# Efficient Multi-User Hybrid Precoding and Combining for mmWave Massive MIMO Systems

Kaijie Zeng, Yuanli Ma, Bin Yan, Zheng Wang, Yili Xia

School of Information Science and Engineering

Southeast University, Nanjing, China

Email: kaijie_zeng@seu.edu.cn, myl@seu.edu.cn, bin_yan@seu.edu.cn, wznuaa@gmail.com, yili_xia@seu.edu.cn

*Abstract*—**For downlink millimeter wave (mmWave) massive multiple-input multiple-output (MIMO) systems, this paper proposes an efficient multi-user hybrid precoding and combining (EMU-HPC) scheme to maximize the sum-rate with low feedback overhead. Specifically, we adopt a non-iterative generalized low rank approximation of matrices (NI-GLRAM) framework in the analog stage to harvest the array gain with minimal feedback overhead. Meanwhile, for each mobile station (MS), a low-complexity truncated singular value decomposition (TSVD) is implemented for efficient analog combining. To further enhance the sum-rate, we propose a joint analog-digital interference suppression strategy, in which the orthogonal matrix projection is applied in analog precoding, followed by the regularized channel diagonalization (RCD) in the digital stage. Numerical results confirm that the proposed EMU-HPC method achieves superior sum-rate performance with low feedback overhead and computational complexity.**

*Index Terms*—**Massive MIMO, multi-user MIMO, hybrid precoding, sum-rate maximization.**

## I. INTRODUCTION

In 5G systems and beyond, the integration of mmWave and massive MIMO compensates for severe path loss and enables high sum-rate via precoding [1]. However, fully digital precoding is impractical due to its prohibitive hardware cost and power consumption from radio frequency (RF) chains [2]. To address this issue, hybrid precoding has emerged as a solution, which connects a digital baseband precoder to an analog RF precoder with far fewer RF chains [3]. In multi-user MIMO (MU-MIMO), the problem of hybrid precoding and combining design is normally divided into analog and digital stages for sum-rate optimization [4]–[7]. Yet achieving higher sum-rate without increasing feedback overhead remains challenging.

In particular, the work in [4] employs exhaustive codebook search and equal gain transmission in the analog stage to obtain array gains, followed by block diagonalization in the digital stage for interference elimination. With the same digital design, [5] exploits low rank approximations to enhance array gains, while [6] maximizes the capacity of baseband channels. Although these methods reduce feedback overhead by separately designing analog precoder and combiners, the interference is managed solely in the digital stage, which limits the sum-rate. To solve this problem, [7] incorporates interference management in the analog stage using orthogonal

vector projections and employs alternating optimization (AO) for array gains, but the increased feedback overhead from successive AOs may hinder its deployment.

To achieve high sum-rate with low feedback overhead, an efficient multi-user hybrid precoding and combining scheme, named EMU-HPC, is proposed for mmWave massive MIMO systems. Specifically, to maximize the array gain, we employ the NI-GLRAM framework in the analog stage with low feedback cost, while the TSVD in analog combining is realized efficiently via the variance-reduced principal component analysis (VR-PCA) in the LazySVD framework. For sum-rate improvement, we propose a joint analog-digital interference management scheme, with the orthogonal matrix projection first applied in analog precoding, followed by RCD in the digital stage. Numerical results demonstrate that EMU-HPC attains superior sum-rate with minimal feedback overhead and reduced computational complexity.

## II. SYSTEM MODEL

We consider a single-cell massive MU-MIMO downlink system. The base station (BS) equipped with $N_t$ transmit antennas and $M_t$ RF chains serves $K$ MSs, and each MS has $N_r$ receive antennas and $M_r$ RF chains to multiplex $N_s$ data streams. For effective communications, the constraints $N_s = M_r \ll N_r$ and $KN_s = M_t \ll N_t$ are assumed. The narrowband mmWave channel for the $k$-th MS (MS-$k$) is modeled by the Saleh–Valenzuela model [4] as

$$\mathbf{H}_k = \sqrt{\frac{\beta_k N_t N_r}{N_c N_{ray}}} \sum_{i=1}^{N_c} \sum_{l=1}^{N_{ray}} \alpha_{il}^k \mathbf{a}_r(\phi_{il}^k) \mathbf{a}_t^H(\theta_{il}^k), \qquad (1)$$

which consists of $N_c$ clusters, each propagating $N_{ray}$ paths, with $\beta_k$ denoting the large-scale fading factor. Uniform linear arrays (ULAs) are employed at the BS and all MSs. For the $(i,l)$-th path in $\mathbf{H}_k$, $\alpha_{il}^k \sim \mathcal{CN}(0,1)$ is the path gain, $\phi_{il}^k$ and $\theta_{il}^k$ are azimuth angles of arrival (AoA) and departure (AoD). $\mathbf{a}_r(\phi_{il}^k)$ and $\mathbf{a}_t(\theta_{il}^k)$ represent the receive and transmit array response vectors, respectively, of the form $\mathbf{a}(\psi) = \frac{1}{\sqrt{N}}[1, e^{j\frac{2\pi}{\lambda} d \sin(\psi)}, \ldots, e^{j(N-1)\frac{2\pi}{\lambda} d \sin(\psi)}]^T$, where $N$ is the number of antennas and $\psi$ can be the AoA or AoD, $\lambda$ is the signal wavelength and $d = \frac{\lambda}{2}$ is the antenna spacing without loss of generality.

At the BS, the aggregated data streams $\mathbf{s} = [\mathbf{s}_1^T, \ldots, \mathbf{s}_K^T]^T$, where $\mathbf{s}_k \in \mathbb{C}^{N_s \times 1}$ contains the data streams for MS-$k$,

is processed by a digital precoder $\mathbf{F}_D \in \mathbb{C}^{M_t \times KN_s}$ and an analog precoder $\mathbf{F}_A \in \mathbb{C}^{N_t \times M_t}$. Thus, the signal transmitted to the mmWave channels is written as $\mathbf{x} = \mathbf{F}_A\mathbf{F}_D\mathbf{s}$. The total power budget $P_t$ is allocated as $\mathbb{E}\{\mathbf{s}\mathbf{s}^H\} = \frac{P_t}{KN_s}\mathbf{I}_{KN_s}$, where $\mathbb{E}\{\cdot\}$ denotes expectation. Thus, to meet the transmit power constraint $\mathbb{E}\{\|\mathbf{x}\|_2^2\} = P_t$, the hybrid precoder should be normalized such that $\|\mathbf{F}_A\mathbf{F}_D\|_F^2 = KN_s$.

At MS-$k$, the received signal is processed by an analog combiner $\mathbf{W}_{A,k} \in \mathbb{C}^{N_r \times M_r}$ and a digital combiner $\mathbf{W}_{D,k} \in \mathbb{C}^{M_r \times N_s}$. In particular, since the analog precoder and combiner are implemented by phase shifters, their entries are subject to constant modulus (CM) constraints, i.e., $|\mathbf{F}_A(i,j)| = \frac{1}{\sqrt{N_t}}, \forall i,j$ and $|\mathbf{W}_{A,k}(i,j)| = \frac{1}{\sqrt{N_r}}, \forall i,j,k$. After the analog precoding and combining, the baseband channel for MS-$k$ is defined as $\widetilde{\mathbf{H}}_k = \mathbf{W}_{A,k}^H\mathbf{H}_k\mathbf{F}_A$.

After digital combining, the signal-to-interference-plus-noise ratio (SINR) of the $i$-th data stream in MS-$k$ is

$$\text{SINR}_{k_i} = \frac{S_{k_i}}{I_{k_i} + N_{k_i}}, \begin{cases} S_{k_i} = P_{k_i}|\mathbf{w}_{k_i}^H\widetilde{\mathbf{H}}_k\mathbf{f}_{k_i}|^2, \\ I_{k_i} = \sum_{\substack{l=1 \\ l \neq i}}^{N_s} P_{k_l}|\mathbf{w}_{k_i}^H\widetilde{\mathbf{H}}_k\mathbf{f}_{k_l}|^2 \\ \quad + \sum_{\substack{m=1 \\ m \neq k}}^{K}\sum_{j=1}^{N_s} P_{m_j}|\mathbf{w}_{k_i}^H\widetilde{\mathbf{H}}_k\mathbf{f}_{m_j}|^2, \\ N_{k_i} = \sigma_n^2\|\mathbf{w}_{k_i}^H\mathbf{W}_{A,k}^H\|_2^2, \end{cases}$$
(2)

where $P_{k_i}$ is the power allocated to the $i$-th data stream of $\mathbf{s}_k$. The first term of $I_{k_i}$ is the intra-user interference and the second term is the multi-user interference (MUI). $\sigma_n^2$ is the power of additive white Gaussian noise, $\mathbf{w}_{k_i}$ is the $i$-th column of $\mathbf{W}_{D,k}$, and $\mathbf{f}_{k_i}$ is the $[(k-1)N_s + i]$-th column of $\mathbf{F}_D$.

Assuming Gaussian data symbols, the design problem of hybrid precoding and combining to maximize the sum-rate is formulated as

$$\max_{\substack{\mathbf{F}_A, \mathbf{F}_D, \\ \{\mathbf{W}_{A,k}, \mathbf{W}_{D,k}\}_{k=1}^K}} R = \sum_{k=1}^{K}\sum_{i=1}^{N_s} \mathbb{E}\{\log_2(1 + \text{SINR}_{k_i})\}, \quad (3a)$$

$$\text{s.t.} \quad \|\mathbf{F}_A\mathbf{F}_D\|_F^2 = KN_s, \quad (3b)$$

$$\mathbf{F}_A \in \mathcal{F}_A, \quad (3c)$$

$$\mathbf{W}_{A,k} \in \mathcal{W}_{A,k}, \ \forall k, \quad (3d)$$

where $\mathcal{F}_A$ and $\mathcal{W}_{A,k}, \forall k$, are the feasible sets for the analog precoder and combiners, respectively, satisfying the CM constraints. Unfortunately, the optimization problem (3) is non-convex, rendering its global optimum intractable. Consequently, existing methods [5], [7] adopt the two-stage design pattern as follows:

   a) **Analog Stage Design**: The analog precoder $\mathbf{F}_A$ and combiners $\{\mathbf{W}_{A,k}\}_{k=1}^K$ are obtained at this stage such that the aggregated baseband channel $\widetilde{\mathbf{H}} = [\widetilde{\mathbf{H}}_1^T, \dots, \widetilde{\mathbf{H}}_K^T]^T$ is a diagonal low-rank approximation of $\mathbf{H} = [\mathbf{H}_1^T, \dots, \mathbf{H}_K^T]^T$. In this way, array gains are maximized and MUI is sufficiently suppressed. Assuming $\{\mathbf{W}_{A,k}\}_{k=1}^K$ and $\mathbf{F}_A$ to be semi-unitary, the analog stage design problem is given as [5]

$$\max_{\mathbf{F}_A, \{\mathbf{W}_{A,k}\}_{k=1}^K} \sum_{k=1}^{K} \|\mathbf{W}_{A,k}^H\mathbf{H}_k\mathbf{F}_A\|_F^2, \quad (4a)$$

$$\text{s.t.} \quad \mathbf{F}_A^H\mathbf{F}_A = \mathbf{I}_{M_t}, \ \mathbf{W}_{A,k}^H\mathbf{W}_{A,k} = \mathbf{I}_{M_r}, \ \forall k, \quad (4b)$$

$$\mathbf{F}_A \in \mathcal{F}_A, \ \mathbf{W}_{A,k} \in \mathcal{W}_{A,k}, \ \forall k. \quad (4c)$$

After relaxing the CM constraints, non-convex problem (4) can be approximately solved by the generalized low rank approximations of matrices [5], which involves AO between $\mathbf{F}_A$ and $\{\mathbf{W}_{A,k}\}_{k=1}^K$, incurring prohibitive feedback overhead.

   b) **Digital Stage Design**: Based on the baseband channel $\{\widetilde{\mathbf{H}}_k\}_{k=1}^K$ formed by $\mathbf{F}_A$ and $\{\mathbf{W}_{A,k}\}_{k=1}^K$, the digital precoder $\mathbf{F}_D$ and combiners $\{\mathbf{W}_{D,k}\}_{k=1}^K$ are derived for interference mitigation.

## III. HYBRID PRECODING AND COMBINING DESIGN

In this section, we elaborate and analyze the proposed EMU-HPC method. In the analog stage, to maximize array gains without AO, we adopt the NI-GLRAM framework [8], which designs the combiners and precoder non-iteratively. Moreover, the analog precoding achieves a balance between array gains and MUI suppression. In the digital stage, RCD [9] is applied to mitigate the residual interference and noise, followed by an analysis of complexity and feedback overhead.

### A. Analog Combining Design

The analog combiner for MS-$k$ is designed to maximize the array gain. Specifically, following NI-GLRAM, we decouple the analog combining design from (4) as

$$\max_{\{\mathbf{W}_{A,k}\}_{k=1}^K} \sum_{k=1}^{K} \|\mathbf{W}_{A,k}^H\mathbf{H}_k\|_F^2, \quad (5a)$$

$$\text{s.t.} \quad \mathbf{W}_{A,k}^H\mathbf{W}_{A,k} = \mathbf{I}_{M_r}, \mathbf{W}_{A,k} \in \mathcal{W}_{A,k}, \ \forall k. \quad (5b)$$

Since (5) remains non-convex, we relax CM constraints and obtain a semi-unitary solution $\widehat{\mathbf{W}}_{A,k}$ for MS-$k, \forall k$, given by the first $M_r$ left singular vectors of $\mathbf{H}_k$, i.e., $\widehat{\mathbf{W}}_{A,k} = \mathbf{U}_k(:, 1:M_r)$, with the singular value decomposition (SVD) $\mathbf{H}_k = \mathbf{U}_k\mathbf{D}_k\mathbf{V}_k^H$ [8]. However, direct TSVD at each MS may be prohibitive under limited computational resources.

To reduce this computational burden, we implement the VR-PCA method [10] within the LazySVD framework [11] to approximate the required TSVD at each MS, which is based on stochastic gradient ascent with guaranteed convergence [10]. Specifically, to estimate the $m$-th left singular vector of $\mathbf{H}_k$ via VR-PCA, the $t$-th iteration within epoch $s$ updates as follows:

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \gamma[\mathbf{h}_{i_t}^{(m)}(\mathbf{h}_{i_t}^{(m)H}\mathbf{w}_{t-1} - \mathbf{h}_{i_t}^{(m)H}\widetilde{\mathbf{w}}_{s-1}) + \boldsymbol{\mu}_s], \quad (6)$$

where $\gamma$ is the step size, $\mathbf{w}_t$ and $\widetilde{\mathbf{w}}_s$ are the per-iteration and per-epoch estimates with unit $\ell_2$ norm. $\boldsymbol{\mu}_s$ is the gradient at $\widetilde{\mathbf{w}}_{s-1}$, and $\mathbf{h}_{i_t}^{(m)}$ is the $i_t$-th column of the deflated channel matrix $\mathbf{H}_k^{(m)}$ with $i_t$ randomly selected. The maximum number of epochs is $N_{\text{epoch}}$, and the number of iterations per epoch is $N_t$ [10]. VR-PCA converges if $\|\mathbf{w}_{N_t} - \widetilde{\mathbf{w}}_{s-1}\|_2 \leq \varepsilon$ is satisfied at the end of an iteration, where $\varepsilon$ is the error tolerance.

In the LazySVD framework, the $m$-th left singular vector $\mathbf{u}_m'$ obtained from VR-PCA is orthonormalized with all previous vectors $\mathbf{U}^{(m-1)} = [\mathbf{u}_1, \dots, \mathbf{u}_{m-1}]$ as

$$\mathbf{u}_m = \frac{(\mathbf{I}_{N_r} - \mathbf{U}^{(m-1)}\mathbf{U}^{(m-1)H})\mathbf{u}_m'}{\left\|(\mathbf{I}_{N_r} - \mathbf{U}^{(m-1)}\mathbf{U}^{(m-1)H})\mathbf{u}_m'\right\|_2}. \quad (7)$$

**Algorithm 1** Analog Combining Scheme for MS-$k$

---

**Input:** $\mathbf{H}_k$, $N_{\text{epoch}}$, $\gamma$, $\varepsilon$

**Output:** $\mathbf{W}_{\text{A},k}$

1: Initialize $\mathbf{U}^{(0)} = \emptyset$, $\mathbf{H}_k^{(1)} = \mathbf{H}_k$
2: **for** $m = 1$ to $M_{\text{r}}$ **do**
3:     Initialize $\widetilde{\mathbf{w}}_0 \sim \mathcal{CN}(\mathbf{0}_{N_{\text{r}} \times 1}, \mathbf{I}_{N_{\text{r}}})$, $\widetilde{\mathbf{w}}_0 \leftarrow \frac{\widetilde{\mathbf{w}}_0}{\|\widetilde{\mathbf{w}}_0\|_2}$
4:     Initialize the error vector $\mathbf{e} = \widetilde{\mathbf{w}}_0$
5:     s = 1
6:     **while** $s \leq N_{\text{epoch}}$ and $\|\mathbf{e}\|_2 > \varepsilon$ **do**
7:         $\boldsymbol{\mu}_s = \frac{1}{N_{\text{t}}} \sum_{i=1}^{N_{\text{t}}} \mathbf{h}_i^{(m)} \left( \mathbf{h}_i^{(m)H} \widetilde{\mathbf{w}}_{s-1} \right)$
8:         $\mathbf{w}_0 = \widetilde{\mathbf{w}}_{s-1}$
9:         **for** $t = 1$ to $N_{\text{t}}$ **do**
10:             Sample $i_t$ randomly from $\{1, 2, \ldots, N_{\text{t}}\}$
11:             Get $\mathbf{h}_{i_t}^{(m)}$ and update $\mathbf{w}_t$ by (6)
12:             Normalize $\mathbf{w}_t$ as $\mathbf{w}_t \leftarrow \frac{\mathbf{w}_t}{\|\mathbf{w}_t\|_2}$
13:         **end for**
14:         Obtain the error vector as $\mathbf{e} = \mathbf{w}_{N_{\text{t}}} - \widetilde{\mathbf{w}}_{s-1}$
15:         $\widetilde{\mathbf{w}}_s = \mathbf{w}_{N_{\text{t}}}$
16:         $s \leftarrow s + 1$
17:     **end while**
18:     Obtain $\mathbf{u}_m^{'}$ as $\mathbf{u}_m^{'} = \widetilde{\mathbf{w}}_{s-1}$ and obtain $\mathbf{u}_m$ by (7)
19:     Obtain $\mathbf{H}_k^{(m+1)}$ by (8)
20:     Enlarge $\mathbf{U}^{(m)}$ by $\mathbf{U}^{(m)} = [\mathbf{U}^{(m-1)}, \mathbf{u}_m]$
21: **end for**
22: Given $\widehat{\mathbf{W}}_{\text{A},k} = \mathbf{U}^{(M_{\text{r}})}$, obtain $\mathbf{W}_{\text{A},k}$ by (9)

---

Based on this, the deflated channel matrix $\mathbf{H}_k^{(m)}$ is updated by projecting it onto the orthogonal complement of $\mathbf{u}_m$ as

$$\mathbf{H}_k^{(m+1)} = (\mathbf{I}_{N_{\text{r}}} - \mathbf{u}_m \mathbf{u}_m^H) \mathbf{H}_k^{(m)}, \tag{8}$$

which is the input for the next round of VR-PCA. The resulting $M_{\text{r}}$ vectors form the semi-unitary analog combiner $\widehat{\mathbf{W}}_{\text{A},k}$. Under the CM constraint, the analog combiner $\mathbf{W}_{\text{A},k}$ approximates $\widehat{\mathbf{W}}_{\text{A},k}$ via phase extraction, given as

$$\mathbf{W}_{\text{A},k} = \frac{1}{\sqrt{N_{\text{r}}}} \exp \left[ \text{j} \angle \left( \widehat{\mathbf{W}}_{\text{A},k} \right) \right], \tag{9}$$

where $\angle(\cdot)$ represents the element-wise phase extraction for a matrix. After that, MS-$k$ computes $\bar{\mathbf{H}}_k = \mathbf{W}_{\text{A},k}^H \mathbf{H}_k$ and feeds it back to the BS for analog precoding. The analog combining scheme for MS-$k$ is summarized in Algorithm 1.

### B. Analog Precoding Design

The analog precoding is designed to balance array gain enhancement and MUI suppression based on $\{\bar{\mathbf{H}}_k\}_{k=1}^K$ from all MSs. Following NI-GLRAM, the analog precoding design based on (4) is given as

$$\max_{\mathbf{F}_{\text{A}}} \sum_{k=1}^K \|\bar{\mathbf{H}}_k \mathbf{F}_{\text{A}}\|_F^2, \tag{10a}$$

$$\text{s.t.} \quad \mathbf{F}_{\text{A}}^H \mathbf{F}_{\text{A}} = \mathbf{I}_{M_{\text{t}}}, \; \mathbf{F}_{\text{A}} \in \mathcal{F}_{\text{A}}. \tag{10b}$$

However, for each MS, (10) fails to distinguish its desired array gain from the MUI it causes to the other MSs, resulting

in maximizing both. This leads to an imbalance between array gain harvesting and MUI suppression. To address this issue, we design a dedicated analog precoder $\mathbf{F}_{\text{A},k} \in \mathbb{C}^{N_{\text{t}} \times M_{\text{r}}}$ for each MS-$k$ and form the overall analog precoder as $\mathbf{F}_{\text{A}} = [\mathbf{F}_{\text{A},1}, \ldots, \mathbf{F}_{\text{A},K}]$. Based on this, the desired array gain for MS-$k$ is $\|\bar{\mathbf{H}}_k \mathbf{F}_{\text{A},k}\|_F^2$, while the MUI from MS-$k$ to the other MSs is $\sum_{p \neq k} \|\bar{\mathbf{H}}_p \mathbf{F}_{\text{A},k}\|_F^2$. However, jointly maximizing the former and minimizing the latter is difficult since the row space of $\bar{\mathbf{H}}_k$ misaligns with the common null space of $\bar{\mathbf{H}}_p$'s, $\forall p \neq k$.

To resolve this, we relax the CM constraint and design the column space of semi-unitary analog precoder $\widehat{\mathbf{F}}_{\text{A},k}$ to align with the row space of $\bar{\mathbf{H}}_k$ after projecting $\bar{\mathbf{H}}_k$ onto the joint orthogonal complement of $\widehat{\mathbf{F}}_{\text{A},1}, \ldots, \widehat{\mathbf{F}}_{\text{A},k-1}$, resulting in the residual channel matrix $\bar{\mathbf{H}}_{\text{res},k}$, given by

$$\bar{\mathbf{H}}_{\text{res},k} = \begin{cases} \bar{\mathbf{H}}_1, & k = 1 \\ \bar{\mathbf{H}}_k \prod_{i=1}^{k-1} \left( \mathbf{I}_{N_{\text{t}}} - \widehat{\mathbf{F}}_{\text{A},i} \widehat{\mathbf{F}}_{\text{A},i}^H \right), & k > 1 \end{cases} \tag{11}$$

if the previously obtained $\widehat{\mathbf{F}}_{\text{A},i}$, $\forall i < k$, are semi-unitary for $k = 2$, or mutually orthogonal for $k \geq 3$. Then, we modify the analog precoding problem for MS-$k$ from (10) to a MUI-aware array gain maximization problem as follows:

$$\max_{\mathbf{F}_{\text{A},k}} \|\bar{\mathbf{H}}_{\text{res},k} \mathbf{F}_{\text{A},k}\|_F^2, \tag{12a}$$

$$\text{s.t.} \quad \mathbf{F}_{\text{A},k}^H \mathbf{F}_{\text{A},k} = \mathbf{I}_{M_{\text{r}}}, \; \mathbf{F}_{\text{A},k} \in \mathcal{F}_{\text{A},k}, \tag{12b}$$

where $\mathcal{F}_{\text{A},k}$ is the feasible set for $\mathbf{F}_{\text{A},k}$. Its semi-unitary solution is $\widehat{\mathbf{F}}_{\text{A},k} = \bar{\mathbf{V}}_k(:, 1 : M_{\text{r}})$ obtained via SVD $\bar{\mathbf{H}}_{\text{res},k} = \bar{\mathbf{U}}_k \bar{\mathbf{D}}_k \bar{\mathbf{V}}_k^H$. Thus, by designing MS-1 to MS-$K$ recursively, each $\widehat{\mathbf{F}}_{\text{A},k}$ satisfies the semi-unitary condition for $k = 1$ and the mutual orthogonality condition for $k \geq 2$, i.e., $\widehat{\mathbf{F}}_{\text{A},p}^H \widehat{\mathbf{F}}_{\text{A},q} = \mathbf{0}_{M_{\text{r}} \times M_{\text{r}}}, \forall p, q \leq k, p \neq q$, thereby validating the computation of $\bar{\mathbf{H}}_{\text{res},k+1}$ via (11).

The elements of $\widehat{\mathbf{F}}_{\text{A},k}$ are then phase-extracted to meet the CM constraint. Therefore, $\mathbf{F}_{\text{A},k}$ is given by

$$\mathbf{F}_{\text{A},k} = \frac{1}{\sqrt{N_{\text{t}}}} \exp \left[ \text{j} \angle \left( \widehat{\mathbf{F}}_{\text{A},k} \right) \right], \tag{13}$$

with the following MUI suppression property.

**Proposition 1.** *The analog precoder in* (13) *minimizes the upper bound of the prior MUI term* $\|\bar{\mathbf{H}}_p \mathbf{F}_{\text{A},k}\|_F^2$, $\forall p < k$.

*Proof.* From the mutual orthogonality condition and (11), it follows that the approximate prior MUI term satisfies $\bar{\mathbf{H}}_p \widehat{\mathbf{F}}_{\text{A},k} = \mathbf{0}_{M_{\text{r}} \times M_{\text{r}}}, \forall p < k$. Leveraging this property, the actual prior MUI term is upper-bounded as

$$\begin{aligned} \left\| \bar{\mathbf{H}}_p \mathbf{F}_{\text{A},k} \right\|_F^2 &= \left\| \bar{\mathbf{H}}_p \mathbf{F}_{\text{A},k} - \bar{\mathbf{H}}_p \widehat{\mathbf{F}}_{\text{A},k} \right\|_F^2 \\ &\leq \left\| \bar{\mathbf{H}}_p \right\|_2^2 \left\| \widehat{\mathbf{F}}_{\text{A},k} - \mathbf{F}_{\text{A},k} \right\|_F^2. \end{aligned} \tag{14}$$

Since $\|\bar{\mathbf{H}}_p\|_2$ is constant, minimizing the upper bound is equivalent to minimizing $\|\widehat{\mathbf{F}}_{\text{A},k} - \mathbf{F}_{\text{A},k}\|_F^2$ to obtain $\mathbf{F}_{\text{A},k}$ under the CM constraint. As shown in [12], the optimal analog precoder achieving this minimization is given by (13), concluding the proof. $\square$

In summary, for each MS, the MUI generated to prior MSs is effectively suppressed, while the MUI to later MSs persists due to the balance with array gain harvesting.

## C. Digital Precoding and Combining Design

In the digital stage, RCD is applied to the baseband channel $\widetilde{\mathbf{H}}_k, \forall k$, to suppress the remaining interference and noise.

The RCD method first addresses the MUI plus noise under minimum mean square error criterion, where the digital precoder for MS-$k$ is $\mathbf{T}_k^{(a)}$. We define $\bar{\bar{\mathbf{H}}}_k = [\widetilde{\mathbf{H}}_1^T, \ldots, \widetilde{\mathbf{H}}_{k-1}^T, \widetilde{\mathbf{H}}_{k+1}^T, \ldots, \widetilde{\mathbf{H}}_K^T]^T$, so that $\|\bar{\bar{\mathbf{H}}}_k \mathbf{T}_k^{(a)}\|_F^2$ quantifies the MUI generated by MS-$k$ to the other MSs, and the design problem of $\{\mathbf{T}_k^{(a)}\}_{k=1}^K$ is given by [7]

$$\min_{\{\mathbf{T}_k^{(a)}\}_{k=1}^K} \sum_{k=1}^K \left( \left\|\bar{\bar{\mathbf{H}}}_k \mathbf{T}_k^{(a)}\right\|_F^2 + \frac{KN_s \sigma_n^2}{P_t} \|\mathbf{T}_k^{(a)}\|_F^2 \right). \quad (15)$$

One of its non-trivial solutions is given by

$$\mathbf{T}_k^{(a)} = \left( \bar{\bar{\mathbf{H}}}_k^H \bar{\bar{\mathbf{H}}}_k + \frac{KN_s \sigma_n^2}{P_t} \mathbf{I}_{M_t} \right)^{-1}, \quad (16)$$

which reduces the MUI to negligible levels. This enables each MS to optimize the data stream transmission in its effective subchannel $\widetilde{\mathbf{H}}_k \mathbf{T}_k^{(a)}$, with SVD $\widetilde{\mathbf{H}}_k \mathbf{T}_k^{(a)} = \mathbf{U}_{b,k} \mathbf{D}_{b,k} \mathbf{V}_{b,k}^H$. For the $k$-th subchannel, its optimal digital precoder $\mathbf{T}_k^{(b)}$ and the digital combiner $\mathbf{W}_{D,k}$, which jointly eliminate the intra-user interference are designed as

$$\mathbf{T}_k^{(b)} = \mathbf{V}_{b,k}(:, 1 : N_s), \quad \mathbf{W}_{D,k} = \mathbf{U}_{b,k}(:, 1 : N_s). \quad (17)$$

Finally, under the total transmit power constraint, the normalized digital precoder for MS-$k$ is given as

$$\mathbf{F}_{D,k} = \sqrt{\frac{KN_s}{\sum_{i=1}^K \|\mathbf{F}_A \mathbf{T}_i^{(a)} \mathbf{T}_i^{(b)}\|_F^2}} \mathbf{T}_k^{(a)} \mathbf{T}_k^{(b)}, \quad (18)$$

and the overall digital precoder is $\mathbf{F}_D = [\mathbf{F}_{D,1}, \ldots, \mathbf{F}_{D,K}]$.

In summary, the proposed EMU-HPC method is shown in Algorithm 2 and its computational complexity is further examined. Specifically, the analog combining has a complexity order of $O(KN_{epoch}N_sN_tN_r)$, dominated by the computation of the $\boldsymbol{\mu}_s$ and $\mathbf{w}_t$ in VR-PCAs. In analog precoding, the complexity order is $O(K^2N_s^2N_t)$, primarily from computing the residual channel matrices. In the digital stage, RCD incurs a complexity of $O(K^4N_s^3)$. Overall, EMU-HPC has a total complexity order of $O(KN_sN_tN_rN_{epoch} + K^2N_s^2N_t + K^4N_s^3)$.

For comparison, Table I summarizes the complexity orders of different hybrid precoding and combining algorithms. It shows that the complexity of EMU-HPC is significantly lower than HyEB [6] and HRCD [9], and is no higher than HyBD [4], HySBD [12], and SROA [7], as verified by simulations.

Additionally, the feedback overhead of EMU-HPC in the analog stage is $KN_sN_t$, since $\{\bar{\mathbf{H}}_k\}_{k=1}^K$ is fed back only once. This is identical to HyBD, HySBD, and HyEB, but much lower than $KN_sN_{iter}(N_t + N_r)$ required by SROA, where $N_{iter}$ is the maximum number of data-feedback rounds among all AOs. These AOs may cause extra communication delays and further hinder practical deployment. Besides, all algorithms have identical feedback overhead in the digital stage.

**Algorithm 2** The Proposed Efficient Multi-User Hybrid Precoding and Combining (EMU-HPC) Algorithm

---

**Input:** $\{\mathbf{H}_k\}_{k=1}^K$, $N_{epoch}$, $\gamma$, $\varepsilon$
**Output:** $\{\mathbf{W}_{A,k}, \mathbf{W}_{D,k}\}_{k=1}^K$, $\mathbf{F}_A$, $\mathbf{F}_D$
1: **for** $k = 1$ to $K$ **do**
2:     Obtain $\mathbf{W}_{A,k}$ by Algorithm 1
3:     Compute $\bar{\mathbf{H}}_k$ by $\bar{\mathbf{H}}_k = \mathbf{W}_{A,k}^H \mathbf{H}_k$
4: **end for**
5: Initialize $\mathbf{F}_A = \emptyset$ and $\bar{\mathbf{H}}_{res,k} = \bar{\mathbf{H}}_k, \forall k$
6: **for** $k = 1$ to $K$ **do**
7:     Perform SVD of $\bar{\mathbf{H}}_{res,k}$ as $\bar{\mathbf{H}}_{res,k} = \bar{\mathbf{U}}_k \bar{\mathbf{D}}_k \bar{\mathbf{V}}_k^H$
8:     Obtain $\widehat{\mathbf{F}}_{A,k}$ by $\widehat{\mathbf{F}}_{A,k} = \bar{\mathbf{V}}_k(:, 1 : M_r)$
9:     Obtain $\mathbf{F}_{A,k}$ by (13)
10:     **for** $p = k + 1$ to $K$ **do**
11:         $\bar{\mathbf{H}}_{res,p} \leftarrow \bar{\mathbf{H}}_{res,p}(\mathbf{I}_{N_t} - \widehat{\mathbf{F}}_{A,k}\widehat{\mathbf{F}}_{A,k}^H)$
12:     **end for**
13:     $\mathbf{F}_A \leftarrow [\mathbf{F}_A, \mathbf{F}_{A,k}]$
14: **end for**
15: **for** $k = 1$ to $K$ **do**
16:     Obtain $\mathbf{T}_k^{(a)}$ by (16)
17:     Obtain $\widetilde{\mathbf{H}}_k \mathbf{T}_k^{(a)}$ with SVD $\widetilde{\mathbf{H}}_k \mathbf{T}_k^{(a)} = \mathbf{U}_{b,k} \mathbf{D}_{b,k} \mathbf{V}_{b,k}^H$
18:     Obtain $\mathbf{T}_k^{(b)}$ and $\mathbf{W}_{D,k}$ by (17)
19: **end for**
20: Obtain the digital precoders by (18), $\forall k$, and form $\mathbf{F}_D$

---

TABLE I
COMPARISON OF COMPUTATIONAL COMPLEXITY ORDERS

| Algorithm | Computational Complexity Order |
|-----------|-------------------------------|
| EMU-HPC | $O\left[KN_tN_s(N_{epoch}N_r + KN_s) + K^4N_s^3\right]$ |
| HyBD | $O\left[KN_t\left(N_r^2 + KN_s^2\right) + K^4N_s^3\right]$ |
| HySBD | $O\left[KN_t\left(N_r^2 + KN_s^2\right) + K^4N_s^3\right]$ |
| HyEB | $O[KN_s^2N_{iter}(N_t^2 + N_r^2) + K^2N_s^4N_t + K^4N_s^3]$ |
| HRCD | $O\left[KN_t\left(N_r^2 + KN_s^2\right) + K^4N_s^3 + KN_sN_cN_{ray}\right]$ |
| SROA | $O\left[KN_tN_s(N_{iter}N_r + KN_r) + K^4N_s^3\right]$ |

## IV. SIMULATION RESULTS

Apart from setting the large-scale fading factor to be $\beta_k = 1, \forall k$, the mmWave channel configuration follows that in [4]. The signal-to-noise ratio (SNR) is defined as $P_t/\sigma_n^2$. The EMU-HPC algorithm is compared with HyBD [4], HySBD [12], MGLRAM [5], HRCD [9], and SROA [7]. The capacity-achieving iterative waterfilling-based dirty paper coding (IWDPC) [13] is used as a benchmark for performance upper bound. For EMU-HPC, the step size of VR-PCA is $\gamma = 0.01$, the maximum number of epochs is 10 and the error tolerance is $\varepsilon = 10^{-3}$ (same as SROA for fairness).

Fig. 1 illustrates the sum-rate performance over SNR in the mmWave channel, showing that EMU-HPC outperforms other hybrid precoding and combining methods, except at very low SNR. This improvement is mainly achieved through joint interference management. Compared with SROA, EMU-HPC requires far less feedback overhead to achieve a superior sum-rate, rendering it more practical for deployment. Fig. 2 further
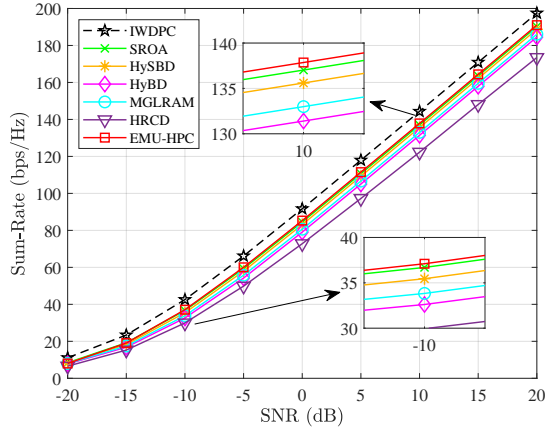
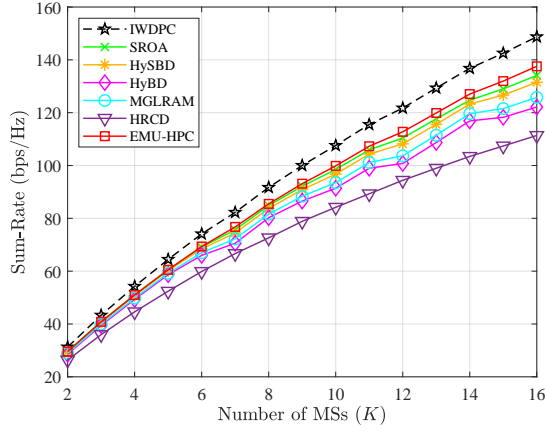Fig. 1. Sum-rate vs. SNR in range of [-20,20] dB for $256 \times 16$ 8-user mmWave massive MIMO systems with $N_\mathrm{s} = 2$.



Fig. 2. Sum-rate vs. the number of MSs in the range of $[2, 16]$ for $256 \times 16$ mmWave massive MIMO systems with $N_\mathrm{s} = 2$ (SNR=0dB).
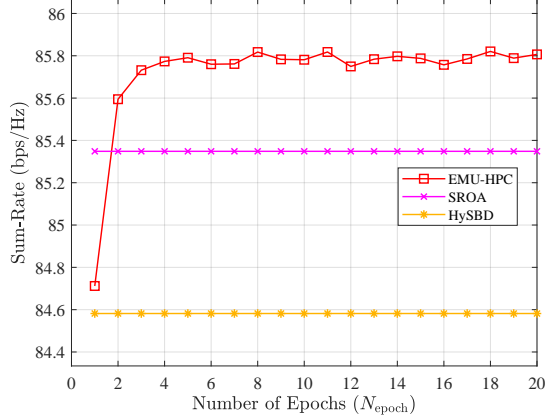


Fig. 3. Sum-rate vs. the maximum number of epochs in the range of $[1, 20]$ for $256 \times 16$ 8-user mmWave massive MIMO systems with $N_\mathrm{s} = 2$ (SNR=0dB).

presents the sum-rate performance versus the increasing numbers of MSs, where EMU-HPC achieves higher sum-rate than other hybrid precoding and combining schemes at medium and high $K$, and remains competitive at low $K$. Compared with MGLRAM, which maximizes both MUI and array gains via (10), EMU-HPC attains better performance by balancing array gains and MUI suppression during analog precoding.

Fig. 3 illustrates the sum-rate performance as $N_\mathrm{epoch}$ in-

creases, indicating that EMU-HPC requires only $N_\mathrm{epoch} = 2$ to surpass SROA, thereby satisfying $N_\mathrm{epoch}N_\mathrm{r} + KN_\mathrm{s} < N_\mathrm{iter}N_\mathrm{r} + KN_\mathrm{r}$ with $N_\mathrm{iter}$ averaging $\bar{N}_\mathrm{iter} = 13.96$. Moreover, $N_\mathrm{epoch} = 1$ is sufficient to outperform HySBD, resulting in $N_\mathrm{s}N_\mathrm{epoch}N_\mathrm{r} \ll N_\mathrm{r}^2$. Hence, EMU-HPC achieves the lowest computational complexity among all hybrid precoding and combining methods in Table I, demonstrating its efficiency.

## V. CONCLUSION

In this paper, we propose the EMU-HPC method for mmWave massive MU-MIMO to improve the sum-rate with low feedback overhead. NI-GLRAM is employed in the analog stage to maximize the array gain with minimal feedback, while joint analog–digital interference suppression is realized through the orthogonal matrix projection in analog precoding and RCD in the digital stage. Numerical results demonstrate that EMU-HPC achieves better performance compared with existing hybrid precoding and combining algorithms.

## REFERENCES

[1] N. Li and P. Fan, "Distributed cell-free massive MIMO versus cellular massive MIMO under UE hardware impairments," *Chinese Journal of Electronics*, vol. 33, no. 5, pp. 1274–1285, 2024.

[2] X. Zhan, Z. Sun, F. Shu, Y. Chen, X. Cheng, Y. Wu, Q. Zhang, Y. Li, and P. Zhang, "Rapid phase ambiguity elimination methods for DOA estimator via hybrid massive MIMO receive array," *Chinese Journal of Electronics*, vol. 33, no. 1, pp. 175–184, 2024.

[3] Q. Zhang and B. Wu, "Characteristic mode analysis for pattern diversity and beamforming: A survey," *Chinese Journal of Electronics*, vol. 33, no. 5, pp. 1117–1126, 2024.

[4] 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.

[5] N. Song, H. Sun, and T. Yang, "Coordinated hybrid beamforming for millimeter wave multi-user massive MIMO systems," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, 2016, pp. 1–6.

[6] C. Hu, J. Liu, X. Liao, Y. Liu, and J. Wang, "A novel equivalent baseband channel of hybrid beamforming in massive multiuser MIMO systems," *IEEE Commun. Lett.*, vol. 22, no. 4, pp. 764–767, Apr. 2018.

[7] G. M. Zilli and W. P. Zhu, "Constrained channel decomposition-based hybrid beamforming for mmwave massive MIMO systems," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 1707–1720, 2020.

[8] J. Liu and S. Chen, "Non-iterative generalized low rank approximation of matrices," *Pattern Recognit. Lett.*, vol. 27, no. 9, pp. 1002–1008, 2006.

[9] 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.

[10] O. Shamir, "A stochastic PCA and SVD algorithm with an exponential convergence rate," in *Proc. 32nd Int. Conf. Mach. Learn.*, 2015, pp. 144–152.

[11] Z. Allen-Zhu and Y. Li, "LazySVD: Even faster SVD decomposition yet without agonizing pain," in *Proc. 30th Int. Conf. on Neural Inf. Process. Syst.*, 2016, pp. 982–990.

[12] X. Wu, D. Liu, and F. Yin, "Hybrid beamforming for multi-user massive MIMO systems," *IEEE Trans. Commun.*, vol. 66, no. 9, pp. 3879–3891, Sep. 2018.

[13] N. Jindal, W. Rhee, S. Vishwanath, S. A. Jafar, and A. Goldsmith, "Sum power iterative water-filling for multi-antenna Gaussian broadcast channels," *IEEE Trans. Inf. Theory*, vol. 51, no. 4, pp. 1570–1580, Apr. 2005.