

# Reconfigurable Intelligent Surface for Green Edge Inference

Sheng Hua, *Graduate Student Member, IEEE*, Yong Zhou<sup>ID</sup>, *Member, IEEE*, Kai Yang<sup>ID</sup>, *Member, IEEE*, Yuanming Shi<sup>ID</sup>, *Senior Member, IEEE*, and Kunlun Wang<sup>ID</sup>, *Member, IEEE*

**Abstract**—Reconfigurable intelligent surface (RIS) as an emerging cost-effective technology can enhance the spectral- and energy-efficiency of wireless networks. In this article, we consider an RIS-aided green edge inference system, where the inference tasks generated from resource-constrained mobile devices (MDs) are uploaded to and cooperatively performed at multiple resource-enhanced base stations (BSs). Taking into account both the computation and uplink/downlink transmit power consumption, we formulate an overall network power consumption minimization problem, which calls for the joint design of the set of tasks performed by each BS, uplink/downlink beamforming vectors of BSs, transmit power of MDs, and uplink/downlink phase-shift matrices at the RIS. However, the resulting combinatorial optimization problem is nonconvex and highly intractable. We tackle the challenge of combinatorial variables by exploiting the group sparsity structure of the beamforming vectors. Moreover, a block-structured optimization with mixed  $\ell_{1,2}$ -norm and difference-of-convex-functions (DC) based three-stage framework is proposed to solve the problem, where the mixed  $\ell_{1,2}$ -norm and DC techniques are adopted to induce the group sparsity structure and handle the nonconvex rank-one constraint, respectively. Simulations demonstrate the supreme performance gain of deploying an RIS and confirm the effectiveness of the proposed algorithm over the baseline algorithms in reducing the overall network power consumption.

**Index Terms**—Reconfigurable intelligent surface, joint uplink and downlink, green edge inference, block-structured optimization, difference-of-convex programming.

Manuscript received October 13, 2020; revised January 16, 2021; accepted February 6, 2021. Date of publication February 11, 2021; date of current version May 20, 2021. The work of Yong Zhou was supported in part by the National Natural Science Foundation of China (NSFC) under Grant 62001294 and Grant 61971286. The work of Kunlun Wang was supported in part by the National Natural Science Foundation of China (NSFC) under Grant 61932014. This article was presented in part at the IEEE Global Communications Conference Workshops, Waikoloa, HI, USA, December 2019. The editor coordinating the review of this article was J. Xu. (*Corresponding author: Yong Zhou.*)

Sheng Hua is with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China, also with the Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences, Shanghai 200050, China, and also with the University of Chinese Academy of Sciences, Beijing 100049, China (e-mail: huasheng@shanghaitech.edu.cn).

Yong Zhou and Yuanming Shi are with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China (e-mail: zhouyong@shanghaitech.edu.cn; shiym@shanghaitech.edu.cn).

Kai Yang is with JD Technology Group, Beijing 100176, China (email: yangkai@shanghaitech.edu.cn).

Kunlun Wang is with the School of Communication and Electronic Engineering, East China Normal University, Shanghai 200241, China (e-mail: klwang@cee.ecnu.edu.cn).

Digital Object Identifier 10.1109/TGCN.2021.3058657

## I. INTRODUCTION

**B**ENEFITING from the availability of big data, recent years have witnessed the prosperity of deep neural network (DNN), which is a branch of artificial intelligence (AI) techniques and has demonstrated its superiority in a variety of intelligent applications (e.g., computer vision and natural language processing). Thanks to its strong representation ability, DNN has been adopted to promote the development of wireless communications and bring convenience to the system implementation in various aspects such as interference management [2]–[3]. On the other hand, wireless communications can promote AI as well, by analyzing the rich data generated by the ever-increasing number of mobile devices (MDs) and thus providing intelligent services for the MDs. With the MDs at the network edge generating 77 exabytes data per month by 2022 [4], the demand of performing inference tasks (e.g., object recognition and machine translation) is anticipated to be ubiquitous, especially in the next generation AI-powered wireless networks [5]. Driven by this trend, it is urgent to push traditional cloud-based DNN models to the network edge so as to unleash the potentials of edge data and in turn provide intelligent services [6]. One possible architecture to perform inference tasks is on-device inference, i.e., running DNN models directly on MDs. While the model compression [7], model selection [8], and hardware acceleration [9] are proposed as promising techniques to help devices run small-sized DNN models, deploying powerful models with millions of parameters is still challenging because of both the memory and battery limitations [10].

By leveraging edge computing [11] and deploying DNN models at the edge base stations (BSs) that have strong computational capacity and large storage resources, edge inference stands out as a promising paradigm to provide intelligent services for MDs [12]. To accomplish the inference tasks, the MDs upload the task-specific data to the BSs and subsequently the BSs deliver the inference results after finishing the inference process. Tailored for latency-critical applications, the authors in [13] and [14] respectively proposed device-edge and edge-cloud synergy frameworks to partition DNN model parameters based on network dynamics to minimize the execution latency. As energy efficiency is a key performance indicator for edge inference systems, the authors in [15] and [16] proposed energy-aware approaches to prune DNN models to minimize the computation power consumption (i.e., power required for the BSs to perform the inference tasks) while maintaining reasonable inference precision. However, the communication power consumption was not considered in [15] and [16]. The authors in [17]–[18]

proposed to minimize the sum of computation and downlink transmit power consumption (i.e., power required for the BSs to deliver inference results to the MDs), while the uplink transmit power (i.e., power required for the MDs to upload data to the BSs) was neglected. However, in edge inference systems, the traffic load in the uplink (e.g., raw images for an object recognition task) is usually comparable to that in the downlink (e.g., labeled images), resulting in high uplink transmit power consumption. Therefore, it is imperative to develop new techniques to reduce both the uplink and downlink transmit power consumption and in turn facilitate an energy-efficient design for edge inference systems.

Recently, a growing line of works focused on an emerging technology named reconfigurable intelligent surface (RIS) [19], which has the potential to significantly reduce the power consumption [1], [20] and improve the energy efficiency [21]. RIS is also envisioned as a key enabler to provide broadband connectivity for the future 6G systems [22]. In particular, an RIS is a low-cost planar array consisting of a large number of passive reflecting elements with reconfigurable phase shifts, each of which can be dynamically tuned via a software controller to reflect the incident signals [23]–[24]. These elements consume negligible energy due to their passive nature. By adaptively adjusting the phase shifts of reflecting elements, an RIS can combine the constructive signals and suppress the interference, thereby greatly enhancing the performance of wireless systems [25]–[26]. By jointly optimizing the beamforming vectors at the BS and the phase-shift matrix at the RIS, deploying an RIS has the potential to reduce the power consumption in various applications, e.g., downlink unicast [20] and broadcast [27] settings, non-orthogonal multiple access [28], and simultaneous wireless information and power transfer [29]. The authors in [30] investigated RIS-aided over-the-air computation to minimize signal distortion. To reduce the complexity of dynamic phase configuration for the RIS, the authors in [31]–[32] proposed deep learning methods to learn the mapping between the locations of transceivers and the optimal phase shifts. In terms of power consumption, all the aforementioned works only considered the downlink transmit power consumed by the BSs. However, in edge inference systems, the computation power consumption is an indispensable component and should be taken into account to accurately characterize the overall network power consumption. In addition, it is essential to optimize both the uplink and downlink phase-shift matrices of the RIS to assist both the uplink and downlink data transmissions. These two key issues make the approaches proposed in the existing works not applicable to RIS-aided edge inference systems.

To guarantee the quality of intelligent services provided for MDs, computation replication [33] allows each inference task to be performed by multiple BSs and creates multiple copies of the inference results at different BSs. These copies enable cooperative downlink transmission among the BSs on delivering the inference results. In terms of the power consumption, however, cooperative transmission and computation replication conflict with each other. Specifically, cooperative transmission reduces the downlink transmit power consumption by exploiting a higher beamforming gain, while computation replication rapidly

increases the computation power consumption by repeatedly running the same DNN model for multiple times. Therefore, it is necessary to strike a balance between the computation and communication power consumption via appropriately selecting inference tasks to be performed by each BS and in turn achieve green edge inference, which motivates this work.

In this article, we consider an RIS-aided green edge inference system with multiple BSs cooperatively performing inference tasks for multiple MDs, taking into account both the uplink and downlink transmit power consumption as well as the computation power consumption. Our objective is to minimize the overall network power consumption subject to prescribed quality-of-service (QoS) requirements, by jointly designing the task selection strategy, transmit/receive beamforming vectors of the BSs, the transmit power of the MDs, and the uplink/downlink phase-shift matrices at the RIS. However, the formulated problem is a combinatorial optimization problem with nonconvex constraints and is highly intractable. The main contributions of this article are summarized as follows

- We propose a joint design of the task selection strategy, transmit/receive beamforming vectors, transmit power, and uplink/downlink phase-shift matrices for an RIS-aided green edge inference system. To the best of our knowledge, this is the first attempt to unify beamforming vectors, transmit power, and phase shifts design in both the uplink and downlink transmission into a general edge inference framework.
- The combinatorial nature of the task selection strategy and the coupled optimization variables stand out as two major challenges. We address the challenge of the combinatorial variables by exploiting the group sparsity structure of the beamforming vectors, and tackle the challenge of the coupled variables by proposing a block-structured optimization (BSO) approach.
- With fixed phase shifts, we adopt the weighted mixed  $\ell_{1,2}$ -norm to induce the group sparsity of beamforming vectors. With fixed beamforming vectors and transmit power, the original problem is transformed to a homogeneous quadratically constrained quadratic programming (QCQP) with a nonconvex rank-one constraint. As the widely adopted semidefinite relaxation (SDR) technique incurs performance degradation when the number of reflecting elements is large, we propose a novel difference-of-convex-functions (DC) representation for this nonconvex constraint, followed by developing an effective DC algorithm. We then propose a BSO with mixed  $\ell_{1,2}$ -norm and DC based three-stage framework to solve the problem.
- Through extensive simulations, we show that the deployment of an RIS can significantly reduce the overall network power consumption of edge inference systems. Furthermore, the proposed BSO with mixed  $\ell_{1,2}$ -norm and DC algorithm achieves a significant performance improvement compared to the BSO with mixed  $\ell_{1,2}$ -norm and SDR algorithm, which demonstrates the effectiveness of DC in yielding the rank-one solutions.

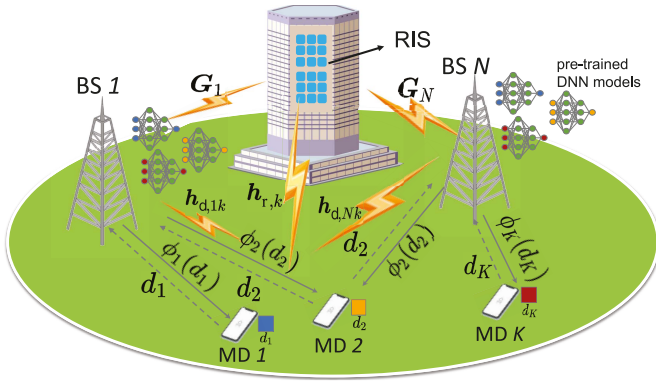


Fig. 1. RIS-aided edge inference system with  $N$  BSs collaboratively serving  $K$  MDs with the assistance of an RIS deployed on the facade of a building.

The remainder of this article is organized as follows. We present the system model and problem formulation in Section II. A BSO approach is developed in Section III to decouple optimization variables. We propose a three-stage framework in Section IV. Simulation results are illustrated in Section V. Finally, Section VI concludes this article.

**Notations:** We use boldface lower-case (e.g.,  $\mathbf{h}$ ) and upper-case letters (e.g.,  $\mathbf{G}$ ) to represent vectors and matrices, respectively. The transpose, conjugate transpose, trace operator and diagonal matrix are denoted as  $(\cdot)^T$ ,  $(\cdot)^H$ ,  $\text{Tr}(\cdot)$  and  $\text{diag}(\cdot)$ , respectively. The symbols  $|\cdot|$  and  $\Re(\cdot)$  denote the modulus and the real component of a complex number. The  $n \times n$  identity matrix is denoted as  $\mathbf{I}_n$ . The complex normal distribution is denoted as  $\mathcal{CN}$ . The inner product of two matrices  $\mathbf{X}$  and  $\mathbf{Y}$  is denoted as  $\langle \mathbf{X}, \mathbf{Y} \rangle$ , which is defined as  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{Tr}(\mathbf{X}^H \mathbf{Y})$ . The  $\ell_2$ -norm of a vector is denoted as  $\|\cdot\|_2$ . The spectral norm and Frobenius norm of a matrix are denoted as  $\|\cdot\|$  and  $\|\cdot\|_F$ , respectively. The  $i$ -th largest singular value of matrix  $\mathbf{X}$  is denoted as  $\sigma_i(\mathbf{X})$ . We use  $\mathbf{1}_{\{\cdot\}}$  to denote the indicator function which outputs 1 if the condition  $\cdot$  is satisfied, and outputs 0 otherwise. In the rest of this article, the superscripts UL and DL refer to uplink and downlink, respectively, and the letters d and r in the subscripts stand for the *direct link* and the *reflected link*, respectively.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

In this section, we describe the system model and the power consumption model for performing inference tasks at the network edge, followed by formulating an overall network power consumption minimization problem for green edge inference systems.

### A. System Model

We consider an RIS-aided edge inference system, where  $N$   $L_n$ -antenna BSs distributed in a small-cell network collaboratively serve  $K$  single-antenna MDs with the assistance of an  $M$ -element RIS deployed on the facade of a building, as shown in Fig. 1. Let  $\mathcal{N} = \{1, \dots, N\}$ ,  $\mathcal{K} = \{1, \dots, K\}$ , and  $\mathcal{M} = \{1, \dots, M\}$  denote the index sets of BSs, MDs, and reflecting elements, respectively. The BSs are resource-enhanced with strong computation and storage capabilities [34]. Each MD has

an inference task (e.g., image recognition) to be processed by a task-specific DNN model (e.g., AlexNet [35]). Specifically, the DNN model denoted as  $\phi_k$  takes MD  $k$ 's local data  $d_k$  (e.g., raw images) as input and generates the inference result  $\phi_k(d_k)$  (e.g., labeled images) as output. As it is impractical to run DNN models on resource-constrained MDs, we in this article propose to upload the inference tasks of the MDs to be performed at the BSs. We assume that all the BSs have downloaded the pre-trained DNN models from cloud servers in advance, therefore they can perform tasks for all the MDs [17].

The overall process of accomplishing the inference tasks in the edge inference system is composed of the following three phases.

- **Uplink Transmission:** The MDs upload the collected input data  $\{d_k, k \in \mathcal{K}\}$  to the BSs.
- **Inference Computation:** The BSs feed data (e.g.,  $d_k$ ) into a specific pre-trained DNN model (e.g.,  $\phi_k$ ) according to the task type and then obtain the inference results (e.g.,  $\phi_k(d_k)$ ).
- **Downlink Transmission:** The BSs deliver the inference results  $\{\phi_k(d_k), k \in \mathcal{K}\}$  to the corresponding MDs.

By exploiting the broadcast nature of wireless channels, each MD's data can be successfully received by multiple BSs in the uplink, enabling computation replication and creating multiple copies of the inference results at different BSs [33]. In the downlink, the BSs performing the same inference task cooperatively transmit the inference results to the corresponding MD. By exploiting the existing channel estimation approaches [36], we assume that the global channel state information (CSI) is available at the BSs as in [27]–[29]. Let  $\mathcal{A}_n \subseteq \mathcal{K}$  denote the set of MD indices whose inference tasks are selectively performed by BS  $n$ , and  $\mathcal{A} = (\mathcal{A}_1, \dots, \mathcal{A}_N)$  denote the task selection strategy. We adopt the time-division duplex (TDD) mode to separate the uplink and downlink transmissions.

1) **Uplink Transmission:** Let  $s_k^{\text{UL}} \in \mathbb{C}$  denote the representative information symbol of input data  $d_k$ , and  $p_k^{\text{UL}} \in \mathbb{R}$  denote the transmit power of MD  $k$ . Without loss of generality,  $\{s_k^{\text{UL}}, k \in \mathcal{K}\}$  are assumed to have zero mean and unit power. The signal received at BS  $n$  can be expressed as

$$\mathbf{y}_n^{\text{UL}} = \sum_{k \in \mathcal{K}} \mathbf{g}_{nk}^{\text{UL}} \sqrt{p_k^{\text{UL}}} s_k^{\text{UL}} + \mathbf{z}_n^{\text{UL}}, \quad \forall n \in \mathcal{N}, \quad (1)$$

where  $\mathbf{g}_{nk}^{\text{UL}} \in \mathbb{C}^{L_n \times 1}$  is the equivalent baseband channel response from MD  $k$  to BS  $n$  and  $\mathbf{z}_n^{\text{UL}} \sim \mathcal{CN}(\mathbf{0}, \sigma_n^2 \mathbf{I}_{L_n})$  is the additive white Gaussian noise (AWGN) at BS  $n$  with  $\sigma_n^2$  being the noise power. With the deployment of an RIS, the equivalent baseband channel from MD  $k$  to BS  $n$  consists of both the direct link and the reflected link, where the reflected link is a concatenation of the MD-RIS link, the phase shifts at the RIS, and the RIS-BS link [20], [27]–[29]. Therefore,  $\mathbf{g}_{nk}^{\text{UL}}$  can be modeled as

$$\mathbf{g}_{nk}^{\text{UL}} = \underbrace{\mathbf{h}_{d,nk}^{\text{UL}}}_{\text{direct link}} + \underbrace{\left(\mathbf{G}_n^{\text{UL}}\right)^H \left(\boldsymbol{\Theta}^{\text{UL}}\right)^H \mathbf{h}_{r,k}^{\text{UL}}}_{\text{reflected link}}, \quad (2)$$

where  $\mathbf{h}_{d,nk}^{\text{UL}} \in \mathbb{C}^{L_n \times 1}$ ,  $\mathbf{h}_{r,k}^{\text{UL}} \in \mathbb{C}^{M \times 1}$ , and  $\mathbf{G}_n^{\text{UL}} \in \mathbb{C}^{M \times L_n}$  denote the channel responses from MD  $k$  to BS  $n$ , from MD  $k$

to the RIS, and from the RIS to BS  $n$ , respectively. In addition,  $\Theta^{\text{UL}} = \beta \text{diag}(\theta_1^{\text{UL}}, \dots, \theta_M^{\text{UL}}) \in \mathbb{C}^{M \times M}$  denotes the diagonal phase-shift matrix for uplink transmission, where  $\beta \in [0, 1]$  is the amplitude reflection coefficient and  $\theta_m^{\text{UL}} = e^{j\varphi_m^{\text{UL}}}$  with  $\varphi_m^{\text{UL}} \in [0, 2\pi)$  being the uplink phase shift of the  $m$ -th reflecting element of the RIS. The reflected link only accounts for one-time reflection, because the power of signals reflected by two or more times is negligible due to the high path loss [20], [27]–[29].

We consider the linear beamforming strategy, and denote the receive beamforming vector of BS  $n$  to decode  $s_k^{\text{UL}}$  as  $\mathbf{v}_{nk}^{\text{UL}} \in \mathbb{C}^{L_n \times 1}$ . BS  $n$  only decodes MD  $k$ 's transmitted symbol  $s_k^{\text{UL}}$  if  $k \in \mathcal{A}_n$ . The estimated symbol at BS  $n$  for MD  $k \in \mathcal{A}_n$ , denoted by  $\hat{s}_{nk}^{\text{UL}} \in \mathbb{C}$ , is given by

$$\hat{s}_{nk}^{\text{UL}} = (\mathbf{v}_{nk}^{\text{UL}})^{\text{H}} \mathbf{y}_n^{\text{UL}} = (\mathbf{v}_{nk}^{\text{UL}})^{\text{H}} \mathbf{g}_{nk}^{\text{UL}} \sqrt{p_k^{\text{UL}}} s_k^{\text{UL}} + (\mathbf{v}_{nk}^{\text{UL}})^{\text{H}} \sum_{l \neq k} \mathbf{g}_{nl}^{\text{UL}} \sqrt{p_l^{\text{UL}}} s_l^{\text{UL}} + (\mathbf{v}_{nk}^{\text{UL}})^{\text{H}} \mathbf{z}_n^{\text{UL}}. \quad (3)$$

The uplink signal-to-interference-plus-noise ratio (SINR) observed at BS  $n$  for MD  $k$  is

$$\text{SINR}_{nk}^{\text{UL}} = \frac{p_k^{\text{UL}} |(\mathbf{v}_{nk}^{\text{UL}})^{\text{H}} \mathbf{g}_{nk}^{\text{UL}}|^2}{\sum_{l \neq k} p_l^{\text{UL}} |(\mathbf{v}_{nk}^{\text{UL}})^{\text{H}} \mathbf{g}_{nl}^{\text{UL}}|^2 + \sigma_k^2 \|\mathbf{v}_{nk}^{\text{UL}}\|_2^2}, \quad \forall k \in \mathcal{A}_n, n \in \mathcal{N}. \quad (4)$$

2) *Downlink Transmission*: After performing the inference tasks, the BSs cooperatively transmit the inference results  $\{\phi_k(d_k), k \in \mathcal{K}\}$  to the corresponding MDs through downlink wireless channels. Let  $s_k^{\text{DL}} \in \mathbb{C}$  denote the representative symbol of  $\phi_k(d_k)$  intended for MD  $k$  and  $\mathbf{v}_{nk}^{\text{DL}}$  denote the downlink beamforming vector from BS  $n$  to MD  $k$ . Without loss of generality,  $\{s_k^{\text{DL}}, k \in \mathcal{K}\}$  are assumed to have zero mean and unit power. The signal transmitted by BS  $n$ , denoted as  $\mathbf{x}_n^{\text{DL}} \in \mathbb{C}^{L_n \times 1}$ , is a summation of beamformed symbols for MD  $k \in \mathcal{A}_n$ , i.e.,  $\mathbf{x}_n^{\text{DL}} = \sum_{k \in \mathcal{A}_n} \mathbf{v}_{nk}^{\text{DL}} s_k^{\text{DL}}$ ,  $\forall n \in \mathcal{N}$ . The signal received by MD  $k$  can be expressed as

$$\begin{aligned} y_k^{\text{DL}} &= \sum_{n \in \mathcal{N}} (\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \mathbf{x}_n^{\text{DL}} + z_k^{\text{DL}} \\ &= \sum_{n \in \mathcal{N}} (\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \\ &\quad \times \left( \mathbf{1}_{\{k \in \mathcal{A}_n\}} \mathbf{v}_{nk}^{\text{DL}} s_k^{\text{DL}} + \sum_{l \in \mathcal{A}_n, l \neq k} \mathbf{v}_{nl}^{\text{DL}} s_l^{\text{DL}} \right) + z_k^{\text{DL}} \\ &= \sum_{n \in \mathcal{N}} \mathbf{1}_{\{k \in \mathcal{A}_n\}} (\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \mathbf{v}_{nk}^{\text{DL}} s_k^{\text{DL}} \\ &\quad + \sum_{l \neq k} \sum_{n \in \mathcal{N}} \mathbf{1}_{\{l \in \mathcal{A}_n\}} (\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \mathbf{v}_{nl}^{\text{DL}} s_l^{\text{DL}} + z_k^{\text{DL}}, \quad \forall k \in \mathcal{K}, \end{aligned} \quad (5)$$

where  $z_k^{\text{DL}} \in \mathbb{C}$  is the AWGN at MD  $k$  with zero mean and power  $\sigma_k^2$ , and  $(\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \in \mathbb{C}^{1 \times L_n}$  is the equivalent downlink channel response from BS  $n$  to MD  $k$ . Similar to the uplink

counterpart,  $(\mathbf{g}_{nk}^{\text{DL}})^{\text{H}}$  can be modeled as

$$(\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} = \underbrace{(\mathbf{h}_{d,nk}^{\text{DL}})^{\text{H}}}_{\text{direct link}} + \underbrace{(\mathbf{h}_{r,k}^{\text{DL}})^{\text{H}} \Theta^{\text{DL}} \mathbf{G}_n^{\text{DL}}}_{\text{reflected link}}, \quad (6)$$

where  $(\mathbf{h}_{d,nk}^{\text{DL}})^{\text{H}} \in \mathbb{C}^{1 \times L_n}$ ,  $(\mathbf{h}_{r,k}^{\text{DL}})^{\text{H}} \in \mathbb{C}^{1 \times M}$ , and  $\mathbf{G}_n^{\text{DL}} \in \mathbb{C}^{M \times L_n}$  denote the channel responses from BS  $n$  to MD  $k$ , from the RIS to MD  $k$ , and from BS  $n$  to the RIS, respectively, and  $\Theta^{\text{DL}} = \beta \text{diag}(\theta_1^{\text{DL}}, \dots, \theta_M^{\text{DL}}) \in \mathbb{C}^{M \times M}$  is the downlink phase-shift matrix with diagonal entries  $\theta_m^{\text{DL}} = e^{j\varphi_m^{\text{DL}}}$  and  $\varphi_m^{\text{DL}} \in [0, 2\pi)$ . Although channel reciprocity is often assumed to hold in TDD systems, we consider a general case, where the channel responses in the uplink and downlink can be different. Based on (5), the SINR observed by MD  $k \in \mathcal{K}$  in the downlink is given by

$$\begin{aligned} \text{SINR}_k^{\text{DL}} &= \frac{\left| \sum_{n \in \mathcal{N}} \mathbf{1}_{\{k \in \mathcal{A}_n\}} (\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \mathbf{v}_{nk}^{\text{DL}} \right|^2}{\sum_{l \neq k} \left| \sum_{n \in \mathcal{N}} \mathbf{1}_{\{l \in \mathcal{A}_n\}} (\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \mathbf{v}_{nl}^{\text{DL}} \right|^2 + \sigma_k^2} \\ &= \frac{\left| \sum_{n \in \mathcal{N}} (\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \mathbf{v}_{nk}^{\text{DL}} \right|^2}{\sum_{l \neq k} \left| \sum_{n \in \mathcal{N}} (\mathbf{g}_{nk}^{\text{DL}})^{\text{H}} \mathbf{v}_{nl}^{\text{DL}} \right|^2 + \sigma_k^2}, \end{aligned} \quad (7)$$

where the second equality holds because BS  $n$  does not transmit data to MD  $k$  by setting  $\mathbf{v}_{nk}^{\text{DL}} = \mathbf{0}$  if  $k \notin \mathcal{A}_n$ .

## B. Power Consumption Model

As running DNN models often incurs high energy consumption due to their high computational complexity [37]–[38] and energy-efficiency is one of the key performance indicators for green communications [39], we in this subsection present the power consumption model of the proposed edge inference system, taking into consideration both the computation power for inference and the communication power for uplink and downlink transmissions.

1) *Computation Power Consumption*: We denote the power consumption of performing MD  $k$ 's inference task at BS  $n$  as  $P_{nk}^{\text{C}}$ . Therefore, the total computation power consumption at all BSs is given by  $P_{\text{comp}}(\mathcal{A}) = \sum_{n \in \mathcal{N}} \sum_{k \in \mathcal{A}_n} P_{nk}^{\text{C}}$ . It is worth noting that the majority of the computation power is consumed for running DNN models, which can be estimated by using the energy estimation methodology proposed in [37]. This methodology provides a layer-wise energy breakdown for arbitrary neural networks. In particular, the DNN configurations (e.g., number of filters, number of input feature maps) are taken as inputs and the normalized layer-wise consumptions energy of the neural network (i.e., normalized by the energy consumption per multiply-and-accumulation (MAC) operation) are generated as outputs [40]. The computation time can be calculated via dividing the number of MAC operations by the average throughput of a CPU chip. Therefore, the power consumption for performing an inference task equals to the total energy consumption divided by the corresponding computation time.

For example, the energy consumption of running AlexNet to process one image on a well-designed energy-efficient Eyeriss chip can be approximated by that of performing  $4 \times 10^9$  MAC

operations, or equivalently 0.45 W when the chip is at core supply voltage 1.2 V [41]. As the typical value of computation power  $P_{nk}^c$  (e.g., 0.45 W) is comparable to the BSs' transmit power (e.g., 1 W [42]), it is necessary to take into account both the computation and transmit power to facilitate the energy-efficient design.

2) *Communication Power Consumption*: The communication power consumption consists of the power consumed by the MDs in the uplink transmission and by the BSs in the downlink transmission. According to (2)-(7), the total uplink transmit power consumption is  $\sum_{k \in \mathcal{K}} p_k^{\text{UL}}$ , while the downlink transmit power consumption of BS  $n$  is given by  $\mathbf{E}[\sum_{k \in \mathcal{A}_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2] = \sum_{k \in \mathcal{A}_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2$ , where  $\mathbf{E}[\cdot]$  denotes the expectation. Therefore, the total communication power consumption for both uplink and downlink transmissions is given by  $P_{\text{comm}}(\mathcal{A}, \{p_k^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\}) = \sum_{k \in \mathcal{K}} p_k^{\text{UL}} + \sum_{n \in \mathcal{N}} \sum_{k \in \mathcal{A}_n} \frac{1}{\eta_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2$ , where  $\eta_n$  is the drain efficiency coefficient of the radio frequency power amplifier of BS  $n$ .

In summary, the overall network power consumption, consisting of both the computation and communication power consumption, can be expressed as

$$\begin{aligned} P_{\text{total}}(\mathcal{A}, \{p_k^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\}) \\ &= P_{\text{comm}}(\mathcal{A}, \{p_k^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\}) + P_{\text{comp}}(\mathcal{A}) \\ &= \sum_{k \in \mathcal{K}} p_k^{\text{UL}} + \sum_{n \in \mathcal{N}} \sum_{k \in \mathcal{A}_n} \frac{1}{\eta_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2 + \sum_{n \in \mathcal{N}} \sum_{k \in \mathcal{A}_n} P_{nk}^c. \end{aligned} \quad (8)$$

It is worth noting that the static power consumption of the BSs (proportional to  $N$ ) and that of the RIS (proportional to  $M$ ) are not included in  $P_{\text{total}}$  and regarded as constants, because they do not vary with the optimization variables and are neglectable compared to computation and communication power consumption.

### C. Problem Formulation and Analysis

In the proposed edge inference system, there exists a fundamental tradeoff between the communication and computation power consumption. Specifically, with computation replication, more BSs performing the same task reduces the downlink transmit power consumption by exploiting a higher cooperative beamforming gain, at the cost of increasing the computation power consumption. Therefore, we propose to achieve green edge inference by minimizing the overall network power consumption via striking a good balance between the communication and computation power consumption.

Let  $\{\gamma_k^{\text{UL}}, k \in \mathcal{K}\}$  and  $\{\gamma_k^{\text{DL}}, k \in \mathcal{K}\}$  denote the SINR thresholds required to successfully receive the input data and inference results in the uplink and downlink, respectively. Given a task selection strategy  $\mathcal{A}$ , the network power consumption minimization problem is formulated as

$$\begin{aligned} \mathcal{P}(\mathcal{A}) : \quad & \underset{\{\mathbf{v}_{nk}^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\}, \{p_k^{\text{UL}}\}}{\text{minimize}} \quad P_{\text{total}}(\mathcal{A}, \{p_k^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\}) \\ & \text{subject to} \quad \text{SINR}_k^{\text{DL}} \geq \gamma_k^{\text{DL}}, \forall k \in \mathcal{K}, \end{aligned} \quad (9a)$$

$$\begin{aligned} \text{SINR}_{nk}^{\text{UL}} &\geq \gamma_k^{\text{UL}}, \\ \forall k \in \mathcal{A}_n, n \in \mathcal{N}, \end{aligned} \quad (9b)$$

$$\begin{aligned} \sum_{k \in \mathcal{A}_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2 &\leq P_{n, \text{max}}^{\text{DL}}, \\ \forall n \in \mathcal{N}, \end{aligned} \quad (9c)$$

$$\begin{aligned} p_k^{\text{UL}} &\leq P_{k, \text{max}}^{\text{UL}}, \forall k \in \mathcal{K}, \end{aligned} \quad (9d)$$

$$\begin{aligned} |\theta_m^{\text{UL}}| &= |\theta_m^{\text{DL}}| = 1, \\ \forall m \in \mathcal{M}, \end{aligned} \quad (9e)$$

$$\begin{aligned} \mathbf{v}_{nk}^{\text{UL}} &= \mathbf{v}_{nk}^{\text{DL}} = \mathbf{0}, \\ \forall k \notin \mathcal{A}_n, n \in \mathcal{N}, \end{aligned} \quad (9f)$$

where  $P_{k, \text{max}}^{\text{UL}}$  and  $P_{n, \text{max}}^{\text{DL}}$  denote the maximum transmit power of MD  $k$  and BS  $n$  in the uplink and downlink, respectively, and (9f) are the group sparsity constraints of beamforming vectors. Specifically, if  $k \notin \mathcal{A}_n$ , BS  $n$  does not decode MD  $k$ 's data in the uplink (i.e.,  $\mathbf{v}_{nk}^{\text{UL}} = \mathbf{0}$ ) and subsequently cannot transmit inference results to MD  $k$  in the downlink (i.e.,  $\mathbf{v}_{nk}^{\text{DL}} = \mathbf{0}$ ).

As  $\mathcal{A}$  is a variable to be designed, we need to search over all possibilities of  $\mathcal{A}$  to obtain the optimal task selection strategy  $\mathcal{A}^*$ . Therefore the overall optimization problem is given by

$$\mathcal{P}_{\text{original}}: \underset{\mathcal{A}_1 \subseteq \mathcal{K}, \dots, \mathcal{A}_N \subseteq \mathcal{K}}{\text{minimize}} \quad p^*(\mathcal{A}) \quad (10)$$

where  $p^*(\mathcal{A})$  is the objective value of problem  $\mathcal{P}(\mathcal{A})$ .

In order to solve problem  $\mathcal{P}_{\text{original}}$ , we are confronted with several main challenges. As set  $\mathcal{K}$  has  $2^K$  different subsets,  $\mathcal{A}$  has a total of  $2^{KN}$  different possibilities, making it apparently impractical to search over all the possibilities. Despite of the troublesome variable  $\mathcal{A}$ , the coupled continuous variables phases shifts and beamforming vectors in constraints (9a)-(9b) pose a unique challenge. Moreover, the unit-modulus constraint (9e) imposed by the phase-shift of each RIS element is nonconvex. In the following, we shall exploit the group sparsity structure of beamforming vectors to get rid of the combinatorial variable  $\mathcal{A}$  in problem  $\mathcal{P}_{\text{original}}$ , thereby facilitating efficient algorithm design.

### III. BLOCK-STRUCTURED OPTIMIZATION APPROACH

In general, the combinatorial optimization problem  $\mathcal{P}_{\text{original}}$  is hard to tackle. Fortunately, based on the key observation that the task selection strategy  $\mathcal{A}$  has an intrinsic connection with the group sparsity structure of beamforming vector  $\mathbf{v}_{nk} = [(\mathbf{v}_{nk}^{\text{UL}})^T, (\mathbf{v}_{nk}^{\text{DL}})^T]^T$ , we can eliminate the troublesome variable  $\mathcal{A}$ . Specifically, all coefficients in the beamforming group  $\mathbf{v}_{nk}$  are zero simultaneously if  $k \notin \mathcal{A}_n$ . In other words, we have  $k \notin \mathcal{A}_n \Leftrightarrow \mathbf{v}_{nk} = \mathbf{0}$  and  $k \in \mathcal{A}_n \Leftrightarrow \mathbf{v}_{nk} \neq \mathbf{0}$ . Therefore, the overall network power consumption given in (8) can be equivalently rewritten as

$$\begin{aligned} P_{\text{total}}(\{p_k^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\}) &= \sum_{k=1}^K p_k^{\text{UL}} + \sum_{n=1}^N \sum_{k=1}^K \frac{1}{\eta_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2 \\ &+ \sum_{n=1}^N \sum_{k=1}^K \mathbf{1}_{\{\mathbf{v}_{nk} \neq \mathbf{0}\}} P_{nk}^c. \end{aligned} \quad (11)$$



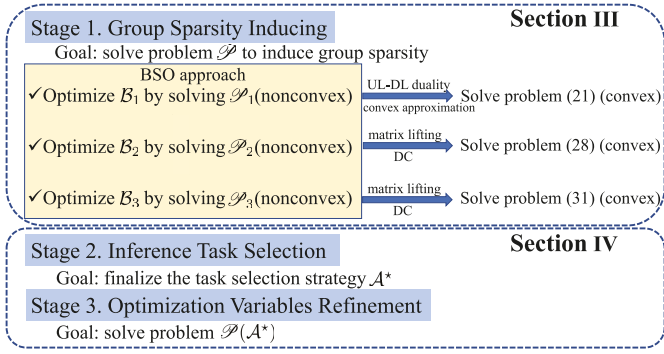


Fig. 2. Overview of the proposed three-stage framework for problem  $\mathcal{P}_{\text{original}}$ .

As multiple tasks may not be performed by a specific BS, the aggregated beamforming vector  $\mathbf{v} \in \mathbb{C}^{K \sum_{n=1}^N L_n}$  defined as  $\mathbf{v} = [\mathbf{v}_{11}^T, \dots, \mathbf{v}_{1K}^T, \dots, \mathbf{v}_{N1}^T, \dots, \mathbf{v}_{NK}^T]^T$  is expected to have the group sparsity structure with only a few non-zero blocks.

The above discussions indicate that we do not have to explicitly optimize the task selection strategy  $\mathcal{A}$ . Instead,  $\mathcal{A}$  can be determined by the group sparsity pattern of the beamforming vectors, i.e.,  $\mathcal{A}_n = \{k | \mathbf{v}_{nk} \neq \mathbf{0}, k \in \mathcal{K}\}, \forall n \in \mathcal{N}$ . Therefore, we propose to solve problem  $\mathcal{P}_{\text{original}}$  based on the following three stages. In *Stage 1*, we presume that each BS performs inference tasks for all the MDs (i.e.,  $\mathcal{A}_n = \mathcal{K}, \forall n \in \mathcal{N}$ ), and solve the following group sparse beamforming problem to obtain the initial group sparse beamforming vectors, which indicate the group sparsity structure of the optimal beamforming vectors

$$\mathcal{P} : \begin{aligned} & \underset{\substack{\{\mathbf{v}_{nk}^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\}, \{p_k^{\text{UL}}\}}} \\ & \underset{\substack{\{\mathbf{v}_{nk}^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\}, \{p_k^{\text{UL}}\}}} \end{aligned} \quad \text{minimize} \quad P_{\text{total}}(\{\mathbf{v}_{nk}^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\})$$

$$\text{subject to } \text{SINR}_k^{\text{DL}} \geq \gamma_k^{\text{DL}}, \forall k \in \mathcal{K}, \quad (12a)$$

$$\text{SINR}_{nk}^{\text{UL}} \geq \gamma_k^{\text{UL}}, \forall n \in \mathcal{N}, k \in \mathcal{K}, \quad (12b)$$

$$\sum_{k=1}^K \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2 \leq P_{n,\max}^{\text{DL}}, \quad \forall n \in \mathcal{N}, \quad (12c)$$

$$p_k^{\text{UL}} \leq P_{k,\max}^{\text{UL}}, \forall k \in \mathcal{K}, \quad (12d)$$

$$|\theta_m^{\text{DL}}| = 1, \forall m \in \mathcal{M}, \quad (12e)$$

$$|\theta_m^{\text{UL}}| = 1, \forall m \in \mathcal{M}. \quad (12f)$$

Based on the initial group sparse beamforming vectors, we in *Stage 2* determine the task selection strategy  $\mathcal{A}^*$ . Finally, we solve problem  $\mathcal{P}(\mathcal{A}^*)$  to refine the optimization variables in *Stage 3*. In the remaining of this section, we propose a BSO approach to solve problem  $\mathcal{P}$ , while the details of *Stage 2* and *Stage 3* are presented in Section IV. We provide a graphic illustration in Fig. 2 to provide a better overview of the proposed three-stage framework.

The main idea of the BSO approach is to first partition the variables into several blocks, and then alternately

optimize one of the blocks in each iteration while keeping the others fixed [43]. Specifically, we partition the five optimization variables into three blocks, denoted as  $\mathcal{B}_1 = (\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_k^{\text{UL}}\})$ ,  $\mathcal{B}_2 = \Theta^{\text{DL}}$ , and  $\mathcal{B}_3 = \Theta^{\text{UL}}$ .

#### A. Optimizing Variables $\{\mathbf{v}_{nk}^{\text{UL}}\}$ , $\{p_k^{\text{UL}}\}$ , and $\{\mathbf{v}_{nk}^{\text{DL}}\}$

When  $\mathcal{B}_2$  and  $\mathcal{B}_3$  are fixed, problem  $\mathcal{P}$  is reduced to the following problem

$$\mathcal{P}_1 : \begin{aligned} & \underset{\substack{\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_k^{\text{UL}}\}}} \end{aligned} \quad \text{minimize} \quad P_{\text{total}}(\{\mathbf{v}_{nk}^{\text{UL}}\}, \{\mathbf{v}_{nk}^{\text{DL}}\})$$

$$\text{subject to } (12a)-(12d).$$

It is observed that optimization variables  $\mathbf{v}_{nk}^{\text{UL}}$  and  $\mathbf{v}_{nk}^{\text{DL}}$  are coupled only in the objective function, but not in the constraints. For the sake of analysis convenience, we temporarily dismiss the indicator term in the objective function and split problem  $\mathcal{P}_1$  into two parts, i.e., the downlink part  $\mathcal{P}_{1-1}$  and the uplink part  $\mathcal{P}_{1-2}$ . Afterwards we will combine the uplink and downlink parts to derive an effective solution to problem  $\mathcal{P}_1$ . The power minimization problem in the downlink part is

$$\mathcal{P}_{1-1} : \begin{aligned} & \underset{\{\mathbf{v}_{nk}^{\text{DL}}\}} \end{aligned} \quad \text{minimize} \quad \sum_{n=1}^N \sum_{k=1}^K \frac{1}{\eta_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2$$

$$\text{subject to } (12a), (12c),$$

which is a celebrated problem formulation in unicast multiple-input single-output systems. Due to the fact that an arbitrary phase rotation of vector  $\mathbf{v}_{nk}^{\text{DL}}$  does not affect the SINR constraints (12a) [44], we can replace (12a) with the following second-order cone (SOC) constraint

$$\sqrt{\sum_{l \neq k} |(\mathbf{g}_k^{\text{DL}})^H \mathbf{v}_l^{\text{DL}}|^2} + \sigma_k^2 \leq \frac{1}{\sqrt{\gamma_k^{\text{DL}}}} \Re((\mathbf{g}_k^{\text{DL}})^H \mathbf{v}_k^{\text{DL}}), \quad (13)$$

where  $\mathbf{g}_k^{\text{DL}} = [(\mathbf{g}_{1k}^{\text{DL}})^T, \dots, (\mathbf{g}_{Nk}^{\text{DL}})^T]^T$  and  $\mathbf{v}_k^{\text{DL}} = [(\mathbf{v}_{1k}^{\text{DL}})^T, \dots, (\mathbf{v}_{Nk}^{\text{DL}})^T]^T$  denote the aggregated channel response vector and transmit beamforming vector with respect to MD  $k$ , respectively. Therefore, problem  $\mathcal{P}_{1-1}$  is recast as a convex second-order cone programming (SOCP), which can be effectively solved by interior-point methods using modern software like CVX [45].

On the other hand, the power minimization problem in the uplink is given by

$$\mathcal{P}_{1-2} : \begin{aligned} & \underset{\substack{\{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_k^{\text{UL}}\}}} \end{aligned} \quad \text{minimize} \quad \sum_{k=1}^K p_k^{\text{UL}} \quad \text{subject to } (12b), (12d).$$

Although the SINR constraints (12b) are similar to those in the downlink counterpart, we cannot convexify them in a similar way because the numerator involves both optimization variables. Moreover, another issue in  $\mathcal{P}_{1-2}$  is that directly optimizing this problem makes  $\mathbf{v}_{nk}^{\text{UL}}$  to be nearly zero, because an arbitrary scaling of  $\mathbf{v}_{nk}^{\text{UL}}$  does not affect the uplink SINR constraints (12b) [46]. Specifically, if  $(\tilde{\mathbf{v}}_{nk}^{\text{UL}}, \tilde{p}_k^{\text{UL}})$  denotes the optimal solution to problem  $\mathcal{P}_{1-2}$ , then we have  $\tilde{\mathbf{v}}_{nk}^{\text{UL}} \approx \mathbf{0}, \forall n \in \mathcal{N}, k \in \mathcal{K}$ . Although the scaling issue does not

violate the group sparsity structure of  $\mathbf{v}_{nk}$ , it indicates that the receive beamforming vectors do not contribute to the task selection. Based on the uplink-downlink duality, in the following, we shall propose a virtual downlink formulation to overcome the scaling issue.

To facilitate an effective algorithm design, we relax problem  $\mathcal{P}_{1-2}$  to the following problem

$$\begin{aligned} & \underset{\{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_{nk}\}}{\text{minimize}} \quad \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K p_{nk} \\ & \text{subject to} \quad p_{nk} \leq P_{k,\max}^{\text{UL}}, \quad \forall n \in \mathcal{N}, k \in \mathcal{K}, \\ & \quad \frac{p_{nk} |(\mathbf{v}_{nk}^{\text{UL}})^H \mathbf{g}_{nk}^{\text{UL}}|^2}{\sum_{l \neq k} p_{nl} |(\mathbf{v}_{nl}^{\text{UL}})^H \mathbf{g}_{nl}^{\text{UL}}|^2 + \sigma_n^2 \|\mathbf{v}_{nk}^{\text{UL}}\|_2^2} \geq \gamma_k^{\text{UL}}, \\ & \quad \forall n \in \mathcal{N}, k \in \mathcal{K}. \end{aligned} \quad (14a)$$

We can easily verify that problem (14) is indeed a relaxation to problem  $\mathcal{P}_{1-2}$ , because given any solution  $(\{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_k^{\text{UL}}\})$  feasible to problem  $\mathcal{P}_{1-2}$ ,  $(\{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_{nk}\})$  is also feasible to problem (14) when  $p_{nk} = p_k^{\text{UL}}, \forall n \in \mathcal{N}$ , and the objective value of problem (14) is no greater than that of  $\mathcal{P}_{1-2}$ . The motivation for this relaxation is that problem (14) can be solved in the virtual downlink formulation so as to overcome the scaling issue.

We first consider an ideal scenario that the MDs have unlimited power budgets (i.e.,  $P_{k,\max}^{\text{UL}} = +\infty, \forall k \in \mathcal{K}$ ). Problem (14) is then equivalent to the following virtual downlink power minimization problem

$$\underset{\{\mathbf{v}_{nk}^{\text{VDL}}\}}{\text{minimize}} \quad \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2 \quad (15a)$$

$$\text{subject to} \quad \text{SINR}_{nk}^{\text{VDL}} \geq \gamma_k^{\text{UL}}, \quad \forall n \in \mathcal{N}, k \in \mathcal{K}, \quad (15b)$$

where  $\mathbf{v}_{nk}^{\text{VDL}} \in \mathbb{C}^{L_n \times 1}$  denotes the virtual downlink transmit beamforming vector from BS  $n$  to MD  $k$ , and  $\text{SINR}_{nk}^{\text{VDL}}$  is the virtual downlink SINR observed by MD  $k$  defined as  $\text{SINR}_{nk}^{\text{VDL}} = \frac{|(\mathbf{g}_{nk}^{\text{UL}})^H \mathbf{v}_{nk}^{\text{VDL}}|^2}{\sum_{l \neq k} |(\mathbf{g}_{nl}^{\text{UL}})^H \mathbf{v}_{nl}^{\text{VDL}}|^2 + \sigma_n^2}$ . Note that in  $\text{SINR}_{nk}^{\text{VDL}}$ , the scaling issue does not exist for  $\{\mathbf{v}_{nk}^{\text{VDL}}, n \in \mathcal{N}, k \in \mathcal{K}\}$ . The rigorous proof of the equivalence of the uplink power minimization problem (14) and the virtual downlink problem (15) can be derived by Lagrangian duality, as shown in the Appendix. The optimal solutions obtained by solving problem (15) have close connections to solutions to problem (14), i.e.,  $\mathbf{v}_{nk}^{\text{VDL}} = \mathbf{v}_{nk}^{\text{UL}}$  and  $\sum_{n=1}^N \sum_{k=1}^K \|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2 = \sum_{n=1}^N \sum_{k=1}^K p_{nk}$ . However, it is worth mentioning that the equivalence between the virtual downlink beamforming power and the uplink transmit power does not necessarily hold, i.e.,  $\|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2 \neq p_{nk}$ . Therefore if  $P_{k,\max}^{\text{UL}} < +\infty$ , we cannot directly rewrite the transmit power constraints (14a) as  $\|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2 \leq P_{k,\max}^{\text{UL}}, \forall n, \forall k$  out of intuition and add them to problem (15). Instead, we consider a sum-power constraint to relax the uplink transmit power constraints (14a), which is given by

$$\sum_{n=1}^N \sum_{k=1}^K \|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2 = \sum_{n=1}^N \sum_{k=1}^K p_{nk} \leq N \sum_{k=1}^K P_{k,\max}^{\text{UL}}. \quad (16)$$

By introducing this mild power control to problem (15), we need to solve the following problem

$$\begin{aligned} & \underset{\{\mathbf{v}_{nk}^{\text{VDL}}\}}{\text{minimize}} \quad \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2 \\ & \text{subject to} \quad (15b), (16). \end{aligned} \quad (17)$$

Similar to (13), constraint (15b) can be replaced with the following SOC constraint

$$\sqrt{\sum_{l \neq k} |(\mathbf{g}_{nl}^{\text{UL}})^H \mathbf{v}_{nl}^{\text{VDL}}|^2 + \sigma_n^2} \leq \frac{1}{\sqrt{\gamma_k^{\text{UL}}}} \Re \left( (\mathbf{g}_{nk}^{\text{UL}})^H \mathbf{v}_{nk}^{\text{VDL}} \right) \quad (18)$$

to make problem (17) a convex SOCP problem.

By exploiting the uplink-downlink duality and transforming the uplink model into a virtual downlink model, we address the challenge of the receive beamforming vectors scaling issue. As  $\{\mathbf{v}_{nk}^{\text{VDL}}\}$  have the same group sparsity pattern as  $\{\mathbf{v}_{nk}^{\text{UL}}\}$ , combining the downlink and virtual downlink parts, we relax  $\mathcal{P}_1$  to the following problem

$$\begin{aligned} & \underset{\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{VDL}}\}}{\text{minimize}} \quad \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2 + \sum_{n=1}^N \sum_{k=1}^K \frac{1}{\eta_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2 \\ & \quad + \sum_{n=1}^N \sum_{k=1}^K \mathbf{1}_{\{\bar{\mathbf{v}}_{nk} \neq \mathbf{0}\}} P_{nk}^{\text{C}} \\ & \text{subject to} \quad (12c), (13), (16), (18), \end{aligned} \quad (19)$$

where  $\bar{\mathbf{v}}_{nk}$  in the objective function is defined as  $\bar{\mathbf{v}}_{nk} = [(\mathbf{v}_{nk}^{\text{VDL}})^T, (\mathbf{v}_{nk}^{\text{DL}})^T]^T$ . It is observed that all the constraints of problem (19) are convex.

Although the feasible set is convex, problem (19) is a nonconvex integer programming problem due to the indicator function in the objective function. We identify that the third term in the objective function of problem (19) is a weighted  $\ell_0$ -norm of vector  $\bar{\mathbf{v}} = [\bar{\mathbf{v}}_{11}^T, \bar{\mathbf{v}}_{12}^T, \dots, \bar{\mathbf{v}}_{NK}^T]^T$  with weights  $P_{nk}^{\text{C}}$ 's, and it is non-convex. As  $\ell_1$ -norm is a well-known convex relaxation to  $\ell_0$ -norm, we relax the weighted  $\ell_0$ -norm as

$$\begin{aligned} \sum_{n=1}^N \sum_{k=1}^K \mathbf{1}_{\{\bar{\mathbf{v}}_{nk} \neq \mathbf{0}\}} P_{nk}^{\text{C}} & \approx \sum_{n=1}^N \sum_{k=1}^K \|\bar{\mathbf{v}}_{nk}\|_2 P_{nk}^{\text{C}} \\ & = \sum_{n=1}^N \sum_{k=1}^K \sqrt{\|\mathbf{v}_{nk}^{\text{DL}}\|_2^2 + \|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2} P_{nk}^{\text{C}}. \end{aligned} \quad (20)$$

Note that (20) is actually the weighted mixed  $\ell_{1,2}$ -norm of vector  $\bar{\mathbf{v}}$ . The mixed  $\ell_{1,2}$ -norm behaves like an  $\ell_1$ -norm on vector  $[\|\bar{\mathbf{v}}_{11}\|_2, \|\bar{\mathbf{v}}_{12}\|_2, \dots, \|\bar{\mathbf{v}}_{NK}\|_2]$ . The outer  $\ell_1$ -norm induces the sparsity structure, while the inner  $\ell_2$ -norm is responsible for forcing all coefficients in the beamforming group  $\bar{\mathbf{v}}_{nk}$  to be zero. By adopting mixed  $\ell_{1,2}$ -norm as the convex relaxation of the indication function term, we can induce the group sparsity structure of beamforming groups  $\{\bar{\mathbf{v}}_{nk}, n \in \mathcal{N}, k \in \mathcal{K}\}$ .

---

**Algorithm 1:** Mixed  $\ell_{1,2}$ -Norm Based Group Sparsity Inducing for Problem  $\mathcal{P}_1$ 


---

**Input:**  $\Theta^{\text{UL}}, \Theta^{\text{DL}}$ 1. Solve the convex relaxation problem (21) to obtain  $\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{VDL}}\}$ 2. Set  $\mathbf{v}_{nk}^{\text{UL}} = \mathbf{v}_{nk}^{\text{VDL}}, \forall n \in \mathcal{N}, k \in \mathcal{K}$ , and then solve problem  $\mathcal{P}_{1-2}$  with fixed  $\{\mathbf{v}_{nk}^{\text{UL}}\}$  to obtain  $\{p_k^{\text{UL}}\}$ **Output:**  $\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_k^{\text{UL}}\}$ 

By replacing the indication function with its convex surrogate, we relax problem (19) as the following convex problem

$$\begin{aligned} & \underset{\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{VDL}}\}}{\text{minimize}} && \sum_{n=1}^N \sum_{k=1}^K \|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2 + \sum_{n=1}^N \sum_{k=1}^K \frac{1}{\eta_n} \|\mathbf{v}_{nk}^{\text{DL}}\|_2^2 \\ & && + \sum_{n=1}^N \sum_{k=1}^K \sqrt{\|\mathbf{v}_{nk}^{\text{DL}}\|_2^2 + \|\mathbf{v}_{nk}^{\text{VDL}}\|_2^2} P_k^c \\ & \text{subject to} && (12c), (13), (16), (18). \end{aligned} \quad (21)$$

Once we solve problem (21) and obtain the solutions  $(\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{VDL}}\})$ , the receive beamforming vector  $\{\mathbf{v}_{nk}^{\text{UL}}\}$  in  $\mathcal{P}_1$  can be obtained by setting  $\mathbf{v}_{nk}^{\text{UL}} = \mathbf{v}_{nk}^{\text{VDL}}, \forall n \in \mathcal{N}, k \in \mathcal{K}$ , and  $\{p_k^{\text{UL}}\}$  can be obtained by solving problem  $\mathcal{P}_{1-2}$  with  $\{\mathbf{v}_{nk}^{\text{UL}}\}$  fixed, which is reduced to a linear program.

The overall algorithm for solving problem  $\mathcal{P}_1$  is summarized in Algorithm 1.

### B. Optimizing Variable $\Theta^{\text{DL}}$

For given  $\mathcal{B}_1$  and  $\mathcal{B}_3$ , we optimize  $\Theta^{\text{DL}}$  to solve the resulting problem, which is termed as  $\mathcal{P}_2$ . Since  $\Theta^{\text{DL}}$  does not appear in the objective function of  $\mathcal{P}$ , problem  $\mathcal{P}_2$  is in fact a downlink feasibility detection problem, which is given by

$$\mathcal{P}_2 : \text{ find } \Theta^{\text{DL}} \text{ subject to } (12a), (12e).$$

According to the SINR<sup>DL</sup> expression given in (7), we have

$$\begin{aligned} & \left| \sum_{n=1}^N \left( \mathbf{g}_{nk}^{\text{DL}} \right)^H \mathbf{v}_{nl}^{\text{DL}} \right|^2 \\ & \stackrel{(a)}{=} \left| \sum_{n=1}^N \left( \left( \mathbf{h}_{d,nk}^{\text{DL}} \right)^H + \left( \mathbf{h}_{r,k}^{\text{DL}} \right)^H \Theta^{\text{DL}} \mathbf{G}_n^{\text{DL}} \right) \mathbf{v}_{nl}^{\text{DL}} \right|^2 \\ & \stackrel{(b)}{=} \left| \left( \mathbf{h}_{d,k}^{\text{DL}} \right)^H \mathbf{v}_l^{\text{DL}} + \beta \left( \mathbf{a}^{\text{DL}} \right)^H \right. \\ & \quad \times \text{diag} \left( \left( \mathbf{h}_{r,k}^{\text{DL}} \right)^H \right) \tilde{\mathbf{G}}^{\text{DL}} \mathbf{v}_l^{\text{DL}} \left. \right|^2, \end{aligned} \quad (22)$$

where (a) follows by substituting (6), and (b) holds by defining  $\mathbf{h}_{d,k}^{\text{DL}} = [(\mathbf{h}_{d,1k}^{\text{DL}})^T, \dots, (\mathbf{h}_{d,Nk}^{\text{DL}})^T]^T$ ,  $\mathbf{v}_l^{\text{DL}} = [(\mathbf{v}_{1l}^{\text{DL}})^T, \dots, (\mathbf{v}_{Nl}^{\text{DL}})^T]^T$ ,  $\mathbf{a}^{\text{DL}} = [\theta_1^{\text{DL}}, \dots, \theta_M^{\text{DL}}]^H$ ,  $\tilde{\mathbf{G}}^{\text{DL}} = [\mathbf{G}_1^{\text{DL}}, \dots, \mathbf{G}_N^{\text{DL}}]$ . Note that in (22) the only term related to phase shifts is  $\mathbf{a}^{\text{DL}}$ . Therefore for notational ease, we define  $\mathbf{w}_{kl}^{\text{DL}} = \beta \text{diag}((\mathbf{h}_{r,k}^{\text{DL}})^H) \tilde{\mathbf{G}}^{\text{DL}} \mathbf{v}_l^{\text{DL}}$ ,  $b_{kl}^{\text{DL}} = (\mathbf{h}_{d,k}^{\text{DL}})^H \mathbf{v}_l^{\text{DL}}$ , and the SINR<sup>DL</sup> expression in (7) can be equivalently rewritten as  $\text{SINR}_k^{\text{DL}} = \frac{|b_{kk}^{\text{DL}} + (\mathbf{a}^{\text{DL}})^H \mathbf{w}_{kk}^{\text{DL}}|^2}{\sum_{l \neq k} |b_{kl}^{\text{DL}} + (\mathbf{a}^{\text{DL}})^H \mathbf{w}_{kl}^{\text{DL}}|^2 + \sigma_k^2}$ , leading to the

following inhomogeneous QCQP problem

$$\begin{aligned} & \text{find } \mathbf{a}^{\text{DL}} \\ & \text{subject to } |a_m^{\text{DL}}|^2 = 1, \forall m \in \mathcal{M}, \end{aligned} \quad (23a)$$

$$\frac{|b_{kk}^{\text{DL}} + (\mathbf{a}^{\text{DL}})^H \mathbf{w}_{kk}^{\text{DL}}|^2}{\sum_{l \neq k} |b_{kl}^{\text{DL}} + (\mathbf{a}^{\text{DL}})^H \mathbf{w}_{kl}^{\text{DL}}|^2 + \sigma_k^2} \geq \gamma_k, \quad \forall k \in \mathcal{K}. \quad (23b)$$

By introducing an auxiliary scalar  $t$ , and defining  $\mathbf{R}_{kl}^{\text{DL}} = \begin{bmatrix} \mathbf{w}_{kl}^{\text{DL}} (\mathbf{w}_{kl}^{\text{DL}})^H & \mathbf{w}_{kl}^{\text{DL}} (b_{kl}^{\text{DL}})^H \\ (\mathbf{w}_{kl}^{\text{DL}})^H b_{kl}^{\text{DL}} & 0 \end{bmatrix}$ , and  $\bar{\mathbf{a}}^{\text{DL}} = \begin{bmatrix} \mathbf{a}^{\text{DL}} \\ t^{\text{DL}} \end{bmatrix}$ , problem (23) can be transformed into the following homogeneous QCQP problem

$$\begin{aligned} & \text{find } \bar{\mathbf{a}}^{\text{DL}} \\ & \text{subject to } |\bar{a}_m^{\text{DL}}|^2 = 1, \text{ for } m = 1, \dots, M+1, \\ & \frac{(\bar{\mathbf{a}}^{\text{DL}})^H \mathbf{R}_{kk}^{\text{DL}} \bar{\mathbf{a}}^{\text{DL}} + |b_{kk}^{\text{DL}}|^2}{\sum_{l \neq k} (\bar{\mathbf{a}}^{\text{DL}})^H \mathbf{R}_{kl}^{\text{DL}} \bar{\mathbf{a}}^{\text{DL}} + |b_{kl}^{\text{DL}}|^2 + \sigma_k^2} \geq \gamma_k, \\ & \forall k \in \mathcal{K}. \end{aligned} \quad (24)$$

The non-convexity of problem (24) lies in the unit-modulus constraints. A common technique used to handle the nonconvex QCQP problems is matrix lifting. For ease of notation, we omit the superscript DL if it does not cause any ambiguity. In problem (24), as  $\bar{\mathbf{a}}^H \mathbf{R}_{kk} \bar{\mathbf{a}} = \text{Tr}(\mathbf{R}_{kk} \bar{\mathbf{a}} \bar{\mathbf{a}}^H)$ , and by introducing a new variable  $\mathbf{A} = \bar{\mathbf{a}} \bar{\mathbf{a}}^H \in \mathbb{C}^{(M+1) \times (M+1)}$ , we rewrite problem (24) as the following feasibility detection problem

$$\begin{aligned} & \text{find } \mathbf{A} \\ & \text{subject to } \text{Tr}(\mathbf{R}_{kk} \mathbf{A}) + |b_{kk}|^2 \geq \gamma_k \sum_{l \neq k} \text{Tr}(\mathbf{R}_{kl} \mathbf{A}) \\ & \quad + \gamma_k \left( \sum_{l \neq k} |b_{kl}|^2 + \sigma_k^2 \right), \forall k \in \mathcal{K}, \quad (25a) \\ & \mathbf{A}_{mm} = 1, \text{ for } m = 1, \dots, M+1, \quad (25b) \\ & \mathbf{A} \succeq \mathbf{0}, \text{ rank}(\mathbf{A}) = 1. \quad (25c) \end{aligned}$$

Here  $\mathbf{A} \succeq \mathbf{0}$  indicates that  $\mathbf{A}$  is a positive semidefinite (PSD) matrix. The challenge in solving problem (25) is the nonconvex rank-one constraint. The SDR technique is widely adopted to tackle the rank-one constraint in QCQP problems [47]. By simply dropping the nonconvex rank-one constraint, SDR relaxes the problem into a convex semidefinite programming (SDP) problem, which can then be solved by CVX. If a feasible  $\mathbf{A}$  with rank one is found, then  $\bar{\mathbf{a}}$  can be obtained by singular value decomposition (SVD) of  $\mathbf{A}$ .

However, such a relaxation may not be tight, i.e., the solution obtained by SDR may not satisfy the rank-one constraint. As pointed out in [48], the performance of SDR degrades sharply as the problem size grows. In our case, when  $M$  and/or  $K$  is large, the probability of returning a rank-one solution is low. If this is the case, additional steps (i.e., Gaussian randomization) are required to construct a rank-one solution from the higher-rank solution obtained by solving problem (25) [20], [47]. However, it is still possible that we



fail to find a feasible solution to problem (24) after a large number of Gaussian randomizations.

In other words, dropping the rank-one constraint cannot accurately detect the feasibility of problem (25). Hence, we propose a novel DC representation for the rank-one constraint, which is guaranteed to obtain a solution satisfying the non-convex rank-one constraint if problem (25) is feasible. Note that for the PSD matrix  $\mathbf{A}$ , the rank-one constraint indicates that  $\sigma_1(\mathbf{A}) > 0$  and  $\sigma_i(\mathbf{A}) = 0, \forall i = 2, \dots, M+1$ , where  $\sigma_i(\mathbf{A})$  is the  $i$ -th largest singular value of  $\mathbf{A}$ . And recall that the trace norm and spectral norm of  $\mathbf{A}$  are defined as  $\text{Tr}(\mathbf{A}) = \sum_{i=1}^{M+1} \sigma_i(\mathbf{A})$  and  $\|\mathbf{A}\| = \sigma_1(\mathbf{M})$ , respectively. Hence, the rank-one constraint can be equivalently rewritten as the difference of these two convex norms, i.e.,

$$\text{rank}(\mathbf{A}) = 1 \Leftrightarrow \text{Tr}(\mathbf{A}) - \|\mathbf{A}\| = 0, \text{Tr}(\mathbf{A}) > 0. \quad (26)$$

Based on the DC representation for the nonconvex rank-one constraint, we solve the following DC program

$$\begin{aligned} & \underset{\mathbf{A} \succeq \mathbf{0}}{\text{minimize}} \quad g(\mathbf{A}) := \text{Tr}(\mathbf{A}) - \|\mathbf{A}\| \\ & \text{subject to} \quad (25a)-(25b). \end{aligned} \quad (27)$$

Note that if the objective value  $g(\mathbf{A})$  is zero, the rank-one constraint is satisfied and we obtain a feasible solution to problem (25). Although problem (27) is still nonconvex due to the concave term  $-\|\mathbf{A}\|$  in the objective, we can adopt successive convex approximation to solve it in an iterative manner. Specifically, by linearizing the concave term, at iteration  $i$  we need to solve the following convex problem

$$\begin{aligned} & \underset{\mathbf{A} \succeq \mathbf{0}}{\text{minimize}} \quad \text{Tr}(\mathbf{A}) - \left\langle \partial \|\mathbf{A}^{[i-1]}\|, \mathbf{A} \right\rangle \\ & \text{subject to} \quad (25a)-(25b). \end{aligned} \quad (28)$$

where  $\mathbf{A}^{[i-1]}$  is the solution obtained at iteration  $i-1$  and  $\partial \|\mathbf{A}^{[i-1]}\|$  denotes the subgradient of spectral norm at point  $\mathbf{A}^{[i-1]}$ . Note that one subgradient of  $\|\mathbf{A}\|$  can be efficiently computed as  $\mathbf{q}_1 \mathbf{q}_1^H$ , where  $\mathbf{q}_1$  is the vector corresponding to the largest singular value  $\sigma_1(\mathbf{A})$  [49]. Given an initial  $\mathbf{A}^{[0]}$  and by iteratively solving (28) until the objective function  $g(\mathbf{A})$  in (27) becomes zero, we obtain an exact rank-one solution according to (26). We design a practical stopping criterion as  $\text{Tr}(\mathbf{A}) - \|\mathbf{A}\| < \epsilon_{DC}$ , where  $\epsilon_{DC}$  is a sufficiently small positive constant.

The convergence characteristic of the iterative DC algorithm for problem (27) is presented in the following proposition.

**Proposition 1:** The generated sequence  $\{g(\mathbf{A}^{[i]})\}$  is strictly decreasing and the sequence  $\{\mathbf{A}^{[i]}\}$  converges to a critical point of  $g$  from an arbitrary initial point  $\mathbf{A}^{[0]}$ .

*Proof:* Please refer to [48, Appendix B] for more details. ■

In fact, the proposed DC algorithm can always find a feasible  $\mathbf{A}$  to problem (25), which guarantees the objective value of problem (27) converges to zero. This is because the feasible region of problem (25) is always non-empty, at least the obtained solution to problem (25) at iteration  $t$  (i.e.,  $\Theta^{(t)}$ ) is still feasible at iteration  $t+1$ . Therefore, the strictly decreasing and non-negative sequence  $\{g(\mathbf{A}^{[i]})\}$  can always converge to zero within finite steps.

After obtaining a feasible  $\mathbf{A}$ , the phase-shift matrix  $\Theta^{\text{DL}}$  in  $\mathcal{P}_2$  can be recovered as follows. By SVD of  $\mathbf{A}$ , we can obtain the solution to problem (24) as  $\bar{\mathbf{a}}^{\text{DL}} = [\mathbf{a}_0^{\text{DL}} \ t_0^{\text{DL}}]^T$ , then the solution to problem (23) can be computed as  $\mathbf{a}^{\text{DL}} = \mathbf{a}_0^{\text{DL}}/t_0^{\text{DL}}$ , and the solution to  $\mathcal{P}_2$  is given as  $\Theta^{\text{DL}} = \beta \text{diag}((\mathbf{a}^{\text{DL}})^H)$ .

### C. Optimizing Variable $\Theta^{\text{UL}}$

As phase-shift matrix  $\Theta^{\text{UL}}$  does not appear in the objective function of  $\mathcal{P}$ , given  $\mathcal{B}_1$  and  $\mathcal{B}_2$ , the resulting problem denoted as  $\mathcal{P}_3$  is also a feasibility detection problem

$$\mathcal{P}_3 : \text{find } \Theta^{\text{UL}} \text{ subject to } (12b), (12f).$$

The same derivation process presented in the last subsection is also applicable to transform the uplink problem  $\mathcal{P}_3$  into a homogeneous QCQP, therefore details are omitted here for brevity. Specifically, by defining

$$\begin{aligned} \mathbf{w}_{nkl}^{\text{UL}} &= \beta \text{diag}\left(\left(\mathbf{h}_{r,l}^{\text{UL}}\right)^H\right) \mathbf{G}_n^{\text{UL}} \mathbf{v}_{nk}^{\text{UL}}, \\ b_{nkl} &= \left(\mathbf{h}_{d,nl}^{\text{UL}}\right)^H \mathbf{v}_{nk}^{\text{UL}}, \\ \mathbf{R}_{nkl}^{\text{UL}} &= \begin{bmatrix} \mathbf{w}_{nkl}^{\text{UL}} (\mathbf{w}_{nkl}^{\text{UL}})^H & \mathbf{w}_{nkl}^{\text{UL}} (b_{nkl}^{\text{UL}})^H \\ (\mathbf{w}_{nkl}^{\text{UL}})^H \mathbf{b}_{nkl}^{\text{UL}} & 0 \end{bmatrix}, \quad \bar{\mathbf{a}}^{\text{UL}} = \begin{bmatrix} \mathbf{a}^{\text{UL}} \\ t^{\text{UL}} \end{bmatrix}, \end{aligned}$$

we have the following uplink homogeneous QCQP problem

$$\begin{aligned} & \text{find } \bar{\mathbf{a}}^{\text{UL}} \\ & \text{subject to} \quad |\bar{a}_m^{\text{UL}}|^2 = 1, \text{ for } m = 1, \dots, M+1, \\ & \quad \frac{p_k^{\text{UL}} \left( (\bar{\mathbf{a}}^{\text{UL}})^H \mathbf{R}_{nkk}^{\text{UL}} \bar{\mathbf{a}}^{\text{UL}} + |b_{nkk}^{\text{UL}}|^2 \right)}{\sum_{l \neq k} p_l^{\text{UL}} \left( (\bar{\mathbf{a}}^{\text{UL}})^H \mathbf{R}_{nkl}^{\text{UL}} \bar{\mathbf{a}}^{\text{UL}} + |b_{nkl}^{\text{UL}}|^2 \right) + \sigma_n^2 \|\mathbf{v}_{nk}^{\text{UL}}\|_2^2} \\ & \quad \geq \gamma_k, \\ & \quad \forall n \in \mathcal{N}, k \in \mathcal{K}. \end{aligned} \quad (29)$$

Similar to the downlink counterpart, we can lift problem (29) to the following feasibility detection problem

$$\begin{aligned} & \text{find } \mathbf{A} \\ & \text{subject to} \quad \text{Tr}(\mathbf{R}_{nkk} \mathbf{A}) + |b_{nkk}|^2 \geq \gamma_k \sum_{l \neq k} \alpha_{kl} \text{Tr}(\mathbf{R}_{nkl} \mathbf{A}) \\ & \quad + \gamma_k \left( \sum_{l \neq k} \alpha_{kl} |b_{nkl}|^2 + c_{nk} \right), \forall n, \forall k, \quad (30a) \\ & \quad \mathbf{A}_{mm} = 1, \text{ for } m = 1, \dots, M+1, \quad (30b) \\ & \quad \mathbf{A} \succeq \mathbf{0}, \text{rank}(\mathbf{A}) = 1, \quad (30c) \end{aligned}$$

where  $\mathbf{A} = \bar{\mathbf{a}}^{\text{UL}} (\bar{\mathbf{a}}^{\text{UL}})^H$ ,  $\alpha_{kl} = p_l^{\text{UL}}/p_k^{\text{UL}}$  and  $c_{nk} = \sigma_n^2 \|\mathbf{v}_{nk}^{\text{UL}}\|_2^2/p_k^{\text{UL}}$ . Similarly, a feasible  $\mathbf{A}$  to problem (30) can be obtained by iteratively solving the following convex problem

$$\begin{aligned} & \underset{\mathbf{A} \succeq \mathbf{0}}{\text{minimize}} \quad \text{Tr}(\mathbf{A}) - \left\langle \partial \|\mathbf{A}^{[i-1]}\|, \mathbf{A} \right\rangle \\ & \text{subject to} \quad (30a)-(30b), \end{aligned} \quad (31)$$

until the stopping criterion  $\text{Tr}(\mathbf{A}) - \|\mathbf{A}\| < \epsilon_{DC}$  is satisfied.

The overall algorithm for solving problem  $\mathcal{P}_2$  (or  $\mathcal{P}_3$ ) is summarized in Algorithm 2.

---

**Algorithm 2:** DC Algorithm for Feasibility Detection Problem  $\mathcal{P}_2$  (or  $\mathcal{P}_3$ )

---

**Input:**  $\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{VDL}}\}, \{p_k^{\text{UL}}\}$ , initial point  $\mathbf{A}^{[0]}$  and set  $i = 0$   
**while**  $\text{Tr}(\mathbf{A}^{[i]}) - \|\mathbf{A}^{[i]}\| < \epsilon_{\text{DC}}$  **do**  
    1.  $i \leftarrow i + 1$ ,  $\partial\|\mathbf{A}^{[i-1]}\| = \mathbf{q}_1 \mathbf{q}_1^H$   
    2. Solve problem (28) (or problem (31))  
**end**  
Decompose  $\mathbf{A}$  as  $\mathbf{A} = \bar{\mathbf{a}} \bar{\mathbf{a}}^H$ ; Denote  $\bar{\mathbf{a}} = [\mathbf{a}_0, \mathbf{t}_0]^T$ ;  
Obtain  $\mathbf{a} = \mathbf{a}_0/\mathbf{t}_0$  and  $\Theta = \beta \text{diag}(\mathbf{a}^H)$   
**Output:**  $\Theta^{\text{UL}}$  (or  $\Theta^{\text{DL}}$ )

---

#### D. Unified BSO Approach

Based on the above discussions, problem  $\mathcal{P}$  is solved by iteratively solving problems  $\mathcal{P}_1$ ,  $\mathcal{P}_2$ , and  $\mathcal{P}_3$  in an alternating manner until convergence. We justify the effectiveness and depict the convergence behavior of the proposed BSO approach in the following proposition.

*Proposition 2:* With the BSO approach, the objective value of  $\mathcal{P}_1$  is non-increasing in the consecutive iterations.

*Proof:* For ease of notation, we denote the objective value of  $\mathcal{P}_1$  as  $f(\mathbf{V}, \Omega)$ , where the first variable  $\mathbf{V}$  is an abstraction of three optimization variables  $\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{UL}}\}$ , and  $\{p_k^{\text{UL}}\}$  in  $\mathcal{P}_1$ , and the second variable  $\Omega$  abstracts phase-shift matrices  $\Theta^{\text{DL}}$  and  $\Theta^{\text{UL}}$  in  $\mathcal{P}_2$  and  $\mathcal{P}_3$ . Assuming that  $(\mathbf{V}^{(t)}, \Omega^{(t)})$  is obtained at iteration  $t$ . If  $\mathcal{P}_2$  and  $\mathcal{P}_3$  are feasible, i.e.,  $(\mathbf{V}^{(t)}, \Omega^{(t+1)})$  exists,  $(\mathbf{V}^{(t)}, \Omega^{(t+1)})$  is also feasible to  $\mathcal{P}_1$ . Therefore,  $(\mathbf{V}^{(t)}, \Omega^{(t)})$  and  $(\mathbf{V}^{(t+1)}, \Omega^{(t+1)})$  are feasible solutions to  $\mathcal{P}_1$  at the consecutive iterations  $t$  and  $t + 1$ , respectively. We have the following inequality  $f(\mathbf{V}^{(t+1)}, \Omega^{(t+1)}) \stackrel{(a)}{\leq} f(\mathbf{V}^{(t)}, \Omega^{(t+1)}) \stackrel{(b)}{=} f(\mathbf{V}^{(t)}, \Omega^{(t)})$ , where (a) holds because  $\mathbf{V}^{(t+1)}$  is the optimal solution to  $\mathcal{P}_1$  for a given  $\Omega^{(t+1)}$  at iteration  $t + 1$ , and (b) holds because the objective value  $P_{\text{total}}$  depends only on the value of  $\mathbf{V}$  and is independent of  $\Omega$ . ■

The obtained beamforming vectors  $\{\mathbf{v}_{nk}^{\text{UL}}\}$  and  $\{\mathbf{v}_{nk}^{\text{DL}}\}$  after solving problem  $\mathcal{P}$  should have group sparse patterns. In the next section, we shall discuss how to design the task selection strategy  $\mathcal{A}$  based on the obtained solutions, and present the holistic three-stage framework for the network power minimization problem  $\mathcal{P}_{\text{original}}$ .

#### IV. THREE-STAGE FRAMEWORK FOR NETWORK POWER MINIMIZATION PROBLEM

In this section, we propose a thorough three-stage framework for problem  $\mathcal{P}_{\text{original}}$ . In *Stage 1*, we adopt the BSO with mixed  $\ell_{1,2}$ -norm and DC algorithm to induce the group sparsity structure of the uplink/downlink beamforming vectors and optimize the phase-shift matrices. The obtained solutions are served as a guideline to determine the inference task selection strategy  $\mathcal{A}$ . In *Stage 2*, an ordering rule is proposed for all tasks according to their priorities, which depend on the structured-sparse beamforming vectors obtained in *Stage 1* as well as several key system parameters (i.e., channel state

information, amplifier efficiency and task computation power). Based on the task ordering, we perform a task selection procedure to finalize  $\mathcal{A}$ . In *Stage 3*, the beamforming vectors and the phase-shift matrices are refined with the finalized task selection strategy  $\mathcal{A}$ .

##### A. Stage 1. Group Sparsity Inducing

With randomly initialized phase-shift matrices  $\Theta^{\text{UL}}$  and  $\Theta^{\text{DL}}$ , we repeatedly solve problems  $\mathcal{P}_1$ ,  $\mathcal{P}_2$ , and  $\mathcal{P}_3$  based on Algorithms 1 and 2 respectively in an alternating manner until the following stopping criterion is satisfied: the relative improvement of objective values of problem  $\mathcal{P}_1$  defined as  $(P_{\text{total}}^{(t)} - P_{\text{total}}^{(t+1)})/P_{\text{total}}^{(t)}$  is below a predefined threshold  $\epsilon$ , where  $P_{\text{total}}^{(t)}$  and  $P_{\text{total}}^{(t+1)}$  represent the objective values obtained in iterations  $t$  and  $t + 1$ , respectively. The yielded beamforming vectors should have the group sparsity structures. It is worth mentioning that the relative improvement is expected to be non-negative, because  $P_{\text{total}}$  is non-increasing as proved in Proposition 2.

##### B. Stage 2. Inference Task Selection

The next question is how to determine the task selection strategy  $\mathcal{A}$ . In fact, it is inappropriate to directly use the beamforming vectors obtained in *Stage 1* to finalize  $\mathcal{A}$ , as the vectors may contain some very small but nonzero coefficients. As illustrated in (11), these nonzero coefficients indicate that the corresponding tasks are performed by the BSs, which may result in unnecessary computation power consumption. To address this issue, we utilize the obtained group-structured beamforming vectors as well as other prior information to provide a guideline to determine  $\mathcal{A}$ .

For ease of exposition, we define the set of all task indices as  $\{(n, k) | n \in \mathcal{N}, k \in \mathcal{K}\}$ . The task indexed by  $(n, k)$  is considered to be *active* if  $k \in \mathcal{A}_n$ , and *inactive* otherwise. Specifically, the priority of task  $(n, k)$  is defined as  $\tau_{nk} = \sqrt{\frac{\|[\mathbf{g}_{nk}^{\text{UL}}, \mathbf{g}_{nk}^{\text{DL}}]\|_2^2 \eta_n}{P_{nk}^c}} \|\mathbf{v}_{nk}\|_2$ . We sort all  $NK$  tasks in a descending order according to their priorities, i.e.,  $\tau_{\pi_1} \geq \tau_{\pi_2} \cdots \geq \tau_{\pi_{NK}}$ , where  $\pi$  is a permutation of task indices  $(n, k)$ 's. Intuitively, if BS  $n$  has a higher power amplifier efficiency, a higher channel gain, and a higher beamforming gain with respect to MD  $k$ , but a lower computation power consumption, task  $(n, k)$  has a higher priority. A higher  $\tau_{nk}$  implies that task  $(n, k)$  is more power-efficient, therefore it is more likely to be activated.

To finalize the set  $\mathcal{A}$ , we need to detect the feasibility of a sequence of problems

$$\begin{aligned} & \text{find } \{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_k^{\text{UL}}\} \\ & \text{subject to (12a)-(12d), } \mathbf{v}_{\pi[j]} = \mathbf{0}, \end{aligned} \quad (32)$$

where  $\pi[j] = \{\pi_{j+1}, \dots, \pi_{NK}\}$  denotes the inactive task set at iteration  $j$ , and  $\mathbf{v}_{\pi[j]} = \mathbf{0}$  represents that all coefficients in those beamforming groups  $\mathbf{v}_{nk}$ 's with index  $(n, k) \in \pi[j]$  are set to zero. Note that the number of active tasks is within  $[K, NK]$ . Starting with  $j = K$ , we terminate the feasibility detection

---

**Algorithm 3:** BSO With Mixed  $\ell_{1,2}$ -Norm and DC Based Three-Stage Framework for Nonconvex Combinatorial Problem  $\mathcal{P}_{\text{original}}$ 


---

**Input:** initial phase-shift matrices  $\Theta^{\text{UL}}$  and  $\Theta^{\text{DL}}$ , and threshold  $\varepsilon$

*Stage 1:* Alternatively optimize  $\mathcal{B}_1$ ,  $\mathcal{B}_2$ , and  $\mathcal{B}_3$

**while** the improvement of the objective in problem  $\mathcal{P}_1$  is greater than  $\varepsilon$  **do**

1. solve  $\mathcal{P}_1$  for  $\mathbf{v}_{nk}^{\text{UL}}, p_k^{\text{UL}}, \mathbf{v}_{nk}^{\text{DL}}$  using Algorithm 1
2. solve  $\mathcal{P}_2$  for  $\Theta^{\text{DL}}$  using Algorithm 2
3. solve  $\mathcal{P}_3$  for  $\Theta^{\text{UL}}$  using Algorithm 2

**end**

*Stage 2:* Determine the inference task selection

1. Compute task priorities and sort all tasks in a descending order based on their priorities
2. Iteratively solve problem (32) until feasible to finalize the task selection strategy  $\mathcal{A}^*$

*Stage 3:* Solve problem  $\mathcal{P}(\mathcal{A}^*)$  to refine variables

$\{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_k^{\text{UL}}\}, \Theta^{\text{DL}}, \Theta^{\text{UL}}$

**Output:**  $\mathcal{A}^*, \{\mathbf{v}_{nk}^{\text{DL}}\}, \{\mathbf{v}_{nk}^{\text{UL}}\}, \{p_k^{\text{UL}}\}, \Theta^{\text{DL}}, \Theta^{\text{UL}}$

---

procedure and return  $\pi^{[j]}$  until problem (32) is feasible. The task selection strategy  $\mathcal{A}^*$  can be easily obtained from  $\pi^{[j]}$ .

Comparing problem (32) to  $\mathcal{P}_1$ , as the added constraint  $\mathbf{v}_{\pi^{[j]}} = \mathbf{0}$  is convex, we can solve problem (32) using Algorithm 1. Details are thus omitted here for brevity.

### C. Stage 3. Optimization Variables Refinement

After determining the task selection strategy  $\mathcal{A}^*$ , we should refine the beamforming vectors to make sure  $\mathbf{v}_{nk} = \mathbf{0}$  if  $k \notin \mathcal{A}_n, \forall n \in \mathcal{N}$ , and the phase-shift matrices at the RIS need to be refined as well. This can be achieved by solving problem  $\mathcal{P}(\mathcal{A}^*)$ . Since the group sparsity constraints (9f) are convex, the BSO with mixed  $\ell_{1,2}$ -norm and DC algorithm to solve problem  $\mathcal{P}$  is also applicable here to obtain the solutions. Details are thus omitted here.

### D. Complexity Discussions

The overall algorithm for green edge inference is summarized in Algorithm 3. The computational complexity of Algorithm 1 is dominated by solving the SOCP problem (21), which is  $\mathcal{O}(L^{3.5}K^{3.5})$  using interior-point methods [50]. The complexity of Algorithm 2 is dominated by iteratively solving the SDP problem (28) (or problem (31)) and computing the SVD of  $\mathbf{A}$ . An SDP problem with an  $a \times a$  semidefinite matrix variable and  $b$  SDP constraints is solved with complexity  $\mathcal{O}(\sqrt{a}(a^3b + a^2b^2 + b^3))$  by interior-point methods [51]. For problems (28) and (31), we have  $a = M + 1$ ,  $b = K + 1$  and  $a = M + 1$ ,  $b = NK + 1$ , respectively. It is also observed in our simulations that the convergence rate of the iterative procedure is fast (less than 10 iterations), therefore, the overall complexity of Algorithm 2 is  $\mathcal{O}(\sqrt{M}(M^3NK + M^2N^2K^2 + N^3K^3))$ . In short, the computational complexity involved in *Stage 1* is  $\mathcal{O}(TR)$ , where  $T$  is the required iterations for the BSO to converge, and

$R = L^{3.5}K^{3.5} + M^{3.5}NK + M^{2.5}N^2K^2 + M^{0.5}N^3K^3$  is the polynomial term.

We need to solve at most  $NK$  SOCP problems in *Stage 2*, and hence the worst-case complexity is  $\mathcal{O}(NK(L^{3.5}K^{3.5}))$ . The complexity involved in *Stage 3* is the same as that in *Stage 1*. Therefore, the overall complexity for the proposed three-stage framework is  $\mathcal{O}(TR + NL^{3.5}K^{4.5})$ . In other words, we propose a polynomial complexity algorithm for the combinatorial problem  $\mathcal{P}_{\text{original}}$ .

## V. SIMULATION RESULTS

In this section, we present the simulation results to verify the effectiveness of our proposed algorithm. We consider a network with five 8-antenna BSs and six MDs uniformly and randomly distributed in a square region of  $300 \text{ m} \times 300 \text{ m}$ . An RIS with 30 reflecting elements is located at the 3-dimensional coordinate (150, 0, 15). In addition, the BSs are with height 30 m (i.e., the coordinates of the BSs are  $(x, y, 30)$ ), while the MDs are with height 0 m (i.e., the coordinates of the MDs are  $(x, y, 0)$ ).

We consider the following distance-dependent path loss model  $L(d) = T_0(\frac{d}{d_0})^{-\alpha}$ , where  $T_0$  is the constant path loss at the reference distance  $d_0 = 1 \text{ m}$ ,  $d$  is the Euclidean distance between the transceivers, and  $\alpha$  is the path loss exponent. Each antenna of the BSs is assumed to have an isotropic radiation pattern (i.e., 0 dBi antenna gain), while each element of the RIS has a 3 dBi gain for fair comparison because it reflects signals only in its front half-space [29], [52]. Moreover, Rician fading is considered as the small-scale fading for the BS-RIS and RIS-MD channels. For example, the BS-RIS channel can be expressed as  $\mathbf{G}_n^x = \sqrt{\frac{\kappa_{\text{BR}}}{1+\kappa_{\text{BR}}}} \mathbf{G}_n^{\text{LOS}} + \sqrt{\frac{1}{1+\kappa_{\text{BR}}}} \mathbf{G}_n^{\text{NLOS}}$ , where  $\kappa_{\text{BR}}$  is the Rician factor representing the ratio of power between the deterministic line-of-sight (LOS) path and the scattered paths,  $\mathbf{G}_n^{\text{LOS}}$  is the LOS component modeled as the product of the steering vectors of the BS-RIS link [53],  $\mathbf{G}_n^{\text{NLOS}}$  is Rayleigh fading components with entries distributed as  $\mathcal{CN}(0, 1)$ , and  $x \in \{\text{UL}, \text{DL}\}$ . The entries in  $\mathbf{G}_n^x$  are then multiplied by the square root of distance-dependent path loss denoted by  $\alpha_{\text{BR}}$ . According to the 3GPP propagation environment from [54, Table B.1.2.1-1], we set  $\alpha_{\text{BR}} = 2.2$ . In addition, other channels are similarly generated with  $\alpha_{\text{BM}} = 3.67$  and  $\kappa_{\text{BM}} = 0$  (i.e., Rayleigh fading to account for rich scattering) for the BS-MD channel,  $\alpha_{\text{RM}} = 1.69$  and  $\kappa_{\text{RM}} = 3$  for the RIS-MD channel.

We consider a system with bandwidth 1 MHz and set  $T_0 = -30 \text{ dB}$ . The effective noise power for the BSs and the MDs are  $\sigma_n^2 = -90 \text{ dBm}$  and  $\sigma_k^2 = -80 \text{ dBm}$ , respectively. Without specified otherwise, other parameters are set as follows:  $P_{n,\text{max}}^{\text{DL}} = 1 \text{ W}$ ,  $P_{n,\text{max}}^{\text{UL}} = 1 \text{ W}$ ,  $P_{nk}^c = 0.45 \text{ W}$ ,  $\eta_n = 0.25$ ,  $\beta = 1$ ,  $\varepsilon = 10^{-2}$ ,  $\gamma_k^{\text{UL}} = \frac{1}{2}\gamma_k^{\text{DL}}$ , and  $\epsilon_{\text{DC}} = 10^{-6}$ .

We compare the proposed BSO with mixed  $\ell_{1,2}$ -norm and DC algorithm (abbreviated as BSO- $\ell_{1,2}$ -DC) with the following benchmarks.

- *Without-RIS:* Without the deployment of an RIS, the equivalent channels in (2) and (6) contain only the direct link, i.e.,  $\mathbf{h}_{r,k}^{\text{UL}} = \mathbf{h}_{r,k}^{\text{DL}} = \mathbf{0}, \forall k$ . As we do not need to

TABLE I  
COMPARISON BETWEEN SYSTEMS WITH AND WITHOUT AN RIS

	Target SINR [dB]	-20	-15	-10	-5	0	5	10
Feasible Probability	Without-RIS	0.98	0.85	0.54	0.12	0.01	0	0
	BSO- $\ell_{1,2}$ -RP	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>0.91</b>	<b>0.49</b>	<b>0.25</b>	<b>0.19</b>
Overall Power Consumption	Without-RIS [W]	3.22	3.32	3.55	5.75	14.35	N/A	N/A
	BSO- $\ell_{1,2}$ -RP [W]	<b>3.22</b>	<b>3.30</b>	<b>3.47</b>	<b>3.99</b>	<b>5.21</b>	<b>9.85</b>	<b>15.32</b>

optimize phase shifts in this case, the alternating process in *Stage 1* is simplified to solve  $\mathcal{P}_1$  only once.

- *BSO with mixed  $\ell_{1,2}$ -norm and Random Phase* (abbreviated as BSO- $\ell_{1,2}$ -RP): In this case, the phase shifts of all reflecting elements in both uplink and downlink transmissions are randomly chosen from  $[0, 2\pi)$  and then used to solve problem  $\mathcal{P}_1$ . We do not solve problems  $\mathcal{P}_2$  and  $\mathcal{P}_3$  in *Stage 1* subsequently to optimize the phase shifts. This benchmark is designed to reveal the necessity of optimizing the phase-shift matrices.
- *BSO with mixed  $\ell_{1,2}$ -norm and SDR* (abbreviated as BSO- $\ell_{1,2}$ -SDR): In this case, the nonconvex rank-one constraints in (25) and (30) are dropped. Gaussian randomization is then adopted to obtain a feasible solution to problems  $\mathcal{P}_2$  and  $\mathcal{P}_3$ . The number of randomly generated vectors for Gaussian randomization is set as 1000. If Gaussian randomization fails to find a feasible solution, we terminate the alternating process in *Stage 1*.
- *Decoupled BSO with mixed  $\ell_{1,2}$ -norm and DC* (abbreviated as BSO- $\ell_{1,2}$ -DC-decoupled): Instead of solving the coupled uplink and downlink power minimization problem (19) in *Stage 1*, in this case we solve  $\mathcal{P}_{1-1}$  and  $\mathcal{P}_{1-2}$  individually to derive the decoupled uplink and downlink designs. *Stage 2* and *Stage 3* are the same as BSO- $\ell_{1,2}$ -DC.

We first compare our RIS-aided communication system with the conventional one without the assistance of an RIS. For fair comparison, we do not explicitly optimize the phase-shift matrices, i.e., we compare the performance of Without-RIS and BSO- $\ell_{1,2}$ -RP.

We study the relationship between the feasible probability of problem  $\mathcal{P}_{\text{original}}$  and the target SINR  $\gamma_k^{\text{DL}}$ . The feasible probability of problem  $\mathcal{P}_{\text{original}}$  is defined as

$$\mathbb{P}\{\mathcal{P}_{\text{original}} \text{ is feasible}\} = \frac{\text{number of simulations } \mathcal{P}_{\text{original}} \text{ is feasible}}{\text{total number of simulations}}.$$

As the target SINR requirements become more stringent, i.e., larger values of  $\gamma_k^{\text{UL}}$  and  $\gamma_k^{\text{DL}}$ , the feasibility probability of problem  $\mathcal{P}_{\text{original}}$  is expected to decline. Results illustrated in Table I are averaged over 1000 independently generated channel realizations. We observe that the conventional system without RIS almost fails to support those settings with a target SINR being higher than 0 dB, while the RIS-aided system can still support with a high probability. In terms of the maximum SINR that the communication systems can support, we observe that there exists at least a 10 dB gain of the RIS-aided system over the system without RIS.

The lower part of Table I illustrates the overall network power consumption of systems with and without RIS. Under

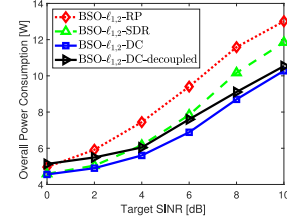


Fig. 3. Overall network power consumption versus target SINR  $\gamma_k^{\text{DL}}$ .

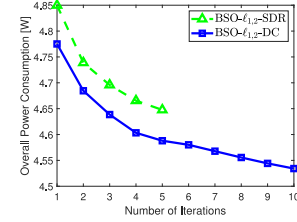


Fig. 4. Convergence behaviors of both BSO- $\ell_{1,2}$ -DC and BSO- $\ell_{1,2}$ -SDR algorithms.

the same SINR requirement, it is observed that BSO- $\ell_{1,2}$ -RP yields significantly lower power consumption. The supreme performance gain demonstrates that the deployment of an RIS in wireless communication systems can greatly boost the SINR and in turn reduce the overall network power consumption.

We also show the superiority of our proposed BSO- $\ell_{1,2}$ -DC algorithm in terms of the overall network power consumption, and the results obtained by averaging over 1000 independent channel realizations are shown in Fig. 3. The first observation is that both the BSO- $\ell_{1,2}$ -DC and BSO- $\ell_{1,2}$ -SDR algorithms significantly outperform BSO- $\ell_{1,2}$ -RP, which demonstrates that dynamically optimizing the phase-shift matrices according to the beamforming vectors can reduce the network power consumption to a large extent. In addition, we observe that the proposed BSO- $\ell_{1,2}$ -DC algorithm yields much lower overall network power consumption and is more energy efficient than the BSO- $\ell_{1,2}$ -SDR algorithm. Given an overall power budget (e.g., 8.5 W), BSO- $\ell_{1,2}$ -DC can achieve 1.5 dB higher SINR for the MDs than BSO- $\ell_{1,2}$ -SDR. Such a performance gain is mainly because BSO- $\ell_{1,2}$ -SDR may early terminate the alternating BSO process in *Stage 1* and cannot further proceed to find feasible solutions to problems  $\mathcal{P}_2$  and/or  $\mathcal{P}_3$ . To make this more explicit, Fig. 4 shows the objective values of problem  $\mathcal{P}_1$  in the first 10 alternating iterations in a specific channel realization. It is observed that as the BSO approach proceeds, the overall network power consumption for both BSO- $\ell_{1,2}$ -DC and BSO- $\ell_{1,2}$ -SDR algorithms are non-increasing, which validates our analysis in Proposition 2. It is also observed that

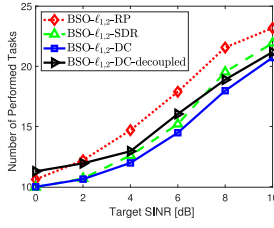


Fig. 5. Number of total performed tasks versus target SINR  $\gamma_k^{\text{DL}}$ .

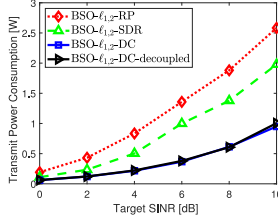


Fig. 6. Transmit power consumption versus target SINR  $\gamma_k^{\text{DL}}$ .

the BSO- $\ell_{1,2}$ -SDR algorithm terminates at the 5th iteration, as SDR fails to obtain feasible solutions to problems  $\mathcal{P}_2$  and/or  $\mathcal{P}_3$  even with Gaussian randomization techniques. In contrast, DC can always yield feasible solutions as we have discussed in Section III-B, therefore BSO- $\ell_{1,2}$ -DC terminates the alternating process only when the consecutive iterations make little progress. Comparing BSO- $\ell_{1,2}$ -DC-decoupled to BSO- $\ell_{1,2}$ -DC, we observe that BSO- $\ell_{1,2}$ -DC always yields lower power consumption than BSO- $\ell_{1,2}$ -DC-decoupled. This is because by coupling uplink and downlink beamforming designs, BSO- $\ell_{1,2}$ -DC provides more accurate information to the task ordering rule and thus arranging all tasks in a more reasonable order. Moreover, the performance gap between BSO- $\ell_{1,2}$ -DC-decoupled and BSO- $\ell_{1,2}$ -DC shrinks as the target SINR grows. This indicates that when the target SINR is high, the impact of task order on reducing the overall power consumption decreases because nearly all tasks should be selected and performed.

Another interesting point worth mentioning in Fig. 3 is that the performance gaps between BSO- $\ell_{1,2}$ -DC and BSO- $\ell_{1,2}$ -SDR are getting larger as the value of the target SINR increases, which indicates that BSO- $\ell_{1,2}$ -DC is especially appealing when high-quality services are required by the MDs. This is because a higher target SINR leads to a narrower feasible region of problems  $\mathcal{P}_2$  and  $\mathcal{P}_3$ , making SDR less likely to find a feasible solution. In short, Fig. 3 shows that the proposed BSO- $\ell_{1,2}$ -DC is able to reduce the overall network power consumption by 10% in the low SINR regime, and by up to 30% in the high SINR regime.

The number of tasks performed by all the BSs and the transmit power consumption versus the target SINR are shown in Fig. 5 and Fig. 6, respectively. As the target SINR increases, both the number of performed tasks and transmit power consumption increase. It is observed in Fig. 5 that BSO- $\ell_{1,2}$ -DC can always perform fewer tasks to satisfy a certain target SINR. In other words, the long-lasting alternating iterations of BSO- $\ell_{1,2}$ -DC shown in Fig. 4 helps in turn promote the group

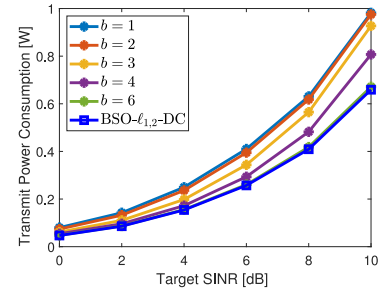


Fig. 7. Discrete and continuous phase shifts.

sparsity structure of beamforming vectors, thereby achieving lower computation power consumption. It is interesting to notice that BSO- $\ell_{1,2}$ -DC-decoupled selects even more tasks than BSO- $\ell_{1,2}$ -RP when the target SINR is 0 dB. This indicates the urgent demands to couple the uplink and downlink beamforming designs especially in the low SINR regime, so as to accurately characterize the tasks ordering rule and select the performed tasks reasonably. In terms of the transmit power consumption depicted in Fig. 6, we make the similar observation that BSO- $\ell_{1,2}$ -DC yields the lowest transmit power consumption. Finally, it is also observed that the performance gaps between BSO- $\ell_{1,2}$ -DC and other algorithms tend to be larger in the high SINR regime.

We compare the performances of discrete and continuous phase shifts in Fig. 7. With a  $b$ -bit resolution, the set of possible discrete phase shifts at each element is given as  $\mathcal{D} = \{0, \Delta\theta, 2\Delta\theta, \dots, (2^b - 1)\Delta\theta\}$ , where  $\Delta\theta = 2\pi/2^b$ . After solving problems  $\mathcal{P}_2$  and  $\mathcal{P}_3$  to obtain continuous phase shift at each reflecting element, we quantize it to its nearest neighbor in  $\mathcal{D}$ . The quantization process takes place in each iteration of *Stage 1*, which influences the subsequent task selection process in *Stage 2*. For fair comparison, in simulations we consider that all  $NK$  tasks are performed and only show the impact of discrete phase shifts on the transmit power consumption. It is observed that the performances of discrete and continuous phase shifts are comparable when  $b = 6$ , which indicates a 6-bit resolution discrete phase shifts is practically sufficient.

## VI. CONCLUSION

In this article, we investigated an RIS-aided edge inference system with multiple BSs cooperatively serving multiple MDs, taking into account both the uplink and downlink transmissions. The design of an energy-efficient edge inference system was formulated as a joint task selection strategy, uplink/downlink beamformers, transmit power, and uplink/downlink phase-shift matrices design problem. A BSO approach was proposed to decouple the optimization variables. For an efficient algorithm design, mixed  $\ell_{1,2}$ -norm was adopted to induce group sparsity of uplink/downlink beamforming vectors, while the matrix lifting and DC techniques were exploited to handle the nonconvex rank-one constraint and in turn solve the phase-shift matrix optimization problems. We also clarified the convergence behavior of the proposed BSO approach. Through extensive simulations,

$$\begin{aligned}
L(\mathbf{v}_{nk}^{\text{VDL}}; \lambda_{nk}) &= \sum_{n=1}^N \sum_{k=1}^K \left\| \mathbf{v}_{nk}^{\text{VDL}} \right\|_2^2 - \sum_{n=1}^N \sum_{k=1}^K \lambda_{nk} \left[ \left| \left( \mathbf{g}_{nk}^{\text{UL}} \right)^H \mathbf{v}_{nk}^{\text{VDL}} \right|^2 / \gamma_k - \sum_{l \neq k} \left| \left( \mathbf{g}_{nk}^{\text{UL}} \right)^H \mathbf{v}_{nl}^{\text{VDL}} \right|^2 - \sigma_n^2 \right] \\
&= \sum_{n=1}^N \sum_{k=1}^K \lambda_{nk} \sigma_n^2 + \sum_{n=1}^N \sum_{k=1}^K \left( \mathbf{v}_{nk}^{\text{VDL}} \right)^H \left[ \mathbf{I} - \frac{\lambda_{nk}}{\gamma_k} \mathbf{g}_{nk}^{\text{UL}} \left( \mathbf{g}_{nk}^{\text{UL}} \right)^H + \sum_{l \neq k} \lambda_{nl} \mathbf{g}_{nl}^{\text{UL}} \left( \mathbf{g}_{nl}^{\text{UL}} \right)^H \right] \mathbf{v}_{nk}^{\text{VDL}} \\
&= \sum_{n=1}^N \sum_{k=1}^K \lambda_{nk} \sigma_n^2 + \sum_{n=1}^N \sum_{k=1}^K \left( \mathbf{v}_{nk}^{\text{VDL}} \right)^H \left[ \mathbf{I} + \sum_{l=1}^K \lambda_{nl} \mathbf{g}_{nl}^{\text{UL}} \left( \mathbf{g}_{nl}^{\text{UL}} \right)^H - \left( 1 + \frac{1}{\gamma_k} \right) \lambda_{nk} \mathbf{g}_{nk}^{\text{UL}} \left( \mathbf{g}_{nk}^{\text{UL}} \right)^H \right] \mathbf{v}_{nk}^{\text{VDL}}
\end{aligned} \tag{33}$$

we demonstrated that deploying an RIS can significantly reduce the overall network power consumption. Furthermore, the effectiveness of the proposed DC algorithm in inducing low-rank solutions was also verified.

We identify the following relevant extensions as our future work. Developing a robust transmission design for RIS-aided edge inference with imperfect CSI is an interesting extension. Taking into account the RIS with coupled reflection amplitude and phase shifts is another interesting extension. In terms of scalability, the statistical CSI and the alternating direction method of multipliers (ADMM)-based parallel convex optimization framework can be utilized to further alleviate the computation burden for large-scale edge inference systems.

#### APPENDIX

We will show that problem (14) with  $P_{k,\max}^{\text{UL}} = +\infty, \forall k \in \mathcal{K}$  is equivalent to problem (15). As the fraction  $\frac{1}{N}$  in the objective does not affect the equivalence, we dismiss the fraction to simplify the presentation.

The Lagrangian of problem (15) is shown in (33) at the top of the page. The dual function is given by  $g(\lambda_{nk}) = \min_{\mathbf{v}_{nk}^{\text{VDL}}} L(\mathbf{v}_{nk}^{\text{VDL}}; \lambda_{nk})$ . The Lagrange dual problem is given as

$$\begin{aligned}
&\text{maximize}_{\{\lambda_{nk}\}} \sum_{n=1}^N \sum_{k=1}^K \lambda_{nk} \sigma_n^2 \\
&\text{subject to } \lambda_{nk} \geq 0, \mathbf{\Lambda}_n \succeq \left( 1 + \frac{1}{\gamma_k} \right) \lambda_{nk} \mathbf{g}_{nk}^{\text{UL}} \left( \mathbf{g}_{nk}^{\text{UL}} \right)^H, \\
&\quad \forall n \in \mathcal{N}, k \in \mathcal{K},
\end{aligned} \tag{34}$$

where  $\mathbf{\Lambda}_n = \mathbf{I} + \sum_l \lambda_{nl} \mathbf{g}_{nl}^{\text{UL}} \left( \mathbf{g}_{nl}^{\text{UL}} \right)^H$ , and strong duality holds.

Now we show that problem (14) is equivalent to problem (34). As  $\mathbf{v}_{nk}^{\text{UL}}$  is not in the objective of problem (14), the optimal  $\mathbf{v}_{nk}^{\text{UL}}$  that maximizes the SINR (14b) is the MMSE receiver given by  $\hat{\mathbf{v}}_{nk}^{\text{UL}} = \mathbf{\Gamma}^{-1} \mathbf{g}_{nk}^{\text{UL}}$ , where  $\mathbf{\Gamma} = \sum_{l=1}^K p_{nl} \mathbf{g}_{nl}^{\text{UL}} \left( \mathbf{g}_{nl}^{\text{UL}} \right)^H + \sigma_n^2 \mathbf{I}$  [55]. Substituting  $\hat{\mathbf{v}}_{nk}^{\text{UL}}$  into (14b), we get

$$\begin{aligned}
&\left( 1 + \frac{1}{\gamma_k} \right) p_{nk} \left( \mathbf{g}_{nk}^{\text{UL}} \right)^H \mathbf{\Gamma}^{-1} \mathbf{g}_{nk}^{\text{UL}} \geq 1 \\
&\stackrel{(a)}{\Leftrightarrow} \mathbf{\Gamma} \preceq \left( 1 + \frac{1}{\gamma_k} \right) p_{nk} \mathbf{g}_{nk}^{\text{UL}} \left( \mathbf{g}_{nk}^{\text{UL}} \right)^H,
\end{aligned}$$

where (a) uses the property  $\mathbf{y}^H \mathbf{X}^{-1} \mathbf{y} \geq 1 \Leftrightarrow \mathbf{X} \preceq \mathbf{y} \mathbf{y}^H$  proved in [55, Lemma 1]. By defining  $p_{nk} = \lambda_{nk} \sigma_n^2$ ,  $\mathbf{\Gamma} =$

$\sigma_n^2 \mathbf{\Lambda}_n$ , the uplink problem (14) is then equivalent to

$$\begin{aligned}
&\text{minimize}_{\{\lambda_{nk}\}} \sum_{n=1}^N \sum_{k=1}^K \lambda_{nk} \sigma_n^2 \\
&\text{subject to } \lambda_{nk} \geq 0, \mathbf{\Lambda}_n \preceq \left( 1 + \frac{1}{\gamma_k} \right) \lambda_{nk} \mathbf{g}_{nk}^{\text{UL}} \left( \mathbf{g}_{nk}^{\text{UL}} \right)^H, \\
&\quad \forall n \in \mathcal{N}, k \in \mathcal{K}.
\end{aligned} \tag{35}$$

The optimal  $\{\lambda_{nk}\}$  in both problems (34) and (35) will meet the SINR constraints with equality [44]. Therefore, we can safely reverse the SINR constraints as well as minimization to maximization in problem (35), which is now exactly the same as problem (34).

#### REFERENCES

- [1] S. Hua and Y. Shi, "Reconfigurable intelligent surface for green edge inference in machine learning," in *Proc. IEEE Globecom Workshops*, Waikoloa, Hawaii, Dec. 2019, pp. 1–6.
- [2] A. Zappone, M. Di Renzo, and M. Debbah, "Wireless networks design in the era of deep learning: Model-based, AI-based, or both?" *IEEE Trans. Commun.*, vol. 67, no. 10, pp. 7331–7376, Oct. 2019.
- [3] H. Sun, X. Chen, Q. Shi, M. Hong, X. Fu, and N. D. Sidiropoulos, "Learning to optimize: Training deep neural networks for interference management," *IEEE Trans. Signal Process.*, vol. 66, no. 20, pp. 5438–5453, Oct. 2018.
- [4] "Cisco visual networking index: Global mobile data traffic forecast update, 2017–2022," Cisco, San Jose, CA, USA, White Paper, Feb. 2019. [Online]. Available: <https://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white-paper-c11-738429.pdf>
- [5] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, and Y.-J. A. Zhang, "The roadmap to 6G: AI empowered wireless networks," *IEEE Commun. Mag.*, vol. 57, no. 8, pp. 84–90, Aug. 2019.
- [6] J. Zhang and K. B. Letaief, "Mobile edge intelligence and computing for the Internet of Vehicles," *Proc. IEEE*, vol. 108, no. 2, pp. 246–261, Feb. 2020.
- [7] S. Liu, Y. Lin, Z. Zhou, K. Nan, H. Liu, and J. Du, "On-demand deep model compression for mobile devices: A usage-driven model selection framework," in *Proc. 16th Annu. Int. Conf. Mobile Syst. Appl. Serv.*, Munich, Germany, Jun. 2018, pp. 389–400.
- [8] B. Taylor, V. S. Marco, W. Wolff, Y. Elkhatib, and Z. Wang, "Adaptive deep learning model selection on embedded systems," *ACM SIGPLAN Notices*, vol. 53, no. 6, pp. 31–43, 2018.
- [9] Z. Du *et al.*, "ShiDianNao: Shifting vision processing closer to the sensor," *ACM SIGARCH Comput. Archit. News*, vol. 43, no. 3, pp. 92–104, 2015.
- [10] J. Chen and X. Ran, "Deep learning with edge computing: A review," *Proc. IEEE*, vol. 107, no. 8, pp. 1655–1674, Aug. 2019.
- [11] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge computing: Vision and challenges," *IEEE Internet Things J.*, vol. 3, no. 5, pp. 637–646, Oct. 2016.
- [12] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge intelligence: Paving the last mile of artificial intelligence with edge computing," *Proc. IEEE*, vol. 107, no. 8, pp. 1738–1762, Aug. 2019.



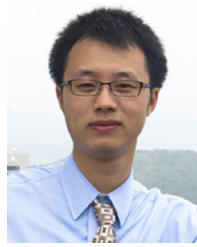
- [13] E. Li, L. Zeng, Z. Zhou, and X. Chen, "Edge AI: On-demand accelerating deep neural network inference via edge computing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 447–457, Jan. 2020.
- [14] C. Hu, W. Bao, D. Wang, and F. Liu, "Dynamic adaptive DNN surgery for inference acceleration on the edge," in *Proc. IEEE INFOCOM*, Paris, France, Apr. 2019, pp. 1423–1431.
- [15] T.-J. Yang, Y.-H. Chen, and V. Sze, "Designing energy-efficient convolutional neural networks using energy-aware pruning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Honolulu, HI, USA, Jul. 2017, pp. 6071–6079.
- [16] C. Louizos, K. Ullrich, and M. Welling, "Bayesian compression for deep learning," in *Proc. Neural Inf. Process. Syst. (NeurIPS)*, Long Beach, CA, USA, Dec. 2017, pp. 3288–3293.
- [17] K. Yang, Y. Shi, W. Yu, and Z. Ding, "Energy-efficient processing and robust wireless cooperative transmission for edge inference," *IEEE Internet Things J.*, vol. 7, no. 10, pp. 9456–9470, Oct. 2020.
- [18] X. Yang, S. Hua, Y. Shi, H. Wang, J. Zhang, and K. B. Letaief, "Sparse optimization for green edge AI inference," *J. Commun. Inf. Netw.*, vol. 5, no. 1, pp. 1–15, 2020.
- [19] E. Basar, M. Di Renzo, J. De Rosny, M. Debbah, M.-S. Alouini, and R. Zhang, "Wireless communications through reconfigurable intelligent surfaces," *IEEE Access*, vol. 7, pp. 116753–116773, 2019.
- [20] Q. Wu and R. Zhang, "Intelligent reflecting surface enhanced wireless network via joint active and passive beamforming," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5394–5409, Nov. 2019.
- [21] C. Huang, A. Zappone, G. C. Alexandropoulos, M. Debbah, and C. Yuen, "Reconfigurable intelligent surfaces for energy efficiency in wireless communication," *IEEE Trans. Wireless Commun.*, vol. 18, no. 8, pp. 4157–4170, Aug. 2019.
- [22] N. Rajatheva et al., "White paper on broadband connectivity in 6G," 2020. [Online]. Available: <https://arxiv.org/abs/2004.14247>.
- [23] Q.-U.-A. Nadeem, A. Kammoun, A. Chaaban, M. Debbah, and M.-S. Alouini, "Asymptotic max-min SINR analysis of reconfigurable intelligent surface assisted MISO systems," *IEEE Trans. Wireless Commun.*, vol. 19, no. 12, pp. 7748–7764, Dec. 2020.
- [24] C. Huang et al., "Holographic mimo surfaces for 6G wireless networks: Opportunities, challenges, and trends," *IEEE Wireless Commun.*, vol. 27, no. 5, pp. 118–125, Oct. 2020.
- [25] M. Di Renzo et al., "Smart radio environments empowered by reconfigurable AI meta-surfaces: An idea whose time has come," *EURASIP J. Wireless Commun. Netw.*, vol. 2019, no. 1, pp. 1–20, Jan. 2019.
- [26] C. Liaskos, S. Nie, A. Tsioliaridou, A. Pitsillides, S. Ioannidis, and I. Akyildiz, "A new wireless communication paradigm through software-controlled metasurfaces," *IEEE Commun. Mag.*, vol. 56, no. 9, pp. 162–169, Sep. 2018.
- [27] H. Han, J. Zhao, D. Niyato, M. D. Renzo, and Q.-V. Pham, "Intelligent reflecting surface aided network: Power control for physical-layer broadcasting," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Dublin, Ireland, Jun. 2020, pp. 1–7.
- [28] M. Fu, Y. Zhou, and Y. Shi, "Intelligent reflecting surface for downlink non-orthogonal multiple access networks," in *Proc. IEEE Globecom Workshops*, Waikoloa, HI, USA, Dec. 2019, pp. 1–6.
- [29] Q. Wu and R. Zhang, "Weighted sum power maximization for intelligent reflecting surface aided SWIPT," *IEEE Wireless Commun. Lett.*, vol. 9, no. 5, pp. 586–590, May 2020.
- [30] Z. Wang, Y. Shi, Y. Zhou, H. Zhou, and N. Zhang, "Wireless-powered over-the-air computation in intelligent reflecting surface-aided IoT networks," *IEEE Internet Things J.*, vol. 8, no. 3, pp. 1585–1598, Feb. 2021.
- [31] C. Huang, G. C. Alexandropoulos, C. Yuen, and M. Debbah, "Indoor signal focusing with deep learning designed reconfigurable intelligent surfaces," in *Proc. 20th IEEE Int. Workshop Signal Process. Adv. Wireless Commun. (SPAWC)*, Cannes, France, Jul. 2019, pp. 1–5.
- [32] G. C. Alexandropoulos, S. Samarakoon, M. Bennis, and M. Debbah, "Phase configuration learning in wireless networks with multiple reconfigurable intelligent surfaces," 2020. [Online]. Available: [arXiv:2010.04376](https://arxiv.org/abs/2010.04376).
- [33] K. Li, M. Tao, and Z. Chen, "Exploiting computation replication for mobile edge computing: A fundamental computation-communication tradeoff study," *IEEE Trans. Wireless Commun.*, vol. 19, no. 7, pp. 4563–4578, Jul. 2020.
- [34] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 4, pp. 2322–2358, 4th Quart., 2017.
- [35] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Int. Conf. Neural Inf. Process. Syst. (NeurIPS)*, Lake Tahoe, NV, USA, Dec. 2012, pp. 1106–1114.
- [36] Q.-U.-A. Nadeem, H. Alwazani, A. Kammoun, A. Chaaban, M. Debbah, and M.-S. Alouini, "Intelligent reflecting surface assisted multi-user MISO communication: Channel estimation and beamforming design," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 661–680, 2020.
- [37] T.-J. Yang, Y.-H. Chen, J. S. Emer, and V. Sze, "A method to estimate the energy consumption of deep neural networks," in *Proc. 51st Asilomar Conf. Signals Syst. Comput. (ASCSSC)*, Pacific Grove, CA, USA, Oct. 2017, pp. 1916–1920.
- [38] X. Xu et al., "Scaling for edge inference of deep neural networks," *Nat. Electron.*, vol. 1, no. 4, pp. 216–222, 2018.
- [39] R. Mahapatra, Y. Nijssure, G. Kaddoum, N. U. Hassan, and C. Yuen, "Energy efficiency tradeoff mechanism towards wireless green communication: A survey," *IEEE Commun. Surveys Tuts.*, vol. 18, no. 1, pp. 686–705, 1st Quart., 2016.
- [40] CNN Energy Estimation Website. [Online]. Available: <http://energyestimation.mit.edu>
- [41] Y.-H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks," *IEEE J. Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, Jan. 2017.
- [42] I. Hwang, B. Song, and S. S. Soliman, "A holistic view on hyper-dense heterogeneous and small cell networks," *IEEE Commun. Mag.*, vol. 51, no. 6, pp. 20–27, Jun. 2013.
- [43] M. Hong, M. Razaviyayn, Z.-Q. Luo, and J.-S. Pang, "A unified algorithmic framework for block-structured optimization involving big data: With applications in machine learning and signal processing," *IEEE Signal Process. Mag.*, vol. 33, no. 1, pp. 57–77, Jan. 2016.
- [44] A. Wiesel, Y. C. Eldar, and S. Shamai, "Linear precoding via conic optimization for fixed MIMO receivers," *IEEE Trans. Signal Process.*, vol. 54, no. 1, pp. 161–176, Jan. 2006.
- [45] M. Grant and S. Boyd. (2014). *CVX: MATLAB Software for Disciplined Convex Programming, Version 2.1*. [Online]. Available: <http://cvxr.com/cvx>
- [46] S. Luo, R. Zhang, and T. J. Lim, "Downlink and uplink energy minimization through user association and beamforming in C-RAN," *IEEE Trans. Wireless Commun.*, vol. 14, no. 1, pp. 494–508, Jan. 2015.
- [47] Z.-Q. Luo, W.-K. Ma, A. M.-C. So, Y. Ye, and S. Zhang, "Semidefinite relaxation of quadratic optimization problems and applications," *IEEE Signal Process. Mag.*, vol. 27, no. 3, pp. 20–34, May 2010.
- [48] K. Yang, T. Jiang, Y. Shi, and Z. Ding, "Federated learning via over-the-air computation," *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 2022–2035, Mar. 2020.
- [49] P. D. Tao and L. T. H. An, "Convex analysis approach to DC programming: Theory, algorithms and applications," *Acta Mathematica Vietnamica*, vol. 22, no. 1, pp. 289–355, Jan. 1997.
- [50] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [51] I. Pólik and T. Terlaky, *Interior Point Methods for Nonlinear Optimization*. Berlin, Germany: Springer, 2010.
- [52] Q. Wu and R. Zhang, "Beamforming optimization for wireless network aided by intelligent reflecting surface with discrete phase shifts," *IEEE Trans. Commun.*, vol. 68, no. 3, pp. 1838–1851, Mar. 2020.
- [53] S. Zhang and R. Zhang, "Capacity characterization for intelligent reflecting surface aided MIMO communication," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 8, pp. 1823–1838, Aug. 2020.
- [54] *Further Advancements for E-UTRA Physical Layer Aspects (Release 9)*, 3GPP Standard TS 36.814, Mar. 2010.
- [55] W. Yu and T. Lan, "Transmitter optimization for the multi-antenna downlink with per-antenna power constraints," *IEEE Trans. Signal Process.*, vol. 55, no. 6, pp. 2646–2660, Jun. 2007.



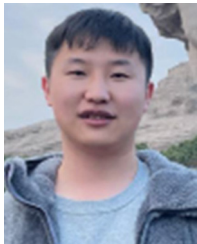
**Sheng Hua** (Graduate Student Member, IEEE) received the B.S. degree from the School of Computer Science, Xidian University, Xi'an, China, in 2018. He is currently pursuing the master's degree with the School of Information Science and Technology, ShanghaiTech University, Shanghai, China. His current research interests include reconfigurable intelligent surface and edge artificial intelligence.



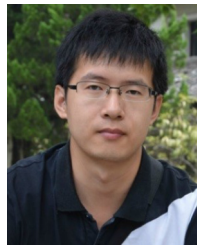
**Yong Zhou** (Member, IEEE) received the B.Sc. and M.Eng. degrees from Shandong University, Jinan, China, in 2008 and 2011, respectively, and the Ph.D. degree from the University of Waterloo, Waterloo, ON, Canada, in 2015. From November 2015 to January 2018, he worked as a Postdoctoral Research Fellow with the Department of Electrical and Computer Engineering, The University of British Columbia, Vancouver, Canada. He is currently an Assistant Professor with the School of Information Science and Technology, ShanghaiTech University, Shanghai, China. His research interests include Internet of Things, edge computing, and reconfigurable intelligent surface.



**Yuanming Shi** (Senior Member, IEEE) received the B.S. degree in electronic engineering from Tsinghua University, Beijing, China, in 2011, and the Ph.D. degree in electronic and computer engineering from the Hong Kong University of Science and Technology in 2015. Since September 2015, he has been with the School of Information Science and Technology, ShanghaiTech University, where he is currently a tenured Associate Professor. He visited the University of California, Berkeley, CA, USA, from October 2016 to February 2017. His research areas include optimization, statistics, machine learning, signal processing, and their applications to 6G, IoT, and AI. She is a recipient of the 2016 IEEE Marconi Prize Paper Award in Wireless Communications, and the 2016 Young Author Best Paper Award by the IEEE Signal Processing Society. He is an Editor of IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS and IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS.



**Kai Yang** (Member, IEEE) received the B.S. degree from Dalian University of Technology in 2015, and the Ph.D. degree from the University of Chinese Academy of Sciences, Beijing, in 2020. He is currently an AI Engineer in JD Technology Group focusing on designing efficient systems and algorithms of machine learning (federated learning in particular) and their applications in financial risk management.



**Kunlun Wang** (Member, IEEE) received the Ph.D. degree in electronic engineering from Shanghai Jiao Tong University, Shanghai, China, in 2016. From 2016 to 2017, he was with Huawei Technologies Company Ltd., where he was involved in energy efficiency algorithm design. From 2017 to 2019, he was with the Key Laboratory of Wireless Sensor Network and Communication, SIMIT, Chinese Academy of Sciences, Shanghai, China. From 2019 to 2020, he was with the School of Information Science and Technology, ShanghaiTech University. Since 2021, he has been a Professor with the School of Information Science and Technology, East China Normal University. His current research interests include energy efficient communications, fog computing networks, resource allocation, and optimization algorithm.