

Multiuser Beam Allocation for Millimeter Wave Massive MIMO Systems

Xuyao Sun and Chenhao Qi
School of Information Science and Engineering
Southeast University, Nanjing 210096, China
Email: qch@seu.edu.cn

Abstract—In this paper multiuser beam allocation for millimeter wave (mmWave) massive MIMO systems is investigated, based on an improved hybrid precoding design framework. The framework includes two stages. In the first stage, an orthogonal pilot (OP) based beam training scheme is proposed, where all users can simultaneously perform the beam training with the base station (BS) and all the RF chains at the BS are fully utilized. In the second stage, a channel estimation method based on the results from the beam training is presented without transmitting any pilot sequences. Note that users close in geographical may share the same beam from the base station (BS) and cause the beam conflict. To mitigate the multiuser interference caused by beam conflicts, a quality of service (QoS) constrained beam allocation scheme is proposed, with the objective to maximize the equivalent channel gain for the users satisfying QoS constraints as well as maximizing the number of users satisfying QoS constraints on the premise of no beam conflict for all users. Simulation results verify the effectiveness of the proposed schemes and show that the QoS constrained beam allocation scheme can achieve higher spectral efficiency than existing schemes.

Index Terms—Millimeter wave communications, massive MIMO, hybrid precoding, beam allocation

I. INTRODUCTION

The combination of millimeter wave (mmWave) communications and massive multi-input multi-output (MIMO) has been regarded as a frontier for future wireless communication systems [1], since it has the potential to dramatically improve wireless access and throughput. Specifically, the mmWave band ranging from 30 GHz to 300 GHz can considerably increase the data rate benefiting from its abundant frequency resource [2]. On the other hand, the transceivers can pack large antenna arrays into small form factors at mmWave frequencies, making it possible to compensate the high path loss caused by high carrier frequency [3].

Since the number of antennas is large and the working frequency is much higher than conventional MIMO systems, hybrid precoding including analog precoding and digital precoding is usually adopted for mmWave massive MIMO communications. In particular, the directional mmWave beam alignment with the channel main path is important for the analog precoding design, as the precision of the beam alignment determines the channel gain and therefore the achievable rate. In [4], in order to fasten the beam alignment, a beam training scheme based on hierarchical codebook is proposed. In [5], a multi-resolution codebook based on beamforming

sequence is proposed. However, in multiuser scenario, aside of efficient beam training and beam alignment, the multiuser beam allocation is critical, as users close in geographical may share the same beam from the base station (BS) and cause the beam conflict. For example, the users in indoor environment may be densely distributed and some users are close to each other. The multiuser interference caused by beam conflicts will result in severe system performance degradation. In [6], in order to reduce the multiuser interference, an user selection algorithm is proposed where only a small subset of users are selected to be served by the BS. In [7], three beam selection algorithms based on three different criterion are proposed for beamspace mmWave massive MIMO systems, where each user served by the BS is equipped with a single omnidirectional antenna. In [8], all users are classified into two user groups including the interference-users (IUs) and non-interference-users (NIUs) and an interference-aware (IA) beam selection algorithm is proposed to mitigate multiuser interference for mmWave massive MIMO systems.

In this paper, we consider multiuser beam allocation for mmWave massive MIMO systems, based on an improved hybrid precoding design framework. The framework includes two stages. In the first stage, we propose an orthogonal pilot (OP) based beam training scheme, where all users can simultaneously perform the beam training with the BS and all the RF chains at the BS are fully utilized. In the second stage, we present a channel estimation method based on the results from the first stage without transmitting any pilot sequences. Note that users close in geographical may share the same beam from the BS and cause the beam conflict. To mitigate the multiuser interference caused by beam conflicts, we propose a quality of service (QoS) constrained beam allocation scheme, with the objective to maximize the equivalent channel gain of the QoS-satisfied users, under the premise that the number of the QoS-satisfied users without beam conflict is maximized.

The notations used in this paper are defined as follows. Symbols for matrices (upper case) and vectors (lower case) are in boldface. According to the convention, a , \mathbf{a} , \mathbf{A} and \mathcal{A} denote a scalar, a vector, a matrix and a set, respectively. $[\mathbf{a}]_i$, $[\mathbf{A}]_{i,:}$, $[\mathbf{A}]_{:,j}$ and $[\mathbf{A}]_{i,j}$ represent the i th entry of \mathbf{a} , the i th-row of \mathbf{A} , the j th-column of \mathbf{A} and the entry on the i th-row and j th-column of \mathbf{A} , respectively. $(\cdot)^T$, $(\cdot)^*$, $(\cdot)^H$, $(\cdot)^{-1}$, $|\cdot|$ and $\|\cdot\|_0$, $\|\cdot\|_F$ denote the transpose, the conjugate, the conjugate transpose (Hermitian), the inverse, the absolute

value, the zero norm and the Frobenius norm, respectively. $\mathbf{0}^K$, \mathbf{I}_K and \emptyset are the zero vector of size K , the identity matrix of size K and the empty set, respectively. $\mathcal{CN}(m, \mathbf{R})$ is the complex Gaussian distribution with the mean of m and the covariance matrix \mathbf{R} . $\mathbb{E}[\cdot]$ denotes the expectation. \mathbb{C} is the set of complex number.

II. PROBLEM FORMULATION

We consider a multiuser mmWave massive MIMO communication system with a single BS and K users. The BS with N_{BS} antennas placed in a uniform linear array (ULA) and N_{RF} RF chains ($N_{\text{BS}} \gg N_{\text{RF}} \geq 1$) employs a hybrid precoding architecture, while each user equipment (UE) with N_{UE} ULA antennas and a single RF chain employs an analog-only combining architecture. The maximum number of users that can be simultaneously served by the BS is restricted by the number of its RF chains, i.e., $K \leq N_{\text{RF}}$.

During the downlink transmission, the signal received by the k ($k = 1, 2, \dots, K$)th user is denoted as

$$\mathbf{y}_k = \mathbf{H}_k^{\text{dl}} \mathbf{F}_{\text{RF}} \mathbf{F}_{\text{BB}} \mathbf{s} + \mathbf{n}_k^{\text{dl}} \quad (1)$$

where $\mathbf{s} \triangleq [s_1, s_2, \dots, s_K]^T$ is the vector of data symbols subjecting to the constraint of total transmit power P_{dl} , i.e., $\mathbb{E}[\mathbf{s}\mathbf{s}^H] = \frac{P_{\text{dl}}}{K} \mathbf{I}_K$. $\mathbf{F}_{\text{BB}} \triangleq [\mathbf{f}_1^{\text{BB}}, \mathbf{f}_2^{\text{BB}}, \dots, \mathbf{f}_K^{\text{BB}}] \in \mathbb{C}^{K \times K}$ and $\mathbf{F}_{\text{RF}} \triangleq [\mathbf{f}_1^{\text{RF}}, \mathbf{f}_2^{\text{RF}}, \dots, \mathbf{f}_K^{\text{RF}}] \in \mathbb{C}^{N_{\text{BS}} \times K}$ are the baseband precoder (digital precoder) and RF precoder (analog precoder), respectively. $\mathbf{H}_k^{\text{dl}} \in \mathbb{C}^{N_{\text{UE}} \times N_{\text{BS}}}$ is the channel matrix between the BS and the k th user. $\mathbf{n}_k^{\text{dl}} \sim \mathcal{CN}(0, \sigma_{\text{dl}}^2 \mathbf{I}_{N_{\text{UE}}})$ denotes the noise term where each entry of \mathbf{n}_k^{dl} independently obeys the complex Gaussian distribution with zero mean and variance of σ_{dl}^2 . After being processed by the RF combiner (analog combiner) \mathbf{w}_k at the k th user, we obtain an estimate of s_k as

$$\begin{aligned} \hat{s}_k &= \mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{F}_{\text{RF}} \mathbf{F}_{\text{BB}} \mathbf{s} + \mathbf{w}_k^H \mathbf{n}_k^{\text{dl}} \\ &= \mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{F}_{\text{RF}} \sum_{n=1}^K \mathbf{f}_n^{\text{BB}} s_n + \mathbf{w}_k^H \mathbf{n}_k^{\text{dl}}. \end{aligned} \quad (2)$$

Note that \mathbf{F}_{RF} and \mathbf{w}_k are implemented using phase shifters. The entries of \mathbf{F}_{RF} and \mathbf{w}_k have constant envelope. Further, the angles of the phase shifters are usually quantized and have a finite set of possible value. With these constraints, we have $[\mathbf{F}_{\text{RF}}]_{m,n} = \frac{1}{\sqrt{N_{\text{BS}}}} e^{j\phi_{m,n}}$, $[\mathbf{w}_k]_m = \frac{1}{\sqrt{N_{\text{UE}}}} e^{j\theta_m}$, where $\phi_{m,n}$ and θ_m are quantized angles. Moreover, we normalize \mathbf{F}_{BB} such that $\|\mathbf{F}_{\text{RF}} \mathbf{f}_k^{\text{BB}}\|_F^2 = 1$, $k = 1, 2, \dots, K$, indicating that the hybrid precoder does not provide power gain.

We adopt the widely used geometric mmWave MIMO channel model with L_k scatterers [2]. The uniform linear arrays (ULAs) are equipped at both the BS and users. Then the channel between the BS and the k th user can be expressed as

$$\mathbf{H}_k^{\text{dl}} = \sqrt{\frac{N_{\text{BS}} N_{\text{UE}}}{L_k}} \sum_{l=1}^{L_k} \alpha_l^k \mathbf{a}_{\text{UE}}(\Theta_l^k) \mathbf{a}_{\text{BS}}^H(\Phi_l^k) \quad (3)$$

where α_l^k is the complex gain of the l th path with $\mathbb{E}[|\alpha_l^k|^2] = \bar{\alpha}$. $\Theta_l^k \triangleq \sin(\theta_l^k)$ and $\Phi_l^k \triangleq \sin(\phi_l^k)$ are the angle of arrival (AoA) and angle of departure (AoD) of the l th path, respectively, where $\theta_l^k \in (-\pi/2, \pi/2)$ and $\phi_l^k \in (-\pi/2, \pi/2)$. The antenna array response vectors of the BS and the k th user are denoted as $\mathbf{a}_{\text{BS}}(\Phi_l^k) = \mathbf{u}(N_{\text{BS}}, \Phi_l^k)$ and $\mathbf{a}_{\text{UE}}(\Theta_l^k) = \mathbf{u}(N_{\text{UE}}, \Theta_l^k)$, respectively. $\mathbf{u}(A, \epsilon)$ is defined as

$$\mathbf{u}(A, \epsilon) \triangleq \frac{1}{\sqrt{A}} [1, e^{j\frac{2\pi}{\lambda} d\epsilon}, \dots, e^{j(A-1)\frac{2\pi}{\lambda} d\epsilon}]^T, \quad (4)$$

where λ is the signal wavelength and d is the distance between two adjacent antennas. Normally we set $d = \lambda/2$.

A commonly adopted performance metric is the sum-rate of the system. Therefore, our objective is to efficiently design the hybrid precoder (analog and digital precoders) at the BS and the analog combiner at each user, so that the sum-rate is maximized.

Based on (2), we can write down the achievable rate of the k th user as

$$R_k = \log_2 \left(1 + \frac{\frac{P}{K} |\mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{F}_{\text{RF}} \mathbf{f}_k^{\text{BB}}|^2}{\frac{P}{K} \sum_{i \neq k} |\mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{F}_{\text{RF}} \mathbf{f}_i^{\text{BB}}|^2 + \sigma_{\text{dl}}^2}} \right). \quad (5)$$

The sum-rate is $R_{\text{sum}} = \sum_{k=1}^K R_k$. We adopt the widely used beamsteering codebooks [9], where the analog precoder and analog combiners are formed by the codewords of the codebook. The codebooks at the BS and the users can be denoted as $\mathbf{F}_c = [\mathbf{f}_c(1), \mathbf{f}_c(2), \dots, \mathbf{f}_c(N_{\text{BS}})]$ and $\mathbf{W}_c = [\mathbf{w}_c(1), \mathbf{w}_c(2), \dots, \mathbf{w}_c(N_{\text{UE}})]$, respectively, where

$$\begin{aligned} \mathbf{f}_c(n) &= \mathbf{u}(N_{\text{BS}}, -1 + (2n-1)/N_{\text{BS}}), \\ \mathbf{w}_c(n) &= \mathbf{u}(N_{\text{UE}}, -1 + (2n-1)/N_{\text{UE}}). \end{aligned} \quad (6)$$

Then the sum-rate maximization problem in terms of \mathbf{F}_{RF} , \mathbf{F}_{BB} and \mathbf{w}_k can be formulated as

$$\begin{aligned} \max_{\mathbf{F}_{\text{RF}}, \mathbf{F}_{\text{BB}}, \mathbf{w}_k} \sum_{k=1}^K \log_2 \left(1 + \frac{\frac{P_{\text{dl}}}{K} |\mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{F}_{\text{RF}} \mathbf{f}_k^{\text{BB}}|^2}{\frac{P_{\text{dl}}}{K} \sum_{i \neq k} |\mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{F}_{\text{RF}} \mathbf{f}_i^{\text{BB}}|^2 + \sigma_{\text{dl}}^2}} \right) \\ \text{s.t.} \quad & [\mathbf{F}_{\text{RF}}]_{:,k} = \mathbf{f}_k^{\text{RF}} \in \mathbf{F}_c, \quad k = 1, 2, \dots, K, \\ & \mathbf{w}_k \in \mathbf{W}_c, \quad k = 1, 2, \dots, K, \\ & \|\mathbf{F}_{\text{RF}} \mathbf{f}_k^{\text{BB}}\|_F^2 = 1, \quad k = 1, 2, \dots, K. \end{aligned} \quad (7)$$

Note that (7) is a mixed integer programming problem, which is difficult to tackle. A typical hybrid precoding design to solve (7) is divided into two stages [9]. The first stage includes the beam training and analog precoding design to determine $\{\mathbf{f}_k^{\text{RF}}\}_{k=1}^K$ and $\{\mathbf{w}_k\}_{k=1}^K$, while neglecting the resulting interference among users. When determining $\{\mathbf{f}_k^{\text{RF}}\}_{k=1}^K$ and $\{\mathbf{w}_k\}_{k=1}^K$, the optimization problem can be written as

$$\begin{aligned} \max_{\{\mathbf{f}_k^{\text{RF}}\}_{k=1}^K, \{\mathbf{w}_k\}_{k=1}^K} \left\{ |\mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{f}_k^{\text{RF}}| \right\}_{k=1}^K, \\ \text{s.t.} \quad \mathbf{w}_k \in \mathbf{W}_c, \quad \mathbf{f}_k^{\text{RF}} \in \mathcal{F}_c, \end{aligned} \quad (8)$$

where $|\mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{f}_k^{\text{RF}}|$ is named as the equivalent channel gain. The second stage includes the channel estimation and digital precoding design to mitigate the multiuser interference, where

\mathbf{F}_{BB} is determined based on the criterion of zero forcing (ZF) or minimum mean square error (MMSE) [6].

III. IMPROVED HYBRID PRECODING DESIGN FRAMEWORK

In this section, we present an improved hybrid precoding design framework. The framework includes two stages. In the first stage, we propose an OP based beam training scheme to improve the beam training efficiency of direct exhaustive search. Then we determine the analog precoder and analog combiners. In the second stage, we present a channel estimation method based on the results from the first stage without transmitting any pilot sequences. After that, we determine the digital precoder.

A. OP-based Beam Training and Analog Precoding Design

In time-division duplex (TDD) system, we have the channel reciprocity, i.e., $\mathbf{H}_k^{\text{dl}} = (\mathbf{H}_k^{\text{ul}})^T$, where the superscript ‘‘ul’’ is short for uplink and \mathbf{H}_k^{ul} denotes the uplink channel matrix between the k th user and the BS. In OP-based beam training scheme, all users transmit mutually orthogonal pilot sequences so that the signal from different users can be distinguished at the BS. The pilot sequences are denoted as $\sqrt{\tau P_{\text{ul}}}\phi_k \in \mathbb{C}^{1 \times \tau}$, $k = 1, 2, \dots, K$, with $\phi_k \phi_k^H = 1$, $\phi_k \phi_j^H = 0$, $k \neq j$, where τ ($\tau \geq K$) is the length of the pilot sequence and P_{ul} is the uplink transmit power of each user.

During the beam training, all users select the same codeword from \mathbf{W}_c , e.g., $\mathbf{w}_c(n)$, $n = 1, 2, \dots, N_{\text{UE}}$, as the analog beamforming vector. Given $\mathbf{w}_c(n)$, the BS sequentially selects each codeword from \mathbf{F}_c as the analog combining vector. In particular, the BS can select N_{RF} different codewords instead of only one codeword each time, since the BS has N_{RF} RF chains. The selected N_{RF} codewords in the m ($m = 1, 2, \dots, N_{\text{BS}}/N_{\text{RF}}$)th selection form an analog combining matrix $\mathbf{F}_{\text{RF}}^{(m)}$. Then the received pilot sequences at the BS are expressed as

$$\mathbf{Y}^{(m,n)} = \sum_{k=1}^K \sqrt{\tau P_{\text{ul}}} (\mathbf{F}_{\text{RF}}^{(m)})^H \mathbf{H}_k^{\text{ul}} \mathbf{w}_c(n) \phi_k + (\mathbf{F}_{\text{RF}}^{(m)})^H \mathbf{N}^{\text{ul}} \quad (9)$$

where \mathbf{N}^{ul} represents the uplink channel noise. Each entry of \mathbf{N}^{ul} independently obeys the complex Gaussian distribution with zero mean and variance of σ_{ul}^2 . For the k ($k = 1, 2, \dots, K$)th user, the BS multiplies $\mathbf{Y}^{(m,n)}$ with the conjugate of ϕ_k on the left, obtaining

$$\begin{aligned} \mathbf{r}_k^{(m,n)} &= \frac{\phi_k^*}{\sqrt{\tau P_{\text{ul}}}} (\mathbf{Y}^{(m,n)})^T \\ &= \mathbf{w}_c^T(n) (\mathbf{H}_k^{\text{ul}})^T (\mathbf{F}_{\text{RF}}^{(m)})^* + \frac{\phi_k^* (\mathbf{N}^{\text{ul}})^T}{\sqrt{\tau P_{\text{ul}}}} (\mathbf{F}_{\text{RF}}^{(m)})^* \\ &= \mathbf{w}_c^T(n) \mathbf{H}_k^{\text{dl}} (\mathbf{F}_{\text{RF}}^{(m)})^* + \frac{\phi_k^* (\mathbf{N}^{\text{ul}})^T}{\sqrt{\tau P_{\text{ul}}}} (\mathbf{F}_{\text{RF}}^{(m)})^*. \quad (10) \end{aligned}$$

It is seen that totally $N_{\text{BS}}N_{\text{UE}}/N_{\text{RF}}$ times of beam training between the k th user and the BS are required for the OP-based beam training scheme. We put together $\mathbf{r}_k^{(m,n)}$, $m =$

$1, 2, \dots, N_{\text{BS}}/N_{\text{RF}}$, $n = 1, 2, \dots, N_{\text{UE}}$ as

$$\mathbf{R}_k = \begin{bmatrix} \mathbf{r}_k^{(1,1)} & \mathbf{r}_k^{(2,1)} & \dots & \mathbf{r}_k^{(N_{\text{BS}}/N_{\text{RF}},1)} \\ \mathbf{r}_k^{(1,2)} & \mathbf{r}_k^{(2,2)} & \dots & \mathbf{r}_k^{(N_{\text{BS}}/N_{\text{RF}},2)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{r}_k^{(1,N_{\text{UE}})} & \mathbf{r}_k^{(2,N_{\text{UE}})} & \dots & \mathbf{r}_k^{(N_{\text{BS}}/N_{\text{RF}},N_{\text{UE}})} \end{bmatrix}. \quad (11)$$

Considering the uplink and downlink channel reciprocity, the absolute value of each entry of \mathbf{R}_k under noiseless condition is essentially the equivalent channel gain. The maximization of the equivalent channel gain in (8) is converted to finding the entry with the largest absolute value from \mathbf{R}_k . Suppose the row index and column index of the entry with the largest absolute value from \mathbf{R}_k are denoted as p_k and q_k , respectively. For uplink transmission of the k th user, the best analog beamforming vector is $\mathbf{w}_c(p_k)$ and the best combining vector at the BS is $\mathbf{f}_c(q_k)$. For the downlink transmission of the k th user, owing to the channel reciprocity in TDD system, the best analog beamforming vector at the BS is $\tilde{\mathbf{f}}_k^{\text{RF}} = (\mathbf{f}_c(q_k))^*$ and the best analog combining vector at the user is $\tilde{\mathbf{w}}_k = (\mathbf{w}_c(p_k))^*$. Then the designed analog precoder at the BS is

$$\tilde{\mathbf{F}}_{\text{RF}} = [\tilde{\mathbf{f}}_1^{\text{RF}}, \tilde{\mathbf{f}}_2^{\text{RF}}, \dots, \tilde{\mathbf{f}}_K^{\text{RF}}], \quad (12)$$

and the designed analog combiner which is fed back by the BS [9] to each user is $\tilde{\mathbf{w}}_k$, $k = 1, 2, \dots, K$.

The advantage of the OP-based beam training scheme can be summarized as follows.

- 1) The beam training efficiency is improved, since all users can simultaneously perform the beam training with the BS while N_{RF} RF chains at the BS are also fully utilized. In the direct exhaustive search, the BS performs the beam training with each user one by one instead of in parallel. Therefore, the total times of beam training for the OP-based scheme and the direct exhaustive search are $KN_{\text{BS}}N_{\text{UE}}$ and $N_{\text{BS}}N_{\text{UE}}/N_{\text{RF}}$, respectively, resulting in a ratio of $1/(KN_{\text{RF}})$, e.g., only 0.39% for $K = N_{\text{RF}} = 16$.
- 2) Compared with the direct exhaustive search, the OP-based beam training scheme is more robust against the noise, according to the knowledge of spread spectrum communications, since a pilot sequence in length of τ instead of a single pilot symbol is used.
- 3) Compared with the beam training schemes based on hierarchical codebook where frequent channel feedback for different layers of codebook is required, the OP-based beam training scheme achieves small overhead of feedback, as the BS only needs to feed back the index of the best user codeword to the corresponding user after the beam training.

B. Channel Estimation and Digital Precoding Design

Note that in this stage we do not transmit any pilot sequences. We make the channel estimation based on the results from the previous beam training stage.

Define $\tilde{\mathbf{H}}$ as

$$\tilde{\mathbf{H}} = \begin{bmatrix} \mathbf{w}_1^H \mathbf{H}_1^{\text{dl}} \mathbf{f}_1^{\text{RF}} & \mathbf{w}_1^H \mathbf{H}_1^{\text{dl}} \mathbf{f}_2^{\text{RF}} & \dots & \mathbf{w}_1^H \mathbf{H}_1^{\text{dl}} \mathbf{f}_K^{\text{RF}} \\ \mathbf{w}_2^H \mathbf{H}_2^{\text{dl}} \mathbf{f}_1^{\text{RF}} & \mathbf{w}_2^H \mathbf{H}_2^{\text{dl}} \mathbf{f}_2^{\text{RF}} & \dots & \mathbf{w}_2^H \mathbf{H}_2^{\text{dl}} \mathbf{f}_K^{\text{RF}} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{w}_K^H \mathbf{H}_K^{\text{dl}} \mathbf{f}_1^{\text{RF}} & \mathbf{w}_K^H \mathbf{H}_K^{\text{dl}} \mathbf{f}_2^{\text{RF}} & \dots & \mathbf{w}_K^H \mathbf{H}_K^{\text{dl}} \mathbf{f}_K^{\text{RF}} \end{bmatrix}. \quad (13)$$

According to (2), we have

$$\hat{\mathbf{s}}_k = [\tilde{\mathbf{H}}]_{k,:} \mathbf{F}_{\text{BB}} \mathbf{s} + \mathbf{w}_k^H \mathbf{n}_k^{\text{dl}}, \quad k = 1, 2, \dots, K. \quad (14)$$

Then we stack $\hat{\mathbf{s}}_k$, $k = 1, 2, \dots, K$ together as $\hat{\mathbf{s}} \triangleq [\hat{\mathbf{s}}_1, \hat{\mathbf{s}}_2, \dots, \hat{\mathbf{s}}_K]^T$, having

$$\hat{\mathbf{s}} = \tilde{\mathbf{H}} \mathbf{F}_{\text{BB}} \mathbf{s} + \mathbf{n}^{\text{dl}} \quad (15)$$

where $\mathbf{n}^{\text{dl}} \triangleq [\mathbf{w}_1^H \mathbf{n}_1^{\text{dl}}, \mathbf{w}_2^H \mathbf{n}_2^{\text{dl}}, \dots, \mathbf{w}_K^H \mathbf{n}_K^{\text{dl}}]^T$. It is seen that the design of \mathbf{F}_{BB} relies on the estimation of $\tilde{\mathbf{H}}$.

In fact, we can derive the estimation of $\tilde{\mathbf{H}}$ based on $\mathbf{R}_1, \mathbf{R}_2, \dots, \mathbf{R}_K$, which has already been obtained in (11). Denote the estimate of $\tilde{\mathbf{H}}$ as $\tilde{\tilde{\mathbf{H}}}$. The entry on the i ($i = 1, 2, \dots, K$)th row and j ($j = 1, 2, \dots, K$)th column of $\tilde{\tilde{\mathbf{H}}}$ can be expressed as

$$[\tilde{\tilde{\mathbf{H}}}]_{i,j} = [\mathbf{R}_i]_{p_i, q_j} \quad (16)$$

where p_i and q_j have already been determined during the beam training.

Given $\tilde{\tilde{\mathbf{H}}}$, the ZF digital precoder and MMSE digital precoder are

$$\mathbf{F}_{\text{BB}}^{\text{ZF}} = \tilde{\tilde{\mathbf{H}}}^H (\tilde{\tilde{\mathbf{H}}} \tilde{\tilde{\mathbf{H}}}^H)^{-1}, \quad (17)$$

and

$$\mathbf{F}_{\text{BB}}^{\text{MMSE}} = \tilde{\tilde{\mathbf{H}}}^H \left(\frac{P}{K} \tilde{\tilde{\mathbf{H}}} \tilde{\tilde{\mathbf{H}}}^H + \sigma_{\text{dl}}^2 \mathbf{I}_K \right)^{-1}, \quad (18)$$

respectively. In order to satisfy the total power constraint, each column of designed digital precoder via (17) or (18), denoted as $\tilde{\mathbf{f}}_k^{\text{BB}}$ should be normalized, i.e., $\tilde{\mathbf{f}}_k^{\text{BB}} = \tilde{\mathbf{f}}_k^{\text{BB}} / \|\tilde{\mathbf{F}}_{\text{RF}} \tilde{\mathbf{f}}_k^{\text{BB}}\|_F$ such that $\|\tilde{\mathbf{F}}_{\text{RF}} \tilde{\mathbf{f}}_k^{\text{BB}}\|_F = 1$, $k = 1, 2, \dots, K$.

Now we have designed the analog precoder at the BS, the digital precoder at the BS and the analog combiner at each user as $\tilde{\mathbf{F}}_{\text{RF}}$, $\{\tilde{\mathbf{f}}_k^{\text{BB}}\}_{k=1}^K$, and $\{\tilde{\mathbf{w}}_k\}_{k=1}^K$, respectively.

IV. MULTIUSER BEAM ALLOCATION

Since the resolution of phase shifters limits the number of available codewords, the BS may assign different users with the same codeword when the number of users increases, leading to the same beam from the BS pointing at different users and making $\tilde{\mathbf{H}}$ low-rank. In this case, no matter how we design the digital precoder \mathbf{F}_{BB} at the BS, the product between \mathbf{F}_{BB} and $\tilde{\mathbf{H}}$ is low-rank and thus can not be diagonalized, which causes severe interference among different users thus reducing the sum-rate. Owing to the multipath property of the channel in (7) where $L_k > 1$, as well as the channel power leakage phenomenon [10], there are several alternative beams with large equivalent channel gain. This inspires us to eliminate multiuser interference through the design of beam allocation after finishing the beam training.

Since the beam allocation is not considered in (8), we treat the following beam allocation problem as

$$\max_{\{\mathbf{f}_k^{\text{RF}}\}_{k=1}^K, \{\mathbf{w}_k\}_{k=1}^K} \left\{ \left| \mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{f}_k^{\text{RF}} \right| \right\}_{k=1}^K, \quad (19)$$

$$\text{s.t. } \mathbf{w}_k \in \mathcal{W}_c, \mathbf{f}_k^{\text{RF}} \in \mathcal{F}_c, \quad (20)$$

$$\mathbf{f}_i^{\text{RF}} \neq \mathbf{f}_j^{\text{RF}}, i, j = 1, 2, \dots, K, i \neq j. \quad (21)$$

It is seen that the objective function expressed in (19) is the same as that of (8). Note that (8) is essentially K independent optimization problems. However, in the constraint expressed in (21), it is required that the beams allocated for different users should be different, implying that there is no beam conflict for different users. Therefore, (21) introduces inner relations among K optimization problems and converts it to be a multi-objective optimization problem [11]. As a result, a set of Pareto optimal solutions instead of a single one are usually obtained. Therefore, optimization preference is needed to determine a proper solution from a set of solutions. Generally, we maximize the number of simultaneously served users by the BS, under the premise that these users satisfy the QoS. Therefore, we set the optimization preference as

$$\max_{\{\tilde{\mathbf{f}}_k^{\text{RF}}\}_{k=1}^K, \{\tilde{\mathbf{w}}_k\}_{k=1}^K} \sum_{k=1}^K \mathcal{I}(|\tilde{\mathbf{w}}_k^H \mathbf{H}_k^{\text{dl}} \tilde{\mathbf{f}}_k^{\text{RF}}|, \gamma_k), \quad (22)$$

where

$$\mathcal{I}(x, y) = u(x - y) \quad (23)$$

is a binary decision function and $u(n)$ is a unit step function. γ_k is a threshold related to the quality-of-service (QoS) for the k th user, meaning that only when the equivalent channel gain is greater than γ_k , the QoS for the k th user can be guaranteed. In practice, different users may have different QoS constraints. For example, an user demanding live video service is constrained by a large γ_k while an user demanding audio service is only constrained by a small γ_k . In this optimization preference, we require that the number of users satisfying QoS constraints, i.e., whose equivalent channel gains are greater than γ_k , is maximized. With this optimization preference, the beam allocation problem can be expressed as a multi-objective bilevel optimization problem [12]

$$\max_{\{\mathbf{f}_k^{\text{RF}}\}_{k=1}^K, \{\mathbf{w}_k\}_{k=1}^K} \left\{ \left| \mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{f}_k^{\text{RF}} \right| \right\}_{k=1}^K \quad (24)$$

$$\text{s.t. } \{\{\mathbf{f}_k^{\text{RF}}\}_{k=1}^K, \{\mathbf{w}_k\}_{k=1}^K\} \in$$

$$\text{argmax}_{\{\tilde{\mathbf{f}}_k^{\text{RF}}\}_{k=1}^K, \{\tilde{\mathbf{w}}_k\}_{k=1}^K} \sum_{k=1}^K \mathcal{I}(|\tilde{\mathbf{w}}_k^H \mathbf{H}_k^{\text{dl}} \tilde{\mathbf{f}}_k^{\text{RF}}|, \gamma_k), \quad (25)$$

$$\mathbf{w}_k \in \mathcal{W}_c, \mathbf{f}_k^{\text{RF}} \in \mathcal{F}_c, \quad (26)$$

$$\mathbf{f}_i^{\text{RF}} \neq \mathbf{f}_j^{\text{RF}}, i, j = 1, 2, \dots, K, i \neq j, \quad (27)$$

where we maximize the equivalent channel gain and the number of users satisfying the QoS in the upper level objectives (24) and lower level objective (25), respectively. Therefore, we aim at maximizing the equivalent channel gain, under the

premise that the number of the QoS-satisfied users without any beam conflict is maximized. When $\gamma_1 = \gamma_2 = \dots = \gamma_K = 0$, (25) can be removed, resulting in the equivalence between the optimization problem expressed by (24)-(27) and the optimization problem expressed by (19)-(21). Note that once an user's QoS cannot be satisfied, it is meaningless to continue to maximize its equivalent channel gain. Therefore, we only further maximize the equivalent channel gain for the users satisfying the QoS constraints. We should narrow the set of all users in (24) to a subset of those users satisfying QoS constraints. Denote

$$\mathcal{T}(x, y) = xu(x - y). \quad (28)$$

where $u(n)$ is a unit step function. Then the optimization problem in (24)-(27) can be expressed as

$$\begin{aligned} & \max_{\{\mathbf{f}_k^{\text{RF}}\}_{k=1}^K, \{\mathbf{w}_k\}_{k=1}^K} \left\{ \mathcal{T}(|\mathbf{w}_k^H \mathbf{H}_k^{\text{dl}} \mathbf{f}_k^{\text{RF}}|, \gamma_k) \right\}_{k=1}^K, \quad (29) \\ & \text{s.t. } \{\{\mathbf{f}_k^{\text{RF}}\}_{k=1}^K, \{\mathbf{w}_k\}_{k=1}^K\} \in \\ & \text{argmax}_{\{\{\tilde{\mathbf{f}}_k^{\text{RF}}\}_{k=1}^K, \{\tilde{\mathbf{w}}_k\}_{k=1}^K\}} \sum_{k=1}^K \mathcal{I}(|\tilde{\mathbf{w}}_k^H \mathbf{H}_k^{\text{dl}} \tilde{\mathbf{f}}_k^{\text{RF}}|, \gamma_k), \quad (30) \end{aligned}$$

$$\mathbf{w}_k \in \mathcal{W}_c, \mathbf{f}_k^{\text{RF}} \in \mathcal{F}_c, \quad (31)$$

$$\mathbf{f}_i^{\text{RF}} \neq \mathbf{f}_j^{\text{RF}}, i, j = 1, 2, \dots, K, i \neq j. \quad (32)$$

However, the aforementioned multi-objective bilevel optimization problem is difficult to handle.

To reduce the computational complexity on solving this problem, we suppose that the beams are sequentially allocated to different users. In this context, the user allocated beam earlier has more choices than that allocated beam later. The user will have fewer candidate beams if the priority of this user is low. Therefore, in order to maximize the number of users satisfying QoS, higher priority should be given to the user with a single candidate beam that satisfies QoS. We start the beam allocation from the user with the largest equivalent channel gain. Only when the beam conflict happens, we give the high priority to the user with a single candidate beam to maximize the number of users satisfying QoS.

Now we propose a QoS constrained (QC) beam allocation scheme, as shown in **Algorithm 1**. Note that the beams are formed by the codewords of \mathbf{F}_c and \mathbf{W}_c . The beam allocation is essentially the codewords allocation. We use two vectors denoted as \mathbf{b}_f and \mathbf{u}_f to store the indices of BS codewords and user codewords we finally allocate to the BS and users, respectively. We initialize both \mathbf{b}_f and \mathbf{u}_f to be zero. The set of indices of users for beam allocation, denoted as \mathcal{K} , is initialized to be $\{1, 2, \dots, K\}$. Note that the size of \mathcal{K} gets smaller as the beams are sequentially allocated to different users.

For each user, instead of only selecting the best pair $(\tilde{\mathbf{f}}_k^{\text{RF}}, \tilde{\mathbf{w}}_k)$ that can maximize the equivalent channel gain, we select several pairs so that we have candidate pairs if the beam conflict happens. Firstly, from the l ($l = 1, 2, \dots, N_{\text{BS}}$)th

column of \mathbf{R}_k , $k = 1, 2, \dots, K$, we select the entry with the largest absolute value, denoted as

$$g_k(l) = \max_{i=1, 2, \dots, N_{\text{UE}}} |[\mathbf{R}]_{i;l}|, \quad l = 1, 2, \dots, N_{\text{BS}}. \quad (33)$$

For each user, we find the largest equivalent channel gain corresponding to each BS codeword. We sort $\{g_k(1), g_k(2), \dots, g_k(N_{\text{BS}})\}$ in descending order, obtaining $\mathbf{g}_k^{\text{sort}}$, where the largest entry of $\mathbf{g}_k^{\text{sort}}$ is $g_k^{\text{sort}}(1)$. Then we update $\mathbf{g}_k^{\text{sort}}$ by

$$\mathbf{g}_k^{\text{sort}} \leftarrow \left\{ g_k^{\text{sort}}(i) \mid g_k^{\text{sort}}(i) \geq \gamma_k, i \in \{1, 2, \dots, N_{\text{BS}}\} \right\}. \quad (34)$$

Suppose the length of $\mathbf{g}_k^{\text{sort}}$ is M_k , i.e., $M_k \leftarrow \|\mathbf{g}_k^{\text{sort}}\|_0$, $k = 1, 2, \dots, K$. We denote the index of the BS codeword corresponding to $g_k^{\text{sort}}(l)$, $l = 1, 2, \dots, M_k$ in \mathbf{F}_c as $b_k(l)$, obtaining \mathbf{b}_k . We also denote the index of the user codeword corresponding to $g_k^{\text{sort}}(l)$, $l = 1, 2, \dots, M_k$ in \mathbf{W}_c as $u_k(l)$, obtaining \mathbf{u}_k . Therefore, for each BS codeword, now we find the user codeword with the largest equivalent channel gain satisfying QoS constraint. These steps are summarized in Step 3.

Then we select the largest entry of $\mathbf{g}_k^{\text{sort}}$, $k \in \mathcal{K}$, forming a set $\mathcal{G} \triangleq \{g_k^{\text{sort}}(1), k \in \mathcal{K}\}$. The index of the largest entry of \mathcal{G} is defined as

$$k_{\max} \triangleq \arg \max_{k \in \mathcal{K}} \{g_k^{\text{sort}}(1)\}, \quad (35)$$

which corresponds to the strongest beam.

If $\|\mathbf{g}_{k_{\max}}^{\text{sort}}\|_0 = 1$ indicating that the k_{\max} th user has only one candidate beam and we cannot allocate this beam to the other users, we set $k_a \leftarrow k_{\max}$ and then go to Step 16, where k_a is defined as the index of the user finally allocated with this beam.

Otherwise, we check if there is beam conflict with the other users. If the conflict happens with some other users who have only one candidate beam, i.e.,

$$\Lambda \triangleq \{k \mid \|\mathbf{g}_k^{\text{sort}}\|_0 = 1, b_k(1) = b_{k_{\max}}(1), k \in \mathcal{K} \setminus \{k_{\max}\}\} \quad (36)$$

where $\Lambda \neq \emptyset$, we obtain the index of the largest entry among these users as

$$k_c \triangleq \arg \max_{k \in \Lambda} g_k^{\text{sort}}(1). \quad (37)$$

Then we set $k_a \leftarrow k_c$. If $\Lambda = \emptyset$ indicating there is no beam conflict with single beam users, we simply set $k_a \leftarrow k_{\max}$.

The indices of BS codeword and the user codeword corresponding to $g_{k_a}^{\text{sort}}(1)$ are $b_{k_a}(1)$ and $u_{k_a}(1)$, respectively. Then we allocate this beam to the k_a th user by writing the indices of the codewords into \mathbf{b}_f and \mathbf{u}_f , i.e., $b_f(k_a) \leftarrow b_{k_a}(1)$, $u_f(k_a) \leftarrow u_{k_a}(1)$.

Once this beam has been allocated to the k_a th user, we delete all the candidate beams of the k_a th user by setting $\mathbf{g}_{k_a}^{\text{sort}}$, \mathbf{b}_{k_a} and \mathbf{u}_{k_a} empty. In addition, we have to delete this beam from the candidate beams of all the other users, as

the other users can no longer be allocated with this beam. Therefore, we update $\mathbf{g}_k^{\text{sort}}$, \mathbf{b}_k , and \mathbf{u}_k , $k \in \mathcal{K}$, as

$$\mathbf{g}_k^{\text{sort}} \leftarrow \begin{cases} \emptyset, & \text{if } k = k_a, \\ \mathbf{g}_k^{\text{sort}} \setminus \{g_k^{\text{sort}}(i) | b_k(i) = b_{k_a}(1), i \in \{1, 2, \dots, M_k\}\}, & \text{else,} \end{cases} \quad (38)$$

$$\mathbf{b}_k \leftarrow \begin{cases} \emptyset, & \text{if } k = k_a, \\ \mathbf{b}_k \setminus \{b_k(i) | b_k(i) = b_{k_a}(1), i \in \{1, 2, \dots, M_k\}\}, & \text{else,} \end{cases} \quad (39)$$

and

$$\mathbf{u}_k \leftarrow \begin{cases} \emptyset, & \text{if } k = k_a, \\ \mathbf{u}_k \setminus \{u_k(i) | b_k(i) = b_{k_a}(1), i \in \{1, 2, \dots, M_k\}\}, & \text{else,} \end{cases} \quad (40)$$

respectively. Since the number of the users for us to allocate beams is decreased by one, we update \mathcal{K} by $\mathcal{K} \leftarrow \mathcal{K} \setminus \{k_a\}$. Meanwhile, we update M_k as the length of $\mathbf{g}_k^{\text{sort}}$ by

$$M_k \leftarrow \|\mathbf{g}_k^{\text{sort}}\|_0, \quad k \in \mathcal{K}. \quad (41)$$

We repeat the above steps until one of the following two conditions is satisfied. 1) We finish the beam allocation to all users, i.e., $\mathcal{K} = \emptyset$. For example, two users share three beams. Once each user is allocated with a beam, it is finished. 2) The set of candidate beams is empty, i.e., $\mathcal{G} = \emptyset$. For example, two users share a beam. Once this beam is allocated to either one of the users, it is finished since there is no candidate beam available. Finally, we output \mathbf{b}_f and \mathbf{u}_f , where the k ($k = 1, 2, \dots, K$)th user is allocated with the BS codeword $\mathbf{f}_c(b_f(k))$ and the user codeword $\mathbf{w}_c(u_f(k))$. If $b_f(k) = u_f(k) = 0$, it means that the BS does not allocate any beam to the k th user, for there is no candidate beam available.

Note that during the beam training described in Section III-A, we find the best analog beamforming vector $\mathbf{w}_c(p_k)$ and the best combining vector $\mathbf{f}_c(q_k)$ for uplink transmission, which does not consider the beam conflict and can now be replaced by **Algorithm 1**.

V. SIMULATION RESULTS

Now we evaluate the performance of the proposed beam allocation scheme. Consider an mmWave massive MIMO system with a BS equipped with N_{BS} antennas serving K users. The number of RF chains at the BS is N_{RF} . The number of antennas at each user is N_{UE} . The number of resolvable multipath in mmWave channel is set as 3 for each user, i.e., $L_k = 3$, while the complex channel gain is set as $\alpha_1^k \sim \mathcal{CN}(0, 1)$ and $\alpha_i^k \sim \mathcal{CN}(0, 0.1)$ for $i \neq 1$. The spectral efficiency illustrated in Fig. 1 and Fig. 2 is defined as the sum-rate averaged over the number of users satisfying QoS constraints. We fix the uplink channel SNR as $\text{SNR}_{\text{ul}} = 10 \log_{10}(\bar{\alpha} P_{\text{ul}} / \sigma_{\text{ul}}^2) = 20$ dB for uplink beam training and channel estimation. The downlink SNR is defined as $\text{SNR}_{\text{dl}} = 10 \log_{10}(\bar{\alpha} P_{\text{dl}} / (\sigma_{\text{dl}}^2 K))$. The OP-based beam training scheme is solely performed in simulations, due to

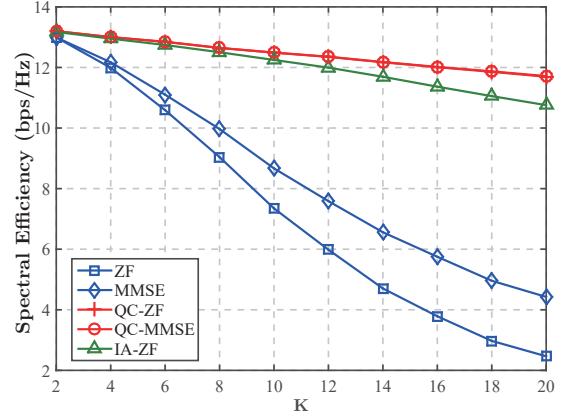


Fig. 1. Comparisons of spectral efficiency for different beam allocation schemes in terms of K .

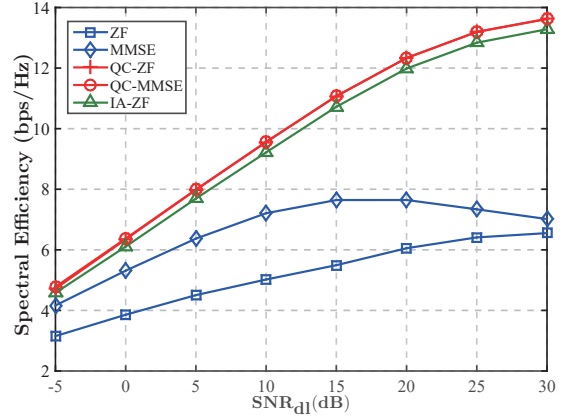


Fig. 2. Comparisons of spectral efficiency for different beam allocation schemes in terms of SNR_{dl} .

its advantage over other schemes. For simplicity, we set $\gamma_1 = \gamma_2 = \dots = \gamma_K = 10\sigma_{\text{dl}}$. Monte Carlo simulations are performed based on 3000 random channel implementations.

As shown in Fig. 1, we compare spectral efficiency for different beam allocation schemes in terms of K . We set $N_{\text{BS}} = 64$, $N_{\text{RF}} = 20$, $N_{\text{UE}} = 16$ and $\text{SNR}_{\text{dl}} = 20$ dB. The curves labeled “ZF” and “MMSE” perform ZF digital precoding as in (17) and MMSE digital precoding as in (18), respectively. Note that the above two curves do not use the beam allocation to solve the problem of beam conflicts, which makes $\bar{\mathbf{H}}$ in (13) low rank and causes the curves to drop rapidly as K increases. Since the MMSE digital precoding can slightly relieve the low rank of $\bar{\mathbf{H}}$, it performs better than the ZF digital precoding. Compared to the curves of “ZF” and “MMSE”, the curves of “QC-ZF” and “QC-MMSE” use the proposed QoS constrained beam allocation scheme in **Algorithm 1**, respectively. As K increases, the beam conflict happens with higher probability. Once the beam conflict happens, the candidate beam with smaller equivalent channel gain is selected for one of the conflicted users, which can effectively mitigate the interference caused by the beam

conflict and therefore stop the curves from fast decreasing like “ZF” and “MMSE”. When $K = 6$, the improvement of spectral efficiency of “QC-ZF” over “ZF” is 21.34%, which verifies the effectiveness of the beam allocation. It is observed that the curves of “QC-ZF” and “QC-MMSE” are almost overlapped. Therefore, once the beam conflict is treated by the beam allocation, the simple ZF digital precoding can be employed. To make comparisons, we also extend the IA beam selection scheme proposed in [8], which is labeled as “IA-ZF”. It is seen that the proposed QC beam allocation scheme outperforms the IA scheme, e.g., 8.92% improvement in spectral efficiency can be achieved when $K = 20$. The reason is that the IA scheme selects the best beam achieving the sum-rate maximization from the group of the interference-users (IUs) at each beam selection, while lacking the overall consideration for the other interference users.

As shown in Fig. 2, we compare spectral efficiency for different beam allocation schemes in terms of SNR_{dl} . We set $N_{\text{BS}} = 64$, $N_{\text{RF}} = 12$, $N_{\text{UE}} = 16$ and $K = 12$. It is seen that the curves labeled “ZF” and “MMSE” climb slowly as SNR_{dl} increases from -5 dB to 15 dB. Since the above two curves do not use beam allocation, where the main factor that affects the system performance is the multiuser interference caused by beam conflicts, the improvement of SNR_{dl} has little contribution to the system performance. Although “MMSE” outperforms “ZF”, the performance gap decreases with the increase of SNR_{dl} from 15 dB to 30 dB, indicating that the effect of noise is reduced and these two estimation methods get close in performance. It is also seen that the curves of “QC-ZF” and “QC-MMSE” using the proposed QoS constrained beam allocation scheme in **Algorithm 1** perform much better than the curves of “ZF” and “MMSE”, since the interference caused by the beam conflict has been effectively mitigated. When $\text{SNR}_{\text{dl}}=15$ dB, the improvement of “QC-ZF” over “MMSE” and over “ZF” in spectral efficiency are 44.95% and 101.71%, respectively. The curves of “QC-ZF” and “QC-MMSE” are almost overlapped, which verifies that the simple ZF digital precoding can be employed once the beam conflict is treated by the beam allocation. Moreover, the proposed QC beam allocation scheme outperforms the existing IA scheme, which has also been illustrated in Fig. 1. In particular, the performance gap between the curves labeled “QC-ZF” and “IA-ZF” keep almost the same as SNR_{dl} increases.

VI. CONCLUSIONS

In this paper, we have presented an improved hybrid precoding design framework including two stages. In the first

stage, we have proposed an OP based beam training scheme. In the second stage, we have presented a channel estimation method based on the results from the beam training without transmitting any pilot sequences. To mitigate the multiuser interference caused by beam conflicts, we have proposed a QoS constrained beam allocation scheme. Simulation results have shown that the proposed beam allocation scheme has higher spectral efficiency than existing schemes. The future work will focus on beam allocation algorithms for multiuser cellular systems taking the out-of-cell interference into consideration.

ACKNOWLEDGMENT

This work is supported in part by the National Natural Science Foundation of China (NSFC) under Grant 61871119 and by the Natural Science Foundation of Jiangsu Province under Grant BK20161428.

REFERENCES

- [1] A. L. Swindlehurst, E. Ayanoglu, P. Heydari, and F. Capolino, “Millimeter-wave massive mimo: The next wireless revolution?” *IEEE Commun. Mag.*, vol. 52, no. 9, pp. 56–62, Sep. 2014.
- [2] R. W. Heath, N. Gonzalez-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, “An overview of signal processing techniques for millimeter wave MIMO systems,” *IEEE J. Sel. Top. Signal Process.*, vol. 10, no. 3, pp. 436–453, Apr. 2016.
- [3] W. Hong, K.-H. Baek, Y. Lee, Y. Kim, and S.-T. Ko, “Study and prototyping of practically large-scale mmWave antenna systems for 5G cellular devices,” *IEEE Commun. Mag.*, vol. 52, no. 9, pp. 63–69, Sep. 2014.
- [4] Z. Xiao, P. Xia, and X.-G. Xia, “Channel Estimation and Hybrid Precoding for Millimeter-Wave MIMO Systems: A Low-Complexity Overall Solution,” *IEEE Access*, vol. 5, pp. 16 100–16 110, July 2017.
- [5] S. Noh, M. D. Zoltowski, and D. J. Love, “Multi-resolution codebook based beamforming sequence design in millimeter-wave systems,” in *2015 IEEE Global Commun. Conf. (GLOBECOM)*, San Diego, CA, USA, Dec. 2014, pp. 1–6.
- [6] W. Yuan, S. M. Armour, and A. Doufexi, “A novel user selection algorithm for multiuser hybrid precoding in mmwave systems,” in *Proc. PIMRC*, Valencia, Spain, Dec. 2016, pp. 1–6.
- [7] P. V. Amadori and C. Masouros, “Low RF-complexity millimeter-wave beamspace-MIMO systems by beam selection,” *IEEE Trans. Commun.*, vol. 63, no. 6, pp. 2212–2223, Jun. 2015.
- [8] X. Gao, L. Dai, Z. Chen, Z. Wang, and Z. Zhang, “Near-optimal beam selection for beamspace mmWave massive MIMO systems,” *IEEE Commun. Lett.*, vol. 20, no. 5, pp. 1054–1057, May 2016.
- [9] A. Alkhateeb, G. Leus, and R. W. Heath, “Limited feedback hybrid precoding for multi-user millimeter wave systems,” *IEEE Trans. Wireless Commun.*, vol. 14, no. 11, pp. 6481–6494, July 2015.
- [10] W. Ma and C. Qi, “Channel Estimation for 3-D Lens Millimeter Wave Massive MIMO System,” *IEEE Commun. Lett.*, vol. 21, no. 9, pp. 2045–2048, Jun. 2017.
- [11] R. T. Marler and J. S. Arora, “Survey of multi-objective optimization methods for engineering,” *Structural and multidisciplinary optimization*, vol. 26, no. 6, pp. 369–395, Apr. 2004.
- [12] A. Sinha, P. Malo, and K. Deb, “A review on bilevel optimization: From classical to evolutionary approaches and applications,” *IEEE Trans. Evol. Comput.*, vol. 22, no. 2, pp. 276–295, Apr. 2018.