

Received October 15, 2021, accepted November 30, 2021, date of publication December 8, 2021, date of current version December 16, 2021.

Digital Object Identifier 10.1109/ACCESS.2021.3134081

ZF-Based Downlink Hybrid Precoding and Combining for Rate Balancing in mmWave Multiuser MIMO Systems

WOO-HEE LIM¹, SUWON JANG², WOOSHYEONG PARK²,
AND JIHOON CHOI¹, (Senior Member, IEEE)

¹School of Electrical and Electronics Engineering, Chung-Ang University, Seoul 06974, South Korea

²School of Electronics and Information Engineering, Korea Aerospace University, Goyang, Gyeonggi-do 10540, South Korea

Corresponding author: Jihoon Choi (jihoon@kau.ac.kr)

This work was supported by the National Research Foundation of Korea (NRF) Grant by the Korean Government through MSIT under Grant 2019R1A2C1006418 and Grant 2021R1A4A2001316.

ABSTRACT This paper proposes a new design strategy for hybrid precoding and combining in the downlink of millimeter wave (mmWave) multiuser multiple-input multiple-output (MU-MIMO) channels. When channel state information is available at the transmitter, the proposed scheme designs the analog precoder and combiners by iteratively factorizing the matrices for fully digital precoding and combining, respectively, using the alternating optimization technique. Then, the digital precoder and combiners are obtained through block diagonalization of effective MU-MIMO channels composed of the analog precoder, MU-MIMO channels, and analog combiners, in order to eliminate inter-user interferences. Moreover, focusing on rate balancing among users, we derive a new power allocation algorithm that exploits a modified gradient descent method. The proposed method iteratively adjusts the power allocated to each user in terms of maximizing the minimum user rate. Through numerical simulations, we verify the convergence of the proposed design procedure for hybrid precoding and combining. Moreover, it is shown that the proposed method outperforms the conventional hybrid precoding and combining methods for rate balancing as well as achieves the minimum user rate close to the performance upper bound.

INDEX TERMS Hybrid precoding, rate balancing, multiuser MIMO, power allocation, zero-forcing.

I. INTRODUCTION

Millimeter wave (mmWave) communication systems have been widely investigated as a means to accommodate rapidly increasing wireless traffic loads. The use of high carrier frequency fundamentally enables the increase in channel capacity [1]–[4], however at the same time, it causes some hurdles to achieving coverage such as huge path loss and rain attenuation [5]–[7]. Since mmWave communication systems can use a large number of antennas and radio frequency (RF) devices by virtue of the reduced wavelength, it is natural to utilize the massive multiple-input multiple-output (MIMO) technology for achieving a huge beamforming gain as a means to mitigate the increase in path loss [8]. Furthermore, it provides a spatial multiplexing gain by simultaneously transmitting data streams using multiple beams.

The associate editor coordinating the review of this manuscript and approving it for publication was Amjad Mehmood¹.

In a mmWave MIMO system, the beamforming and spatial multiplexing gains are achieved by the fully digital precoding, fully analog precoding, and hybrid precoding methods. The fully digital precoding can attain the theoretical channel capacity by arbitrarily adjusting the precoding coefficients, however the implementation cost for large-scale MIMO transceivers is excessive because the number of RF chains is equal to the number of active transmit antennas [6], [7]. In contrast, the fully analog beamforming enables low-complexity implementation with minimum number of RF chains, yet reduces the beamforming gain due to the limited resolution of phase shifters [9]–[11]. As a compromise of the analog and digital precoding methods, the hybrid precoding technique has been studied in many literatures [12]–[24]. The hybrid beamforming scheme interconnects a small number of digital data streams to a large number of RF antennas through two-stage architectures composed of digital processing and analog beamforming,

providing a good tradeoff between performance and complexity. For example, the compressive sensing theory was used to design sparse hybrid precoding and combining matrices in [12], and the theoretical performance of hybrid precoding was analyzed in [13] considering the relationship between the number of RF chains and the number of data streams. Moreover, the hybrid precoder and/or combiner were designed by matrix factorization techniques based on the alternating minimization algorithm [14]–[17], and practical implementation issues were considered such as the limited resolution of phase shifters in analog beamforming [18]–[21] and the limited feedback in closed-loop hybrid precoding [22]–[24].

When the fully digital precoding is used in the downlink of multiuser MIMO (MU-MIMO) systems, the precoding matrix can be designed based on two criteria – one is to maximize the sum rate for total throughput optimization [25]–[29] and the other is to maximize the minimum user rate for fairness [30], [31]. As an extension of point-to-point mmWave communication systems, the hybrid precoding and combining architectures are also considered in the downlink of mmWave MU-MIMO systems [24], [32]–[37]. Several hybrid precoding methods have been developed for MU-MIMO systems using the fully analog combining [32], employing the codebook-based hybrid precoding and fully analog combining [24], and utilizing the phase-shifting analog precoding in combination with digital precoding [33]–[37]. When the analog precoder is implemented by phase shifters, the digital precoder can be designed by the zero-forcing (ZF) [33], [34], regularized channel diagonalization [35], and minimum mean square error (MMSE) techniques [36], [37]. Moreover, the MU-MIMO precoding techniques have been extended to wideband mmWave communication systems with a common analog precoder and frequency-specific digital precoders [38], [39].

This paper proposes a new hybrid precoding method for the downlink of mmWave MU-MIMO systems based on the block diagonalization (BD) and power allocation, when analog precoder and combiners are implemented by phase shifters. Different from the point-to-point communications, the simultaneous transmission of multiple data streams causes inter-stream and/or inter-user interferences in a MU-MIMO channel. When the channel state information (CSI) of all users is known to the transmitter, the inter-stream and inter-user interferences can be completely removed by the BD technique which is a sort of ZF methods, and the MU-MIMO channel can be separated into multiple independent channels. Because the number of RF chains are limited in mmWave MIMO communications, the analog precoder is designed prior to the BD processing. To this end, we first compute the precoder for fully digital precoding and combining exploiting the BD technique of MU-MIMO channels, and then the proposed analog precoder is obtained through matrix factorization of the fully digital precoder. The digital precoder is designed through BD of the effective channels including the analog processing and real channels.

In order to guarantee the fairness among users, the proposed method conducts power allocation to data streams in terms of maximizing the user-wise achievable rate. The main contributions of this paper are summarized as follows.

- Considering the rate balancing among users, we formulate an optimization problem to design the hybrid precoder and combiners for a mmWave MU-MIMO system. To solve the optimal problem, we first propose an iterative algorithm that decomposes the fully digital precoder and combiners into corresponding hybrid precoders and combiners, respectively. The proposed method is applicable to the design of the analog precoder on the transmitter side and the analog combiners on the user side. In the proposed algorithm, the analog precoder is obtained by iteratively updating the factorized matrices in the direction that the Frobenius norm of factorization error matrix is minimized, and the analog combiners are computed in a similar manner.
- The effective MU-MIMO channels are defined by applying the designed analog precoder and combiners to the original channels. Then, the effective channels are converted to independent interference-free channels via the BD technique. A new power allocation problem is formulated under the total power constraint for rate balancing, and it is shown that the objective is quasiconcave. An iterative power allocation algorithm is proposed for maximizing the minimum user rate through some modification of the gradient descent method. The proposed power allocation scheme ensures that all users have the same achievable rate irrespective of channel conditions.
- The complexity of the proposed method is compared with those of existing precoding and combining schemes. Through computer simulations, we evaluate the performance of the proposed method in terms of the minimum user rate, and show that the proposed design method is beneficial compared to conventional hybrid precoding and combining techniques. Moreover, the effect of imperfect CSI is presented through numerical simulations.

The remainder of this paper is organized as follows. Section II introduces the downlink of a mmWave MU-MIMO system with hybrid precoding and combining, and briefly explains the conventional BD technique. Section III describes the proposed ZF-based design method of hybrid precoder and combiners as well as the proposed power allocation algorithm for rate balancing. Complexity analysis and simulation results are provided in Sections IV and V, and conclusions are presented in Section VI.

Notations: Superscripts T , H , $*$, and -1 denote transposition, Hermitian transposition, complex conjugate, and inversion, respectively, for any scalar, vector, or matrix. $|x|$ means the absolute value of x ; the notations $|X|$, $\|X\|$, and $\|X\|_F$ denote the determinant, ℓ_2 -norm, and Frobenius-norm of matrix X , respectively; \mathbf{I}_m represents an $m \times m$ identity matrix; $\mathbf{0}_{m \times n}$ and $\mathbf{1}_{m \times n}$ denote the $m \times n$ zero matrix and

all-ones matrix, respectively; $\text{tr}(\mathbf{A})$ is the trace operation of matrix \mathbf{A} ; $(\mathbf{A})_{m,n}$ means the (m,n) th entry of \mathbf{A} ; $\text{diag}(\mathbf{x})$ returns a diagonal matrix whose main diagonal elements are equal to \mathbf{x} ; $\text{blkdiag}(\{\mathbf{A}\}_{m=1}^M)$ denotes a block-diagonal matrix composed of $\mathbf{A}_1, \dots, \mathbf{A}_M$; $\mathbf{A} \circ \mathbf{B}$ represents the Hadamard product of matrices \mathbf{A} and \mathbf{B} ; and $x \sim \mathcal{CN}(0, \sigma^2)$ means that a complex random variable x conforms to a complex normal distribution with zero mean and variance σ^2 . $\mathbb{E}[x]$ stands for the expectation of random variable x .

II. SYSTEM MODEL AND PREVIOUS WORK

This section introduces the system model for the downlink of a mmWave MU-MIMO system using hybrid precoding and combining, and briefly explains the previous work related to the BD approach for ZF-based transmission.

A. SYSTEM MODEL FOR HYBRID PRECODING AND COMBINING

Fig. 1 describes the system model for the downlink MU-MIMO system using fully-connected hybrid precoding and combining composed of M_{RF} transmit RF chains, M transmit antennas, K users with N antennas each, and N_{RF} receive RF chains. For notational convenience, we assume that the users receive the same number of data streams using the same number of antennas, i.e. each user receives L data streams using N antenna elements. Notice that the proposed method derived under this assumption can be applied to a general case with different numbers of antennas at users, if the number of transmit antennas is equal to or greater than the number of total data streams.

The modulated symbol vector $\mathbf{s} \in \mathbb{C}^{KL \times 1}$ satisfying $\mathbb{E}[\mathbf{s}\mathbf{s}^H] = \mathbf{I}_{LK}$ is transmitted to K users via analog/digital hybrid precoding. The transmitted signal is expressed as

$$\mathbf{x} = \mathbf{F}_A \mathbf{F}_D \mathbf{s} \quad (1)$$

where $\mathbf{F}_D \in \mathbb{C}^{M_{RF} \times KL}$ is the baseband digital precoding matrix for adjusting the magnitudes and phases, $\mathbf{F}_A \in \mathbb{C}^{M \times M_{RF}}$ is the RF analog precoding matrix with unit magnitude and phase shifters. Denote that the total transmit power as P_t , and then it holds that $\|\mathbf{F}_A \mathbf{F}_D\|_F^2 \leq P_t$.

Let us denote the channel between the transmitter and user k as $\mathbf{H}_k \in \mathbb{C}^{N \times M}$, i.e. a MIMO flat fading channel. Suppose that the CSIs for all users $\{\mathbf{H}_k; 1 \leq k \leq K\}$ are known to the transmitter. User k conducts the analog/digital hybrid combining to detect the modulated symbol vector $\mathbf{s}_k \in \mathbb{C}^{L \times 1}$, where $\mathbf{s} = [\mathbf{s}_1^T, \mathbf{s}_2^T, \dots, \mathbf{s}_K^T]^T$. At user k , the received signal after hybrid combining is denoted as

$$\begin{aligned} \mathbf{r}_k &= \mathbf{W}_{D,k}^H \mathbf{W}_{A,k}^H (\mathbf{H}_k \mathbf{x} + \mathbf{n}_k) \\ &= \mathbf{W}_{D,k}^H \mathbf{W}_{A,k}^H \mathbf{H}_k \mathbf{F}_A \mathbf{F}_D \mathbf{s} + \mathbf{W}_{D,k}^H \mathbf{W}_{A,k}^H \mathbf{n}_k, \end{aligned} \quad (2)$$

where $\mathbf{W}_{A,k} \in \mathbb{C}^{N \times N_{RF}}$ is the RF analog precoding matrix with unit magnitude and phase shifters, $\mathbf{W}_{D,k} \in \mathbb{C}^{N_{RF} \times L}$ is the baseband digital precoding matrix for adjusting magnitudes and phases, and $\mathbf{n}_k \in \mathbb{C}^{N \times 1}$ is the noise vector whose elements are independent and identically distributed (i.i.d.)

Gaussian random variables with zero mean and variance σ_k^2 , i.e. $\mathbf{n}_k \sim \mathcal{CN}(\mathbf{0}, \sigma_k^2 \mathbf{I}_L)$. Using the channels $\{\mathbf{H}_k\}$, we design \mathbf{F}_A , \mathbf{F}_D , $\{\mathbf{W}_{A,k}\}$, and $\{\mathbf{W}_{D,k}\}$ for hybrid precoding and combining.

B. BLOCK DIAGONALIZATION FOR MU-MIMO SYSTEM

When fully digital precoding and combining are used in a MU-MIMO system, the BD technique in [25] is introduced based on the generalized channel inversion. Specifically, when the number of transmit antennas is equal to or greater than the number of total data streams, i.e. $M \geq KL$, the BD-based precoding completely removes the inter-user and inter-stream interferences. When a linear precoding is used at the transmitter, the received signal at user k is expressed as

$$\mathbf{y}_k = \mathbf{H}_k \mathbf{F}_{FD} \mathbf{s} + \mathbf{n}_k \quad (3)$$

where $\mathbf{F}_{FD} \in \mathbb{C}^{M \times KL}$ is the fully digital precoding matrix, and $\mathbf{y}_k \in \mathbb{C}^{N \times 1}$ is the received signal vector at user k . Note that $M_{RF} = M$, $\mathbf{F}_D = \mathbf{F}_{FD}$, and $\mathbf{F}_A = \mathbf{I}_M$, because the fully digital precoding method is used. By separating the desired signal from the inter-user interferences, we can rewrite as

$$\begin{aligned} \mathbf{y}_k &= \mathbf{H}_k \mathbf{F}_k \mathbf{s}_k + \mathbf{H}_k \sum_{\ell=1, \ell \neq k}^K \mathbf{F}_\ell \mathbf{s}_\ell + \mathbf{n}_k \\ &= \mathbf{H}_k \mathbf{F}_k \mathbf{s}_k + \mathbf{H}_k \tilde{\mathbf{F}}_k \tilde{\mathbf{s}}_k + \mathbf{n}_k, \end{aligned} \quad (4)$$

where $\mathbf{F}_k \in \mathbb{C}^{M \times L}$ is the precoder for user k , i.e. $\mathbf{F}_{FD} = [\mathbf{F}_1 \mathbf{F}_2 \dots \mathbf{F}_K]$, and $\tilde{\mathbf{F}}_k$ and $\tilde{\mathbf{s}}_k$ are the precoding matrix and transmit symbol vector corresponding to all users except user k , defined as

$$\tilde{\mathbf{F}}_k = [\mathbf{F}_1 \dots \mathbf{F}_{k-1} \mathbf{F}_{k+1} \dots \mathbf{F}_K] \quad (5)$$

$$\tilde{\mathbf{s}}_k = [\mathbf{s}_1^T \dots \mathbf{s}_{k-1}^T \mathbf{s}_{k+1}^T \dots \mathbf{s}_K^T]^T. \quad (6)$$

When the fully digital combining is used, $\mathbf{W}_{A,k} = \mathbf{I}_N$ for all k and the first-stage digital combiner for user k , $\mathbf{W}_k \in \mathbb{C}^{N \times L}$, is composed of the left singular vectors corresponding to the L largest singular values of \mathbf{H}_k through singular value decomposition (SVD). After the first-stage digital combining, we have

$$\begin{aligned} \mathbf{r}_k &= \mathbf{W}_k^H \mathbf{y}_k \\ &= \tilde{\mathbf{H}}_k \mathbf{F}_k \mathbf{s}_k + \tilde{\mathbf{H}}_k \tilde{\mathbf{F}}_k \tilde{\mathbf{s}}_k + \tilde{\mathbf{n}}_k, \end{aligned} \quad (7)$$

where $\tilde{\mathbf{H}}_k = \mathbf{W}_k^H \mathbf{H}_k$ and $\tilde{\mathbf{n}}_k = \mathbf{W}_k^H \mathbf{n}_k$.

The purpose of BD is to make the interference term removed in (7) as well as to maximize the user-wise achievable rate. For notational convenience, we define the interference channel as

$$\tilde{\mathbf{H}}_k = [\tilde{\mathbf{H}}_1^T \dots \tilde{\mathbf{H}}_{k-1}^T \tilde{\mathbf{H}}_{k+1}^T \dots \tilde{\mathbf{H}}_K^T]^T. \quad (8)$$

Suppose that \tilde{L}_k is the rank of $\tilde{\mathbf{H}}_k$, i.e. $\tilde{L}_k = \text{rank}(\tilde{\mathbf{H}}_k)$. Then, we perform SVD for $\tilde{\mathbf{H}}_k$ to remove the inter-user interference as follows:

$$\tilde{\mathbf{H}}_k = \tilde{\mathbf{U}}_k \tilde{\mathbf{\Sigma}}_k \left[\tilde{\mathbf{V}}_k^{(1)} \tilde{\mathbf{V}}_k^{(0)} \right]^H, \quad (9)$$

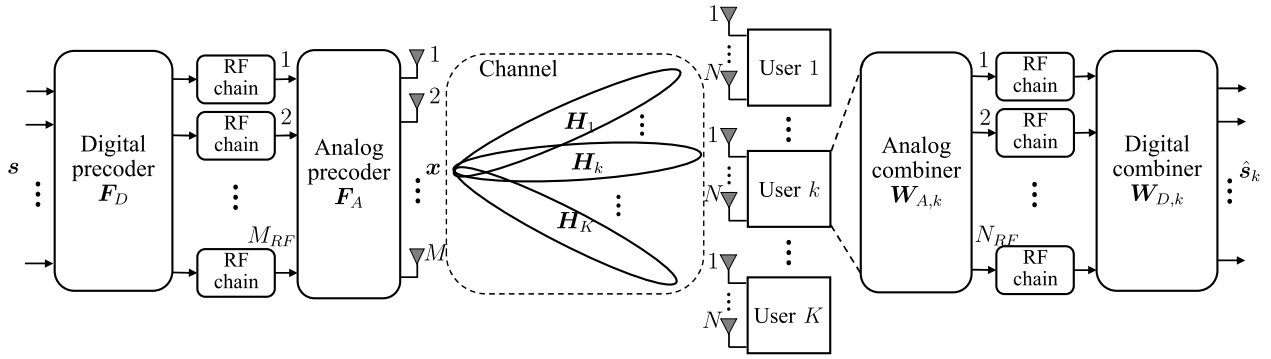


FIGURE 1. Downlink of the MU-MIMO system using hybrid precoding and combining composed of M_{RF} transmit RF chains, M transmit antennas, K users with N antennas each, and N_{RF} receive RF chains.

where $\tilde{U}_k \in \mathbb{C}^{L(K-1) \times L(K-1)}$ is a unitary matrix, $\tilde{\Sigma}_k \in \mathbb{R}^{L(K-1) \times M}$ is a rectangular diagonal matrix with nonnegative diagonal elements, $\tilde{V}_k^{(1)} \in \mathbb{C}^{M \times L_k}$ and $\tilde{V}_k^{(0)} \in \mathbb{C}^{M \times (M-L_k)}$ are the right singular vectors corresponding to nonzero and zero singular values, respectively. Suppose that $L_k = \text{rank}(\tilde{H}_k \tilde{V}_k^{(0)})$. Since $M \geq KL$, it holds $L_k = L$. Applying the SVD again to $\tilde{H}_k \tilde{V}_k^{(0)}$, we have

$$\tilde{H}_k \tilde{V}_k^{(0)} = U_k [\Sigma_k \mathbf{0}] \begin{bmatrix} \mathbf{V}_k^{(1)} & \mathbf{V}_k^{(0)} \end{bmatrix}^H, \quad (10)$$

where $U_k \in \mathbb{C}^{L \times L}$ is a unitary matrix, $\Sigma_k \in \mathbb{R}^{L \times L}$ is a diagonal matrix with nonnegative elements, $\mathbf{V}_k^{(1)} \in \mathbb{C}^{(M-L_k) \times L}$ and $\mathbf{V}_k^{(0)} \in \mathbb{C}^{(M-L_k) \times (M-L_k)}$ are the right singular vectors corresponding to nonzero and zero singular values, respectively. Using the right singular vectors $\tilde{V}_k^{(0)}$ and $\mathbf{V}_k^{(1)}$ in (9) and (10), we design the ZF-based precoder as follows:

$$\mathbf{F}_{FD} = \begin{bmatrix} \tilde{V}_1^{(0)} \mathbf{V}_1^{(1)} & \tilde{V}_2^{(0)} \mathbf{V}_2^{(1)} & \dots & \tilde{V}_K^{(0)} \mathbf{V}_K^{(1)} \end{bmatrix} \mathbf{P}^{\frac{1}{2}} \quad (11)$$

where $\mathbf{P} \in \mathbb{R}^{KL \times KL}$ is a diagonal matrix with nonnegative elements representing the power allocation to individual data streams. From the total transmit power constraint, it holds $\text{tr}(\mathbf{P}) \leq P_t$. Moreover, U_k in (10) is the second-stage digital combiner for the effective channel $\tilde{H}_k \tilde{V}_k^{(0)}$, and the fully digital combiner for user k is denoted as

$$\mathbf{W}_{FD,k} = \mathbf{W}_k U_k, \text{ for } 1 \leq k \leq K. \quad (12)$$

III. PROPOSED HYBRID PRECODING AND COMBINING METHOD

Fig. 2 shows the overall procedure of the proposed design method for hybrid precoding and combining. For the ZF-based design, it is assumed that $M \geq KL$ and $M_{RF} \geq KN_{RF}$. In the following subsections, we propose a design procedure of the analog/digital precoders and combiners for hybrid processing in the downlink MU-MIMO system, and then derive a new power allocation method for rate balancing among users.

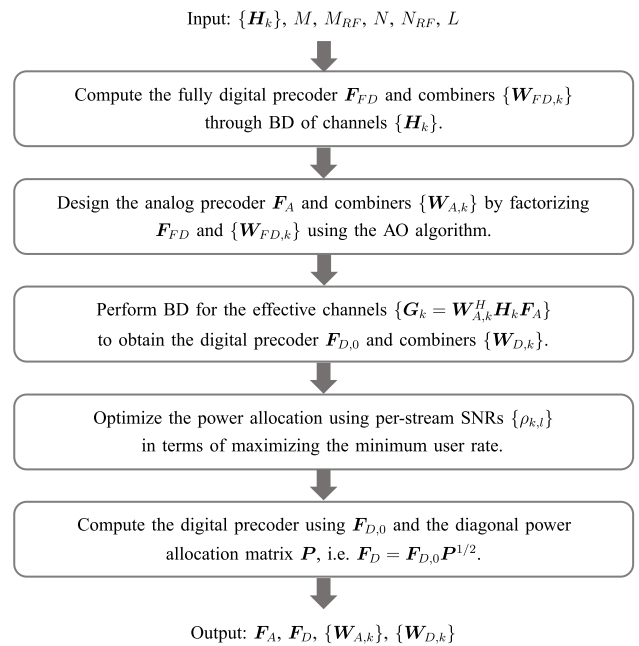


FIGURE 2. Overall procedure of the proposed design method for hybrid precoding and combining.

A. DESIGN OF HYBRID PRECODER AND COMBINERS

As the first step, the fully digital precoder and combiners are evaluated using (7)–(12). Here, the power allocation is not considered, i.e. $\mathbf{P} = \frac{P_t}{KL} \mathbf{I}_{KL}$, because the fully digital precoder \mathbf{F}_{FD} in (11) is utilized only for the design of the analog precoder \mathbf{F}_A . In the proposed design procedure, the analog precoder \mathbf{F}_A and analog combiners $\{\mathbf{W}_{A,k}\}$ are determined by factorizing the fully digital precoder \mathbf{F}_{FD} and combiners $\{\mathbf{W}_{FD,l}\}$, respectively, in the MMSE sense through the alternating optimization (AO) algorithm in [14]. The MMSE matrix factorization problem is formulated as

$$P1: \mathbf{F}_A = \arg \min_{\mathbf{F}_A, \mathbf{F}_D} \|\mathbf{F}_{FD} - \mathbf{F}_A \mathbf{F}_D\|_F^2 \quad (13a)$$

$$s.t. |(\mathbf{F}_A)_{m,n}| = 1, \text{ for } \forall m, n \quad (13b)$$

$$\|\mathbf{F}_A \mathbf{F}_D\|_F^2 = P_t \quad (13c)$$

where (13b) denotes that the analog precoder is implemented by phase shifters and (13c) means the transmit power constraint. In (13c), we use the equality constraint, because the downlink achievable rate is proportional to the transmit power. When \mathbf{F}_A is fixed, the digital precoder minimizing the cost function (13a) is given by

$$\mathbf{F}_D = c\mathbf{F}_A^\dagger \mathbf{F}_{FD}, \quad (14)$$

where \mathbf{A}^\dagger means the pseudo-inverse of \mathbf{A} and c is a scaling factor to meet the transmit power constraint (13c). On the other hand, when \mathbf{F}_D is fixed, an iterative algorithm is derived to find the analog precoder minimizing the cost function (13a) which exploits the conjugate gradient method combined with the Riemannian gradient for the unit modulus solution. When denoting $J(\mathbf{X}(i)) = \|\mathbf{F}_{FD} - \mathbf{X}(i)\mathbf{F}_D\|_F^2$, the Euclidean gradient of $J(\cdot)$ is given by

$$\nabla J(\mathbf{X}(i)) = (\mathbf{F}_{FD} - \mathbf{X}(i)\mathbf{F}_D)(-\mathbf{F}_D^H), \quad (15)$$

where $\mathbf{X}(i) \in \mathbb{C}^{M \times M_{RF}}$ is the analog precoder at i th iteration. To find the solution conforming to the unit modulus constraint, the Riemannian gradient is computed through orthogonal projection of $\nabla J(\mathbf{X})$ onto the tangent space of the complex circle manifold:

$$\text{Proj}(\nabla J(\mathbf{X}(i))) = \nabla J(\mathbf{X}(i)) - \text{Re}(\nabla J(\mathbf{X}(i)) \circ \mathbf{X}(i)^*) \circ \mathbf{X}(i). \quad (16)$$

Suppose that $\mathbf{D}(i) \in \mathbb{C}^{M \times M_{RF}}$ is the conjugate direction at i th iteration. From the Riemannian gradient in (16), the conjugate direction is updated as

$$\mathbf{D}(i+1) = -\text{Proj}(\nabla J(\mathbf{X}(i))) + \beta(i)\text{Proj}(\mathbf{D}(i)) \quad (17)$$

where $\text{Proj}(\mathbf{D}(i))$ is obtained by replacing $\nabla J(\mathbf{X}(i))$ with $\mathbf{D}(i)$ in (16) and $\beta(i)$ is a step-size parameter. Finally, the analog precoder is updated as below:

$$(\mathbf{X}(i+1))_{m,n} = \frac{(\mathbf{X}(i) + \alpha(i)\mathbf{D}(i+1))_{m,n}}{|(\mathbf{X}(i) + \alpha(i)\mathbf{D}(i+1))_{m,n}|} \quad (18)$$

where $\alpha(i)$ is a step-size parameter, $1 \leq m \leq M$, and $1 \leq n \leq M_{RF}$. By repeating (15)–(18) until $\mathbf{X}(i)$ converges, we have the analog precoder minimizing $J(\cdot)$:

$$\mathbf{F}_A = \mathbf{X}(\text{end}) \quad (19)$$

where $\mathbf{X}(\text{end})$ is the finally updated matrix satisfying the convergence criterion. Similarly, the analog combiners are designed by solving the MMSE matrix factorization problem for each user as follows:

$$P2: \mathbf{W}_{A,k} = \arg \min_{\mathbf{W}_{A,k}, \mathbf{W}_{D,k}} \|\mathbf{W}_{FD,k} - \mathbf{W}_{A,k}\mathbf{W}_{D,k}\|_F^2 \quad (20a)$$

$$s.t. |\mathbf{W}_{A,k}|_{m,n}| = 1, \text{ for } \forall m, n. \quad (20b)$$

The problem (P2) is the same as (P1) except the transmit power constraint, thus the analog combiner $\mathbf{W}_{A,k}$ is determined by performing the conjugate gradient method in a similar manner to (15)–(18). In this case, we set $c = 1$ in (14)

because the scaling for the transmit power constraint is not required.

Now, we design the digital precoder and combiners. By utilizing the analog precoder and combiners designed from (P1) and (P2), the effective channel for user k , $\mathbf{G}_k \in \mathbb{C}^{N_{RF} \times M_{RF}}$, is expressed as

$$\mathbf{G}_k = \mathbf{W}_{A,k}^H \mathbf{H}_k \mathbf{F}_A, \text{ for } 1 \leq k \leq K. \quad (21)$$

By replacing $\{\mathbf{H}_k\}$ with $\{\mathbf{G}_k\}$, we carry out the BD procedure for MU-MIMO downlink channels in (8)–(12). Then, from (11), the digital precoder prior to power allocation, $\mathbf{F}_{D,0} \in \mathbb{C}^{M_{RF} \times KL}$, is obtained as

$$\mathbf{F}_{D,0} = \left[\tilde{\mathbf{V}}_1^{(0)} \mathbf{V}_1^{(1)} \tilde{\mathbf{V}}_2^{(0)} \mathbf{V}_2^{(1)} \cdots \tilde{\mathbf{V}}_K^{(0)} \mathbf{V}_K^{(1)} \right], \quad (22)$$

and the digital combiner for user k , $\mathbf{W}_{D,k} \in \mathbb{C}^{N_{RF} \times L}$, is determined from (12) as below:

$$\mathbf{W}_{D,k} = \mathbf{W}_k \mathbf{U}_k. \quad (23)$$

Overall, the proposed algorithm for designing \mathbf{F}_A , $\mathbf{F}_{D,0}$, $\{\mathbf{W}_{A,k}\}$, and $\{\mathbf{W}_{D,k}\}$ is summarized as Algorithm 1.

B. OPTIMIZATION OF POWER ALLOCATION

From (11) and (22), the digital precoder is given by

$$\mathbf{F}_D = \mathbf{F}_{D,0} \mathbf{P}^{1/2}. \quad (24)$$

This subsection proposes a power allocation scheme taking into account the rate balancing among users, i.e. we design \mathbf{P} maximizing the minimum user rate. Substituting \mathbf{F}_A , $\mathbf{F}_{D,0}$, $\{\mathbf{W}_{A,k}\}$, and $\{\mathbf{W}_{D,k}\}$ designed in Section III-A to the received signal in (2), we have

$$\mathbf{r}_k = \mathbf{\Sigma}_k \mathbf{P}_k^{1/2} \mathbf{s}_k + \mathbf{W}_{H,k}^H \mathbf{n}_k, \quad (25)$$

where $\mathbf{\Sigma}_k$ is obtained from the BD of $\{\mathbf{G}_k\}$ shown in (10), $\mathbf{P}_k \in \mathbb{R}^{L \times L}$ is a submatrix of \mathbf{P} satisfying $\mathbf{P} = \text{blkdiag}(\{\mathbf{P}_k\}_{k=1}^K)$, and $\mathbf{W}_{H,k} = \mathbf{W}_{A,k} \mathbf{W}_{D,k}$. Since $\mathbf{\Sigma}_k$ and \mathbf{P}_k are diagonal matrices, the nominal signal-to-noise ratio (SNR) of l th data stream at user k is computed as

$$\gamma_{k,\ell} = \frac{\rho_{k,\ell}^2}{\sigma_k^2 (\mathbf{W}_{H,k}^H \mathbf{W}_{H,k})_{\ell,\ell}} \quad (26)$$

where $\mathbf{\Sigma}_k = \text{diag}(\rho_{k,1}, \rho_{k,2}, \dots, \rho_{k,L})$. Suppose that $p_{k,\ell}$ is the power allocated to the l th data stream of user k , i.e. $\mathbf{P}_k = \text{diag}(p_{k,1}, p_{k,2}, \dots, p_{k,L})$. Then, the achievable rate for user k is expressed as

$$f_k(p_k) = \max_{\{p_{k,\ell}\}} \sum_{\ell=1}^L \log_2(1 + p_{k,\ell} \gamma_{k,\ell}) \quad (27a)$$

$$s.t. \sum_{\ell=1}^L p_{k,\ell} = p_k, p_{k,\ell} \geq 0 \text{ for } \forall \ell. \quad (27b)$$

Here, notice that the achievable rate is a function of p_k describing the power allocated to user k . $f_k(p_k)$ is obtained by the water-filling algorithm that optimally assigns the power

Algorithm 1 Proposed Algorithm for Designing Hybrid Precoders and Combiners

1. **Input:** $F_{FD}, \{W_{FD,k}; 1 \leq k \leq K\}, M_{RF}, N_{RF}, L$
2. Initialize: $j = -1$ and $(F_A(0))_{m,n} = e^{j2\pi\theta_{m,n}}$ for $\forall m, n$ where $\theta_{m,n}$ is a random variable uniformly distributed over $[0, 1)$.
3. **repeat**
4. $j = j + 1$.
5. Compute the digital precoder: $F_D(j) = F_A^\dagger(j)F_{FD}$.
6. Given $F_D(j)$, update the analog precoder until convergence utilizing (15)–(18).
7. Get the new analog precoder $F_A(j + 1)$ from (19).
8. Evaluate the cost function $J(F_A(j + 1))$.
9. **until** $\|J(F_A(j + 1)) - J(F_A(j))\| < \epsilon_F$, where ϵ_F is the precoder tolerance for termination.
10. **for** $k = 1 : K$ **do**
11. Initialize: $j = -1$ and $(W_{A,k}(0))_{m,n} = e^{j2\pi\phi_{m,n}}$ for $\forall m, n$ where $\phi_{m,n}$ is a random variable uniformly distributed over $[0, 1)$.
12. **repeat**
13. $j = j + 1$.
14. Compute the digital combiner: $W_{D,k}(j) = W_{A,k}^\dagger(j)W_{FD,k}$.
15. Update the analog combiner until convergence utilizing (15)–(18), by replacing F_{FD} and F_D with $W_{FD,k}$ and $W_{D,k}(j)$, respectively.
16. Get the new analog combiner: $W_{A,k}(j + 1) = X(end)$.
17. Evaluate the cost function: $J(W_{A,k}(j + 1))$.
18. **until** $\|J(W_{A,k}(j + 1)) - J(W_{A,k}(j))\| < \epsilon_W$, where ϵ_W is the combiner tolerance for termination.
19. **end for**
20. Obtain the effective channels $\{G_k\}$ from (21), and perform the BD procedure in (8)–(12) using $\{G_k\}$ instead of $\{H_k\}$.
21. Compute $F_{D,0}$ and $\{W_{D,k}\}$ from (22) and (23).
22. **Output:** $F_A, F_{D,0}, \{W_{A,k}\}$, and $\{W_{D,k}\}$.

p_k into L independent channels with $\{\gamma_{k,1}, \dots, \gamma_{k,L}\}$. Specifically, the optimal power allocation is expressed as

$$p_{k,\ell} = \left[\lambda - \frac{1}{\gamma_{k,\ell}} \right]_+, \quad (28)$$

where λ is the Lagrange multiplier denoting the L level conforming to the power constraint (27b), and $[x]_+ = \max(0, x)$.

To ensure fairness of achievable rates among users, we formulate the following max-min optimization problem:

$$P3 : p_o = \arg \max_p \min\{f_1(p_1), f_2(p_2), \dots, f_K(p_K)\} \quad (29a)$$

$$s.t. \sum_{k=1}^K p_k = P_t, p_k \geq 0 \text{ for } \forall k, \quad (29b)$$

Algorithm 2: Proposed Power Allocation Algorithm for Rate Balancing

1. **Input:** $\{\rho_{k,\ell}; 1 \leq k \leq K, 1 \leq \ell \leq L\}, \{\sigma_k; 1 \leq k \leq K\}, \{W_{H,k}; 1 \leq k \leq K\}, P_t$
2. Initialize $i = -1$ and $p_k(0) = \frac{P_t}{K}$.
3. Compute the per-stream SNRs $\{\gamma_{k,\ell}\}$ using (26).
4. **repeat**
5. $i = i + 1$.
6. Compute the user rates $\{f_k(p_k(i))\}$ in (27a) using the water-filling solution in (28).
7. Evaluate the numerical gradient of $f_k(p_k)$ using (30).
8. Estimate the optimal max-min user rate f_o using (33).
9. Compute $\Delta p(i)$ from (32).
10. Adjust the parameter $\alpha_p(i)$ using (35).
11. Update the power allocation vector using (34).
12. **until** $\|\Delta p(i)\| < \epsilon_p$, where ϵ_p is the tolerance for termination.
13. **Output:** $p_o = p(i)$.

where $\mathbf{p} = [p_1, p_2, \dots, p_K]^T$. It is difficult to find a closed-form solution for (P3). Instead, we derive an iterative algorithm based on the gradient descent method and the water-filling approach. When $\{\gamma_{k,\ell}\}$ in (26) are fixed, the achievable rate $f_k(p_k)$ is a monotonically increasing function with respect to (w.r.t.) p_k and the objective of (P3) is a quasiconcave function from Property 1. In other words, (P3) has a unique solution that can be found by a convex optimization technique.

Property 1: Let us define

$$g_K(\mathbf{p}) \triangleq \min\{f_1(p_1), f_2(p_2), \dots, f_K(p_K)\}.$$

Given $\{\gamma_{k,\ell}; 1 \leq k \leq K, 1 \leq \ell \leq L\}$, $g_K(\mathbf{p})$ is a quasiconcave function w.r.t. \mathbf{p} under the constraints $p_1 + p_2 + \dots + p_K = P_t$ and $p_k \geq 0$ for $1 \leq k \leq K$.

Proof: See Appendix A. ■

To find the optimal solution of (P3), we evaluate the numerical gradient of $f_k(p_k)$:

$$\nabla f_k(p_k) \cong \frac{f_k(p_k + \Delta p) - f_k(p_k - \Delta p)}{2\Delta p} \quad (30)$$

where Δp is a small positive constant. At the optimal point, it holds that $f_o \triangleq f_1(p_{o,1}) = f_2(p_{o,2}) = \dots = f_K(p_{o,K})$. Let us define the power allocation vector at i th iteration as $\mathbf{p}(i) = [p_1(i), p_2(i), \dots, p_K(i)]^T$. With linear approximation around $p_k(i)$, we may write the difference between the optimal user rate f_o and the actual user rate $f_k(p_k(i))$ as follows:

$$f_o - f_k(p_k(i)) \approx \nabla f_k(p_k(i)) \Delta p_k(i). \quad (31)$$

Also, we can rewrite (31) as

$$\Delta p_k(i) = \frac{f_o - f_k(p_k(i))}{\nabla f_k(p_k(i))}. \quad (32)$$

Due to the transmit power constraint, the total transmit power is not changed, i.e. $\sum_{k=1}^K \Delta p_k(i) = 0$. By applying this condition

to (32), we have

$$f_o = \left(\sum_{k=1}^K \frac{1}{\nabla f_k(p_k(i))} \right)^{-1} \sum_{k=1}^K \frac{f_k(p_k(i))}{\nabla f_k(p_k(i))}. \quad (33)$$

Again, by substituting (33) into (32), we can compute $\Delta p_k(i)$ for all k , and the power allocation vector can be updated in the direction of increasing $g_K(\mathbf{p})$ by employing the gradient descent method:

$$\mathbf{p}(i+1) = \mathbf{p}(i) + \alpha_p(i)\Delta\mathbf{p}(i), \quad (34)$$

where $\Delta\mathbf{p}(i) = [\Delta p_1(i), \Delta p_2(i), \dots, \Delta p_K(i)]^T$ and $\alpha_p(i)$ is given by

$$\alpha_p(i) = \beta_p \min_{1 \leq k \leq K} \left\{ \frac{\min(p_k(i), 1 - p_k(i))}{|\Delta p_k(i)|} \right\} \quad (35)$$

where $0 < \beta_p \leq 1$ is the step-size parameter. The proposed power allocation method is summarized as Algorithm 2. Finally, the optimal power allocation matrix is given by $\mathbf{P}_o = \text{diag}(\mathbf{p}_o)$, and we design the digital precoder \mathbf{F}_D by substituting \mathbf{P}_o into (24).

IV. COMPLEXITY ANALYSIS

We compare the complexity of the proposed hybrid precoding and combining method with those of the fully digital processing and the correlation-based scheme in [37]. We define $Q_1 = M - (K - 1)L$ for the BD in (8)–(12) using $\{\mathbf{H}_k\}$ and $Q_2 = M_{RF} - (K - 1)L$ for the BD using $\{\mathbf{G}_k\}$. In the proposed method, J_1 denotes the number of iterations for the AO method, and the complexity of the AO algorithm for obtaining $\{\mathbf{W}_{A,k}\}$ is neglected, because finding \mathbf{F}_A by the AO requires much more computations than obtaining $\{\mathbf{W}_{A,k}\}$. In the correlation-based scheme, we define $J_2 = C_t + C_r$, where C_t and C_r are the codebook sizes for the analog precoder and analog combiners, respectively.

The number of operations is presented in Table 1. Whereas the fully digital processing requires only the BD of $\{\mathbf{H}_k\}$, the proposed method necessitates additional computations for the AO algorithm designing \mathbf{F}_A , obtaining the effective channels $\{\mathbf{G}_k\}$, and the BD of $\{\mathbf{G}_k\}$. The correlation-based scheme requires slightly less computational load than the proposed method, because the complexity of channel correlations is smaller than that of the AO step, i.e. $O(J_2MM_{RF}N) < O(J_1M^2M_{RF})$. Notice that the proposed method and the correlation-based scheme use only M_{RF} transmit and N_{RF} receive RF chains, while the fully digital processing utilizes M and N RF chains that require excessive hardware implementation cost.

V. SIMULATION RESULTS

In this section, we evaluate the performance of the proposed hybrid precoding and combining method through numerical simulations, and compare the proposed scheme with existing hybrid processing techniques for MU-MIMO systems. Specifically, we consider four hybrid precoding meth-

TABLE 1. Number of operations for various precoding and combining methods.

Step	Fully Digital	Proposed Method	Corr.-Based [37]
BD of $\{\mathbf{H}_k\}$	$O(M^2Q_1K)$	$O(M^2Q_1K)$	$O(M^2Q_1K)$
AO in (14)–(18)	–	$O(J_1M^2M_{RF})$	–
Channel Corr.	–	–	$O(J_2MM_{RF}N)$
$\{\mathbf{G}_k\}$ in (21)	–	$O(MM_{RF}N)$	$O(MM_{RF}N)$
BD of $\{\mathbf{G}_k\}$	–	$O(M_{RF}^2Q_2K)$	$O(M_{RF}^2Q_2K)$

ods and three joint hybrid precoding/combining schemes as follows:

- Fully digital precoding: the fully digital precoder is used at the transmitter (i.e. $M = M_{RF}$, $\mathbf{F}_A = \mathbf{I}_M$, $\mathbf{F}_{D,0} = \mathbf{F}_{FD}$); the fully digital combiners are used at the receiver (i.e. $N_{RF} = N$, $\mathbf{W}_{A,k} = \mathbf{I}_N$, and $\mathbf{W}_{D,k} = \mathbf{W}_{FD,k}$); and the conventional BD technique in Section II-B is applied. This method presents the performance upper bound of the ZF-based multiuser transmission.
- Proposed hybrid precoding: the fully digital combiners $\{\mathbf{W}_{FD,k}\}$ are used at the receivers, and the hybrid precoders \mathbf{F}_A and $\mathbf{F}_{D,0}$ are designed by Algorithm 1 considering the fully digital combining.
- Corr.-based hybrid precoding [37]: following the approach in [37], the analog precoder is constructed by selecting the M_{RF} beamforming vectors having the largest correlations with $\{\mathbf{H}_k\}$, when the fully digital combiners are used at the receivers. The digital precoder is designed via the BD of effective channels.
- Random analog precoding: the analog precoder is defined as an arbitrary matrix with unit modulus via random phase shifting; the fully digital combiners are used at the receivers; and the BD technique is used to determine the digital precoder.
- Proposed hybrid precoding & combining: the hybrid precoders \mathbf{F}_A and $\mathbf{F}_{D,0}$ are jointly designed with the hybrid combiners $\{\mathbf{W}_{A,k}\}$ and $\{\mathbf{W}_{D,k}\}$ by Algorithm 1.
- Corr.-based hybrid precoding & combining [37]: following the approach in [37], the analog precoder and combiners are jointly constructed by sequentially selecting the beamformers and combiners having the largest correlations with $\{\mathbf{H}_k\}$. The BD technique is used to determine the digital precoder and combiners.
- Random analog precoding & combining: the analog precoder and combiners, \mathbf{F}_A and $\{\mathbf{W}_{A,k}\}$, are defined as random phase shifting, and the digital precoder and combiners are designed by BD of effective channels. This scheme denotes the performance lower bound of the ZF-based hybrid precoding and combining.

Algorithm 2 provides the power allocation for rate balancing (RB) that maximizes the minimum user rate. For comparison, the power allocation for sum rate maximization (SRM) is also used based on the water-filling algorithm. Notice that we put the keywords *RB* and *SRM* in the legend to denote the power allocation methods. The keywords are dropped if no confusion arises between SRM and RB criteria.

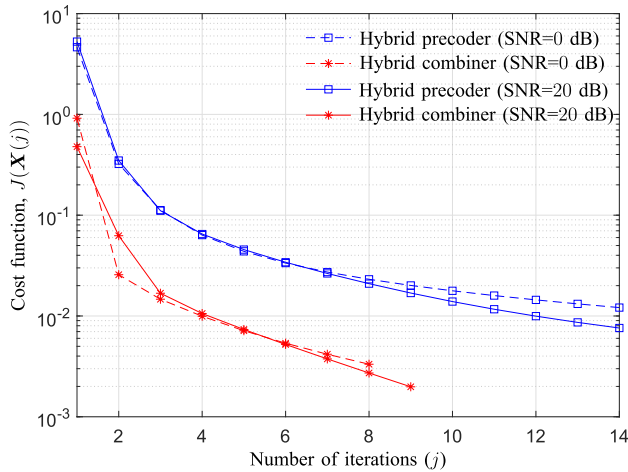


FIGURE 3. Convergence behavior of Algorithm 1 when $K = 4$, $M = 32$, and SNR=0 or 20 dB.

In the simulation, we used the following parameters: $N = 4$, $N_{RF} = 3$, $L = 2$, and $M_{RF} = KN_{RF}$ for receivers; $\epsilon_F = \epsilon_W = 0.001$ for Algorithm 1; and $\epsilon_p = 10^{-5}$, $\Delta p = 10^{-5}$, and $\beta_p = 0.5$ for Algorithm 2, unless otherwise stated. In (17) and (18), $\alpha(i)$ and $\beta(i)$ are adjusted by the backtracking line search [41, Ch 9.2] from the initial values $\alpha(0) = 1$ and $\beta(0) = 0.5$. For the correlation-based method, we use 128 independent phase shifting vectors designed by the phase quantization approach in [12] to determine the analog precoding matrix and the analog combining matrices, respectively. As in [12]–[24], the Saleh-Valenzuela channel model is used to generate mmwave MU-MIMO channels with the following parameters: the carrier frequency is 28 GHz; the number of clusters is 3; the number of subpaths per cluster is 8; the angle-of-departure (AoD) and angle-of-arrival (AoA) for each cluster are uniformly distributed from $-\pi$ to π in the azimuth direction and from -0.5π to 0.5π in the elevation direction, respectively; the subpath angular spread for azimuth and elevation directions is $\pi/64$ at the transmitter and $\pi/16$ at the receiver, respectively; and the inter-element spacing is equal to half wavelength at both the transmitter and receiver. The average channel gains are asymmetrically configured to reflect the distance variation between the transmitter and receiver, i.e. $E[\|\mathbf{H}_K\|_F^2] = 0.1E[\|\mathbf{H}_1\|_F^2]$, and $E[\|\mathbf{H}_k\|_F^2] = \zeta_k E[\|\mathbf{H}_1\|_F^2]$ for $2 \leq k \leq K - 1$ where ζ_k is a random variable uniformly distributed in the range of (0.1, 1.0). For simplicity, we set $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_K^2$, and SNR is defined as P_t/σ_k^2 . Every point denoting the achievable rate was obtained by averaging the simulation results over more than 200 independent channel realizations, except the convergence analysis in Figs. 3 and 4.

A. PERFORMANCE EVALUATION UNDER PERFECT CSI

This subsection presents the simulation results when the perfect CSI is available at the transmitter. Figs. 3 and 4 show the convergence behaviors of Algorithms 1 and 2, respectively, when $K = 4$, $M = 32$, and SNR = 0 or 20 dB. In Fig. 3,

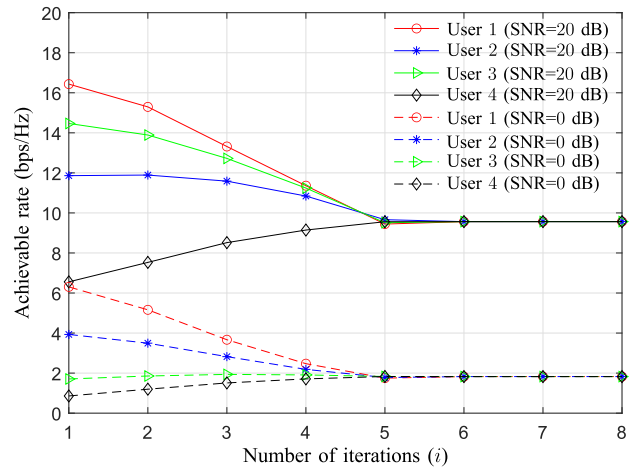


FIGURE 4. Convergence behavior of Algorithm 2 when $K = 4$, $M = 32$, and SNR=0 or 20 dB.

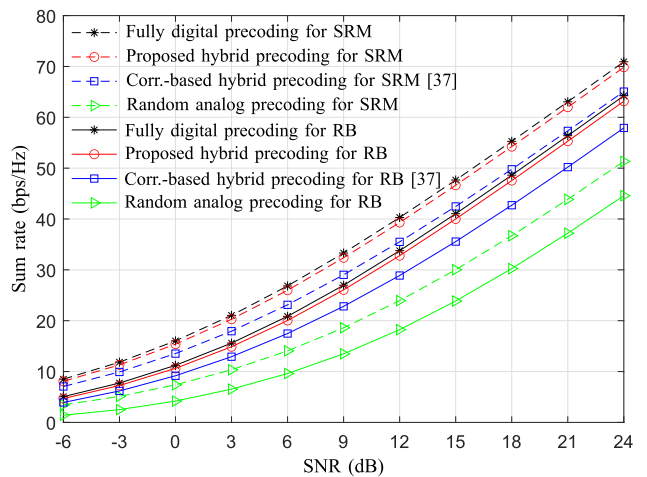


FIGURE 5. Sum rate of various hybrid precoding methods across SNR when $K = 4$ and $M = 32$.

F_A is a 32×12 matrix and $F_{D,0}$ is a 12×8 matrix, while $W_{A,k}$ is a 4×3 matrix and $W_{D,k}$ is a 3×2 matrix. Due to higher dimensions of factorized matrices, the hybrid precoder design requires more number of iterations than the design of hybrid combiners. Also, the convergence speed varies depending on the characteristics of the fully digital combiners $\{W_{FD,k}\}$ when designing the hybrid combiners for users, and Algorithm 1 converges slightly faster in SNR=20 dB than SNR=0 dB. Even in the worst case, the matrix factorization procedure in Algorithm 1 is completed within 20 iterations. In Fig. 4, whereas the user rate difference is very large before the power allocation ($i = 1$), the user rates rapidly converge to the same value through Algorithm 2 and the minimum user rate is maximized. The power allocation procedure requires only 5 ~ 8 iterations, and the required number of iterations is almost similar irrespective of initial achievable rates and SNR.

Figs. 5 and 6 compare the sum rate and the minimum user rate of various hybrid precoding methods, respectively, when

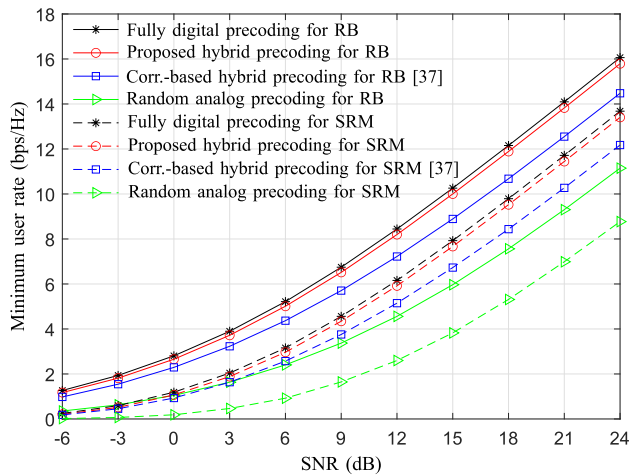


FIGURE 6. Minimum user rate of various hybrid precoding methods across SNR when $K = 4$ and $M = 32$.

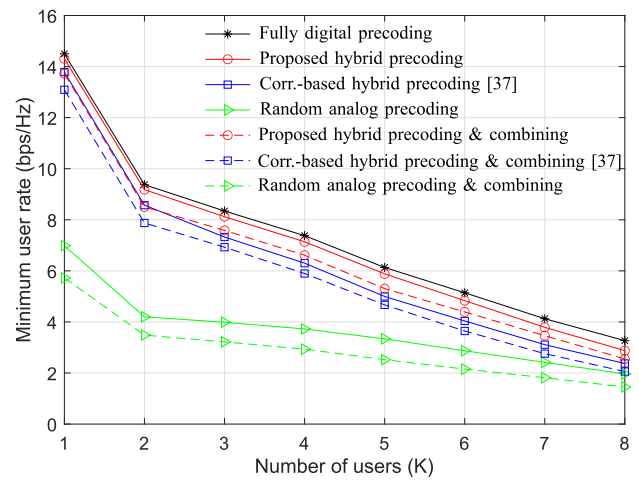


FIGURE 8. Minimum user rate across the number of users when $M = 32$ and $\text{SNR} = 10$ dB.

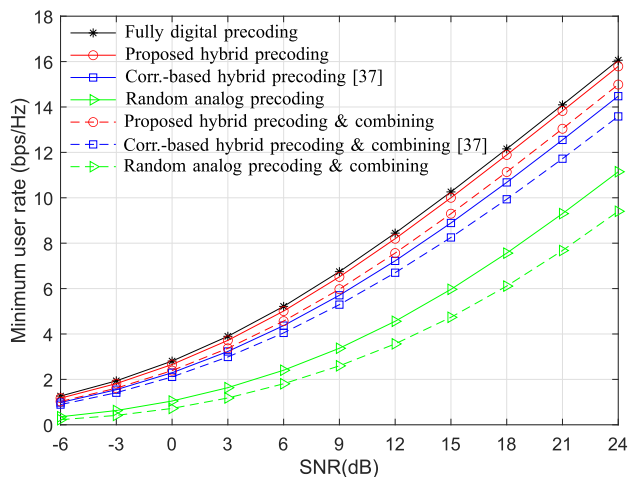


FIGURE 7. Minimum user rate according to SNR for various hybrid precoding and combining methods when $K = 4$ and $M = 32$.

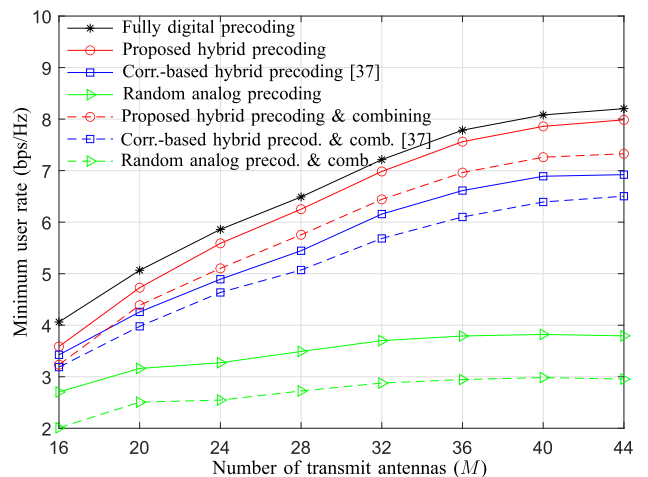


FIGURE 9. Minimum user rate across the number of transmit antennas (M) when $K = 4$ and $\text{SNR} = 10$ dB.

$K = 4$ and $M = 32$. For SRM, the power allocation is conducted by the water-filling algorithm across multiple data streams, thereby the sum rate of all users is maximized at the cost of the rate imbalance among users. For RB, the power allocation is performed by Algorithm 2, thus the minimum user rate is maximized with some sum rate loss. In Fig. 5, the proposed hybrid precoding for SRM outperforms the Corr.-based hybrid precoding for SRM and the random analog precoding for SRM, and presents the sum rate performance comparable to the fully digital precoding for SRM. Similarly, the proposed hybrid precoding for RB significantly improves the minimum user rate compared to the corr.-based hybrid precoding for RB and the random analog precoding for RB, and also performs very close to the fully digital precoding for RB which is the upper bound. As expected, the RB-based precoding methods achieve better minimum user rate than the SRM-based schemes.

In Figs. 7–9, we compare the minimum user rate of various transmission techniques, when the power allocation is con-

ducted by Algorithm 2 for RB. Fig. 7 shows the minimum user rate across SNR when $K = 4$ and $M = 32$. In addition, Figs. 8 and 9 present the change of the minimum user rate with increment of the number of users and the number of transmit antennas, respectively, when $\text{SNR} = 10$ dB, $M = 32$ for Fig. 8, and $K = 4$ for Fig. 9. For all cases, the proposed hybrid precoding scheme outperforms the corr.-based hybrid precoding and the random analog precoding. Also, the proposed hybrid precoding & combining method performs better than the existing counterparts such as the corr.-based hybrid precoding & combining and the random analog precoding & combining. The proposed hybrid precoding method exhibits reasonable performance loss compared to the fully digital precoding scheme. The performance gap between the proposed hybrid precoding and the proposed hybrid precoding & combining is higher than that between the fully digital precoding and the proposed hybrid precoding, because the nullity of $\tilde{\mathbf{H}}_k$ in (9) (or the rank of $\tilde{\mathbf{V}}_k^{(0)}$) is reduced in BD by the use of hybrid combining. In Fig. 8, the analog

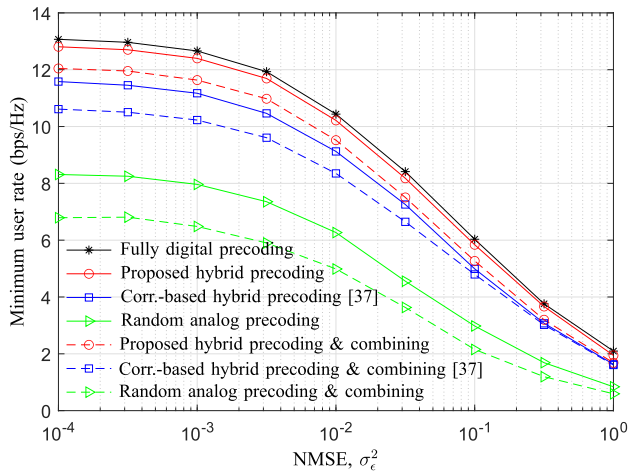


FIGURE 10. Minimum user rate according to NMSE of the channel when $K = 4$, $M = 32$, and $\text{SNR} = 20$ dB.

precoding gain decreases as the number of users increases, thus the rate loss of the random analog precoding is reduced with increment of K . In contrast, the performance loss of the random analog precoding grows as the number of transmit antennas M increases, because the analog precoding gain is proportional to M . In Fig. 9, the minimum user rates for all methods grow logarithmically as M increases, because a MIMO channel capacity is proportional to $\log_2(M)$ with a large M and fixed N .

B. PERFORMANCE EVALUATION UNDER CSI UNCERTAINTY

In this subsection, we evaluate the performance of the proposed method when both the transmitter and receiver design the hybrid precoder and hybrid combiners with some CSI errors, respectively. The CSI error is denoted as $\mathbf{E}_k \in \mathbb{C}^{N \times M}$ whose elements are i.i.d. Gaussian noises with zero mean, and the channel matrix with CSI uncertainty is expressed as

$$\tilde{\mathbf{H}}_k = \mathbf{H}_k + \mathbf{E}_k \tag{36}$$

where $k = 1, 2, \dots, K$. We define the normalized mean square error (NMSE) to describe the average power of the CSI error relative to the mean channel power as follows:

$$\sigma_\epsilon^2 = \frac{E[\|\mathbf{E}_k\|_F^2]}{E[\|\mathbf{H}_k\|_F^2]} \tag{37}$$

In practical systems, the precoder and combiners are designed using $\tilde{\mathbf{H}}_k$ instead of \mathbf{H}_k , resulting to imperfect cancellation of inter-user interferences in ZF-based transmission. Denote the hybrid precoders and combiners obtained from $\tilde{\mathbf{H}}_k$ as $\tilde{\mathbf{F}}_A, \tilde{\mathbf{F}}_D, \tilde{\mathbf{W}}_{A,k}$, and $\tilde{\mathbf{W}}_{D,k}$. In this case, the achievable rate for user k is computed as

$$f_k(\tilde{\mathbf{F}}_H, \{\tilde{\mathbf{W}}_{H,k}\}) = \log_2 |\mathbf{C}_{0,k}| - \log_2 |\mathbf{C}_{1,k}|, \tag{38}$$

where $\tilde{\mathbf{F}}_H = \tilde{\mathbf{F}}_A \tilde{\mathbf{F}}_D$ is the entire matrix for hybrid precoding, $\tilde{\mathbf{W}}_{H,k} = \tilde{\mathbf{W}}_{A,k} \tilde{\mathbf{W}}_{D,k}$ is the entire matrix for hybrid combin-

ing, and the matrices $\mathbf{C}_{0,k}$ and $\mathbf{C}_{1,k}$ are given by

$$\mathbf{C}_{0,k} = \tilde{\mathbf{W}}_{H,k}^H \mathbf{H}_k \tilde{\mathbf{F}}_H \tilde{\mathbf{F}}_H^H \mathbf{H}_k^H \tilde{\mathbf{W}}_{H,k} + \sigma_k^2 \tilde{\mathbf{W}}_{H,k}^H \tilde{\mathbf{W}}_{H,k} \tag{39a}$$

$$\mathbf{C}_{1,k} = \mathbf{C}_{0,k} - \tilde{\mathbf{W}}_{H,k}^H \mathbf{H}_k \tilde{\mathbf{F}}_{H,k} \tilde{\mathbf{F}}_{H,k}^H \mathbf{H}_k^H \tilde{\mathbf{W}}_{H,k}. \tag{39b}$$

Here, $\tilde{\mathbf{F}}_{H,k} \in \mathbb{C}^{M \times L}$ is the precoding matrix for the data streams transferred to user k that satisfies $\tilde{\mathbf{F}}_H = [\tilde{\mathbf{F}}_{H,1} \tilde{\mathbf{F}}_{H,2} \dots \tilde{\mathbf{F}}_{H,K}]$.

Fig. 10 denotes the minimum user rate according to the NMSE, σ_ϵ^2 , when $K = 4$, $M = 32$, and $\text{SNR} = 20$ dB. It is assumed that the NMSE is identical to all users. As the NMSE increases, the minimum user rate gradually decreases in all hybrid precoding methods. The proposed hybrid precoding scheme performs better than the corr.-based hybrid precoding and the random analog precoding methods, irrespective of the NMSE, and the performance difference is reduced with increment of the NMSE. As in the case of only hybrid precoding, the proposed approach achieves higher minimum user rate than the existing methods when the hybrid precoding and combining are jointly used. As before, the proposed method performs very close to the fully digital precoding in the entire NMSE region.

VI. CONCLUSION

In this paper, we have proposed a new design procedure for hybrid precoding and combining in mmWave MU-MIMO systems considering the rate balancing among users. In the proposed scheme, the analog precoder and combiners are determined by factorizing the fully digital precoder and combiners in the least squares sense, while the digital precoder and combiners are designed by BD of effective channels. The proposed ZF-based hybrid processing ensures interference-free transmission in the downlink under the perfect CSI. Moreover, the proposed power allocation algorithm maximizes the minimum user rate given the hybrid precoders and combiners. Numerical simulation results show that the proposed approach is more beneficial than the existing hybrid precoding and combining schemes in terms of the minimum user rate.

The proposed method is applicable to the transceiver design of cellular base stations and Wi-Fi APs equipped with large-scale transmit antenna elements operating in mmWave bands. Also, the proposed design scheme can be utilized in the uplink of mmWave MU-MIMO systems, by exploiting the duality between the downlink and uplink. It is a good future research topic to design hybrid precoders and combiners for rate balancing when the MMSE-based precoding and combining are used.

APPENDIX A PROOF OF PROPERTY 1

Given $p_k \geq 0$, the achievable rate for user k , $f_k(p_k)$, is obtained by the water-filling algorithm. When $\gamma_{k,\ell} > 0$ for all ℓ , $f_k(p_k)$ is a monotonically increasing function. When the number of users is two ($K = 2$), the constraints are given by $p_1 + p_2 = P_t$ and $p_1 \geq 0, p_2 \geq 0$. For the

variable p_2 , $f_1(p_1) = f_1(P_t - p_2)$ is a monotonically decreasing function while $f_2(p_2)$ is a monotonically increasing function. Therefore, from [40], $g_2(p_1, p_2) = \min(f_1(p_1), f_2(p_2))$ is a quasiconcave function, and $g_2(p_1, p_2)$ subject to $p_1 + p_2 = y$, $p_1 \geq 0$, and $p_2 \geq 0$ increases as y increases. In other words, $g_2(p_1, y - p_1)$ is a nondecreasing function w.r.t. y .

Similarly, $p_1 + p_2 + p_3 = P_t$ and $p_1 \geq 0, p_2 \geq 0, p_3 \geq 0$, when $K = 3$. For the variable p_3 , $g_2(p_1, P_t - p_1 - p_3)$ is a nonincreasing function while $f_3(p_3)$ is a monotonically increasing function. Again, from [40], $g_3(p_1, p_2, p_3) = \min(f_1(p_1), f_2(p_2), f_3(p_3)) = \min(g_2(p_1, p_2), f_3(p_3))$ is a quasiconcave function, and $g_3(p_1, p_2, y - (p_1 + p_2))$ is a nondecreasing function w.r.t. y .

By repeating this procedure, it is shown that $g_K(\mathbf{p}) = \min(f_1(p_1), f_2(p_2), \dots, f_K(p_K))$ subject to $p_1 + p_2 + \dots + p_K = P_t$ and $p_k \geq 0$ is a quasiconcave function. This completes the proof. ■

REFERENCES

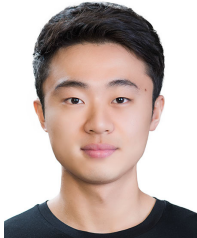
- [1] O. El Ayach, R. W. Heath, Jr., S. Abu-Surra, S. Rajagopal, and Z. Pi, "The capacity optimality of beam steering in large millimeter wave MIMO systems," in *Proc. IEEE 13rd Int. Workshop Signal Process. Adv. Wireless Commun. (SPAWC)*, Jun. 2012, pp. 100–104.
- [2] T. S. Rappaport, S. Sun, R. Mayzus, H. Zhao, Y. Azar, K. Wang, G. N. Wong, J. K. Schulz, M. Samimi, and F. Gutierrez, "Millimeter wave mobile communications for 5G cellular: It will work!" *IEEE Access*, vol. 1, pp. 335–349, 2013.
- [3] 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.
- [4] I. A. Hemadeh, K. Satyanarayana, M. El-Hajjar, and L. Hanzo, "Millimeter-wave communications: Physical channel models, design considerations, antenna constructions, and link-budget," *IEEE Commun. Surveys Tuts.*, vol. 20, no. 2, pp. 870–913, 2nd Quart., 2018.
- [5] M. R. Akdeniz, Y. Liu, M. K. Samimi, S. Sun, S. Rangan, T. S. Rappaport, and E. Erkip, "Millimeter wave channel modeling and cellular capacity evaluation," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1164–1179, Jun. 2014.
- [6] R. W. Heath, Jr., N. González-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, "An overview of signal processing techniques for millimeter wave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 436–453, Apr. 2016.
- [7] S. A. Busari, K. M. S. Huq, S. Mumtaz, L. Dai, and J. Rodriguez, "Millimeter-wave massive MIMO communication for future wireless systems: A survey," *IEEE Commun. Surveys Tuts.*, vol. 20, no. 2, pp. 836–869, 2nd Quart., 2018.
- [8] L. Lu, G. Y. Li, A. L. Swindlehurst, A. Ashikhmin, and R. Zhang, "An overview of massive MIMO: Benefits and challenges," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 742–758, Oct. 2014.
- [9] P. Xia, R. W. Heath, Jr., and N. Gonzalez-Prelcic, "Robust analog precoding designs for millimeter wave MIMO transceivers with frequency and time division duplexing," *IEEE Trans. Commun.*, vol. 64, no. 11, pp. 4622–4634, Nov. 2016.
- [10] J. Choi, "Analog beamforming for low-complexity multiuser detection in mm-wave systems," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 6747–6752, Aug. 2016.
- [11] Y. Wang, W. Zou, and Y. Tao, "Analog precoding designs for millimeter wave communication systems," *IEEE Trans. Veh. Technol.*, vol. 67, no. 12, pp. 11733–11745, Dec. 2018.
- [12] O. El Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. W. Heath, Jr., "Spatially sparse precoding in millimeter wave MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 13, no. 3, pp. 1499–1513, Jan. 2014.
- [13] F. Sotirani and W. Yu, "Hybrid digital and analog beamforming design for large-scale antenna arrays," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 501–513, Apr. 2016.
- [14] X. Yu, J.-C. Shen, J. Zhang, and K. B. Letaief, "Alternating minimization algorithms for hybrid precoding in millimeter wave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 485–500, Apr. 2016.
- [15] X. Qiao, Y. Zhang, M. Zhou, and L. Yang, "Alternating optimization based hybrid precoding strategies for millimeter wave MIMO systems," *IEEE Access*, vol. 8, pp. 113078–113089, 2020.
- [16] W. Ni, X. Dong, and W. S. Lu, "Near-optimal hybrid processing for massive MIMO systems via matrix decomposition," *IEEE Trans. Signal Process.*, vol. 65, no. 15, pp. 3922–3933, Aug. 2017.
- [17] J. Jin, Y. R. Zheng, W. Chen, and C. Xiao, "Hybrid precoding for millimeter wave MIMO systems: A matrix factorization approach," *IEEE Trans. Wireless Commun.*, vol. 17, no. 5, pp. 3327–3339, May 2018.
- [18] R. Méndez-Rial, C. Rusu, N. González-Prelcic, A. Alkhateeb, and R. W. Heath, Jr., "Hybrid MIMO architectures for millimeter wave communications: Phase shifters or switches?" *IEEE Access*, vol. 4, pp. 247–267, 2016.
- [19] Z. Wang, M. Li, Q. Liu, and A. L. Swindlehurst, "Hybrid precoder and combiner design with low-resolution phase shifters in mmWave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 2, pp. 256–269, May 2018.
- [20] F. Dong, W. Wang, and Z. Wei, "Low-complexity hybrid precoding for multi-user mmWave systems with low-resolution phase shifters," *IEEE Trans. Veh. Technol.*, vol. 68, no. 10, pp. 9774–9784, Oct. 2019.
- [21] A. N. Uwaechia, N. M. Mahyuddin, M. F. Ain, N. M. A. Latiff, and N. F. Za'bah, "On the spectral-efficiency of low-complexity and resolution hybrid precoding and combining transceivers for mmWave MIMO systems," *IEEE Access*, vol. 7, pp. 109259–109277, 2019.
- [22] A. W. Shaban, O. Damen, Y. Xin, and E. Au, "Statistically-aided codebook-based hybrid precoding for millimeter wave channels," *IEEE Access*, vol. 8, pp. 101500–101513, 2020.
- [23] S. He, J. Wang, Y. Huang, B. Ottersten, and W. Hong, "Codebook-based hybrid precoding for millimeter wave multiuser systems," *IEEE Trans. Signal Process.*, vol. 65, no. 20, pp. 5289–5304, Oct. 2017.
- [24] A. Alkhateeb, G. Leus, and R. W. Heath, Jr., "Limited feedback hybrid precoding for multi-user millimeter wave systems," *IEEE Trans. Wireless Commun.*, vol. 14, no. 11, pp. 6481–6494, Nov. 2015.
- [25] Q. H. Spencer, A. L. Swindlehurst, and M. Haardt, "Zero-forcing methods for downlink spatial multiplexing in multiuser MIMO channels," *IEEE Trans. Signal Process.*, vol. 52, no. 2, pp. 461–471, Feb. 2004.
- [26] R. Zhang, "Cooperative multi-cell block diagonalization with per-base-station power constraints," *IEEE J. Sel. Areas Commun.*, vol. 28, no. 9, pp. 1435–1445, Dec. 2010.
- [27] J. Choi, S. Han, and J. Joung, "Low-complexity multiuser MIMO precoder design under per-antenna power constraints," *IEEE Trans. Veh. Technol.*, vol. 67, no. 9, pp. 9011–9015, Sep. 2018.
- [28] J. Joung and Y. H. Lee, "Regularized channel diagonalization for multiuser MIMO downlink using a modified MMSE criterion," *IEEE Trans. Signal Process.*, vol. 55, no. 4, pp. 1573–1579, Apr. 2007.
- [29] V. Stankovic and M. Haardt, "Generalized design of multi-user MIMO precoding matrices," *IEEE Trans. Wireless Commun.*, vol. 7, no. 3, pp. 953–961, Mar. 2008.
- [30] S. Shi, M. Schubert, and H. Boche, "Downlink MMSE transceiver optimization for multiuser MIMO systems: MMSE balancing," *IEEE Trans. Signal Process.*, vol. 56, no. 8, pp. 3702–3712, Aug. 2008.
- [31] I. Ghamnia, D. Slock, and Y. Yuan-Wu, "Rate balancing for multiuser MIMO systems," in *Proc. IEEE 20th Int. Workshop Signal Process. Adv. Wireless Commun. (SPAWC)*, Jul. 2019, pp. 1–5.
- [32] A. Alkhateeb, R. W. Heath, Jr., and G. Leus, "Achievable rates of multi-user millimeter wave systems with hybrid precoding," in *Proc. IEEE Int. Conf. Commun.*, Jun. 2015, pp. 1232–1237.
- [33] W. Ni and X. Dong, "Hybrid block diagonalization for massive multiuser MIMO systems," *IEEE Trans. Commun.*, vol. 64, no. 1, pp. 201–211, Jan. 2016.
- [34] K. Duan, H. Du, and Z. Wu, "Hybrid alternating precoding and combining design for mmWave multi-user MIMO systems," in *Proc. IEEE/CIC Int. Conf. Commun. China (ICCC)*, Aug. 2018, pp. 217–221.
- [35] F. Khalid, "Hybrid beamforming for millimeter wave massive multiuser MIMO systems using regularized channel diagonalization," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3, pp. 705–708, Jun. 2019.
- [36] D. H. N. Nguyen, L. B. Le, and T. Le-Ngoc, "Hybrid MMSE precoding for mmWave multiuser MIMO systems," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2016, pp. 1–6.
- [37] D. H. Nguyen, L. B. Le, T. Le-Ngoc, and R. W. Heath, Jr., "Hybrid MMSE precoding and combining designs for mmWave multiuser systems," *IEEE Access*, vol. 5, pp. 19167–19181, 2017.

- [38] J. P. González-Coma, J. Rodríguez-Fernández, N. González-Prelcic, L. Castedo, and R. W. Heath, Jr., "Channel estimation and hybrid precoding for frequency selective multiuser mmWave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 2, pp. 353–367, May 2018.
- [39] J. P. Gonzalez-Coma, W. Utschick, and L. Castedo, "Hybrid LISA for wideband multiuser millimeter-wave communication systems under beam squint," *IEEE Trans. Wireless Commun.*, vol. 18, no. 2, pp. 1277–1288, Feb. 2019.
- [40] K. C. Kiwiel, "Convergence and efficiency of subgradient methods for quasiconvex minimization," *Math. Program.*, vol. 90, no. 1, pp. 1–25, Mar. 2001.
- [41] S. Boyd and L. Vandenberghe, *Convex Optimization*. New York, NY, USA: Cambridge Univ. Press, 2004.



WOO-HEE LIM received the B.S. degree from Korea Aerospace University (KAU), Goyang-si, South Korea, in 2021. He is currently pursuing the master's degree.

From 2020 to 2021, he was with the Intelligent Signal Processing Laboratory (ISPL), School of Electronics and Information Engineering, KAU, as an Undergraduate Research Assistant, where he conducted research on signal processing for MIMO and mmWave communications, and compressive sensing algorithms. Since 2021, he has been with the RF Circuits and Systems, Antennas/Application Laboratory, School of Electrical and Electronics Engineering, Chung-Ang University, Seoul, South Korea. His research interests include MIMO communication techniques and signal processing algorithms, mmWave antenna design, phased array antennas, smart antennas, antenna-in-package, and RFIC design.



SUWON JANG is currently pursuing the B.S. degree with KAU, Goyang-si, South Korea.

Since 2020, he has been with the Intelligent Signal Processing Laboratory (ISPL), School of Electronics and Information Engineering, KAU, as an Undergraduate Research Assistant, where he performed research on MIMO communications, mmWave transmission and signal processing algorithms, mmWave channel modeling and estimation, and intelligent reflecting surface (IRS) aided communications. His research interests include signal processing for IRS aided mmWave communications, unmanned aerial vehicle (UAV) communications, compressive sensing algorithms, and future cellular networks.



WOOHYEONG PARK is currently pursuing the B.S. degree with KAU, Goyang-si, South Korea.

In 2020, he joined the Intelligent Signal Processing Laboratory (ISPL), School of Electronics and Information Engineering, KAU, as an Undergraduate Research Assistant, where he performed research on signal processing for MIMO and IRS aided communications, compressive sensing algorithms, and transceiver design for mmWave communications. His research interests include mobile communication techniques, satellite communications, signal processing for IRS aided communication networks, and transceiver design for next generation cellular networks.



JIHOON CHOI (Senior Member, IEEE) received the B.S., M.S., and Ph.D. degrees from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 1997, 1999, and 2003, respectively.

From 2003 to 2004, he was with the Department of Electrical and Computer Engineering, The University of Texas at Austin, where he performed research on MIMO-OFDM systems as a Postdoctoral Fellow. From 2004 to 2008, he was with Samsung Electronics, South Korea, where he worked on developments of radio access stations for M-WiMAX and base stations for CDMA 1xEV-DO Rev.A/B. In 2008, he joined KAU, Goyang, South Korea, as a Faculty Member. He is currently a Professor with the School of Electronics and Information Engineering, KAU, where he is also the Chief Investigator of the Intelligent Signal Processing Laboratory (ISPL). His research interests include MIMO communications and signal processing algorithms, secure transmission in the physical layer, radar signal processing, UAV trajectory optimization, mobile edge computing, and modem design for future cellular networks, wireless LANs, the IoT devices, and digital broadcasting systems.

• • •