# On the Impacts of Hybrid Beamforming on Millimeter Wave Cellular Network Optimization

Jia Liu[†]

[†]Department of Electrical and Computer Engineering, The Ohio State University

## Abstract

Millimeter wave communication (mmWave) has recently emerged as a key technology for building 5G wireless networks and beyond. To reconcile the conflict between the large antenna arrays and the limited amount of RF chains in mmWave systems, hybrid beamforming becomes a promising beamforming architecture. However, existing research on hybrid beamforming focused mostly on the physical layer or signal processing aspects. So far, there is a lack of theoretical understanding on how hybrid beamforming could affect mmWave *network optimization*. In this paper, we study the impacts of hybrid beamforming on utility-optimality and queueing delay in mmWave cellular networks. Our contributions in this paper are three-fold: i) we develop a joint hybrid beamforming and congestion control algorithmic framework for mmWave network utility maximization; ii) we reveal a pseudoconvexity structure in the hybrid beamforming scheduling problem, which leads to simplified analog beamforming protocol design; and iii) we theoretically characterize the scalings of utility-optimality and delay with respect to channel state information (CSI) accuracy in digital beamforming.

## 1 Introduction

In recent years, millimeter wave communication (mmWave) has emerged as a promising technology for building 5G wireless networks and beyond. The excitements of mmWave communications are primarily due to: i) the rich unlicensed spectrum resources in 60 GHz bands; ii) the ease of packing large antenna arrays into small form factors (a consequence of the short wavelengths); and iii) a much simplified interference management thanks to the highly directional "pencil-beam-like" mmWave signals. Moreover, recent field tests (see, e.g., [1, 2], etc.) have shown that the large directivity gains of mmWave transceivers can offset the high atmospheric attenuation in mmWave bands, dispelling the common concern that mmWave is not suitable for outdoor communications. The potential of mmWave networks has also stimulated many standardization activities (e.g., IEEE 802.15.3 wireless personal area networks, 802.11ad wireless local area networks, and fast-growing interests in mmWave cellular networks [3]).

1

However, the highly directional propagation of mmWave signals and the special mmWave hardware requirements also introduce several unique technical challenges for network systems. One major problem in mmWave networking is its vulnerability to blockage, which is due to the weak diffraction ability of mmWave communications [3, 4]. Mitigating blockage in mobile cellular networks requires a frequent search for new unblocked directed spatial paths, which entails a large communication overhead and complicates the scheduling and congestion control algorithmic designs at higher layers.

Another main technical challenge is the beamforming architecture design, which lies at the heart of mmWave directional networking. Although large antenna arrays can be easily deployed in mmWave systems, the high power consumption of mixed mmWave signal components significantly limits the number of radio frequency hardware chains (RF chains), rendering full digital beamforming (requiring one RF chain per antenna) impractical [5]. Moreover, most of the digital beamforming schemes in traditional MIMO systems require full channel state information (CSI), which is difficult to acquire in mmWave systems due to the fast fading in mmWave spectrum and the low signal-to-noise ratio (SNR) before beamforming [6]. Because of the RF chain limitations in mmWave systems, analog beamforming approaches have been proposed (see, e.g., [7, 8]). The basic idea of analog beamforming is to control the phase shifters of antenna elements, so that the energy of the transmitted data stream is concentrated in a single direction to obtain a high directivity gain. Compared to digital beamforming, analog beamforming can be achieved by only one RF chain without requiring any CSI at the transmitter. However, analog beamforming can only transmit in a single beam direction and cannot leverage any spatial multiplexing capability of the large mmWave antenna array.

In light of the limitations of analog and digital beamformings, there is a growing consensus that the more suitable architecture for mmWave cellular networks is the *hybrid beamforming* architecture, which exploits the large mmWave antenna arrays and yet only requires a limited number of RF chains [6, 9–12]. Hybrid beamforming enjoys the best of both worlds: On one hand, it uses analog beamforming to offer spatial division and directivity gains to combat large mmWave channel attenuations. On the other hand, digital beamforming provides multiplexing gains for the lower dimensional *effective channels*, for which the CSI is relatively easier to acquire. It has been shown in [6, 13] that hybrid beamforming achieves a comparable data rate performance compared to full digital beamforming with 8 to 16 times fewer RF chains.

So far, however, the existing works on mmWave hybrid beamforming are mostly concerned with problems at the physical layer or signal processing aspects. To date, there remains a lack of theoretical understanding on how hybrid beamforming could affect the performances of mmWave network control, scheduling, and resource optimization algorithms. In this paper, our goal is to fill this gap by conducting an in-depth study on the impacts of hybrid beamforming on throughput and delay performances in mmWave cellular network optimization.

Specficaly, in this paper, we focus on the algorithmic design and the throughput-delay analysis for the celebrated queue-length-based congestion control and scheduling framework (QCS) (see, e.g., [14, 15], and [16] for a survey) in hybrid-beamforming-based mmWave cellular networks. Our main results and technical contributions are as follows:

- We develop an accurate analytical model that captures the essence of hybrid beamforming in mmWave cellular networks, while being tractable enough to enable network-level understanding and analysis. Based on this analytical model, we formulate the problem of joint hybrid beamforming and congestion control for network utility maximization. We show that the joint hybrid beamforming and congestion control optimization is non-convex by nature, which creates challenges for the algorithmic designs in the MaxWeight scheduling component in the QCS framework.

- By exploiting the special problem structure of the mmWave MaxWeight scheduling component, we show that the non-convex scheduling subproblem admits a *pseudoconvex approximation* under a wide range of hybrid beamforming parameters of practical interests. Moreover, our analysis reveals that, to solve the scheduling subproblem, one only needs to adjust the analog beamwidth at the base station (BS), while the analog beamwidth adjustment at the mobile station (MS) side is unnecessary. This insight greatly simplifies the analog beamforming training protocol design.

- We investigate the impact of CSI inaccuracy on network performance with hybrid beamforming, where we assume that the true CSI is quantized by $Q$ bits. We reveal a pair of interesting phase transition phenomena in utility-optimality and delay in the following sense: There exists a critical value $Q^\sharp$ such that: i) if $0 < Q < Q^\sharp$, then the deviations of steady-state queue-length grows linearly and congestion control rate is bounded by a constant; ii) If $Q \geq Q^\sharp$, the deviations of queue-lengths and congestion control rates have the same scaling laws as in the full CSI case.

Collectively, these results not only deepen our theoretical understanding of mmWave network optimization with hybrid beamforming, but also provide insights for low-complexity analog beam training and effective CSI quantization in practice. The remainder of this paper is organized as follows: In Section 2, we introduce network models and the problem formulation. Section 3 presents the mmWave congestion control and scheduling framework, as well as the algorithmic design for analog beam training. Section 4 studies the impacts of inaccurate CSI on digital beamforming. Section 5 provides numerical results and Section 6 concludes this paper.

## 2    Network Model and Problem Formulation

**Notation:**    We use boldface to denote matrices/vectors. $\mathbf{A}^\dagger$ denotes the conjugate transpose of $\mathbf{A}$. We use $\|\cdot\|$ and $\|\cdot\|_1$ to denote $L^2$- and $L^1$-norms, respectively. We let $\mathbf{I}$ denote the identity
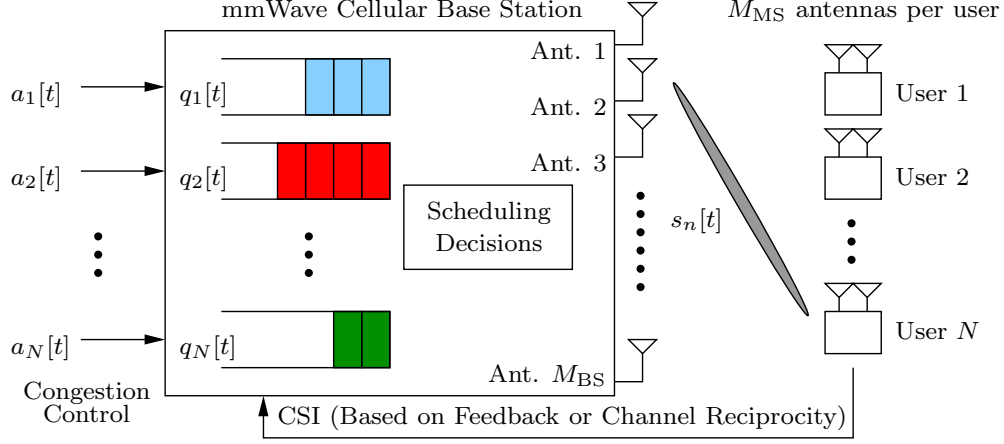
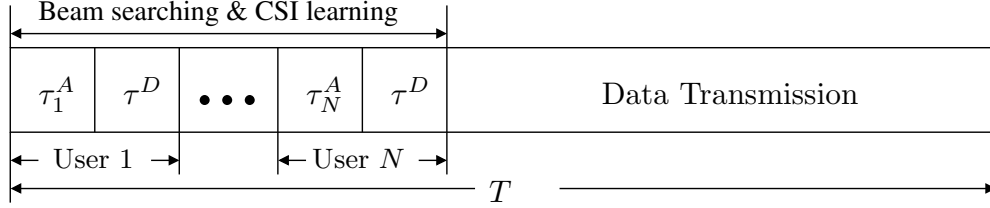Figure 1: A mmWave cellular downlink with a $M_{\mathrm{BS}}$-antenna base station and $N$ $M_{\mathrm{MS}}$-antenna users.



Figure 2: Frame structure of a time-slot in mmWave cellular networks with hybrid beamforming.

matrix, whose dimension is conformal to the context. We let $\mathbb{R}$ and $\mathbb{C}$ denote real and complex spaces, respectively.

**1) Hybrid-Beamforming-Based mmWave Downlink:** As shown in Fig. 1, we consider a mmWave cellular downlink system with $N$ users. The BS and each user have $M_{\mathrm{BS}}$ and $M_{\mathrm{MS}}$ antennas, respectively. The mmWave downlink adopts a hybrid beamforming architecture with $M_{\mathrm{RF}}^{\mathrm{B}}$ and $M_{\mathrm{RF}}^{\mathrm{M}}$ RF chains at the BS and each user's MS, respectively (see Fig. 3). The system operates in a time-slotted mode. The time-slots are indexed by $t \in \{0, 1, 2, \ldots\}$. As shown in Fig. 2, each time-slot is of period $T$ and contains two phases. The first phase is further divided into $N$ mini-slots corresponding to the $N$ users. Each mini-slot contains two parts $\tau_n^A$ and $\tau_n^D$. In $\tau_n^A$, both the BS and user $n$ perform analog beam search to refresh their beam directions to mitigate link breakage caused by user $n$'s movements. In $\tau_n^D$, the BS estimates the CSI of user $n$ for digital beamforming. In the data transmission phase, based on the analog beam and digital CSI training results, the BS picks one of the $N$ users and steers an analog beam to this user. Likewise, the scheduled user also steers an analog beam toward the BS. Further, by leveraging the learned CSI to perform spatial multiplexing, the BS and the scheduled user communicate via $K$ data streams. For mmWave systems in practice, we usually have: i) $K \leq M_{\mathrm{RF}}^{\mathrm{B}} \leq M_{\mathrm{BS}}$; ii) $K \leq M_{\mathrm{RF}}^{\mathrm{M}} \leq M_{\mathrm{MS}}$; iii) $M_{\mathrm{RF}}^{\mathrm{M}} \leq M_{\mathrm{RF}}^{\mathrm{B}}$; and iv) $M_{\mathrm{MS}} \leq M_{\mathrm{BS}}$.

**2) Hybrid Beamforming Achievable Rate Region:** Suppose that user $n$ is scheduled
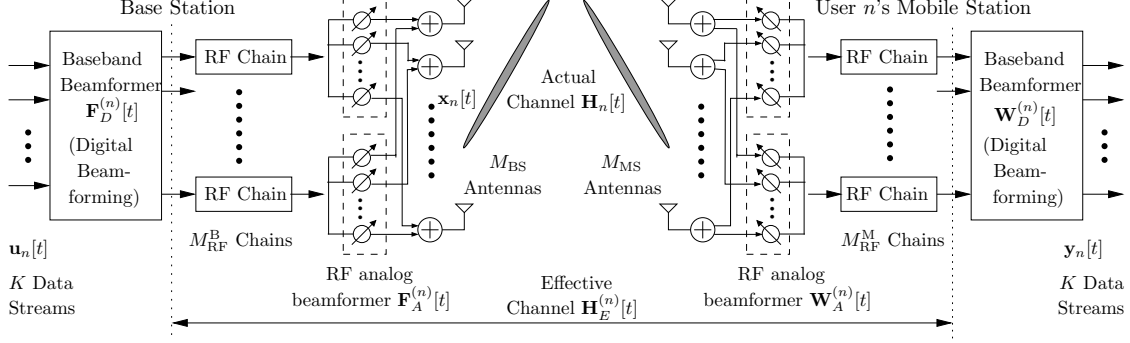
Figure 3: Block diagram of a mmWave cellular network with hybrid beamforming.

in time-slot $t$. As shown in Fig. 3, the BS applies a digital baseband beamformer $\mathbf{F}_D^{(n)}[t] \in \mathbb{C}^{M_{\mathrm{RF}}^{\mathrm{B}} \times K}$, which is followed by an analog RF beamformer $\mathbf{F}_A^{(n)}[t] \in \mathbb{C}^{M_{\mathrm{BS}} \times M_{\mathrm{RF}}^{\mathrm{B}}}$. Here, $\mathbf{F}_A^{(n)}[t]$ is implemented *only* by RF phase shifters. Then, the transmitted signal $\mathbf{x}_n[t]$ can be written as $\mathbf{x}_n[t] = \mathbf{F}_A^{(n)}[t]\mathbf{F}_D^{(n)}[t]\mathbf{u}_n[t]$, where $\mathbf{u}_n[t] \in \mathbb{C}^K$ is the transmitted symbol of user $n$ and satisfies $\mathbb{E}\{\mathbf{u}_n[t]\mathbf{u}_n^\dagger[t]\} = (P_{\mathrm{tot}}/K)\mathbf{I}$, where $P_{\mathrm{tot}}$ denotes the total transmit power.

We let $\mathbf{H}_n[t] \in \mathbb{C}^{M_{\mathrm{MS}} \times M_{\mathrm{BS}}}$ denote the channel matrix in time-slot $t$ between the BS and user $n$. We assume independent quasi-static block fading, i.e., each entry in $\mathbf{H}_n[t]$ is constant in one time-slot and independently varies in the next time-slot. Then, the received signal $\mathbf{r}_n[t]$ at user $n$ can be written as: $\mathbf{z}_n[t] = \mathbf{H}_n[t]\mathbf{F}_A^{(n)}[t]\mathbf{F}_D^{(n)}[t]\mathbf{u}_n[t] + \mathbf{n}[t]$, where $\mathbf{n}[t]$ is the Gaussian noise at the receiver. As seen in Fig. 3, in time-slot $t$, user $n$ applies analog and digital beamformers $\mathbf{W}_A^{(n)}[t]$ and $\mathbf{W}_D^{(n)}[t]$ to process $\mathbf{z}_n[t]$ to yield the baseband signal:

$$\mathbf{y}_n[t] = \mathbf{W}_D^{(n)\dagger}[t]\mathbf{H}_E^{(n)}[t]\mathbf{F}_D^{(n)}[t]\mathbf{u}_n[t] + \tilde{\mathbf{n}}[t], \tag{1}$$

where $\mathbf{H}_E^{(n)}[t] \triangleq \mathbf{W}_A^{(n)\dagger}[t]\mathbf{H}_n[t]\mathbf{F}_A^{(n)}[t] \in \mathbb{C}^{M_{\mathrm{RF}}^{\mathrm{M}} \times M_{\mathrm{RF}}^{\mathrm{B}}}$ denotes the *effective channel* after analog beamforming and $\tilde{\mathbf{n}}[t] \triangleq \mathbf{W}_D^{(n)\dagger}[t]\mathbf{W}_A^{(n)\dagger}[t]\mathbf{n}[t]$ denotes the effective noise. It is clear from (1) that the hybrid beamforming achievable rate depends on the choices of analog and digital beamformers $\mathbf{F}_A^{(n)}[t]$, $\mathbf{F}_D^{(n)}[t]$, $\mathbf{W}_A^{(n)}[t]$, and $\mathbf{W}_D^{(n)}[t]$. In what follows, we will present the models for analog and digital beamforming.

*a) Analog beamforming process:* In time-slot $t$, $\mathbf{F}_A^{(n)}[t]$ and $\mathbf{W}_A^{(n)}[t]$ are determined by a beam training process, during which the BS and user $n$ search over all possible direction combinations within their corresponding sectors[1], as shown in Fig. 4 (this exhaustive beam training process is compatible with IEEE 802.11ad and IEEE 802.15.3c).

Let $T_p$ denote the time required for transmitting and receiving a pilot symbol. Let $\psi^{\mathrm{B}}$ and $\psi_n^{\mathrm{M}}$ denote the sector-level beamwidth at the BS and user $n$, respectively. Also, let $\theta_B[t]$ and $\theta_n[t]$

---

[1]In this paper, we assume that both the BS and user know the sectors of each other's location in each time-slot. This is a reasonable assumption because the sector information can be inferred with high accuracy from the beam direction in the previous time-slot and the mobility/trajectory information of the user.
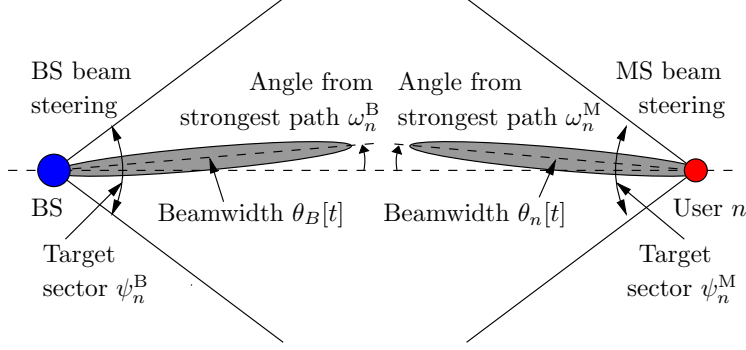
Figure 4: The analog beamforming training procedure.

denote the beam-level beamwidth at the BS and user $n$, respectively. Then, the beam search time $\tau_n^A$ can be computed as: $\tau_n^A = \frac{\psi_n^B}{\theta_B[t]} \frac{\psi_n^M}{\theta_n[t]} T_p$.

In this paper, we adopt a widely used sectored antenna pattern model (see, e.g., [17–19]): We assume that the gains are a constant for all angles within the main lobe and equal to a smaller constant in the side lobes. As shown in Fig. 4, we let $\omega_n^B$ and $\omega_n^M$ represent the angles deviating from the strongest path between the BS and user $n$, respectively (the strongest path needs not be line-of-sight and Fig. 4 is only for illustrative purposes). We let $g_n^B(\omega_n^B, \theta_B[t])$ and $g_n^M(\omega_n^M, \theta_n[t])$ denote the transmission and reception gains at the BS and user $n$, which are modeled as follows [17–19]:

$$g_n^B(\omega_n^B, \theta_B[t]) = \begin{cases} \frac{2\pi - (2\pi - \theta_B[t])\eta}{\theta_B[t]}, & \text{if } |\omega_n^B| \leq \frac{\theta_B[t]}{2}, \\ \eta, & \text{otherwise,} \end{cases} \tag{2}$$

$$g_n^M(\omega_n^M, \theta_n[t]) = \begin{cases} \frac{2\pi - (2\pi - \theta_n[t])\eta}{\theta_n[t]}, & \text{if } |\omega_n^M| \leq \frac{\theta_n[t]}{2}, \\ \eta, & \text{otherwise,} \end{cases} \tag{3}$$

where $\eta \in [0, 1)$ is the side lobe gain. In practice, $\eta \ll 1$ for narrow beams (i.e., $\theta_B[t]$ and $\theta_n[t]$ are small). This model captures the essential features of antenna patterns (e.g., directivity gains, front-to-back ratio, and half-power beamwidth, etc.). The beam training is finished when the BS and the user's beams are aligned with the strongest path, i.e., the conditions $|\omega_n^B| \leq \frac{\theta_B[t]}{2}$ in (2) and $|\omega_n^M| \leq \frac{\theta_n[t]}{2}$ in (3) are satisfied.

*b) Digital beamforming process:* Once the analog beam search is completed, the analog beam-formers $\mathbf{F}_A^{(n)}[t]$ and $\mathbf{W}_D^{(n)}[t]$ are known. Therefore, we can estimate the CSI of the effective channel $\mathbf{H}_E^{(n)}[t]$, which is assumed to take $\tau^D = \beta T_p$ amount of time (cf. Fig. 2), where $\beta > 0$ is some constant. With the learned CSI, the BS and user $n$ jointly choose baseband beamformers $\mathbf{F}_D^{(n)}[t]$ and $\mathbf{W}_D^{(n)}[t]$ by some digital beamforming strategies. For example, we can use singular value decomposition (SVD), where the columns of $\mathbf{F}_D^{(n)}[t]$ and $\mathbf{W}_D^{(n)}[t]$ are set to $\mathbf{H}_E^{(n)}[t]$'s right and left singular vectors, respectively. The SVD scheme is well known to be capacity-achieving for arbitrary $\mathbf{H}_E^{(n)}[t]$ [20]. As another example, if $M_{RF}^B \geq M_{RF}^M$, one can let $\mathbf{W}_D^{(n)}[t] = \mathbf{I}$ and let $\mathbf{F}_D^{(n)}[t]$

be the pseudoinverse of $\mathbf{H}_E^{(n)}[t]$. This scheme is known as "zero-forcing beamforming" and it is asymptotically capacity-achieving under high SNR regime [21].

One particularly interesting case arises when $M_{\mathrm{RF}}^{\mathrm{B}} \gg M_{\mathrm{RF}}^{\mathrm{M}}$. In this case, the row vectors in the effective channel $\mathbf{H}_E^{(n)}[t]$ are asymptotically orthogonal to each other as $M_{\mathrm{RF}}^{\mathrm{B}}$ gets large. Thanks to this nice property, one can let $\mathbf{F}_D^{(n)}[t] = \mathbf{H}_E^{(n)}[t]^\dagger$ and $\mathbf{W}_D^{(n)}[t] = \mathbf{I}$, which is called "conjugate beamforming." It has been shown that conjugate beamforming is also asymptotically capacity-achieving in the high SNR regime [22]. We will further discuss conjugate beamforming in Section 4.

Regardless the choice of digital beamforming schemes, the digital beamforming process converts $\mathbf{H}_E^{(n)}[t]$ into $K \leq \min\{M_{\mathrm{RF}}^{\mathrm{B}}, M_{\mathrm{RF}}^{\mathrm{M}}\}$ spatial channels (depending on the rank of $\mathbf{H}_E^{(n)}[t]$). We let $g_n^{(k)}[t]$ denote the effective gain of the $k$-th spatial channel. Based on the models of hybrid analog/digital beamforming, we have that the hybrid beamforming achievable rate of user $n$ can be computed as:[2]

$$r_n(\theta_B[t], \theta_n[t]) = \left(1 - \frac{\tau^A + N\tau^D}{T}\right) \sum_{k=1}^K \log_2\left(1 + \frac{P_{\max}}{KN_0} g_n^{\mathrm{B}}(\omega_n^{\mathrm{B}}, \theta_B[t]) g_n^{\mathrm{M}}(\omega_n^{\mathrm{M}}, \theta_n[t]) g_n^{(k)}[t]\right), \quad (4)$$

where $\tau^A \triangleq \sum_{n=1}^N \tau_n^A$ and $P_{\max}$ denotes the maximum transmission power at the BS. Then, for a given channel state in time-slot $t$, we let $\mathcal{C}_n[t]$ denote the instantaneous achievable rate region under a chosen digital beamforming scheme:

$$\mathcal{C}_n[t] \triangleq \left\{ r_n(\theta_B[t], \theta_n[t]) \,\middle|\, \begin{array}{l} \theta_B[t] \in (0, \psi_n^{\mathrm{B}}], \\ \theta_n[t] \in (0, \psi_n^{\mathrm{M}}]. \end{array} \right\}. \quad (5)$$

It can be seen from (4) that the beamwidths $\theta_B[t]$ and $\theta_n[t]$ need to be chosen judiciously: On one hand, from (2) and (3), $g_n^{\mathrm{B}}(\omega_n^{\mathrm{B}}, \theta_B[t])$ and $g_n^{\mathrm{M}}(\omega_n^{\mathrm{M}}, \theta_n[t])$ increase as $\theta_B[t]$ and $\theta_n[t]$ decrease, leading to a higher SNR and hence a higher data rate. However, the smaller the beamwidths $\theta_B[t]$ and $\theta_n[t]$, the shorter the transmission phase, i.e., there exists a trade-off between data rate and transmission time.

**3) Queueing Model:** As shown in Fig. 1, the BS maintains a separate queue for each user. Let $a_n[t]$ denote the number of packets injected into queue $n$ in time-slot $t$. The arrival processes $\{a_n[t]\}$, $\forall n$, are controlled by a congestion controller. We assume that there exists a finite constant $A^{\max}$ such that $a_n[t] \leq A^{\max}$, $\forall n, t$. Let $\mathbf{s}[t] \triangleq [s_1[t], \ldots, s_N[t]]^\top$ denote the scheduled service rate vector in time-slot $t$ (the scheduling algorithm that determines $\mathbf{s}[t]$ will be presented in Section 3). Then, the queue-length of user $n$ evolves as: $q_n[t+1] = \left(q_n[t] - s_n[t] + a_n[t]\right)^+$, $\forall n$, where $(\cdot)^+ \triangleq \max(0, \cdot)$. Let $\mathbf{q}[t] = [q_1[t]], \ldots, q_N[t]]^\top$. In this paper, we adopt the following notion of queue-stability (same as in [14, 15]): We say that a network is *stable* if the steady-state total queue-length is finite, i.e., $\limsup_{t\to\infty} \mathbb{E}\{\|\mathbf{q}[t]\|_1\} < \infty$.

---

[2]In this paper, equal power allocation is used to trade for lower rate evaluation complexity in the effective MIMO channel, although water-filling is the optimal power allocation solution for MIMO channels [20]. Equal power allocation is desirable in practice because it has been shown that the rate loss due to equal power allocation is negligible under high SNR and that equal power allocation is asymptotically capacity-achieving as SNR increases [21].

**5) Problem Formulation:** Let $\bar{a}_n \triangleq \lim_{T\to\infty} \frac{1}{T}\sum_{t=0}^{T-1} a_n[t]$ denote the average controlled arrival rate of user $n$. We associate each user $n$ with a utility function $U_n(\bar{a}_n)$, which is assumed to be strongly concave, increasing, and twice continuously differentiable. $U_n(\bar{a}_n)$ represents the utility gained by user $n$ when data is injected at rate $\bar{a}_n$. Then, the joint congestion control and scheduling (JCS) optimization problem for network utility maximization can be written as:

$$\textbf{JCS: } \text{Maximize} \quad \sum_{n=1}^{N} U_n(\bar{a}_n)$$
$$\text{subject to Queue-length stability constraints,}$$
$$s_n[t] \in \mathcal{C}_n[t], \ a_n[t] \in [0, A^{\max}] \ \forall n, t.$$

In Section 3, we will first consider the algorithmic design for solving Problem JCS under perfect CSI. Then, in Section 4, we will conduct an in-depth investigation on the impacts of CSI inaccuracy on throughput and delay.

# 3 Algorithmic Design under Perfect CSI

Because of the utility maximization formulation, Problem JCS can be solved by the well-known queue-length-based congestion control and scheduling (QCS) framework (see, e.g., [14–16]) in the following sense: The congestion control rate $\bar{\mathbf{a}}$ achieves an optimality gap $O(1/K)$ at the price of an $O(K)$ queue-length, where $K > 0$ is a system parameter. Hence, the utility-optimality gap can be made arbitrarily small by increasing $K$. Now, consider the following QCS algorithm specialized for hybrid beamforming in mmWave networks:

## 3.1 The QCS Algorithm Specialized for Hybrid Beamforming

**Algorithm 1:** Queue-Length-Based Congestion Control and Scheduling in mmWave Cellular Network.

**Initialization:** Choose parameters $K > 0$. Set $t = 0$.

**Main Loop:**

1. *MaxWeight Scheduler:* In time $t \geq 1$, given queue-lengths $\mathbf{q}[t]$ and CSI $\mathbf{H}[t]$, the scheduler chooses a service rate vector $\mathbf{s}[t]$ from $\mathcal{C}_n[t]$ by hybrid beamforming such that:

$$\mathbf{s}[t] = \underset{r_n \in \mathcal{C}_n[t], \forall n}{\arg\max} \left\{ \sum_{n=1}^{N} q_n[t] r_n \right\}, \tag{6}$$

   where $\mathcal{C}_n[t]$ is defined in (5).

2. *Congestion Controller:* Given queue-lengths $\mathbf{q}[t]$, the congestion controller chooses data injection rates $a_n[t]$, $\forall n$, which are integer-valued random variables satisfying:

$$\mathbb{E}\{a_n[t]|q_n[t]\} = \min\left\{ U_n'^{-1}\left(\frac{q_n[t]}{K}\right), A^{\max} \right\}, \tag{7}$$

$$\mathbb{E}\{a_n^2[t]|q_n[t]\} \leq A_2^{\max} < \infty, \quad \forall q_n[t], \tag{8}$$

where $U_n^{'-1}(\cdot)$ represents the inverse function of first-order derivative of $U_n(\cdot)$. In (7) and (8), $A^{\max}$ and $A_2^{\max}$ are some predefined sufficiently large positive constants.

3. *Queue-Length Updates:* Update the queue-lengths as $q_n[t+1] = (q_n[t] - s_n[t] + a_n[t])^+$, $\forall n$. Let $t = t + 1$. Go to Step 1 and repeat the process.

Although being optimal, the QCS framework has a major limitation in that the MaxWeight scheduling problem is difficult to solve in general and could be NP-Hard in many wireless networks [16]. Surprisingly, in what follows, we will show that the physical layer properties of hybrid beamforming in certain settings imply several special mathematical structures that leads to efficient solution for the MaxWeight subproblem.

## 3.2   The MaxWeight Scheduling Subproblem

To solve the MaxWeight scheduling subproblem, we start by examining the properties of the set of instantaneous hybrid beamforming achievable rates $\{r_n \in \mathcal{C}_n[t], \forall n\}$. First, we note that the BS forms only one beam to one of the $N$ users in each time-slot, say user $n$. This implies that $r_{n'} = 0$, $\forall n' \neq n$. Hence, the MaxWeight problem in (6) can be simplified as:

$$\max_{r_n \in \mathcal{C}_n[t], \forall n} \left\{ \sum_{n=1}^N q_n[t] r_n \right\} = \max_{n \in \{1,\dots,N\}} \left\{ q_n[t] r_n \big| r_n \in \mathcal{C}_n[t] \right\}$$

$$\overset{(a)}{=} \max_{n \in \{1,\dots,N\}} \left\{ q_n[t] \left[ \max_{\theta_B[t], \theta_n[t]} \left\{ r_n(\theta_B[t], \theta_n[t], ) \right\} \right] \right\}, \tag{9}$$

where $(a)$ follows from the fact that $r_n(\theta_B[t], \theta_n[t])$ does not depend on $q_n[t]$. As a result, solving the MaxWeight scheduling problem in (6) boils down to first solving the inner rate maximization problem in (9) for each user, and then choosing the user who has the largest rate-queue-length product. Toward this end, we explicitly write down the inner maximization problem in (9) for each user $n$ as follows:

$$\text{Maximize} \qquad \left( 1 - \frac{\tau^A + \tau^D}{T} \right) \sum_{k=1}^K \log_2 \left( 1 + \frac{P_{\max}}{K N_0} g_n^B(\omega_n^B, \theta_B[t]) \right.$$

$$\left. g_n^M(\omega_n^M, \theta_n[t]) g_n^{(k)}[t] \right) \tag{10}$$

$$\text{subject to} \qquad \theta_B[t] \in (0, \psi_n^B], \ \theta_n[t] \in (0, \psi_n^M].$$

Unfortunately, due to the multiplication between the time fraction and the rate in the objective function, Problem (10) falls into the class of polynomial programming problems, which is non-convex and NP-Hard [23]. In this paper, we consider a slightly modified and yet practically relevant homogenous setting: the BS can coordinate all users in each time-slot to ensure $\theta_n[t] = \theta_M[t]$, $\forall n$, where $\theta_M[t]$ denotes the common beamwidth of all users in time $t$. Under the homogenous setting, Problem (10) *remains* a non-convex polynomial program. However, it turns out that if the side lobe gain $\eta$ is small, Problem (10) can be approximated by a univariate pseudoconvex problem. We state this result as follows:

**Lemma 1** (Univariate approximation)**.** *If the side lobe gain satisfies $\eta \ll \frac{1}{3}$, then Problem (10) can be approximated by the following univariate optimization problem:*

$$\text{Maximize} \qquad \left(b_0 - \frac{b_1}{\tilde{\theta}[t]}\right) \sum_{k=1}^{K} \log_2\left(1 + \frac{4\pi^2 c_n^{(k)}}{\tilde{\theta}[t]}\right) \qquad (11)$$

$$\text{subject to} \qquad \tilde{\theta}[t] \in \left[\frac{b_1}{b_0}, \ \psi_n^{\text{B}}\psi_n^{\text{M}}\right],$$

*where $b_0 \triangleq 1 - (N\beta T_p/T)$, $b_1 \triangleq \frac{T_p}{T}\sum_{n=1}^{N}\psi_n^{\text{B}}\psi_n^{\text{M}}$, and $c_n^{(k)} \triangleq (P_{\max}/KN_0)g_n^{(k)}[t]$ are constants.*

Lemma 1 can be proved by substituting the antenna radiation pattern model in (2) and (3) into the objective function of Problem (10) and then exploiting the condition $\eta \ll \frac{1}{3}$ to simplify. We relegate the proof details to the appendix. With Lemma 1, we state the first main result of this paper as follows:

**Theorem 1** (Pseudoconvexity[3] of the approximation)**.** *Problem (11) is a pseudoconvex optimization problem. Moreover, if $T_p \ll T$, then Problem (11) has one unique maximum achieved in the interior of $[b_1/b_0, \psi_n^{\text{B}}\psi_n^{\text{M}}]$.*

Since Problem (11) is a maximization problem with one simple box constraint, showing its pseudoconvexity is equivalent to showing the pseudoconcavity of the objective function. We refer readers to the appendix for proof details.

**Remark 1.** Three remarks of Lemma 1 and Theorem 1 are in order: i) In practice, the conditions $T_p \ll T$ and $\eta \ll \frac{1}{3}$ can usually be satisfied because a pilot symbol is much shorter compared to a time-slot and the mmWave beams are sharp; ii) The pseudoconvex (which further implies strictly quasiconvex) and univariate properties suggest that Problem (11) can be solved by simple one-dimensional line search methods [23, Theorem 8.1.1] (e.g., the bisection or the golden section methods); iii) It can be seen from the proof of Lemma 1 that we have defined $\tilde{\theta}[t] \triangleq \theta_B[t]\theta_M[t]$. Note that the optimal objective value of Problem (11) is only a function of $\tilde{\theta}^*[t]$ and does not depend on the specific values of $\theta_B[t]$ and $\theta_M[t]$, as long as their product is equal to $\tilde{\theta}^*[t]$. This implies that we can simply set $\theta_M[t]$ to some appropriate fixed value and only adjust $\theta_B[t]$ at the BS side. In other words, there is no need to jointly adjust $\theta_B[t]$ and $\theta_M[t]$. This insight greatly simplifies the protocol designs in the analog beamforming phase. □

Collectively, the results in this section provide an algorithmic solution to Problem JCS assuming that the CSI learned in $\tau^D$ (hence the digital beamforming gains $g_n^{(k)}[t]$) is accurate. However, it remains unclear how the network utility and delay performance of Algorithm 1 will be affected if the CSI is inaccurate. This problem will be addressed in the next section.

---

[3]In convex analysis, a function $f : S \in \mathbb{R}^N \to \mathbb{R}$ is said to be pseudoconvex if for each $\mathbf{x}_1, \mathbf{x}_2 \in S$, $\nabla f(\mathbf{x}_1)^\top(\mathbf{x}_2 - \mathbf{x}_1) \geq 0$ implies $f(\mathbf{x}_2) \geq f(\mathbf{x}_1)$ or equivalently $\nabla f(\mathbf{x}_2)^\top(\mathbf{x}_2 - \mathbf{x}_1) \geq 0$. The function $f$ is said to be pseudoconcave if $-f$ is pseudoconvex.

# 4 The Impacts of Inaccurate CSI on the QCS Algorithm with Hybrid Beamforming

Generally speaking, in traditional multi-antenna networks, CSI is measured at each MS based on pilot symbol training and then fed back to the BS. However, due to the short coherence time of mmWave channels (around an order of magnitude lower than that of microwave bands since Doppler shifts scale linearly with frequencies [3]), this traditional CSI feedback approach is not suitable for mmWave-based cellular networks. Another CSI acquisition method is to have the system run in time-division duplexing (TDD) mode. Based on the channel reciprocity, the uplink CSI measured at the BS will be used for downlink transmissions. However, the limited transmit power at the MSs and the lack of beamforming gains for the uplink pilot symbols limit the accuracy of TDD-based CSI estimation. Further, the short coherence time of mmWave bands implies that the channel reciprocity assumption is only valid for low-mobility scenarios. Given these CSI estimation challenges in mmWave cellular systems, it is likely that the CSI learned during the $\tau^D$ period (cf. Fig. 2) is inaccurate.

In studying the impacts of CSI inaccuracy, we are interested in the case where the number of RF chains at the BS is much greater than that at the MSs, (e.g., tens of times larger). This setting is particularly interesting because it is the relevant case for mmWave cellular networks in practice: First, as mentioned in Section 1, large antenna arrays can be easily deployed at the BS due to the short wavelengths of mmWave bands. Also, because of the physical size and power constraints, the MSs usually accommodate much fewer RF chains compared to that at the BS side. Moreover, as noted in many studies [22], CSI acquisition is one of the most fundamental limiting factors in the system designs of large-scale antenna cellular systems. In what follows, we start with the digital beamforming for effective mmWave channels with a large number of RF chains at the BS and its operations under a limited CSI model.

**1) Digital Beamforming for Effective Channels with $M_{\mathrm{RF}}^{\mathrm{B}} \gg M_{\mathrm{RF}}^{\mathrm{M}}$:** As mentioned in Section 2, due to the near orthogonality between the rows in the effective channel in this case, the simple conjugate digital beamforming technique can be used. Recall that the received signal of user $n$ can be written as: $\mathbf{y}[t] = \mathbf{W}_D^{(n)\dagger}[t]\mathbf{H}_E^{(n)}[t]\mathbf{F}_D^{(n)}[t]\mathbf{s}_n[t] + \tilde{\mathbf{n}}[t]$, where $\mathbf{H}_E^{(n)}[t] \in \mathbb{C}^{M_{\mathrm{RF}}^{\mathrm{M}} \times M_{\mathrm{RF}}^{\mathrm{B}}}$ is the effective channel by taking into account the effects of analog beamforming; and $\mathbf{F}_D^{(n)}[t]$ and $\mathbf{W}_D^{(n)\dagger}[t]$ are the transmit and receive digital beamformers, respectively. Under conjugate beamforming, we let $\mathbf{W}_D^{(n)\dagger}[t] = \mathbf{I}$ and $\mathbf{F}_D^{(n)}[t] = \mathbf{H}_E^{(n)}[t]^{\dagger}$, i.e., the conjugate transpose of $\mathbf{H}_E^{(n)}[t]$. Also, we assume that the effective channel $\mathbf{H}_E^{(n)}[t]$ is of full row rank so that we can let $K = M_{\mathrm{RF}}^{\mathrm{M}}$ (i.e., all receiver RF chains are utilized). Then, the achievable rate under digital conjugate beamforming can be computed as:

$$r_n[t] \approx \left(1 - \frac{\tau^A + \tau^D}{T}\right)\sum_{k=1}^{K}\log_2\left(1 + \frac{P_{\max}}{KN_0}\|\mathbf{h}_{E,k}^{(n)}[t]\|^2\right), \tag{12}$$

where $\mathbf{h}_{E,k}^{(n)}[t]$ denotes the $k$-th row of $\mathbf{H}_E^{(n)}[t]$. In (12), the approximation holds because the rows of $\mathbf{H}_E^{(n)}[t]$ are nearly orthogonal as $M_{\mathrm{RF}}^{\mathrm{B}}$ gets large. We note here that (12) is equivalent to (4) since $\mathbf{h}_{E,k}^{(n)}[t]$ has absorbed the gains $g_n^{\mathrm{B}}(\omega_n^{\mathrm{B}}, \theta_B[t])$ and $g_n^{\mathrm{M}}(\omega_n^{\mathrm{M}}, \theta_n[t])$ achieved by analog beamforming.

**2) CSI Inaccuracy Modeling:** Given the inevitable CSI errors and to alleviate the CSI estimation burden for digital beamforming, we adopt the limited CSI model in the literature (see, e.g., [21] and references therein). Such limited CSI can be obtained by $Q$ bits of feedback from each user. Alternatively, based on the channel reciprocity, the BS could use $Q$ bits to rapidly quantize the uplink CSI (see Fig. 1). In either case, the value of $Q$ depends on the CSI learning time $\tau^D$ and efficiency of the specific CSI learning algorithm. The $Q$-bit limited CSI for each RF chain $k$ can be determined by a vector quantization codebook $\mathcal{B}_k \triangleq \{\mathbf{c}_k^1, \ldots, \mathbf{c}_k^{2^Q}\}$, where $\mathbf{c}_k^i \in \mathbb{C}^{M_{\mathrm{RF}}^{\mathrm{B}}}$, $i = 1, \ldots, 2^Q$, represents a codeword. Given an effective channel $\mathbf{H}_E^{(n)}[t]$, the codeword for its $k$-th row vector $\mathbf{h}_{E,k}^{(n)}[t]$ is chosen by picking the one that is closest to $\mathbf{h}_{E,k}^{(n)}[t]$ in the following sense [21]: $i_k^*[t] = \arg\min_{j \in \{1,\ldots,2^Q\}} \sin^2(\angle(\mathbf{h}_{E,k}^{(n)}[t], \mathbf{c}_n^j))$, where $i_k^*[t]$ denotes the index of the chosen codeword in time-slot $t$. Let $\widehat{\mathbf{H}}_E^{(n)}[t] \in \mathbb{C}^{M_{\mathrm{RF}}^{\mathrm{M}} \times M_{\mathrm{RF}}^{\mathrm{B}}}$ denote the estimated channel matrix by collecting all codewords $i_k^*[t], \forall k$. Then, based on $\widehat{\mathbf{H}}_E^{(n)}[t]$, the BS performs conjugate beamforming to construct $K$ spatial channels. However, due to the errors in $\widehat{\mathbf{H}}_E^{(n)}[t]$, inter-channel interference is not negligible under conjugate beamforming, and the amount of interference depends on the codebook size $2^Q$ and the choice of the quantization scheme.

Let $\widehat{r}_n^Q[t]$ denote the actual conjugate beamforming achievable rate under the true CSI $\mathbf{H}_E^{(n)}[t]$ while the system is treating the $Q$-bit limited CSI $\widehat{\mathbf{H}}_E^{(n)}[t]$ as if it is accurate. Also, let $\widehat{\mathbf{H}}_{E,1}^{(n)}[t]$ and $\widehat{\mathbf{H}}_{E,2}^{(n)}[t]$ represent two estimated CSI values obtained by using $Q_1$ and $Q_2$ bits, respectively. Then, we can show that the following monotonicity result of the conjugate beamforming achievable rate holds under limited CSI, which will be used in our subsequent analysis (the proof is relegated to Appendix C):

**Lemma 2** (Monotonicity of beamforming achievable rate)**.** *If $Q_1 \leq Q_2$, then there exists a CSI quantization scheme under which $\widehat{r}_n^{Q_1}[t] \leq \widehat{r}_n^{Q_2}[t]$. Further, $\widehat{r}_n^Q[t] \uparrow r_n[t]$ as $Q \to \infty$.*

**3) Algorithmic Changes to the QCS Framework:** Due to the use of $Q$-bit limited CSI in mmWave hybrid beamforming, the QCS algorithmic framework in Algorithm 1 also needs to be modified accordingly as follows:

---

**Algorithm 2:** Queue-Length-Based Congestion Control and Scheduling in mmWave Cellular Network with $Q$-Bit CSI.

---

**Initialization:** Choose parameters $K > 0$. Set $t = 0$.

**Main Loop:**

1. *MaxWeight Scheduler:* In time-slot $t \geq 1$, given the queue-length vector $\mathbf{q}[t]$ and the $Q$-bit estimated CSI $\widehat{\mathbf{H}}_E^{(n)}[t]$, $\forall n$, we let $\tilde{r}_n[t]$ be the believed conjugate beamforming achievable rate under

$\widehat{\mathbf{H}}_E^{(n)}[t]$, $\forall n$. Then, the scheduler chooses a user $n$ such that $n = \arg\max_{n' \in \{1,\dots,N\}} \{q_{n'}[t]\tilde{r}_{n'}[t]\}$. As a result, the actual achievable service rates are $s_{Q,n}[t] = \widehat{r}_n^Q[t]$ and $s_{Q,n'}[t] = 0$, $\forall n' \neq n$.

2. *Congestion Controller:* Same as in Algorithm 1.

3. *Queue-Length Updates:* Same as in Algorithm 1.

**4) Performance Analysis:** To describe our main theoretical results, we first need the following deterministic problem, where we assume that the channel state process is not random and fixed at its mean. We let $\bar{\mathcal{C}}^Q \triangleq \{r_n^Q, \forall n : r_n^Q = \mathbb{E}\{\widehat{r}_n^Q[t]\}\}$ denote the mean achievable rate region. Also, the congestion control and scheduling variables are time-invariant and denoted as $a_n$ and $s_{Q,n}$, $\forall n$, respectively. Then, the deterministic congestion control and scheduling problem can be written as:

$$\text{Maximize} \left\{ K \sum_{n=1}^N U_n(a_n) \,\middle|\, \begin{array}{l} a_n - s_{Q,n} \leq 0, \forall n, \\ s_{Q,n} \in \bar{\mathcal{C}}^Q, \forall n, \\ a_n \in [0, a^{\max}], \ \forall n. \end{array} \right\}. \tag{13}$$

Based on the convex approximation argument in Theorem 1, it is clear that Problem (13) is (approximately) convex. Thus, there exists a unique optimal solution. Associating dual variables $q_{Q,n} \geq 0$, $\forall n$ with the constraints $a_n - s_{Q,n} \leq 0$, $\forall n$, we obtain the Lagrangian as follows:

$$\Theta_K(\mathbf{q}_Q) \triangleq \max_{\mathbf{a}, \mathbf{s}_Q \in \bar{\mathcal{C}}^Q} \left\{ K \sum_{n=1}^N U_n(a_n) + \sum_{n=1}^N q_{Q,n}(s_{Q,n} - a_n) \right\}, \tag{14}$$

where the notation $\Theta_K(\cdot)$ signifies the Lagrangian's dependence on $K$ and the vector $\mathbf{q}_Q \triangleq [q_{Q,1}, \dots, q_{Q,N}]^\top \in \mathbb{R}_+^N$ contains all dual variables. Then, the Lagrangian dual problem of Problem (13) can be written as:

$$\text{Minimize} \qquad \Theta_K(\mathbf{q}_Q), \text{ subject to} \qquad \mathbf{q}_Q \in \mathbb{R}_+^N. \tag{15}$$

It is easy to show that Problem (15) satisfies the Slater condition [23]. Therefore, the optimal value of Problem (13) is equal to that of the dual problem in (15). Let $(\mathbf{a}_Q^*, \mathbf{s}_Q^*)$ and $\mathbf{q}_{Q,(K)}^*$ be the optimal primal and dual solutions to Problem (13) and Problem (15), respectively. Then, it can be shown that $\mathbf{q}_{Q,(K)}^*$ satisfies the following properties:

**Lemma 3** (Dual solution). $\mathbf{q}_{Q,(K)}^* = K\mathbf{q}_{Q,(1)}^*$, *i.e.,, $\mathbf{q}_{Q,(K)}^*$ grows linearly and the slope depends on $\mathbf{q}_{Q,(1)}^*$. Further, if $Q_1 \leq Q_2$, then the slopes satisfy $\mathbf{q}_{Q_1,(1)}^* \geq \mathbf{q}_{Q_2,(1)}^*$.*

Lemma 3 can be proved by using the Karush-Kuhn-Tucker (KKT) conditions [23] (the proof is relegated to Appendix D). Note that $K$ is only a scaling factor in the objective function in (13). Also, by contradiction, we must have $\mathbf{a}_Q^* = \mathbf{s}_Q^*$ (otherwise, $\mathbf{a}_Q^*$ can be further increased, contradicting to the fact that $\mathbf{a}_Q^*$ is optimal). Thus, we have the following result holding true:

**Lemma 4** (Primal solution). *The congestion control solution $\mathbf{a}_Q^*$ is independent of $K$ and equal to the service rate $\mathbf{s}_Q^*$.*

13

With Lemmas 3 and 4, we are now ready to present the main results in this section. Our first result says that the steady-state queue-length vector $\mathbf{q}^\infty$ lies in a bounded neighborhood of the dual solution $\mathbf{q}^*_{Q,(K)}$ of Problem (15). Further, the size of the neighborhood manifests a phase-transition phenomenon.

**Theorem 2** (Queueing delay phase transition). *Under Algorithm 2 with any given $Q$-bit CSI quantization scheme, there exists a critical value $Q^\sharp$ independent of the performance control parameter $K$ of Algorithm 2, such that:*

- *If $0 < Q < Q^\sharp$, then $\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}^*_{Q,(K)}\|\} = O(C_{(Q)}K)$, where $C_{(Q)} \geq 0$ is a constant depending on the quantization codebook, and $C_{(Q)}$ decreases as $Q$ increases;*
- *If $Q \geq Q^\sharp$, then $\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}^*_{Q,(K)}\|\} = O(\sqrt{K})$.*

**Remark 2.** Theorem 2 and Lemma 3 characterize the steady-state queue-length scalings: If $Q$ is larger than the critical value $Q^\sharp$, the steady state queue-length deviation grows sublinearly, which is much slower compared to the linear growth when $Q \leq Q^\sharp$. Also, the slope of mean queue-length $\mathbf{q}^*_{Q,(K)}$ depends on $Q$: the smaller the value of $Q$ (i.e., poorer CSI accuracy), the steeper the slope. Note that the $O(\sqrt{K})$-scaling of queue-length deviation when $Q \geq Q^\sharp$ is the same as that under the full CSI case [15]. This shows a somewhat unexpected insight that full CSI is not necessary to produce (in order sense) the original QCS queue-length scaling behavior. □

Now, let $a^\infty_{Q,n} \triangleq \mathbb{E}\{\min\{U'^{-1}_n(q^\infty_n/K), a^{\max}\}\}$, $\forall n$, be the steady-state congestion control rates under a given $Q$-bit CSI quantization scheme and let $\mathbf{a}^\infty_Q \triangleq [a^\infty_{Q,1}, \ldots, a^\infty_{Q,N}]^\top$. The next main result characterizes the phase transition of the scaling of $\mathbf{a}^\infty_Q$'s deviation from the solution $\mathbf{a}^*_Q$ of Problem (13):

**Theorem 3** (Congestion control phase transition). *Under Algorithm 2 with any $Q$-bit CSI quantization scheme, there exists a critical value $Q^\sharp$, same as in Theorem 2, such that:*

- *If $0 < Q < Q^\sharp$, then $\|\mathbf{a}^\infty_Q - \mathbf{a}^*_Q\| = O(C_{(Q)})$, where $C_{(Q)} \geq 0$ is the same constant as defined in Theorem 2;*
- *If $Q \geq Q^\sharp$, then $\|\mathbf{a}^\infty_Q - \mathbf{a}^*_Q\| = O(1/\sqrt{K})$.*

**Remark 3.** Theorem 3 also indicates a phase transition for $\mathbf{a}^\infty_Q$: When $Q < Q^\sharp$, the performance control parameter $K$ of Algorithm 2 has no effect on the deviation $\|\mathbf{a}^\infty_Q - \mathbf{a}^*_Q\|$. However, when $Q \geq Q^\sharp$, $\mathbf{a}^\infty_Q$'s deviation from $\mathbf{a}^*_Q$ grows as $O(1/\sqrt{K})$, which is the same as the full CSI case [14,15]. Another way to interpret this phase transition phenomenon is that $Q^\sharp$ represents the minimum codebook size for a CSI quantization scheme, such that it can resurrect the performance tuning capability of parameter $K$ in the QCS algorithm. Both Theorems 2 and 3 can be proved by Lyapunov stability analysis, and the details are relegated to the appendix. □

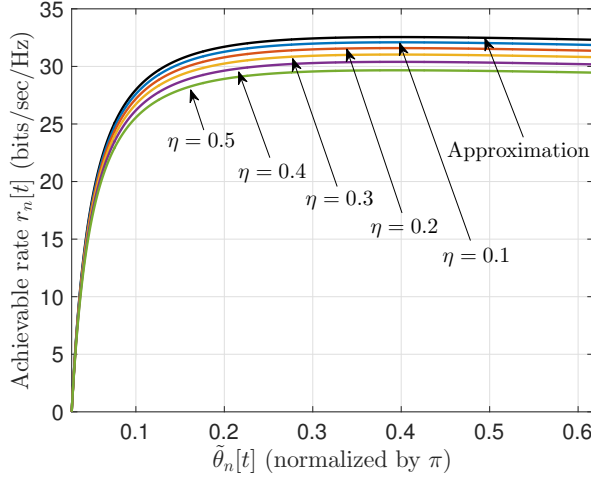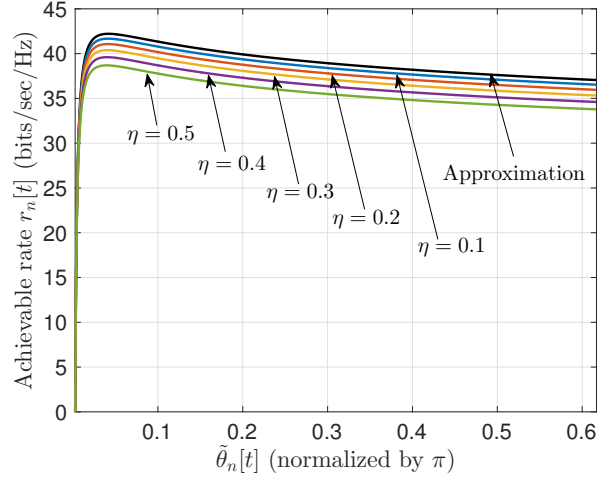Figure 5: The approximation gaps of Problem (11) under different analog slide lobe gains $\eta$ ($T_p/T = 0.01$).

Figure 6: The approximation gaps of Problem (11) under different analog slide lobe gains $\eta$ ($T_p/T = 0.001$).

## 5 Numerical Results

In this section, we conduct simulations to demonstrate the theoretical results in Sections 3 and 4. We first verify the approximation accuracy and the pseudoconvexity of Problem (11). We set SNR to 30 dB and set the $T_p/T$ ratios to 0.01 and 0.001. We vary the side lobe gain $\eta$ from 0.1 to 0.5 and the results are shown in Figs. 5 and 6. We can see that, under both $T_p/T$ ratios, the approximation gaps shrinks as $\eta$ decreases. In these examples, the gaps under $\eta = 0.1$ are almost negligible. Moreover, we note that the approximation function is indeed pseudoconcave, as predicted by Theorem 1.

Next, we examine the impacts of $Q$ on the queue-lengths and the results are shown in Fig. 7. In our simulations, we suppose that the BS and each MS have 128 and 2 RF chains, respectively. The total SNR is 40 dB. We use $\log(\cdot)$ as the utility function for each user (i.e., the proportional fairness metric [16]) and adopt random vector quantization (RVQ) as our $Q$-bit CSI quantization codebook [21]. We set the value of $Q$ to be 1, 2, 4, 8, 16, 32, 48, and 64. We also draw an accompanying dash line to show the scaling trend of each curve in Fig. 7. For small $Q$ values, we can see that the mean queue-length deviation increases faster than the square root law, roughly displaying a linear scaling with respect to $K$ as indicated in Theorem 2. For this example, the critical value of $Q$ turns out to be 8. Once $Q \geq 8$, the queue-length deviations scale as $O(\sqrt{K})$, also confirming Theorem 2.

Lastly, we study the impacts of $Q$-bit CSI on the congestion control rates and the results are illustrated in Fig. 8. For small $Q$ values, we can see that $\mathbf{a}_Q^\infty$ is only affected by $Q$ and is a constant
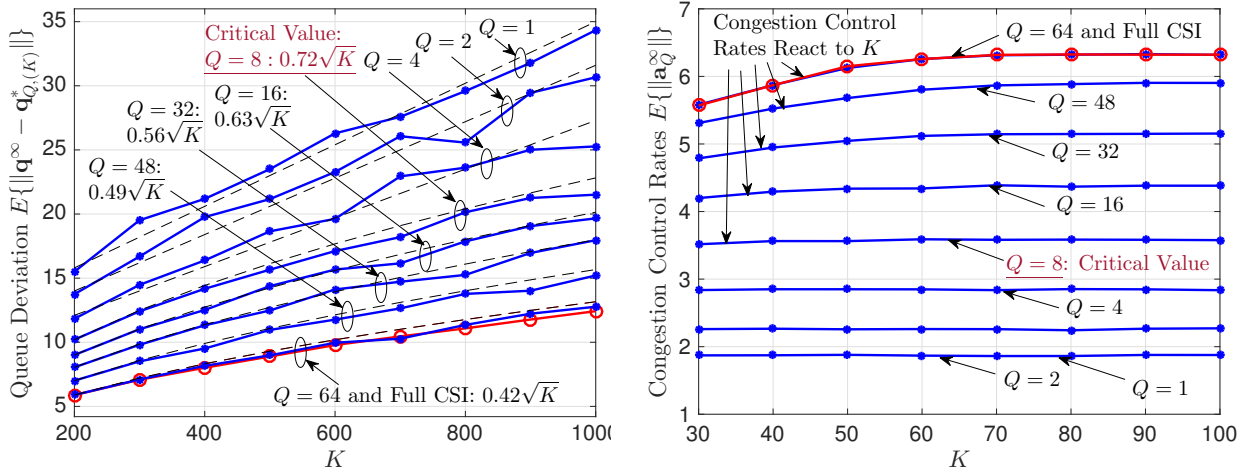
Figure 7: Average queue-length deviation with respect to $K$ for $Q = 1, 2, 4, 8, 16, 32, 48, 64$ bits.

Figure 8: The congestion control rates with respect to $K$ for $Q = 1, 2, 4, 8, 16, 32, 48, 64$ bits.

independent of $K$. Also, $\mathbf{a}_Q^\infty$'s gap to the full CSI case shrinks as $Q$ increases, which confirms Lemma 4 and Theorem 3. Again, we can observe that the critical value of $Q$ is 8: When $Q \geq 8$, $\mathbf{a}_Q^\infty$ displays an $O(1/\sqrt{K})$ diminishing gap to $\mathbf{a}_Q^*$, which agrees with Theorem 3.

## 6  Conclusion

In this paper, we studied the impacts of hybrid beamforming on the delay and network utility performance in mmWave cellular network optimization. We proposed a queue-length-based hybrid beamforming scheduling and congestion control framework for mmWave network utility maximization. We first showed that the hybrid beamforming scheduling subproblem in this framework enjoys a hidden pseudoconvexity structure, which leads to simplified analog beam training design. We then characterized two phase transition phenomena in throughput and delay with respect to CSI accuracy in digital beamforming. Collectively, these results deepen our understanding of mmWave networking performances. Hybrid beamforming in mmWave networking is an exciting and under-explored research area. Our future directions include, e.g., multi-cell mmWave networks with hybrid beamforming, the impacts of CSI inaccuracy on limited RF chains at the BS side, etc.

## A  Proof of Lemma 1

For simplicity, we let $r_n[t]$ denote the objective function of Problem (10). Substituting (2) and (3) into $r_n[[t]$ and using the defined constants, we can rewrite the objective function of Problem (10)

as:

$$r_n[t] = \left(b_0 - \frac{b_1}{\theta_B[t]\theta_M[t]}\right) \sum_{k=1}^{K} \log_2\left(1 + \underbrace{\frac{2\pi - (2\pi - \theta_B[t])\eta}{\theta_B[t]} \cdot \frac{2\pi - (2\pi - \theta_M[t])\eta}{\theta_M[t]}}_{(\Delta)} c_n^{(k)}\right). \quad (16)$$

Note that the term $(\Delta)$ can be further written as:

$$(\Delta) = \left(\frac{2\pi(1-\eta)}{\theta_B[t]} + \eta\right)\left(\frac{2\pi(1-\eta)}{\theta_M[t]} + \eta\right)$$

$$\overset{(a)}{=} \left(\frac{4\pi^2(1-\eta)^2}{\tilde{\theta}[t]} + \frac{2\pi(1-\eta)(\theta_B[t] + \theta_M[t])}{\tilde{\theta}[t]} + \eta^2\right),$$

where in $(a)$ we define $\tilde{\theta}[t] \triangleq \theta_B[t]\theta_M[t]$. Now, we claim that

$$\frac{4\pi^2(1-\eta)^2}{\tilde{\theta}[t]} \gg \frac{2\pi(1-\eta)(\theta_B[t] + \theta_M[t])}{\tilde{\theta}[t]} \quad (17)$$

is true if $\eta \ll \frac{1}{3}$. To see this, we first note that $\eta \ll \frac{1}{3}$ implies $4\pi \ll \frac{2\pi(1-\eta)}{\eta}$. Also, since $\theta_B[t], \theta_M[t] \in (0, 2\pi]$, we have

$$\theta_B[t] + \theta_M[t] \le 4\pi \ll \frac{2\pi(1-\eta)}{\eta},$$

which implies that (17) is true. Hence, it follows from (17) and $\eta \ll \frac{1}{3} < 1$ that $(\Delta) \approx (\frac{4\pi^2}{\tilde{\theta}[t]} + \eta^2)$, which further implies

$$r_n[t] = (16) \approx \left(b_0 - \frac{b_1}{\tilde{\theta}[t]}\right) \sum_{k=1}^{K} \log_2\left(1 + \frac{4\pi^2 c_n^{(k)}}{\tilde{\theta}[t]}\right),$$

i.e., the objective function in (11). This completes the proof.

## B  Proof of Theorem 1

As mentioned earlier, verifying the pseudoconvexity of Problem (11) means verifying the pseudo-concavity of the objective function. Toward this end, we let $f(\tilde{\theta}[t])$ denote the *negative* objective function and our goal is to show that $f(\tilde{\theta}[t])$ is pseudoconvex, which means that for any $\tilde{\theta}_1[t]$ and $\tilde{\theta}_2[t]$ in the feasible interval, if $f'(\tilde{\theta}_1[t])(\tilde{\theta}_2[t] - \tilde{\theta}_1[t]) \ge 0$, we must also have $f'(\tilde{\theta}_2[t])(\tilde{\theta}_2[t] - \tilde{\theta}_1[t]) \ge 0$.

First, let us consider the case where $\tilde{\theta}_2[t] \ge \tilde{\theta}_1[t]$. Then, showing $f'(\tilde{\theta}_2[t])(\tilde{\theta}_2[t] - \tilde{\theta}_1[t]) \ge 0$ is equivalent to showing $f'(\tilde{\theta}_2[t]) \ge 0$. Note that, in this case, the condition $f'(\tilde{\theta}_1[t])(\tilde{\theta}_2[t] - \tilde{\theta}_1[t]) \ge 0$ simply means $f'(\tilde{\theta}_1[t]) \ge 0$, i.e.,

$$f'(\tilde{\theta}_1[t]) = \sum_{k=1}^{K} \frac{1}{\tilde{\theta}_1^2} \left[ \underbrace{\frac{4\pi^2 c_n^{(k)}}{\ln(2)} \cdot \frac{b_0\tilde{\theta}_1[t] - b_1}{\tilde{\theta}_1[t] + 4\pi^2 c_n^{(k)}}}_{(P1)} \underbrace{-b_1 \log_2\left(1 + \frac{4\pi^2 c_n^{(k)}}{\tilde{\theta}_1[t]}\right)}_{(P2)} \right] \ge 0. \quad (18)$$

17

It is obvious that the term $(P2)$ is an increasing function of $\tilde{\theta}[t]$. Now, consider the fractional term $\frac{b_0\tilde{\theta}_1[t]-b_1}{\tilde{\theta}_1[t]+4\pi^2 c_n^{(k)}}$ in $(P1)$, which is negative-valued according to the definitions of $b_0$, $b_1$, and the feasible interval. Also, from the definition of $b_0$, we have $b_0 < 1$, implying that the absolute value of the nominator is increasing at a slower rate than that of the denominator. This means that $(P1)$ is also an increasing function of $\tilde{\theta}[t]$. Hence, $f'(\tilde{\theta}[t])$ is increasing since both $(P1)$ and $(P2)$ are increasing. As a result, $f'(\tilde{\theta}_1[t]) \geq 0$ and $\tilde{\theta}_2[t] \geq \tilde{\theta}_1[t]$ imply $f'(\tilde{\theta}_2[t]) \geq 0$ and thus the case of $\tilde{\theta}_2[t] \geq \tilde{\theta}_1[t]$ is proved. The other case where $\tilde{\theta}_2[t] \leq \tilde{\theta}_1[t]$ can also be proved by similar arguments and we omit the details in here for brevity.

To show that the optimal solution is unique and achieved in the interior of the feasible interval, it suffices to show that $\frac{\partial r_n[t]}{\partial \tilde{\theta}[t]}$'s values at two end points of the interval have opposite signs. Then, from the decreasing derivative property of $r_n[t]$ $(r_n[t] = -f(\tilde{\theta}[t]))$, $\frac{\partial r_n[t]}{\partial \tilde{\theta}[t]}$ must have exactly one zero-crossing point in the interior of the feasible interval. Also, the pseudoconcavity of $r_n[t]$ means that the zero-crossing point is the global maximum. First, if $\tilde{\theta}[t] = \frac{b_1}{b_0}$, we have $(P1) = 0$. Hence, $\frac{\partial r_n[t]}{\partial \tilde{\theta}[t]} > 0$ since $-(P2) > 0$ (because $b_1$, $\tilde{\theta}^2[t]$, and the $\log(\cdot)$ rate expressions are positive). On the other hand, when $\tilde{\theta}[t] \uparrow \psi_n^{\mathrm{B}}\psi_n^{\mathrm{M}}$, it follows from $T_p \ll T$ that

$$r_n[t] = \left[1 - \left(N\beta + \frac{\sum_{n'=1}^{N}\psi_{n'}^{\mathrm{B}}\psi_{n'}^{\mathrm{M}}}{\psi_n^{\mathrm{B}}\psi_n^{\mathrm{M}}}\right)\frac{T_p}{T}\right] \times \sum_{k=1}^{K}\log_2\left(1 + \frac{4\pi^2 c_n^{(k)}}{\tilde{\theta}[t]}\right) \approx \sum_{k=1}^{K}\log_2\left(1 + \frac{4\pi^2 c_n^{(k)}}{\tilde{\theta}[t]}\right),$$

which is decreasing in $\tilde{\theta}[t]$ and must have a negative derivative at $\tilde{\theta}[t] = \psi_n^{\mathrm{B}}\psi_n^{\mathrm{M}}$. This completes the proof.

## C   Proof of Lemma 2

We first show the second part of Lemma 2. Since the BS performs conjugate beamforming with equal power allocation by treating $\widehat{\mathbf{H}}_E^{(n)}[t]$ as if it is the accurate CSI, the received signal can be written as $y_{n,k}[t] = u_{n,k}[t](P_{\max}/K)\mathbf{h}_{E,k}^{(n)\dagger}[t]\widehat{\mathbf{w}}_{n,k}[t] + \sum_{k'=1,\neq k}^{K}u_{n,k'}[t](P_{\max}/K)\mathbf{h}_{E,k'}^{(n)\dagger}[t]\widehat{\mathbf{w}}_{j,k}[t] + v_n[t]$, where $\widehat{\mathbf{w}}_{n,k}[t] = \widehat{\mathbf{h}}_{E,k}^{(n)}[t]$, $1 \leq k \leq K$, i.e., the $k$-th row of $\widehat{\mathbf{H}}_E^{(n)}[t]$. Hence, the conjugate beamforming rates $s_{Q,n}[t]$ achieved under $\mathbf{H}_E^{(n)}[t]$ based on the belief that the CSI is $Q$-bit CSI $\widehat{\mathbf{H}}_E^{(n)}[t]$ can be computed as:

$$\widehat{r}_n^Q[t] = \sum_{k=1}^{K}\log_2\left(1 + \frac{(P_{\max}/K)\big|\mathbf{h}_{E,k}^{(n)\dagger}[t]\widehat{\mathbf{h}}_{E,k}^{(n)}[t]\big|^2}{N_0 + \sum_{k'=1,\neq k}^{K}(P_{\max}/K)\big|\mathbf{h}_{E,k'}^{(n)\dagger}[t]\widehat{\mathbf{h}}_{E,k'}^{(n)}[t]\big|^2}\right)$$

$$< \sum_{k=1}^{K}\log_2\left(1 + \frac{P_{\max}}{KN_0}\big\|\mathbf{h}_{E,k}^{(n)}[t]\big\|^2\right) = r_n[t], \ \forall n, \tag{19}$$

where the inequality in (19) holds because $\big|\mathbf{h}_{E,k}^{(n)\dagger}[t]\widehat{\mathbf{h}}_{E,k}^{(n)}[t]\big|^2 \leq \|\mathbf{h}_{E,k}^{(n)}[t]\|^2$ and $|\mathbf{h}_{E,k'}^{(n)\dagger}[t]\widehat{\mathbf{h}}_{E,k'}^{(n)}[t]|^2 \geq 0$. Thus, for every rate point $\widehat{\mathbf{r}}^Q[t] = [\widehat{r}_n^Q[t], \ldots, \widehat{r}_n^Q[t]]^T$, equal power allocation achieves a rate point

$\mathbf{r}[t] = [r_1[t], \ldots, r_N[t]]^\top$ that dominates $\widehat{\mathbf{r}}^Q[t]$ in every coordinate. Also, as $Q \to \infty$, $\widehat{\mathbf{H}}_E^{(n)}[t] \to \mathbf{H}_E^{(n)}[t]$. It thus follows from (19) that $\widehat{\mathbf{r}}_n^Q[t] \uparrow \mathbf{r}[t]$.

Next, we prove that the first part of Lemma 2 is true. Let $\mathcal{B}_n^1$ and $\mathcal{B}_n^2$ denote the vector quantization codebooks corresponding to $Q_1$ and $Q_2$ bits, respectively. Since $Q_1 \leq Q_2$, it follows that the codebook sizes $|\mathcal{B}_n^1| \leq |\mathcal{B}_n^2|$. Hence, given codebook $\mathcal{B}_n^1$, one can construct $\mathcal{B}_n^2$ by simply retaining all codewords in $\mathcal{B}_n^1$ and adding new code words that are not in $\mathcal{B}_n^1$, which implies $\mathcal{B}_n^1 \subset \mathcal{B}_n^2$. As a result, for any given CSI $\mathbf{h}_{E,k}^{(n)}[t]$, one can always find a codeword in $\mathcal{B}_n^2$ whose distance to $\mathbf{h}_{E,k}^{(n)}[t]$ is not larger than that from $\mathcal{B}_n^1$. Hence, the SINR term in (19) becomes larger under $\mathcal{B}_n^2$, implying $\widehat{r}_n^{Q_1}[t] \leq \widehat{r}_n^{Q_2}[t]$. This completes the proof.

## D  Proof of Lemma 3

Dividing $K$ on both sides of (14), we have

$$\frac{1}{K}\Theta_K(\mathbf{q}_Q) = \max_{\mathbf{a}, \mathbf{s}_Q \in \bar{\mathcal{C}}^Q} \left\{ \sum_{n=1}^{N} U_n(a_n) + \sum_{n=1}^{N} \widehat{q}_{Q,n}(s_{Q,n} - a_n) \right\},$$

where $\widehat{q}_{Q,n} = q_{Q,n}/K$. Note that the right hand side is precisely $\Theta_1(\mathbf{q}_Q)$, for which the maximizer is $\widehat{\mathbf{q}} = \mathbf{q}_{Q,(1)}^*$. Hence, we have $\Theta_K(\mathbf{q})$ is maximized at $K\mathbf{q}_{Q,(1)}^*$. This proves the first part of Lemma 3.

To show the second part of Lemma 3, we first have from the KKT complementary slackness condition and the monotonicity of $U_n(\cdot)$ that, at optimality, $a_n^* = s_{Q,n}^*$, $\forall n$. We let $a_n^*(Q_1)$ and $a_n^*(Q_2)$ denote the optimal congestion control rates under $Q_1$ and $Q_2$, respectively. If $Q_1 \leq Q_2$, we have from Lemma 2 that $s_{Q_1,n}^* \leq s_{Q_2,n}^*$, which further implies $a_n^*(Q_1) \leq a_n^*(Q_2)$. On the other hand, from the KKT stationarity condition, we have $U_n'(a_n^*(Q)) - q_{(Q),n}^* = 0$. Since $a_n^*(Q_1) \leq a_n^*(Q_2)$, it follows from the concavity of $U_n(\cdot)$ that $q_{Q_1,n}^* \geq q_{Q_2,n}^*$. This completes the proof.

## E  Proof of Theorem 2

To prove Theorem 2, we first show the existence of steady-state by proving a positive Harris-recurrence result of the queue-length process. This result implies the existence of steady-state, which lays the foundation for proving Theorems 2. We let $\mathbb{1}_{\mathcal{A}}(\mathbf{x})$ denote the indicator function, which takes value 1 if $\mathbf{x} \in \mathcal{A}$ and 0 otherwise. We state the queue-length positive Harris-recurrence result as follows:

**Proposition 1** (Queue-Length Positive Recurrence). *Consider a Lyapunov function* $V(\mathbf{q}[t]) \triangleq \frac{1}{2K}\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^2$ *for a given $K$. For the scheduler (6) and congestion controller (7)–(8), there exist constants $\delta, \eta > 0$, both independent of $K$, such that the queue-length process $\{\mathbf{q}[t]\}_{t=0}^{\infty}$ satisfies*

*the following conditional mean drift condition:*

$$\mathbb{E}\{\Delta V(\mathbf{q}[t]|\mathbf{q}[t])\} \triangleq \mathbb{E}\{V(\mathbf{q}[t+1]) - V(\mathbf{q}[t])|\mathbf{q}[t]\}$$

$$\leq -\frac{\delta}{\Phi K}\|\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)}\|\mathbb{1}_{\mathcal{Q}^c}(\mathbf{q}[t]) + \eta\mathbb{1}_{\mathcal{Q}}(\mathbf{q}[t]), \qquad (20)$$

*where* $\mathcal{Q} \triangleq \{\mathbf{q} \in \mathbb{Z}_+^N |\|\mathbf{q} - \mathbf{q}^*_{Q,(K)}\| \leq \gamma K\}$ *for some constant* $\gamma > 0$ *and* $\mathcal{B}^c$ *denotes the complement of* $\mathcal{Q}$ *in* $\mathbb{Z}_+^N$.

*Proof.* Consider the quadratic Lyapunov function defined in Proposition 1: $V(\mathbf{q}[t]) = \frac{1}{2K}\|\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)}\|^2$, where $\mathbf{q}[t]$ represents the queue-length vector in time-slot $t$ under parameters $K$ and $Q$; and $\mathbf{q}^*_{Q,(K)}$ denotes the optimal dual solution for the static version of Problem JCS under parameter $K$. Then, the one-slot mean Lyapunov drift of $V(\mathbf{q}[t])$, which can computed as:

$$\mathbb{E}\{V(\mathbf{q}[t+1]) - V(\mathbf{q}[t])|\mathbf{q}[t]\}$$

$$= \mathbb{E}\left\{\frac{1}{2K}\|\mathbf{q}[t+1] - \mathbf{q}^*_{(K)}\|^2 - \frac{1}{2K}\|\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)}\|^2 \Big|\mathbf{q}[t]\right\}$$

$$= \frac{1}{2K}\mathbb{E}\left\{(\mathbf{q}[t+1] - \mathbf{q}[t])^\top(\mathbf{q}[t+1] + \mathbf{q}[t] - 2\mathbf{q}^*_{Q,(K)})\Big|\mathbf{q}[t]\right\}$$

$$\overset{(a)}{\leq} \frac{1}{2K}\mathbb{E}\left\{(-\mathbf{s}_Q[t] + \mathbf{a}[t])^\top(2\mathbf{q}[t] - 2\mathbf{q}^*_{Q,(K)} - \mathbf{s}_Q[t] + \mathbf{a}[t])\Big|\mathbf{q}[t]\right\}$$

$$= \frac{1}{2K}\mathbb{E}\left\{\|-\mathbf{s}_Q[t] + \mathbf{a}[t]\|^2 + 2(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top(-\mathbf{s}^Q[t] + \mathbf{a}[t])\Big|\mathbf{q}[t]\right\}$$

$$= \frac{1}{K}(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top(-\mathbf{s}_Q[t] + \mathbf{a}[t]) + \frac{1}{2K}\mathbb{E}\left\{\|-\mathbf{s}_Q[t] + \mathbf{a}[t]\|^2\right\},$$

where $(a)$ follows from the non-expansive property of the $\max\{0, \cdot\}$ operation. Note that, from the definition of Algorithm 1, we have $\mathbb{E}\{\|\mathbf{a}[t]\|^2|\mathbf{q}[t]\} < A_2^{\max}N$. Also, since $s_{Q,n}[t]$ falls in a bounded instantaneous capacity region $\mathcal{C}_{\widehat{\mathbf{H}}[t]}$, $\forall n$, we must have $s_{Q,n}[t] \leq s^{\max}$ for some $s^{\max} > 0$. Hence, by defining $D_0 \triangleq \frac{N}{2}(A_2^{\max} + (s^{\max})^2)$, we have

$$\mathbb{E}\{\Delta V(\mathbf{q}[t])|\mathbf{q}[t]\} \leq \frac{1}{K}(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top\mathbb{E}\{\mathbf{a}[t] - \mathbf{s}_Q[t]\} + \frac{D_0}{K}$$

$$\overset{(a)}{=} \frac{1}{K}(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top(\mathbb{E}\{\mathbf{a}[t]|\mathbf{q}[t]\} - \mathbf{s}^*_Q) +$$

$$\frac{1}{K}\mathbb{E}\{(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top(\mathbf{s}^*_Q - \mathbf{s}_Q[t])|\mathbf{q}[t]\} + \frac{D_0}{K},$$

$$\overset{(b)}{\leq} \frac{1}{K}(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top(\mathbb{E}\{\mathbf{a}[t]|\mathbf{q}[t]\} - \mathbf{s}^*_Q) +$$

$$\frac{1}{K}\|\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top\| \times \mathbb{E}\{\|\mathbf{s}^*_Q - \mathbf{s}_Q[t]\||\mathbf{q}[t]\} + \frac{D_0}{K}, \qquad (21)$$

where $\mathbf{s}^*_Q$ is such that $(\mathbf{s}^*_Q, \mathbf{q}^*_{Q,(K)})$ is a pair of optimal primal and dual solutions to Problem (15) under parameter $K$. In (21), $(a)$ follows from adding and subtracting $\mathbf{s}^*_Q$ as well as the fact that $\mathbf{a}[t]$ is independent of the channel state and determined solely by $\mathbf{q}[t]$; and $(b)$ follows from Cauchy-Schwarz inequality.

Note from Lemma 4 that $\mathbf{s}_Q^*$ is independent of $K$ and $s_{Q,n}[t] \in \mathcal{C}_{\widehat{\mathbf{H}}[t]}$ is upper-bounded. Thus, we have

$$\mathbb{E}\{\|\mathbf{s}_Q^* - \mathbf{s}_Q[t]\|\|\mathbf{q}[t]\} \leq C_{(Q)} \triangleq \max_{\mathbf{q}:\|\mathbf{q}\|=1} \mathbb{E}\{\|\mathbf{s}_Q^* - \mathbf{s}_Q\|\mathbf{q}\}, \tag{22}$$

where $C_{(Q)}$ signifies that its value depends on $Q$. Hence, we can further upper bound (21) as:

$$\mathbb{E}\{\Delta V(\mathbf{q}[t])|\mathbf{q}[t]\} \leq \frac{1}{K}(\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*)^\top (\mathbb{E}\{\mathbf{a}[t]|\mathbf{q}[t]\} - \mathbf{s}_Q^*) + \frac{1}{K}\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^\top\|C_{(Q)} + \frac{D_0}{K}, \tag{23}$$

Now, let us consider the first term on the right hand side in (23), i.e., $\frac{1}{K}(\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*)^\top (\mathbb{E}\{\mathbf{a}[t]|\mathbf{q}[t]\} - \mathbf{s}^*)$. Since $U_n(\cdot)$ is concave and increasing, $\forall n$, we have

$$\left(q_n[t] - q_{Q,(K),n}^*\right)^\top \left[U_n'^{-1}\left(\frac{q_n[t]}{K}\right) - U_n'^{-1}\left(\frac{q_{Q,(K),n}^*}{K}\right)\right] \leq 0.$$

Thus, by Cauchy-Schwatz inequality, we have:

$$(\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*)^\top (\mathbb{E}\{\mathbf{a}[t]|\mathbf{q}[t]\} - \mathbf{s}_Q^*) = \sum_{n=1}^{N} \left(q_n[t] - q_{Q,(K),n}^*\right)^\top$$

$$\times \left[U_n'^{-1}\left(\frac{q_n[t]}{K}\right) - U_n'^{-1}\left(\frac{q_{Q,(K),n}^*}{K}\right)\right] \leq -\sum_{n=1}^{N} |q_n[t] -$$

$$q_{Q,(K),n}^*|\left|U_n'^{-1}\left(\frac{q_n[t]}{K}\right) - U_n'^{-1}\left(\frac{q_{Q,(K),n}^*}{K}\right)\right|. \tag{24}$$

By the strong convexity of $-U_n(\cdot)$ and the Lipschitz continuity of $U_n'(\cdot)$, we have

$$\left|U_n'\left(a_{n,1}\right) - U_n'\left(a_{n,2}\right)\right| \leq \Phi\left|a_{n,1} - a_{n,2}\right|.$$

Therefore, by the inverse function lemma, we have

$$\frac{1}{\Phi}\left|\frac{q_n[t]}{K} - \frac{q_{Q,(K),n}^*}{K}\right| \leq \left|U_n'^{-1}\left(\frac{q_n[t]}{K}\right) - U_n'^{-1}\left(\frac{q_{Q,(K),n}^*}{K}\right)\right|.$$

Hence, we can further upper-bound (24) as:

$$(\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*)^\top (\mathbb{E}\{\mathbf{a}[t]|\mathbf{q}[t]\} - \mathbf{s}_Q^*) \leq -\frac{1}{\Phi K}\sum_{n=1}^{N} \left(q_n[t] - q_{Q,(K),n}^*\right)^2 = -\frac{1}{\Phi K}\left\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\right\|^2. \tag{25}$$

Substituting (25) into (23), we have

$$\mathbb{E}\{\Delta V(\mathbf{q}[t])|\mathbf{q}[t]\} \quad \leq \quad -\frac{1}{\Phi K^2}\left\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\right\|^2 \quad + \quad \frac{1}{K}\|\mathbf{q}[t] \quad - \quad \mathbf{q}_{Q,(K)}^*\|^\top\|C_{(Q)} \quad + \quad \frac{D_0}{K}. \tag{26}$$

Now, suppose that $\left\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\right\| \geq \gamma_1 K$, where $\beta_1$ will be specified shortly. Note also that $K \geq 1$, we have

$$\frac{1}{\left\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\right\|} \leq \frac{1}{\gamma_1 K} \leq \frac{1}{\gamma_1}.$$

21

It then follows that (26) can be further upper bounded as:

$$
\begin{aligned}
\mathbb{E}\{\Delta V(\mathbf{q}[t])|\mathbf{q}[t]\} &= -\frac{1}{\Phi K}\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\big\|\cdot\frac{\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\big\|}{K}\\
&\quad+\frac{1}{K}\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)})^\top\big\|D_1+\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\big\|\frac{D_0}{\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\big\|K}\\
&\leq-\frac{1}{\Phi K}\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\big\|\Big(\gamma_1-C_{(Q)}\Phi-\frac{D_0\Phi}{\gamma_1}\Big).
\end{aligned}
\tag{27}
$$

By choosing $\gamma_1$ such that $\gamma_1-D_1\Phi-\frac{D_0\Phi}{\gamma_1}>0$, we have

$$
\mathbb{E}\{\Delta V(\mathbf{q}[t])|\mathbf{q}[t]\}\leq-\frac{\hat{\delta}_1}{\Phi K}\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\big\|
\tag{28}
$$

where $\hat{\delta}_1=\gamma_1-C_{(Q)}\Phi-\frac{D_0\Phi}{\gamma_1}$. Solving $\beta_1-C_{(Q)}\Phi-\frac{D_0\Phi}{\beta_1}=0$ and plugging in the obtained $\gamma_1$ to define a ball $\mathcal{B}_1\triangleq\{\mathbf{q}:\big\|\mathbf{q}-\mathbf{q}^*_{Q,(K)}\big\|\leq\frac{K}{2}[(C_{(Q)}\Phi)+\sqrt{(C_{(Q)}\Phi)^2+4D_0\Phi}]\}$, we have

$$
\mathbb{E}\{\Delta V(\mathbf{q}[t])|\mathbf{q}[t]\}\leq-\frac{\delta_1}{K}\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\big\|,\text{ if }\mathbf{q}[t]\in\mathcal{B}_1^c,
\tag{29}
$$

where $\delta_1\triangleq\frac{\hat{\delta}_1}{\Phi}$. On the other hand, when $\mathbf{q}[t]\in\mathcal{B}_1$, it is clearly true that $\mathbb{E}\{\Delta V(\mathbf{q}[t])|\mathbf{q}[t]\}\leq\eta_1$ for some $\eta_1>0$. Combining these facts yields the following:

$$
\mathbb{E}\{\Delta V(\mathbf{q}[t])|\mathbf{q}[t]=\mathbf{q}\}\leq-\frac{\delta_1}{K}\big\|\mathbf{q}-\mathbf{q}^*_{Q,(K)}\big\|\mathbb{1}_{\mathcal{B}_1^c}(\mathbf{q})+\eta_1\mathbb{1}_{\mathcal{B}_1}(\mathbf{q}).
$$

This completes the proof of Proposition 1. $\qquad\square$

The inequality in (20) suggests that the conditional mean drift is negative when the deviation of the queue-length vector $\mathbf{q}[t]$ away from $\mathbf{q}^*_{Q,(K)}$ is sufficiently large. Since (20) is just the Foster-Lyapunove criterion [24, Proposition I.5.3], $\{\mathbf{q}[t]\}_{t=0}^\infty$ is positive recurrent, we have that a steady-state distribution of queue-lengths exists. Thus, we let $\mathbf{q}^\infty$ denote the queue-length vector in steady-state. With Proposition 1, we are now in a position to prove Theorem 2.

Next, to prove Theorem 2, we use an $\alpha$-parameterized quadratic Lyapunov function: $V_\alpha(\mathbf{q}[t])=\frac{1}{2K^\alpha}\big\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\big\|^2$, where the parameter $\alpha\in\{0,1\}$ and its value will be specified later. Following similar steps in the proof of Proposition 1, we can bound the conditional mean Lyapunov drift as follows:

$$\mathbb{E}\{V_\alpha(\mathbf{q}[t+1]) - V_\alpha(\mathbf{q}[t])|\mathbf{q}[t]\}$$

$$\stackrel{(a)}{\leq} \frac{1}{K^\alpha}(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top (\mathbb{E}\{\mathbf{a}[t]|\mathbf{q}[t]\} - \mathbf{s}^*_Q) +$$

$$\frac{1}{K^\alpha}\mathbb{E}\{(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top(\mathbf{s}^*_Q - \mathbf{s}_Q[t])|\mathbf{q}[t]\} + \frac{D_0}{K^\alpha},$$

$$\stackrel{(b)}{\leq} \frac{1}{K^\alpha}\Big[-\frac{1}{\Phi K}\|\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)}\|^2 + D_0\Big] +$$

$$\frac{1}{K^\alpha}\mathbb{E}\Big\{(\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)})^\top(\mathbf{s}^*_Q - \mathbf{s}_Q[t])\big|\mathbf{q}[t]\Big\}$$

$$\stackrel{(c)}{\leq} \frac{1}{K^\alpha}\Big[-\frac{1}{\Phi K}\|\mathbf{q}[t] - \mathbf{q}^*_{Q,(K)}\|^2 + D_0\Big] +$$

$$\frac{1}{K^\alpha}\mathbb{E}\Big\{(\mathbf{q}[t])^\top(\mathbf{s}^* - \mathbf{s}_Q[t])\big|\mathbf{q}[t]\Big\}, \tag{30}$$

where $D_0 \triangleq \frac{N}{2}(A_2^{\max} + (s^{\max})^2)$ and $\mathbf{s}^* \triangleq \lim_{Q\to\infty}\mathbf{s}^*_Q$. In (30), (a) follows from adding and subtracting $\mathbf{s}^*_Q$; (b) follows from (25); and (c) follows from $\mathbf{s}^*_Q \leq \mathbf{s}^*$ (by Lemma 2) and the scheduler design, which implies $(\mathbf{q}^*_{Q,(K)})^\top \mathbf{s}_Q[t] \leq (\mathbf{q}^*_{Q,(K)})^\top \mathbf{s}^*_Q$. Next, consider the $T$-step conditional mean Lyapunov drift. For any $\mathbf{q}[0] \geq \mathbf{0}$, we have that

$$\mathbb{E}\{V_\alpha(\mathbf{q}[T])|\mathbf{q}[0]\} - V_\alpha(\mathbf{q}[0]) = \sum_{t=0}^{T-1}\mathbb{E}\{V(\mathbf{q}[t+1]) - V(\mathbf{q}[t])|\mathbf{q}[0]\}$$

$$\stackrel{(a)}{=} \sum_{t=0}^{T-1}\sum_{\mathbf{q}\in\mathcal{Z}_+^N}\Big[\Pr(\mathbf{q}[t]=\mathbf{q}|\mathbf{q}[0])\mathbb{E}\{V_\alpha(\mathbf{q}[t+1]) - V_\alpha(\mathbf{q}[t])|\mathbf{q}[t]=\mathbf{q}\}\Big]$$

$$\stackrel{(b)}{\leq} \sum_{t=0}^{T-1}\sum_{\mathbf{q}\in\mathcal{Z}_+^N}\Pr(\mathbf{q}[t]=\mathbf{q}|\mathbf{q}[0])\Big\{\frac{1}{K^\alpha}\Big[\frac{-1}{\Phi K}\|\mathbf{q}[t]-\mathbf{q}^*_{Q,(K)}\|^2 + D_0\Big]\Big\}$$

$$+\sum_{t=0}^{T-1}\sum_{\mathbf{q}\in\mathcal{Z}_+^N}\Pr(\mathbf{q}[t]=\mathbf{q}|\mathbf{q}[0])\Big\{\frac{1}{K^\alpha}\mathbb{E}\Big\{\mathbf{q}^\top(\mathbf{s}^*-\mathbf{s}_Q[t])\Big\}\Big\}, \tag{31}$$

where (a) follows from the fact that $\mathbf{q}[t]$ is a discrete state Markov chain in $\mathbb{Z}_+^N$ and (b) follows from (30). Note that for any $\mathbf{q}[t] \in \mathbb{Z}_+^N$, $\lim_{T\to\infty}\frac{1}{T}\sum_{t=0}^{T-1}\Pr(\mathbf{q}[t]=\mathbf{q}|\mathbf{q}[0]) = \pi_{\mathbf{q}}^\infty$, where $\pi_{\mathbf{q}}^\infty$ denotes the stationary distribution of the Markov chain $\mathbf{q}[t]$. Moving $V(\mathbf{q}[0])$ to the right hand side, dividing both sides by $T$, and letting $T \to \infty$ yields:

$$0 \leq J + \sum_{\mathbf{q}\in\mathbb{Z}_+^N}\pi_{\mathbf{q}}^\infty(\mathbf{q})^\top(\mathbf{s}^* - \mathbf{s}_B^\infty) = J + \mathbb{E}\{(\mathbf{q}^\infty)^\top(\mathbf{s}^* - \mathbf{s}_B^\infty\}, \tag{32}$$

where $J \triangleq \lim_{T\to\infty}\frac{1}{T}\sum_{t=0}^{T-1}\sum_{\mathbf{q}\in\mathcal{Z}_+^N}\Pr(\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0])\{\frac{1}{K^\alpha}[\frac{-1}{\Phi K}\|\mathbf{q}[t] - \mathbf{q}^*_{(K)}\|^2 + D_0]\}$, and $\mathbf{s}_Q^\infty \triangleq \arg\max_{\mathbf{x}\in\mathcal{C}_{\mathbf{H}[\infty]|\widehat{\mathbf{H}}[\infty]}}(\mathbf{q}^\infty)^\top\mathbf{x}$ represents the steady-state service rates with $Q$-bit CSI.

Next, consider the term $\mathbb{E}\{(\mathbf{q}^\infty)^\top(\mathbf{s}^* - \mathbf{s}_Q^\infty)\}$ in (32). For any given realization of $\mathbf{q}^\infty$ in the steady-state, from the design of the MaxWeight scheduler in (6), we have that

$$(\mathbf{q}^\infty)^\top\mathbf{s}^* \leq \max_{\mathbf{x}\in\mathcal{C}_{\mathbf{H}[\infty]}}(\mathbf{q}^\infty)^\top\mathbf{x} = (\mathbf{q}^\infty)^\top\mathbf{s}^\infty. \tag{33}$$

where $\mathbf{s}^\infty \triangleq \lim_{Q \to \infty} \mathbf{s}_Q^\infty$ and $\mathbf{H}[\infty]$ represent the full CSI in the steady state. Hence, for any realization of $\mathbf{q}^\infty$ such that $\mathbf{q}^\infty \neq \rho \mathbf{s}^*$ for some $\rho \in \mathbb{R}$, if $Q$ is sufficiently large, we must have $(\mathbf{q}^\infty)^\top \mathbf{s}^* - (\mathbf{q}^\infty)^\top \mathbf{s}_Q^\infty \leq 0$. Hence, there exists a critical value $Q^\sharp$ such that for all $Q > Q^\sharp$, the average value of $(\mathbf{q}^\infty)^\top \mathbf{s}^* - (\mathbf{q}^\infty)^\top \mathbf{s}_Q^\infty$ can be made non-positive, i.e., $\mathbb{E}\{(\mathbf{q}^\infty)^\top (\mathbf{s}^* - \mathbf{s}_Q^\infty)\} \leq 0$. Hence, we consider two cases based on the positivity of $\mathbb{E}\{(\mathbf{q}^\infty)^\top (\mathbf{s}^* - \mathbf{s}_Q^\infty)\}$ as follows:

_Case I):_ $Q \geq Q^\sharp$ such that $\mathbb{E}\{(\mathbf{q}^\infty)^\top (\mathbf{s}^* - \mathbf{s}_Q^\infty)\} \leq 0$: In this case, it follows from (32) that

$$0 \leq \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{\mathbf{q} \in \mathcal{Z}_+^N} \Pr(\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0]) \left\{ \frac{1}{K^\alpha} \left[ -\frac{1}{\Phi K} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^2 + D_0 \right] \right\}. \quad (34)$$

We now consider the term in the second line in (34) by setting $\alpha = 0$. Similar to the proof of Proposition 1, suppose that $\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\| \geq \gamma \sqrt{K}$, where $\gamma$ will be specified shortly. This implies that $\frac{1}{\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|} \leq \frac{1}{\gamma}$. Then, the second line in (34) can be upper bounded as:

$$-\frac{1}{\Phi K} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^2 + D_0 = -\frac{1}{\Phi \sqrt{K}} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\| \times$$
$$\left( \frac{\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|}{\sqrt{K}} + \frac{D_0 \Phi \sqrt{K}}{\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|} \right)$$
$$\leq -\frac{1}{\Phi \sqrt{K}} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\| \left( \gamma - \frac{D_0 \Phi}{\gamma} \right). \quad (35)$$

Hence, by choosing $\gamma > \sqrt{D_0 \Phi}$, we have

$$-\frac{1}{\Phi K} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^2 + D_0 \leq -\frac{\hat{\delta}}{\Phi \sqrt{K}} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|, \quad (36)$$

where $\hat{\delta} = \gamma - \frac{D_0 \Phi}{\gamma} > 0$. Plugging in $\gamma > \sqrt{D_0 \Phi}$ to define a ball $\mathcal{B} \triangleq \{\mathbf{q} : \|\mathbf{q} - \mathbf{q}_{Q,(K)}^*\| \leq \sqrt{D_0 \Phi K}\}$, we have

$$-\frac{1}{\Phi K} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^2 + D_0 \leq -\frac{\delta}{\sqrt{K}} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|, \text{ if } \mathbf{q}[t] \in \mathcal{B}^c.$$

On the other hand, when $\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\| \leq \sqrt{D_0 \Phi K}$, it is clear that $-(1/\Phi K) \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^2 + D_0 \leq \eta$ for some $\eta > 0$. Combining these facts, we have

$$-\frac{1}{\Phi K} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^2 + D_0 \leq -\frac{\delta}{K} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\| \mathbb{1}_{\mathcal{B}^c}(\mathbf{q}[t]) + \eta \mathbb{1}_{\mathcal{B}}(\mathbf{q}[t]). \quad (37)$$

Substituting (37) into (34) yields:

$$0 \leq \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{\mathbf{q} \in \mathcal{Z}_+^N} \Pr(\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0]) \times$$
$$\left( -\frac{\delta}{K} \|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\| \mathbb{1}_{\mathcal{B}^c}(\mathbf{q}) + \eta \mathbb{1}_{\mathcal{B}}(\mathbf{q}) \right)$$
$$= \eta \sum_{\mathbf{q} \in \mathcal{B}} \pi_{\mathbf{q}}^\infty - \frac{\delta}{\sqrt{K}} \sum_{\mathbf{q} \in \mathcal{B}^c} \|\mathbf{q} - \mathbf{q}_{Q,(K)}^*\| \pi_{\mathbf{q}}^\infty. \quad (38)$$

where we use the fact that, $\forall \mathbf{q} \in \mathcal{Z}_+^N$, $\lim_{T\to\infty} \frac{1}{T} \sum_{t=0}^{T-1} \Pr\{\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0]\} = \pi_{\mathbf{q}}^\infty$. Re-arranging the terms and with some manipulations, the above inequality can be written as:

$$\frac{\delta}{\sqrt{K}} \sum_{\mathbf{q}\in\mathcal{Z}_+^N} \|\mathbf{q}-\mathbf{q}_{Q,(K)}^*\| \pi_{\mathbf{q}}^\infty \leq \sum_{\mathbf{q}\in\mathcal{B}} \left( \eta + \frac{\delta}{\sqrt{K}}\|\mathbf{q} - \mathbf{q}_{Q,(K)}^*\| \right) \pi_{\mathbf{q}}^\infty$$

$$\leq (\eta + \delta\gamma) \sum_{\mathbf{q}\in\mathcal{B}} \pi_{\mathbf{q}}^\infty \leq (\eta + \delta\gamma), \tag{39}$$

where the second inequality follows from the definition of $\mathcal{B}$. Note here that the left-hand-side is precisely $\frac{\delta}{\sqrt{K}}\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|\}$. Thus, multiplying both sides by $\sqrt{K}/\delta$, we have:

$$\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|\} \leq \left( \gamma + \frac{\eta}{\delta} \right) \sqrt{K} = O(\sqrt{K}). \tag{40}$$

$\underline{Case\ II):}$ $Q \leq Q^\sharp$ such that $\mathbb{E}\{(\mathbf{q}^\infty)^\top(\mathbf{s}^* - \mathbf{s}_Q^\infty)\} > 0$: In this case, we set $\alpha = 1$. It thus follows from (30) that:

$$\mathbb{E}\{\Delta V_1(\mathbf{q}[t])|\mathbf{q}[t]\} \leq -\frac{1}{\Phi K^2} \left\| \mathbf{q}[t] - \mathbf{q}_{Q,(K)}^* \right\|^2 + \frac{1}{K}\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*)^\top\|C_{(Q)} + \frac{D_0}{K}, \tag{41}$$

where $C_{(Q)}$ is defined in the proof of Proposition 1 (cf. Eq. (22)). Note that (41) is identical to (26). Then, following exactly the same steps as in the proof of Proposition 1, we have:

$$\mathbb{E}\{\Delta V_1(\mathbf{q}[t])|\mathbf{q}[t]=\mathbf{q}\} \leq -\frac{\delta_1}{K}\|\mathbf{q}-\mathbf{q}_{Q,(K)}^*\|\mathbb{1}_{\mathcal{B}_1^c}(\mathbf{q})+\eta_1\mathbb{1}_{\mathcal{B}_1}(\mathbf{q}).$$

where $\delta_1$, $\eta_1$, and $\mathcal{B}_1$ are the same as in the proof of Proposition 1. Then, it follows from (31) that

$$\mathbb{E}\{V_1(\mathbf{q}[T]|\mathbf{q}[0])\} - V_1(\mathbf{q}[0]) \leq \eta_1 \sum_{\mathbf{q}\in\mathcal{B}_1} \sum_{t=0}^{T-1} \Pr\{\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0]\}$$

$$- \frac{\delta_1}{K} \sum_{\mathbf{q}\in\mathcal{B}_1^c} \|\mathbf{q} - \mathbf{q}_{(K)}^*\| \sum_{t=0}^{T-1} \Pr\{\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0]\}. \tag{42}$$

Following similar steps as in Case I to divide $T$ on both sides on (42) and let $T \to \infty$, we have $0 \leq \eta_1 \sum_{\mathbf{q}\in\mathcal{B}_1} \pi_{\mathbf{q}}^\infty - \frac{\delta_1}{K} \sum_{\mathbf{q}\in\mathcal{B}_1^c} \|\mathbf{q}-\mathbf{q}_{Q,(K)}^*\|\pi_{\mathbf{q}}^\infty$. Re-arranging the terms and with some manipulations, the above inequality can be written as:

$$\frac{\delta_1}{K} \sum_{\mathbf{q}\in\mathcal{Z}_+^N} \|\mathbf{q}-\mathbf{q}_{Q,(K)}^*\|\pi_{\mathbf{q}}^\infty \leq \sum_{\mathbf{q}\in\mathcal{B}_1} \left( \eta_1+\frac{\delta_1}{K}\|\mathbf{q}-\mathbf{q}_{Q,(K)}^*\| \right) \pi_{\mathbf{q}}^\infty$$

$$\leq (\eta_1 + \delta_1\gamma_1) \sum_{\mathbf{q}\in\mathcal{B}} \pi_{\mathbf{q}}^\infty \leq (\eta_1 + \delta_1\gamma_1),$$

where $\gamma_1$ is the same as in the proof of Proposition 1. Note that the left-hand-side is $\frac{\delta_1}{K}\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|\}$. Multiplying both sides by $\frac{K}{\delta_1}$, we have:

$$\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|\} \leq \left( \gamma_1 + \frac{\eta_1}{\delta_1} \right) K$$

$$= \left( \left[ (C_{(Q)}\Phi) + \sqrt{(C_{(Q)}\Phi)^2 + 4D_0\Phi} \right] + \frac{\eta}{\delta} \right) K = O(C_{(Q)}K).$$

This completes the proof of Theorem 2.

# F    Proof of Theorem 3

To show the results in Theorem 3, we first note that $\mathbb{E}\{a_n[t]|q_n[t]\} = \min\{U_n'^{-1}(\frac{q_n[t]}{K}, A^{\max})\}$ and $a_n^* = U_n'^{-1}(\frac{q_n^*}{K})$, $\forall n$. Thus, we have:

$$
\begin{aligned}
\|\mathbf{a}_Q^\infty - \mathbf{a}_Q^*\| &\le \|\mathbf{a}_Q^\infty - \mathbf{a}_Q^*\|_1 \\
&= \sum_{n=1}^{N} \left| \mathbb{E}\Big\{ \min\Big\{ U_n'^{-1}\Big(\frac{q_n^\infty}{K}, A^{\max}\Big)\Big\}\Big\} - U_n'^{-1}\Big(\frac{q_{Q,(K),n}^*}{K}\Big)\right| \\
&\overset{(a)}{\le} \sum_{n=1}^{N} \mathbb{E}\Big\{ \Big| \min\Big\{ U_n'^{-1}\Big(\frac{q_n^\infty}{K}, A^{\max}\Big)\Big\} - U_n'^{-1}\Big(\frac{q_{Q,(K),n}^*}{K}\Big)\Big|\Big\} \\
&\overset{(b)}{\le} \sum_{n=1}^{N} \mathbb{E}\Big\{ \Big| U_n'^{-1}\Big(\frac{q_n^\infty}{K}\Big) - U_n'^{-1}\Big(\frac{q_{Q,(K),n}^*}{K}\Big)\Big|\Big\} \\
&\overset{(c)}{=} \sum_{n=1}^{N} \mathbb{E}\Big\{ \Big| \big[U_n'^{-1}\big(\tfrac{\tilde{q}_n}{K}\big)\big]' \Big(\frac{q_n^\infty}{K} - \frac{q_{Q,(K),n}^*}{K}\Big)\Big|\Big\} \\
&\overset{(d)}{\le} \sum_{n=1}^{N} \mathbb{E}\Big\{ \Big| \frac{1}{U_n''(\frac{\tilde{q}_n}{K})}\Big| \Big|\frac{q_n^\infty}{K} - \frac{q_{Q,(K),n}^*}{K}\Big|\Big\} \\
&\le \sum_{n=1}^{N} \mathbb{E}\Big\{ \frac{1}{\phi K}|q_n^\infty - q_{Q,(K),n}^*|\Big\} = \frac{1}{\phi K}\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|_1\} \\
&\le \frac{\sqrt{N}}{\phi K}\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|\},
\end{aligned}
\tag{43}
$$

where $(a)$ follows from Jensen's inequality and the convexity of the $L^1$-norm; $(b)$ follows from relaxing the projection onto $[0, A^{\max}]$; $(c)$ follows from the mean value theorem; and $(d)$ follows from the inverse function lemma. Recall in the proof of Theorem 2 (cf. (32)), we have $0 \le J + \sum_{\mathbf{q} \in \mathbb{Z}_+^N} \pi_{\mathbf{q}}^\infty(\mathbf{q})^\top(\mathbf{s}^* - \mathbf{s}_Q^\infty) = J + \mathbb{E}\{(\mathbf{q}^\infty)^\top(\mathbf{s}^* - \mathbf{s}_Q^\infty)\}$. Again, based on the positivity of the term $\mathbb{E}\{(\mathbf{q}^\infty)^\top(\mathbf{s}^* - \mathbf{s}_Q^\infty)\}$, we consider two cases:

_Case I):_ $Q > Q^\sharp$ such that $\mathbb{E}\{(\mathbf{q}^\infty)^\top(\mathbf{s}^* - \mathbf{s}_Q^\infty)\} \le 0$: In this case, we can again discard $\mathbb{E}\{(\mathbf{q}^\infty)^\top(\mathbf{s}^* - \mathbf{s}_Q^\infty)\}$ in (32) and let $\alpha = 0$ to obtain:

$$
0 \le \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{\mathbf{q} \in \mathcal{Z}_+^N} \Pr(\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0])\Big\{ -\frac{1}{\Phi K}\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\|^2\Big\} + D_0.
$$

By re-arranging, multiplying both sides by $\Phi K$, and noting that $\lim_{T \to \infty} \frac{1}{T}\sum_{t=0}^{T-1} \Pr\{\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0]\} = \pi_{\mathbf{q}}^\infty$, we have

$$
\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|^2\} \le D_0 \Phi K.
\tag{44}
$$

It then follows from (43) that

$$\|\mathbf{a}_Q^\infty - \mathbf{a}_Q^*\|^2 \leq \left(\frac{\sqrt{N}}{\phi K}\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|\}\right)^2$$

$$\overset{(a)}{\leq} \frac{N}{\phi^2 K^2}\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|^2\} \overset{(b)}{\leq} \frac{N}{\phi^2 K^2}D_0\Phi K = \frac{ND_0}{\phi^2 K}, \tag{45}$$

where $(a)$ follows from Jensen's inequality; and $(b)$ follows from (44). Taking square root on both sides of (45) yields $\|\mathbf{a}_Q^\infty - \mathbf{a}_Q^*\| = O(1/\sqrt{K})$.

   <u>Case II):</u> $Q \leq Q^\sharp$ such that $\mathbb{E}\{(\mathbf{q}^\infty)^\top(\mathbf{s}^* - \mathbf{s}_Q^\infty)\} > 0$: In this case, we set $\alpha = 1$ and it follows from (30) that:

$$\mathbb{E}\{\Delta V_1(\mathbf{q}[t])|\mathbf{q}[t]\} \leq -\frac{1}{\Phi K^2}\left\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\right\|^2 +$$

$$\frac{C_{(Q)}}{K}\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*)^\top\| + \frac{D_0}{K}$$

$$= -\frac{1}{\Phi K^2}\left(\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\| - \frac{C_{(Q)}\Phi K}{2}\right)^2 + D, \tag{46}$$

where $C_{(Q)}$ is defined in the proof of Proposition 1 (cf. Eq. (22)) and $D \triangleq \frac{C_{(Q)}}{4} + \frac{D_0}{\Phi K}$. Telescoping the inequality in (46) from $t = 0$ to $T - 1$ yields:

$$\mathbb{E}\{V_1(\mathbf{q}[T]|\mathbf{q}[0])\} - V_1(\mathbf{q}[0]) \leq -\frac{1}{\Phi K^2}\sum_{t=0}^{T-1}\sum_{\mathbf{q}\in\mathbb{Z}_+^N}\Pr\{\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0]\}$$

$$\times \left(\|\mathbf{q}[t] - \mathbf{q}_{Q,(K)}^*\| - \frac{C_{(Q)}\Phi K}{2}\right)^2 + DT. \tag{47}$$

Dividing both sides of (47) by $\frac{T}{K^2}$, letting $T \to \infty$, and noting that $\lim_{T\to\infty}\frac{1}{T}\sum_{t=0}^{T-1}\Pr\{\mathbf{q}[t] = \mathbf{q}|\mathbf{q}[0]\} = \pi_\mathbf{q}^\infty, \forall \mathbf{q} \in \mathcal{Z}_+^N$, we have that:

$$\mathbb{E}\left\{\left(\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\| - \frac{C_{(Q)}\Phi K}{2}\right)^2\right\} \leq D\Phi K^2.$$

Taking square root on both sides yields:

$$\left[\mathbb{E}\left\{\left(\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\| - \frac{C_{(Q)}\Phi K}{2}\right)^2\right\}\right]^{\frac{1}{2}} \leq K\sqrt{D\Phi}. \tag{48}$$

Moreover, examining the left-hand-side of (48), we have

$$\left[\mathbb{E}\left\{\left(\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\| - \frac{C_{(Q)}\Phi K}{2}\right)^2\right\}\right]^{\frac{1}{2}}$$

$$\overset{(a)}{\geq} \mathbb{E}\left\{\left[\left(\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\| - \frac{C_{(Q)}\Phi K}{2}\right)^2\right]^{\frac{1}{2}}\right\}$$

$$= \mathbb{E}\left\{\left|\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\| - \frac{C_{(Q)}\Phi K}{2}\right|\right\}$$

$$\geq \mathbb{E}\left\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\| - \frac{C_{(Q)}\Phi K}{2}\right\}$$

$$= \mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|\} - \frac{C_{(Q)}\Phi K}{2}, \tag{49}$$

where $(a)$ follows from Jensen's inequality. Combining (43), (48), and (49) yields:

$$\|\mathbf{a}_Q^\infty - \mathbf{a}_Q^*\| \leq \frac{\sqrt{N}}{\phi K}\mathbb{E}\{\|\mathbf{q}^\infty - \mathbf{q}_{Q,(K)}^*\|\} = \frac{\sqrt{N}}{\phi K}\left(\frac{C_{(Q)}\Phi K}{2} + K\sqrt{D\Phi}\right) = O(C_{(Q)}).$$

Note that Cases I and II are exactly the same results as stated in Theorem 3. This completes the proof.

# References

[1] Y. Zhu, Z. Zhang, Z. Marzi, C. Nelson, U. Madhow, B. Y. Zhao, and H. Zheng, "Demystifying 60 GHz outdoor picocells," in *Proc. ACM MobiCom*, Maui, HI, September 2014, pp. 5 – 16.

[2] S. Sur, V. Venkateswaran, X. Zhang, and P. Ramanathan, "60 GHz indoor networking through flexible beams: A link-level profiling," in *Proc. ACM SIGMETRICS*, Portland, OR, June 2015, pp. 71–84.

[3] S. Rangan, T. S. Rappaport, and E. Erkip, "Millimeter-wave cellular wireless networks: Potentials and challenges," *Proceedings of the IEEE*, vol. 102, no. 3, pp. 366–385, mar 2014.

[4] H. Zhao, R. Mayzus, S. Sun, M. Samimi, J. K. Schulz, Y. Azar, K. Wang, G. N. Wong, F. Gutierrez, and T. S. Rappaport, "28 GHz millimeter wave cellular communication measurements for reflection and penetration loss in and around buildings in New York City," *IEEE Commun. Mag.*, vol. 51, no. 7, pp. 154–160, July 2013.

[5] Z. Pi and F. Khan, "An introduction to millimeter-wave mobile broadband systems," *IEEE Commun. Mag.*, vol. 49, no. 6, pp. 101–107, June 2011.

[6] A. Alkhateeb, O. E. Ayach, GeertLeus, and R. W. Heath, "Channel estimation and hybrid precoding for millimeter wave cellular systems," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 831–846, October 2014.

[7] J. Wang et al., "Beam codebook based beamforming protocol formulti-Gbps millimeter-waveWPAN systems," *IEEE J. Sel. Areas Commun.*, vol. 27, no. 8, pp. 1390–1399, August 2009.

[8] S. Hur, T.Kim, D.Love, J. Krogmeier, T. Thomas, and A. Ghosh, "Millimeter wave beamforming for wireless backhaul and access in small cell networks," *IEEE Trans. Commun.*, vol. 61, no. 10, pp. 4391–4403, October 2013.

[9] S. Han, I. Chih-Lin, Z. Xu, and C. Rowell, "Large-scale antenna systems with hybrid analog and digital beamforming for millimeter wave 5G," *IEEE Commun. Mag.*, vol. 53, no. 1, pp. 186–194, Jan. 2015.

[10] X. Zhang, A. Molisch, and S. Kung, "Variable-phase-shift-based RF-baseband codesign for MIMO antenna selection," *IEEE Trans. Signal Process.*, vol. 53, no. 11, pp. 4091–4103, Nov. 2005.

[11] V. Venkateswaran and A. van der Veen, "Analog beamforming in MIMO communications with phase shift networks and online channel estimation," *IEEE Trans. Signal Process*, vol. 58, no. 8, pp. 4131–4143, Aug. 2010.

[12] O. E. Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. W. Heath, "Spatially sparse precoding in millimeter wave MIMO systems," *IEEE Trans. Wireless Commun*, vol. 13, no. 3, pp. 1499–1513, Mar. 2013.

[13] T. Bogale, L. B. Le, A. Haghighat, and L. Vandendorpe, "On the number of RF chains and phase shifters, and scheduling design with hybrid analog-digital beamforming," *IEEE Trans. Wireless Commun.*, vol. 15, no. 5, pp. 3311–3326, May 2015.

[14] M. J. Neely, E. Modiano, and C.-P. Li, "Faireness and optimal stochastic control for heterogeneous networks," *IEEE/ACM Trans. Netw.*, vol. 16, no. 2, pp. 396–409, Apr. 2008.

[15] A. Eryilmaz and R. Srikant, "Fair resource allocation in wireless networks using queue-length-based scheduling and congestion control," *IEEE/ACM Trans. Netw.*, vol. 15, no. 6, pp. 1333–1344, Dec. 2007.

[16] X. Lin, N. B. Shroff, and R. Srikant, "A tutorial on cross-layer optimization in wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 8, pp. 1452–1463, Aug. 2006.

[17] A. M. Hunter, J. G. Andrews, and S. Weber, "Transmission capacity of ad hoc networks with spatial diversity," *IEEE Trans. Wireless Commun.*, vol. 7, no. 12, pp. 5058–5071, Dec. 2008.

[18] J. Wildman, P. H. Nardelli, M. Latva-aho, and S. Weber, "On the joint impact of beamwidth and orientation error on throughput in wireless directional poisson networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 12, pp. 7072–7085, Jun. 2014.

[19] H. Shokri-Ghadikolaei, L. Gkatzikis, and C. Fischione, "Beam-searching and transmission scheduling in millimeter wave communications," in *Proc. IEEE ICC*, London, UK, Jun. 2015, pp. 1292 – 1297.

[20] I. E. Telatar, "Capacity of multi-antenna Gaussian channels," *European Trans. Telecomm.*, vol. 10, no. 6, pp. 585–596, Nov. 1999.

[21] N. Jindal, "MIMO broadcast channels with finite rate feedback," *IEEE Trans. Inf. Theory*, vol. 52, no. 11, pp. 5045–5059, Nov. 2006.

[22] E. G. Larsson, O. Edfors, and T. L. Marzetta, "Massive MIMO for next generation wireless systems," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 186–195, Feb. 2014.

[23] M. S. Bazaraa, H. D. Sherali, and C. M. Shetty, *Nonlinear Programming: Theory and Algorithms*, 3rd ed. New York, NY: John Wiley & Sons Inc., 2006.

[24] S. P. Meyn and R. L. Tweedie, *Markov Chains and Stochastic Stability*, 2nd ed. Cambridge, UK: Cambridge University Press, 2009.