# Cell-Free Massive MIMO in O-RAN: Energy-Aware Joint Orchestration of Cloud, Fronthaul, and Radio Resources

Özlem Tuğfe Demir, *Member, IEEE*, Meysam Masoudi, *Member, IEEE*, Emil Björnson, *Fellow, IEEE*, and Cicek Cavdar, *Member, IEEE*

*Abstract*—For the energy-efficient deployment of cell-free massive MIMO functionality in a practical wireless network, the end-to-end (from radio site to the cloud) energy-aware operation is essential. In line with the cloudification and virtualization in the open radio access networks (O-RAN), it is indisputable to envision prospective cell-free infrastructure on top of the O-RAN architecture. In this paper, we explore the performance and power consumption of cell-free massive MIMO technology in comparison with traditional small-cell systems, in the virtualized O-RAN architecture. We compare two different functional split options and different resource orchestration mechanisms. In the end-to-end orchestration scheme, we aim to minimize the end-to-end power consumption by jointly allocating the radio, optical fronthaul, and virtualized cloud processing resources. We compare end-to-end orchestration with two other schemes: i) "radio-only" where radio resources are optimized independently from the cloud and ii) "local cloud coordination" where orchestration is only allowed among a local cluster of radio units. We develop several algorithms to solve the end-to-end power minimization and sum spectral efficiency maximization problems. The numerical results demonstrate that end-to-end resource allocation with fully virtualized fronthaul and cloud resources provides a substantial additional power saving than the other resource orchestration schemes.

*Index Terms*—Cell-free massive MIMO, virtualized O-RAN, joint transmission, end-to-end resource allocation, joint network orchestration.

## I. INTRODUCTION

The number of mobile user equipments (UEs) and their capacity requests are anticipated to increase continuously during the next decade [1]. It is expected that by 2030 ubiquitous and limitless connectivity will be needed everywhere [2]. New networking and new air interface solutions are the two key enablers to support next-generation wireless systems [1]. The networking technologies include the softwarization, virtualization, and open radio access networks (O-RAN), whereas massive MIMO (multiple-input multiple-output) and cell-free massive MIMO are among the main air interface technologies. This paper studies cell-free massive MIMO in the O-RAN

architecture, which brings together two inherently compatible networking and air interface entities of beyond 5G networks.

Cell-free massive MIMO has been proposed as a physical-layer technology that combines ultra-dense networks with joint transmission/reception (JT), and the low-complexity linear processing schemes from massive MIMO [3]–[5]. By taking advantage of both joint processing and macro diversity, cell-free massive MIMO reduces the large data rate variations across the coverage area, which solves one of the main drawbacks of the current cellular networks. For joint processing of UEs' signals, it is required to centralize parts of the baseband processing, which constructs the inherent connection between cell-free massive MIMO and the centralized RAN (C-RAN) architecture [6]. The separation of software from hardware not only allows for joint processing in cell-free operation but also creates new energy-saving opportunities with green and agile virtualization [7]. To meet the coordination and signaling requirements of the envisioned cell-free massive MIMO network, O-RAN is envisaged as a promising architecture by providing substantial adaptability [8].

In this paper, in a cell-free massive MIMO network with O-RAN architecture, we jointly allocate radio, optical transport network, and cloud processing resources given the number of UEs and their performance requirements in an area. We derive the required processing resources for each operation in the cloud and determine the required optical transport resources based on different functional split options in a cell-free massive MIMO system. The joint orchestration of end-to-end resources enabled by O-RAN architecture is critical to fully benefit from energy-saving mechanisms in a cell-free network. Indeed, O-RAN enables the joint resource allocation by the real-time and near-real-time softwarized controllers [9]. Thanks to this joint orchestration, the number of active radio units (RUs) can scale down together with the amount of optical and processing resources in the transport network and the cloud, respectively, following the UE demand, to minimize the total end-to-end power consumption of the network.

### A. Evolution of RAN architecture towards O-RAN

The distributed RAN (D-RAN) is the most widely deployed legacy network architecture, in which all baseband processing functions for each base station (BS) are co-located with the RU at the cell site [10]. This architecture is not scalable, cost-efficient, and, most importantly, is not capable of supporting heterogeneous services efficiently in terms of energy

consumption and throughput [11]. The C-RAN architecture was developed to improve the network's energy efficiency and resource utilization. In conventional C-RAN, the baseband processing functions are detached from the RUs at the cell site and moved into a centralized resource pool at a central cloud (CC). Although this architecture is more energy-efficient than D-RAN due to centralized cooling, each RU has its dedicated processing unit making the processing resources under-utilized when the traffic is unbalanced between cells during peak hours.

Virtualized C-RAN emerged as an architecture that decouples the hardware and software by virtualizing network functions [10], [12]–[17]. In this architecture, the deployed processing units are no longer dedicated to one specific RU, but network functionalities are implemented in software and run on general-purpose processors (GPPs) [18]. In this way, processing resources can be shared between various loaded cells, further improving resource utilization and reducing network energy consumption. The advantages of implementing virtualized C-RAN are i) simplified network management; ii) enabled resource pooling; and iii) improved coordination of radio resources required for cell-free operations. Although the virtualized C-RAN architecture advancement is promising, full centralization of physical-layer processing significantly upscales the fronthaul signal capacity, especially when technologies like massive MIMO are employed. Therefore, a more convenient and potentially flexible architecture needed to be further investigated to keep the scalability benefits of C-RAN while resolving the bandwidth congestion and allowing for effective coordinated multipoint (CoMP) coordination for cell-edge UEs.

Recently, a consensus has been reached between network vendors and operators to support an O-RAN architecture and standards [19]. The O-RAN Alliance [20] is an industry-wide standardization for RAN interfaces that complements the 3GPP standards and covers RAN disaggregation, RAN automation, and RAN virtualization. O-RAN architecture, proposed by the O-RAN Alliance, is a virtualized C-RAN with an open, interoperable interface and virtualization, allowing multiple vendor products to work together in one network. O-RAN and following standardization efforts enable building the virtualized C-RAN on open hardware and cloud, and allowing full exploitation of virtualization and sharing in virtualized C-RAN. The three key elements of O-RAN are i) cloudification; ii) intelligence and automation; and iii) open internal RAN interface. The primary mission of the O-RAN is to reshape the RAN industry towards open, virtualized, and fully interoperable mobile networks [21].

### B. Cell-free massive MIMO in the O-RAN architecture

In the previous work, mostly the physical-layer aspects of cell-free massive MIMO have been studied and only the radio site power consumption has been considered [22]. In [23], the power consumption of the fronthaul transport is also taken into account, but the authors assume all RUs serve all UEs, which is not power efficient. In [10], a cell-free network architecture is presented by optimizing the user-centric formation of soft BSs defined as joint allocation of spectrum, optical wavelength,

and cloud processing resources together with a set of RUs considering JT. End-to-end power consumption [10] and network throughput [12] are optimized considering radio, optical fronthaul, and cloud processing resources. However, massive MIMO is not considered in these cell-free networks. In [24], RU selection for JT under fronthaul constraints is studied, but the cloud processing power consumption is simplified as a fixed parameter. In [25], the processing requirements in the cloud are taken into consideration, but only radio site power consumption is minimized.

Recently, the authors of [26] have refined the cell-free massive MIMO terminology according to the O-RAN architecture and discussed several implementation options of cell-free functionality in the current or future O-RAN generations. The works [8], [27] also studied the performance of the cell-free massive MIMO on top of the O-RAN architecture. In [28], the placement of the central processing unit and allocation of radio bandwidth have been studied in terms of throughput. To the best of the authors' knowledge, all the related work on cell-free massive MIMO in O-RAN either considered architectural high-level views or focused on spectral efficiency (SE) or throughput performance.

### C. Contributions

In the conference version [29] of this paper, the end-to-end power consumption modeling and minimization were considered only with fully-centralized functional splitting option in virtualized C-RAN architecture. In this paper, we follow a holistic approach by studying different resource orchestration schemes for the cell-free massive MIMO and small-cell systems in the O-RAN architecture. The considered end-to-end network power consumption for cell-free massive MIMO involves the impact of radio, optical fronthaul, and cloud processing resources. Extending the conference version, we additionally consider intra-physical-layer functional splitting option by modifying the power consumption accordingly. We derive the cloud processing requirements of a cell-free massive MIMO OFDM system, given the required system performance. Based on our developed end-to-end power consumption model, we cast two optimization problems to jointly allocate the radio, fronthaul, and cloud resources. The first problem, which was the only problem considered in [29], minimizes the end-to-end power consumption by joint allocation of transmit powers, optical fronthaul resources, cloud processing resources, and the set of O-RUs (in line with the O-RAN terminology) serving the UEs to meet their quality of service (QoS) requirements. We cast the problem in a mixed binary second-order cone programming form, which can be optimally solved. Different from the conference paper, an approximated version of the original power minimization problem is obtained via $l_0$ norm, and a concave-convex programming (CCP)-based algorithm is proposed to solve this problem in a more manageable form to gain further insights into larger cell-free setups. Moreover, a joint sum SE maximization and total network power minimization problem is proposed. After novel transformations, a proper approximated form of this problem is solved via the same CCP-based algorithm. Using the found solutions, we compare

the performance of the fully virtualized end-to-end, local cloud coordination-based, and radio-only resource allocation, where only the end-to-end scheme was considered in [29]. Through numerical simulations, we show how much power saving is achieved by the virtualized end-to-end resource allocation compared to the case of fixed fronthaul resources and partial resource sharing in the cloud (local coordination), and the cloud-unaware radio-only scheme. Moreover, the SE improvement provided by the cell-free massive MIMO over conventional small-cell networks, where each UE is served by only one O-RU, is quantified for different scenarios. The effect of different functional splits is discussed.

### D. Paper Outline

The remainder of this paper is organized as follows. Section II overviews the O-RAN architecture for cell-free massive MIMO functionality. Section III introduces the channel model, channel estimation, and downlink operation in a cell-free massive MIMO system. In Section IV, the end-to-end power consumption modeling together with the analysis of processing complexity is elaborated. The details of the end-to-end power minimization problem and the respective algorithms are provided in Section V. Section VI extends the developed optimization methodology to the joint sum SE maximization and power minimization problem. The performance of the proposed end-to-end energy-aware algorithms is quantified and compared with the partial energy-saving mechanisms in Section VII. Finally, Section VIII concludes the paper.

**Reproducible research:** All the simulation results can be reproduced using the Matlab code available at: https://github.com/ozlemtugfedemir/O-RAN-cell-free

## II. Architecture Overview of O-RAN for Cell-free Massive MIMO

We consider a cell-free massive MIMO system that is built on the top of the O-RAN architecture in line with the next-generation virtualized C-RAN ecosystem as shown in Fig. 1 [19]. There are $L$ O-RUs and $K$ UEs that are arbitrarily distributed in the coverage area. All UEs have a single antenna while each O-RU is equipped with $N$ antennas. All the O-RUs are connected to the O-Cloud with virtualization and processing resource sharing capabilities [10], via fronthaul connections. O-Cloud consists of two main units, which are O-CU (centralized unit) and O-DU (distributed unit). According to the O-RAN specification, the O-DU is responsible for the lower network layer operations (RCL, MAC, and PHY) whereas the O-CU implements the higher layer operations as illustrated in Fig. 1. O-RAN also has logical nodes known as near real-time RAN intelligent controller (near-RT RIC) and non-RT RIC. The near-RT-RIC is responsible for near real-time intelligent optimization of RAN resources. On the other hand, non-RT RIC is located in the service management and orchestration (SMO) unit, which is responsible for non-real-time intelligent orchestration. These two logical nodes enable fully virtualized end-to-end resource optimization, which we consider in this paper. O-RAN has multiple deployment options, in some of which O-DU and O-CU are co-located,

and in some of which, they are separated logically and geographically.

In this study, we consider Scenario A as the O-RAN deployment setup, where O-CU and O-DU are bundled and co-located with the logical node near-RT RIC communicating through E2 interface defined by O-RAN [30]. One key advantage of this scenario is the minimized delay while the deployment cost might be larger compared to the other scenarios [31]. Hierarchical deployment scenarios in O-RAN are kept as future work since they will not affect the key findings of our study. A set of O-RUs serves each UE coherently to improve SE. Let $x_{k,l} \in \{0, 1\}$ be the binary variable denoting whether UE $k$ is served by O-RU $l$ or not. It is one if UE $k$ is served by O-RU $l$, and zero otherwise. For example, in Fig. 1, the colored circular regions for each UE indicate the O-RUs that are serving them. In this paper, we will optimize these subsets to satisfy several optimization metrics that take the SE of each UE and the end-to-end network power consumption into account.

In the O-Cloud, there are $W$ stacks of general-purpose processors (GPPs) in line with the cloudification framework of O-RAN [19]. These pooled GPPs are used for baseband processing due to their processing capabilities and programmability, which allows virtualization. The workload of each GPP is routed via a dispatcher that is controlled by a global cloud controller [16], which is near-RT-RIC in the considered O-RAN architecture. To perform cell-free joint processing, the respective data and control signals for a particular UE should first exist in the same GPP even though the locally precoded signals can be computed and sent to the O-RUs from different GPPs (thanks to distributed operation and computational sharing and virtualization among multiple GPPs). In Fig. 1, the same colors are used to show which O-RUs are connected to which GPP. UE 1 is served by O-RU 1 and O-RU 2, which are all connected to GPP 1. In this case, UE 1 can be served using potentially only GPP 1 without any need for data exchange among GPPs. However, since the O-RUs that serve UE 3 are connected to different GPPs, the connection should be activated between GPP 1 and GPP 2 for the sharing of UE 3 data and control signal. These are all orchestrated by the near-RT-RIC.

We consider the recently released evolved CPRI (eCPRI) specification for the fronthaul/midhaul transmission. A time- and wavelength-division multiplexed passive optical network (TWDM-PON) is utilized as the fronthaul transport network to carry eCPRI packets to meet the high capacity requirements of the fronthaul transmission in a cell-free network [10], [13]. As shown in Fig. 1, each O-RU is connected to one optical network unit (ONU) that is assigned to one of the multiple wavelengths in the fiber communication. Each wavelength can be shared by more than one O-RU using time-division multiplexing. In the O-Cloud, there exists an optical line terminal (OLT) with a WDM multiplexer (WDM MUX) and multiple line cards (LCs), each of which is connected to one GPP. Each LC serves only one wavelength and, thus, each O-RU's signals are received at or transmitted from the GPP that uses the same wavelength. For example, GPP 1 is responsible for the fronthaul transport of the first four O-RUs'
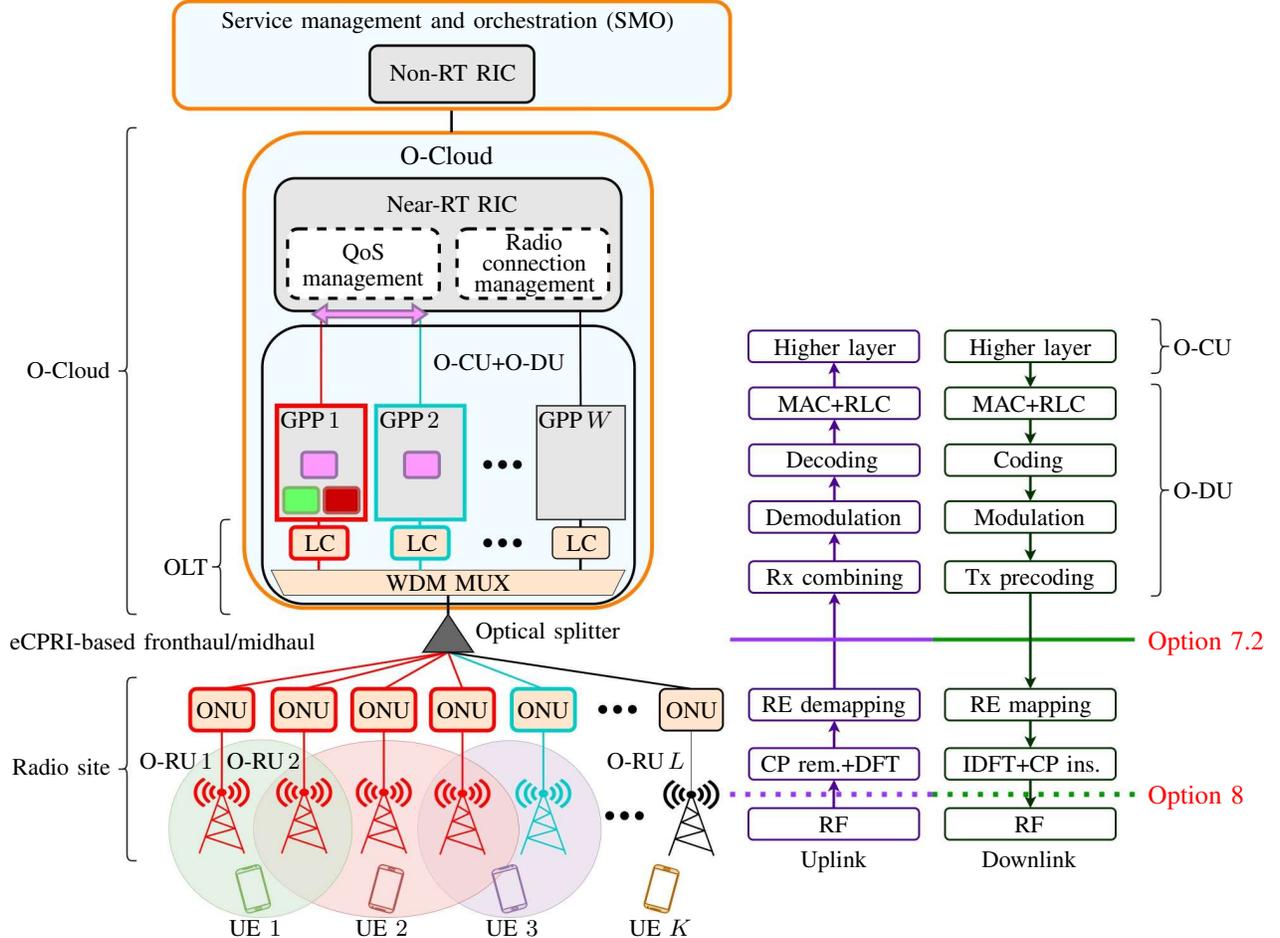
Fig. 1. O-RAN architecture for cell-free massive MIMO with functional splitting options 8 and 7.2 in uplink and downlink.

data, whereas O-RU 5 is connected to GPP 2.

In the considered O-RAN architecture, two different functional split options from 3GPP specification [32] are considered. The first one, which is called *Option 8* and shown in Fig. 1, is the physical layer (PHY)-radio frequency (RF) functional split to fully benefit from the efficient processing resources and housing facilities in the O-Cloud.[1] According to this split, the O-RUs only perform RF processing and transmit (receive) the pure sampled and quantized baseband signals to (from) the O-Cloud via fronthaul links. All the remaining processing is done in the O-Cloud. Hence, Option 8 has the advantage of the lowest RU complexity [34]. In particular, with functional split option 8 in the uplink training, all the operations below the dashed purple line in Fig. 1 are implemented at the O-RUs whereas the remaining physical-layer, MAC and RLC operations are carried out in the O-DU. Then, higher-layer functions are implemented in O-CU. Similarly, the dashed green line determines the same functional split for the downlink operation. With PHY-RF functional split option 8, the required fronthaul data rate for each AP is given

as [35]

$$R_{\mathrm{fronthaul}}^{(8)} = 2 f_s N_{\mathrm{bits}} N, \quad (1)$$

where $f_s$ and $N_{\mathrm{bits}}$ are the sampling frequency and the number of bits to quantize the signal samples, respectively. Due to the limited capacity of each wavelength in TWDM-PON, which is denoted by $R_{\max}$, we can assign at most $W_{\max} \triangleq \lfloor R_{\max}/R_{\mathrm{fronthaul}}^{(8)} \rfloor$ O-RUs to each wavelength and, hence, to each GPP $w$, for $w = 1, \ldots, W$.

The other functional split is *Option 7.2* and is more specifically considered in O-RAN [34], [36]. This option splits the network functions so that some of the PHY operations close to RF processing (low-PHY functions) are implemented in the O-RUs whereas the remaining PHY (high-PHY) and higher layer operations are moved to the O-Cloud as shown by the solid purple and green lines in Fig. 1.[2] This functional splitting option generally lowers the fronthaul data rate compared to Option 8, by still allowing cell-free JT processing and low O-RU complexity. The uplink receive combining and transmit precoding are considered in the O-Cloud for coherent JT. The

[1]The functional split option 7.2 is the one that is mainly supported by O-RAN. However, functional split option 8 is also important to consider not only due to its energy-saving potential but also the high experimentation interest by the leading researchers in the field according to the survey results in [33, Fig. 1].

[2]In the literature, some authors use the term *midhaul* to refer to the link between O-RU and O-DU when the functional split option 7.2 is adopted. In this paper, we use *fronthaul* to refer to the corresponding link for all the split options without loss of generality.

corresponding fronthaul data rate requirement is given as [35]

$$R_{\text{fronthaul}}^{(7.2)} = \frac{2N_{\text{bits}}N_{\text{used}}N}{T_s}, \tag{2}$$

where $N_{\text{used}}$ is the number of used subcarriers and $T_s$ is the OFDM symbol duration. We expect Option 7.2 is more advantageous since it usually holds that $R_{\text{fronthaul}}^{(7.2)} < R_{\text{fronthaul}}^{(8)}$ so that we can assign more O-RUs to each GPP, i.e., $W_{\max} \triangleq \lfloor R_{\max}/R_{\text{fronthaul}}^{(7.2)} \rfloor$ O-RUs to each wavelength. To quantify this, let us consider a setup with $R_{\max} = 10\,\text{GBps}$, $N_{\text{bits}} = 12$, and conventional LTE parameters $f_s = 30.72\,\text{MHz}$, $T_s = 71.4\,\mu\text{s}$ $N_{\text{used}} = 1200$. Further assuming $N = 4$, each GPP, and, thus each wavelength, can support $W_{\max} = 3$ O-RUs with Option 8, while $W_{\max} = 6$ O-RUs can be connected to each GPP with Option 7.2. When the number of O-RU antennas is greater than 13, the maximum capacity of each TWDM-PON wavelength is not big enough to provide a fronthaul connection between the O-RUs and the O-Cloud with Option 8. However, with Option 7.2, still, up to 24-antenna O-RUs can be supported by each wavelength due to the reduced fronthaul rate requirement.

Apart from fronthaul transport limitations, the relation of the processing requirements in the O-Cloud and how the GPPs can handle those based on their limited capabilities determine the number of active GPPs. Shutting down the unused active GPPs and the corresponding LCs, we can save power by minimizing the active idle power of the network [16]. However, there is a trade-off between minimizing the number of active GPPs and the SE performance of UEs. In this paper, we will derive the related processing requirements of each operation in the O-Cloud for a cell-free massive MIMO OFDM system. Before we introduce our optimization problems, we will present the physical-layer foundations of cell-free massive MIMO system in the next section.

## III. CELL-FREE MASSIVE MIMO SYSTEM

A cell-free massive MIMO system with time-division duplex and OFDM is considered. The carrier and sampling frequencies are $f_c$ and $f_s$, respectively. We assume a block-fading channel model as illustrated in Fig. 2. The total number of subcarriers is $N_{\text{DFT}}$ across the total bandwidth of $B$ Hz. $N_{\text{DFT}}$ is also the dimension of the discrete Fourier transform (DFT), while the number of used subcarriers is $N_{\text{used}} \leq N_{\text{DFT}}$ shown as red rectangles. The light blue rectangles represent the null subcarriers. Each OFDM symbol has a duration of $T_s$ seconds. According to the block fading channel modeling, the channels are constant time-invariant and frequency-flat in each coherence block that consists of $N_{\text{smooth}}$ consecutive OFDM subcarriers and $N_{\text{slot}}$ OFDM symbols [37], [5, Remark 2.1]. The channel is constant across $\tau_c = N_{\text{smooth}}N_{\text{slot}}$ channel uses, which is the number of useful samples in each coherence block, and takes independent realizations between different blocks. Although in practice each OFDM sample (channel use) has a unique channel response, the block-fading assumption can be accurately applied with a proper selection of $N_{\text{slot}}$ and $N_{\text{smooth}}$ since the channels do not change abruptly across consecutive time and frequency samples or there is a known mapping between them [5, Remark 2.1]. Hence, the channels

can be assumed to be approximately constant (smooth) across $\tau_c$ time-frequency samples in each coherence block without loss of generality.

There are mainly two types of cell-free network operation. The first one is the centralized operation, where all the processing regarding channel estimation and payload data detection/precoding are performed in the central processing unit (O-Cloud in the considered architecture) [5, Sec. 5.1 and 6.1]. On the other hand, in the distributed operation [5, Sec. 5.2 and 6.2], each O-RU first estimates the local channels and these estimates are used for local combining/precoding during data transmission. In this paper, we consider distributed downlink operation where the channel estimates corresponding to each O-RU are used for distributed per-O-RU precoding [5, Sec. 6.2]. In this way, we can divide the UE-related processing tasks into small independent blocks that can be virtualized in the O-Cloud to minimize the number of active GPPs.[3] Note that the fronthaul requirements scale with the number of O-RU antennas with both functional split Options 8 and 7.2. In addition, it is proportional to the number of used subcarriers if Option 7.2 is adopted. In both options, they do not depend on the number of UEs that an O-RU can serve.

Since we focus on the downlink operation, each coherence block is divided into two phases: i) uplink training with $\tau_p$ samples and ii) downlink payload data transmission with $\tau_d = \tau_c - \tau_p$ samples. All the UEs are served on the same time-frequency resources using spatial multiplexing. Moreover, channel estimation and precoding are implemented in each coherence block in the same way. Hence, we will focus on an arbitrary time-frequency resource block as in [5].

### A. Channel Model

We let $\mathbf{h}_{kl} \in \mathbb{C}^N$ denote the frequency-domain channel from UE $k$ to O-RU $l$ in an arbitrary coherence block. The channels are modeled using correlated Rayleigh fading, i.e., $\mathbf{h}_{kl} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}_N, \mathbf{R}_{kl})$ and they are independent for different UEs and O-RUs. The correlation matrix $\mathbf{R}_{kl} \in \mathbb{C}^{N \times N}$ is determined by the spatial correlation of the channel $\mathbf{h}_{kl}$ between the antennas of O-RU $l$ and the corresponding average channel gain, which is denoted by $\beta_{kl} = \text{tr}(\mathbf{R}_{kl})/N$. The channel gain is dependent on large-scale effects such as geometric attenuation and shadowing. The spatial correlation matrices are fixed throughout the communication and they are known in accordance with the related literature [5], [38].

### B. Uplink Training: Channel Estimation

To perform coherent transmit processing, the channels need to be estimated in each coherence block. In a large cell-free network with many UEs, there will not be enough pilot resources to assign orthogonal sequences to all UEs. Hence, we consider a set of $\tau_p$ mutually orthogonal pilot sequences

---

[3]This is not possible in the centralized cell-free operation where higher-dimensional precoding should be computed jointly for all the serving O-RUs, based on pooling of channel estimates, and then applied to a particular UE in the same GPP.
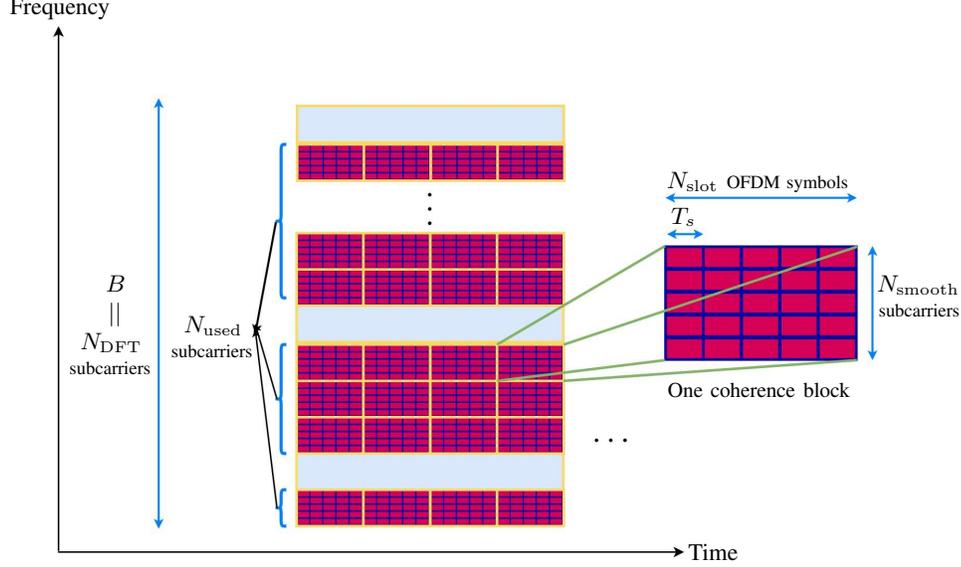
Fig. 2. In each coherence block of $\tau_c = N_{\text{smooth}} N_{\text{slot}}$ complex samples, the channel is modeled as time-invariant and frequency-flat according to the block-fading model. Out of total number of $N_{\text{DFT}}$ subcarriers, $N_{\text{used}} \leq N_{\text{DFT}}$ subcarriers are utilized.

$\phi_1, \ldots, \phi_{\tau_p} \in \mathbb{C}^{\tau_p}$ that are assigned to the UEs and reused by multiple UEs. The sequences satisfy

$$\phi_{t_1}^{\text{H}} \phi_{t_2} = \begin{cases} \tau_p, & t_1 = t_2, \\ 0, & t_1 \neq t_2. \end{cases} \quad (3)$$

The pilot sequences are assigned to the UEs in a deterministic way and $t_k$ denotes the index of the pilot assigned to UE $k$ as $t_k \in \{1, \ldots, \tau_p\}$. The set of UEs that share the same pilot with UE $k$ is defined as

$$\mathcal{P}_k = \{i : t_i = t_k, \ i = 1, \ldots, K\} \subset \{1, \ldots, K\}. \quad (4)$$

As shown in Fig. 1, either after initial RF processing (with Option 8) or after resource demapping (with Option 7.2), the received signals are transmitted to the O-Cloud and the channel estimation is performed there for each coherence block. The frequency-domain received signal at O-RU $l$ and a particular coherence block during the entire pilot transmission is given by

$$\mathbf{Y}_l^{\text{pilot}} = \sum_{i=1}^{K} \sqrt{\eta_i} \mathbf{h}_{il} \phi_{t_i}^{\text{T}} + \mathbf{N}_l \quad (5)$$

where $\eta_i \geq 0$ is the pilot transmit power of UE $i$ and $\mathbf{N}_l \in \mathbb{C}^{N \times \tau_p}$ is the receiver noise with independent and identically distributed (i.i.d.) elements as $\mathcal{N}_{\mathbb{C}}(0, \sigma^2)$. Using all the received pilots signals, the minimum mean-squared error (MMSE) channel estimate $\widehat{\mathbf{h}}_{kl}$ of the channel $\mathbf{h}_{kl}$ is given as follows [5, Corol. 4.1]:

$$\widehat{\mathbf{h}}_{kl} = \sqrt{\eta_k} \mathbf{R}_{kl} \left( \sum_{i \in \mathcal{P}_k} \eta_i \tau_p \mathbf{R}_{il} + \sigma^2 \mathbf{I}_N \right)^{-1} \mathbf{Y}_l^{\text{pilot}} \phi_{t_k}^*. \quad (6)$$

*C. Downlink Data Transmission*

Let $\varsigma_i \in \mathbb{C}$ denote the downlink data signal of UE $i$ with $\mathbb{E}\{|\varsigma_i|^2\} = 1$. Let the normalized (in terms of average power) precoding vector and transmit power corresponding to UE $i$ and O-RU $l$ for $x_{i,l} = 1$ be $\mathbf{w}_{il} \in \mathbb{C}^N$ and $p_{il} \geq 0$,

respectively. In the O-Cloud, the frequency-domain precoded signal to be transmitted from O-RU $l$ is constructed as

$$\mathbf{x}_l = \sum_{i=1}^{K} \sqrt{p_{il}} x_{i,l} \mathbf{w}_{il} \varsigma_i \in \mathbb{C}^N. \quad (7)$$

In selecting the precoding vectors (local partial MMSE (LP-MMSE) precoding [5, Sec. 6.2.2]), the MMSE channel estimates $\{\widehat{\mathbf{h}}_{il}\}$ from (6) are used. The received frequency-domain downlink signal at UE $k$ is[4]

$$y_k^{\text{dl}} = \sum_{l=1}^{L} \mathbf{h}_{kl}^{\text{T}} \mathbf{x}_l + n_k = \sum_{l=1}^{L} \sum_{i=1}^{K} \sqrt{p_{il}} x_{i,l} \mathbf{h}_{kl}^{\text{T}} \mathbf{w}_{il} \varsigma_i + n_k \quad (8)$$

where $n_k \sim \mathcal{N}_{\mathbb{C}}(0, \sigma^2)$ is the receiver noise. The downlink achievable SE of UE $k$ at each resource block can be computed using [5, Corr. 6.3 and Sec. 7.1.2] as follows:

$$\mathsf{SE}_k = \frac{\tau_d}{\tau_c} \log_2 (1 + \mathsf{SINR}_k) \quad \text{bit/s/Hz} \quad (9)$$

with the effective signal-to-interference-plus-noise ratio (SINR) given by

$$\mathsf{SINR}_k = \frac{|\mathbf{b}_k^{\text{T}} \boldsymbol{\rho}_k|^2}{\sum_{i=1}^{K} \boldsymbol{\rho}_i^{\text{T}} \mathbf{C}_{ki} \boldsymbol{\rho}_i + \sigma^2} \quad (10)$$

where

$$\boldsymbol{\rho}_k = [\sqrt{p_{k1}} x_{k,1} \ \cdots \ \sqrt{p_{kL}} x_{k,L}]^{\text{T}} \in \mathbb{R}_{\geq 0}^L, \quad (11)$$

$$\mathbf{b}_k \in \mathbb{R}_{\geq 0}^L, \quad [\mathbf{b}_k]_l = \mathbb{E}\{\mathbf{h}_{kl}^{\text{T}} \mathbf{w}_{kl}\} \quad (12)$$

$$\mathbf{C}_{ki} \in \mathbb{C}^{L \times L},$$

$$[\mathbf{C}_{ki}]_{lr} = \begin{cases} \mathbb{E}\{\mathbf{h}_{kl}^{\text{T}} \mathbf{w}_{kl} \mathbf{w}_{kr}^{\text{H}} \mathbf{h}_{kr}^*\} - [\mathbf{b}_k]_l [\mathbf{b}_k]_r^*, & i = k, \\ \mathbb{E}\{\mathbf{h}_{kl}^{\text{T}} \mathbf{w}_{il} \mathbf{w}_{ir}^{\text{H}} \mathbf{h}_{kr}^*\}, & i \neq k. \end{cases} \quad (13)$$

---

[4]Note that an additional conjugation on the channel vectors is not introduced different from [5], so the Hermitian transpose on the channel vectors in [5] are replaced by tranpose.

In this paper, we use a slightly modified version of the LP-MMSE precoding [5, Sec. 6.2.2], where $\mathbf{w}_{kl} = \overline{\mathbf{w}}_{kl}/\sqrt{\mathbb{E}\{\|\overline{\mathbf{w}}_{kl}\|^2\}}$ with

$$\overline{\mathbf{w}}_{kl} = \eta_k \left( \sum_{i \in \overline{\mathcal{P}}_l} \eta_i \left( \widehat{\mathbf{h}}_{il} \widehat{\mathbf{h}}_{il}^{\text{H}} + \mathbf{Z}_{il} \right) + \sigma^2 \mathbf{I}_N \right)^{-1} \widehat{\mathbf{h}}_{kl}. \quad (14)$$

Here, $\mathbf{Z}_{il}$ is the correlation matrix of the channel estimation error $\mathbf{h}_{il} - \widehat{\mathbf{h}}_{il}$ and $\overline{\mathcal{P}}_l$ only includes the UE indices with the strongest channel per each pilot at O-RU $l$.

## IV. POWER CONSUMPTION MODELING

For the considered virtualized O-RAN architecture, we do not stick to specific hardware but instead adopt a generic power model that can be applied to the different technologies. The network power consumption consists of two main components: i) the radio site power consumption that includes the O-RU power consumption $P_{\text{RU},l}$, for $l = 1, \ldots, L$ and the ONU power consumption, $P_{\text{ONU}}$; and ii) the O-Cloud power consumption, $P_{\text{Cloud}}$ [14]. The total power consumption is given as

$$P_{\text{tot}} = \sum_{l=1}^{L} P_{\text{RU},l} + \sum_{l=1}^{L} z_l P_{\text{ONU}} + P_{\text{Cloud}} \quad (15)$$

where the binary variable $z_l$ indicates whether O-RU $l$ is active or not. If it is active, then $z_l = 1$ and otherwise $z_l = 0$. The O-RU power consumption for a particular O-RU $l$, i.e., $P_{\text{RU},l}$ is given by [14], [39]

$$
\begin{aligned}
P_{\text{RU},l} = {} & \\
& z_l \left( P_{\text{RU},0} + \Delta^{\text{tr}} \sum_{k=1}^{K} x_{k,l} p_{kl} + \mathbb{I}_s \left( P_{\text{RU},0}^{\text{proc}} + \Delta_{\text{RU}}^{\text{proc}} \frac{C_{\text{RU},l}}{C_{\text{RU}}^{\text{max}}} \right) \right) \\
= {} & z_l \left( P_{\text{RU},0} + \mathbb{I}_s P_{\text{RU},0}^{\text{proc}} + \Delta_{\text{RU}}^{\text{proc}} \frac{C_{\text{RU},l}}{C_{\text{RU}}^{\text{max}}} \right) + \Delta^{\text{tr}} \sum_{k=1}^{K} x_{k,l} p_{kl},
\end{aligned}
$$
$$(16)$$

where $P_{\text{RU},0}$ is the static power consumption of each O-RU when there is no transmission at the active mode and $\sum_{k=1}^{K} x_{k,l} p_{kl}$ is the transmit power of O-RU $l$. The load-dependent power consumption is modeled by the slope $\Delta^{\text{tr}}$. The binary indicator $\mathbb{I}_s \in \{0, 1\}$ specifies which functional splitting option is used. When it is one, Option 7.2 is utilized. In this case, the low-PHY functions between the dashed and solid lines in Fig. 1 are implemented in the processing units of the active O-RUs. Otherwise, when $\mathbb{I}_s = 0$, Option 8 is adopted and all the baseband processing is done in the O-Cloud. When $\mathbb{I}_s = 1$, the term $P_{\text{RU},0}^{\text{proc}} + \Delta_{\text{RU}}^{\text{proc}} \frac{C_{\text{RU},l}}{C_{\text{RU}}^{\text{max}}}$ in the above expression represents the processing power consumption using a load-dependent power consumption model [16]. $P_{\text{RU},0}^{\text{proc}}$ is the idle mode processing power consumption of each active processing unit at the O-RU corresponding to zero utilization. $\Delta_{\text{RU}}^{\text{proc}}$ is the slope of the load-dependent processing power consumption The maximum processing capacity at each O-RU is given by $C_{\text{RU}}^{\text{max}}$ in giga-operations per second (GOPS) depending on the used technology. The processing utilization at O-RU $l$ is given by $0 \leq C_{\text{RU},l} \leq C_{\text{RU}}^{\text{max}}$ in GOPS. Later,

we will compute the required GOPS. In the second equality in (16), we have simplified the equation by noting that when $z_l$ is zero, $x_{k,l}, \forall k$ become zero by definition. Similarly, when $\mathbb{I}_s = 0$, irrespective of whether O-RU $l$ is active or not, there is no processing power consumption and $C_{\text{RU},l}$ is automatically zero. In this case, the processing power is not included. It is worth mentioning that if an O-RU is not active, i.e., $z_l = 0$, we turn it off to save power.

The total power consumption in the O-Cloud is computed using the load-dependent GPP power consumption model from [16] as

$$
\begin{aligned}
P_{\text{Cloud}} = {} & P_{\text{fixed}} + \frac{1}{\sigma_{\text{cool}}} \left( P_{\text{OLT}} \sum_{w=1}^{W} w \ell_w \right. \\
& \left. + P_{\text{GPP},0}^{\text{proc}} \sum_{w=1}^{W} w d_w + \Delta_{\text{GPP}}^{\text{proc}} \frac{C_{\text{GPP}}}{C_{\text{GPP}}^{\text{max}}} \right), \quad (17)
\end{aligned}
$$

where $P_{\text{fixed}}$ is the load-independent fixed power consumption that includes the power consumption of the O-Cloud dispatcher, housing facilities, etc. The cooling efficiency is $0 < \sigma_{\text{cool}} \leq 1$. $P_{\text{OLT}}$ is the power consumption of each OLT module per GPP [14]. The binary variable $\ell_w \in \{0, 1\}$ is one if the LCs of $w$ number of GPPs are active and zero otherwise. Similarly, $d_w \in \{0, 1\}$ is one if $w$ number of GPPs are active for either processing or fronthaul connection to the O-RUs, and zero otherwise. This particular definition of the binary variables is to ensure that the constraints can be written as linear functions of the binary variables in the optimization problem we will consider in the next section. It is worth mentioning that the LC of an active GPP may be inactive if the corresponding GPP participates only in the processing that is redirected to it coming from other GPPs. $P_{\text{GPP},0}^{\text{proc}}$ is the idle mode processing power consumption of each active GPP in the O-Cloud corresponding to zero utilization. $\Delta_{\text{GPP}}^{\text{proc}}$ is the slope of the load-dependent processing power consumption of each GPP. For each GPP, the maximum processing capacity is given by $C_{\text{GPP}}^{\text{max}}$ in GOPS. The total processing utilization is given by $0 \leq C_{\text{GPP}} \leq W C_{\text{GPP}}^{\text{max}}$ in GOPS.

### A. GOPS Analysis of Digital Operations at the Radio Site and O-Cloud

In this section, we will analyze the GOPS for digital signal processing operations of a cell-free massive MIMO system. In the uplink, if Option 8 is used, after RF processing at the O-RUs, the quantized baseband signals are directly sent to the cloud. Then, at the GPPs, cyclic prefix (CP) removal and $N_{\text{DFT}}$-point DFT are performed to obtain frequency-domain signals. After resource element (RE) demapping, the remaining PHY and higher-layer operations are implemented for a particular UE. Similarly, in the downlink, after the higher-layer functions, the precoded signals are obtained. In the sequel, RE mapping, inverse DFT, and CP insertion are realized. Then, the time-domain signals are sent to the O-RUs and RF transmission takes place. If Option 7.2 is used, then the order of operations remains the same while the frequency-domain signals are sent after RE demapping to the O-Cloud in the uplink and the precoded frequency-domain signals are sent

from the O-Cloud to the O-RUs for the remaining low-PHY operations in the downlink as shown in Fig. 1.

To compute the total GOPS in O-RU $l$, $C_{\text{RU},l}$ and in the O-Cloud, $C_{\text{GPP}}$, we will mainly use the results from cellular massive MIMO [40], [41]. In [41], a factor two of overhead is taken in arithmetic operation calculations to account for memory operations. In the following GOPS calculations, we will also consider this by including a multiplication by two in each arithmetic operation. The first operation after RF processing is baseband filtering. Considering 10 taps with a polyphase filtering implementation, the corresponding complexity per O-RU is given in GOPS as $C_{\text{filter}} = 40Nf_s/10^9$ [41]. The next operation is DFT in the uplink and inverse DFT in the downlink, which has the complexity with fast Fourier transform (FFT) as $C_{\text{DFT}} = 8NN_{\text{DFT}}\log_2\left(N_{\text{DFT}}\right)/\left(T_s10^9\right)$, which is obtained by dividing the number of required real operations by the OFDM symbol duration $T_s$ [40].

The GOPS of the sample-based arithmetic operations such as precoding scale with $N_{\text{used}}/T_s$ [40]. An additional multiplying factor $\tau_d/\tau_c$ should be taken into account in precoding the downlink data since $\tau_d$ samples are precoded in each coherence block of length $\tau_c$. For the channel estimation, reciprocity calibration, and precoding computation, it scales with $N_{\text{used}}/(T_s\tau_c)$ since the corresponding operations are common for each sample in a coherence block of length $\tau_c$. To this end, from [5, Sec. 6.2.2], for LP-MMSE transmit precoding together with the required channel estimation of the served UEs by O-RU $l$ (and of the strongest UE per pilot), the GOPS (in terms of real multiplications) is computed as

$$
\begin{aligned}
C_{\text{prec},l} = &\underbrace{\frac{N_{\text{used}}}{T_s\tau_c10^9}\left(8N\tau_p^2 + 8N^2\left(\tau_p + \sum_{i=1}^{K}x_{i,l}\right)\right)}_{\text{Channel estimation}} \\
&+ \underbrace{\frac{N_{\text{used}}\tau_d}{T_s\tau_c10^9}8N\sum_{i=1}^{K}x_{i,l}}_{\text{Precoding}} + \underbrace{\frac{N_{\text{used}}}{T_s\tau_c10^9}8N\sum_{i=1}^{K}x_{i,l}}_{\text{Reciprocity calibration}} \\
&+ \underbrace{\frac{N_{\text{used}}}{T_s\tau_c10^9}\left((4N^2 + 4N)\tau_p + 8N^2\sum_{i=1}^{K}x_{i,l} + \frac{8\left(N^3 - N\right)}{3}\right)}_{\text{Precoding computation}}
\end{aligned}
$$

$$(18)$$

where we have also included the complexity of applying precoding and reciprocity calibration from [40]. It is worth mentioning that local precoding computation $C_{\text{prec},l}$ in (18) can be implemented at any other GPP $w$ different than the connected GPP to O-RU $l$ (benefiting from cloud sharing and virtualization via GPP dispatcher).

In addition to precoding, there are other GOPS regarding OFDM modulation/demodulation, mapping/demapping, channel coding, higher-layer control and network operations. These can be computed using the flexible power modeling in [42]. Let $C_{\text{other,O−RU}}$ and $C_{\text{other,UE}}$ denote the GOPS for the other operations which scale with the number of active O-RUs and the number of UEs that are served by each O-RU, respectively. There is also a fixed GOPS for UE operations, which are independent of the number of serving O-RUs. This

is represented by $\mathcal{F}$. The total GOPS in O-RU $l$ is given by

$$C_{\text{RU},l} = \mathbb{I}_s\underbrace{\left(C_{\text{filter}} + C_{\text{DFT}}\right)}_{\triangleq\mathcal{S}} = \mathbb{I}_s\mathcal{S}. \qquad (19)$$

If Option 8 is used, then $\mathbb{I}_s = 0$ and there is no digital signal processing implemented in the O-RUs. In this case, we include the corresponding low-PHY GOPS in the total GOPS of the O-Cloud as

$$
\begin{aligned}
C_{\text{GPP}} = \\
\sum_{l=1}^{L} z_l\left((1 - \mathbb{I}_s)\left(C_{\text{filter}} + C_{\text{DFT}}\right) + C_{\text{prec},l} + C_{\text{other,O−RU}}\right) \\
+ \sum_{l=1}^{L}\sum_{k=1}^{K} x_{k,l}C_{\text{other,UE}} + \mathcal{F} \\
= \mathcal{Z}\sum_{l=1}^{L} z_l + \mathcal{X}\sum_{l=1}^{L}\sum_{k=1}^{K} x_{k,l} + \mathcal{F}
\end{aligned}
\qquad (20)
$$

where the constant parameters $\mathcal{Z}$ and $\mathcal{X}$ are defined for ease of notation in the optimization problem. Note that when both the O-RU $l$ is active ($z_l = 1$) and Option 8 is used for functional splitting ($\mathbb{I}_s = 0$), the corresponding low-PHY GOPS is included in $C_{\text{GPP}}$.

## V. POWER-EFFICIENT O-RU SELECTION, GPP AND POWER ALLOCATION

In this section, we will introduce the proposed optimization problem that minimizes total power consumption. The aim is to decide which O-RUs serve which UEs, i.e., the binary variables $x_{k,l} \in \{0, 1\}$, the transmit powers allocated to the UEs, i.e., $p_{kl}$, which O-RUs are active and connected to the O-Cloud, i.e., $z_l \in \{0, 1\}$, and the number of active LCs and GPPs, i.e., $\ell_w, d_w \in \{0, 1\}$ in the O-Cloud. We note that the considered optimization problem is a mixed integer program since the O-RU selection together with power allocation and DU allocation is considered, which has a combinatorial nature. To express both the objective function and the constraints in a mixed binary linear or conic form, we introduce the following additional optimization variables:

$$\boldsymbol{\rho}_k = \left[\sqrt{p_{k1}}x_{k,1} \ \ldots \ \sqrt{p_{kL}}x_{k,L}\right]^{\text{T}} = \left[\rho_{k,1} \ \ldots \ \rho_{k,L}\right]^{\text{T}}, \quad (21)$$
$$\boldsymbol{\rho}'_l = \left[\sqrt{p_{1l}}x_{1,l} \ \ldots \ \sqrt{p_{Kl}}x_{K,l}\right]^{\text{T}} = \left[\rho_{1,l} \ \ldots \ \rho_{K,l}\right]^{\text{T}}. \qquad (22)$$

Due to the limited capacity of each wavelength in TWDM-PON, which is denoted by $R_{\max}$, we can assign at most

$$W_{\max} = \left\lfloor \frac{R_{\max}}{\mathbb{I}_s R_{\text{fronthaul}}^{(7.2)} + (1 - \mathbb{I}_s)R_{\text{fronthaul}}^{(8)}} \right\rfloor \qquad (23)$$

O-RUs to each wavelength and, hence, to each GPP $w$, for $w = 1, \ldots, W$. When $\mathbb{I}_s = 1$, the functional split option 7.2 is used with the corresponding required fronthaul data rate $R_{\text{fronthaul}}^{(7.2)}$. On the other hand, when $\mathbb{I}_s = 0$, the functional split option 8 is used with the required fronthaul data rate $R_{\text{fronthaul}}^{(8)}$.

In the considered network power consumption minimization problem, we assume each UE $k$ has a SE request with the corresponding minimum SINR requirement $\gamma_k$. Hence, we

have QoS constraints in the form of $\mathsf{SINR}_k \geq \gamma_k$ for each UE $k$. The optimization problem can be cast using the introduced optimization variables as

$$
\begin{aligned}
\underset{\substack{z_l, x_{k,l}, \forall k, \forall l \\ \ell_w, d_w, \forall w \\ \boldsymbol{\rho}_k, \forall k}}{\text{minimize}} \quad & P_{\text{fixed}} + (P_{\text{RU},0} + P_{\text{ONU}}) \sum_{l=1}^{L} z_l
\end{aligned}
$$

$$
+ \left( \mathbb{I}_s \left( P_{\text{RU},0}^{\text{proc}} + \frac{\Delta_{\text{RU}}^{\text{proc}} \mathcal{S}}{C_{\text{RU}}^{\text{max}}} \right) + \frac{\Delta_{\text{GPP}}^{\text{proc}} \mathcal{Z}}{\sigma_{\text{cool}} C_{\text{GPP}}^{\text{max}}} \right) \sum_{l=1}^{L} z_l
$$

$$
+ \Delta^{\text{tr}} \sum_{l=1}^{L} \sum_{k=1}^{K} \rho_{k,l}^2 + \frac{P_{\text{OLT}}}{\sigma_{\text{cool}}} \sum_{w=1}^{W} w\ell_w + \frac{P_{\text{GPP},0}^{\text{proc}}}{\sigma_{\text{cool}}} \sum_{w=1}^{W} wd_w
$$

$$
+ \frac{\Delta_{\text{GPP}}^{\text{proc}} \mathcal{X}}{\sigma_{\text{cool}} C_{\text{GPP}}^{\text{max}}} \sum_{l=1}^{L} \sum_{k=1}^{K} x_{k,l} + \frac{\Delta_{\text{GPP}}^{\text{proc}} \mathcal{F}}{\sigma_{\text{cool}} C_{\text{GPP}}^{\text{max}}} \tag{24a}
$$

subject to:

$$
\frac{|\mathbf{b}_k^{\text{T}} \boldsymbol{\rho}_k|^2}{\sum_{i=1}^{K} \boldsymbol{\rho}_i^{\text{T}} \mathbf{C}_{ki} \boldsymbol{\rho}_i + \sigma^2} \geq \gamma_k, \quad \forall k, \tag{24b}
$$

$$
\sum_{l=1}^{L} z_l \leq W_{\text{max}} W, \tag{24c}
$$

$$
\frac{\sum_{k=1}^{K} x_{k,l}}{K} \leq z_l \leq \sum_{k=1}^{K} x_{k,l}, \quad \forall l, \tag{24d}
$$

$$
\sum_{w=1}^{W} (w-1)\ell_w \leq \frac{\sum_{l=1}^{L} z_l}{W_{\text{max}}} \leq \sum_{w=1}^{W} w\ell_w, \tag{24e}
$$

$$
\mathcal{Z} \sum_{l=1}^{L} z_l + \mathcal{X} \sum_{l=1}^{L} \sum_{k=1}^{K} x_{k,l} + \mathcal{F} \leq C_{\text{GPP}}^{\text{max}} \sum_{w=1}^{W} wd_w, \tag{24f}
$$

$$
\sum_{w=1}^{W} \ell_w = 1, \quad \sum_{w=1}^{W} d_w = 1, \tag{24g}
$$

$$
\sum_{w=1}^{W} w\ell_w \leq \sum_{w=1}^{W} wd_w, \tag{24h}
$$

$$
0 \leq \rho_{k,l} \leq \sqrt{p_{\text{max}}} x_{k,l}, \quad \forall k, \forall l, \tag{24i}
$$

$$
\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\text{max}}} z_l, \quad \forall l, \tag{24j}
$$

$$
z_l, \ell_w, d_w, x_{k,l} \in \{0,1\}, \quad \forall k, \forall l, \forall w. \tag{24k}
$$

The constraints in (24b) are to guarantee that each UE's minimum SINR requirement is satisfied. The constraint in (24c) guarantees that the number of active O-RUs is determined not to exceed the maximum allowable number determined by the fronthaul limitations of the given functional splitting option. The constraints in (24d) relate the binary variables $x_{k,l}$ and $z_l$, i.e., an O-RU is active if and only if it serves at least one UE. The constraint in (24e) connects the number of active O-RUs to the the number of required active LCs. The constraint in (24f) guarantees that the total GOPS does not exceed the processing capability of active GPPs in the O-Cloud. The constraints in (24g) are to satisfy that $\ell_w$ and $d_w$ are only one for one value of $w$ since these binary variables are one when the number of active LCs or GPPs is equal to their sub-index. The constraint in (24h) is to ensure that the number of active GPPs is always greater than or equal to the number of

active LCs. The constraints in (24i) guarantee that the square root of the power coefficient for UE $k$ and O-RU $l$ is zero if UE $k$ is not served by O-RU $l$. Here, $p_{\text{max}}$ is the maximum transmit power budget of each O-RU and when $x_{k,l} = 1$, this constraint does not limit $\rho_{k,l}$. (24j) represents the per-O-RU transmit power constraints. Finally, the constraints in (24k) specify which variables are binary.

Note that the SINR constraints in (24b) can be re-written in second-order cone form [5, Sec. 7.1.2]. As a result, the overall optimization problem is a mixed binary second-order cone programming problem, which has a convex structure except for the binary constraints. Hence, the global optimum can be obtained by the branch-and-bound algorithm [43]. It is known that the complexity grows exponentially with the number of discrete variables, which in our case means the number of O-RUs, UEs, and GPPs. In the next part, we will develop a lower complexity solution by approximating the problem and using concave programming.

### A. Approximate Optimization Problem Formulation Using $l_0$ Norm Minimization

In devising a lower complexity algorithm, the main step is to eliminate the binary variables in the optimization problem. To this end, we will keep the SINR constraints and some other few constraints as they are while reflecting the other constraints in the objective function by using novel approximations that allow elimination of several binary variables and obtaining a concave objective function with convex constraints at the end. Let us go through all the constraints in (24b)-(24k) in the sequel. First, we express the SINR constraints in (24b) in second-order cone form as

$$
\left\| \begin{bmatrix} \sqrt{\gamma_k} \mathbf{C}_{k1}^{\frac{1}{2}} \boldsymbol{\rho}_1 \\ \vdots \\ \sqrt{\gamma_k} \mathbf{C}_{kK}^{\frac{1}{2}} \boldsymbol{\rho}_K \\ \sqrt{\gamma_k} \sigma \end{bmatrix} \right\| \leq \mathbf{b}_k^{\text{T}} \boldsymbol{\rho}_k, \quad \forall k. \tag{25}
$$

To simplify the problem and make it manageable, we will not consider the constraints in (24c)-(24h) and the binary constraints in (24k). Hence, the variables $\ell_w$ and $d_w$ do not exist in the modified problem. To reflect their impact on the objective function, we will approximate the values of $\sum_{w=1}^{W} w\ell_w$ and $\sum_{w=1}^{W} wd_w$ in the objective function (24a) using the constraints in (24e)-(24h) as follows:

$$
\sum_{w=1}^{W} w\ell_w \approx \frac{\sum_{l=1}^{L} z_l}{W_{\text{max}}} = \frac{\|\mathbf{z}\|_0}{W_{\text{max}}}, \tag{26}
$$

$$
\sum_{w=1}^{W} wd_w \approx \max \left( \sum_{w=1}^{W} w\ell_w, \frac{\mathcal{Z}\|\mathbf{z}\|_0 + \mathcal{X}\sum_{k=1}^{K} \|\boldsymbol{\rho}_k\|_0 + \mathcal{F}}{C_{\text{GPP}}^{\text{max}}} \right) \tag{27}
$$

where $\mathbf{z} = [z_1 \ \dots \ z_L]^{\text{T}}$ and we have used the fact that $\sum_{l=1}^{L} x_{k,l} = \|\boldsymbol{\rho}_k\|_0$. To write $\sum_{w=1}^{W} wd_w$ in the objective function of the approximate modified problem, we can eliminate the maximum operation and approximate $\sum_{w=1}^{W} wd_w$ by

using the upper bound of (27) as

$$\sum_{w=1}^{W} w d_w \approx \sum_{w=1}^{W} w \ell_w + \frac{\mathcal{Z}\|\mathbf{z}\|_0 + \mathcal{X}\sum_{k=1}^{K}\|\boldsymbol{\rho}_k\|_0 + \mathcal{F}}{C_{\mathrm{GPP}}^{\max}}$$

$$\approx \frac{\|\mathbf{z}\|_0}{W_{\max}} + \frac{\mathcal{Z}\|\mathbf{z}\|_0 + \mathcal{X}\sum_{k=1}^{K}\|\boldsymbol{\rho}_k\|_0 + \mathcal{F}}{C_{\mathrm{GPP}}^{\max}}. \quad (28)$$

Keeping the constraints (24j) and the left-side constraint in (24i), and neglecting the fixed part of the objective function, the approximated optimization problem in terms of reduced number of variables can be expressed as

$$\underset{\mathbf{z},\boldsymbol{\rho}_k,\forall k}{\text{minimize}} \quad (P_{\mathrm{RU},0} + P_{\mathrm{ONU}})\|\mathbf{z}\|_0$$

$$+ \left(\mathbb{I}_s\left(P_{\mathrm{RU},0}^{\mathrm{proc}} + \frac{\Delta_{\mathrm{RU}}^{\mathrm{proc}}\mathcal{S}}{C_{\mathrm{RU}}^{\max}}\right) + \frac{\Delta_{\mathrm{GPP}}^{\mathrm{proc}}\mathcal{Z}}{\sigma_{\mathrm{cool}}C_{\mathrm{GPP}}^{\max}}\right)\|\mathbf{z}\|_0$$

$$+ \Delta^{\mathrm{tr}}\sum_{l=1}^{L}\sum_{k=1}^{K}\rho_{k,l}^2 + \frac{P_{\mathrm{OLT}} + P_{\mathrm{GPP},0}^{\mathrm{proc}}}{\sigma_{\mathrm{cool}}}\frac{\|\mathbf{z}\|_0}{W_{\max}}$$

$$+ \frac{P_{\mathrm{GPP},0}^{\mathrm{proc}}}{\sigma_{\mathrm{cool}}C_{\mathrm{GPP}}^{\max}}\left(\mathcal{Z}\|\mathbf{z}\|_0 + \mathcal{X}\sum_{k=1}^{K}\|\boldsymbol{\rho}_k\|_0\right)$$

$$+ \frac{\Delta_{\mathrm{GPP}}^{\mathrm{proc}}\mathcal{X}}{\sigma_{\mathrm{cool}}C_{\mathrm{GPP}}^{\max}}\sum_{k=1}^{K}\|\boldsymbol{\rho}_k\|_0 \quad (29a)$$

subject to:

$$\left\|\begin{bmatrix}\sqrt{\gamma_k}\mathbf{C}_{k1}^{\frac{1}{2}}\boldsymbol{\rho}_1 \\ \vdots \\ \sqrt{\gamma_k}\mathbf{C}_{kK}^{\frac{1}{2}}\boldsymbol{\rho}_K \\ \sqrt{\gamma_k}\sigma\end{bmatrix}\right\| \leq \mathbf{b}_k^{\mathrm{T}}\boldsymbol{\rho}_k, \quad \forall k, \quad (29b)$$

$$\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\max}}z_l, \quad \forall l \quad (29c)$$

$$\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\max}}, \quad \forall l, \quad \rho_{k,l} \geq 0, \ \forall k, \forall l \quad (29d)$$

where we have included $\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\max}}$ in (29d) to guarantee the per-O-RU power constraints are satisfied even if $z_l > 1$ for some $l$.

Note that $l_0$ norm is not a real norm and $\|.\|_0$ is not a convex function. Hence, the optimization problem in (29) is not convex. One way to solve this problem in an efficient way is to replace the $l_0$ norm by a more tractable function. After evaluating several methods, we have observed that the following approach leads to decent performance with nice convergence properties in the considered iterative algorithm. The following continuously differentiable concave function can be used in place of $\|\mathbf{z}\|_0$:

$$f(\mathbf{z}) = \sum_{l=1}^{L}(1 - e^{-\alpha z_l}) \quad (30)$$

with $\alpha > 0$ [44]. It can be seen that as $\alpha \to \infty$, $f(\mathbf{z}) \to \|\mathbf{z}\|_0$. Using the concave function $f(\mathbf{z})$, the non-convex problem in (29) can be approximated as

$$\underset{\mathbf{z},\boldsymbol{\rho}_k,\forall k}{\text{minimize}} \quad \mathcal{A}_z f(\mathbf{z}) + \mathcal{A}_\rho \sum_{k=1}^{K} f(\boldsymbol{\rho}_k) + \Delta^{\mathrm{tr}}\sum_{l=1}^{L}\sum_{k=1}^{K}\rho_{k,l}^2$$

$$(31a)$$

subject to $(29b) - (29d)$ $\quad (31b)$

---

**Algorithm 1** CCP algorithm for solving the problem in (31).

1: **Initialization:**
  - Set $\mathbf{z}^{(0)}$ and $\boldsymbol{\rho}_k^{(0)}$, for $k = 1, \dots, K$, with arbitrary positive entries and the solution accuracy $\varepsilon > 0$.
  - Set the iteration counter $t = 0$.
2: **repeat**
3:     $t \leftarrow t + 1$.
4:     Solve the convex problem

$$\underset{\mathbf{z},\boldsymbol{\rho}_k,\forall k}{\text{minimize}} \quad \Delta^{\mathrm{tr}}\sum_{l=1}^{L}\sum_{k=1}^{K}\rho_{k,l}^2 + \mathcal{A}_z\nabla f\left(\mathbf{z}^{(t-1)}\right)^{\mathrm{T}}\mathbf{z}$$

$$+ \mathcal{A}_\rho\sum_{k=1}^{K}\nabla f\left(\boldsymbol{\rho}_k^{(t-1)}\right)^{\mathrm{T}}\boldsymbol{\rho}_k \quad (32a)$$

    subject to $(29b) - (29d)$. $\quad (32b)$

5:     Set $\mathbf{z}^{(t)}, \boldsymbol{\rho}_1^{(t)}, \dots, \boldsymbol{\rho}_K^{(t)}$ as the solution of (32).
6: **until** the normalized squared error difference between the current and previous objective functions in (31a) is less than $\varepsilon$.
7: **Output:** $\mathbf{z}^{(t)}, \boldsymbol{\rho}_1^{(t)}, \dots, \boldsymbol{\rho}_K^{(t)}$.

---

where the constants $\mathcal{A}_i$, for $i \in \{z, \rho\}$ in the objective function are obtained by summing the terms that multiply the corresponding $\|\cdot\|_0$ terms in (29a). The above problem is in concave programming form. It is not convex, but the so-called *concave-convex procedure (CCP)* outlined in Algorithm 1 can be used, where the concave objective function is convexified around the previous solution and a convex problem is solved at each iteration. From [45, Thm. 4], Algorithm 1 converges to a stationary point of the problem in (31) under suitable constraint qualification when (31) is feasible.

The solution found by Algorithm 1 can further be improved by enforcing more sparsity, and thus less number of active O-RUs, via the refinement steps outlined in Algorithm 2. First, the power coefficients $\rho_{k,l}$, whose values normalized by the maximum power coefficient are smaller than the threshold $0 < \zeta \ll 1$, are set to zero. Then, the number of active O-RUs is determined by checking the non-zero power coefficients. Next, the problem of minimizing radio site power consumption under the SINR and per-O-RU power constraints given in (34) is solved. The number of active O-RUs is iteratively reduced according to the O-RU transmit powers found previously until this problem is infeasible. Once infeasibility is detected, the lastly found power coefficients are returned as the output of Algorithm 2. Here, the threshold $\zeta$ should be sufficiently small to ensure that the problem in (34) is feasible at the first iteration, but also fine-tuned in a way to eliminate unnecessarily small power coefficients that do not affect the feasibility of the power allocation problem.

**Algorithm 2** Refinement algorithm for improving the solution found in Algorithm 1.

1: **Initialization:**
- Set $\boldsymbol{\rho}_k^{(0)}$, for $k = 1, \ldots, K$, as the output of Algorithm 1.
- Compute $\overline{\rho} = \max\limits_{l=1,\ldots,L, k=1,\ldots,K} \rho_{k,l}^{(0)}$.
- Set the power coefficients $\rho_{k,l}^{(0)}$, which are sufficiently small so that $\rho_{k,l}^{(0)}/\overline{\rho} \leq \zeta$, to zero, where $0 < \zeta \ll 1$ is the threshold parameter.
- Determine the number of active O-RUs, denoted by $L_{\text{active}}^{(0)}$, as the number of O-RUs $l$ so that

$$\sum_{k=1}^{K} \left( \rho_{k,l}^{(0)} \right)^2 > 0. \tag{33}$$

- Set the iteration counter $t = 0$.

2: **repeat**
3:      $t \leftarrow t + 1$.
4:      Solve the convex problem

$$\underset{\boldsymbol{\rho}_k, \forall k}{\text{minimize}} \quad \Delta^{\text{tr}} \sum_{l=1}^{L} \sum_{k=1}^{K} \rho_{k,l}^2 \tag{34a}$$

subject to:

$$(29b) \tag{34b}$$

$$\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\max}}, \quad \forall l, \quad \rho_{k,l} \geq 0, \quad \forall k, \quad \forall l \tag{34c}$$

$$\rho_{k,l} = 0, \quad \text{for } \rho_{k,l}^{(t-1)} = 0, \quad \forall k, \quad \forall l. \tag{34d}$$

5:
- If the problem is feasible, set $\boldsymbol{\rho}_1^{(t)}, \ldots, \boldsymbol{\rho}_K^{(t)}$ as the solution of (34).
- Reduce the number of active O-RUs by one, i.e., $L_{\text{active}}^{(t)} = L_{\text{active}}^{(t-1)} - 1$.
- Set $\rho_{k,l}^{(t)}$ to zero for the $L - L_{\text{active}}^{(t)}$ O-RUs that have the least total transmit powers $\sum_{k=1}^{K} \left( \rho_{k,l}^{(t)} \right)^2$.

6: **until** An infeasible solution is obtained or $L_{\text{active}}^{(t)} = 1$.
7: **Output:** $\boldsymbol{\rho}_1^{(t)}, \ldots, \boldsymbol{\rho}_K^{(t)}$.

## VI. JOINT SUM-SE MAXIMIZATION AND POWER CONSUMPTION MINIMIZATION

In this section, we will consider a multiobjective optimization problem that aims to jointly minimize the total end-to-end power consumption and maximize the sum-SE of all the UEs. To this end, we remove the SINR constraints in (29b) from the problem (31) and include the sum-SE with a certain weight to the objective function. The considered problem can be cast as

$$\underset{\mathbf{z}, \boldsymbol{\rho}_k, \forall k}{\text{minimize}} \quad \Delta^{\text{tr}} \sum_{l=1}^{L} \sum_{k=1}^{K} \rho_{k,l}^2 + \mathcal{A}_z f(\mathbf{z}) + \mathcal{A}_\rho \sum_{k=1}^{K} f(\boldsymbol{\rho}_k)$$

$$- \lambda \sum_{k=1}^{K} \ln \left( 1 + \frac{|\mathbf{b}_k^{\text{T}} \boldsymbol{\rho}_k|^2}{\sum_{i=1}^{K} \boldsymbol{\rho}_i^{\text{T}} \mathbf{C}_{ki} \boldsymbol{\rho}_i + \sigma^2} \right) \tag{35a}$$

subject to:

$$\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\max}} z_l, \quad \forall l \tag{35b}$$

$$\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\max}}, \forall l, \quad \rho_{k,l} \geq 0, \forall k, \forall l \tag{35c}$$

where the parameter $\lambda > 0$ determines the weighting of the sum-SE term in the multi-objective function. We have used $\ln(\cdot)$ instead of $\log_2(\cdot)$ for simplicity without loss of generality since only a constant factor differs between them. Moreover, the pre-log factor in (9) is not considered since its effect can be absorbed into the penalty parameter $\lambda$. Weighted MMSE is a common technique to find local optima of sum-SE maximization problems under power constraints [46]. However, here our problem differs from the traditional sum-SE maximization structure since it also includes the total power consumption in the objective function, which is in a non-convex form. In this paper, we will devise an alternative algorithm to take advantage of CCP that we have already constructed for the power minimization problem in the previous section. We first utilize the transform of $\ln(1 + \text{SINR}_k)$ to an equivalent form from [47, Thm. 3]. Using the newly defined optimization variables $\chi_k$, for $k = 1, \ldots, K$, the optimization problem in (35) can be equivalently (in terms of the optimal solution) expressed as

$$\underset{\mathbf{z}, \boldsymbol{\rho}_k, \chi_k, \forall k}{\text{minimize}} \quad \Delta^{\text{tr}} \sum_{l=1}^{L} \sum_{k=1}^{K} \rho_{k,l}^2 + \mathcal{A}_z f(\mathbf{z}) + \mathcal{A}_\rho \sum_{k=1}^{K} f(\boldsymbol{\rho}_k)$$

$$- \lambda \left( \sum_{k=1}^{K} \ln(1 + \chi_k) - \sum_{k=1}^{K} \chi_k \right)$$

$$- \lambda \sum_{k=1}^{K} (1 + \chi_k) \frac{|\mathbf{b}_k^{\text{T}} \boldsymbol{\rho}_k|^2}{\sum_{i=1}^{K} \boldsymbol{\rho}_i^{\text{T}} \mathbf{C}_{ki} \boldsymbol{\rho}_i + |\mathbf{b}_k^{\text{T}} \boldsymbol{\rho}_k|^2 + \sigma^2} \tag{36a}$$

subject to:

$$\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\max}} z_l, \quad \forall l \tag{36b}$$

$$\|\boldsymbol{\rho}_l'\| \leq \sqrt{p_{\max}}, \forall l, \quad \rho_{k,l} \geq 0, \forall k, \forall l \tag{36c}$$

where the equivalency can be shown by taking the derivative of the objective function with respect to $\chi_k$, for $k = 1, \ldots, K$, equating them to zero, and inserting the optimal $\chi_k$ to the objective function. Next, we will develop a further novel transformation and integrate CCP to solve the resulting problem. The first step is to introduce the optimization variables $u_k$ and $r_k$, for $k = 1, \ldots, K$, to represent upper bounds to $(1 + \chi_k)^{-1}$ and $\left( \sum_{i=1}^{K} \boldsymbol{\rho}_i^{\text{T}} \mathbf{C}_{ki} \boldsymbol{\rho}_i + |\mathbf{b}_k^{\text{T}} \boldsymbol{\rho}_k|^2 + \sigma^2 \right) / (1 + \chi_k)$

in (36a), respectively. We then re-cast the problem in (36) as

$$\underset{\mathbf{z},\boldsymbol{\rho}_k,\chi_k,u_k,r_k,\forall k}{\text{minimize}} \quad \Delta^{\text{tr}} \sum_{l=1}^{L} \sum_{k=1}^{K} \rho_{k,l}^2 + \mathcal{A}_z f(\mathbf{z}) + \mathcal{A}_\rho \sum_{k=1}^{K} f(\boldsymbol{\rho}_k)$$

$$+ \lambda \left( \sum_{k=1}^{K} \ln(u_k) + \sum_{k=1}^{K} \chi_k - \sum_{k=1}^{K} \frac{(\mathbf{b}_k^{\text{T}}\boldsymbol{\rho}_k)^2}{r_k} \right) \tag{37a}$$

subject to:

$$\|\boldsymbol{\rho}_l'\| \le \sqrt{p_{\max}} z_l, \quad \forall l \tag{37b}$$

$$\|\boldsymbol{\rho}_l'\| \le \sqrt{p_{\max}}, \ \forall l, \quad \rho_{k,l} \ge 0, \ \forall k, \forall l \tag{37c}$$

$$\left\| \begin{bmatrix} \sqrt{2}\mathbf{C}_{k1}^{\frac{1}{2}}\boldsymbol{\rho}_1 \\ \vdots \\ \sqrt{2}\mathbf{C}_{kK}^{\frac{1}{2}}\boldsymbol{\rho}_K \\ \sqrt{2}\mathbf{b}_k^{\text{T}}\boldsymbol{\rho}_k \\ \sqrt{2}\sigma \\ (1+\chi_k) \\ r_k \end{bmatrix}^{\text{T}} \right\| \le 1 + \chi_k + r_k, \quad \forall k \tag{37d}$$

$$\left\| \begin{bmatrix} (1+\chi_k) \\ u_k \\ \sqrt{2} \end{bmatrix} \right\| \le 1 + \chi_k + u_k, \quad \forall k \tag{37e}$$

where we have used second-order cone constraints to construct the inequalities $(1+\chi_k)^{-1} \le u_k$ and $\left( \sum_{i=1}^{K} \boldsymbol{\rho}_i^{\text{T}} \mathbf{C}_{ki}\boldsymbol{\rho}_i + |\mathbf{b}_k^{\text{T}}\boldsymbol{\rho}_k|^2 + \sigma^2 \right)/(1+\chi_k) \le r_k$. The following lemma demonstrates the equivalency between the problems (36) and (37) in terms of the optimal solutions.

**Lemma 1.** *The optimal values of* $\{\mathbf{z}, \boldsymbol{\rho}_k, \forall k\}$ *are the same for the problems* (36) *and* (37).

*Proof.* To prove the claim, it is enough to show that the constraints in (37d)-(37e) are satisfied with equality at the optimal solution. We can prove this by contradiction. Assume that for the optimal solution of (37), at least one of the constraints in (37d) is satisfied with strict inequality. Then, we can reduce the value of respective $r_k$ until the corresponding constraint is satisfied with equality without violating any other constraint. In this way, we also improve the objective value, which contradicts that the initial $r_k$ is optimum. Hence, all the constraints in (37d) are satisfied with equality. Moreover, assume that at least one of the constraints in (37e) is satisfied with strict inequality, which leads to $(1+\chi_k)^{-1} < u_k$, for the corresponding $k$. Then, we can reduce the value of $u_k$ until $(1+\chi_k)^{-1} = u_k$ without violating any constraints, and improving the objective function. This contradicts that $(1+\chi_k)^{-1} < u_k$ for at least one $k$. Hence, all the constraints in (37e) are satisfied with equality. Then, by inserting the value of $r_k$ and $(1+\chi_k)^{-1} = u_k$ into the objective function, the problems (36) and (37) can be shown in identical form. $\square$

We note that the objective function in (37a) is a summation of a convex and concave function explained as follows. The first part of the objective function, i.e., $\Delta^{\text{tr}} \sum_{l=1}^{L} \sum_{k=1}^{K} \rho_{k,l}^2$ is a convex function. The term $\mathcal{A}_z f(\mathbf{z}) + \mathcal{A}_\rho \sum_{k=1}^{K} f(\boldsymbol{\rho}_k)$ is a concave function. Let's go through the remaining terms one by one. We note that $\ln(u_k)$ is a concave function, which

can be shown by the sign of its second derivative. The last term $-\frac{(\mathbf{b}_k^{\text{T}}\boldsymbol{\rho}_k)^2}{r_k}$ is a concave function of $\boldsymbol{\rho}_k$ and $r_k$ since $g(t_k, r_k) = \frac{t_k^2}{r_k}$, which is a quadratic-over-linear function, is convex in terms of $t_k, r_k$, where $t_k = \mathbf{b}_k^{\text{T}}\boldsymbol{\rho}_k$. We also note that the constraints in (37d)-(37e) are convex, and, thus, we can use CCP by solving a convexified problem at each iteration. Furthermore, the convexified problem has a quadratic objective function with second-order cone constraints, which enables the use of efficient algorithms thanks to our unique reformulation. The steps of the respective algorithm are similar to those of Algorithm 1, so we do not repeat it for brevity.

### A. End-to-end vs. local coordination-based vs. radio-only resource orchestration

The optimization problems that have been considered so far provide the user-centric O-RU clusters, active O-RUs, and the respective transmit power allocation. Based on that, the processing requirements in terms of GOPS and the required number of active LCs (TWDM-PON wavelengths) and the number of GPPs are computed for three different resource orchestration schemes.

(i) Fully virtualized end-to-end resource orchestration: The fronthaul resources and connections between the O-RUs, LCs, and GPPs are assumed to be fully virtualized and jointly orchestrated so that the minimal number of LCs and GPPs are activated. Number of LCs is set as the smallest integer that is greater than or equal to the number of active O-RUs divided by the number of O-RUs that each LC can serve. The number of active GPPs, which affect the GPP idle power, is also determined as the smallest possible value in the O-Cloud.

(ii) Local coordination-based resource orchestration: This is not a fully virtualized system, where the fronthaul resources for each O-RU are fixed. The number of active LCs is determined based on the dedicated optical wavelengths to the active O-RUs. Although full virtualization is assumed for UE-specific operations in the GPPs, the O-RU-specific operations are assumed to be executed in the dedicated GPPs to the fixed optical LC and wavelength serving a particular O-RU. The respective processing resources are virtualized among the GPPs allocated to the O-RUs that share the same optical wavelength. Hence, an increase in the fronthaul and processing idle power consumption is expected in the local-coordination-based resource orchestration. In such a system, the only power-saving mechanism is shutting down unused O-RUs and LCs, and partial virtualization within GPPs connected to each LC.

(iii) Radio-only scheme: The fronthaul and cloud resources are fixed. Although the load-dependent and O-RU powers facilitate the same sparsity induced by the proposed algorithms, the number of active LCs is selected according to the fixed fronthaul connections, and the number of active GPPs is selected under the peak traffic assumption in which all O-RUs are active. The only power-saving mechanism is shutting down unused O-RUs and LCs, which explains why we call this scheme "radio-only" resource allocation.

It is worth pointing out that in all three resource allocation schemes, the GPP and O-RU load-dependent power consumption is the same, where the difference lies in the GPP idle

power and fronthaul power consumption. The second and third schemes are developed as benchmarks.

## VII. NUMERICAL RESULTS AND DISCUSSION

In this part, we quantify the end-to-end power consumption of a cell-free massive MIMO system in the O-RAN architecture to gain an understanding of the impact of joint radio, fronthaul, and cloud resource allocation, the selection of the functional splitting option, and the most power-consuming network components. Three resource orchestration schemes, which are described in Section VI.A, are considered: i) fully virtualized end-to-end; ii) local coordination-based; and iii) radio-only resource allocation.

The simulation parameters are outlined in Table I, and they are mainly set from the works [13], [14], [16], [29], [39], [40], [48]. In particular, we consider pico-cell power parameters from [39]. The GOPS/Watt for each of $W$ GPPs in the cloud and the processing unit of each O-RU is 2.434 according to 2x Intel Xeon E5-2683 v4 processor from [48, Tab. 1]. The idle power $P_{\mathrm{GPP},0}^{\mathrm{proc}} = P_{\mathrm{RU},0}^{\mathrm{proc}}$ and the slope $\Delta_{\mathrm{GPP}}^{\mathrm{proc}} = \Delta_{\mathrm{RU}}^{\mathrm{proc}}$ are scaled linearly such that each GPP in the cloud and O-RU processing unit has $C_{\mathrm{GPP}}^{\max} = C_{\mathrm{RU}}^{\max} = 180\,\mathrm{GOPS}$ as in [16]. The deployment and radio site parameters are as in the running example of [5, Sec. 5.3] with $f_c = 2\,\mathrm{GHz}$ except for the parameters that are listed in Table I. The functional split options 8 and 7.2 are denoted as FS-8 and FS-7.2 in the simulations.

The provided parameters are treated as constants in the optimization problems (24), (31), and (35). Once the solutions to the optimization problems are found, the set of active O-RUs, the O-RU cluster for each UE, and transmit power allocation are determined by checking the non-zero power coefficients $\rho_{k,l}$ found by each method. The required GOPS can then be computed using the expression (20) according to the active O-RUs, clusters, and functional splitting option. The number of active LCs and GPPs are set differently for each resource orchestration scheme as explained in detail in Section VI.A. In the end, the total power consumption in (24a) is computed using the fixed parameters, the output of the optimization problems, and the virtualization level of the considered resource orchestration.

In the CCP-based algorithms, five random initializations are considered and the best solution among them is noted. For the non-fully virtualized resource allocation schemes, for each setup, five random fixed TWDM-PON wavelength and LC connections are assumed and the averages of the power consumption values are presented in the figures. For the solution of (31), the $\alpha$ parameter in (30) is selected as $7 - \mathrm{SE_{req}}$, where $\mathrm{SE_{req}}$ is the required SE value for all the UEs and this selection is empirically based on the several experiments and the observation that larger $\alpha$ is needed for smaller SE requirement. For the problem in (35), $\alpha = 3$ is selected. For the algorithms, a minimum and a maximum iteration number of 10 and 50 are set, respectively. The solution accuracy parameter in Algorithm 1 and the threshold in Algorithm 2 are selected as $\varepsilon = 10^{-5}$ and $\zeta = 10^{-3}$, respectively.

First, we consider a setup with $L = 16$ O-RUs, $K = 8$ UEs, and $W = 4$ GPPs with FS-8 and find the optimal end-to-end resource allocation by solving the mixed-binary second-order cone programming problem in (24). In addition to three resource allocation schemes with cell-free massive MIMO, we also consider a conventional small-cell system where each UE is only served by one O-RU with end-to-end resource allocation. To this end, we consider the same virtualized O-RAN architecture and obtain the power-optimal resource allocation by solving (24) for the small-cell system by adding an additional constraint $\sum_{l=1}^{L} x_{k,l} = 1, \forall k$ to guarantee that only one O-RU is transmitting data to each UE. Hence, the problem for small-cell is more restrictive than its cell-free counterpart.

We consider 30 random O-RU and UE locations. For each random setup, we take the average of five random permutations of the fixed fronthaul wavelength assignments in the local coordination-based and radio-only resource allocation schemes. Fig. 3 shows the average total power consumption for a given SE requirement that is assumed to be the same for every UE. We consider the same random setups for all the methods. However, due to the SINR constraints in (24b), the optimization problem is not guaranteed to be feasible. In the figure, the average is taken out of all feasible setups at each point and it is only plotted when the feasibility ratio is greater than 50%. As the plot shows, when the SE requirement is greater than 1.75, the small-cell system cannot guarantee reliable performance due to infeasibility. On the other hand, the cell-free system benefits from user-centric JT to support the UEs with much higher SEs. For small SE values, the power consumption is almost the same for both systems. The reason is that for some setups, the cell-free optimization problem results in small-cell solution making their power consumption the same. However, as the SE increases, more O-RUs and DUs are activated to serve UEs when using the small-cell system. This results in increased power consumption compared to cell-free massive MIMO, as shown in the figure for the SE range $[0.5, 1.75]$. The maximum power saving is 14% when the SE requirement is 1.25 bit/s/Hz. In conclusion, cell-free massive MIMO results in less or equal power consumption for a small-cell system to guarantee a certain SE requirement. Based on our simulations, we observe that cell-free massive MIMO provides around 1.7 times more rate to the UEs with almost the same energy per bit in comparison to the small-cell system.

When we compare the total power consumption of the virtualized cell-free massive MIMO system with the other two orchestration schemes, we see a consistent power saving achieved by the former one for all the considered SE values. 30% and 17% saving are obtained when the SE requirement is 1.25 bit/s/Hz over radio-only and local coordination-based resource allocation, respectively. To better understand what contributes most to the achieved power saving, we show the average power consumption breakdown for the particular SE requirement of 1.25 bit/s/Hz in Fig. 4. Virtualized cell-free massive MIMO with end-to-end resource allocation provides reduced power consumption in comparison to the small-cell system due to the reduced RAN, fronthaul, and cloud pro-

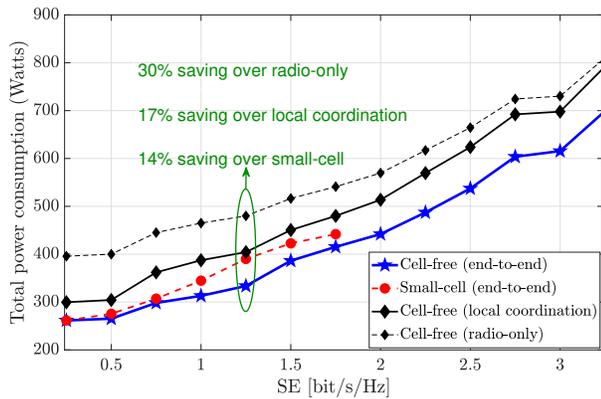| $N$ | 4 | $f_s, B$ | 30.72 MHz, 20 MHz | $N_{\text{DFT}}, N_{\text{used}}$ | 2048, 1200 |
|---|---|---|---|---|---|
| $T_s$ | 71.4 μs | $N_{\text{smooth}}, N_{\text{slot}}$ | 12, 16 | $\tau_c, \tau_p$ | 192, 8 |
| Size of coverage area | 1 km × 1 km | $P_{\text{RU},0}, \Delta^{\text{tr}}$ | 6.8N W, 4 | Pilot power, $p_{\max}$ | 100 mW, 1 W |
| $P_{\text{fixed}}, \sigma_{\text{cool}}$ | 120 W, 0.9 | $P_{\text{ONU}}, P_{\text{OLT}}$ | 7.7 W, 20 W | $P_{\text{GPP},0}^{\text{proc}}, P_{\text{RU},0}^{\text{proc}}$ | 20.8 W |
| $\Delta_{\text{GPP}}^{\text{proc}}, \Delta_{\text{RU}}^{\text{proc}}$ | 74 W | $C_{\text{GPP}}^{\max}, C_{\text{RU}}^{\max}$ | 180 GOPS | $R_{\max}, N_{\text{bits}}$ | 10 Gbps, 12 |



Fig. 3. The total power versus the SE requirement per UE for $L = 16$ and $K = 8$.
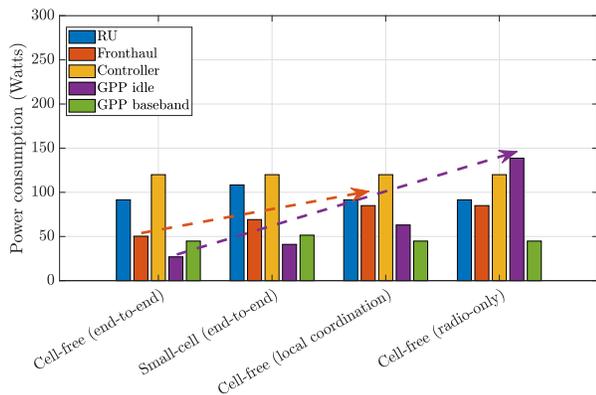


Fig. 4. Power consumption breakdown for the SE requirement of 1.25 bit/s/Hz, $L = 16$, and $K = 8$.

cessing power consumption when cell-free operation is used thanks to the less number of activated network components. When cell-free operation is used with local coordination, the frounthaul and GPP idle power consumption increases due to the fixedly assigned fronthaul resources and the partial intra-LC cloud resource sharing mechanism. The GPP idle power further increases when the radio-only scheme is considered due to the increased number of activated LCs and GPPs.

In Table II, we list the total power consumption of fully virtualized cell-free massive MIMO for different SE requirements from Fig. 3, which corresponds to the optimal solution found by solving the mixed-binary second-order cone programming problem in (24). We compare these optimal values with those obtained by solving the proposed approximate problem in (31) via the CCP approach outlined in Algorithm 1 and the

refinement method in Algorithm 2. As the table shows, there is a slight power consumption increase when solving the lower-complexity approximate problem. However, the resulting increase is at most 8%, which showcases the effectiveness of the proposed method.

In the remainder of the simulations, we consider a setup with $L = 36$ O-RUs. We do not specify the maximum number of GPPs and LCs, just assign the required number of them. First, we solve the proposed problem in (31) by the CCP approach outlined in Algorithm 1 and the refinement method in Algorithm 2 for different numbers of UEs, where the SE requirement for each UE is 2 bit/s/Hz. As seen in Fig. 5, for both FS-8 and FS-7.2, the end-to-end resource allocation results in smaller total power consumption thanks to the pooled cloud resources and sharing among the GPPs. Compared to radio-only orchestration, the fully virtualized end-to-end resource allocation provides 39% power saving for $K = 8$. On the other hand, when the intra-PHY split FS-7.2 is used, a processing unit is activated for each active O-RU leading to inefficient resource utilization with increased power. This reduces the cloud processing requirements under the assumption that all O-RUs are connected to the O-Cloud in the radio-only scheme. Moreover, each LC can now serve two times more O-RUs (6/3) than FS-8 allows, thanks to the corresponding significantly reduced fronthaul rate requirements. Hence, the power saving obtained by the fully virtualized orchestration over the radio-only one becomes smaller, i.e., 26%. This can also be observed by the power consumption breakdown illustrated in Fig. 6. When FS-8 is used, there is a substantial saving opportunity regarding the GPP idle power. On the other hand, this opportunity is less apparent when FS-7.2 is used since RU baseband processing power is the same for all three orchestration schemes.

In Figs. 7, 8, and Table III, we present the results of the joint sum SE maximization and total power consumption minimization problem in (37), which is solved by the proposed CCP algorithmic framework similar to Algorithm 1. To obtain the solution of the small-cell system, we select the O-RU that has the largest assigned power by the cell-free solution for each UE. Two values of the penalty parameter $\lambda$ for the sum SE in the multi-objective function in (37) are considered: i) $\lambda = 5$, which is called low-weight scenario and ii) $\lambda = 50$, which is called medium-weight scenario according to our observations and trials. As another benchmark, we consider only the sum SE maximization approach, where the power consumption is not included in the objective function. In Fig. 7, the cumulative distribution function (CDF) of the SE per UE of all the considered three methods for both cell-free and small-cell systems with FS-8. The SE values are very

TABLE II
COMPARISON OF THE TOTAL POWER CONSUMPTION (W) BETWEEN THE OPTIMAL SOLUTION AND THE CCP APPROACH.

| SE [bit/s/Hz] | 0.25 | 0.5 | 0.75 | 1 | 1.25 | 1.5 | 1.75 | 2 | 2.25 | 2.5 | 2.75 | 3 | 3.25 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Optimal solution | 262 | 265 | 298 | 313 | 334 | 386 | 415 | 442 | 487 | 537 | 604 | 615 | 699 |
| CCP solution | 266 | 282 | 304 | 333 | 361 | 400 | 432 | 473 | 517 | 571 | 646 | 666 | 756 |



Fig. 5. The total power versus the number of UEs for $L = 36$ and the SE requirement of 2 bit/s/Hz.



Fig. 6. Power consumption breakdown for $K = 8$, $L = 36$, and the SE requirement of 2 bit/s/Hz.

close to the ones in this figure when FS-7.2 is selected; thus, we skip those results. As shown in this figure and Table III, the main benefit of cell-free operation is the significantly increased SE for the most unfortunate UEs with the lowest SEs. The so-called 90%-likely SE, which can be provided to 90% of all the UEs, increases by more than two-fold when cell-free massive MIMO transmission is adopted. The sum SE also improves with the cell-free massive MIMO, but its impact is less significant considering the increased end-to-end power consumption. The increased power consumption is mainly compensated for by the higher SE guaranteed to all the UEs provided by cell-free massive MIMO. As the weight given to the sum SE maximization increases in the problem, this worst-case SE improves but with a cost of higher total power consumption.

In Fig. 8, we plot the total power consumption for the three resource allocation schemes and for different weights given to the sum SE and functional splits. It is worth noting that for a given weight, all the resource allocation schemes provide the same SE values. The difference lies in the power saving achieved by full virtualization and resource sharing in the cloud. As shown in the figure, the power saving is higher when the weight is low, and hence, the number of active O-RUs and the respective radio power are in a small value range. As observed before, the power saving is higher when FS-8 is utilized. As also demonstrated in Table III, the total power consumption is consistently larger with the FS-7.2.

## VIII. CONCLUSIONS

In this paper, we have modeled the end-to-end network power consumption of a cell-free massive MIMO system in the O-RAN architecture with fully centralized and intra-PHY functional splitting options. We have solved the two end-to-
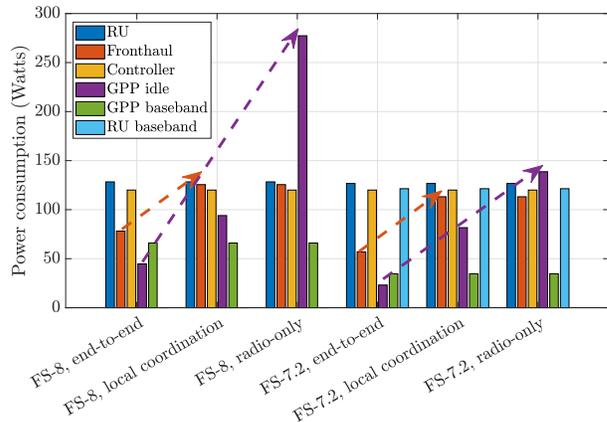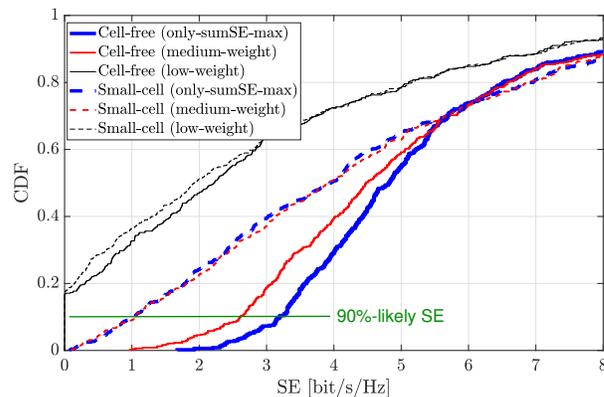


Fig. 7. The CDF of SE per UE for the joint sum SE maximization and power consumption minimization problem. For the sake of presentation, the CDF values are only shown for the SE range $[0 - 8]$ bit/s/Hz.

end resource allocation problems to find the power-efficient O-RU selection, O-RU/UE association, and the respective virtualized optical fronthaul and cloud resources. The first problem minimizes end-to-end consumption under the SE requests and network resource constraints. The second problem jointly maximizes sum SE and minimizes total power consumption to construct a balanced trade-off between SE performance and power cost.

The proposed fully virtualized end-to-end resource allocation achieves up to 39% and 19% power saving compared to the radio-only and local coordination-based orchestration, which only benefit from the O-RU turn-off mechanism and/or partial resource sharing among the computational tasks of the O-RUs that share the same fixed fronthaul wavelength. The key enabler of this power saving is the reduced fronthaul and GPP

TABLE III
SUM SE, 90%-LIKELY SE, AND TOTAL POWER CONSUMPTION FOR THE JOINT SUM SE MAXIMIZATION AND POWER CONSUMPTION MINIMIZATION PROBLEM.

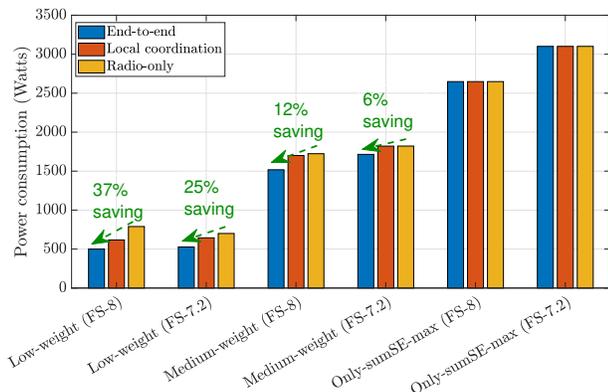| | Sum SE (bit/s/Hz) | | | 90%-likely SE (bit/s/Hz) | | |
|---|---|---|---|---|---|---|
| | Low-weight | Medium-weight | Only-sumSE-max | Low-weight | Medium-weight | Only-sumSE-max |
| Cell-free | 46.9 | 80.5 | 84.0 | 0 | 2.6 | 3.2 |
| Small-cell | 45.3 | 71.9 | 71.3 | 0 | 1.0 | 1.0 |
| | Total power consumption (Watts) | | | | | |
| | FS-8 | | | FS-7.2 | | |
| | Low-weight | Medium-weight | Only-sumSE-max | Low-weight | Medium-weight | Only-sumSE-max |
| Cell-free | 500 | 1518 | 2649 | 527 | 1715 | 3101 |
| Small-cell | 490 | 1062 | 1189 | 515 | 1222 | 1383 |



Fig. 8. The total power consumption for the joint sum SE maximization and power consumption minimization problem.

idle power consumption thanks to the less number of activated fronthaul wavelengths (LCs) and GPPs in the O-Cloud.

The considered cell-free system is advantageous over conventional small-cells in the same O-RAN architecture for both maximum rate and minimum power consumption. When the first optimization problem with the SE request constraints is considered, up to $14\%$ power saving is possible. On the other hand, for very small SE requests, the performance of cell-free and small-cell systems are the same, since the cell-free functionality is not activated. Cell-free massive MIMO increases the maximum provided rate by 1.7 with less energy per bit in comparison to small cells. When the required SE is high (more than 1.75 bit/s/Hz), the small-cell scenario may not be feasible, while by activating cell-free massive MIMO, we can obtain a feasible and power-efficient solution. When it comes to sum SE maximization, the main benefit of cell-free operation is the significantly increased 90%-likely SE (more than two-fold) over a small-cell system. This comes with a cost of increased power consumption, which is mainly compensated for by the higher guaranteed SE to all UEs rather than a relatively smaller sum SE improvement.

If the power consumption is prioritized over sum SE by setting the penalty parameter to a low value ($\lambda = 5$), we can save power up to 69% in comparison to the medium-weight setup ($\lambda = 50$), but the SE performance is inferior in terms of both sum SE and 90%-likely SE. Tuning the penalty parameter as in the medium-weight setup, we can save power up to %45 by only sacrificing the sum rate by %4.

## REFERENCES

[1] W. Jiang, B. Han, M. A. Habibi, and H. D. Schotten, "The road towards 6G: A comprehensive survey," *IEEE Open Journal of the Communications Society*, vol. 2, pp. 334–366, 2021.

[2] C. Han, Y. Wu, Z. Chen *et al.*, "Network 2030 a blueprint of technology applications and market drivers towards the year 2030 and beyond," 2018.

[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] G. Interdonato, E. Björnson, H. Quoc Ngo, P. Frenger, and E. G. Larsson, "Ubiquitous cell-free massive mimo communications," *EURASIP Journal on Wireless Communications and Networking*, vol. 2019, no. 1, pp. 1–13, 2019.

[5] Ö. T. Demir, E. Björnson, and L. Sanguinetti, "Foundations of user-centric cell-free massive MIMO," *Foundations and Trends® in Signal Processing*, vol. 14, no. 3-4, pp. 162–472, 2021.

[6] D. Wang, C. Zhang, Y. Du, J. Zhao, M. Jiang, and X. You, "Implementation of a cloud-based cell-free distributed massive MIMO system," *IEEE Communications Magazine*, vol. 58, no. 8, pp. 61–67, 2020.

[7] A. L. Imoize, H. I. Obakhena, F. I. Anyasi, and S. N. Sur, "A review of energy efficiency and power control schemes in ultra-dense cell-free massive MIMO systems for sustainable 6G wireless communication," *Sustainability*, vol. 14, no. 17, p. 11100, 2022.

[8] J. S. Vardakas, K. Ramantas, E. Vinogradov, M. A. Rahman, A. Girycki, S. Pollin, S. Pryor, P. Chanclou, and C. Verikoukis, "Machine learning-based cell-free support in the O-RAN architecture: An innovative converged optical-wireless solution toward 6G networks," *IEEE Wireless Communications*, vol. 29, no. 5, pp. 20–26, 2022.

[9] F. Malandrino, E. Chiaramello, M. Parazzini, and C. F. Chiasserini, "Performance and EMF exposure trade-offs in human-centric cell-free networks," in *2022 20th WiOpt Optimization in Mobile, Ad hoc, and Wireless Networks (WiOpt)*. IEEE, 2022, pp. 377–382.

[10] X. Wang, C. Cavdar, L. Wang, M. Tornatore, H. S. Chung, H. H. Lee, S. M. Park, and B. Mukherjee, "Virtualized cloud radio access network for 5G transport," *IEEE Communications Magazine*, vol. 55, no. 9, pp. 202–209, 2017.

[11] N. Alliance, "NGMN 5G P1 requirements & architecture work stream end-to-end architecture description of network slicing concept," 2016.

[12] X. Wang, C. Cavdar, L. Wang, M. Tornatore, Y. Zhao, H. S. Chung, H. H. Lee, S. Park, and B. Mukherjee, "Joint allocation of radio and optical resources in virtualized cloud RAN with CoMP," in *2016 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2016, pp. 1–6.

[13] X. Wang, S. Thota, M. Tornatore, H. S. Chung, H. H. Lee, S. Park, and B. Mukherjee, "Energy-efficient virtual base station formation in optical-access-enabled cloud-RAN," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 5, pp. 1130–1139, 2016.

[14] M. Masoudi, S. S. Lisi, and C. Cavdar, "Cost-effective migration toward virtualized C-RAN with scalable fronthaul design," *IEEE Systems Journal*, vol. 14, no. 4, pp. 5100–5110, 2020.

[15] A. Alabbasi, X. Wang, and C. Cavdar, "Optimal processing allocation to minimize energy and bandwidth consumption in hybrid CRAN," *IEEE Transactions on Green Communications and Networking*, vol. 2, no. 2, pp. 545–555, 2018.

[16] T. Sigwele, A. S. Alam, P. Pillai, and Y. F. Hu, "Energy-efficient cloud radio access networks by cloud based workload consolidation for 5G," *Journal of Network and Computer Applications*, vol. 78, pp. 1–8, 2017.

[17] M. Masoudi, M. G. Khafagy, A. Conte, A. El-Amine, B. Françoise, C. Nadjahi, F. E. Salem, W. Labidi, A. Süral, A. Gati *et al.*, "Green mobile networks for 5G and beyond," *IEEE Access*, vol. 7, pp. 107 270–107 299, 2019.

[18] L. Bonati, M. Polese, S. D'Oro, S. Basagni, and T. Melodia, "Open, programmable, and virtualized 5G networks: State-of-the-art and the road ahead," *Computer Networks*, vol. 182, p. 107516, 2020.

[19] A. Garcia-Saavedra and X. Costa-Perez, "O-RAN: Disrupting the virtualized RAN ecosystem," *IEEE Communications Standards Magazine*, 2021.

[20] O. Alliance, "O-ran whitepaper-building the next generation ran," *O-RAN Alliance, Tech. Rep., Oct*, 2019.

[21] M. Polese, L. Bonati, S. D'Oro, S. Basagni, and T. Melodia, "Understanding O-RAN: Architecture, interfaces, algorithms, security, and research challenges," *arXiv preprint arXiv:2202.01032*, 2022.

[22] T. Van Chien, E. Björnson, and E. G. Larsson, "Joint power allocation and load balancing optimization for energy-efficient cell-free massive MIMO networks," *IEEE Transactions on Wireless Communications*, vol. 19, no. 10, pp. 6798–6812, 2020.

[23] P. Agheli, M. J. Emadi, and H. Beyranvand, "Designing cost-and energy-efficient cell-free Massive MIMO network with Fiber and FSO fronthaul links," *CoRR*, vol. abs/2011.08511, 2021. [Online]. Available: http://arxiv.org/abs/2011.08511

[24] C. Pan, H. Zhu, N. J. Gomes, and J. Wang, "Joint precoding and RRH selection for user-centric green MIMO C-RAN," *IEEE Transactions on Wireless Communications*, vol. 16, no. 5, pp. 2891–2906, 2017.

[25] V. N. Ha and L. B. Le, "Computation capacity constrained joint transmission design for C-RANs," in *2016 IEEE Wireless Communications and Networking Conference*. IEEE, 2016, pp. 1–6.

[26] V. Ranjbar, A. Girycki, M. A. Rahman, S. Pollin, M. Moonen, and E. Vinogradov, "Cell-free mMIMO support in the O-RAN architecture: A PHY layer perspective for 5G and beyond networks," *IEEE Communications Standards Magazine*, vol. 6, no. 1, pp. 28–34, 2022.

[27] J. S. Vardakas, K. Ramantas, E. Datsika, M. Payaró, S. Pollin, E. Vinogradov, M. Varvarigos, P. Kokkinos, R. González-Sánchez, J. J. V. Olmos *et al.*, "Towards machine-learning-based 5G and beyond intelligent networks: The MARSAL project vision," in *2021 IEEE MeditCom Networking (MeditCom)*. IEEE, 2021, pp. 488–493.

[28] T. Murakami, N. Aihara, A. Ikami, Y. Tsukamoto, and H. Shinbo, "Analysis of CPU placement of cell-free massive MIMO for user-centric RAN," in *NOMS 2022-2022 IEEE/IFIP Network Operations and Management Symposium*. IEEE, 2022, pp. 1–7.

[29] Ö. T. Demir, M. Masoudi, E. Björnson, and C. Cavdar, "Cell-free massive MIMO in virtualized CRAN: How to minimize the total network power?" in *ICC 2022 - IEEE International Conference on Communications*, 2022, pp. 159–164.

[30] "O-RAN Use Cases and Deployment Scenarios Towards Open and Smart RAN," O-RAN Alliance, White Paper, February 2020.

[31] M. Dryjański, "O-RAN deployment scenarios," May 2023. [Online]. Available: https://rimedolabs.com/blog/o-ran-deployment-scenarios/

[32] 3GPP, 3rd Generation Partnership Project, "Study on new radio access technology: Radio access architecture and interfaces," 3GPP TR 38.801 V14.0, Tech. Rep., Mar. 2017.

[33] A. S. Abdalla, P. S. Upadhyaya, V. K. Shah, and V. Marojevic, "Toward next generation open radio access networks–what O-RAN can and cannot do!" *IEEE Network*, 2022.

[34] V. Q. Rodriguez, F. Guillemin, A. Ferrieux, and L. Thomas, "Cloud-RAN functional split for an efficient fronthaul network," in *2020 International Wireless Communications and Mobile Computing (IWCMC)*. IEEE, 2020, pp. 245–250.

[35] G. O. Pérez, J. A. Hernández, and D. Larrabeiti, "Fronthaul network modeling and dimensioning meeting ultra-low latency requirements for 5G," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 10, no. 6, pp. 573–581, 2018.

[36] O-RAN Fronthaul Working Group, "Control, user and synchronization plane specification," O-RAN, Specification, Tech. Rep., 2019.

[37] T. L. Marzetta, E. G. Larsson, H. Yang, and H. Q. Ngo, *Fundamentals of Massive MIMO*. Cambridge University Press, 2016.

[38] L. Sanguinetti, E. Björnson, and J. Hoydis, "Toward Massive MIMO 2.0: Understanding spatial correlation, interference suppression, and pilot contamination," *IEEE Transactions on Communications*, vol. 68, no. 1, pp. 232–257, 2020.

[39] G. Auer, V. Giannini, C. Desset, I. Godor, P. Skillermark, M. Olsson, M. A. Imran, D. Sabella, M. J. Gonzalez, O. Blume, and A. Fehske, "How much energy is needed to run a wireless network?" *IEEE Wireless Commun.*, vol. 18, no. 5, pp. 40–49, October 2011.

[40] S. Malkowsky, J. Vieira, L. Liu, P. Harris, K. Nieman, N. Kundargi, I. C. Wong, F. Tufvesson, V. Öwall, and O. Edfors, "The world's first real-time testbed for massive MIMO: Design, implementation, and validation," *IEEE Access*, vol. 5, pp. 9073–9088, 2017.

[41] C. Desset and B. Debaillie, "Massive MIMO for energy-efficient communications," in *2016 46th European Microwave Conference (EuMC)*. IEEE, 2016, pp. 138–141.

[42] B. Debaillie, C. Desset, and F. Louagie, "A flexible and future-proof power model for cellular base stations," in *VTC Spring*, 2015.

[43] Gurobi Optimization, LLC, "Gurobi Optimizer Reference Manual," 2021. [Online]. Available: https://www.gurobi.com

[44] F. Rinaldi, F. Schoen, and M. Sciandrone, "Concave programming for minimizing the zero-norm over polyhedral sets," *Computational Optimization and Applications*, vol. 46, no. 3, pp. 467–486, 2010.

[45] B. K. Sriperumbudur and G. R. Lanckriet, "On the convergence of the concave-convex procedure." in *Nips*, vol. 9. Citeseer, 2009, pp. 1759–1767.

[46] 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, 2011.

[47] K. Shen and W. Yu, "Fractional programming for communication systems—part II: Uplink scheduling via matching," *IEEE Transactions on Signal Processing*, vol. 66, no. 10, pp. 2631–2644, 2018.

[48] D. Simeonidou, "Dynamic softwarised RAN function placement in optical data centre networks," in *International IFIP Conference on Optical Network Design and Modeling*, vol. 11616. Springer Nature, 2020.