# Resource Allocation and Scheduling in Non-coherent User-centric Cell-free MIMO

Hussein A. Ammar*, Raviraj Adve*, Shahram Shahbazpanahi†*, Gary Boudreau‡, and Kothapalli Srinivas‡

*University of Toronto, Dep. of Elec. and Comp. Eng., Toronto, Canada
†University of Ontario Institute of Technology, Dep. of Elec. and Comp. Eng., Oshawa, Canada
‡Ericsson Canada, Ottawa, Canada

*Abstract*—We study the problem of user-scheduling and resource allocation in distributed multi-user, multiple-input multiple-output (MIMO) networks implementing user-centric clustering and non-coherent transmission. We formulate a weighted sum-rate maximization problem which can provide user proportional fairness. As in this setup, users can be served by many transmitters, user scheduling is particularly difficult. To solve this issue, we use block coordinate descent, fractional programming, and compressive sensing to construct an algorithm that performs user-scheduling and beamforming. Our results show that the proposed framework provides an 8- to 10-fold gain in the long-term user spectral efficiency compared to benchmark schemes such as round-robin scheduling. Furthermore, we quantify the performance loss due to imperfect channel state information and pilot training overhead using a defined *area-based* pilot-reuse factor.

*Index Terms*—User-centric clustering, cell-free, user-scheduling, resource allocation, distributed MIMO, distributed antennas system, fairness, imperfect CSI.

## I. INTRODUCTION

Deploying user-centric clustering in distributed multiple-input multiple-output (MIMO) networks enhances the performance of the conventional cell-edge users by placing each user at the effective center of its serving cluster [1], [2]. User-centric clustering can outperform general cell-free networks that assume all the remote radio heads (RRHs) can serve the users [3]. Recently, resource allocation under cell-free MIMO has attracted significant attention. The studies in [4], [5] optimize beamforming design by minimizing a weighted sum mean square error (MSE) utility, which is easier to tackle than weighted sum rate (WSR) maximization problems but suffers a penalty in terms of sum-rate [4].

The work in [3] considers optimizing power allocation to maximize lower bounds for sum-rate and minimum rate. Similarly, [6] optimizes the beamforming to maximize the minimum rate. Note that, max-min rate solutions do not provide flexibility to control the fairness. Moreover, the authors in [7] consider a near-optimal power control algorithm using zero-forcing (ZF) and conjugate beamforming that is simpler than the max–min power approach for cell-free massive MIMO networks. Furthermore, the work in [8] optimizes the beamforming by using a lower-bound for the logarithm function of the rate to obtain a local optimum.

A crucial component in optimizing the WSR is user-scheduling. In conventional networks, techniques that have been investigated include dual decomposition and the gradient method [9], where the scheduling variables are relaxed from being binary. Notably, this relaxation is optimal for a large number of subcarriers [10]. Furthermore, the investigations in [11], [12] use fractional programming to perform resource allocation in conventional networks, where the user-scheduling part is performed using a combinatorial search.

In summary, the main limitation of these studies is either not specifically addressing the user-centric clustering scheme, or ignoring the user-scheduling step for the users, that is, the scheduled users are assumed to be *preselected*. In this paper, we optimize user-scheduling and resource allocation in a user-centric cell-free MIMO network through formulating a WSR problem. We study the non-coherent transmission mode, which does not require the RRHs to strictly synchronize their transmissions, but it prevents from directly using the weighted minimum mean square error (WMMSE) [13]. The scheduling part of the problem cannot be solved efficiently using a combinatorial search algorithm because each user can be served by many RRHs with overlapping serving clusters. To tackle this, we employ tools from block coordinate descent, fractional programming, and compressive sensing, which allow the construction of an algorithm that guarantees convergence of the network sum-rate through a smooth non-decreasing pattern. In summary, the contributions of this paper are:

- Formulating the WSR problem for the non-coherent cell-free multiuser MIMO setting
- Using fractional programming to optimize beamforming and employing compressive sensing to solve the scheduling problem
- Developing and implementing robust beamforming to account for channel estimation errors

The rest of the paper is organized as follows. Section II presents our system model, while Section III formulates the optimization problem. Section IV presents our proposed resource allocation algorithm. Finally, Sections V and VI report our simulation results and conclusion, respectively.

## II. SYSTEM MODEL

### A. Network Model

As shown in Fig. 1, we consider the downlink of a time-division duplex (TDD) system comprising several RRHs, represented by the set $\mathcal{B}$, each equipped with $M$ antennas and jointly serving the active users, represented by the set $\mathcal{U}$. Both
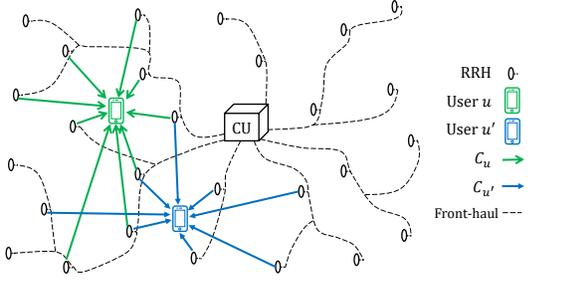
Fig. 1: Serving cluster using user-centric clustering.

RRHs and users are randomly located in 2D space. The RRHs are controlled by a single control unit (CU), and as in [3], we assume a relaxed front-haul constraint, which can be realized through technologies like the radio stripes system [14].

For each user $u \in \mathcal{U}$, we define a cluster $\mathcal{C}_u$ that includes the RRHs that *potentially can be selected* to serve the user according to user-centric clustering. Specifically, $\mathcal{C}_u$ comprises the RRHs with strong average channels, i.e., $\mathcal{C}_u = \{r \mid (\psi_{ru}L(d_{ru})) \geq \rho\}$, where $\psi_{ru}$ denotes the shadowing, $L(d_{ru})$ accounts for the path loss; here, $d_{ru}$ is the distance between RRH $r$ and user $u$. If no RRH meets this criterion, the $\mathcal{C}_u$ for the user comprises the RRH with largest $(\psi_{ru}L(d_{ru}))$. Finally, we represent the users that may be served by RRH $r$ as $\mathcal{E}_r$.

### B. Channel Estimation

Channel estimation is performed through an uplink pilot-training phase of length $\tau_p$. During this phase, we can write the signal $\mathbf{Y}_r \in \mathbb{C}^{M \times \tau_p}$ received at RRH $r$ as

$$\mathbf{Y}_r = \sum_{u \in \mathcal{U}} \sqrt{p_u} \mathbf{h}_{ru} \mathbf{\Phi}_u + \mathbf{Z}_r, \tag{1}$$

where $\mathbf{\Phi}_u \in \mathbb{C}^{1 \times \tau_p}$ is the unit norm ($\mathbf{\Phi}_u \mathbf{\Phi}_u^H = 1$) pilot sequence used by user $u$, $p_u$ is the transmit power of the user, and $\mathbf{Z}_r$ is the additive white Gaussian noise (AWGN) with entries $\sim \mathcal{CN}\left(0, \sigma_Z^2\right)$; $\mathbf{h}_{ru} \in \mathbb{C}^{M \times 1}$, the channel between RRH $r$ and user $u$ is modeled as $\mathbf{h}_{ru} \triangleq \sqrt{\psi_{ru}L(d_{ru})}\mathbf{g}_{ru}$, where $\mathbf{g}_{ru} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_M)$ accounts for small-scale fading.

As in [2], we assume knowledge of the users' transmit powers and large-scale fading. Hence, using $\mathbf{\check{y}}_r = \text{vec}\{\mathbf{Y}_r\} \in \mathbb{C}^{M\tau_p \times 1}$ and linear MMSE, the channel estimate $\mathbf{\hat{h}}_{ru}, \forall u \in \mathcal{E}_r$ can be obtained as $\mathbf{\hat{h}}_{ru} = \mathbf{R}_{ru}\mathbf{R}_r^{-1}\mathbf{\check{y}}_r$, with $\mathbf{R}_{ru} = \sqrt{p_u}\psi_{ru}L(d_{ru})\left(\mathbf{\Phi}_u^* \otimes \mathbf{I}_M\right)$ and

$$\mathbf{R}_r = \sum_{u \in \mathcal{U}} p_u \psi_{ru}L(d_{ru})\left(\mathbf{\Phi}_u^T \mathbf{\Phi}_u^* \otimes \mathbf{I}_M\right) + \sigma_z^2 \mathbf{I}_{M\tau_p}.$$

When the number of users $|\mathcal{U}| \geq \tau_p$, the available pilot sequences need to be reused by the users, adding pilot contamination. This results in the estimated channel $\mathbf{\hat{h}}_{ru} \sim \mathcal{CN}\left(\mathbf{0}, \mathbf{\Psi}_{ru}\right)$, with the error covariance matrix given by [15]

$$\mathbf{\Psi}_{ru} \triangleq \mathbf{D}_{ru}\left(\sum_{u' \in \mathcal{U}_u} \mathbf{D}_{ru'} + \frac{\sigma_Z^2}{p_u}\mathbf{I}_M\right)^{-1}\mathbf{D}_{ru}, \tag{2}$$

where $\mathbf{D}_{ru} \in \mathbb{C}^{M \times M}$ is a diagonal matrix with diagonal entries $[\mathbf{D}_{ru}]_{mm} \triangleq \psi_{ru}L(d_{ru})$, and $\mathcal{U}_u$ is the set of users

employing the same pilot sequence as user $u$ (including user $u$). It is known from MMSE that the channel estimation error $\mathbf{e}_{ru} = \mathbf{h}_{ru} - \mathbf{\hat{h}}_{ru}$ is uncorrelated with $\mathbf{\hat{h}}_{ru}$ and can be modeled as $\mathbf{e}_{ru} \sim \mathcal{CN}\left(\mathbf{0}, \mathbf{\Theta}_{ru}\right)$, where $\mathbf{\Theta}_{ru} \triangleq \mathbf{D}_{ru} - \mathbf{\Psi}_{ru}$.

The downlink signal received at user $u$ can be modeled as

$$y_u = \sum_{r \in \mathcal{C}_u} \sqrt{s_{ru}}\left(\mathbf{\hat{h}}_{ru}^H + \mathbf{e}_{ru}^H\right)\mathbf{w}_{ru}x_{ru}$$

$$+ \sum_{r' \in \mathcal{B}} \sum_{u' \in \mathcal{E}_{r'}, u' \neq u} \sqrt{s_{r'u'}}\left(\mathbf{\hat{h}}_{r'u}^H + \mathbf{e}_{r'u}^H\right)\mathbf{w}_{r'u'}x_{r'u'} + z_u \quad (3)$$

where $\{x_{ru} : r \in \mathcal{C}_u\}$ are the symbols transmitted by the serving RRHs for user $u$ with $\mathbb{E}\{|x_{ru}|^2\} = 1$, $\mathbf{w}_{ru} \in \mathbb{C}^{M \times 1}$ is the precoding vector used by RRH $r$ to serve user $u$, and $z_u \sim \mathcal{CN}(0, \sigma_z^2)$ is the AWGN.

### C. Pilot Assignment (PA) Policy

Properly assigning the pilots to the users is clearly pivotal to decrease pilot contamination.

**Proposition 1.** We propose to use a heuristic low-overhead location-based PA policy. Our policy assigns non-orthogonal pilots for users that are far from each other by using the hierarchical agglomerative clustering (HAC) algorithm [16]. The HAC creates a tree to cluster the users into many groups each containing a number of users less than or equal to the number of available orthogonal pilot sequences. We then assign each group the available orthogonal sequences. The algorithm can be constructed as follows:

1) Treat each active user as a cluster head.
2) Combine the two nearest clusters into one using an average linkage, e.g., Ward's minimum variance criterion.
3) Repeat Step 2 until you reach the root of the tree where all the users are in the same cluster.
4) While backtracking the tree starting from the root, define each cluster when its number of users is less than or equal $\tau_p$.
5) Assign the orthogonal pilots to each cluster randomly.

The HAC algorithm is more consistent than the K-means and Gaussian mixture models, and it is not sensitive to the choice of the used distance-metric [16]. Also, as this algorithm does not require selecting the number of clusters needed, it allows us to easily define the cluster based on an upper limit of the number of users belonging to it, i.e., relate it to $\tau_p$.

### III. PROBLEM FORMULATION

#### A. Problem Definition

To decode the data streams from the RRHs, the users employ successive interference cancellation (SIC). Under the assumption of perfect SIC, the effective achievable rate[1] for

---

[1] This expression is based on using the famous Jensen's Inequality to write down a lower-bound for the data rate with an expectation over the unknown instantaneous channel state information (CSI) error $\{\mathbf{e}_{ru} : r \in \mathcal{B}, u \in \mathcal{E}_r\}$, i.e., $\mathbb{E}_{\mathbf{e}}\{\log(1 + \widetilde{\gamma}_u)\} \geq \log\left(1 + 1/\mathbb{E}_{\mathbf{e}}\{\widetilde{\gamma}_u^{-1}\}\right)$, then using SIC to decode the received data streams at the user. Note that this expression is only used to perform the resource allocation, however, when we plot the performance, we use the actual achievable rate using the actual channels.

user $u$ can be modeled by the CU as [13]

$$R_u = \frac{(\tau_d - \tau_p)}{\tau_d} \log\left(1 + \frac{\sum_{r \in \mathcal{C}_u} s_{ru} |\hat{\mathbf{h}}_{ru}^H \mathbf{w}_{ru}|^2}{A_u(\mathcal{S}, \mathcal{W})}\right), \quad (4)$$

$$\text{with} \quad A_u(\mathcal{S}, \mathcal{W}) = \sum_{r' \in \mathcal{B}} \sum_{u' \in \mathcal{E}_{r'}, u' \neq u} s_{r'u'} \left|\hat{\mathbf{h}}_{r'u}^H \mathbf{w}_{r'u'}\right|^2$$
$$+ \sum_{r' \in \mathcal{B}} \sum_{u' \in \mathcal{E}_{r'}} s_{r'u'} \mathbf{w}_{r'u'}^H \boldsymbol{\Theta}_{r'u} \mathbf{w}_{r'u'} + \sigma_z^2 \quad (5)$$

where $\tau_d$ is the channel coherence time, $\mathcal{S} = \{\mathbf{s}_1, \ldots, \mathbf{s}_{|\mathcal{B}|}\}$ is the set of binary scheduling variables at the RRHs with $\mathbf{s}_r = [s_{ru_1} \ \ldots \ s_{ru_{|\mathcal{E}_r|}}]^T \in \mathbb{B}^{|\mathcal{E}_r| \times 1}$, i.e, if $s_{ru} = 1$, user $u$ is scheduled by RRH $r$, else it is not. Similarly, the set of the beamformers is $\mathcal{W} = \{\mathbf{W}_1, \ldots, \mathbf{W}_{|\mathcal{B}|}\}$ with $\mathbf{W}_r = [\mathbf{w}_{ru_1}, \ldots, \mathbf{w}_{ru_{|\mathcal{E}_r|}}] \in \mathbb{C}^{M \times |\mathcal{E}_r|}$. The term $\boldsymbol{\Theta}_{ru'} = \mathbb{E}\{\mathbf{e}_{ru'}\mathbf{e}_{ru'}^H\}$ is the covariance of the estimation error of the channel between RRH $r$ and user $u'$, and including it in the model allows to construct a robust beamforming.

We formulate the following WSR problem on the CU

$$\text{(P1)} \quad \max_{\mathcal{S}, \mathcal{W}} \quad \sum_{u \in \mathcal{U}} \delta_u \log(1 + \gamma_u) \quad (6a)$$

$$\text{s.t.} \quad \sum_{u \in \mathcal{E}_r} s_{ru} \leq M, \quad r \in \mathcal{B} \quad (6b)$$

$$\sum_{u \in \mathcal{E}_r} \|\mathbf{w}_{ru}\|_2^2 \leq p, \quad r \in \mathcal{B} \quad (6c)$$

$$\gamma_u = \frac{\sum_{r \in \mathcal{C}_u} s_{ru} \left|\hat{\mathbf{h}}_{ru}^H \mathbf{w}_{ru}\right|^2}{A_u(\mathcal{S}, \mathcal{W})}, \quad u \in \mathcal{U} \quad (6d)$$

$$s_{ru} \in \{0, 1\} \quad r \in \mathcal{B}, u \in \mathcal{E}_r \quad (6e)$$

where $\delta_u$ denotes the proportional fair weights for user $u$. The term $A_u$ is defined in (5). Problem (6) optimizes the decision variables $\mathcal{S}$ and $\mathcal{W}$ which determine the user-scheduling and beamforming weight vectors, respectively, such that the total utility in (6a) is maximized. Constraints (6b) prevent the RRHs from simultaneously serving more than $M$ users on the same channel. Constraints (6c) satisfy the power budget of the RRHs, and (6e) show that a user $u$ can be scheduled or not. Constraints (6d) define the effective signal to interference and noise ratio (SINR) as an auxiliary variable.

Problem (6) is a mixed-integer non-convex problem and obtaining a global optimum is mathematically prohibitive.

*B. Problem Analysis*

The beamforming vectors are constructed for users that are actually scheduled on the channel and hence

$$s_{ru} = \mathbb{1}\{\|\mathbf{w}_{ru}\|_2^2\} = \left\|\|\mathbf{w}_{ru}\|_2^2\right\|_0 \quad (7)$$

where $\|\cdot\|_0$ is the $\ell_0$-norm. Using the literature of compressive sensing, the $\ell_0$-norm of a vector $\mathbf{x}$ can be approximated as a weighted convex $\ell_1$-norm $\|\mathbf{x}\|_0 \simeq \sum_m \alpha_m |x_m| = \|\boldsymbol{\alpha}\mathbf{x}\|_1$ [17], where $\alpha_m$ are positive weights that penalize the nonzero coefficients $x_m$, and $\boldsymbol{\alpha} = \mathbf{diag}\{\alpha_1, \alpha_2, \ldots\}$ is a diagonal matrix. For our case, $\mathbf{x} = \|\mathbf{w}_{ru}\|_2$ which is scalar.

We can construct an iterative process to find these weights at each iteration $i$ as suggested in [17]

$$\alpha_{ru}^{(i+1)} = \frac{1}{\left\|\mathbf{w}_{ru}^{(i)}\right\|_2^2 + \epsilon}, \quad (8)$$

where $\epsilon > 0$ provides stability and ensures that a zero-valued component in $\|\mathbf{w}_{ru}\|_2^2$ does not strictly prohibit a nonzero estimate at the update in the next iteration.

As a result, our problem can be formulated as follows

$$\text{(P2)} \quad \max_{\mathcal{W}} \quad \sum_{u \in \mathcal{U}} \delta_u \log(1 + \gamma_u) \quad (9a)$$

$$\text{s.t.} \quad \sum_{u \in \mathcal{E}_r} \alpha_{ru} \|\mathbf{w}_{ru}\|_2^2 \leq M, \quad r \in \mathcal{B} \quad (9b)$$

$$\sum_{u \in \mathcal{E}_r} \|\mathbf{w}_{ru}\|_2^2 \leq p, \quad r \in \mathcal{B} \quad (9c)$$

$$\gamma_u = \frac{\sum_{r \in \mathcal{C}_u} \left|\hat{\mathbf{h}}_{ru}^H \mathbf{w}_{ru}\right|^2}{B_u(\mathcal{W})}, \quad u \in \mathcal{U} \quad (9d)$$

where $B_u(\mathcal{W}) = A_u(\mathbf{1}, \mathcal{W})$. We use the Lagrangian for the equality constraints in (9d)

$$\mathcal{L}(\mathcal{W}, \boldsymbol{\gamma}, \boldsymbol{\nu}) = \sum_{u \in \mathcal{U}} \delta_u \log(1 + \gamma_u)$$
$$- \sum_{u \in \mathcal{U}} \nu_u \left(\gamma_u - \frac{\sum_{r \in \mathcal{C}_u} \left|\hat{\mathbf{h}}_{ru}^H \mathbf{w}_{ru}\right|^2}{B_u(\mathcal{W})}\right) \quad (10)$$

When $\mathcal{W}$ is fixed, we evaluate the first optimality condition of the SINR auxiliary variable $\gamma_u$ by setting the derivative of (10) with respect to $\gamma_u$ to zero, which results in a value for $\nu_u$ that satisfies this optimality. Substituting $\nu_u$ back into (10):

$$f_1(\mathcal{W}, \boldsymbol{\gamma}) = \sum_{u \in \mathcal{U}} \delta_u (\log(1 + \gamma_u) - \gamma_u)$$
$$+ \sum_{u \in \mathcal{U}} \delta_u \left(\frac{(1 + \gamma_u) \sum_{r \in \mathcal{C}_u} \left|\hat{\mathbf{h}}_{ru}^H \mathbf{w}_{ru}\right|^2}{\sum_{r' \in \mathcal{B}} \sum_{u' \in \mathcal{E}_{r'}} \mathbf{w}_{r'u'}^H \left(\hat{\mathbf{h}}_{r'u}\hat{\mathbf{h}}_{r'u}^H + \boldsymbol{\Theta}_{r'u}\right) \mathbf{w}_{r'u'} + \sigma_z^2}\right) \quad (11)$$

Setting the derivative of (11) to zero, we obtain the expected optimal formula for $\gamma_u$ in (9d), which means they are equivalent.

Hence, our new reformulated problem can be written as

$$\text{(P3)} \quad \max_{\mathcal{W}, \boldsymbol{\gamma}} \quad f_1(\mathcal{W}, \boldsymbol{\gamma}) \quad (12)$$
$$\text{s.t.} \quad \text{(9b) and (9c)}$$

Note that we are not writing the dual problem here, but rather we are introducing SINR auxiliary variables $\boldsymbol{\gamma}$ that act as a proxy to account for the changes of the other variables.

**Proposition 2.** Maximizing the second term in the objective function in (12) is equivalent to maximizing the resulting $|\mathcal{C}_u|$ terms if we expand the numerator, where $|\mathcal{C}_u|$ is the size of the serving cluster for user $u$, i.e., the number of possible serving RRHs. If we decouple these terms and reorganize them with respect to each RRH $r$, we can restructure (11) as

$$f_1(\mathcal{W}, \boldsymbol{\gamma}) = \sum_{u \in \mathcal{U}} \delta_u (\log(1 + \gamma_u) - \gamma_u) + \sum_{r \in \mathcal{B}} f_2(r; \mathcal{W}, \boldsymbol{\gamma}),$$
$$(13)$$

where for each RRH $r$ we have

$$f_2(r; \mathcal{W}, \boldsymbol{\gamma}) =$$

$$\sum_{u \in \mathcal{E}_r} \delta_u \left( \frac{(1+\gamma_u) \left| \hat{\mathbf{h}}_{ru}^H \mathbf{w}_{ru} \right|^2}{\sum_{r' \in \mathcal{B}} \sum_{u' \in \mathcal{E}_{r'}} \mathbf{w}_{r'u'}^H \left( \hat{\mathbf{h}}_{r'u} \hat{\mathbf{h}}_{r'u}^H + \boldsymbol{\Theta}_{r'u} \right) \mathbf{w}_{r'u'} + \sigma_z^2} \right) \quad (14)$$

This restructuring follows from the fact that $\sum_{u \in \mathcal{U}} \left( \frac{a_u \sum_{r \in \mathcal{C}_u} A_{ru}}{B_u} \right) = \sum_{r \in \mathcal{B}} \sum_{u \in \mathcal{E}_r} \left( \frac{a_u A_{ru}}{B_u} \right)$, where each term in the summation in (14) is the fraction of the useful signal received at user $u$ from RRH $r$ over the total signals received at this user (including the useful signals).

Using fractional programming [11, Corollary 1] over (14), we can define the following function.

$$f_3(r; \mathcal{W}, \boldsymbol{\gamma}, \boldsymbol{\beta}_r) = \sum_{u \in \mathcal{E}_r} \left( 2\text{Re} \left\{ \beta_{ru}^* \sqrt{\delta_u (1+\gamma_u)} \mathbf{w}_{ru}^H \hat{\mathbf{h}}_{ru} \right\} \right.$$

$$\left. - |\beta_{ru}|^2 \left( \sum_{r' \in \mathcal{B}} \sum_{u' \in \mathcal{E}_{r'}} \mathbf{w}_{r'u'}^H \left( \hat{\mathbf{h}}_{r'u} \hat{\mathbf{h}}_{r'u}^H + \boldsymbol{\Theta}_{r'u} \right) \mathbf{w}_{r'u'} + \sigma_z^2 \right) \right) \quad (15)$$

where vector $\boldsymbol{\beta}_r \in \mathbb{C}^{|\mathcal{E}_r|}$ is introduced as a new auxiliary variable, and $\text{Re}\{\cdot\}$ is the real part. The function (15) is concave in $\boldsymbol{\beta}_r$. Also, it can be shown to be equivalent to (14) in the same way as was done with (11), i.e., by setting the partial derivative with respect to $\beta_{ru}^*$ to zero, then substituting the value of $\beta_{ru}$ in (15) which yields (14).

Then, our objective function in (12) can be written as

$$f_4(\mathcal{W}, \boldsymbol{\gamma}, \boldsymbol{\beta}) = \sum_{u \in \mathcal{U}} \delta_u \left( \log(1+\gamma_u) - \gamma_u \right) + \sum_{r \in \mathcal{B}} f_3(r; \mathcal{W}, \boldsymbol{\gamma}, \boldsymbol{\beta}_r), \quad (16)$$

where $\boldsymbol{\beta} = \left[ (\boldsymbol{\beta}_1)^T \dots (\boldsymbol{\beta}_{|\mathcal{B}|})^T \right]^T$ is the concatenation of the auxiliary variables $\boldsymbol{\beta}_r \in \mathbb{C}^{|\mathcal{E}_r|}$ introduced in (15) for each RRH $r$. $\qquad \square$

## IV. RESOURCE ALLOCATION

### A. Optimal Expressions

When the variables other than $\boldsymbol{\beta}_r$ are fixed, the optimal value of the auxiliary variable $\beta_{ru}$ can be obtained from its corresponding first-order optimality condition from (16) as

$$\beta_{ru} = \frac{\sqrt{\delta_u (1+\gamma_u)} \mathbf{w}_{ru}^H \hat{\mathbf{h}}_{ru}}{\sum_{r' \in \mathcal{B}} \sum_{u' \in \mathcal{E}_{r'}} \mathbf{w}_{r'u'}^H \left( \hat{\mathbf{h}}_{r'u} \hat{\mathbf{h}}_{r'u}^H + \boldsymbol{\Theta}_{r'u} \right) \mathbf{w}_{r'u'} + \sigma_z^2} \quad (17)$$

Similarly for the beamformers $\mathbf{w}_{ru}$, we can write the Lagrangian formulation using the new objective function (16) and the constraints in (12), then evaluating the corresponding first-order optimality condition to write $\mathbf{w}_{ru}$ as

$$\mathbf{w}_{ru} = \sqrt{\delta_u (1+\gamma_u)} \beta_{ru}^* \left( \sum_{r' \in \mathcal{B}} \sum_{u' \in \mathcal{E}_{r'}} |\beta_{r'u'}|^2 \left( \hat{\mathbf{h}}_{ru'} \hat{\mathbf{h}}_{ru'}^H + \boldsymbol{\Theta}_{ru'} \right) \right.$$

$$\left. + (\mu_r + \lambda_r \alpha_{ru}) \mathbf{I}_M \right)^{-1} \mathbf{h}_{ru} \quad (18)$$

where the Lagrangian multipliers $\lambda_r \geq 0$ and $\mu_r \geq 0$ correspond to the capacity (9b) and power (9c) constraints. Importantly, both these constraints relate to the power used at RRH $r$, i.e., both cannot be tight simultaneously. From complementary slackness, therefore, one of these Lagrange multipliers, both corresponding to RRH $r$, must be zero.

Unfortunately, we do not know a priori which constraint will remain tight. As we will see in our algorithm section, we propose a heuristic that, at each iteration of the algorithm, checks for whether the capacity constraint is satisfied (allowing $\lambda_r = 0$); if it is not satisfied, we update set $\lambda_r$ to a small value and update $\mu_r$ using a bisection search to meet the power constraint. Our results show that after a few iterations, $\lambda_r$ always converges to zero; we will comment on this in the results section.

### B. Optimization Algorithm

---
**Algorithm 1:** User-scheduling and resource allocation
---
1 Initialize $\mathcal{W}$ and weights $\alpha_{ru}$ for *all* users.
2 **while** *NOT converged* **do**
3 $\quad$ Update $\boldsymbol{\gamma}$ using (9d).
4 $\quad$ Update $\boldsymbol{\beta}$ using (17).
5 $\quad$ Update $\mathcal{W}$ using (18).
6 $\quad$ Update $\{\mu_r, \lambda_r : r \in \mathcal{B}\}$ as described using complementary slackness.
7 $\quad$ Update weights $\boldsymbol{\alpha}$ using (8).
8 **end**
---

We construct Algorithm 1 to allocate the resources for the users. The algorithm initializes some variables (Step 1) (e.g., conjugate beamforming to initialize $\mathbf{w}_{ru}$). Then, it updates the variables $\boldsymbol{\gamma}$, $\boldsymbol{\beta}$, $\mathcal{W}$, and $\boldsymbol{\alpha}$ iteratively one at a time until convergence.

The complexity of updating $\boldsymbol{\gamma}$, $\boldsymbol{\beta}$, and $\boldsymbol{\alpha}$ is $\mathcal{O}(|\mathcal{U}|)$, $\mathcal{O}(|\mathcal{U}|C_{\text{avg}})$, and $\mathcal{O}(|\mathcal{U}|C_{\text{avg}})$, respectively, where $C_{\text{avg}}$ is the average cluster size per user, and it is affected by both the density of the active users and the large scale fading threshold $\rho$. The complexity of the beamforming using a weighted MMSE [18] is $\mathcal{O}(|\mathcal{U}_s|^2 M^2 + |\mathcal{U}_s| M^3)$, where $\mathcal{U}_s$ is the set of scheduled users. Hence, leading to a total algorithm complexity of at most $\mathcal{O}(M^3|\mathcal{B}|^2 + M^4|\mathcal{B}| + |\mathcal{U}|C_{\text{avg}})$, where the number of the scheduled users is at most $|\mathcal{U}_s| \leq M|\mathcal{B}|$.

## V. NUMERICAL RESULTS AND ANALYSIS

To eliminate network borders, we consider a wrap-around structure consisting of $Q = 7$ hexagonal *virtual* cells[2] each having an inner radius $500$ m and containing $N$ RRHs that are uniformly distributed in each *virtual* cell. Users are randomly distributed with a density $\lambda_{\text{users}}$ and a circular exclusion region of $20$ m around each RRH. We average our results using Monte Carlo simulations over both network realizations and time slots (TSs), and we include the effect of the users fairness

---
[2]We create these virtual cells to allow for wrap-around; the cells have no physical meaning.

TABLE I: Simulation parameters.

| | *Parameter* | *Value* |
|---|---|---|
| Cell config. | $Q$, $N$, $M$, $\lambda_{\text{users}}$ | 7, 10, 8, 200 users/km$^2$ |
| Power, Imperfect CSI | $p$, $\tau_d$, $(\tau_p)$, $p_u$ | 30 dBm, 200, (16, 32, 64), 20 dBm |
| Noise spectral density, Noise figure | $S_z$, $F_z$, Bandwidth | $-174$dBm/Hz, 8 dBm, 180 KHz |
| Others | $\sigma_{\text{shadowing}}$, $\rho$, $\eta$, $\epsilon$ | 4 dB, $L(0.4)$, 0.2, $\frac{0.9p}{M}$ |



(a) Allocated power for the users' beamformers on a typical RRH.

(b) Convergence plot for many channel realizations.

Fig. 2: Evolution of the algorithm.

by simulating 100 TSs and averaging the results over the last allocated 50 TSs, representing steady state performance[3].

We use the COST231 Walfish-Ikegami [19] to model the path loss at 1800 MHz, resulting in $L_{|\text{dB}}(d_{ru}) = -112.4271 - 38\log_{10}(d_{ru})$ where $d_{ru}$ is in km. In Table I, we summarize the parameters used.

The proportional fairness weight, $\delta_u$, for user $u$ is the inverse of the achieved long-term average rate over an exponentially decaying window; in time slot $t$ we set $\delta_u$ as [20]

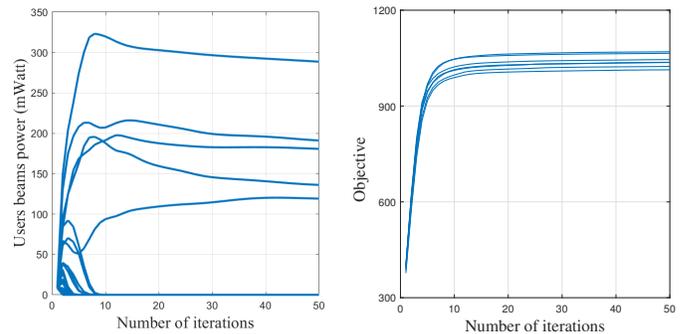$$\delta_u^{(t)} = \frac{1}{\bar{R}_u^{(t)}}, \tag{19}$$

where $\delta_u^{(t)}$ is the value of $\delta_u$ at time slot $t$, and $\bar{R}_u^{(t)}$ is the user exponentially weighted rate averaged over previous time slots, and it is updated as $\bar{R}_u^{(t)} = \eta R_u^{(t)} + (1-\eta)\bar{R}_u^{(t-1)}$ with a forgetting factor $0 \le \eta \le 1$, where $R_u^{(t)}$ is the rate achieved by user $u$ at time $t$, and it can be defined as [13]

$$R_u^{(t)} = \frac{(\tau_d - \tau_p)}{\tau_d}\log\left(1 + \frac{\sum_{r \in \mathcal{C}_u} s_{ru}\left|\mathbf{h}_{ru}^H\mathbf{w}_{ru}\right|^2}{\sum_{r' \in \mathcal{B}}\sum_{u' \in \mathcal{E}_{r'}, u' \neq u} s_{r'u'}\left|\mathbf{h}_{r'u}^H\mathbf{w}_{r'u'}\right|^2 + \sigma_z^2}\right) \tag{20}$$

In Fig. 2(a), we plot the evolution of the allocated power of the beamformer's weights for a typical RRH in a typical network as a function of the algorithm iterations. It is clear that after a few iterations, the power constraint (9c) is tighter than the capacity constraint (9b) which becomes deactivated, as previously discussed. In Fig. 2(b), we illustrate the convergence of the algorithm for several channel realizations.

In Fig. 3, we plot the long-term performance results under ideal CSI, i.e., the channels are known and there is no pilot training overhead. Fig. 3(a) shows the long-term network sum of spectral efficiency (SE) as a function of the algorithm iterations. The evolution of the curve shows that the algorithm converges smoothly with a non-decreasing fashion. Also, the results show a huge performance gain from using our approach compared to the ZF and the conjugate beamforming schemes with a round-robin scheduling. Compared to these schemes respectively, we obtain about a $9.1$-fold and $10.6$-fold improvement. Additionally, we plot the resulting network sum SE when using the ZF beamforming scheme with the optimized user-scheduling obtained from our proposed approach. The
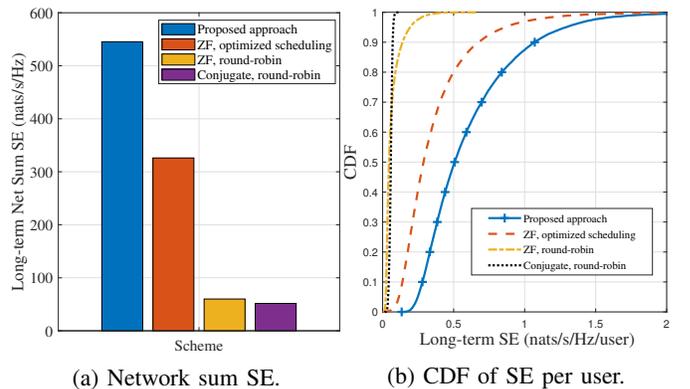
---

[3]We emphasize that plotting the results from allocating a single TS would definitely give much higher performance, because the users with the best channel's conditions would be served, i.e., fairness is equal for all the users. Nonetheless, we are interested in studying the effect of the scheme on the long-term.

results still show a $1.67$-fold improvement. Clearly, this gap is due to the fact that the ZF beamformers are constructed at each RRH using only the channels of the served users, and each user is being allocated equal power irrespective of the channel conditions. Moreover, to quantify the effect of optimized user-scheduling, we compare the round-robin scheduling with that of the optimized one for the ZF beamforming. The result highlights the importance of optimized scheduling, where a $5.44$-fold improvement is achieved.



(a) Network sum SE.

(b) CDF of SE per user.

Fig. 3: Long-term results, $N = 10$.

In Fig. 3(b), we plot the cumulative density function (CDF) of the long-term SE of the users, where an $8$ to $10$-fold gain is observed in the median long-term user spectral efficiency for our approach compared to round-robin scheduling. We emphasize that this plot presents the long-term average rate that accounts for the user scheduling. Users may not be scheduled in every time slot; this is determined by their channels and their weights as defined in (19). The gains for the $10^{th}$-percentile rate, is clear. The proposed approach results in about $7$-fold and $2$-fold improvement in the cell-edge long-term rate compared to round-robin scheduling and ZF beamforming with optimized scheduling respectively.

To quantify the performance of imperfect CSI, we compare the following cases:

- PI: Our *proposed* approach using *ideal* channels, where no channel estimation phase is accounted for.

- PEAR: Our *proposed* approach when the algorithm is using the *estimated* channel and using *robust* beamforming, i.e., accounting for the estimation error. However, when plotting the results, we plot the *actual* network performance, i.e., using (20).

Since we use user-centric clustering, we define an *area-based* pilot-reuse factor (not cell-based) as $\xi_p \triangleq \tau_p/\lambda_{\text{users}}$. For example, $\xi_p = 0.25$, means that on-average one-quarter of the users found in an area of $1 \times 1$ km$^2$ are using orthogonal pilots. Under the user density specified in Table I, the pilot sequence lengths $\tau_p = 64,\ 32,\ 16$ produce on-average $\xi_p = 0.32,\ 0.16,\ 0.08$ respectively.

In Fig. 4, we plot the long-term network sum SE of the different studied cases using, but this time using $N = 5$ RRHs per virtual cell. The results show a drop of performance for PEAR by 39.21, 37, and 43.21 percent compared to the ideal channel case (PI). If we are to quantify the performance drop due to only imperfect CSI, we obtain 10.6, 25, and 38.27 percent drop in the performance compared to the ideal case. From the results, using $\tau_p = 32$ provides the highest sum SE, i.e., it is a good compromise between the pilot contamination and the pilot-training overhead.
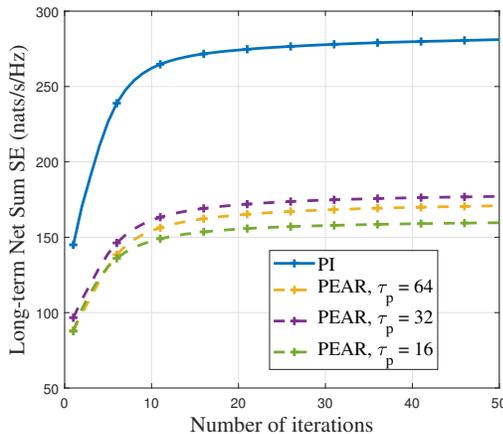


Fig. 4: Evolution of the Long-term sum of SE, $N = 5$.

## VI. Conclusion

This paper optimized user-scheduling and resource allocation in a distributed cell-free MIMO system under the user-centric clustering scheme and the non-coherent transmission mode using a weighted sum rate problem formulation. We used tools from block coordinate descent, fractional programming, and compressive sensing to provide closed-form expressions for the optimized variables, while keeping the other variables fixed. This allowed us to construct an iterative optimization algorithm that converges smoothly in non-decreasing fashion. Our key contribution is optimized user-scheduling, which is neglected in most of the literature. The numerical results show that our optimized resource allocation boosts network performance, both in terms of sum-rate and long-term proportional fair rates, compared to conventional round-robin schemes, where an 8 to 10-fold gain in the long-term user spectral efficiency is observed.

## References

[1] H. A. Ammar and R. Adve, "Power delay profile in coordinated distributed networks: User-centric v/s disjoint clustering," in *2019 IEEE Global Conf. on Signal and Inf. Processing*, pp. 1–5, Nov 2019.

[2] H. Q. Ngo, A. Ashikhmin, H. Yang, E. G. Larsson, and T. L. Marzetta, "Cell-free massive MIMO versus small cells," *IEEE Transactions on Wireless Communications*, vol. 16, no. 3, pp. 1834–1850, 2017.

[3] S. Buzzi, C. D'Andrea, A. Zappone, and C. D'Elia, "User-centric 5G cellular networks: Resource allocation and comparison with the cell-free massive MIMO approach," *IEEE Transactions on Wireless Communications*, vol. 19, pp. 1250–1264, Feb 2020.

[4] I. Atzeni, B. Gouda, and A. Tölli, "Distributed precoding design via over-the-air signaling for cell-free massive MIMO," *IEEE Transactions on Wireless Communications*, pp. 1–1, 2020.

[5] B. Dai and W. Yu, "Sparse beamforming and user-centric clustering for downlink cloud radio access network," *IEEE Access*, vol. 2, pp. 1326–1339, 2014.

[6] M. Bashar, K. Cumanan, A. G. Burr, H. Q. Ngo, M. Debbah, and P. Xiao, "Max–min rate of cell-free massive MIMO uplink with optimal uniform quantization," *IEEE Transactions on Communications*, vol. 67, no. 10, pp. 6796–6815, 2019.

[7] E. Nayebi, A. Ashikhmin, T. L. Marzetta, H. Yang, and B. D. Rao, "Precoding and power optimization in cell-free massive MIMO systems," *IEEE Trans. on Wireless Comm.*, vol. 16, pp. 4445–4459, July 2017.

[8] Q. D. Vu, L. N. Tran, and M. Juntti, "Noncoherent joint transmission beamforming for dense small cell networks: Global optimality, efficient solution and distributed implementation," *IEEE Transactions on Wireless Communications*, vol. 19, no. 9, pp. 5891–5907, 2020.

[9] D. W. K. Ng, E. S. Lo, and R. Schober, "Dynamic resource allocation in MIMO-OFDMA systems with full-duplex and hybrid relaying," *IEEE Transactions on Communications*, vol. 60, no. 5, pp. 1291–1304, 2012.

[10] W. Yu and R. Lui, "Dual methods for nonconvex spectrum optimization of multicarrier systems," *IEEE Transactions on Communications*, vol. 54, no. 7, pp. 1310–1322, 2006.

[11] K. Shen and W. Yu, "Fractional programming for communication systems—part II: Uplink scheduling via matching," *IEEE Transactions on Signal Processing*, vol. 66, pp. 2631–2644, May 2018.

[12] A. A. Khan, R. S. Adve, and W. Yu, "Optimizing downlink resource allocation in multiuser MIMO networks via fractional programming and the hungarian algorithm," *IEEE Transactions on Wireless Communications*, vol. 19, no. 8, pp. 5162–5175, 2020.

[13] C. Pan, H. Ren, M. Elkashlan, A. Nallanathan, and L. Hanzo, "The non-coherent ultra-dense C-RAN is capable of outperforming its coherent counterpart at a limited fronthaul capacity," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 11, pp. 2549–2560, 2018.

[14] P. Frenger, J. Hederen, M. Hessler, and G. Interdonato, "Antenna arrangement for distributed massive MIMO," Nov. 28 2019. US Patent App. 16/435,054.

[15] S. M. Kay, *Fundamentals of Statistical Signal Processing, vol. 1: Estimation Theory*. Prentice Hall PTR, 1993.

[16] M. G. Karypis, V. Kumar, and M. Steinbach, "A comparison of document clustering techniques," in *TextMining Workshop at KDD2000 (May 2000)*, 2000.

[17] E. J. Candes, M. B. Wakin, and S. P. Boyd, "Enhancing sparsity by reweighted $\ell_1$ minimization," *Journal of Fourier analysis and applications*, vol. 14, no. 5-6, pp. 877–905, 2008.

[18] Q. Shi, M. Razaviyayn, Z. Luo, and C. He, "An iteratively weighted MMSE approach to distributed sum-utility maximization for a MIMO interfering broadcast channel," *IEEE Transactions on Signal Processing*, vol. 59, pp. 4331–4340, Sep. 2011.

[19] J. Walfisch and H. L. Bertoni, "A theoretical model of UHF propagation in urban environments," *IEEE Transactions on Antennas and Propagation*, vol. 36, pp. 1788–1796, Dec 1988.

[20] W. Yu, T. Kwon, and C. Shin, *Adaptive resource allocation in cooperative cellular networks*, ch. 9, pp. 233–256. Cambridge Univ. P., 2011.