

A New Look at Cell-Free Massive MIMO: Making It Practical With Dynamic Cooperation

Emil Björnson*, Luca Sanguinetti†

*Department of Electrical Engineering (ISY), Linköping University, Linköping, Sweden (emil.bjornson@liu.se)

†Dipartimento di Ingegneria dell'Informazione, University of Pisa, 56122 Pisa, Italy (luca.sanguinetti@unipi.it)

Abstract—This paper takes a new look at Cell-free Massive MIMO (multiple-input multiple-output) through the lens of the *dynamic cooperation cluster* framework from the Network MIMO literature. The purpose is to identify and address scalability issues that appear in prior work. We provide distributed algorithms for initial access, pilot assignment, cluster formation, precoding, and combining that are scalable in the sense of being implementable with arbitrarily many users. Interestingly, the suggested precoding and combining outperform conjugate beamforming and matched filtering, respectively, while also being fully distributed.

Index Terms—Cell-free Massive MIMO, dynamic cooperation clustering, scalability, combining and precoding.

I. INTRODUCTION

By transmitting a signal coherently from multiple antennas, the received power can be increased without increasing the total transmit power [1]. This is the phenomenon utilized by classic beamforming from co-located antenna arrays but can be also utilized when transmitting coherently from multiple access points (APs) [2]. Even if the APs have different channel gains to the receiver, the benefit of coherent transmission makes it better to spread out the transmit power over multiple APs than transmitting only from the AP with the best channel [3]. Such coherent joint transmission from multiple APs has many different names, including Network MIMO [4].

The early Network MIMO papers assumed all APs have network-wide channel state information (CSI) and transmit to all user equipments (UEs). These are two preferable but impractical/unscalable assumptions that lead to immense backhaul signaling for CSI and data sharing, respectively. Fortunately, [5] proved that Network MIMO can operate without CSI sharing, by sacrificing the ability for APs to jointly cancel interference. Moreover, to limit data sharing, each UE can be served only by a subset of the APs. Initially, a *network-centric* approach was taken by dividing the APs into non-overlapping cooperation clusters in which the APs are sharing data to serve only UEs residing in the joint coverage area. This approach was considered in LTE but provides small gains in practice [6], partially due to substantial interference between clusters. The alternative is *dynamic cooperation clusters* (DCC) [7], which is a *user-centric* approach where each UE is served by the AP subset providing the best channel conditions. DCC didn't gain much attention at the time, since Massive MIMO (mMIMO)

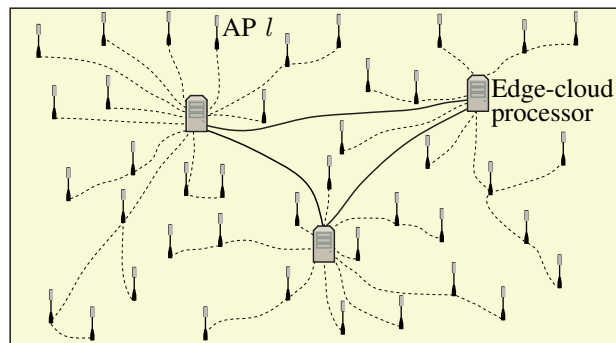


Fig. 1: Illustration of a Cell-free mMIMO network with many distributed APs connected to edge-cloud processors. The APs are jointly serving all the UEs in the coverage area.

was simultaneously proposed [8] and rightfully gained the spotlight, but it was implemented in the pCell technology [9].

Now that mMIMO is a rather mature technology [1], the research focus is shifting back to Network MIMO, but under the new name of *Cell-free mMIMO* [3], [10]. The key novelty is the rigorous ergodic spectral efficiency (SE) analysis with imperfect CSI, but conceptually, it is a special case of Network MIMO. In fact, it was initially a step backward in terms of implementation feasibility since all APs were assumed to serve all UEs and emphasis was put on developing network-wide power control algorithms [3], [10]–[12]. The user-centric approach was reintroduced for Cell-free mMIMO in [13] but without making connections to DCC or other implementation-related aspects that had already been considered in the Network MIMO literature and summarized in the textbook [14].

Contributions: In this paper, we first expose the potential scalability issues of Cell-free mMIMO and then prove that Cell-free mMIMO is a special case of the DCC framework in [7], [14]. We utilize this perspective to present new distributed and scalable algorithms for initial access, pilot assignment, and cooperation cluster formation. We derive downlink and uplink SEs with multi-antenna APs and propose new scalable forms of signal-to-leakage-and-noise ratio (SLNR) precoding and regularized zero-forcing (RZF) combining. These methods are fully distributed and, importantly, outperform the standard conjugate beamforming and matched filtering methods.

II. SYSTEM MODEL AND SCALABILITY

We consider a cell-free network consisting of K single-antenna UEs and L APs, each equipped with N antennas. The

E. Björnson was supported by ELLIIT and the Wallenberg AI, Autonomous Systems and Software Program (WASP). L. Sanguinetti was supported by the University of Pisa under the PRA 2018-2019 Research Project CONCEPT.

APs are connected to edge-cloud processors [9], [15], [16], as illustrated in Fig. 1, which enables coherent joint transmission and reception to the UEs in the entire coverage area.

The channel between AP l and UE k is denoted $\mathbf{h}_{kl} \in \mathbb{C}^N$ and the collective channel from all APs is $\mathbf{h}_k = [\mathbf{h}_{k1}^T \dots \mathbf{h}_{kL}^T]^T \in \mathbb{C}^M$, where $M = NL$. The network operates according to a time-division duplex (TDD) protocol with a data transmission phase and a pilot phase for channel estimation. We consider the standard TDD protocol [1] in which each coherence block is divided into τ_p channel uses for uplink pilots, τ_u for uplink data, and τ_d for downlink data with $\tau_c = \tau_p + \tau_u + \tau_d$. In each block, an independent flat-fading realization is drawn using correlated Rayleigh fading:

$$\mathbf{h}_{kl} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{R}_{kl}) \quad (1)$$

where the spatial correlation matrix $\mathbf{R}_{kl} \in \mathbb{C}^{N \times N}$ describes geometric attenuation, shadowing, and spatial properties.

A. Original Cell-free mMIMO Model

To motivate the approach taken in this paper, we first review the system model in the original papers on Cell-free mMIMO [3], [10], where network-wide downlink transmission from all APs to all the UEs is considered. Let $\mathbf{w}_{il} \in \mathbb{C}^N$ denote the precoding vector that AP l assigns to UE k , then the received downlink signal at UE k is

$$y_k^{\text{dl}} = \sum_{l=1}^L \sum_{i=1}^K \mathbf{h}_{kl}^T \mathbf{w}_{il} s_i + n_k = \sum_{i=1}^K \mathbf{h}_k^T \mathbf{w}_i s_i + n_k \quad (2)$$

where $s_i \in \mathbb{C}$ is the independent unit-power data signal intended for UE i , $\mathbf{w}_k = [\mathbf{w}_{k1}^T \dots \mathbf{w}_{kL}^T]^T \in \mathbb{C}^M$ is the collective precoding vector, and $n_k \sim \mathcal{N}_{\mathbb{C}}(0, \sigma^2)$ is the receiver noise.

The collective channel is distributed as $\mathbf{h}_k \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{R}_k)$ where $\mathbf{R}_k = \text{diag}(\mathbf{R}_{k1}, \dots, \mathbf{R}_{kL}) \in \mathbb{C}^{M \times M}$ is the block-diagonal spatial correlation matrix. The system model (2) is equivalent to a single-cell downlink mMIMO system with correlated fading. The achievable SEs in Cell-free mMIMO, thus, follow easily from the literature on mMIMO with correlated fading, recently summarized in [1]. The key difference from that literature is which precoding vectors can be selected, since these should satisfy per-AP power constraints and (preferably) use only local CSI. Network-wide downlink power optimization methods were developed in [3], [10], among others.

Similarly, during uplink data transmission, the received signal $\mathbf{y}_l^{\text{ul}} \in \mathbb{C}^N$ at AP l is

$$\mathbf{y}_l^{\text{ul}} = \sum_{i=1}^K \mathbf{h}_{il} s_i + \mathbf{n}_l \quad (3)$$

where $s_i \in \mathbb{C}$ is the signal transmitted from UE i with power p_i and $\mathbf{n}_l \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$. Network-wide uplink decoding was considered in the original papers on Cell-free mMIMO [3], [11]. In that case, AP l selects a receive combining vector \mathbf{v}_{kl} for UE k and computes $\mathbf{v}_{kl}^H \mathbf{y}_l^{\text{ul}}$ locally. The network then estimates s_k by computing the summation

$$\hat{s}_k = \sum_{l=1}^L \mathbf{v}_{kl}^H \mathbf{y}_l^{\text{ul}}. \quad (4)$$

Note that (4) is equivalent to an uplink single-cell mMIMO system model with correlated fading, thus the achievable SEs easily follow from that literature [1]. The difference is which combining vectors can be used, since these should (preferably) use only local CSI. Network-wide uplink power optimization methods were developed in [3], [11], [12], among others.

B. Scalability Issues

Although the network-wide processing in the original Cell-free mMIMO papers is appealing, it is not practical for large-scale network deployments with many UEs. To determine if the processing is scalable or not, it is helpful to let $K \rightarrow \infty$ and see which of the following operations are implementable.

- 1) Precoding and combining: AP l computes K precoding vectors (\mathbf{w}_{lk} for all k) and K combining vectors (\mathbf{v}_{lk} for all k). The complexity becomes infinite as $K \rightarrow \infty$.
- 2) Estimation: AP l must compute channel estimates for all K UEs, with infinite complexity as $K \rightarrow \infty$.
- 3) Fronthaul signaling: AP l needs to receive K downlink data signals over the fronthaul network and forward K received signals $\mathbf{v}_{kl}^H \mathbf{y}_l^{\text{dl}}$ over the fronthaul network.
- 4) Power optimization: Any network-wide power optimization has a complexity that goes to infinity as $K \rightarrow \infty$.

The original form of Cell-free mMIMO is clearly not scalable.

Definition 1. A Cell-free mMIMO network is said to be *scalable* if none of four above-listed issues appears.

In the remainder of this paper, we outline a scalable implementation framework according to Definition 1. We start from the DCC framework for Network MIMO in [7], [14], which was claimed to be scalable but we fill in many missing details.

C. Dynamic Cooperation Clusters

The DCC framework was proposed in [7], [14] to enable “*unified analysis of anything from interference channels to ideal network MIMO*”. To this end, the diagonal matrix $\mathbf{D}_{il} \in \mathbb{C}^{N \times N}$ was defined, where the j th diagonal element is 1 if the j th antenna of AP l is allowed to transmit to and decode signals from UE i and 0 otherwise. By modifying (2), the received downlink signal at UE k becomes

$$y_k^{\text{dl}} = \sum_{l=1}^L \sum_{i=1}^K \mathbf{h}_{kl}^T \mathbf{D}_{il} \mathbf{w}_{il} s_i + n_k = \sum_{i=1}^K \mathbf{h}_k^T \mathbf{D}_i \mathbf{w}_i s_i + n_k \quad (5)$$

where $\mathbf{D}_i = \text{diag}(\mathbf{D}_{i1}, \dots, \mathbf{D}_{iL}) \in \mathbb{C}^{M \times M}$ is block-diagonal. By selecting $\mathbf{D}_1, \dots, \mathbf{D}_K$ in different ways, (5) can be used to model many different types of multi-AP networks; see [14].

The original Cell-free mMIMO in (2) is obtained from (5) in the special case of $\mathbf{D}_i = \mathbf{I}_M \forall i$, where all antennas serve all UEs. The user-centric approach to Cell-free mMIMO described in [13] is also an instance of the DCC framework. In [13], $\mathcal{M}(k) \subset \{1, \dots, L\}$ denotes the subset of APs that communicate with UE k , which corresponds to setting

$$\mathbf{D}_{kl} = \begin{cases} \mathbf{I}_N & \text{if } l \in \mathcal{M}(k), \\ \mathbf{0}_N & \text{if } l \notin \mathcal{M}(k). \end{cases} \quad (6)$$

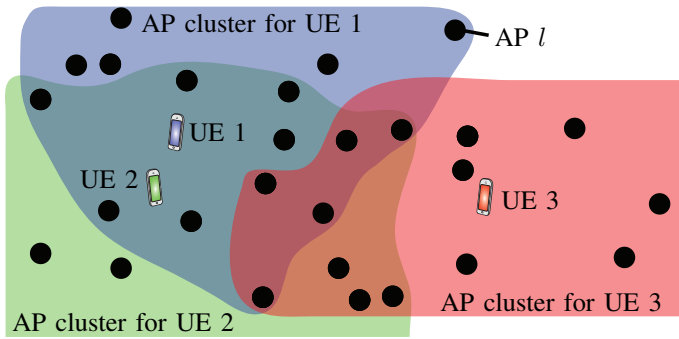


Fig. 2: Example of dynamic cooperation clusters for three UEs.

This is exactly the same setup as considered in [7].

The DCC framework does not change the received uplink signal in (3), but the uplink data estimate in (4) changes to

$$\hat{s}_k = \sum_{l=1}^L \mathbf{v}_{kl}^H \mathbf{D}_{kl} \mathbf{y}_l^{\text{ul}} = \sum_{l \in \mathcal{M}(k)} \mathbf{v}_{kl}^H \mathbf{y}_l^{\text{ul}} \quad (7)$$

where the second equality only holds when using (6).

Fig. 2 illustrates a network with three UEs that are served by a large number of APs. The colored regions illustrate which clusters of APs are transmitting to which UEs. The fact that the clusters are partially overlapping is a core feature of DCCs, and also demonstrates that this is a cell-free network.

The DCC framework was proposed in [7] to achieve scalability in Network MIMO, for example, in the following way.

Lemma 1. The UEs served by AP l have indices in the set

$$\mathcal{D}_l = \{i : \text{tr}(\mathbf{D}_{il}) \geq 1, i \in \{1, \dots, K\}\}. \quad (8)$$

If the cardinality $|\mathcal{D}_l|$ is constant as $K \rightarrow \infty$, the precoding/combining complexity and fronthaul signaling parts of Definition 1 are satisfied and thus scalable.

Proof: AP l only needs to compute precoding and combining vectors for $|\mathcal{D}_l|$ UEs, and it only needs to send/receive data related to these UEs over the fronthaul network. ■

With this result in mind, the practically important question is how to select the sets $\{\mathcal{D}_l : \forall l\}$ in a scalable way, while guaranteeing service to all UEs. This challenge is tackled below.

III. DISTRIBUTED AND SCALABLE IMPLEMENTATION

In this section, we propose a scalable, distributed implementation of Cell-free mMIMO. It is inspired by the guidelines for distributed Network MIMO in [14, Sec. 4.3, 4.7], but is far more detailed and also focused on resource allocation. We begin by making the following key assumption.

Assumption 1. Each AP serves at most one UE per pilot and uses all its antennas to serve these UEs. This implies $|\mathcal{D}_l| \leq \tau_p$ and $\mathbf{D}_{il} = \mathbf{I}_N$ for all $i \in \mathcal{D}_l$, $l = 1, \dots, L$. Hence, the scalability requirement in Lemma 1 is satisfied.

The rationale for Assumption 1 is: a) pilot contamination makes the channel estimates of pilot-sharing UEs similar (or

identical) so the AP will cause strong interference if it transmits to more than one such UE; b) the signal processing complexity becomes fixed and scalable, although all N antennas are used; c) the fronthaul capacity can be dimensioned to support τ_p parallel uplink/downlink data signals per AP.

A. Pilot Transmission and Channel Estimation

There are τ_p mutually orthogonal τ_p -length pilot signals that are assigned to the UEs. Note that the pilot transmission protocol is scalable since τ_p is a constant, independent of K . An algorithm for pilot assignment is provided in Section III-B, but for now we let $\mathcal{S}_t \subset \{1, \dots, K\}$ denote the subset of UEs assigned to pilot t . When these UEs transmit their pilot, the received pilot signal $\mathbf{y}_{tl}^p \in \mathbb{C}^N$ at AP l is

$$\mathbf{y}_{tl}^p = \sum_{i \in \mathcal{S}_t} \sqrt{\tau_p p_i} \mathbf{h}_{il} + \mathbf{n}_{tl} \quad (9)$$

where p_i is the transmit power, τ_p is the processing gain, and $\mathbf{n}_{tl} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$ is noise. Using standard results [1, Sec. 3], the minimum mean-squared error (MMSE) estimate of \mathbf{h}_{kl} for $k \in \mathcal{S}_t$ is

$$\hat{\mathbf{h}}_{kl} = \sqrt{p_k \tau_p} \mathbf{R}_{kl} \mathbf{\Psi}_{tl}^{-1} \mathbf{y}_{tl}^p \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, p_k \tau_p \mathbf{R}_{kl} \mathbf{\Psi}_{tl}^{-1} \mathbf{R}_{kl}) \quad (10)$$

where the correlation matrix $\mathbf{\Psi}_{tl} = \mathbb{E}\{\mathbf{y}_{tl}^p (\mathbf{y}_{tl}^p)^H\}$ is given by

$$\mathbf{\Psi}_{tl} = \sum_{i \in \mathcal{S}_t} \tau_p p_i \mathbf{R}_{il} + \sigma^2 \mathbf{I}_N. \quad (11)$$

The AP uses these estimates for receiving the uplink data and for precoding the downlink data. AP l only needs to compute estimates $\hat{\mathbf{h}}_{kl}$ for $k \in \mathcal{D}_l$, which under Assumption 1 is at most one UE per pilot. Since the complexity per AP is independent of K , the pilot transmission is scalable when $K \rightarrow \infty$.

B. Initial Access and Pilot Assignment

When a new UE wants to access the network, it needs to be assigned a pilot and make it into the set \mathcal{D}_l of at least one AP. This must be done in a distributed fashion, which has the risk that the UE is inadvertently dropped from service since no AP decides to transmit to it. To avoid that, each UE appoints a *Master AP* that is required to transmit to it and coordinate the decoding of the uplink data [14]. Let $K+1$ be the index of the connecting UE, then the proposed access procedure is:

- 1) The UE measures $\beta_l = \text{tr}(\mathbf{R}_{(K+1)l})/N$ for all nearby APs, using periodically broadcasted synchronization signals, and appoints $\text{AP } \ell = \arg \max_l \beta_l$ as its Master AP. The UE also uses this signal to synchronize to the AP.
- 2) The UE contacts its Master AP via a standard random access procedure. The AP responds by assigning pilot $\tau = \arg \min_t \text{tr}(\mathbf{\Psi}_{tl})$ to the UE, with $\mathbf{\Psi}_{tl}$ given in (11).
- 3) The Master AP informs a limited set of neighboring APs that it is now serving UE $K+1$ on pilot τ . These APs independently decide if they will also serve the UE.

In summary, the UE appoints the AP with the strongest channel as its Master AP and it is assigned to the pilot that this AP

observes the least pilot power on.¹ When other APs decide whether to also serve the new UE, Assumption 1 must be enforced. To limit interference (and pilot contamination), it is reasonable for an AP to switch to serving the new UE if it has a better channel to it than to the UE it currently serves on that pilot *and* it is not the Master AP of the current UE. When a UE moves around, the proposed access procedure can be redone when needed; the UE then acts as if it is connecting and appoints a new Master AP, which might assign a new pilot. The old Master AP transfers its status to the new Master AP.

C. Downlink SE and Distributed Precoding

Next, we derive and analyze a general achievable downlink SE expression (i.e., a lower bound on the ergodic capacity) for the DCC system model in (5). AP l selects its precoding vectors \mathbf{w}_{kl} for $k \in \mathcal{D}_l$ as a function of $\{\hat{\mathbf{h}}_{kl} : k \in \mathcal{D}_l\}$, while $\mathbf{D}_{il}\mathbf{w}_{il} = \mathbf{0}$ for $i \notin \mathcal{D}_l$ in all expressions so these precoding vectors need not be selected. The precoding vectors of the UEs that the AP serves must also satisfy the power constraint

$$\sum_{k \in \mathcal{D}_l} \mathbb{E}\{\|\mathbf{w}_{kl}\|^2\} \leq \rho \quad (12)$$

where ρ is the total transmit power of an AP.

We use the *hardening bound* that is widely used in the mMIMO literature to compute SEs [1, Th. 4.6]; it has also been used in [3], [10] for Cell-free mMIMO with $N = 1$, $\mathbf{D}_i = \mathbf{I}_M \forall i$, for specific choices of precoding schemes.

Proposition 1. An achievable downlink SE R_k^{dl} [bit/s/Hz] for UE k is

$$\frac{\tau_d}{\tau_c} \log_2 \left(1 + \frac{|\mathbb{E}\{\mathbf{h}_k^T \mathbf{D}_k \mathbf{w}_k\}|^2}{\sum_{i=1}^K \mathbb{E}\{|\mathbf{h}_k^T \mathbf{D}_i \mathbf{w}_i|^2\} - |\mathbb{E}\{\mathbf{h}_k^T \mathbf{D}_k \mathbf{w}_k\}|^2 + \sigma^2} \right). \quad (13)$$

Proof: This is proved by following the same approach as in [1, Th. 4.6], but for the system model in (5). ■

The expectations in (13) can be computed by Monte-Carlo simulations for any choice of precoding vectors. The precoding at AP l should only depend on $\{\hat{\mathbf{h}}_{il} : i \in \mathcal{D}_l\}$ to achieve a scalable implementation [5]. Without loss of generality, we set

$$\mathbf{w}_{il} = \sqrt{\frac{\rho_{il}}{\mathbb{E}\{\|\bar{\mathbf{w}}_{il}\|^2\}}} \bar{\mathbf{w}}_{il} \quad \forall i \in \mathcal{D}_l \quad (14)$$

where $\rho_{il} \geq 0$ is the transmit power and $\bar{\mathbf{w}}_{il} \in \mathbb{C}^N$ gives the precoding direction. Two schemes that satisfy the scalability requirement are maximum ratio (MR) and SLNR [5]:

$$\bar{\mathbf{w}}_{kl} = \begin{cases} \hat{\mathbf{h}}_{kl}^* & \text{with MR,} \\ \left(\sum_{i \in \mathcal{D}_l} \rho_{il} \hat{\mathbf{h}}_{il}^* \hat{\mathbf{h}}_{il}^T + \sigma^2 \mathbf{I}_N \right)^{-1} \hat{\mathbf{h}}_{kl}^* & \text{with SLNR,} \end{cases} \quad (15)$$

¹This should be a pilot on which the AP is not currently serving a UE as being its Master AP, since that role has higher priority. Each AP can only be the Master AP of up to τ_p UEs in the proposed framework, but in the unlikely event that this cannot be satisfied, multiple UEs can be assigned to the same pilot but multiplexed in time and/or frequency instead. Alternatively, the second strongest AP can be appointed the Master AP.

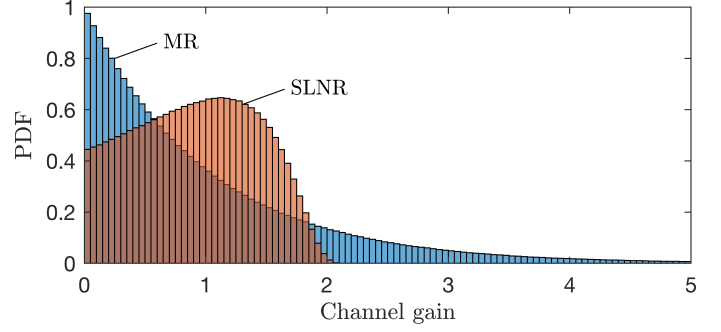


Fig. 3: For $N = M = K = \rho = \sigma^2 = 1$ and perfect CSI, the channel gain is $|h|^2$ with MR and $\frac{|h|^2}{(|h|^2+1)^2} / \mathbb{E}\{\frac{|h|^2}{(|h|^2+1)^2}\}$ with SLNR, where $x \sim \mathcal{N}_{\mathbb{C}}(0, 1)$. Their PDFs are widely different, particularly only SLNR has bounded support.

for $k \in \mathcal{D}_l$. MR is also known as conjugate beamforming and is the standard scheme in the Cell-free mMIMO literature. The scalable SLNR precoding in (15) is new since previous expressions consider all UEs in the network [14]. The benefit of SLNR over MR is two-fold: 1) it suppresses interference spatially if $N > 1$ since $\bar{\mathbf{w}}_{kl}$ maximizes the ratio between desired signal power and interference caused to the other UEs served by the AP; and 2) it reduces variations in the effective gains $\mathbf{h}_{kl}^T \mathbf{w}_{il}$ of desired and interfering channels for any N .

The latter is a non-trivial phenomenon that appears even with $N = M = K = 1$ and perfect CSI. Fig. 3 shows the probability density function (PDF) of the channel gains with MR and SLNR in that case. The channel gains have identical mean values, but MR gives an exponential distribution with an infinite tail while SLNR has small and compact support. This behavior will lead to higher SE when using SLNR.

The “only” benefit of MR is that the SE can be computed in closed form, following the same approach as in [1, Cor. 4.7].

Corollary 1. With MR, the expectations in (13) become

$$\mathbb{E}\{\mathbf{h}_k^T \mathbf{D}_k \mathbf{w}_k\} = \sum_{l \in \mathcal{D}_k} \sqrt{\rho_{il} p_k \tau_p} \text{tr}(\mathbf{D}_{kl} \mathbf{R}_{kl} \Psi_{t_k}^{-1} \mathbf{R}_{kl}) \quad (16)$$

$$\mathbb{E}\{|\mathbf{h}_k^T \mathbf{D}_i \mathbf{w}_i|^2\} = \sum_{l=1}^L \rho_{il} \frac{\text{tr}(\mathbf{D}_{il} \mathbf{R}_{il} \Psi_{t_l}^{-1} \mathbf{R}_{il} \mathbf{D}_{il} \mathbf{R}_{kl})}{\text{tr}(\mathbf{R}_{il} \Psi_{t_l}^{-1} \mathbf{R}_{il})} + \begin{cases} \left| \sum_{l=1}^L \sqrt{\rho_{il} p_k \tau_p} \frac{\text{tr}(\mathbf{D}_{il} \mathbf{R}_{il} \Psi_{t_l}^{-1} \mathbf{R}_{kl})}{\sqrt{\text{tr}(\mathbf{R}_{il} \Psi_{t_l}^{-1} \mathbf{R}_{il})}} \right|^2 & \text{if } t_i = t_k \\ 0 & \text{otherwise,} \end{cases} \quad (17)$$

where t_i is the index of the pilot assigned to UE i .

AP l needs to select the transmit powers ρ_{il} for $i \in \mathcal{D}_l$. Network-wide optimization algorithms, considered in [3], [10], [14], are not scalable as $K \rightarrow \infty$ since the number of optimization variables grows with K .² Since each AP is (at

²It is possible to implement network-wide optimization problems in an iterative semi-distributed way, for example, using dual decomposition theory [14, Sec. 4.3]. However, these approaches converge slowly, require even more optimization variables, and need a lot of backhaul signaling. Hence, this approach is neither practical nor scalable.

$$R_k^{\text{ul}} = \frac{\tau_u}{\tau_c} \log_2 \left(1 + \frac{p_k \left| \sum_{l=1}^L \mathbb{E} \{ \mathbf{v}_{kl}^H \mathbf{D}_{kl} \mathbf{h}_{kl} \} \right|^2}{\sum_{i=1}^K p_i \mathbb{E} \left\{ \left| \sum_{l=1}^L \mathbf{v}_{kl}^H \mathbf{D}_{kl} \mathbf{h}_{il} \right|^2 \right\} - p_k \left| \sum_{l=1}^L \mathbb{E} \{ \mathbf{v}_{kl}^H \mathbf{D}_{kl} \mathbf{h}_{kl} \} \right|^2 + \sigma^2 \sum_{l=1}^L \mathbb{E} \{ \|\mathbf{D}_{kl}^H \mathbf{v}_{kl}\|^2 \}} \right)} \right) \quad (20)$$

least partially) unaware of the power allocation decisions made at other APs, only heuristic solutions are scalable. There are plenty of such schemes in the literature; some examples are found in [5], [7], [10], [16], [14, Sec. 3.4.4]. Since evaluation and comparison of heuristic schemes require extensive simulations, which is outside the scope of this paper, we consider only equal power allocation at each AP:

$$\rho_{il} = \begin{cases} \frac{\rho}{|\mathcal{D}_l|} & \text{if } i \in \mathcal{D}_l, \\ 0 & \text{otherwise.} \end{cases} \quad (18)$$

Since each UE is guaranteed to be served by at least one AP (i.e., its Master AP), each UE will be assigned non-zero transmit power when using (18) and, thus, get a non-zero SE.

D. Uplink SE and Distributed Combining

Next, we derive and analyze an achievable uplink SE expression based on the combined uplink signal in (7) from the serving APs. We use the *use-and-then-forget bound* that is widely used in the mMIMO literature [1, Th. 4.4], and also used in [11], [12] for Cell-free mMIMO with $N = 1$, $\mathbf{D}_i = \mathbf{I}_M \forall i$, for specific choices of combining schemes.

Proposition 2. An achievable uplink SE R_k^{ul} [bit/s/Hz] for UE k is given in (20) on the top of this page.

Proof: This is proved by following the same approach as in [1, Th. 4.4], but for the received signal in (7). ■

A key difference between the uplink and downlink is that, in the uplink, only the APs that serve the UE affects its SE.

AP l need to select its combining vectors \mathbf{v}_{il} for $i \in \mathcal{D}_l$ as a function of $\{\hat{\mathbf{h}}_{il} : i \in \mathcal{D}_l\}$. Two schemes that satisfy this requirement are MR (i.e., matched filtering) and RZF:

$$\mathbf{v}_{kl} = \begin{cases} \hat{\mathbf{h}}_{kl} & \text{with MR,} \\ p_k \left(\sum_{i \in \mathcal{D}_l} p_i \hat{\mathbf{h}}_{il} \hat{\mathbf{h}}_{il}^H + \sigma^2 \mathbf{I}_N \right)^{-1} \hat{\mathbf{h}}_{kl} & \text{with RZF,} \end{cases} \quad (21)$$

for $k \in \mathcal{D}_l$. This new variant of RZF only includes the channel estimates of the UEs that the AP is serving, thus making it a novel contribution and different from the non-scalable L-MMSE scheme recently proposed in [17]. The expectations in (20) can be computed by Monte Carlo simulations for RZF, while a closed form expression similar to Corollary 1 can be obtained for MR, but we omit this part for space limitations.

The uplink transmit powers $\{p_k : \forall k\}$ need to be selected and we assume that each UE has a maximum power of P . The network-wide power optimization methods in [3], [11], [12] are not scalable, thus a heuristic solution is needed. Since transmission at full power provided good performance for both strong and weak UEs in [17], we use that power control policy:

$$p_k = P, \quad k = 1, \dots, K. \quad (22)$$

E. Network Topology, Signal Encoding and Decoding

The proposed algorithms are transparent to the network topology since only neighboring APs cooperate. One option is to have local processing at each AP, as in classic cellular networks, and backhaul connections to the core network. Another option is to divide the APs into disjunct sets and connect each one via fronthaul to an edge-cloud processor [9], [15], [16] for centralized processing, as illustrated in Fig. 1. Many other cloud-RAN implementations are possible.

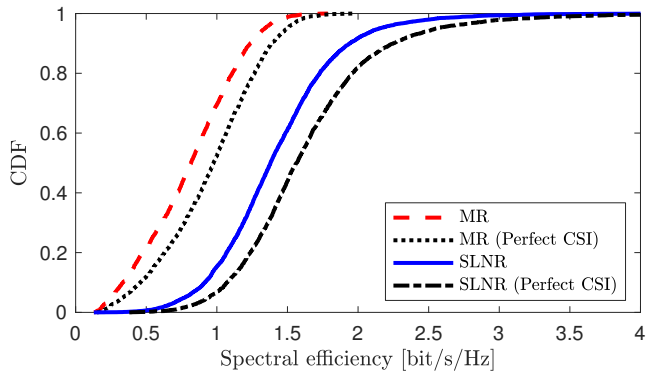
The most complicated part to implement is the encoding of downlink data and decoding of uplink data. We propose that the Master AP is carrying out these tasks, either locally or by delegating the task to a nearby edge-cloud processor. In the downlink, the data is shared to the neighboring APs that also transmit to the UE. In the uplink, the neighboring APs make soft estimates $\mathbf{v}_{kl}^H \mathbf{y}_l^{\text{ul}}$ of the data, which are sent to the Master AP (or a nearby edge-cloud processor) for final decoding.

IV. NUMERICAL RESULTS

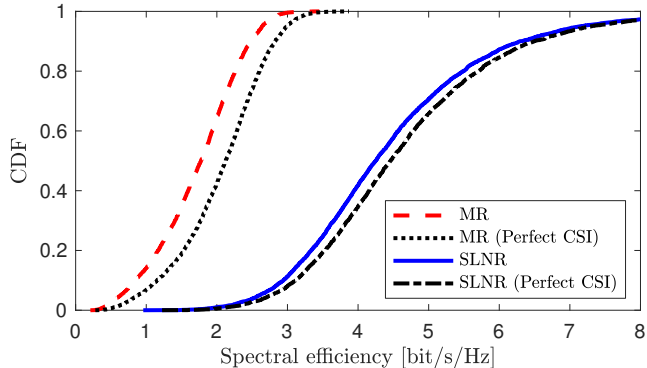
In this section, we compare the SEs achieved with standard MR and the proposed SLNR precoding and RZF combining. We consider a setup with $M = 400$ APs and $K = 100$ UEs independently and uniformly distributed in a 2×2 km square. By using the wrap-around technique, we approximate an infinitely large network with 100 APs/km² and 25 UEs/km². The UEs connect to the network as described in Section III-B, starting with τ_p UEs with different pilots and then letting UEs connect one after the other. We use the propagation model from [1, Sec. 4.1.3] with correlated fading, with the difference that the APs are deployed 10 m above the UEs, which gives a minimum distance. We have $\tau_c = 200$, $\tau_p = 10$, $\rho = p_k = 100$ mW, and 20 MHz bandwidth. We use $\tau_d = 190$ and $\tau_u = 190$ when evaluating downlink and uplink, respectively.

Fig. 4 shows the cumulative distribution function (CDF) of the downlink SE per UE in (13), where the randomness considers different AP and UE locations. We compare MR or SLNR precoding. Fig. 4(a) considers $N = 1$ and we notice that SLNR achieves 80% higher average SE than MR; the improvement is largest for the UEs with the best channel conditions. This means that the phenomenon illustrated in Fig. 3 has a huge impact on performance: MR achieves slightly higher signal power than SLNR for every UE, but also a much higher interference for most UEs due to the larger variations.

Fig. 4(b) considers $N = 4$ and the multiple AP antennas improve the SE of all the UEs, since the precoding gain can increase the received signal power with up to $4 \times$. The gain is particularly large with SLNR since each AP can now also suppress interference spatially. SLNR now achieves 155% higher average SE than MR. The genie-aided case with perfect



(a) $N = 1$ antennas per AP.



(b) $N = 4$ antennas per AP.

Fig. 4: Downlink SE per UE for different precoding schemes.

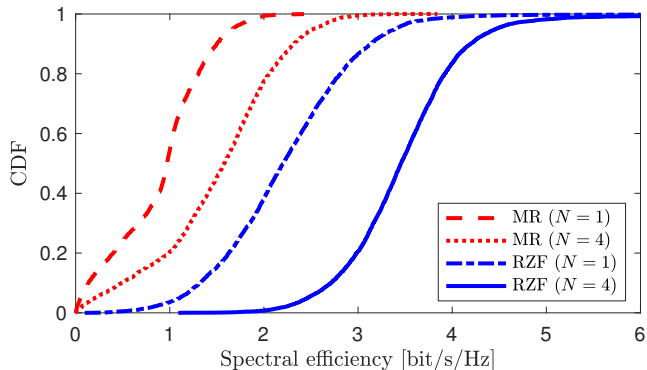


Fig. 5: Uplink SE per UE for different combining schemes.

CSI at the UEs is shown as a reference. We notice that the SEs are close to these upper bounds, thus the difference between MR and SLNR is not an artifact of the hardening bound.

Fig. 5 shows the CDF of the uplink SE per UE in (20) with $N = 1$ or $N = 4$. We notice that RZF outperforms MR, which is expected in light of the downlink results—RZF is basically the uplink counterpart of SLNR. The difference is so large that SLNR with $N = 1$ gives 48% higher average SE than MR with $N = 4$. In addition to the results shown in Fig. 5, we compared with the SE in the genie-aided case when perfect CSI is available in the decoding. RZF gives SE close to that genie-bound, while MR does not. Hence, it is only when using MR that the lack of channel hardening in Cell-free mMIMO (which was proved in [18]) make the SE expressions loose.

V. CONCLUSIONS

This paper has proposed a scalable, distributed implementation of Cell-free mMIMO by exploiting the DCC framework from the Network MIMO literature. We considered initial access, pilot assignment, cooperation cluster formation, precoding, and combining. We have demonstrated that MR is outperformed by the proposed SLNR (RZF) in the downlink (uplink), even for single-antenna APs. We stress that the new SLNR and RZF are fully distributed schemes, in contrast to the network-wide schemes considered in [10], [11], [17]. Hence, SLNR and RZF are respectively the new state-of-the-art in distributed precoding and combining for Cell-free mMIMO.

REFERENCES

- [1] E. Björnson, J. Hoydis, and L. Sanguinetti, “Massive MIMO networks: Spectral, energy, and hardware efficiency,” *Foundations and Trends® in Signal Processing*, vol. 11, no. 3-4, pp. 154–655, 2017.
- [2] S. Shamai and B. M. Zaidel, “Enhancing the cellular downlink capacity via co-processing at the transmitting end,” in *IEEE Vehicular Technology Conference (VTC Spring)*, vol. 3, 2001, pp. 1745–1749.
- [3] H. Q. Ngo, A. Ashikhmin, H. Yang, E. G. Larsson, and T. L. Marzetta, “Cell-free Massive MIMO versus small cells,” *IEEE Trans. Wireless Commun.*, vol. 16, no. 3, pp. 1834–1850, 2017.
- [4] S. Venkatesan, A. Lozano, and R. Valenzuela, “Network MIMO: Overcoming intercell interference in indoor wireless systems,” in *Asilomar Conference on Signals, Systems and Computers*, 2007, pp. 83–87.
- [5] E. Björnson, R. Zakhour, D. Gesbert, and B. Ottersten, “Cooperative multicell precoding: Rate region characterization and distributed strategies with instantaneous and statistical CSI,” *IEEE Trans. Signal Process.*, vol. 58, no. 8, pp. 4298–4310, 2010.
- [6] A. Osseiran, J. F. Monserrat, and P. Marsch, *5G Mobile and Wireless Communications Technology*. Cambridge: Cambridge University Press, 2016, ch. 9, Coordinated Multi-Point Transmission in 5G.
- [7] E. Björnson, N. Jaldén, M. Bengtsson, and B. Ottersten, “Optimality properties, distributed strategies, and measurement-based evaluation of coordinated multicell OFDMA transmission,” *IEEE Trans. Signal Process.*, vol. 59, no. 12, pp. 6086–6101, 2011.
- [8] T. L. Marzetta, “Noncooperative cellular wireless with unlimited numbers of base station antennas,” *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, 2010.
- [9] S. Perlman and A. Forenza, “An introduction to pCell,” 2015, Artemis Networks LLC, White paper. [Online]. Available: <http://www.rearden.com/artemis/An-Introduction-to-pCell-White-Paper-150224.pdf>
- [10] E. Nayebe, A. Ashikhmin, T. L. Marzetta, H. Yang, and B. D. Rao, “Precoding and power optimization in cell-free Massive MIMO systems,” *IEEE Trans. Wireless Commun.*, vol. 16, no. 7, 2017.
- [11] E. Nayebe, A. Ashikhmin, T. L. Marzetta, and B. D. Rao, “Performance of cell-free massive MIMO systems with MMSE and LSFD receivers,” in *Asilomar Conference on Signals, Systems and Computers*, 2016.
- [12] M. Bashar, K. Cumanan, A. G. Burr, M. Debbah, and H. Q. Ngo, “On the uplink max-min SINR of cell-free massive MIMO systems,” *IEEE Trans. Wireless Commun.*, vol. 18, no. 4, pp. 2021–2036, 2019.
- [13] S. Buzzi and C. D’Andrea, “Cell-free massive MIMO: User-centric approach,” *IEEE Commun. Lett.*, vol. 6, no. 6, pp. 706–709, 2017.
- [14] E. Björnson and E. Jorswieck, “Optimal resource allocation in coordinated multi-cell systems,” *Foundations and Trends® in Communications and Information Theory*, vol. 9, no. 2-3, pp. 113–381, 2013.
- [15] A. G. Burr, M. Bashar, and D. Maryopi, “Ultra-dense radio access networks for smart cities: Cloud-RAN, fog-RAN and “cell-free” massive MIMO,” in *PIMRC, International Workshop of CorNer*, 2018.
- [16] G. Interdonato, P. Frenger, and E. G. Larsson, “Scalability aspects of cell-free massive MIMO,” in *IEEE International Conference on Communication (ICC)*, 2019.
- [17] E. Björnson and L. Sanguinetti, “Making cell-free massive MIMO competitive with MMSE processing and centralized implementation,” *CoRR*, vol. abs/1903.10611, 2019. [Online]. Available: <http://arxiv.org/abs/1903.10611>
- [18] Z. Chen and E. Björnson, “Channel hardening and favorable propagation in cell-free Massive MIMO with stochastic geometry,” *IEEE Trans. Commun.*, vol. 17, no. 11, pp. 5205–5219, 2018.