

Movable Antenna-Aided Hybrid Beamforming for Multi-User Communications

Yichi Zhang ¹, Graduate Student Member, IEEE,
Yuchen Zhang ¹, Member, IEEE, Lipeng Zhu ², Member, IEEE,
Sa Xiao ¹, Wanbin Tang ¹, Yonina C. Eldar ³, Fellow, IEEE,
and Rui Zhang ¹, Fellow, IEEE

Abstract—In this correspondence, we propose a movable antenna (MA)-aided multi-user hybrid beamforming scheme with a sub-connected structure, where multiple movable sub-arrays can independently change their positions within different local regions. To maximize the system sum rate, we jointly optimize the digital beamformer, analog beamformer, and positions of sub-arrays, under the constraints of unit modulus, finite movable regions, and power budget. Due to the non-concave/non-convex objective function/constraints, as well as the highly coupled variables, the formulated problem is challenging to solve. By employing fractional programming, we develop an alternating optimization framework to solve the problem via a combination of Lagrange multipliers, penalty method, and gradient descent. Numerical results reveal that the proposed MA-aided hybrid beamforming scheme significantly improves the sum rate compared to its fixed-position antenna (FPA) counterpart. Moreover, with sufficiently large movable regions, the proposed scheme with sub-connected MA arrays even outperforms the fully-connected FPA array.

Index Terms—Movable antenna, hybrid beamforming, multi-user communication, sum-rate maximization.

I. INTRODUCTION

Utilizing high-frequency millimeter-wave and/or terahertz bands for wireless communication has been recognized as an essential trend for advancing beyond 5 G systems. This is attributed to the enhanced data rate available by exploiting their ultra-broad bandwidths. In order to compensate for the severe path loss in these bands, large antenna arrays are usually needed to provide high beamforming gains. Unfortunately, the heavy power consumption and high cost of radio frequency (RF) chains make fully-digital beamforming impractical. Hybrid beamforming [1], [2], which leverages a small number of RF chains and an analog

front end consisting of phase shifters (PSs), has been considered as a promising technique to achieve a good trade-off between hardware cost and communication performance.

Based on the fabrication of PS enabled front end, the hybrid beamformer can be categorized into fully-connected and sub-connected structures [3]. The fully-connected structure, with all RF chains connected to all antennas through PSs, allows more design degrees of freedom (DoFs). However, the excessive use of PSs entails a large cost. In comparison, by connecting each RF chain only with a portion of antennas, the sub-connected structure is more energy efficient and easier to implement [4]. To improve the communication performance with a limited number of RF chains, different algorithms have been proposed [1], [2]. In [1], the hybrid beamformer was optimized by approximating the fully-digital beamformer in an iterative manner. In [2], the authors developed a general framework for hybrid beamforming design under various hardware architectures. However, both existing fully-connected and sub-connected hybrid beamforming schemes rely on fixed-position antennas (FPAs), which limit their communication performance because the spatial variation of wireless channels is not fully exploited at the transmitter/receiver.

Recently, the movable antenna (MA), also known as fluid antenna system, has been proposed as a promising technology to enhance wireless communication performance [5], [6], [7], [8]. Specifically, an MA can flexibly tune its position and/or rotation to reconfigure the wireless channel towards a more favorable condition. Several prior works have investigated MA-aided systems for improving communication performance [9], [10], [11], [12], [13]. The authors in [9] initially introduced a field-response based channel model for MA-aided systems. Based on this model, the work [10] characterized the capacity of an MA-aided multi-input multi-output (MIMO) system. Compared to antenna selection, MA systems can yield better communication performance using fewer antennas because continuous movement of antennas in a large region can fully exploit the channel spatial variation, thereby reducing antenna costs. In [11], [12], [13], the MA was generalized to multi-user systems. The authors in [12] extended MAs to secure communication by jointly optimizing the MA array geometry and beamforming vector. Most existing works considered that each MA is connected to an RF chain and a motor (for antenna movement) [8], [9], [10], [11], [12], [13]. However, implementing a fully-digital structure of an MA array with a large number of motors brings additional power consumption and significant challenges in system implementation, such as the intertwining of antenna connecting wires during movement.

The main contributions of this correspondence are as follows:

- Instead of independently moving each antenna element, each sub-array is driven by a single motor and moves within a restricted region such that the antennas therein can move collectively, which significantly reduces the hardware costs and simplifies the implementation complexity.
- Utilizing the proposed MA-aided sub-connected structure, we consider downlink multi-user communications with the aim of maximizing the system sum rate, which is non-convex and thus difficult to be optimally solved. Facilitated by fractional programming (FP), we propose an alternating optimization (AO) framework to solve the problem, by jointly applying Lagrange multipliers, penalty method, and gradient descent.
- Numerical results demonstrate that our proposed MA-aided hybrid beamforming scheme outperforms its FPA counterpart. Moreover, under certain practical conditions, the proposed scheme with sub-connected MA arrays yields a higher sum rate compared to the fully-connected FPA array.

Received 2 April 2024; revised 15 September 2024; accepted 4 January 2025. Date of publication 23 January 2025; date of current version 20 June 2025. This work was supported by the National Natural Science Foundation of China under Grant 62301117, Grant 62001094, and Grant U19B2014. The review of this article was coordinated by Prof. Hsiao-Feng Lu. (Corresponding authors: Lipeng Zhu; Yuchen Zhang.)

Yichi Zhang, Yuchen Zhang, and Wanbin Tang are with the National Key Laboratory of Wireless Communications, University of Electronic Science and Technology of China, Chengdu 611731, China (e-mail: yczhang@std.uestc.edu.cn; yc_zhang@std.uestc.edu.cn; wbtang@uestc.edu.cn).

Lipeng Zhu is with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 117583 (e-mail: zhulp@nus.edu.sg).

Sa Xiao is with the National Key Laboratory of Science and Technology on Communications, University of Electronic Science and Technology of China, Chengdu 611731, China, and also with the Kash Institute of Electronics and Information Industry, Kash 844000, China (e-mail: xiaosa@uestc.edu.cn).

Yonina C. Eldar is with the Faculty of Mathematics and Computer Science, Weizmann Institute of Science, Rehovot 7610001, Israel (e-mail: yonina.eldar@weizmann.ac.il).

Rui Zhang is with the The Chinese University of Hong Kong, Shenzhen, and Shenzhen Research Institute of Big Data, Shenzhen 518172, China, and also with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 117583 (e-mail: rzhang@cuhk.edu.cn).

Digital Object Identifier 10.1109/TVT.2025.3533078

B. Problem Formulation

We aim to maximize the sum rate of all users jointly optimizing the central positions of UPAs \mathbf{c} , the analog beamformer \mathbf{W}_A , and the digital beamformer \mathbf{W}_D . The optimization problem is formulated as

$$\max_{\mathbf{W}_D, \mathbf{W}_A, \mathbf{c}} \sum_{k=1}^K R_k \quad (7a)$$

$$\text{s.t. } \mathbf{c}_{m,n} \in \mathcal{C}_{m,n}, \forall m, n, \quad (7b)$$

$$\|\mathbf{W}_A \mathbf{W}_D\|_F^2 \leq P_{\max}, \quad (7c)$$

$$|p_{m,n}^{i,j}| = 1, \forall m, n, i, j, \quad (7d)$$

where P_{\max} denotes the power budget. Due to the non-concave/non-convex objective function/constraints, as well as the coupled variables, (7) is difficult to solve. Specifically, in contrast to the existing FPA-based hybrid beamforming designs [1], [2], [15], the movable positions of the UPAs are implicitly contained in the channel vectors and introduce additional interdependencies among variables in both the numerator and denominator of the achievable rate for each user, thus posing further challenges.

III. PROPOSED SOLUTION

In this section, we develop a low-complexity algorithm to obtain a suboptimal solution for (7). Specifically, we first invoke the FP framework [16], a iterative optimization framework shared similar inspiration with the well-known weighted minimum mean-square error (WMMSE) framework [17], to recast (7) into a more tractable form. Next, based on AO, we proceed to decompose the FP problem into three subproblems. Then, the penalty method is employed to address the unit modulus constraint. Finally, we propose a gradient decent method to optimize the positions of the UPAs.

A. FP-Based Reformulation

Let $\gamma = [\gamma_1, \gamma_2, \dots, \gamma_K]$ and $\omega = [\omega_1, \omega_2, \dots, \omega_K]$ be slack variables. To address the fractional forms in the objective function, we employ the FP framework in [16], through which (7) is equivalently reformulated as

$$\max_{\mathbf{W}_D, \mathbf{W}_A, \mathbf{c}, \gamma, \omega} \mathcal{L} = \sum_{k=1}^K \log(1 + \gamma_k) - \gamma_k + (1 + \gamma_k) \left(2\mathcal{R}\{\omega_k^* a_k\} - |\omega_k|^2 b_k \right) \quad (8a)$$

$$\text{s.t. } \gamma_k > 0, \omega_k \in \mathbb{C}, \forall k \in \mathcal{K}, \quad (8b)$$

$$(7b), (7c), (7d), \quad (8c)$$

where, for fixed \mathbf{W}_D , \mathbf{W}_A , and \mathbf{c} , the optimal values of γ_k and ω_k are given by

$$\gamma_k^* = \frac{|a_k|^2}{\sum_{k' \neq k} |\mathbf{h}_k(\mathbf{c})^H \mathbf{W}_A \mathbf{w}_{k'}|^2 + \sigma_k^2} \quad (9)$$

and

$$\omega_k^* = \frac{a_k}{b_k}, \quad (10)$$

respectively, with $a_k = \mathbf{h}_k(\mathbf{c})^H \mathbf{W}_A \mathbf{w}_k$ and $b_k = \sigma_k^2 + \sum_{k'=1}^K |\mathbf{h}_k(\mathbf{c})^H \mathbf{W}_A \mathbf{w}_{k'}|^2$.

To decouple the variables $\{\mathbf{W}_D, \mathbf{W}_A, \mathbf{c}, \gamma, \omega\}$ in (8), we employ the AO framework to update each variable iteratively with other variables being fixed. Since the closed-form expressions of variables γ and ω are provided in (9) and (10), respectively, we only need to update the other variables with given $\{\gamma, \omega\}$ in each iteration. Define $\alpha_k = \log(1 + \gamma_k) - \gamma_k - (1 + \gamma_k)|\omega_k|^2 \sigma_k^2$. Since α_k is only relevant to variables γ and ω , we omit it in the rest of this correspondence.

B. Digital Beamformer Design

Given $\{\mathbf{W}_A, \mathbf{c}, \gamma, \omega\}$, we seek the optimal value of the digital beamformer \mathbf{W}_D . We recast the power constraint as [1]

$$\|\mathbf{W}_A \mathbf{W}_D\|_F^2 = N_v N_h \|\mathbf{W}_D\|_F^2 \leq P_{\max}. \quad (11)$$

Thus, \mathbf{W}_A and \mathbf{W}_D are decoupled in (7c). The subproblem of digital beamformer design is then formulated as

$$\max_{\mathbf{W}_D} - \sum_{k=1}^K \mu_k \sum_{k'=1}^K |\xi_k^H \mathbf{w}_{k'}|^2 - 2\mathcal{R}\{\beta_k^H \mathbf{w}_k\} \quad (12a)$$

$$\text{s.t. } \|\mathbf{W}_D\|_F^2 \leq \frac{P_{\max}}{N_v N_h}, \quad (12b)$$

where $\beta_k^H = (1 + \gamma_k)\omega_k^* \xi_k^H$, $\xi_k^H = \mathbf{h}_k(\mathbf{c})^H \mathbf{W}_A$, and $\mu_k = (1 + \gamma_k)|\omega_k|^2$. Let $\Xi = \sum_{k=1}^K \mu_k \xi_k \xi_k^H$. According to the Lagrangian multipliers, the convex quadratic optimization problem admits the closed-form solution given by

$$\mathbf{w}_k^*(\lambda) = \left(\Xi + \lambda \mathbf{I}_{N_{\text{RF}}^h N_{\text{RF}}^v} \right)^{-1} \beta_k, \quad (13)$$

where λ is a Lagrange multiplier such that $\lambda(\|\mathbf{W}_D\|_F^2 - \frac{P_{\max}}{N_v N_h}) = 0$. That is, if $\|\mathbf{W}_D^*\|_F^2 = \sum_{k=1}^K \|\mathbf{w}_k^*(0)\|^2 \leq \frac{P_{\max}}{N_v N_h}$, $\lambda = 0$. Otherwise, $\lambda > 0$ can be found by a bisection search.

C. Analog Beamformer Design

With given $\{\mathbf{W}_D, \mathbf{c}, \gamma, \omega\}$, we set out to optimize \mathbf{W}_A . Define $\tilde{\mathbf{p}} = [\mathbf{p}_{1,1}^T, \mathbf{p}_{1,2}^T, \dots, \mathbf{p}_{N_{\text{RF}}^h, N_{\text{RF}}^v}^T]^T$, $\tilde{\mathbf{h}}_{k,k'} = \mathbf{h}_k(\mathbf{c})^H (\text{diag}(\mathbf{w}_{k'}) \otimes \mathbf{I}_{N_h N_v})$, and $\tilde{\beta}_k^H = (1 + \gamma_k)\omega_k^* \xi_k^H$. Based on the penalty method, the subproblem of analog beamformer design is formulated as

$$\max_{\tilde{\mathbf{p}}, \phi} \sum_{k=1}^K \left(-\mu_k \sum_{k'=1}^K |\tilde{\mathbf{h}}_{k,k'}^H \phi|^2 + 2\mathcal{R}\{\tilde{\beta}_k^H \phi\} \right) - \eta \|\phi - \tilde{\mathbf{p}}\|_2^2 \quad (14a)$$

$$\text{s.t. } |\tilde{\mathbf{p}}_s| = 1, \forall s = 1, 2, \dots, N_{\text{RF}}^h N_{\text{RF}}^v N_h N_v, \quad (14b)$$

where ϕ is the introduced continuous variable and $\eta > 0$ is the penalty parameter. Problem (14) can be solved via updating $\tilde{\mathbf{p}}$ and ϕ iteratively. With given $\tilde{\mathbf{p}}$, the optimal solution of ϕ is

$$\phi^* = [\tilde{\Xi}]^{-1} \tilde{\beta}, \quad (15)$$

where $\tilde{\Xi} = \sum_{k=1}^K \sum_{k'=1}^K (\mu_k \tilde{\mathbf{h}}_{k,k'} \tilde{\mathbf{h}}_{k,k'}^H) + \eta \mathbf{I}_N$ and $\tilde{\beta} = \sum_{k=1}^K \tilde{\beta}_k + \eta \tilde{\mathbf{p}}$. Given ϕ , the optimization of $\tilde{\mathbf{p}}$ can be extracted from (14) as

$$\min_{\tilde{\mathbf{p}}} \|\phi - \tilde{\mathbf{p}}\|_2^2 \quad (16a)$$

$$\text{s.t. } |\tilde{\mathbf{p}}_s| = 1, \forall s = 1, 2, \dots, N_{\text{RF}}^h N_{\text{RF}}^v N_h N_v. \quad (16b)$$

The optimal phase of $\tilde{\mathbf{p}}$ is given by

$$\arg\{\tilde{\mathbf{p}}\} = \arg\{\phi\}. \quad (17)$$

D. MA Position Design

Given $\{\mathbf{W}_D, \mathbf{W}_A, \gamma, \omega\}$, the subproblem of MA position design can be decomposed into iteratively optimizing $\mathbf{c}_{m,n}$ with the other

central points of UPAs fixed. The optimization problem is formulated as

$$\max_{\mathbf{c}_{m,n} \in \mathcal{C}_{m,n}} \mathcal{L}(\mathbf{c}_{m,n}) = \sum_{k=1}^K 2\mathcal{R} \left\{ \mathbf{h}_k(\mathbf{t}_{m,n})^H \bar{\beta}_{k,m,n} \right\} - \mu_k \sum_{k'=1}^K \left| \sum_{m=1}^{N_{\text{RF}}^h} \sum_{n=1}^{N_{\text{RF}}^v} \mathbf{h}_k(\mathbf{t}_{m,n})^H \bar{\xi}_{k',m,n} \right|^2, \quad (18)$$

$\bar{h}_{k,k'}(\mathbf{t}_{m,n})$

where $\bar{\beta}_{k,m,n} = (1 + \gamma_k) \omega_k^* \bar{\xi}_{k,m,n}$, $\bar{\xi}_{k,m,n} = w_{k,m,n} \mathbf{p}_{m,n}$, with $w_{k,m,n}$ denoting the $((m-1)N_{\text{RF}}^h + n)$ -th element of \mathbf{w}_k . Due to the complicated expression of $\mathcal{L}(\mathbf{c}_{m,n})$, we propose a gradient descent method to optimize $\mathbf{c}_{m,n}$. In each iteration, the gradient can be calculated as

$$\nabla_{\mathbf{c}_{m,n}} \mathcal{L}(\mathbf{c}_{m,n}) = 2 \sum_{k=1}^K \mathcal{R} \left\{ \frac{\partial \mathbf{h}_k(\mathbf{t}_{m,n})}{\partial \mathbf{c}_{m,n}} \bar{\beta}_{k,m,n}^* - \mu_k \sum_{k'=1}^K \frac{\partial \mathbf{h}_k(\mathbf{t}_{m,n})}{\partial \mathbf{c}_{m,n}} \bar{\xi}_{k',m,n}^* \bar{h}_{k,k'}(\mathbf{t}_{m,n}) \right\}, \quad (19)$$

where the expression of partial derivation is given by

$$\left[\frac{\partial \mathbf{h}_k(\mathbf{t}_{m,n})}{\partial \mathbf{c}_{m,n}} \right]_{:, (i-1)N_v + j} = \sum_{l=1}^{L_k} -j \frac{2\pi}{\lambda} \sigma_{k,l} e^{-j \frac{2\pi}{\lambda} (\mathbf{t}_{m,n}^{i,j})^T \boldsymbol{\rho}_{k,l}} \boldsymbol{\rho}_{k,l}, \quad \forall i, j. \quad (20)$$

with $\sigma_{k,l} = [\mathbf{\Sigma}_k \mathbf{1}_{L_k^r}]_l$. Then, the (m, n) -th central point in the t -th iteration is updated by moving along the gradient direction and checking whether the new point is still located within the feasible region, i.e.,

$$\mathbf{c}_{m,n}^{(t+1)} = \begin{cases} \tilde{\mathbf{c}}_{m,n}, & \text{if } \tilde{\mathbf{c}}_{m,n} \in \mathcal{C}_{m,n} \text{ and } \mathcal{L}(\tilde{\mathbf{c}}_{m,n}) \geq \mathcal{L}(\mathbf{c}_{m,n}^{(t)}), \\ \mathbf{c}_{m,n}^{(t)}, & \text{otherwise,} \end{cases} \quad (21)$$

where $\tilde{\mathbf{c}}_{m,n} = \mathbf{c}_{m,n}^{(t)} + \kappa \nabla_{\mathbf{c}_{m,n}} \mathcal{L}(\mathbf{c}_{m,n}^{(t)})$ with κ being the step size. Note that the performance of gradient descent heavily depends on the step size. Thus, for each iteration, we initialize κ with a large positive number and update $\kappa \leftarrow \frac{\kappa}{2}$ until finding a feasible point such that $\mathcal{L}(\tilde{\mathbf{c}}_{m,n}) \geq \mathcal{L}(\mathbf{c}_{m,n}^{(t)})$ or achieve the stopping criterion $\kappa < \kappa_{\min}$. It can be observed that the sequence $\{\mathcal{L}(\mathbf{c}^{(t)})\}$ keeps non-decreasing based on update guidelines (21), leading to the convergence of the sequence. Therefore, the gradient descent method is applicable to any objective function with a specific expression.

Algorithm 1 summarizes the workflow for solving problem (7). We next provide a brief proof for the convergence of the proposed algorithm. For analysis convenience, let $L(\mathbf{W}_{\text{D}}^t, \mathbf{W}_{\text{A}}^t, \mathbf{c}^t)$ denote the value of the penalty objective function (14a) over t -th iteration. Since the penalty term $\eta \|\phi - \tilde{\mathbf{p}}\|_2^2$ is irrelevant to digital beamformer and central positions of UPAs, Algorithm 1 generates a non-decreasing sequence $\{L(\mathbf{W}_{\text{D}}^t, \mathbf{W}_{\text{A}}^t, \mathbf{c}^t)\}$, which has a finite upper bound, thus ensuring the convergence of the penalty objective value sequence [18]. Moreover, the penalty factor can be set such that the penalty term tends to zero, causing the penalized objective function (14a) to reduce to the original problem (8a), thus guaranteeing the convergence of Algorithm 1. The computational complexities of solving the variables $\{\mathbf{W}_{\text{D}}, \mathbf{W}_{\text{A}}, \mathbf{c}, \gamma, \omega\}$ are scaled with $\mathcal{O}((N_{\text{RF}}^h N_{\text{RF}}^v)^3)$, $\mathcal{O}(N^3)$, $\mathcal{O}(K N N_{\text{RF}}^h N_{\text{RF}}^v \sum_{k=1}^K L_k)$, $\mathcal{O}(K N)$, and $\mathcal{O}(K N)$, respectively. Furthermore, we denote the number of iterations needed for the convergence of the AO algorithm as I_{AO} . Considering the dominant computational steps, the maximum computational complexity of Algorithm 1 can be expressed as $\mathcal{O}(I_{\text{AO}}(N^3 + K N N_{\text{RF}}^h N_{\text{RF}}^v \sum_{k=1}^K L_k))$.

Algorithm 1: Proposed Algorithm for Solving Problem (7).

Input: $\mathbf{W}_{\text{D}}^{(0)}, \mathbf{W}_{\text{A}}^{(0)}, \mathbf{c}^{(0)}, \gamma^{(0)}, \omega^{(0)}, \{R_k^{(0)}\}, \tilde{\kappa}, \kappa_{\min}, \varepsilon, I$.

Output: $\mathbf{W}_{\text{D}}^*, \mathbf{W}_{\text{A}}^*, \mathbf{c}^*$.

```

1: Set iteration index  $t = 1$ .
2: repeat
3:   Update  $\gamma^{(t)}$  and  $\omega^{(t)}$  via (9) and (10), respectively.
4:   Obtain  $\mathbf{W}_{\text{D}}^{(t)}$  via (13).
5:   Obtain  $\phi^{(t)}$  via (15).
6:   Obtain  $\mathbf{W}_{\text{A}}^{(t)}$  via (5) and (17).
7:   for all  $m = 1 : N_{\text{RF}}^h$  do
8:     for all  $n = 1 : N_{\text{RF}}^v$  do
9:       Calculate the gradient  $\nabla_{\mathbf{c}_{m,n}} \mathcal{L}(\mathbf{c}_{m,n}^{(t)})$  via (19).
10:      Initialize the step size  $\kappa = \tilde{\kappa}$ .
11:      repeat
12:        Compute  $\tilde{\mathbf{c}}_{m,n} = \mathbf{c}_{m,n}^{(t)} + \kappa \nabla_{\mathbf{c}_{m,n}} \mathcal{L}(\mathbf{c}_{m,n}^{(t)})$ .
13:        Shrink the step size  $\kappa \leftarrow \frac{\kappa}{2}$ .
14:        Update  $\mathbf{c}_{m,n}^{(t+1)}$  according to (21).
15:      until  $\tilde{\mathbf{c}}_{m,n} \in \mathcal{C}_{m,n}$  and  $\mathcal{L}(\tilde{\mathbf{c}}_{m,n}) \geq \mathcal{L}(\mathbf{c}_{m,n}^{(t)})$  or
         $\kappa < \kappa_{\min}$ .
16:      end for
17:    end for
18:    Calculate the sum rate  $\sum_{k=1}^K R_k^{(t+1)}$  according to (6).
19:    Update  $t = t + 1$ .
20: until  $|\sum_{k=1}^K R_k^{(t+1)} - R_k^{(t)}| < \varepsilon$  or the maximum iteration
    number  $I$  is reached.

```

IV. NUMERICAL RESULTS

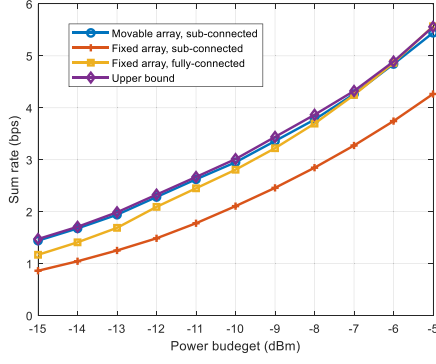
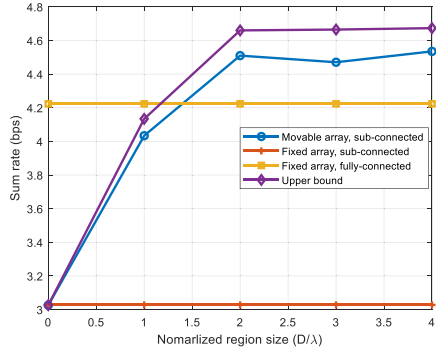
In this section, numerical results are provided to evaluate the performance of our proposed scheme. Unless stated otherwise, we consider a scenario where a 16-antenna BS operating at $f_c = 30$ GHz serves $K = 16$ users. Besides, $N_{\text{RF}}^h = N_{\text{RF}}^v = 4$, $N_h = N_v = 2$. The central point of each UPA remains in the x - z plane, where the inter-element spacing is $\frac{\lambda}{2}$ to avoid antenna coupling and $\lambda = 0.01$ m represents the wavelength. The distance between the k -th user and the BS is uniformly distributed from 20 to 100m, i.e., $d_k \sim \mathcal{U}[20, 100]$. For each user, we assume an equal number of transmit and receive channel paths $L_k^t = L_k^r = L = 6$. The PRM is defined as a diagonal matrix $\Sigma_k = \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_L\}$, where $\sigma_{k,l} \sim \mathcal{CN}(0, (\rho_0 d_k^{-\alpha})^2 / L)$, $\rho_0 = -40$ dB is the channel power gain at the reference distance 1m, and $\alpha = 2.8$ denotes the corresponding path loss exponent. The elevation and azimuth angles are assumed to follow the joint probability density

function $f(\theta_{k,l}^x, \phi_{k,l}^x) = \frac{\cos \theta_{k,l}^x}{2\pi} [9]$, $\chi \in \{t, r\}$. For each UPA, the size of movable frame is $D \times D$. In addition, $P_{\max} = 10$ dBm, $\sigma_k^2 = -80$ dBm, $D = 2\lambda$, $\tilde{\kappa} = 10$, $\varepsilon = 10^{-3}$, $\eta = 10$, $\kappa_{\min} = 10^{-3}$, and $I = 200$.

To verify the effectiveness of the proposed scheme with Algorithm 1, we consider two baselines for comparison: 1) Fixed array, sub-connected [3]. 2) Fixed array, fully-connected [3]. Besides, by finding the optimal positions of all UPAs under the sub-connected structure within the whole region based on exhaustive search, the upper bound on the sum rate is also obtained for comparison.

Fig. 2 depicts the sum-rate versus the power budget for different schemes. It is observed that with the same power budget, our proposed algorithm achieves a higher sum-rate compared to its FPA counterpart. This suggests that the performance gain from MAs can compensate the performance loss of applying sub-connected structure. Moreover, when transmit power is sufficiently low, the MA-aided sub-connected structure can even outperform the fully-connected FPAs. This indicates that with lower inter-user interference, the MAs' ability to enhance received signal power gain significantly contributes to performance improvement.

Fig. 3 demonstrates the sum-rate versus D/λ . Note that the sum-rate increases as the normalized region size increases until reaching a

Fig. 2. Sum-rate versus P_{max} with $D = 2\lambda$.Fig. 3. Sum-rate versus the normalized region size with $P_{max} = -10$ dBm.

constant value. The reason for it is that increasing region size provides more DoFs for the MAs to exploit the channel spatial variation. However, for a sufficiently large region size, the channel gain may exhibit periodic character in the spatial domain, which indicates a strong spatial correlation [9]. For this case, the effective channel gain for each user cannot be further increased by antenna position optimization, while the multiuser interference has already approached the lower bound. Therefore, further increasing the region size offers negligible improvement, which indicates that a moderate size of the spatial region is sufficient for deploying MAs. Moreover, Fig. 3 also shows that with sufficiently large region size, our proposed scheme is superior to all baselines. However, the performance gap from the upper bound is still substantial, due to that the gradient method is prone to converging to local optimum.

V. CONCLUSION

This correspondence investigated MA-aided multi-user hybrid beamforming under the sub-connected structure. We studied the sum rate maximization problem by jointly designing the digital beamformer, analog beamformer, and movable UPAs' positions. To solve the non-convex optimization problem, we employed FP to transform it into a more tractable form and then developed an AO-based algorithm by applying the techniques of Lagrange multiplier, penalty method and gradient descent. Numerical results demonstrated the superiority of

the proposed MA-aided sub-connected structure compared to the FPA-based system. Moreover, under certain conditions, the proposed scheme with sub-connected MA arrays even outperforms the fully-connected FPA array.

REFERENCES

- [1] 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.
- [2] S. S. Ioushua and Y. C. Eldar, "A family of hybrid analog–digital beamforming methods for massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 67, no. 12, pp. 3243–3257, Jun. 2019.
- [3] J. Du, W. Xu, H. Shen, X. Dong, and C. Zhao, "Hybrid precoding architecture for massive multiuser MIMO with dissipation: Sub-connected or fully connected structures?," *IEEE Trans. Wireless Commun.*, vol. 17, no. 8, pp. 5465–5479, Aug. 2018.
- [4] O. E. Ayach, R. W. Heath, S. Rajagopal, and Z. Pi, "Multimode precoding in millimeter wave MIMO transmitters with multiple antenna sub-arrays," in *Proc. IEEE Glob. Commun. Conf.*, Dec. 2013, pp. 3476–3480.
- [5] L. Zhu and K. K. Wong, "Historical review of fluid antenna and movable antenna," 2024, *arXiv:2401.02362*.
- [6] K. -K. Wong, W. K. New, X. Hao, K. -F. Tong, and C. -B. Chae, "Fluid antenna system—Part I: Preliminaries," *IEEE Commun. Lett.*, vol. 27, no. 8, pp. 1919–1923, Aug. 2023.
- [7] K. -K. Wong, K. -F. Tong, and C. -B. Chae, "Fluid antenna system—Part II: Research opportunities," *IEEE Commun. Lett.*, vol. 27, no. 8, pp. 1924–1928, Aug. 2023.
- [8] L. Zhu, W. Ma, and R. Zhang, "Movable antennas for wireless communication: Opportunities and challenges," *IEEE Commun. Mag.*, vol. 62, no. 6, pp. 114–120, Jun. 2024.
- [9] L. Zhu, W. Ma, and R. Zhang, "Modeling and performance analysis for movable antenna enabled wireless communications," *IEEE Trans. Wireless Commun.*, vol. 23, no. 6, pp. 6234–6250, Jun. 2024.
- [10] W. Ma, L. Zhu, and R. Zhang, "MIMO capacity characterization for movable antenna systems," *IEEE Trans. Wireless Commun.*, vol. 23, no. 4, pp. 3392–3407, Apr. 2024.
- [11] L. Zhu, W. Ma, B. Ning, and R. Zhang, "Movable-antenna enhanced multiuser communication via antenna position optimization," *IEEE Trans. Wireless Commun.*, vol. 23, no. 7, pp. 7214–7229, Jul. 2024.
- [12] Y. Wu, D. Xu, D. W. K. Ng, W. Gerstacker, and R. Schober, "Movable antenna-enhanced multiuser communication: Optimal discrete antenna positioning and beamforming," in *Proc. IEEE Glob. Commun. Conf.*, Dec. 2023, pp. 7508–7513.
- [13] G. Hu, Q. Wu, K. Xu, J. Si, and N. Al-Dhahir, "Secure wireless communication via movable-antenna array," *IEEE Signal Process. Lett.*, vol. 31, pp. 516–520, 2024.
- [14] W. Ma, L. Zhu, and R. Zhang, "Compressed sensing based channel estimation for movable antenna communications," *IEEE Commun. Lett.*, vol. 27, no. 10, pp. 2747–2751, Oct. 2023.
- [15] W. Zhu, H. D. Tuan, E. Dutkiewicz, H. V. Poor, and L. Hanzo, "Max-min rate optimization of low-complexity hybrid multi-user beamforming maintaining rate-fairness," *IEEE Trans. Wireless Commun.*, vol. 23, no. 6, pp. 5648–5662, Jun. 2024.
- [16] K. Shen and W. Yu, "Fractional programming for communication systems—Part II: Uplink scheduling via matching," *IEEE Trans. Signal Process.*, vol. 66, no. 10, pp. 2631–2644, May 2018.
- [17] Q. Shi, M. Razaviyayn, Z.-Q. Luo, and C. He, "An iteratively weighted MMSE approach to distributed sum-utility maximization for a MIMO interfering broadcast channel," *IEEE Trans. Signal Process.*, vol. 59, no. 9, pp. 4331–4340, Sep. 2011.
- [18] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.