiple subchannels. If the packet is successfully transmitted to the BS over the i th subchannel assigned to the m th sensor, then $\mathbf{1}_{m,i}^u = 1$. Otherwise, $\mathbf{1}_{m,i}^u = 0$. From (5), we have

$$\Pr\{\mathbf{1}_{m,i}^u = 0\} = \int_0^\infty e_{m,i}^u f_g(x) dx, \quad (6)$$

where $f_g(x)$ is the distribution of instantaneous channel gain. If the elements of $\mathbf{h}_{m,i}^u$ are complex Gaussian distributed, then $f_g(x) = \frac{1}{(N_t-1)!} x^{N_t-1} e^{-x}$. Since each packet is transmitted over N_m^u subchannels, the packet loss probability is given by

$$\begin{aligned} \Pr \left\{ \bigcap_{i=1}^{N_m^u} (\mathbf{1}_{m,i}^u = 0) \right\} &= \prod_{i=1}^{N_m^u} \Pr \{ \mathbf{1}_{m,i}^u = 0 \} \\ &= \left[\int_0^\infty e_{m,i}^u f_g(x) dx \right]^{N_m^u}, \end{aligned} \quad (7)$$

which should be no large than ε^u to guarantee the UL reliability.

Optimizing resource allocation under the constraint on the packet loss probability in (7) is very difficult, because the expression of $e_{m,i}^u$ in (5) is too complicated to obtain any useful

insights. To simplify the analysis, we consider an upper bound of it. Since $f_Q(\cdot)$ is a decreasing function, $e_{m,i}^u$ decreases as $g_{m,i}^u$ increases. Then, an upper bound of $e_{m,i}^u$ can be obtained from

$$e_{m,i}^u \leq \begin{cases} e_m^{u,\text{th}}, & \text{if } g_{m,i}^u \geq g_m^{u,\text{th}}, \\ 1, & \text{if } g_{m,i}^u < g_m^{u,\text{th}}, \end{cases} \quad (8)$$

where $g_m^{u,\text{th}}$ can be obtained by substituting $e_m^{u,\text{th}}$ into $R_{m,i}^u(D^u - 2T_f)/T_f = b$ as

$$g_m^{u,\text{th}} \approx \frac{\phi N_0 B_m^u N_m^u}{\alpha_m^u P_{\max}^u} \left\{ \exp \left[\frac{b \ln 2}{(D^u - 2T_f) B_m^u} + \sqrt{\frac{1}{(D^u - 2T_f) B_m^u}} f_Q^{-1}(e_m^{u,\text{th}}) \right] - 1 \right\}. \quad (9)$$

The upper bound in (8) means that if the instantaneous channel gain $g_{m,i}^u$ is higher than a threshold $g_m^{u,\text{th}}$ such that $R_{m,i}^u(D^u - 2T_f)/T_f \geq b$, then a packet with size b can be transmitted successfully with probability $1 - e_m^{u,\text{th}}$ over the i th subchannel. Otherwise, the packet cannot be transmitted successfully over the i th subchannel.

The upper bound of decoding error probability in (8) is different from outage probability, defined as the probability that the SNR or SINR is lower than a threshold in [15] and [30]. When the channel gain exceeds the threshold, the outage probability is zero. As shown in (8), however, even when the channel gain is higher than the threshold $g_m^{u,\text{th}}$, $e_m^{u,\text{th}}$ is not zero. The relation between $e_m^{u,\text{th}}$ and $g_m^{u,\text{th}}$ is shown in (9).

From (8), the packet loss probability can be bounded by

$$\Pr \left\{ \bigcap_{i=1}^{N_m^u} (\mathbf{1}_{m,i}^u = 0) \right\} \leq \left(\Pr \{ g_{m,i}^u < g_m^{u,\text{th}} \} + e_m^{u,\text{th}} \right)^{N_m^u}. \quad (10)$$

Since $\Pr \{ g_{m,i}^u < g_m^{u,\text{th}} \} = \int_0^{g_m^{u,\text{th}}} f_g(x) dx$, the constraint on UL packet loss probability is

$$f_m^u(N_m^u, B_m^u, e_m^{u,\text{th}}) \triangleq \left[\int_0^{g_m^{u,\text{th}}} f_g(x) dx + e_m^{u,\text{th}} \right]^{N_m^u} \leq \varepsilon^u, \quad m = 1, \dots, M. \quad (11)$$

Remark 1: The upper bound in (8) is not accurate when $g_{m,i}^u$ is smaller than the threshold $g_m^{u,\text{th}}$. In this case, the decoding error probability ranges from $e_m^{u,\text{th}}$ to 1. Since $e_m^{u,\text{th}}$ is smaller than 1, setting $e_{m,i}^u = 1$ when $g_{m,i}^u < g_m^{u,\text{th}}$ leads to conservative resource allocation. We will show the impact of the loose upper bound on bandwidth allocation with numerical results.

2) *Constraint on DL Transmission Packet Loss Probability:* Since the frame duration is much shorter than the E2E delay, it is possible to adjust transmission duration of each packet D^d . To achieve a constant rate of E_B^+ packets per frame, the number of packets that are transmitted simultaneously is $\frac{D^d}{T_f} E_B^+$. Without CSIT, the maximal transmit power of a BS P_{\max}^d is equally allocated among $\frac{D^d}{T_f} E_B^+ N^d$ active subchannels. Denote the average channel gain of the k th user as α_k^d .

⁵The impact of channel estimation errors on data rate can be equivalent to a SNR loss, which depends on the velocity of sensors [31]. Velocity of devices ranges from 0 to 500 km/h [19]. For sensors with slow and median velocity, ϕ is close to 1.

⁶If the transmission duration of each block is less than the channel coherence time, then the channel is referred to as quasi-static channel in [10].

Similar to UL transmission, we can derive the threshold in DL transmission as follows:

$$g_k^{\text{d,th}} \approx \frac{\phi N_0 B^{\text{d}} D^{\text{d}} E_{\text{B}}^+ N^{\text{d}} N_{\text{t}}}{\alpha_k^{\text{d}} P_{\text{max}}^{\text{d}} T_{\text{f}}} \times \left\{ \exp \left[\frac{b \ln 2}{D^{\text{d}} B^{\text{d}}} + \sqrt{\frac{1}{D^{\text{d}} B^{\text{d}}} f_{\text{Q}}^{-1} \left(e_k^{\text{d,th}} \right)} \right] - 1 \right\}, \quad (12)$$

where $e_k^{\text{d,th}}$ is the block error probability when the instantaneous channel gain is $g_k^{\text{d,th}}$. Then, the probability that the packet is not successfully transmitted to the k th user is bounded by

$$f_k^{\text{d}}(N^{\text{d}}, B^{\text{d}}, e_k^{\text{d,th}}) \triangleq \left[\int_0^{g_k^{\text{d,th}}} f_{\text{g}}(x) dx + e_k^{\text{d,th}} \right]^{N^{\text{d}}}. \quad (13)$$

As shown in (12), $g_k^{\text{d,th}}$ decreases with α_k^{d} . Moreover, $f_k^{\text{d}}(N^{\text{d}}, B^{\text{d}}, e_k^{\text{d,th}})$ increases with $g_k^{\text{d,th}}$. Therefore, the user with the lowest average channel gain has the highest packet loss probability. To study DL transmission reliability for all users, we only need to consider the user with the index $k_{\min} = \arg \min_k \alpha_k^{\text{d}}$. Then, the constraint on DL packet loss probability is,

$$f_{k_{\min}}^{\text{d}}(N^{\text{d}}, B^{\text{d}}, e_{k_{\min}}^{\text{d,th}}) \leq \varepsilon^{\text{d}}. \quad (14)$$

Remark 2: Since $V_{m,i}^{\text{u}} \leq 1$, by substituting $V_{m,i}^{\text{u}} = 1$ into (4), we can obtain a lower bound of the achievable rate. As validated in [12], V is very close to 1 when the SNR is higher than 10 dB, which is the typical SNR at the edge of a cell [32]. On the other hand, to guarantee the QoS requirement of typical applications in URLLC [19], the required SNR should be high, which can be supported by equipping multiple antennas at the BS. To simplify the analysis, we use the lower bounds in the rest of this work, i.e., $V_{m,i}^{\text{u}} = 1$. For DL transmission, we can also obtain a lower bound of achievable rate in this way.

B. Total Bandwidth of the System

To ensure D^{q} and ε^{q} , the packet rate for DL transmission should be E_{B}^+ packets per frame. To ensure downlink transmission delay of each packet D^{d} with frame duration T_{f} , the number of packets that are transmitted simultaneously should be $\frac{D^{\text{d}}}{T_{\text{f}}} E_{\text{B}}^+$. Therefore, the bandwidth required for DL transmission in each cell is $\frac{D^{\text{d}}}{T_{\text{f}}} E_{\text{B}}^+ N^{\text{d}} B^{\text{d}}$. Denote the number of active sensors in one frame as M_{a} , and the indices of these sensors as set \mathcal{M}_{a} . Then, the total bandwidth for UL transmission is $\sum_{m \in \mathcal{M}_{\text{a}}} N_m^{\text{u}} B_m^{\text{u}}$. Thus, the required total bandwidth of the system is given by

$$\sum_{m \in \mathcal{M}_{\text{a}}} N_m^{\text{u}} B_m^{\text{u}} + F_{\text{R}}^{-1} \frac{D^{\text{d}}}{T_{\text{f}}} E_{\text{B}}^+ N^{\text{d}} B^{\text{d}}, \quad (15)$$

where F_{R} is the frequency-reuse factor.

Since M_{a} is a random variable, and the resource allocation changes with M_{a} , the BSs need to solve the optimization problem when the number of active sensors changes

(e.g., every millisecond). To reduce computational complexity for solving optimization problem, we introduced an upper bound of the number of active sensors.

Denote $\mathbf{1}_m$ as an indicator function. If the m th sensor is active, $\mathbf{1}_m = 1$. Otherwise, $\mathbf{1}_m = 0$. Then, $\mathbb{E}(\mathbf{1}_m)$ can be expressed as a function of p_m , which is the probability that the m th sensor has a transmission request in each frame. In particular, if there is a request, the sensor will stay active in $(D^{\text{u}} - 2T_{\text{f}})/T_{\text{f}}$ frames, and $\mathbb{E}(\mathbf{1}_m) = (D^{\text{u}} - 2T_{\text{f}})p_m/T_{\text{f}}$. Since all the sensors could be active at the same time, a simple upper bound of (15) can be obtained by setting $M_{\text{a}} = M$. However, such an bound will lead to very conservative bandwidth assignment for UL transmission if the number of sensors is large, where the probability that all the sensors are active is extremely small. In what follows, we provide a threshold of M_{a} , which is an upper bound of M_{a} with high probability. Denote the threshold as M_{a}^{th} . With probability ε_{M} , M_{a} is higher than M_{a}^{th} , i.e., $\varepsilon_{\text{M}} \triangleq \Pr\{M_{\text{a}} > M_{\text{a}}^{\text{th}}\}$. Since M_{a} is the sum of M Bernoulli process, it can be approximated as a Poisson process with parameter $\frac{D^{\text{u}} - 2T_{\text{f}}}{T_{\text{f}}} \sum_{m=1}^M p_m$. Hence, it is not hard to obtain M_{a}^{th} with given ε_{M} . Then, the bandwidth for UL transmission is bounded by $\frac{M_{\text{a}}^{\text{th}}}{M} \sum_{m=1}^M N_m^{\text{u}} B_m^{\text{u}}$ with high probability. When $M_{\text{a}} > M_{\text{a}}^{\text{th}}$, some packets may be lost due to insufficient bandwidth, which however has little impact on the overall reliability if $\varepsilon_{\text{M}} \ll \varepsilon_{\text{max}}$. Then, an upper bound of the required total bandwidth of the system can be obtained as $\frac{M_{\text{a}}^{\text{th}}}{M} \sum_{m=1}^M N_m^{\text{u}} B_m^{\text{u}} + F_{\text{R}}^{-1} \frac{D^{\text{d}}}{T_{\text{f}}} E_{\text{B}}^+ N^{\text{d}} B^{\text{d}} \triangleq B_{\text{tot}}$. With this upper bound, the BSs only need to solve the optimization problem and update resource allocation when the large-scale channel gains change (e.g., every second).

C. Optimization Problem

The optimal UL and DL transmission delays, queueing delay and subchannel assignment that minimize the upper bound of the required total bandwidth to ensure the QoS can be obtained from the following problem:

$$\min_{\substack{D^{\text{u}}, D^{\text{d}}, D^{\text{q}}, N^{\text{d}}, B^{\text{d}}, \\ N_m^{\text{u}}, B_m^{\text{u}}, \varepsilon_m^{\text{u,th}}, e_{k_{\min}}^{\text{d,th}}, \\ m=1, \dots, M}} \frac{M_{\text{a}}^{\text{th}}}{M} \sum_{m=1}^M N_m^{\text{u}} B_m^{\text{u}} + F_{\text{R}}^{-1} \frac{D^{\text{d}}}{T_{\text{f}}} E_{\text{B}}^+ N^{\text{d}} B^{\text{d}} \quad (16)$$

$$\text{s.t. } E_{\text{B}}^+ = \left\{ \frac{\ln(1/\varepsilon^{\text{q}})}{D^{\text{q}} \ln \left[\frac{T_{\text{f}} \ln(1/\varepsilon^{\text{q}})}{\lambda D^{\text{q}}} + 1 \right]} \right\}^+, \quad (16a)$$

$$0 < B_m^{\text{u}} \leq W_{\text{c}}, 0 < B^{\text{d}} \leq W_{\text{c}}, \quad (16b)$$

$$B_m^{\text{u}}, B^{\text{d}} \in \{z B_0, z \in \mathbb{Z}\} \quad (16c)$$

$$0 < N_m^{\text{u}}, 0 < N^{\text{d}}, N_m^{\text{u}}, N^{\text{d}} \in \mathbb{Z}, \quad (16c)$$

$$D^{\text{u}} \in \{3T_{\text{f}}, 4T_{\text{f}}, \dots, D_{\text{max}} - D^{\text{b}} - 2T_{\text{f}}\},$$

$$D^{\text{d}}, D^{\text{q}} \in \{T_{\text{f}}, 2T_{\text{f}}, \dots, D_{\text{max}} - D^{\text{b}} 4T_{\text{f}}\},$$

$$(1), (11), \text{ and } (14),$$

where constraint in (16a) is the required DL packet service rate for ensuring queueing delay and queueing delay bound

violation probability, constraint (16b) ensures the bandwidth of each subchannel less than the coherence bandwidth such that each copy of a packet is transmitted over a flat fading channel,⁷ (1) is the constraint on the E2E delay, and the constraints in (11) and (14) ensure the transmission packet loss probabilities in UL and DL, respectively. Because the upper bounds in (11) and (14) are not unique and depend on $e_m^{u,th}$ and $e_{k_{min}}^{d,th}$, the values of $e_m^{u,th}$ and $e_{k_{min}}^{d,th}$ affect the optimal solution and the total bandwidth. To minimize the required total bandwidth, we adjust $e_m^{u,th}$ and $e_{k_{min}}^{d,th}$ in the upper bounds in an optimal manner.

Remark 3: Similar to the delay components, the system can also adjust packet loss components in queueing and UL and DL transmissions. With different values of ε^u , ε^q and ε^d , the required total bandwidth is different. In problem (16), the values of ε^u , ε^q and ε^d are given. If they are optimization variables, then problem (16) will become intractable. Owing to the following reason, we can provide a simple but reasonable way to divide the constraint on the overall packet loss probability into constraints on different components. If $\varepsilon^u \rightarrow 0$, then according to (11), $N_m^u \rightarrow \infty$, which means that the required bandwidth tends to infinite. Similarly, if $\varepsilon^d \rightarrow 0$, $N^d \rightarrow \infty$. Moreover, by substituting (3) into (15), the relation between B_{tot} and ε^q can be expressed as $B_{tot} = C_1 + C_2 \ln(1/\varepsilon^q) / \ln[C_3 \ln(1/\varepsilon^q) + 1]$, where C_1 , C_2 and C_3 are parameters that do not change with ε^q . If $\varepsilon^q \rightarrow 0$, then B_{tot} approaches infinite. This suggests that all the values of ε^u , ε^q and ε^d cannot be ignored. Here we set them as one third of ε_{max} , i.e., $\varepsilon^u = \varepsilon^q = \varepsilon^d = \varepsilon_{max}/3$. In simulation part, we will show that the total bandwidth with optimal values of ε^u , ε^q and ε^d is almost the same as that with $\varepsilon^u = \varepsilon^q = \varepsilon^d = \varepsilon_{max}/3$.

V. JOINT UL AND DL RESOURCE CONFIGURATION

In this section, we provide a two-step method to find the optimal solution of problem (16). In the first step, we find the optimal subchannel assignment that minimizes total bandwidth with give delay components. In the second step, we find the optimal delay components that minimize the total bandwidth. Finally, we prove that the two-step method can provide the optimal solution.

A. Bandwidth Assignment Optimization

In this subsection, we fix the values of D^u , D^d and D^q , and optimize the values of N_m^u , B_m^u , N^d and B^d . Because the UL and DL bandwidth assignments can be decoupled given the delay components and packet loss components, we first optimize UL subchannel assignment, and then consider DL subchannel assignment.

1) *UL subchannel assignment:* We optimize the values of N_m^u , B_m^u , and $e_m^{u,th}$ in (8) to minimize the UL bandwidth under the constraint on ε^u with given D^u from the following

problem:

$$\begin{aligned} \min_{N_m^u, B_m^u, e_m^{u,th}} \quad & \frac{M_a^{th}}{M} \sum_{m=1}^M N_m^u B_m^u \\ \text{s.t.} \quad & 0 < B_m^u \leq W_c, 0 < N_m^u, N_m \in \mathbb{Z}, \\ & \text{and (11),} \end{aligned} \quad (17)$$

where B_m^u is relaxed to continuous variable. After obtaining the solution of problem (17), the discrete value of B_m^u can be directly obtained from $(B_m^u/B_0)^+ B_0$.

In the sequel, we propose an algorithm to find the optimal solution of problem (17). Since the constraints for each sensor do not depend on those of the other sensors, problem (17) can be further equivalently decomposed into M single-sensor problems as follows:

$$\begin{aligned} \min_{N_m^u, B_m^u, e_m^{u,th}} \quad & N_m^u B_m^u \\ \text{s.t.} \quad & (17a), \text{ and (11).} \end{aligned} \quad (18)$$

To solve problem (18), we need some properties of $f_m^u(N_m^u, B_m^u, e_m^{u,th})$.

Property 1: Given the values of N_m^u and $e_m^{u,th}$, we can find a unique solution of $B_m^{u,min}$ that minimizes $f_m^u(N_m^u, B_m^u, e_m^{u,th})$. Moreover, $f_m^u(N_m^u, B_m^u, e_m^{u,th})$ strictly decreases with B_m^u in the region $B_m^u \in [0, B_m^{u,min}]$ and strictly increases with B_m^u in the region $B_m^u \in [B_m^{u,min}, \infty)$.

Proof: See proof in Appendix B. \square

According to the numerical results in [34], $B_m^{u,min}$ is larger than W_c in typical scenarios of URLLC. In what follows, we propose an algorithm that can find the global optimal solution when $B_m^{u,min} \geq W_c$. For the case $B_m^{u,min} < W_c$, a local optimal solution can be obtained.

Based on Property 1, we can obtain the following property.

Property 2: Given the value of N_m^u , when $B_m^{u,min} \geq W_c$, $f_m^u(N_m^u, B_m^u, e_m^{u,th*})$ strictly decreases with B_m^u in the region $[0, W_c]$, where $e_m^{u,th*}$ is the optimal value that minimizes $f_m^u(N_m^u, B_m^u, e_m^{u,th})$ with given N_m^u and B_m^u .

Proof: See proof in Appendix C. \square

If $f_m^u(N_m^u, W_c, e_m^{u,th*}) > \varepsilon^u$, the reliability can not be guaranteed, and more subchannels are needed. If $f_m^u(N_m^u, W_c, e_m^{u,th*}) \leq \varepsilon^u$, the minimal value of B_m^u that satisfies (11) can be obtained when $f_m^u(N_m^u, B_m^u, e_m^{u,th*}) = \varepsilon^u$, and can be obtained via the binary search method [35]. The search algorithm needs to compute the value of $f_m^u(N_m^u, B_m^u, e_m^{u,th*})$, and hence needs to find $e_m^{u,th*}$ with given B_m^u . To show when $e_m^{u,th*}$ can be obtained with a low complexity method, we provide the following property.

Property 3: Given the values of N_m^u and B_m^u , $f_m^u(N_m^u, B_m^u, e_m^{u,th})$ is convex in $e_m^{u,th}$ when $g_m^{u,th} < N_t - 1$.

Proof: See proof in Appendix D. \square

If $f_m^u(N_m^u, B_m^u, e_m^{u,th})$ is convex in $e_m^{u,th}$, then $e_m^{u,th*}$ can be obtained by the exact linear search method [35]. Otherwise, to obtain $e_m^{u,th*}$, the exhaustive search method should be used. Note that to ensure ultra-high reliability in (11), $g_m^{u,th}$ cannot be too large. For example, when $N_t \geq 2$ and $\varepsilon_{max} \leq 10^{-5}$ (which is typical for URLLC), we have $g_m^{u,th} < N_t - 1$ under

⁷The value of W_c depends on propagation environment, which is not hard to obtain before a system is configured [33].

TABLE I
ALGORITHM TO SOLVE PROBLEM (18)

Input: N_{\max}^u , T_f , b , N_0 , N_t , α_m^u , P_{\max}^u , and accuracy requirement of binary search method δ_b .
Output: N_m^{u*} , B_m^{u*} , and $e_m^{u,th*}$.

```

1:  $N_m^u := 1$ 
2: while  $N_m^u \leq N_{\max}^u$  do
3:   Set  $B_{lb} := 0$ ,  $B_{ub} := W_c$ ,  $B_{bs}^u := 0.5(B_{lb} + B_{ub})$ .
4:   while  $B_{ub} - B_{lb} > \delta_b$  do
5:     Apply exact linear search method to find  $e_{bs}^u$  that minimizes
        $f_m^u(N_m^u, B_{bs}^u, e_{bs}^u)$ .
6:     if  $f_m^u(N_m^u, B_{bs}^u, e_{bs}^u) > \varepsilon^u$  then
7:        $B_{lb} := B_{bs}$ ,  $B_{bs} := 0.5(B_{lb} + B_{ub})$ .
8:     else
9:        $B_{ub} := B_{bs}$ ,  $B_{bs} := 0.5(B_{lb} + B_{ub})$ .
10:    end if
11:  end while
12:  if  $f_m^u(N_m^u, B_{bs}^u, e_{bs}^u) \leq \varepsilon^u$  then
13:     $B_m^u(N_m^u) := B_{bs}^u$  and  $e_m^{u,th}(N_m^u) := e_{bs}^u$ .
14:  else
15:     $B_m^u(N_m^u) := \text{NaN}$  and  $e_m^{u,th}(N_m^u) := \text{NaN}$ .
16:  end if
17:   $N_m^u := N_m^u + 1$ .
18: end while
19:  $N_m^{u*} := \arg \min_{N_m^u} [B_m^u(N_m^u) / B_0]^+ B_0$ .
20:  $B_m^{u*} := [B_m^u(N_m^{u*}) / B_0]^+ B_0$ ,  $e_m^{u,th*} := e_m^{u,th}(N_m^{u*})$ .
21: return  $N_m^{u*}$ ,  $B_m^{u*}$ ,  $e_m^{u,th*}$ .
```

constraint (11) in the cases where $N_m^u \leq 10$. Since large N_m^u results in large bandwidth, and our goal is to minimize the total bandwidth, N_m^u will not be too large. In the proposed search algorithm, we find the optimal solution of problem (17) in the region $0 < N_m^u \leq N_{\max}^u$, where N_{\max}^u is the maximal number of subchannels that can be assigned to each sensor. We will validate that the optimal value of N_m^u is not large with numerical results.

Given the value of N_m^u , according to Property 2 and Property 3, the optimal values of B_m^u and $e_m^{u,th}$ that minimize (18) can be found via binary search method and exact linear search method, respectively. By searching B_m^u and $e_m^{u,th}$ with different values of $N_m^u \in \{1, \dots, N_{\max}^u\}$, the optimal solution of problem (18) can be obtained, denoted as $\{N_m^{u*}, B_m^{u*}, e_m^{u,th*}\}$. To solve problem (17), we need to solve problem (18) for M sensors. Hence, the complexity of the proposed algorithm is $O(MN_{\max}^u)$. The details of the algorithm are provided in Table I.

2) *DL Subchannel Assignment*: The optimal DL subchannel assignment that minimizes the required DL bandwidth can be obtained by solving the following problem:

$$\min_{N^d, B^d, e_{k_{\min}}^{d,th}} \frac{D^d}{T_f} E_B^+ N^d B^d \quad (19)$$

$$\text{s.t. } 0 < B^d \leq W_c, \quad (19a)$$

$$0 < N^d, N^d \in \mathbb{Z}, \quad (19b)$$

$$\text{and (14),}$$

where F_R^{-1} is removed from the objective function since it does not change the optimal solution. Similar to $f_m^u(N_m^u, B_m^u, e_m^{u,th})$ in (11), we can prove that

$f_{k_{\min}}^d(N^d, B^d, e_{k_{\min}}^{d,th})$ in (14) also satisfies Property 1, Property 2 and Property 3. The proofs are similar to that in Appendices B, C and D, and hence are omitted for conciseness. Therefore, the solution of problem (19) can also be found with the algorithm in Table I, and is denoted as $\{N_m^{d*}, B_m^{d*}, e_m^{d,th*}\}$.

B. Delay Components Optimization

To show how to optimize the delay components and to reduce complexity in optimization, we first analyze the relations between the required bandwidth and the delay components.

1) *Increasing Queueing Delay Bound*: From (3), we can directly obtain that the required DL service rate decreases with D^q . As a result, the DL bandwidth decreases with D^q .

2) *Increasing UL Transmission Delay*: In order to show the relationship between UL bandwidth and UL transmission delay, we compare two systems with different UL transmission delays $\hat{D}^u < \tilde{D}^u$. With $2T_f$ delay caused by control signaling, the transmission delay is $\hat{D}^u - 2T_f$ (or $\tilde{D}^u - 2T_f$). Denote $\hat{\mathbf{1}}_m$ and $\tilde{\mathbf{1}}_m$ as the indicator functions that indicate whether the m th sensor is active in the first and the second systems, respectively. Given the probability that the m th sensor requests to transmit a packet in a certain frame as p_m , the probabilities that the m th sensor is active can be expressed as $\mathbb{E}(\hat{\mathbf{1}}_m) = \frac{\hat{D}^u - 2T_f}{T_f} p_m$ and $\mathbb{E}(\tilde{\mathbf{1}}_m) = \frac{\tilde{D}^u - 2T_f}{T_f} p_m$, respectively. Then,

$$\frac{T_f \mathbb{E}(\hat{\mathbf{1}}_m)}{\hat{D}^u - 2T_f} = p_m = \frac{T_f \mathbb{E}(\tilde{\mathbf{1}}_m)}{\tilde{D}^u - 2T_f}. \quad (20)$$

The bandwidth for UL transmission in (17) is hard to analyze since M_a^{th} has no closed-form expression. To obtain some useful insights, we study the average bandwidth for UL transmission.

Proposition 2: Increasing D^u can reduce the required average bandwidth for UL transmission.

Proof: See proof in Appendix E. \square

Although the average bandwidth cannot reflect the bandwidth requirement directly, the simulations in the next section validate that the minimal UL bandwidth decreases with D^u .

3) *Increasing DL Transmission Time*: Based on the following proposition for DL transmission, we can obtain a different conclusion from UL transmission.

Proposition 3: If constraint (19a) is inactive (i.e., $B^d < W_c$), the required minimal bandwidth for DL transmission does not change with D^d .

Proof: See proof in Appendix F. \square

Given the number of subchannels for each packet transmission, B^d increases as D^d decreases. Whether constraint (19a) is inactive or not depends on the number of antennas at each BS, the radius of each cell and communication environment. If N_t is large or the radius of a cell is small, then $B^d < W_c$ even when $D^d = T_f$. We will provide related numerical results in the next section.

4) *Joint Optimization of the Three Delay Components*: The above analysis shows that given E2E delay, the tradeoff between delay components lead to a tradeoff between UL

bandwidth and DL bandwidth. To minimize the total bandwidth, we need to optimize the delay components.

For any given D^u , D^d and D^q , by solving problem (17) and problem (19), we can obtain the optimal bandwidth assignment policy and the minimized total bandwidth, which are denoted as $\Phi^*(D^u, D^d, D^q) \triangleq (N_m^{u*}, B_m^{u*}, e_m^{u,th*}, N_m^{d*}, B_m^{d*}, e_m^{d,th*})$ and $B_{tot}(\Phi^*(D^u, D^d, D^q))$, respectively. Since the optimal bandwidth assignment policy depends on the delay components, $\Phi^*(\cdot)$ and $B_{tot}(\cdot)$ are functions of the delay components. To obtain the optimal delay components, we search the values of D^u , D^d , and D^q . Since the possible values of D^u , D^d , and D^q in problem (16) are finite, it is not hard to obtain D^{u*} , D^{d*} , and D^{q*} that minimize $B_{tot}(\Phi^*(D^u, D^d, D^q))$ with the exhaustive search method.

For large N_t or small cells, $B^d < W_c$. According to Proposition 3, $D^{d*} = T_f$. We only need to search D^u and D^q under the constraint $D^u + D^q \leq D_{max} - D^b - T_f$. Further considering that the bandwidth is minimized when $D^u + D^q = D_{max} - D^b - T_f$, we only need to search D^u in $(0, D_{max} - T_f)$, which is one-dimensional searching, and hence the complexity is not high.

C. Optimality of the Two-Step Method

To show that D^{u*} , D^{d*} , D^{q*} and $\Phi^*(D^{u*}, D^{d*}, D^{q*})$ is the optimal solution of problem (16), we need the following proposition.

Proposition 4: For an arbitrary solution of problem (16), \tilde{D}^u , \tilde{D}^d , \tilde{D}^q and $\tilde{\Phi}(\tilde{D}^u, \tilde{D}^d, \tilde{D}^q)$, we always have $B_{tot}(\Phi^*(D^{u*}, D^{d*}, D^{q*})) \leq B_{tot}(\tilde{\Phi}(\tilde{D}^u, \tilde{D}^d, \tilde{D}^q))$.

Proof: See proof in Appendix G. \square

This suggests that if both the solution for bandwidth assignment policy and the solution for delay components are global optimal, then the two-step method gives rise to a global optimal solution of problem (16). The global optimal bandwidth assignment policy can be obtained by the algorithm in Table I when $W_c \leq B_m^{u,min}$ and $g_m^{u,th} < N_t - 1$. Otherwise, we need to use exhaustive searching to find the global optimal solution.

VI. SIMULATION AND NUMERICAL RESULTS

In this section, we validate the analyses and demonstrate the required bandwidth to support URLLC. With simulation results, we show the impact of minimizing the upper bound of the required total bandwidth and illustrate the performance gain with jointly UL and DL configuration, where arrivals of packet at each sensor are generated by simulation. With numerical results, we show the impact of the upper bound in (8) on the bandwidth required by each sensor and illustrate the optimal number of subchannels allocated to each sensor.

Simulation Parameters are listed in Table II. Since broadcast is used in DL transmission, the required bandwidth does not change with the number of users. To ensure the QoS requirement of all users, we only consider the users with the worst large-scale channel gains, i.e., the users located at the edge of each cell. The path loss model is $-10\lg(a_m) = 35.3 + 37.6\lg(d_m)$, where d_m (m) is the distance between sensors and the BSs they associated with. d_m is uniformly

TABLE II
SIMULATION PARAMETERS [4], [26], [27], [33]

Number of BSs	3
Number of sensors M	3000
Maximal transmit power of a sensor P_{max}^u	23 dBm
Maximal transmit power of a BS P_{max}^d	46 dBm
Overall packet loss probability ε_{max}	10^{-7}
Latency in radio access network $D_{max} - D^b$	1 ms
Frame duration T_f	0.1ms
Backhaul delay D^b	0.1 ms
Single-sided noise spectral density N_0	-174 dBm/Hz
Coherence bandwidth W_c	0.5 MHz
Maximal number of subchannels allocated to each packet N_{max}^u and N_{max}^d	10
Average packet rate generated by each sensor	100 packets/s

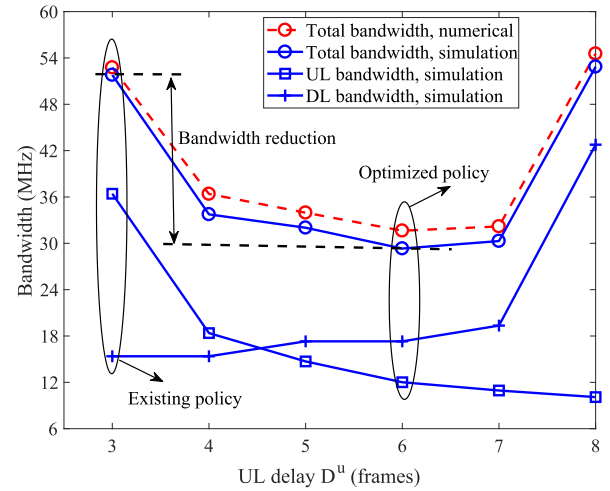


Fig. 6. Required bandwidth versus delay components, where $\varepsilon^u = \varepsilon^d = \varepsilon^q = \varepsilon_{max}/3$, and $N_t = 8$.

distributed in $[50, r]$ m, where r is the radius of each cell. We only consider the scenarios where $d_m > 50$ m, because the large-scale channel gain α_m^u decreases with d_m , and more resources are needed to guarantee QoS requirement with larger sensor-BS distance.

We solve problem (19) with different values of N_t and cell size to show when the required bandwidth of each subchannel is less than the channel coherence bandwidth. Our results show that if the radius of each cell is 100 m, then $B^{d*} < W_c$ when $N_t \geq 4$, and if the radius is 250 m, then $B^{d*} < W_c$ when $N_t \geq 8$. In the rest of this section, the radius of each cell is set to be 250 m and $N_t \geq 8$. As a result, Proposition 3 holds. In other words, required minimal bandwidth for DL transmission does not change with D^d . The required total bandwidth is minimized when $D^{d*} = T_f$.

Figure 6 shows the required bandwidth with different delay components. The numerical results are obtained by solving problem (16) with different values of D^u and D^q , where $D^{d*} = T_f$ and $\varepsilon_M = \Pr\{M_a > M_a^{th}\} = 10^{-15}$ (i.e., $\varepsilon_M \ll \varepsilon_{max}$). M_a can be approximated by a Poisson process

TABLE III
TOTAL BANDWIDTH WITH DIFFERENT PACKET LOSS COMPONENTS

	$N_t = 8$	$N_t = 16$	$N_t = 32$
Optimal ε^u , ε^d and ε^q	28.6 MHz	19.9 MHz	16.8 MHz
$\varepsilon^u = \varepsilon^d = \varepsilon^q = \varepsilon_{\max}/3$	29.3 MHz	20.2 MHz	17.0 MHz

with mean $\lambda(D^u - 2T_f)/T_f$, where λ average number of packets generated by M sensors in one frame. Based on this distribution, we can obtain M_a^{th} in (16) from $\Pr\{M_a > M_a^{\text{th}}\} = 10^{-15}$. To show the impact of minimizing the upper bound of the required total bandwidth, we also provide the simulation results. To obtain the simulation result, we first compute the total bandwidth in (15) achieved by the optimal transmission policy during 10^6 frames. The maximal total bandwidth in the 10^6 frames is the required total bandwidth to ensure the QoS requirement and is shown in Figure 6. We can see that the total bandwidth obtained via numerical results is higher than that obtained via simulation results, and the gap between them is small. This means that the objective function (16) is a tight upper bound of the total bandwidth in (15). To show the gain of jointly optimizing the delay components, we provide the results with an existing policy, where the UL data transmission finished in each frame [1]. With the existing policy, the required bandwidth is show in Fig. 6 when $D^u = 3 T_f$ (two frames are occupied by control signaling). We can see that nearly half of the total bandwidth can be saved by optimizing the delay components. The results also indicate that the maximal bandwidth for UL transmission decreases with D^u , which agrees with Proposition 2.

Simulation results in Table III show the impact of packet loss probabilities on the required bandwidth. The values of D^u , D^d and D^q are set as the optimal values that minimize the total bandwidth in Fig. 6. The optimal values of ε^u , ε^d and ε^q are obtained by exhaustive search in the region $[0, \varepsilon_{\max}]$. To reduce complexity, the accuracy is set to be $0.05\varepsilon_{\max}$. For any given values of ε^u , ε^d and ε^q , the bandwidth assignment is obtained by solving problem (18) and problem (19). With the bandwidth assignment, the total bandwidth is obtained via simulation, i.e., the maximal total bandwidth in 10^6 frames. The results show that the total bandwidth with optimal values of ε^u , ε^d and ε^q is very close to that with $\varepsilon^u = \varepsilon^d = \varepsilon^q = \varepsilon_{\max}/3$. This validates Remark 3.

The upper bound of packet loss probability in (8) is used to formulate problem (16), and hence the bandwidth allocation is conservative. The numerical results in Fig. 7 show the impact of the upper bound on the required bandwidth, where UL transmission is considered. For DL transmission, the results are similar. Given the required decoding error probability, the minimal bandwidth with accurate model is obtained by exhaustive searching under constraint $\int_0^\infty e_{m,i}^u f_g(x) dx \leq \varepsilon^u$, where $e_{m,i}^u$ is given in (5). The results show that the gap between the minimal bandwidth obtained via the upper bound and that with the accurate model decreases with the number of antennas and increases with the sensor-BS distance. This means that the upper bound has little impact on the resource allocation for macro BSs with a large number of antennas

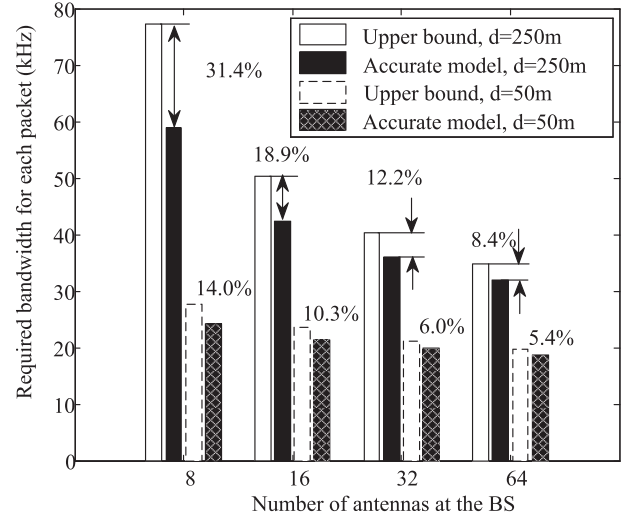


Fig. 7. Impact of the upper bound in (8) on the required bandwidth, where $\varepsilon^u = \varepsilon_{\max}/3$, $D^u = 6T_f$, and $N_m^u = 1$.

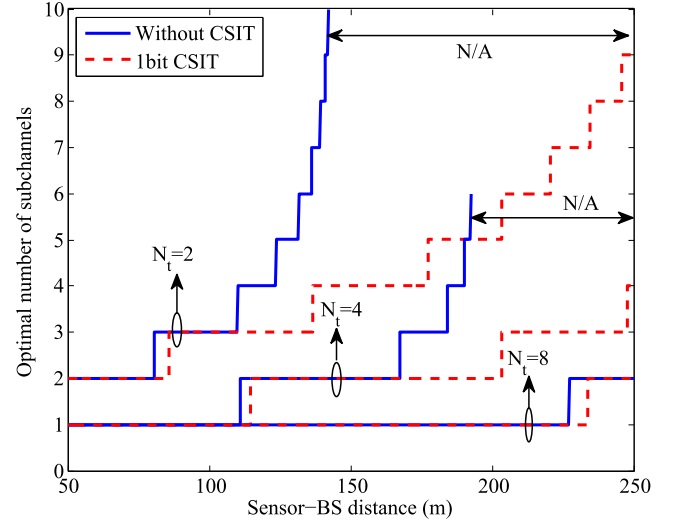


Fig. 8. Optimal number of links for frequency diversity.

and small BSs with short sensor-BS distance. However, when the number of antennas is small and the sensor-BS distance is large, the upper bound leads to conservative resource allocation. This is because when $g_{m,i}^u < g_{m,i}^{u,\text{th}}$, the decoding error probability is much smaller than 1, setting $e_{m,i}^u = 1$ leads to conservative resource allocation.

The optimal number of subchannels assigned to each sensor for UL transmission is illustrated by the numerical results in Fig. 8. We compare two kinds of policies. The first policy does not exploit CSIT, and equally allocates transmit power on the subchannels. With the second policy considered in [1], one bit information of channel gain on each subchannel is available at the transmitter for resource allocation. With the one bit information, the transmitter knows whether the channel gain on a subchannel is above a threshold that a packet can be decoded successfully with the required probability if the maximal transmit power is allocated to the subchannel. Then, the packet is transmitted on one of the subchannels with good channel. When $N_t \geq 16$, $N_m^* = 1$ for all the sensors with any of these two policies (which are not shown in the figure).

TABLE IV
AVAILABILITY WHEN SHADOWING IS CONSIDERED

N_t	16	32	64	128
$D^u = 3T_f$	0.941	0.985	0.9961	0.9991
$D^u = 4T_f$	0.988	0.9981	0.9997	0.99995
$D^u = 5T_f$	0.995	0.9993	0.99991	0.999989
$D^u = 6T_f$	0.997	0.9996	0.99996	0.999996

The results show that CSIT is not helpful for saving bandwidth when the sensor-BS distance is short or the number of active antennas is large.

Availability is another key performance metric for the systems supporting URLLC except the reliability and latency, which is the probability that a system can provide the required QoS (i.e., D_{\max} and ε_{\max}) for all users [3]. Considering that the availability highly depends on shadowing, in the following we provide simulation with shadowing, which follows a lognormal distribution with zero mean and 8 dB standard deviation [33]. Since the transmit power of each sensor is limited, we provide the results for UL transmission in Table IV. To obtain the simulation results, we generate the locations and shadowing of 3000 devices randomly with 10^4 times. The results show that to guarantee the availability of 0.99999 with a single wireless link, $N_t = 128$ and $D^u = 0.6$ ms. To further improve availability, macro-diversity is an option [36].

VII. CONCLUSION

In this paper, we studied joint UL and DL resource configuration to minimize the total bandwidth under strict E2E delay and packet loss probability requirements. A packet delivery mechanism was proposed. In UL transmission, bandwidth is only assigned to the active sensors, and broadcast is used in DL transmission. Channel state information is not available at sensors in UL transmission and not available at the BS in DL transmission. To reduce the required packet rate for ensuring queueing delay at the buffers of the BSs, a statistical multiplexing queueing mode was considered. The total bandwidth required by the mechanism to ensure the E2E delay and overall reliability was minimized by jointly optimizing UL and DL transmission delays, queueing delay and bandwidth assignment. A two-step method was proposed to find the optimal solution of the problem. We first optimized the bandwidth assignment with given delay components and packet loss components. Then, the UL and DL transmission delays and queueing delay were optimized given the E2E delay requirement. Analysis showed that there is a tradeoff between UL and DL bandwidth and it is necessary to optimize the delay components in order to minimize the total bandwidth. Simulation and numerical results validated our analysis and showed that the joint resource configuration can save half of the total bandwidth comparing with an existing policy, where UL and DL transmission delays are not optimized.

APPENDIX A PROOF OF PROPOSITION 1

Proof: For Poisson arrival process with λ , the required minimal constant service rate is provided in (3). From (3),

the required minimal constant service rate of Poisson arrival process with average rate $\tilde{\lambda} = \lambda/L$ is

$$\tilde{E}_B = \frac{T_f \ln(1/\varepsilon^q)}{D^q \ln \left[L \frac{T_f \ln(1/\varepsilon^q)}{\lambda D^q} + 1 \right]} \text{ packets/frame.} \quad (\text{A.1})$$

To prove $L\tilde{E}_B > E_B$ ($L = 2, 3, \dots, K$), we only need to prove

$$L \ln \left[\frac{T_f \ln(1/\varepsilon^q)}{\lambda D^q} + 1 \right] > \ln \left[L \frac{T_f \ln(1/\varepsilon^q)}{\lambda D^q} + 1 \right], \quad (\text{A.2})$$

which can be obtained by substituting (3) and (A.1) into $L\tilde{E}_B > E_B$. Denote $x = \frac{T_f \ln(1/\varepsilon^q)}{\lambda D^q}$. Then, (A.2) can be equivalently rewritten as $f_L(x) \triangleq L \ln(x+1) - \ln(Lx+1) > 0, \forall x > 0$. It is easy to show that $f_L(0) = 0$, and $f'_L(x) = \frac{L}{x+1} - \frac{L}{Lx+1}$, which is positive for $L > 1$. Therefore, $f_L(x) > 0, \forall x > 0$. This completes the proof. \square

APPENDIX B PROOF OF PROPERTY 1

Proof: By substituting $f_g(x) = \frac{1}{(N_t-1)!} x^{N_t-1} e^{-x}$ into (11), we have

$$\begin{aligned} & f_m^u(N_m^u, B_m^u, e_m^{u,\text{th}}) \\ &= \left[\int_0^{g_m^{u,\text{th}}} \frac{1}{(N_t-1)!} x^{N_t-1} e^{-x} dx + e_m^{u,\text{th}} \right]^{N_m^u}. \end{aligned} \quad (\text{B.1})$$

Denote $f_e = \int_0^{g_m^{u,\text{th}}} \frac{1}{(N_t-1)!} x^{N_t-1} e^{-x} dx$. To prove property 1, we only need to prove that f_e first strictly decreases with B_m and then increases with B_m . To this end, we first prove that f_e strictly increases with $g_m^{u,\text{th}}$ and then prove that $g_m^{u,\text{th}}$ first strictly decreases with B_m^u and then increases with B_m .

From f_e , we can obtain that $\frac{\partial f_e}{\partial g_m^{u,\text{th}}} = \frac{(g_m^{u,\text{th}})^{N_t-1} e^{-g_m^{u,\text{th}}}}{(N_t-1)!} > 0$. As a result, f_e strictly increases with $g_m^{u,\text{th}}$. For notation simplicity, (9) can be rewritten as follows,

$$g_m^{u,\text{th}} = C_1 B_m^u \left[\exp \left(\frac{C_2}{B_m^u} + \frac{C_3}{\sqrt{B_m^u}} \right) - 1 \right], \quad (\text{B.2})$$

where $C_1 = \frac{\phi N_0}{\alpha_m^u T_{\max}^u} > 0$, $C_2 = \frac{b \ln 2}{D^u - 2T_f} > 0$, and $C_3 = \sqrt{\frac{1}{D^u - 2T_f}} f_Q^{-1}(e_m^{u,\text{th}}) > 0$.

It is not hard to see that (B.2) is the same as (19) in [34]. According to the proof in Appendix B in [34], $g_m^{u,\text{th}}$ first strictly decreases with B_m^u and then strictly increases with B_m^u , and there is a unique solution of B_m^u that minimizes $g_m^{u,\text{th}}$. This completes the proof. \square

APPENDIX C PROOF OF PROPERTY 2

Proof: To prove that $f_m^u(N_m^u, B_m^u, e_m^{u,\text{th}*})$ decreases with B_m^u in the region $B_m^u \in [0, W_c]$, we show that for any $W_m^u < \tilde{W}_m^u \leq W_c$, $f_m^u(N_m^u, W_m^u, e_m^{u,\text{th}*}) > f_m^u(N_m^u, \tilde{W}_m^u, \tilde{e}_m^{u,\text{th}*})$, where $\tilde{e}_m^{u,\text{th}*}$ is the optimal value of $e_m^{u,\text{th}}$ that minimizes $f_m^u(N_m^u, \tilde{W}_m^u, e_m^{u,\text{th}})$ with given N_m^u and \tilde{W}_m^u . According to Property 1, given N_m^u and $e_m^{u,\text{th}*}$, we have

$$f_m^u(N_m^u, W_m^u, e_m^{u,\text{th}*}) > f_m^u(N_m^u, \tilde{W}_m^u, e_m^{u,\text{th}*}). \quad (\text{C.1})$$

Since $\tilde{e}_m^{u,th*}$ is the optimal value of $e_m^{u,th}$ that minimizes $f_m^u(N_m^u, \tilde{W}_m^u, e_m^{u,th})$, we have

$$f_m^u(N_m^u, \tilde{W}_m^u, e_m^{u,th*}) \geq f_m^u(N_m^u, \tilde{W}_m^u, \tilde{e}_m^{u,th*}). \quad (C.2)$$

From (C.1) and (C.2), we can obtain that $f_m^u(N_m^u, \tilde{W}_m^u, e_m^{u,th*}) > f_m^u(N_m^u, \tilde{W}_m^u, \tilde{e}_m^{u,th*})$. The proof follows. \square

APPENDIX D PROOF OF PROPERTY 3

Proof: According to (B.1), to study the convexity of $f_m^u(N_m^u, B_m^u, e_m^{u,th})$, we only need to study the convexity of f_e . To this end, we first prove that $g_m^{u,th}$ in (9) is convex in $e_m^{u,th}$. Then, we show that f_e is an increasing and convex function of $g_m^{u,th}$ when $g_m^{u,th} < N_t - 1$.

For the Q-function $f_Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty \exp\left(-\frac{\tau^2}{2}\right) d\tau$, we have $f_Q'(x) \triangleq -\frac{1}{\sqrt{2\pi}} e^{-x^2/2} < 0$, and $f_Q''(x) = \frac{x}{\sqrt{2\pi}} e^{-x^2/2} > 0$ when $x > 0$. Thus, when $x > 0$, $f_Q(x)$ is a decreasing and convex function. Moreover, $f_Q(x) < 0.5, \forall x > 0$. Because $e_m^{u,th} < \varepsilon_{\max} < 0.5$ that is true for URLLC applications, and the inverse function of a decreasing and convex function is also convex [35], $f_Q^{-1}(e_m^{u,th})$ is convex in $e_m^{u,th}$, $\forall e_m^{u,th} < 0.5$. Denote $z = f_Q^{-1}(e_m^{u,th})$. Then, $g_m^{u,th}$ in (9) can be rewritten as $g_m^{u,th} = C_4 [\exp(C_5 + C_6 z) - 1]$, where $C_4 = \frac{\phi N_0 B_m^u}{\alpha_m^u P_{\max}^u} > 0$, $C_5 = \frac{b \ln 2}{(D^u - 2T_f) B_m^u} > 0$ and $C_6 = \sqrt{\frac{1}{(D^u - 2T_f) B_m^u}} > 0$. It is easy to see that $g_m^{u,th}$ is an increasing and convex function of z . According to the composition rules, $g_m^{u,th}$ is convex in $e_m^{u,th}$ [35].

It is not hard to derive that $\frac{\partial^2 f_e}{\partial (g_m^{u,th})^2} = \frac{(g_m^{u,th})^{N_t-2} e^{-g_m^{u,th}}}{(N_t-1)!} (N_t - 1 - g_m^{u,th})$. When $N_t - 1 \geq g_m^{u,th}$, f_e is increasing and convex in $g_m^{u,th}$. According to the composition rules, f_e is convex in $e_m^{u,th}$, when $N_t - 1 \geq g_m^{u,th}$. This completes the proof. \square

APPENDIX E PROOF OF PROPOSITION 2

Proof: Denote the bandwidth for the m th sensor in two systems as \hat{B}_m^u and \tilde{B}_m^u , respectively. To keep the decoding error probability identical, the values of $N_m^u, g_m^{u,th}, e_m^{u,th}$ in (11) are fixed in the two systems. According to (9), the relationship between \hat{B}_m^u and \tilde{B}_m^u can be obtained from

$$\begin{aligned} \hat{B}_m^u & \left\{ \exp \left[\frac{b \ln 2}{(\hat{D}^u - 2T_f) \hat{B}_m^u} + \frac{f_Q^{-1}(e_m^{u,th})}{\sqrt{(\hat{D}^u - 2T_f) \hat{B}_m^u}} \right] - 1 \right\} \\ & = \tilde{B}_m^u \left\{ \exp \left[\frac{b \ln 2}{(\tilde{D}^u - 2T_f) \tilde{B}_m^u} + \frac{f_Q^{-1}(e_m^{u,th})}{\sqrt{(\tilde{D}^u - 2T_f) \tilde{B}_m^u}} \right] - 1 \right\}. \end{aligned} \quad (E.1)$$

Since in typical scenarios $B_m^{u,min} \geq W_c$, $g_m^{u,th}$ in (9) strictly decreases with B_m^u in the region $[0, W_c]$. Hence, both left and right hand sides of (E.1) decrease with B_m^u . Moreover,

by substituting $\tilde{B}_m^u = \frac{\hat{D}^u - 2T_f}{\tilde{D}^u - 2T_f} \hat{B}_m^u$ into (E.1), the left hand side of (E.1) is larger than the right hand side of it. Therefore, to satisfy (E.1), $\tilde{B}_m^u < \frac{\hat{D}^u - 2T_f}{\tilde{D}^u - 2T_f} \hat{B}_m^u$. From $\tilde{B}_m^u < \frac{\hat{D}^u - 2T_f}{\tilde{D}^u - 2T_f} \hat{B}_m^u$ and (20), we have

$$\begin{aligned} \mathbb{E} \left(\sum_{m=1}^M \tilde{\mathbf{I}}_m N_m^u \tilde{B}_m^u \right) & < \sum_{m=1}^M \frac{\tilde{D}^u - 2T_f}{\hat{D}^u - 2T_f} \mathbb{E}(\hat{\mathbf{I}}_m) N_m^u \frac{\hat{D}^u - 2T_f}{\tilde{D}^u - 2T_f} \hat{B}_m^u \\ & = \mathbb{E} \left(\sum_{m=1}^M \hat{\mathbf{I}}_m N_m^u \hat{B}_m^u \right). \end{aligned}$$

The proof follows. \square

APPENDIX F PROOF OF PROPOSITION 3

Proof: To prove Proposition 3, we need to prove that the minimal DL bandwidth obtained by solving problem (19) does not change with D^d when constraint (19a) is inactive. We consider two systems with different DL transmission time, i.e., $\hat{D}^d \neq \tilde{D}^d$. We refer to problem (19) with \hat{D}^d and \tilde{D}^d as Problem A and Problem B, respectively. Denote the optimal solutions of Problem A and Problem B as $\{\hat{N}^{d*}, \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$ and $\{\tilde{N}^{d*}, \tilde{B}^{d*}, \tilde{e}_{k_{\min}}^{d,th*}\}$, respectively. Given transmission duration \hat{D}^d and $\{\hat{N}^{d*}, \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$, the threshold in (12) is denoted as $\hat{g}_{k_{\min}}^{d,th}$.

To prove $\frac{\hat{D}^d}{T_f} E_B^+ \hat{N}^{d*} \hat{B}^{d*} = \frac{\tilde{D}^d}{T_f} E_B^+ \tilde{N}^{d*} \tilde{B}^{d*}$, we assume they are not equal, and find contradiction. Without loss of generality, we assume $\frac{\hat{D}^d}{T_f} E_B^+ \hat{N}^{d*} \hat{B}^{d*} < \frac{\tilde{D}^d}{T_f} E_B^+ \tilde{N}^{d*} \tilde{B}^{d*}$. To this end, we first validate that $\{\hat{N}^{d*}, \frac{\hat{D}^d}{\tilde{D}^d} \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$ is a feasible solution of problem B. Since \hat{N}^{d*} is a solution of problem A, constraint (19b) is satisfied. Then, we only need to validate that (14) is satisfied. \hat{N}^{d*} and $\hat{e}_{k_{\min}}^{d,th*}$ are the same in $\{\hat{N}^{d*}, \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$ and $\{\hat{N}^{d*}, \frac{\hat{D}^d}{\tilde{D}^d} \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$. Since $\{\hat{N}^{d*}, \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$ is a solution of problem A, constraint (14) is satisfied with $\hat{N}^{d*}, \hat{e}_{k_{\min}}^{d,th*}$ and $\hat{g}_{k_{\min}}^{d,th}$. If $\hat{g}_{k_{\min}}^{d,th}$ with transmission duration \tilde{D}^d and solution $\{\hat{N}^{d*}, \frac{\hat{D}^d}{\tilde{D}^d} \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$ is the same as $\hat{g}_{k_{\min}}^{d,th}$, then constraint (14) is satisfied. Substituting $\frac{\hat{D}^d}{\tilde{D}^d} \hat{B}^{d*}$ and $\hat{e}_{k_{\min}}^{d,th*}$ into (12), we have

$$\begin{aligned} g_{k_{\min}}^{d,th} & = \frac{\tilde{D}^d \frac{\hat{D}^d}{\tilde{D}^d} \hat{B}^{d*}}{T_f} \\ & \times \left\{ \exp \left[\frac{b \ln 2}{\tilde{D}^d \frac{\hat{D}^d}{\tilde{D}^d} \hat{B}^{d*}} + \sqrt{\frac{1}{\tilde{D}^d \frac{\hat{D}^d}{\tilde{D}^d} \hat{B}^{d*}}} f_Q^{-1}(\hat{e}_{k_{\min}}^{d,th*}) \right] - 1 \right\} \\ & = \frac{\hat{D}^d \hat{B}^d}{T_f} \left\{ \exp \left[\frac{b \ln 2}{\hat{D}^d \hat{B}^d} + \sqrt{\frac{1}{\hat{D}^d \hat{B}^d}} f_Q^{-1}(\hat{e}_{k_{\min}}^{d,th*}) \right] - 1 \right\} \\ & = \hat{g}_{k_{\min}}^{d,th}. \end{aligned} \quad (F.1)$$

Therefore, $\{\hat{N}^{d*}, \frac{\hat{D}^d}{\tilde{D}^d} \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$ is a feasible solution of problem B.

Given the transmission duration \tilde{D}^d , the number of packets that are transmitted simultaneously is $\frac{\tilde{D}^d}{T_f} E_B^+$. Therefore, the required bandwidth for DL transmission

with $\{\hat{N}^{d*}, \frac{\hat{D}^d}{\hat{D}^d} \hat{B}^{d*}, \hat{e}_{k_{\min}}^{d,th*}\}$ satisfies $\frac{\hat{D}^d}{T_f} E_B^+ \hat{N}^{d*} \frac{\hat{D}^d}{\hat{D}^d} \hat{B}^{d*} = \frac{\hat{D}^d}{T_f} E_B^+ \hat{N}^{d*} \hat{B}^{d*} < \frac{\hat{D}^d}{T_f} E_B^+ \tilde{N}^{d*} \tilde{B}^{d*}$, which contradicts with the assumption that $\{\tilde{N}^{d*}, \tilde{B}^{d*}, \tilde{e}_{k_{\min}}^{d,th*}\}$ is the optimal solution of problem B. The proof follows. \square

APPENDIX G PROOF OF PROPOSITION 4

Proof: Given the delay components $\tilde{D}^u, \tilde{D}^d, \tilde{D}^q$, from the first step of the two-step method the optimal subchannel assignment policy is $\Phi^*(\tilde{D}^u, \tilde{D}^d, \tilde{D}^q)$, i.e.,

$$B_{\text{tot}}(\Phi^*(\tilde{D}^u, \tilde{D}^d, \tilde{D}^q)) \leq B_{\text{tot}}(\tilde{\Phi}(\tilde{D}^u, \tilde{D}^d, \tilde{D}^q)). \quad (\text{G.1})$$

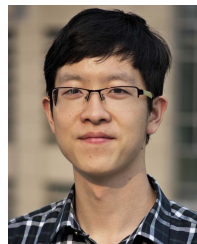
According to the second step of the two-step method, the optimal delay components that minimizes $B_{\text{tot}}(\Phi^*(D^u, D^d, D^q))$ are D^{u*}, D^{d*} and D^{q*} , and hence

$$B_{\text{tot}}(\Phi^*(D^{u*}, D^{d*}, D^{q*})) \leq B_{\text{tot}}(\Phi^*(\tilde{D}^u, \tilde{D}^d, \tilde{D}^q)). \quad (\text{G.2})$$

From (G.1) and (G.2), we have Proposition 4. This completes the proof. \square

REFERENCES

- [1] C. She, C. Yang, and T. Q. S. Quek, "Uplink transmission design with massive machine type devices in tactile Internet," in *Proc. IEEE Global Commun. Conf. (GlobeCom) Workshops*, Dec. 2016, pp. 1–6.
- [2] 3GPP, "Study on scenarios and requirements for next generation access technologies," Tech. Specification Group Radio Access Netw., 3GPP, Valbonne, France, Tech. Rep. 38.913, Release 14, Jun. 2017.
- [3] P. Popovski *et al.*, *Deliverable d6.3 Intermediate System Evaluation Results*, document ICT-317669-METIS/D6.3, 2014. [Online]. Available: https://www.metis2020.com/wp-content/uploads/deliverables/METIS_D6.3_v1.pdf
- [4] G. P. Fettweis, "The tactile Internet: Applications and challenges," *IEEE Veh. Technol. Mag.*, vol. 9, no. 1, pp. 64–70, Mar. 2014.
- [5] M. Simsek, A. Aijaz, M. Dohler, J. Sachs, and G. Fettweis, "5G-enabled tactile Internet," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 3, pp. 460–473, Mar. 2016.
- [6] *Further Advancements for E-UTRA Physical Layer Aspects*, document TSG RAN TR 36.814 v9.0.0, 3GPP, Mar. 2010.
- [7] C. She and C. Yang, "Energy efficiency and delay in wireless systems: Is their relation always a tradeoff?" *IEEE Trans. Wireless Commun.*, vol. 15, no. 11, pp. 7215–7228, Nov. 2016.
- [8] G. Durisi, T. Koch, and P. Popovski, "Toward massive, ultrareliable, and low-latency wireless communication with short packets," *Proc. IEEE*, vol. 104, no. 9, pp. 1711–1726, Sep. 2016.
- [9] Y. Polyanskiy, H. V. Poor, and S. Verdú, "Channel coding rate in the finite blocklength regime," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2307–2359, May 2010.
- [10] W. Yang, G. Durisi, T. Koch, and Y. Polyanskiy, "Quasi-static multiple-antenna fading channels at finite blocklength," *IEEE Trans. Inf. Theory*, vol. 60, no. 7, pp. 4232–4265, Jul. 2014.
- [11] S. Xu, T.-H. Chang, S.-C. Lin, C. Shen, and G. Zhu, "Energy-efficient packet scheduling with finite blocklength codes: Convexity analysis and efficient algorithms," *IEEE Trans. Wireless Commun.*, vol. 15, no. 8, pp. 5527–5540, Aug. 2016.
- [12] S. Schiessl, J. Gross, and H. Al-Zubaidy, "Delay analysis for wireless fading channels with finite blocklength channel coding," in *Proc. ACM MSWiM*, 2015, pp. 13–22.
- [13] Y. Hu, A. Schmeink, and J. Gross, "Blocklength-limited performance of relaying under quasi-static Rayleigh channels," *IEEE Trans. Wireless Commun.*, vol. 15, no. 7, pp. 4548–4558, Jul. 2016.
- [14] C. She, C. Yang, and T. Q. S. Quek, (2017). "Cross-layer optimization for ultra-reliable and low-latency radio access networks." [Online]. Available: <https://arxiv.org/pdf/1703.09575.pdf>
- [15] J. Jia, Y. Deng, J. Chen, A.-H. Aghvami, and A. Nallanathan, "Availability analysis and optimization in CoMP and CA-enabled hetnets," *IEEE Trans. Commun.*, vol. 65, no. 6, pp. 2438–2450, Jun. 2017.
- [16] A. Aijaz, "Towards 5G-enabled tactile Internet: Radio resource allocation for haptic communications," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, Apr. 2016, pp. 1–6.
- [17] 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, pp. 72–78, Jun. 2017.
- [18] M. Condoluci, T. Mahmoodi, E. Steinbach, and M. Dohler, "Soft resource reservation for low-delayed teleoperation over mobile networks," *IEEE Access*, vol. 5, pp. 10445–10455, May 2017.
- [19] P. Schulz *et al.*, "Latency critical IoT applications in 5G: Perspective on the design of radio interface and network architecture," *IEEE Commun. Mag.*, vol. 55, no. 2, pp. 70–78, Feb. 2017.
- [20] M. Harchol-Balter, *Performance Modeling and Design of Computer Systems: Queueing Theory in Action*. Cambridge, U.K.: Cambridge Univ. Press, 2013.
- [21] G. Zhang, T. Q. S. Quek, M. Kountouris, A. Huang, and H. Shan, "Fundamentals of heterogeneous backhaul design—Analysis and optimization," *IEEE Trans. Commun.*, vol. 64, no. 2, pp. 876–889, Feb. 2016.
- [22] S. A. Ashraf, F. Lindqvist, R. Baldemair, and B. Lindoff, "Control channel design trade-offs for ultra-reliable and low-latency communication system," in *Proc. IEEE Global Commun. Conf. (GlobeCom) Workshops*, Dec. 2015, pp. 1–6.
- [23] S. A. Ashraf, I. Aktas, E. Eriksson, K. W. Helmersson, and J. Ansari, "Ultra-reliable and low-latency communication for wireless factory automation: From LTE to 5G," in *Proc. IEEE 21st Int. Conf. Emerg. Technol. Factory Autom. (ETFA)*, Sep. 2016, pp. 1–8.
- [24] H. A. Omar, W. Zhuang, A. Abdrabou, and L. Li, "Performance evaluation of VeMAC supporting safety applications in vehicular networks," *IEEE Trans. Emerg. Topics Comput.*, vol. 1, no. 1, pp. 69–83, Jun. 2013.
- [25] R. Abbas, M. Shirvanimoghaddam, Y. Li, and B. Vucetic, "Random access for M2M communications with QoS guarantees," *IEEE Trans. Commun.*, vol. 65, no. 7, pp. 2889–2903, Jul. 2017.
- [26] *Analysis on Traffic Model and Characteristics for MTC and Text Proposal*, document G. R1-120056, TSG-RAN Meeting WG1#68, Dresden, Germany, 2012.
- [27] M. Khabazian, S. Aissa, and M. Mehmet-Ali, "Performance modeling of safety messages broadcast in vehicular ad hoc networks," *IEEE Trans. Intell. Trans. Syst.*, vol. 14, no. 1, pp. 380–387, Mar. 2013.
- [28] G. L. Choudhury, D. M. Lucantoni, and W. Whitt, "Squeezing the most out of ATM," *IEEE Trans. Commun.*, vol. 44, no. 2, pp. 203–217, Feb. 1996.
- [29] D. Öhmann, M. Simsek, and G. P. Fettweis, "Achieving high availability in wireless networks by an optimal number of Rayleigh-fading links," in *Proc. IEEE Global Commun. Conf. (GlobeCom) Workshops*, Dec. 2014, pp. 1402–1407.
- [30] G. Poccovi, B. Soret, M. Lauridsen, K. I. Pedersen, and P. Mogensen, "Signal quality outage analysis for ultra-reliable communications in cellular networks," in *Proc. IEEE Global Commun. Conf. (GlobeCom) Workshops*, Dec. 2015, pp. 1–6.
- [31] X. Liu, S. Han, and C. Yang, "Energy-efficient training-assisted transmission strategies for closed-loop MISO systems," *IEEE Trans. Veh. Technol.*, vol. 64, no. 7, pp. 2846–2860, Jul. 2015.
- [32] D. Tse, *Fundamentals of Wireless Communication*. Cambridge, U.K.: Cambridge Univ. Press, 2005.
- [33] *Evolved Universal Terrestrial Radio Access (EUTRA); Further Advancements for E-UTRA Physical Layer Aspects*, document 3GPP TR 36.814, Release 9, 3GPP, 2010.
- [34] C. Sun, C. She, and C. Yang, (2017). "Energy-efficient resource allocation for ultra-reliable and low-latency communications." [Online]. Available: <https://arxiv.org/pdf/1707.09720.pdf>
- [35] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [36] C. She and C. Yang, "Available range of different transmission modes for ultra-reliable and low-latency communications," in *Proc. IEEE Veh. Tech. Conf. (VTC Spring)*, Jun. 2017, pp. 1–5.



Changyang She (S'12–M'17) received the B. Eng degree from the Honors College, Beihang University (BUAA), Beijing, China, in 2012, and the Ph.D. degree from the School of Electronics and Information Engineering, BUAA, in 2017. Since 2017, he has become a Post-Doctoral Research Fellow with the Singapore University of Technology and Design. His research interests include ultra-reliable and low-latency communications, Tactile Internet, big data for resource allocation in wireless networks, and energy efficient transmission in 5G communication systems.



Chenyang Yang (M'99–SM'08) received the Ph.D. degree from Beihang University, Beijing, China, in 1997. Since 1999, she has been a Full Professor with the School of Electronic and Information Engineering, Beihang University. She has authored over 200 international journal and conference papers and filed over 70 patents in the fields of energy efficient transmission, coordinated multi-point, interference management, and cognitive radio and relay. Her recent research interests include local caching, URLLC, and wireless big data. She was nominated as an Outstanding Young Professor of Beijing in 1995 and was supported by the First Teaching and Research Award Program for Outstanding Young Teachers of Higher Education Institutions by the Ministry of Education from 1999 to 2004. She was the Chair of the Beijing chapter of the IEEE Communications Society from 2008 to 2012 and the Membership Development Committee Chair of the Asia Pacific Board of IEEE Communications Society from 2011 to 2013. She has served as a Technical Program Committee Member for numerous IEEE conferences. She has ever served as an Associate Editor for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS and the Guest Editor for the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS and the IEEE JOURNAL ON SELECTED TOPICS IN SIGNAL PROCESSING. She is currently an Associate Editor-in-Chief of the *Chinese Journal of Communications* and the *Chinese Journal of Signal Processing*.



Tony Q. S. Quek (S'98–M'08–SM'12–F'18) received the B.E. and M.E. degrees in electrical and electronics engineering from the Tokyo Institute of Technology, and the Ph.D. degree in electrical engineering and computer science from MIT. He is currently a tenured Associate Professor with the Singapore University of Technology and Design (SUTD). He also serves as the Associate Head of ISTD Pillar and the Deputy Director of the SUTD-ZJU IDEA. His current research topics include wireless communications and networking, security, big data processing, network intelligence, and Internet-of-Things.

Dr. Quek has been actively involved in organizing and chairing sessions, and has served as a member of the technical program committee and symposium chairs in a number of international conferences. He is currently an elected member of IEEE Signal Processing Society SPCOM Technical Committee.. He is a co-author of the book *Small Cell Networks: Deployment, PHY Techniques, and Resource Allocation* (Cambridge University Press, 2013) and the book *Cloud Radio Access Networks: Principles, Technologies, and Applications* (Cambridge University Press, 2017). He was an Executive Editorial Committee Member for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, an Editor for the IEEE TRANSACTIONS ON COMMUNICATIONS and the IEEE WIRELESS COMMUNICATIONS LETTERS.

Dr. Quek was a recipient of the 2008 Philip Yeo Prize for Outstanding Achievement in Research, the IEEE GLOBECOM 2010 Best Paper Award, the 2012 IEEE William R. Bennett Prize, the 2015 SUTD Outstanding Education Awards – Excellence in Research, the 2016 IEEE Signal Processing Society Young Author Best Paper Award, the 2017 CTTC Early Achievement Award, the 2017 IEEE ComSoc AP Outstanding Paper Award, the 2017 Clarivate Analytics Highly Cited Researcher, and the 2017 Communication Theory Technical Committee Early Achievement Award. He is currently a Distinguished Lecturer of the IEEE Communications Society IEEE.