

# Joint Mode Selection and Transceiver Design for Device-to-Device Communications Underlying Multi-User MIMO Cellular Networks

Jingran Lin, Qingjiang Shi, Qiang Li, and Dongmei Zhao

**Abstract**—Consider a network consisting of one multi-antenna base station (BS) and multiple pairs of multi-antenna user equipments (UEs). For each UE pair, the communication between transmitter and receiver is established either through BS or via device-to-device (D2D) link. We assume that D2D transmission and cellular transmission are equally prioritized and share the same resources. To improve the network throughput, we maximize the sum rate by jointly optimizing the transmission mode of each UE pair and the associated transceivers. Due to the NP-hardness of this problem, we first perform some efficient approximation to it and then design an iterative algorithm, which is guaranteed to converge to a stationary solution by solving a series of weighted minimum mean square error (WMMSE) problems. The proposed algorithm has two distinguishing features. First, it only solves the WMMSE problem inexactly in each iteration, which thereby has a simplified algorithm structure and accelerated convergence behavior than the classical one. Second, we further fit the WMMSE problem into the alternating direction method of multipliers (ADMM) framework, making it amenable to parallel and distributed computation. Finally, the approximated problem is solved efficiently and distributively, with simple closed-form solutions in each step.

**Index Terms**—Device-to-device (D2D) communication, mode selection, transceiver design, weighted minimum mean square error (WMMSE), alternating direction method of multipliers (ADMM).

## I. INTRODUCTION

Device-to-device (D2D) communication, which allows two user equipments (UEs) to exchange information directly without the intervention of a base station (BS), has been widely accepted as an important key technology in the fifth-generation (5G) mobile networks [1], [2]. By exploiting the proximity gain, reuse gain, hop gain, and pairing gain, D2D can remarkably improve the network performance in terms of spectral and energy efficiencies [3]–[5]. However, to achieve these potential benefits, D2D communications should be appropriately

managed. Many important yet challenging issues, such as D2D discovery, D2D synchronization, transmission mode selection, wireless resource allocation, power control, interference mitigation, etc., need to be carefully addressed [6], [7]. In this paper, we focus on the joint transmission mode selection and interference mitigation problem for the D2D communication underlying cellular networks. Although such a problem has been intensively investigated, most current studies impose some constraints onto the setting of D2D communications, e.g., the UE status [4], [8], the UE/antenna number [8], [9], the resource sharing strategy [7], the priority of D2D transmission [4], etc. As a consequence, the D2D performance gains are not fully utilized. Naturally, to maximally explore the D2D potentials, a more general and flexible formulation of the joint mode selection and interference mitigation problem should be considered.

### A. State of the Art

So far, various schemes have been proposed for joint mode selection and interference mitigation in D2D communications. Basically, the literature can be categorized according to the network size and the way of algorithm execution (centralized or distributed).

1) *Network Size*: Most early studies formulated this problem in single-antenna networks which consist of only one cellular UE and one pair of D2D UEs [8]–[10]. They usually selected the transmission mode (e.g., cellular mode or underlay/overlay D2D mode) for the D2D UE pair by exhaustive search, and then optimized the power allocation to achieve the best performance in throughput and/or power efficiency. Recently, to accommodate the explosive demand for wireless data, the network size (e.g., the UE number and/or the antenna number) has been increasing sharply and inevitably. The approach of exhaustive search is thereby computationally prohibitive, and some advanced techniques, such as game theory [11], graph theory [12], and nonlinear optimization [13], have been introduced to help the mode selection and interference mitigation. Specifically, as far as the multi-antenna D2D communication is concerned, beamforming is widely utilized to manage the interference. For instance, many studies applied zero-forcing (ZF) beamforming to avoid the interference between cellular UEs and D2D UEs [14]–[16]. A more general beamforming scheme was addressed in [17], where the transceivers were designed to maximize the entire D2D and cellular transmission rates, and the sequential quadratic programming method was

J. Lin and Q. Li are with the School of Information and Communication Engineering, University of Electronic Science and Technology of China, Chengdu 611731, China (e-mail: jingranlin@uestc.edu.cn, and lq@uestc.edu.cn). Q. Shi is with the School of Software Engineering, Tongji University, Shanghai 200092, China (e-mail: qing.j.shi@gmail.com). D. Zhao is with the ZTE Corporation, Shenzhen 518063, China (e-mail: zdmzpz@sina.com).

The work of J. Lin was supported in part by the natural Science Foundation of China (NSFC) under Grant 61671120, and in part by the Sichuan Science and Technology Program under Grant 2018JY0147. The work of Q. Shi was supported in part by NSFC under Grant 61671411, in part by National Key Research and Development Project under grant 2017YFE0119300, and in part by the Fundamental Research Funds for the Central Universities under grant 22120180113. The work of Q. Li was supported in part by NSFC under Grant 61531009.

Part of this work has been reported in IEEE ICASSP 2016 [1].

adopted to solve this problem. This approach outperforms the ZF approach, but does not involve the issue of mode selection. The authors in [18] considered the joint mode selection and transceiver design for rate maximization in multi-input multi-output (MIMO) D2D networks, and then solved the problem by applying the approach of successive convex approximation (SCA). In this work, however, each UE pair may perform cellular and D2D transmissions simultaneously, thus rendering heavy operational burdens. Obviously, a more practical mode selection strategy should be exclusive, i.e., the communication within each UE pair is established either through base station (BS) or via D2D link.

2) *Algorithm Execution*: In addition to the centralized algorithms, many distributed algorithms have been proposed since it may not be easy to solve the high-dimensional D2D communication problems as the network size increases. For instance, several heuristic algorithms were developed to distributively solve the joint mode selection and power control problem [19], [20]. In [21], some energy-splitting variables were introduced such that the mode selection and the resource allocation could be decoupled and optimized independently. Another popular approach worth mentioning is game theory, which has been successfully applied to tackle many distributed selfish optimization problems about joint mode selection and resource allocation in D2D communications [22], [23].

In summary, despite the differences in problem formulation and algorithm design, most current studies divided the UEs into cellular UEs and D2D UEs in advance, and fixed their statuses during the entire communication process. Moreover, they usually took cellular UEs as the primary users, and thus their quality of service (QoS) requests were delivered with priority. In this circumstance, one important task of D2D management was to limit the interference from D2D UEs to cellular UEs. To guarantee this, cellular UEs were set to communicate with BS constantly, while only D2D UEs were involved in mode selection. In addition, to better protect cellular transmissions from D2D interference in underlay D2D communications, D2D UEs usually reused the cellular resources in uplink only, while keeping idle in downlink [7], [24], [25].

Obviously, this has significantly affected the flexibility of D2D management. Thus, the performance gains of D2D are not fully utilized. Actually, there are still continuing debates in the Third Generation Partnership Project (3GPP) on the mode selection between cellular and D2D communications and its priority, although a lot of progress has been made in the standardization of D2D discovery, channel models, deployment scenarios, evaluation methodologies, and so on [24]–[26]. In this paper, we relax these constraints and formulate the D2D communication problem from a more general perspective. A key point is that we equally prioritize the cellular transmission and the D2D transmission. Let us explain the reasonability of this setting. The idea that the D2D UEs exchange data by reusing the cellular resources is quite similar to the concept of the secondary user (SU) introduced in cognitive radio (CR) networks [5]. This is one possible reason why a lot of studies have prioritized cellular UEs over D2D UEs, and imposed some constraints on D2D communications. However, there are some essential differences between D2D and CR. In particular,

the D2D connection utilizes licensed spectrum bands and is supervised by a central entity (e.g., the cellular BS), whereas in CR the SU is not controlled by the primary user (PU) networks. In brief, the involvement of the cellular network in the control plane is the key difference between D2D and CR [2], [5], [25]. Since D2D communications are fully controlled by the cellular network, to maximally utilize the performance gains of D2D, it may be more reasonable to assume that cellular UEs and D2D UEs have the same priority. A similar idea has been mentioned in [4], [7], [19], but no details have been provided on how to fulfill it.

## B. Contributions and Organization

In this paper, we consider the D2D communication underlaying multi-user MIMO networks. By jointly optimizing the transmission mode and the associated transceiver for each UE pair, we aim to maximize the network throughput. Departing from most current studies, we formulate this problem from a more general perspective. Specifically, we prioritize the cellular transmission and D2D transmission equally. They share the same resources in both uplink and downlink. The status of each UE pair (cellular or D2D) is not specified in advance, and all the UE pairs are involved in mode selection. The transmission mode of each UE pair can freely switch between cellular and D2D, depending on its contribution to the total throughput and its interference to other UE pairs. This setting endows more flexibilities to D2D management, but yields a challenging NP-hard problem [27]. As a compromise, we pursue some efficient approximate solutions with manageable complexity.

Specifically, by performing the weighted minimum mean square error (WMMSE) reformulation [28], [29], we develop an algorithm which handles this problem by iteratively solving a series of WMMSE problems. In particular, the WMMSE problem in each iteration is only solved inexactly to simplify the algorithm structure and accelerate the convergence. For this reason, the algorithm is referred to as the iterative algorithm based on inexact WMMSE (IA-IWMMSE). We show that IA-IWMMSE is guaranteed to converge to a stationary solution. Next, to facilitate the algorithm execution, we further fit the WMMSE problem into the framework of alternating direction method of multipliers (ADMM) [30], so that IA-IWMMSE can be executed in a distributed manner. Finally, the problem can be solved distributively and efficiently, with a simple closed-form solution in each step.

The rest of this paper is organized as follows. The system model and problem formulation are given in Section II. In Section III, we propose IA-IWMMSE to solve the problem and show its convergence. In Section IV, we next show how to distributively solve the WMMSE problem in IA-WMMSE, with the help of ADMM. Simulation results are provided in Section V. Section VI concludes this paper.

*Notations*: Boldface capital and little letters denote matrices and vectors, respectively. Italic letters denote scalars. For a given matrix  $\mathbf{X}$ , we denote its transpose, Hermitian, and inverse by  $\mathbf{X}^T$ ,  $\mathbf{X}^\dagger$ , and  $\mathbf{X}^{-1}$ , respectively. Similarly, we denote the transpose and Hermitian of vector  $\mathbf{x}$  by  $\mathbf{x}^T$  and  $\mathbf{x}^\dagger$ . We use  $\|\cdot\|_p$  to denote the  $l_p$ -norm,  $p = 0, 1, 2$ . Little subscripts “ $m$ ”

TABLE I  
SUMMARY OF THE VARIABLES USED IN THIS PAPER

Uplink (UL)	$\mathbf{G}_m^U$	Cellular UL channel between TX_UE <sub>m</sub> and BS
	$\mathbf{F}_{m,n}^U$	D2D UL channel between TX_UE <sub>n</sub> and RX_UE <sub>m</sub>
	$\mathbf{v}_m^U$	UL TX beamformer of TX_UE <sub>m</sub>
	$\mathbf{u}_{d,m}^U$	D2D UL RX beamformer of RX_UE <sub>m</sub> for TX_UE <sub>m</sub>
	$\mathbf{u}_{c,m}^U$	Cellular UL RX beamformer of BS for TX_UE <sub>m</sub>
Downlink (DL)	$\mathbf{G}_m^D$	Cellular DL channel between BS and RX_UE <sub>m</sub>
	$\mathbf{F}_{m,n}^D$	D2D DL channel between TX_UE <sub>n</sub> and RX_UE <sub>m</sub>
	$\mathbf{v}_{d,m}^D$	D2D DL TX beamformer of TX_UE <sub>m</sub> for RX_UE <sub>m</sub>
	$\mathbf{v}_{c,m}^D$	Cellular DL TX beamformer of BS for RX_UE <sub>m</sub>
	$\mathbf{u}_{d,m}^D$	D2D DL RX beamformer of RX_UE <sub>m</sub> for TX_UE <sub>m</sub>
	$\mathbf{u}_{c,m}^D$	Cellular DL RX beamformer of RX_UE <sub>m</sub> for BS

and “ $n$ ” denote the UE index; little subscripts “ $c$ ” and “ $d$ ” denote cellular mode and D2D mode; capital superscripts “U” and “D” denote uplink and downlink. For instance,  $\mathbf{x}_{d,m}^U$  is some vector variable associated with UE  $m$  in the uplink D2D transmission, while  $\mathbf{x}_{c,n}^D$  is some vector variable associated with UE  $n$  in the downlink cellular transmission.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

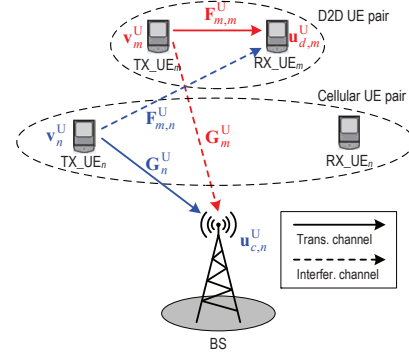
### A. System Model

As shown in Fig. 1, we consider a network that consists of one BS and a set  $\mathcal{M} = \{1, 2, \dots, M\}$  of UE pairs, where each UE pair includes a transmitter UE (TX\_UE) and a receiver UE (RX\_UE). The BS is equipped with  $N_b$  antennas, while each UE with  $N_u$  antennas. By employing the frequency division duplexing (FDD) strategy, each UE pair can communicate either in cellular mode or in D2D mode. In cellular mode, the BS adopts the decode-and-forward strategy to transmit data from one TX\_UE to its intended RX\_UE through orthogonal uplink and downlink frequency bands. In D2D mode, one TX\_UE sends data to its intended RX\_UE directly by reusing the cellular resources in both uplink and downlink. We assume the interfering broadcast channel (IBC) model [28], i.e., each transmitter generates interference to all the other receivers.

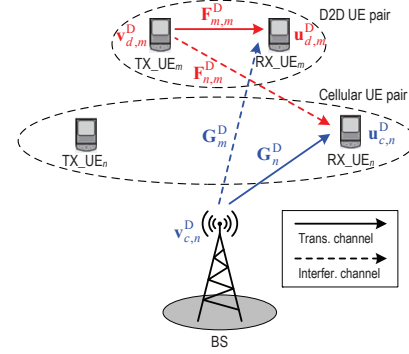
The channel and transceiver variables used in uplink and downlink are summarized in Table I. Without loss of generality, let us assume in Fig. 1 that the  $m$ th UE pair (the transmitter and receiver are denoted as TX\_UE<sub>m</sub> and RX\_UE<sub>m</sub>) selects the D2D mode, while the  $n$ th UE pair (the transmitter and receiver are denoted as TX\_UE<sub>n</sub> and RX\_UE<sub>n</sub>) selects the cellular mode. Therefore, the system depicted in Fig. 1 works as follows.

The uplink transmissions are shown in Fig. 1(a). In the uplink band, TX\_UE<sub>m</sub> sends data to RX\_UE<sub>m</sub> directly through the uplink D2D channel  $\mathbf{F}_{m,m}^U \in \mathbb{C}^{N_u \times N_u}$ , using the TX beamformer  $\mathbf{v}_m^U \in \mathbb{C}^{N_u \times 1}$ , while TX\_UE<sub>n</sub> sends data to BS via the uplink cellular channel  $\mathbf{G}_n^U \in \mathbb{C}^{N_b \times N_u}$ , using the TX beamformer  $\mathbf{v}_n^U \in \mathbb{C}^{N_u \times 1}$ . Then, RX\_UE<sub>m</sub> decodes the data from TX\_UE<sub>m</sub>, using the RX beamformer  $\mathbf{u}_{d,m}^U \in \mathbb{C}^{N_u \times 1}$ , while BS decodes the data from TX\_UE<sub>n</sub>, using the RX beamformer  $\mathbf{u}_{c,n}^U \in \mathbb{C}^{N_b \times 1}$ .

The downlink transmissions are shown in Fig. 1(b). In the downlink band, TX\_UE<sub>m</sub> sends data to RX\_UE<sub>m</sub> directly



(a) Transmissions in uplink band



(b) Transmissions in downlink band

Fig. 1. D2D communications underlying multi-user MIMO cellular networks. The solid lines denote the intended transmit channels, and the dashed lines denote the interfering channels.

through the downlink D2D channel  $\mathbf{F}_{m,m}^D \in \mathbb{C}^{N_u \times N_u}$ , using the TX beamformer  $\mathbf{v}_{d,m}^D \in \mathbb{C}^{N_u \times 1}$ , while BS forwards the data from TX\_UE<sub>n</sub> (decoded in uplink) to RX\_UE<sub>n</sub> via the downlink cellular channel  $\mathbf{G}_n^D \in \mathbb{C}^{N_u \times N_b}$ , using the TX beamformer  $\mathbf{v}_{c,n}^D \in \mathbb{C}^{N_b \times 1}$ . Then, RX\_UE<sub>m</sub> decodes the data from TX\_UE<sub>m</sub>, using the RX beamformer  $\mathbf{u}_{d,m}^D \in \mathbb{C}^{N_u \times 1}$ , while RX\_UE<sub>n</sub> decodes the data from BS, using the RX beamformer  $\mathbf{u}_{c,n}^D \in \mathbb{C}^{N_u \times 1}$ .

Compared with many current schemes, our setting endows more flexibilities to the D2D management, thereby making it possible to fully utilize the network resources and maximize the system throughput. On the other hand, this flexible setting also introduces extra interferences, and then complicates the network management. To achieve the potential benefits of this flexible setting, we should appropriately select the transmission mode for each UE pair and carefully design the associated transceivers to mitigate the interference.

### B. Problem Statement

We first define the D2D set  $\mathcal{D}$  and the cellular set  $\mathcal{D}^\perp$ , with  $\mathcal{D} \cap \mathcal{D}^\perp = \emptyset$  and  $\mathcal{D} \cup \mathcal{D}^\perp = \mathcal{M}$ . If UE pair  $m$  works in D2D mode, we have  $m \in \mathcal{D}$  and  $\mathbf{u}_{c,m}^U = \mathbf{0}$ ,  $\mathbf{v}_{c,m}^D = \mathbf{0}$ ,  $\mathbf{u}_{c,m}^D = \mathbf{0}$ ; otherwise, we have  $m \in \mathcal{D}^\perp$  and  $\mathbf{u}_{d,m}^U = \mathbf{0}$ ,  $\mathbf{v}_{d,m}^D = \mathbf{0}$ ,  $\mathbf{u}_{d,m}^D = \mathbf{0}$ . Next, we formulate the signal-to-interference-plus-noise-ratio (SINR) and rate terms of each UE pair. As shown in Fig. 1, the uplink D2D transmission between TX\_UE<sub>m</sub> and

RX\_UE<sub>m</sub> via channel  $\mathbf{F}_{m,m}^U$  is interfered by all the other TX\_UE<sub>n</sub> via channel  $\mathbf{F}_{m,n}^U$ ,  $\forall n \in \mathcal{M}$ ,  $n \neq m$ , while the downlink D2D transmission between TX\_UE<sub>m</sub> and RX\_UE<sub>m</sub> via channel  $\mathbf{F}_{m,m}^D$  is interfered by all the other D2D transmitter TX\_UE<sub>n</sub> via  $\mathbf{F}_{m,n}^D$ ,  $\forall n \in \mathcal{D}$ ,  $n \neq m$ , and by the BS via channel  $\mathbf{G}_m^D$ . Thus, in D2D mode, the uplink and downlink SINRs of UE pair  $m$  are computed as

$$\text{SINR}_{d,m}^U = \frac{|(\mathbf{u}_{d,m}^U)^\dagger \mathbf{F}_{m,m}^U \mathbf{v}_m^U|^2}{\sigma^2 \|\mathbf{u}_{d,m}^U\|_2^2 + \sum_{n \in \mathcal{M}, n \neq m} |(\mathbf{u}_{d,m}^U)^\dagger \mathbf{F}_{m,n}^U \mathbf{v}_n^U|^2}, \quad (1a)$$

$$\text{SINR}_{d,m}^D = \frac{|(\mathbf{u}_{d,m}^D)^\dagger \mathbf{F}_{m,m}^D \mathbf{v}_{d,m}^D|^2}{\left( \sigma^2 \|\mathbf{u}_{d,m}^D\|_2^2 + \sum_{n \in \mathcal{D}, n \neq m} |(\mathbf{u}_{d,m}^D)^\dagger \mathbf{F}_{m,n}^D \mathbf{v}_{d,n}^D|^2 + |(\mathbf{u}_{d,m}^D)^\dagger \mathbf{G}_m^D \mathbf{v}_{c,m}^D|^2 \right)}, \quad (1b)$$

where  $\sigma^2$  is the noise power. Let  $B^U$  and  $B^D$  denote the uplink bandwidth and the downlink bandwidth, respectively. Then, the total rate of UE pair  $m$  in D2D mode can be expressed as

$$R_{d,m} = B^U \log(1 + \text{SINR}_{d,m}^U) + B^D \log(1 + \text{SINR}_{d,m}^D). \quad (2)$$

Similarly, in cellular mode, the uplink and downlink SINRs of UE pair  $m$  are computed as

$$\text{SINR}_{c,m}^U = \frac{|(\mathbf{u}_{c,m}^U)^\dagger \mathbf{G}_m^U \mathbf{v}_m^U|^2}{\sigma^2 \|\mathbf{u}_{c,m}^U\|_2^2 + \sum_{n \in \mathcal{M}, n \neq m} |(\mathbf{u}_{c,m}^U)^\dagger \mathbf{G}_n^U \mathbf{v}_n^U|^2}, \quad (3a)$$

$$\text{SINR}_{c,m}^D = \frac{|(\mathbf{u}_{c,m}^D)^\dagger \mathbf{G}_m^D \mathbf{v}_{c,m}^D|^2}{\left( \sigma^2 \|\mathbf{u}_{c,m}^D\|_2^2 + \sum_{n \in \mathcal{D}, n \neq m} |(\mathbf{u}_{c,m}^D)^\dagger \mathbf{G}_n^D \mathbf{v}_{c,n}^D|^2 + \sum_{n \in \mathcal{D}} |(\mathbf{u}_{c,m}^D)^\dagger \mathbf{F}_{m,n}^D \mathbf{v}_{d,n}^D|^2 \right)}. \quad (3b)$$

The rate of UE pair  $m$  in cellular mode is actually determined by the smaller one of the uplink rate and the downlink rate, i.e.,

$$R_{c,m} = \min\{B^U \log(1 + \text{SINR}_{c,m}^U), B^D \log(1 + \text{SINR}_{c,m}^D)\}. \quad (4)$$

Therefore, the sum-rate maximization problem based on joint mode selection and transceiver design is formulated as

$$\begin{aligned} (\text{PA}) : \quad & \max_{\{\mathcal{D}, \mathcal{D}^\perp, \mathbf{V}, \mathbf{U}\}} \sum_{m \in \mathcal{D}} R_{d,m} + \sum_{m \in \mathcal{D}^\perp} R_{c,m} \\ \text{s.t.} \quad & \sum_{m \in \mathcal{M}} \|\mathbf{v}_{c,m}^D\|_2^2 \leq P_B, \end{aligned} \quad (5a)$$

$$\|\mathbf{v}_m^U\|_2^2 \leq p_m^U, \|\mathbf{v}_{d,m}^D\|_2^2 \leq p_m^D, \forall m \in \mathcal{M}, \quad (5b)$$

$$\mathcal{D} \cap \mathcal{D}^\perp = \emptyset, \mathcal{D} \cup \mathcal{D}^\perp = \mathcal{M}, \quad (5c)$$

where  $\mathbf{V}$  and  $\mathbf{U}$  denote  $\{\mathbf{v}_m^U, \mathbf{v}_{d,m}^D, \mathbf{v}_{c,m}^D\}_{m \in \mathcal{M}}$  and  $\{\mathbf{u}_{d,m}^U, \mathbf{u}_{c,m}^U, \mathbf{u}_{d,m}^D, \mathbf{u}_{c,m}^D\}_{m \in \mathcal{M}}$ , respectively;  $p_m^U$  and  $p_m^D$  are the transmit power budgets of TX\_UE<sub>m</sub> in uplink and downlink,  $\forall m \in \mathcal{M}$ ;  $P_B$  is the BS power budget.

### C. Problem Reformulation

In this subsection, we reformulate (PA) to avoid optimizing  $\mathcal{D}$  and  $\mathcal{D}^\perp$  directly. First, we give some more general forms of the downlink SINRs in D2D and cellular modes. Replacing  $n \in \mathcal{D}$  and  $n \in \mathcal{D}^\perp$  in (1b) and (3b) simply by  $n \in \mathcal{M}$ , we

get

$$\overline{\text{SINR}}_{d,m}^D = \frac{|(\mathbf{u}_{d,m}^D)^\dagger \mathbf{F}_{m,m}^D \mathbf{v}_{d,m}^D|^2}{\left( \sigma^2 \|\mathbf{u}_{d,m}^D\|_2^2 + \sum_{n \in \mathcal{M}, n \neq m} |(\mathbf{u}_{d,m}^D)^\dagger \mathbf{F}_{m,n}^D \mathbf{v}_{d,n}^D|^2 + |(\mathbf{u}_{d,m}^D)^\dagger \mathbf{G}_m^D \mathbf{v}_{c,m}^D|^2 \right)}, \quad (6a)$$

$$\overline{\text{SINR}}_{c,m}^D = \frac{|(\mathbf{u}_{c,m}^D)^\dagger \mathbf{G}_m^D \mathbf{v}_{c,m}^D|^2}{\left( \sigma^2 \|\mathbf{u}_{c,m}^D\|_2^2 + \sum_{n \in \mathcal{M}, n \neq m} |(\mathbf{u}_{c,m}^D)^\dagger \mathbf{G}_n^D \mathbf{v}_{c,n}^D|^2 + \sum_{n \in \mathcal{M}} |(\mathbf{u}_{c,m}^D)^\dagger \mathbf{F}_{m,n}^D \mathbf{v}_{d,n}^D|^2 \right)}. \quad (6b)$$

The rates of UE pair  $m$  in D2D and cellular modes can then be computed as

$$\bar{R}_{d,m} = B^U \log(1 + \text{SINR}_{d,m}^U) + B^D \log(1 + \overline{\text{SINR}}_{d,m}^D), \quad (7a)$$

$$\bar{R}_{c,m} = \min\{B^U \log(1 + \text{SINR}_{c,m}^U), B^D \log(1 + \overline{\text{SINR}}_{c,m}^D)\}. \quad (7b)$$

Second, in our original formulation, the mode selection strategy is exclusive; i.e., the communication within a UE pair should be established either through BS or via D2D link. To achieve this, we further introduce the binary vector  $\mathbf{d} = [d_1, d_2, \dots, d_M]$  to indicate the transmission mode of each UE pair. If  $d_m = 1$ , UE pair  $m$  works in D2D mode; if  $d_m = 0$ , UE pair  $m$  works in cellular mode. Then, the total achievable rate of UE pair  $m$  is computed as

$$R_m = d_m \bar{R}_{d,m} + (1 - d_m) \bar{R}_{c,m}. \quad (8)$$

*Proposition 1:* (PA) can be equivalently reformulated as the following problem (PB),

$$\begin{aligned} (\text{PB}) : \quad & \max_{\{\mathbf{V}, \mathbf{U}, \mathbf{d}\}} \sum_{m \in \mathcal{M}} R_m \\ \text{s.t.} \quad & (5a) \text{ and } (5b) \text{ satisfied,} \\ & d_m \in \{0, 1\}, \forall m \in \mathcal{M}. \end{aligned} \quad (9)$$

*Proof:* See Appendix A for the details of proof.  $\square$

It has been well known that the sum-rate maximization problem in multi-user MIMO-IBC networks is NP-hard [27]. Compared with the conventional sum-rate maximization problem, (PB) is obviously more challenging. This motivates us to seek some approximate solutions with manageable complexity.

## III. ITERATIVE ALGORITHM BASED ON INEXACT WMMSE

### A. Brief Review of WMMSE

The WMMSE approach [28], [29] transforms the sum-rate maximization problem to a weighted mean square error (MSE) minimization problem, and the two problems are equivalent in the sense that they have the same global optimal solutions.

We consider a MIMO-IBC network consisting of  $M$  pairs of UEs, where each UE pair contains one TX\_UE and one RX\_UE. Denote  $\mathbf{H}_{n,m}$  as the channel between TX\_UE<sub>m</sub> and RX\_UE<sub>n</sub>. Let  $\mathbf{v}_m$  and  $\mathbf{u}_m$  denote the TX beamformer of TX\_UE<sub>m</sub> and the RX beamformer of RX\_UE<sub>m</sub>, respectively. We can express the SINR and rate of RX\_UE<sub>m</sub> as

$$\begin{cases} \text{SINR}_m = \frac{|\mathbf{u}_m^\dagger \mathbf{H}_{m,m} \mathbf{v}_m|^2}{\sigma^2 \|\mathbf{u}_m\|_2^2 + \sum_{n \neq m} |\mathbf{u}_m^\dagger \mathbf{H}_{m,n} \mathbf{v}_n|^2}, \\ R_m = \log(1 + \text{SINR}_m) \end{cases} \quad (10)$$

Utilizing the well-known relation between SINR and MSE, we have  $\max_{\{\mathbf{u}_m\}} (1 + \text{SINR}_m) = \max_{\{\mathbf{u}_m\}} e_m^{-1}$ , where  $e_m$  is the MSE at RX\_UE<sub>m</sub>, which is given by

$$e_m = |1 - \mathbf{u}_m^\dagger \mathbf{H}_{m,m} \mathbf{v}_m|^2 + \sum_{n \neq m} |\mathbf{u}_m^\dagger \mathbf{H}_{m,n} \mathbf{v}_n|^2 + \sigma^2 \|\mathbf{u}_m\|_2^2. \quad (11)$$

Note that  $\log(e_m) = \min_{\{w_m > 0\}} w_m e_m - \log(w_m) + 1$ , then the following sum-rate maximization problem (PC) and the weighted sum-MSE minimization problem (PD) are equivalent

$$\begin{aligned} \text{(PC)} : \quad & \max_{\{\mathbf{v}_m, \mathbf{u}_m\}} \sum_m R_m \\ \text{s.t.} \quad & \|\mathbf{v}_m\|_2^2 \leq P_m, \forall m, \end{aligned}$$

$$\begin{aligned} \text{(PD)} : \quad & \min_{\{\mathbf{v}_m, \mathbf{u}_m, w_m\}} \sum_m [w_m e_m - \log(w_m) + 1] \\ \text{s.t.} \quad & \|\mathbf{v}_m\|_2^2 \leq P_m, \forall m, \end{aligned}$$

where  $w_m$  is the weight of  $e_m$ , and  $P_m$  is the transmit power budget of TX\_UE<sub>m</sub>.

The weighted sum-MSE minimization problem (PD) can be solved by the block coordinate descent (BCD) technique which updates  $\{\mathbf{v}_m\}_{m=1}^M$ ,  $\{\mathbf{u}_m\}_{m=1}^M$  and  $\{w_m\}_{m=1}^M$  alternately, i.e.,

$$\mathbf{v}_m = \left( \sum_{n=1}^M w_n \mathbf{H}_{n,m}^\dagger \mathbf{u}_n \mathbf{u}_n^\dagger \mathbf{H}_{n,m} + \zeta_m \mathbf{I} \right)^{-1} \mathbf{H}_{m,m}^\dagger \mathbf{u}_m w_m, \quad (12a)$$

$$\mathbf{u}_m = \left( \sum_{n=1}^M \mathbf{H}_{m,n} \mathbf{v}_n \mathbf{v}_n^\dagger \mathbf{H}_{m,n}^\dagger + \sigma^2 \mathbf{I} \right)^{-1} \mathbf{H}_{m,m} \mathbf{v}_m, \quad (12b)$$

$$w_m = (1 - \mathbf{u}_m^\dagger \mathbf{H}_{m,m} \mathbf{v}_m)^{-1}, \quad (12c)$$

where  $\zeta_m \geq 0$  is the Lagrangian multiplier associated with  $\|\mathbf{v}_m\|_2^2 \leq P_m$ , and should be chosen carefully such that the Karush-Kuhn-Tucker (KKT) complementarity conditions [31] are satisfied;  $\mathbf{u}_m$  is the MMSE receiver such that the MSE, i.e.,  $e_m$  defined in (11), can be minimized;  $w_m$  is updated by the reciprocal of the minimum  $e_m$ .

According to *Theorem 3* in [28], the WMMSE algorithm<sup>1</sup> iteratively computes a stationary point of (PC) by repeating the steps in (12).

### B. Problem Approximation

In this subsection, we formulate an efficient approximation of (PB) for which the WMMSE approach is applicable.

To this end, we first relax the binary indicator  $d_m \in \{0, 1\}$  by  $0 \leq d_m \leq 1$ . After that, we deal with the product terms in  $R_m = d_m \bar{R}_{d,m} + (1 - d_m) \bar{R}_{c,m}$ . Specifically, we introduce some auxiliary variables  $\{s_{d,m}^U, s_{d,m}^D, s_{c,m}\}_{m \in \mathcal{M}}$ , and bound them by  $d_m B^U \log(1 + \text{SINR}_{d,m}^U) \geq (s_{d,m}^U)^2$ ,  $d_m B^D \log(1 +$

$\text{SINR}_{d,m}^D) \geq (s_{d,m}^D)^2$ ,  $(1 - d_m) B^U \log(1 + \text{SINR}_{c,m}^U) \geq s_{c,m}^2$ , and  $(1 - d_m) B^D \log(1 + \text{SINR}_{c,m}^D) \geq s_{c,m}^2$ ,  $m \in \mathcal{M}$ . Then, the term of  $R_m$  can be replaced by  $(s_{d,m}^U)^2 + (s_{d,m}^D)^2 + s_{c,m}^2$  in the objective of (PB). Considering that the ‘‘quadratic-over-linear’’ function  $h(\omega, \nu) \triangleq \frac{\omega^2}{\nu}$  is convex as  $\nu > 0$  [31], we can recast  $\varrho \nu \geq \omega^2$  as a convex constraint  $\varrho \geq \frac{\omega^2}{\nu}$  for  $\nu > 0$ . Following this idea, we get an approximation of (PB) as

$$\begin{aligned} \text{(PE)} : \quad & \max_{\{\mathbf{V}, \mathbf{U}, \mathbf{d}, \mathbf{s}\}} \sum_{m \in \mathcal{M}} (s_{d,m}^U)^2 + (s_{d,m}^D)^2 + s_{c,m}^2 \\ \text{s.t.} \quad & \text{(5a) and (5b) satisfied,} \\ & 0 \leq d_m \leq 1, \forall m \in \mathcal{M}, \quad (13a) \\ & B^U \log(1 + \text{SINR}_{d,m}^U) \geq \frac{(s_{d,m}^U)^2}{d_m + \epsilon}, \forall m \in \mathcal{M}, \quad (13b) \\ & B^D \log(1 + \text{SINR}_{d,m}^D) \geq \frac{(s_{d,m}^D)^2}{d_m + \epsilon}, \forall m \in \mathcal{M}, \quad (13c) \\ & \min \left\{ B^U \log(1 + \text{SINR}_{c,m}^U), B^D \log(1 + \text{SINR}_{c,m}^D) \right\} \geq \frac{s_{c,m}^2}{1 - d_m + \epsilon}, \\ & \forall m \in \mathcal{M}, \quad (13d) \end{aligned}$$

where  $\mathbf{s}$  denotes  $\{s_{d,m}^U, s_{d,m}^D, s_{c,m}\}_{m \in \mathcal{M}}$ , and  $\epsilon$  is a small positive number introduced to avoid the numerical problems of zero denominator.

We further apply the WMMSE reformulation to handle the  $\log(\cdot)$  terms in (13b) – (13d), thereby generating (PF). From *Theorem 1* in [28] and *Lemma 3* in [29], (PE) and (PF) are equivalent in the sense that they have the same optimal solutions. Specifically, (PF) is expressed as

$$\begin{aligned} \text{(PF)} : \quad & \min_{\{\mathbf{V}, \mathbf{U}, \mathbf{W}, \mathbf{d}, \mathbf{s}\}} \sum_{m \in \mathcal{M}} (s_{d,m}^U)^2 + (s_{d,m}^D)^2 + s_{c,m}^2 \\ \text{s.t.} \quad & \text{(5a), (5b) and (13a) satisfied,} \\ & \log(w_{d,m}^U) - w_{d,m}^U e_{d,m}^U(\mathbf{V}^U, \mathbf{u}_{d,m}^U) + 1 \geq \frac{(s_{d,m}^U)^2}{B^U(d_m + \epsilon)}, \\ & \quad \forall m \in \mathcal{M}, \quad (14a) \\ & \log(w_{d,m}^D) - w_{d,m}^D e_{d,m}^D(\mathbf{V}^D, \mathbf{u}_{d,m}^D) + 1 \geq \frac{(s_{d,m}^D)^2}{B^D(d_m + \epsilon)}, \\ & \quad \forall m \in \mathcal{M}, \quad (14b) \\ & \log(w_{c,m}^U) - w_{c,m}^U e_{c,m}^U(\mathbf{V}^U, \mathbf{u}_{c,m}^U) + 1 \geq \frac{s_{c,m}^2}{B^U(1 - d_m + \epsilon)}, \\ & \quad \forall m \in \mathcal{M}, \quad (14c) \\ & \log(w_{c,m}^D) - w_{c,m}^D e_{c,m}^D(\mathbf{V}^D, \mathbf{u}_{c,m}^D) + 1 \geq \frac{s_{c,m}^2}{B^D(1 - d_m + \epsilon)}, \\ & \quad \forall m \in \mathcal{M}, \quad (14d) \end{aligned}$$

where  $\mathbf{W}$  is the collection of weighting factors  $\{w_{d,m}^U, w_{c,m}^U, w_{d,m}^D, w_{c,m}^D\}_{m \in \mathcal{M}}$ ;  $\mathbf{V}^U$  and  $\mathbf{V}^D$  denote the uplink TX beamformers  $\{\mathbf{v}_m^U\}_{m \in \mathcal{M}}$  and the downlink TX beamformers  $\{\mathbf{v}_m^D\}_{m \in \mathcal{M}}$ ;  $e_{d,m}^U(\mathbf{V}^U, \mathbf{u}_{d,m}^U)$  and  $e_{d,m}^D(\mathbf{V}^D, \mathbf{u}_{d,m}^D)$  are the uplink and downlink D2D MSE values of UE pair  $m$ ;  $e_{c,m}^U(\mathbf{V}^U, \mathbf{u}_{c,m}^U)$  and  $e_{c,m}^D(\mathbf{V}^D, \mathbf{u}_{c,m}^D)$  are the uplink and downlink cellular MSE values of UE pair  $m$ . Similar as (11), these

<sup>1</sup>There are key differences between (PB) and the conventional sum-rate maximization problem (PC). First, we consider the underlay D2D communication scenario which involves both D2D uplink/downlink transmissions and cellular uplink/downlink transmissions, while (PC) only considers the cellular downlink transmission. Second, the objective of (PB) adopts a more complicated form involving the product terms of  $d_m \bar{R}_{c,m}$  and  $d_m \bar{R}_{d,m}$ . Therefore, the WMMSE algorithm cannot be directly applied to solve (PB).

MSE values are defined in (15),

$$e_{d,m}^U(\mathbf{V}^U, \mathbf{u}_{d,m}^U) = \sigma^2 \|\mathbf{u}_{d,m}^U\|_2^2 + |1 - (\mathbf{u}_{d,m}^U)^\dagger \mathbf{F}_{m,m}^U \mathbf{v}_m^U|^2 + \sum_{\substack{n \in \mathcal{M} \\ n \neq m}} |(\mathbf{u}_{d,m}^U)^\dagger \mathbf{F}_{m,n}^U \mathbf{v}_n^U|^2, \quad (15a)$$

$$e_{d,m}^D(\mathbf{V}^D, \mathbf{u}_{d,m}^D) = \sigma^2 \|\mathbf{u}_{d,m}^D\|_2^2 + |1 - (\mathbf{u}_{d,m}^D)^\dagger \mathbf{F}_{m,m}^D \mathbf{v}_{d,m}^D|^2 + \sum_{\substack{n \in \mathcal{M} \\ n \neq m}} |(\mathbf{u}_{d,m}^D)^\dagger \mathbf{F}_{m,n}^D \mathbf{v}_{d,n}^D|^2 + \sum_{n \in \mathcal{M}} |(\mathbf{u}_{d,m}^D)^\dagger \mathbf{G}_m^D \mathbf{v}_{c,n}^D|^2, \quad (15b)$$

$$e_{c,m}^U(\mathbf{V}^U, \mathbf{u}_{c,m}^U) = \sigma^2 \|\mathbf{u}_{c,m}^U\|_2^2 + |1 - (\mathbf{u}_{c,m}^U)^\dagger \mathbf{G}_m^U \mathbf{v}_m^U|^2 + \sum_{\substack{n \in \mathcal{M} \\ n \neq m}} |(\mathbf{u}_{c,m}^U)^\dagger \mathbf{G}_n^U \mathbf{v}_n^U|^2, \quad (15c)$$

$$e_{c,m}^D(\mathbf{V}^D, \mathbf{u}_{c,m}^D) = \sigma^2 \|\mathbf{u}_{c,m}^D\|_2^2 + |1 - (\mathbf{u}_{c,m}^D)^\dagger \mathbf{G}_m^D \mathbf{v}_{c,m}^D|^2 + \sum_{n \in \mathcal{M}} |(\mathbf{u}_{c,m}^D)^\dagger \mathbf{F}_{m,n}^D \mathbf{v}_{d,n}^D|^2 + \sum_{\substack{n \in \mathcal{M} \\ n \neq m}} |(\mathbf{u}_{c,m}^D)^\dagger \mathbf{G}_m^D \mathbf{v}_{c,n}^D|^2. \quad (15d)$$

### C. Iterative Algorithm Based on Inexact WMMSE

Due to its non-convex objective, we handle (PF) by iteratively solving its sequential convex approximations (SCA) based on the idea of difference-of-convex (DC) programming [32]. Specifically, in each iteration, we perform the first-order approximation to  $[(s_{d,m}^U)^2 + (s_{d,m}^D)^2 + s_{c,m}^2]$ , and then solve the resultant problem

$$\begin{aligned} \text{(PG)} : \min_{\{\mathbf{V}, \mathbf{U}, \mathbf{W}, \mathbf{d}, \mathbf{s}\}} & -2 \sum_{m \in \mathcal{M}} \hat{s}_{d,m}^U s_{d,m}^U + \hat{s}_{d,m}^D s_{d,m}^D + \hat{s}_{c,m} s_{c,m} \\ \text{s.t.} & \text{ (5a), (5b), (13a) and (14a) – (14d) satisfied,} \end{aligned}$$

with  $\hat{s}_{d,m}^U$ ,  $\hat{s}_{d,m}^D$  and  $\hat{s}_{c,m}$  being the iterates of  $s_{d,m}^U$ ,  $s_{d,m}^D$  and  $s_{c,m}$  in the previous iteration.

Following the WMMSE framework of (12), we divide the variables  $\{\mathbf{V}, \mathbf{U}, \mathbf{W}, \mathbf{d}, \mathbf{s}\}$  of (PG) into three blocks, i.e.,  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$ ,  $\mathbf{U}$ , and  $\mathbf{W}$ , and then apply the BCD method to solve it iteratively. First, fixing  $\mathbf{U}$  and  $\mathbf{W}$ , we solve the following convex problem to update  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$ ,

$$\begin{aligned} \text{(PH)} : \min_{\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}} & -2 \sum_{m \in \mathcal{M}} \hat{s}_{d,m}^U s_{d,m}^U + \hat{s}_{d,m}^D s_{d,m}^D + \hat{s}_{c,m} s_{c,m} \\ \text{s.t.} & \text{ (5a), (5b), (13a), and (14a) – (14d) satisfied.} \end{aligned}$$

Problem (PH) is convex can be optimally solved by CVX [33]. Define the mapping function  $\Omega(\mathbf{U}, \mathbf{W}; \hat{\mathbf{s}})$ , where  $\hat{\mathbf{s}}$  denotes  $\{\hat{s}_{d,m}^U, \hat{s}_{d,m}^D, \hat{s}_{c,m}\}_{m \in \mathcal{M}}$ , such that every element in the range of the map is an optimal solution to (PH). Then, the update of  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$  can be simply written as

$$\{\mathbf{V}, \mathbf{d}, \mathbf{s}\} \in \Omega(\mathbf{U}, \mathbf{W}; \hat{\mathbf{s}}). \quad (16)$$

Next, fixing  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$  and  $\mathbf{W}$ , we update  $\mathbf{U}$  by solving the following problem,

$$\begin{aligned} \text{(PI)} : \min_{\mathbf{U}} & -2 \sum_{m \in \mathcal{M}} \hat{s}_{d,m}^U s_{d,m}^U + \hat{s}_{d,m}^D s_{d,m}^D + \hat{s}_{c,m} s_{c,m} \\ \text{s.t.} & \text{ (14a) – (14d) satisfied.} \end{aligned}$$

In this subproblem,  $\mathbf{U}$  only appears in the MSE terms in (14), i.e.,  $\{e_{d,m}^U, e_{d,m}^D, e_{c,m}^U, e_{c,m}^D\}_{m \in \mathcal{M}}$ . Since we consider a sum-rate maximization problem here,  $\mathbf{U}$  is selected to minimize the MSE terms so that the rates of the UE pairs can be maximally improved in this step. In addition, updating  $\mathbf{U}$  as the optimal MMSE receiver also maximizes the right-hand side of (14), thereby maximally enlarging the searching space of  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$  in the next iteration. This helps improve the system throughput quickly. Therefore, the optimal solution of  $\mathbf{U}$  is given by

$$\mathbf{u}_{d,m}^U = \left[ \sum_{n \in \mathcal{M}} \mathbf{F}_{m,n}^U \mathbf{v}_n^U (\mathbf{v}_n^U)^\dagger (\mathbf{F}_{m,n}^U)^\dagger + \sigma^2 \mathbf{I} \right]^{-1} \mathbf{F}_{m,m}^U \mathbf{v}_m^U \quad (17a)$$

$$\mathbf{u}_{c,m}^U = \left[ \sum_{n \in \mathcal{M}} \mathbf{G}_n^U \mathbf{v}_n^U (\mathbf{v}_n^U)^\dagger (\mathbf{G}_n^U)^\dagger + \sigma^2 \mathbf{I} \right]^{-1} \mathbf{G}_m^U \mathbf{v}_m^U \quad (17b)$$

$$\mathbf{u}_{d,m}^D = \left[ \sum_{n \in \mathcal{M}} \left[ \mathbf{F}_{m,n}^D \mathbf{v}_{d,n}^D (\mathbf{v}_{d,n}^D)^\dagger (\mathbf{F}_{m,n}^D)^\dagger + \mathbf{G}_m^D \mathbf{v}_{c,n}^D (\mathbf{v}_{c,n}^D)^\dagger (\mathbf{G}_m^D)^\dagger \right] + \sigma^2 \mathbf{I} \right]^{-1} \mathbf{F}_{m,m}^D \mathbf{v}_{d,m}^D \quad (17c)$$

$$\mathbf{u}_{c,m}^D = \left[ \sum_{n \in \mathcal{M}} \left[ \mathbf{F}_{m,n}^D \mathbf{v}_{d,n}^D (\mathbf{v}_{d,n}^D)^\dagger (\mathbf{F}_{m,n}^D)^\dagger + \mathbf{G}_m^D \mathbf{v}_{c,n}^D (\mathbf{v}_{c,n}^D)^\dagger (\mathbf{G}_m^D)^\dagger \right] + \sigma^2 \mathbf{I} \right]^{-1} \mathbf{G}_m^D \mathbf{v}_{c,m}^D. \quad (17d)$$

For the sake of brevity, we simplify the formulas of updating  $\mathbf{U}$  as

$$\mathbf{U} = \Gamma(\mathbf{V}). \quad (18)$$

Lastly, fixing  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$  and  $\mathbf{U}$ , we update  $\mathbf{W}$  by solving the following problem,

$$\begin{aligned} \text{(PJ)} : \min_{\mathbf{W}} & -2 \sum_{m \in \mathcal{M}} \hat{s}_{d,m}^U s_{d,m}^U + \hat{s}_{d,m}^D s_{d,m}^D + \hat{s}_{c,m} s_{c,m} \\ \text{s.t.} & \text{ (14a) – (14d) satisfied.} \end{aligned}$$

Similarly,  $\mathbf{W}$  is chosen to maximize the right-hand side of (14), which is actually the reciprocal of the minimum MSE, i.e.,

$$w_{d,m}^U = [1 - (\mathbf{u}_{d,m}^U)^\dagger \mathbf{F}_{m,m}^D \mathbf{v}_m^U]^{-1} \quad (19a)$$

$$w_{c,m}^U = [1 - (\mathbf{u}_{c,m}^U)^\dagger \mathbf{G}_m^U \mathbf{v}_m^U]^{-1} \quad (19b)$$

$$w_{d,m}^D = [1 - (\mathbf{u}_{d,m}^D)^\dagger \mathbf{F}_{m,m}^D \mathbf{v}_{d,m}^D]^{-1} \quad (19c)$$

$$w_{c,m}^D = [1 - (\mathbf{u}_{c,m}^D)^\dagger \mathbf{G}_m^D \mathbf{v}_{c,m}^D]^{-1} \quad (19d)$$

where  $\mathbf{u}_{d,m}^U$ ,  $\mathbf{u}_{c,m}^U$ ,  $\mathbf{u}_{d,m}^D$  and  $\mathbf{u}_{c,m}^D$  are computed as in (17). Again, for the sake of brevity, we simplify the formulas of updating  $\mathbf{W}$  as

$$\mathbf{W} = \Upsilon(\mathbf{V}, \mathbf{U}). \quad (20)$$

By repeating the steps of (16), (18) and (20), the WMMSE algorithm iteratively solves (PG) with some given  $\hat{\mathbf{s}}$ . Then, we update  $\hat{\mathbf{s}} = \mathbf{s}$ , and repeat the above WMMSE loop. Finally, a two-layer (including the outer DC programming layer and the inner WMMSE layer) algorithm was developed to iteratively solve (PF) or the equivalent (PE) with stationary convergence guarantee, which is referred to as the iterative algorithm based on exact WMMSE (IA-EWMMSE) in this paper, since the update of  $\hat{\mathbf{s}}$  is based on the exact solution of the inner WMMSE problem (PG). We summarize IA-EWMMSE in Table II

TABLE II  
SUMMARY OF IA-EWMMSE

1.	Initialize $\{\mathbf{U}, \mathbf{W}, \hat{\mathbf{s}}\}$ ;
2.	<b>Repeat</b> (outer DC programming layer)
3.	<b>Repeat</b> (inner WMMSE layer)
4.	$\{\mathbf{V}, \mathbf{d}, \mathbf{s}\} \in \Omega(\mathbf{U}, \mathbf{W}; \hat{\mathbf{s}})$ ;
5.	$\mathbf{U} = \Gamma(\mathbf{V})$ ;
6.	$\mathbf{W} = \Upsilon(\mathbf{V}, \mathbf{U})$ ;
7.	<b>Until</b> some stopping criterion is satisfied;
8.	$\hat{\mathbf{s}} = \mathbf{s}$ ;
9.	<b>Until</b> some stopping criterion is satisfied;

TABLE III  
SUMMARY OF IA-IWMMSE

1.	Initialize $\{\mathbf{U}, \mathbf{W}, \hat{\mathbf{s}}\}$ ;
2.	<b>Repeat</b>
3.	$\{\mathbf{V}, \mathbf{d}, \mathbf{s}\} \in \Omega(\mathbf{U}, \mathbf{W}; \hat{\mathbf{s}})$ ;
4.	$\mathbf{U} = \Gamma(\mathbf{V})$ ;
5.	$\mathbf{W} = \Upsilon(\mathbf{V}, \mathbf{U})$ ;
6.	$\hat{\mathbf{s}} = \mathbf{s}$ ;
7.	<b>Until</b> some stopping criterion is satisfied;

TABLE IV  
SUMMARY OF THE INTRODUCED AUXILIARY VARIABLES

Variable Name	Description
$\mathbf{T}_d^U(m)$	$\mathbf{T}_d^U(m) = \mathbf{V}^U$ ; $\mathbf{T}_d^U(m) = \{\mathbf{t}_{d,n}^U(m)\}_{n \in \mathcal{M}}$ , $\mathbf{V}^U = \{\mathbf{v}_n^U\}_{n \in \mathcal{M}}$ .
$\mathbf{T}_c^U(m)$	$\mathbf{T}_c^U(m) = \mathbf{V}^U$ ; $\mathbf{T}_c^U(m) = \{\mathbf{t}_{c,n}^U(m)\}_{n \in \mathcal{M}}$ , $\mathbf{V}^U = \{\mathbf{v}_n^U\}_{n \in \mathcal{M}}$ .
$\mathbf{T}_d^D(m)$	$\mathbf{T}_d^D(m) = \mathbf{V}^D$ ; $\mathbf{T}_d^D(m) = \{\mathbf{t}_{d,n}^D(m), \mathbf{t}_{d,c,n}^D(m)\}_{n \in \mathcal{M}}$ , $\mathbf{V}^D = \{\mathbf{v}_{d,n}^D, \mathbf{v}_{c,n}^D\}_{n \in \mathcal{M}}$ .
$\mathbf{T}_c^D(m)$	$\mathbf{T}_c^D(m) = \mathbf{V}^D$ ; $\mathbf{T}_c^D(m) = \{\mathbf{t}_{c,n}^D(m), \mathbf{t}_{c,d,n}^D(m)\}_{n \in \mathcal{M}}$ , $\mathbf{V}^D = \{\mathbf{v}_{d,n}^D, \mathbf{v}_{c,n}^D\}_{n \in \mathcal{M}}$ .
$x_{d,m}^U, x_{c,m}^U, x_{d,m}^D, x_{c,m}^D$ $y_{d,m}^U, y_{d,m}^D, y_{c,m}$	$B^U x_{d,m}^U = y_{d,m}^U, B^U x_{c,m}^U = y_{c,m}$ , $B^D x_{d,m}^D = y_{d,m}^D, B^D x_{c,m}^D = y_{c,m}$ .
$z_{d,m}^U, z_{d,m}^D, z_{c,m}$	$z_{d,m}^U = d_m + \epsilon, z_{d,m}^D = d_m + \epsilon$ , $z_{c,m} = 1 - d_m + \epsilon$ .

(more details can be found in our previous conference paper [1]), where the stopping criterion is satisfied as the difference between the iterates of two adjacent iterations falls below some predefined threshold (same for Tables III and IV).

However, it may be difficult to implement IA-EWMMSE in practice due to the complicated two-layer loops. To simplify the algorithm structure, it is natural to consider designing some algorithm, which updates  $\hat{\mathbf{s}}$  based on the inexact solution of (PG) while with convergence guarantee. Motivated by this idea, we propose an algorithm referred to as the iterative algorithm based on inexact WMMSE (IA-IWMMSE). The main steps of IA-IWMMSE are summarized in Table III, where (16), (18), and (20) are executed only once in each iteration to update  $\hat{\mathbf{s}}$ . Interestingly, although  $\hat{\mathbf{s}}$  is updated in an inexact manner, we still have the following convergence result, i.e., Proposition 2, for IA-IWMMSE<sup>2</sup>. As shown later in the simulation results of Fig. 4, in addition to having the simpler single-layer structure, IA-IWMMSE accelerates the convergence, and thereby has lower total complexity than IA-EWMMSE.

*Proposition 2:* Every accumulation point of the iterates generated by IA-IWMMSE in Table III is a stationary solution of (PE).

*Proof:* See Appendix B for the details of proof.  $\square$

Before closing this section, we summarize our methodology for problem (PA). First, to avoid optimizing the set variables  $\mathcal{D}$  and  $\mathcal{D}^\perp$  directly, we introduce the binary indicators  $\{d_m\}_{m \in \mathcal{M}}$ , and obtain the equivalent problem (PB). Then, we relax the binary indicators and approximate (PB) by (PE). Since (PE) is still NP-hard, we next perform the WMMSE reformulation to get (PF), which is equivalent to (PE) in the sense that they have the same optimal solutions. However, the non-convex objective of (PF) prevents us from applying the

WMMSE algorithm. To tackle this, we iteratively solve its sequential convex approximations, i.e., (PG), by performing the first-order approximation. In each iteration, we further use BCD to divide (PG) into three convex subproblems, i.e., (PH), (PI), and (PJ). Finally, we get a stationary solution of (PF), which is also a stationary solution of (PE). In summary, the original problem (PA) is challenging to solve, and even its approximation (PE) is still NP-hard. As a compromise, we pursue a stationary solution of (PE) with manageable complexity.

#### IV. DISTRIBUTED ADMM ALGORITHM FOR (PH)

Except for updating  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$ , IA-IWMMSE can be executed distributively. Unfortunately, updating  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$  in the centralized way, e.g., by the interior-point (IP) method [31], requires a per-iteration complexity of  $\mathcal{O}(M^3(N_u + N_d)^3)$ , which dominates the total computational cost. Obviously, a fully distributed algorithm is preferred in this circumstance, which may compute  $\{\mathbf{v}_m^U, \mathbf{v}_{d,m}^D, \mathbf{v}_{c,m}^D\}_{m \in \mathcal{M}}$ ,  $\{d_m\}_{m \in \mathcal{M}}$  and  $\{s_{d,m}^U, s_{d,m}^D, s_{c,m}\}_{m \in \mathcal{M}}$  independently. To this end, we further recast (PH) so that it fits into the ADMM framework [30], and then devise a distributed algorithm to solve it, with a simple closed-form solution in each step.

##### A. ADMM Reformulation of (PH)

To decouple the variables in (PH), we first introduce three groups of auxiliary variables, which are summarized in Table IV.

1) We introduce  $M$  copies of the TX beamformers, i.e.,  $\{\mathbf{T}_d^U(m), \mathbf{T}_c^U(m), \mathbf{T}_d^D(m), \mathbf{T}_c^D(m)\}_{m \in \mathcal{M}}$ , where  $\mathbf{T}_d^U(m) = \mathbf{T}_c^U(m) = \mathbf{V}^U$ , and  $\mathbf{T}_d^D(m) = \mathbf{T}_c^D(m) = \mathbf{V}^D$ ,  $\forall m \in \mathcal{M}$ . Then, we replace the  $\mathbf{V}^U$  in (14a) and (14c) by  $\mathbf{T}_d^U(m)$  and  $\mathbf{T}_c^U(m)$ , respectively, to decouple the  $2M$  uplink constraints of  $\mathbf{V}^U$ . Similarly, we replace the  $\mathbf{V}^D$  in (14b) and (14d) by  $\mathbf{T}_d^D(m)$  and  $\mathbf{T}_c^D(m)$ , respectively, to decouple the  $2M$  downlink constraints of  $\mathbf{V}^D$ .

<sup>2</sup>After IA-IWMMSE converges, we quantize  $\{d_m\}_{m \in \mathcal{M}}$  to binary values to identify the cellular UEs and the D2D UEs. Then, we optimize the transceivers by solving two independent (i.e., uplink and downlink) sum-rate problems with known transmission modes. Since these problems have been intensively studied [28], we omit the details due to the space limitation.

2) We introduce two series of variables:  $\{x_{d,m}^U, x_{c,m}^U, x_{d,m}^D, x_{c,m}^D\}_{m \in \mathcal{M}}$  and  $\{y_{d,m}^U, y_{d,m}^D, y_{c,m}^U, y_{c,m}^D\}_{m \in \mathcal{M}}$ , where  $B^U x_{d,m}^U = y_{d,m}^U$ ,  $B^U x_{c,m}^U = y_{c,m}^U$ ,  $B^D x_{d,m}^D = y_{d,m}^D$ , and  $B^D x_{c,m}^D = y_{c,m}^D$ ,  $\forall m \in \mathcal{M}$ .

3) We also introduce a series of indicators:  $\{z_{d,m}^U, z_{d,m}^D, z_{c,m}^U, z_{c,m}^D\}_{m \in \mathcal{M}}$ , where  $z_{d,m}^U = d_m + \epsilon$ ,  $z_{d,m}^D = d_m + \epsilon$ , and  $z_{c,m} = 1 - d_m + \epsilon$ ,  $\forall m \in \mathcal{M}$ .

Inserting these variables into (PH), we equivalently recast it as

$$\begin{aligned}
 (\text{PK}) : \quad & \min_{\{\mathbf{V}, \mathbf{d}, \mathbf{s}, \mathbf{T}, \mathbf{x}, \mathbf{y}, \mathbf{z}\}} -2 \sum_{m \in \mathcal{M}} \left( \hat{s}_{d,m}^U s_{d,m}^U + \hat{s}_{d,m}^D s_{d,m}^D \right) \\
 & \text{s.t.} \quad (5a), (5b) \text{ and } (13a) \text{ satisfied,} \\
 & \log(w_{d,m}^U) - w_{d,m}^U e_{d,m}^U(\mathbf{T}_d^U(m), \mathbf{u}_{d,m}^U) + 1 \geq x_{d,m}^U, \\
 & \quad \quad \quad \forall m \in \mathcal{M}, \quad (21a) \\
 & \log(w_{d,m}^D) - w_{d,m}^D e_{d,m}^D(\mathbf{T}_d^D(m), \mathbf{u}_{d,m}^D) + 1 \geq x_{d,m}^D, \\
 & \quad \quad \quad \forall m \in \mathcal{M}, \quad (21b) \\
 & \log(w_{c,m}^U) - w_{c,m}^U e_{c,m}^U(\mathbf{T}_c^U(m), \mathbf{u}_{c,m}^U) + 1 \geq x_{c,m}^U, \\
 & \quad \quad \quad \forall m \in \mathcal{M}, \quad (21c) \\
 & \log(w_{c,m}^D) - w_{c,m}^D e_{c,m}^D(\mathbf{T}_c^D(m), \mathbf{u}_{c,m}^D) + 1 \geq x_{c,m}^D, \\
 & \quad \quad \quad \forall m \in \mathcal{M}, \quad (21d) \\
 & y_{d,m}^U \geq \frac{(s_{d,m}^U)^2}{z_{d,m}^U}, \quad y_{d,m}^D \geq \frac{(s_{d,m}^D)^2}{z_{d,m}^D}, \quad y_{c,m} \geq \frac{s_{c,m}^2}{z_{c,m}}, \\
 & \quad \quad \quad \forall m \in \mathcal{M}, \quad (21e) \\
 & z_{d,m}^U \geq \epsilon, \quad z_{d,m}^D \geq \epsilon, \quad z_{c,m} \geq \epsilon, \quad \forall m \in \mathcal{M}, \quad (21f) \\
 & B^U x_{d,m}^U = y_{d,m}^U, \quad B^D x_{d,m}^D = y_{d,m}^D, \quad \forall m \in \mathcal{M}, \quad (21g) \\
 & B^U x_{c,m}^U = y_{c,m}^U, \quad B^D x_{c,m}^D = y_{c,m}^D, \quad \forall m \in \mathcal{M}, \quad (21h) \\
 & z_{d,m}^U = d_m + \epsilon, \quad z_{d,m}^D = d_m + \epsilon, \quad z_{c,m} = 1 - d_m + \epsilon, \\
 & \quad \quad \quad \forall m \in \mathcal{M}, \quad (21i) \\
 & \mathbf{t}_{d,n}^U(m) = \mathbf{v}_n^U, \quad \mathbf{t}_{c,n}^U(m) = \mathbf{v}_n^U, \quad \forall m, n \in \mathcal{M}, \quad (21j) \\
 & \mathbf{t}_{d,n}^D(m) = \mathbf{t}_{c,n}^D(m) = \mathbf{v}_{d,n}^D, \quad \forall m, n \in \mathcal{M}, \quad (21k) \\
 & \mathbf{t}_{c,n}^D(m) = \mathbf{t}_{c,n}^D(m) = \mathbf{v}_{c,n}^D, \quad \forall m, n \in \mathcal{M}. \quad (21l)
 \end{aligned}$$

The partial *augmented Lagrangian function* [30] of (PK) is defined as (22) at the top of next page, where  $c > 0$  is the penalty parameter;  $\Phi, \psi$ , and  $\theta$  are the associated Lagrangian multipliers. Dividing the variables  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}, \mathbf{T}, \mathbf{x}, \mathbf{y}, \mathbf{z}\}$  into two blocks of  $\{\mathbf{V}, \mathbf{s}, \mathbf{y}, \mathbf{z}\}$  and  $\{\mathbf{T}, \mathbf{d}, \mathbf{x}\}$ , we perform the two-block ADMM framework shown in Table V to iteratively solve (PK), with global convergence guarantee [30].

*Remark 1:* In ADMM framework, (PK) is divided into two simple subproblems of  $\{\mathbf{V}, \mathbf{s}, \mathbf{y}, \mathbf{z}\}$  and  $\{\mathbf{T}, \mathbf{d}, \mathbf{x}\}$ . Moreover, due to their separable structures, these two subproblems can be further separated into smaller problems. More interestingly, all these smaller problems have closed-form solutions. Finally, (PK), or the equivalent (PH), can be solved distributively and efficiently. We will elaborate more on this in the following subsections.

### B. Update $\{\mathbf{V}, \mathbf{s}, \mathbf{y}, \mathbf{z}\}$

The problem of  $\{\mathbf{V}, \mathbf{s}, \mathbf{y}, \mathbf{z}\}$  is totally separable among  $\mathbf{v}_m^U$ ,  $\mathbf{v}_m^D$ ,  $\{s_{d,m}^U, y_{d,m}^U, z_{d,m}^U\}$ ,  $\{s_{d,m}^D, y_{d,m}^D, z_{d,m}^D\}$ ,  $\{s_{c,m}, y_{c,m}, z_{c,m}\}$

TABLE V  
TWO-BLOCK ADMM FRAMEWORK FOR (PK)

1.	<b>Repeat</b>
2.	Update $\{\mathbf{V}, \mathbf{s}, \mathbf{y}, \mathbf{z}\}$ with other variables fixed: $\{\mathbf{V}, \mathbf{s}, \mathbf{y}, \mathbf{z}\} = \underset{\text{s.t.}}{\text{argmin}}_{\{\mathbf{V}, \mathbf{s}, \mathbf{y}, \mathbf{z}\}} \mathcal{L}_c(\cdot)$ (5a), (5b), (21e) and (21f),
3.	Update $\{\mathbf{T}, \mathbf{d}, \mathbf{x}\}$ with other variables fixed: $\{\mathbf{T}, \mathbf{d}, \mathbf{x}\} = \underset{\text{s.t.}}{\text{argmin}}_{\{\mathbf{T}, \mathbf{d}, \mathbf{x}\}} \mathcal{L}_c(\cdot)$ (13a) and (21a) – (21d),
4.	Update Lagrangian variables $\{\Phi, \psi, \theta\}$ [30].
5.	<b>Until</b> some stopping criterion is satisfied

$z_{c,m}\}$  for  $m \in \mathcal{M}$ , and  $\{\mathbf{v}_{c,m}^D\}_{m \in \mathcal{M}}$ . Therefore, we can update the  $(5M + 1)$  subproblems independently and in parallel.

1) *Update  $\mathbf{v}_m^U$ ,  $\mathbf{v}_m^D$ , and  $\{\mathbf{v}_{c,m}^D\}_{m \in \mathcal{M}}$ :* The problem of  $\mathbf{v}_m^U$  is expressed as

$$\begin{aligned}
 (\text{PL}) : \quad & \min_{\mathbf{v}_m^U} \frac{c}{2} \sum_{n \in \mathcal{M}} \left( \|\mathbf{v}_m^U - \mathbf{t}_{d,m}^U(n) - \frac{1}{c} \phi_{d,m}^U(n)\|_2^2 \right. \\
 & \quad \left. + \|\mathbf{v}_m^U - \mathbf{t}_{c,m}^U(n) - \frac{1}{c} \phi_{c,m}^U(n)\|_2^2 \right) \\
 & \text{s.t.} \quad \|\mathbf{v}_m^U\|_2^2 \leq p_m^U, \quad (23)
 \end{aligned}$$

which is solved as

$$\begin{cases} \mathbf{v}_m^U = \frac{\sum_{n \in \mathcal{M}} [\phi_{d,m}^U(n) + \phi_{c,m}^U(n) + c(\mathbf{t}_{d,m}^U(n) + \mathbf{t}_{c,m}^U(n))]}{2(cM + \alpha_m^U)} \\ \alpha_m^U = \left[ \frac{\|\sum_{n \in \mathcal{M}} [\phi_{d,m}^U(n) + \phi_{c,m}^U(n) + c(\mathbf{t}_{d,m}^U(n) + \mathbf{t}_{c,m}^U(n))]\|_2}{2\sqrt{p_m^U}} - cM \right]^+ \end{cases} \quad (24)$$

where  $\alpha_m^U$  is the Lagrangian multiplier for  $\|\mathbf{v}_m^U\|_2^2 \leq p_m^U$ , and  $[\cdot]^+ = \max\{0, \cdot\}$ .

The problems of  $\mathbf{v}_m^D$  and  $\{\mathbf{v}_{c,m}^D\}_{m \in \mathcal{M}}$  can be similarly solved. We omit the details for brevity.

2) *Update  $\{s_{d,m}^U, y_{d,m}^U, z_{d,m}^U\}$ ,  $\{s_{d,m}^D, y_{d,m}^D, z_{d,m}^D\}$ , and  $\{s_{c,m}, y_{c,m}, z_{c,m}\}$ :* The problem related to  $\{s_{d,m}^U, y_{d,m}^U, z_{d,m}^U\}$  is expressed as

$$\begin{aligned}
 (\text{PM}) : \quad & \min_{\{s_{d,m}^U, y_{d,m}^U, z_{d,m}^U\}} \left[ \frac{c}{2} (y_{d,m}^U - B^U x_{d,m}^U - \frac{1}{c} \psi_{d,m}^U)^2 \right. \\
 & \quad \left. + \frac{c}{2} (z_{d,m}^U - d_m - \epsilon + \frac{1}{c} \theta_{d,m}^U)^2 \right. \\
 & \quad \left. - 2\hat{s}_{d,m}^U s_{d,m}^U \right] \\
 & \text{s.t.} \quad y_{d,m}^U \geq \frac{(s_{d,m}^U)^2}{z_{d,m}^U}, \quad z_{d,m}^U \geq \epsilon. \quad (25)
 \end{aligned}$$

Exploring the first-order optimality conditions, we obtain

$$\begin{cases} s_{d,m}^U = \frac{s_{d,m}^U z_{d,m}^U}{\beta_{d,m}^U}, \\ y_{d,m}^U = \frac{c B^U x_{d,m}^U + \psi_{d,m}^U + \beta_{d,m}^U}{c}, \\ z_{d,m}^U = \max \left\{ \left[ d_m + \epsilon + \frac{(s_{d,m}^U)^2}{c \beta_{d,m}^U} - \frac{\theta_{d,m}^U}{c} \right], \epsilon \right\}, \end{cases} \quad (26)$$

where  $\beta_{d,m}^U \geq 0$  is the Lagrangian multiplier of  $y_{d,m}^U \geq \frac{(s_{d,m}^U)^2}{z_{d,m}^U}$ , and should be chosen properly such that the Karush-Kuhn-Tucker (KKT) complementarity conditions [31] are satisfied. Specifically, we define  $\beta_{d,m}^* \triangleq \frac{(s_{d,m}^U)^2}{\theta_{d,m}^U - c d_m}$ . If  $\beta_{d,m}^* \geq 0$  and  $y_{d,m}^U \leq \frac{(s_{d,m}^U)^2}{z_{d,m}^U} |_{\beta_{d,m}^U = \beta_{d,m}^*}$ , we solve the cubic equation

$$\frac{c B^U x_{d,m}^U + \psi_{d,m}^U + \beta_{d,m}^U}{c} = \frac{(s_{d,m}^U)^2 \epsilon}{(\beta_{d,m}^U)^2} \quad (27)$$



$$\begin{aligned}
\mathcal{L}_c(\mathbf{V}, \mathbf{d}, \mathbf{s}, \mathbf{T}, \mathbf{x}, \mathbf{y}, \mathbf{z}, \Phi, \psi, \theta) = & -2 \sum_{m \in \mathcal{M}} (\hat{s}_{d,m}^U s_{d,m}^U + \hat{s}_{d,m}^D s_{d,m}^D + \hat{s}_{c,m} s_{c,m}) \\
& + \sum_{m \in \mathcal{M}} \sum_{n \in \mathcal{M}} \left\{ \text{Re} \left[ (\phi_{d,n}^U(m))^\dagger (\mathbf{t}_{d,n}^U(m) - \mathbf{v}_n^U) + (\phi_{c,n}^U(m))^\dagger (\mathbf{t}_{c,n}^U(m) - \mathbf{v}_n^U) \right] \right. \\
& \quad \left. + \frac{c}{2} [\|\mathbf{t}_{d,n}^U(m) - \mathbf{v}_n^U\|_2^2 + \|\mathbf{t}_{c,n}^U(m) - \mathbf{v}_n^U\|_2^2] \right\} \\
& + \sum_{m \in \mathcal{M}} \sum_{n \in \mathcal{M}} \left\{ \text{Re} \left[ (\phi_{d,n}^D(m))^\dagger (\mathbf{t}_{d,n}^D(m) - \mathbf{v}_{d,n}^D) + (\phi_{c,n}^D(m))^\dagger (\mathbf{t}_{c,n}^D(m) - \mathbf{v}_{c,n}^D) \right] \right. \\
& \quad \left. + \frac{c}{2} [\|\mathbf{t}_{d,n}^D(m) - \mathbf{v}_{d,n}^D\|_2^2 + \|\mathbf{t}_{c,n}^D(m) - \mathbf{v}_{c,n}^D\|_2^2] \right\} \\
& + \sum_{m \in \mathcal{M}} \sum_{n \in \mathcal{M}} \left\{ \text{Re} \left[ (\phi_{c,n}^D(m))^\dagger (\mathbf{t}_{c,n}^D(m) - \mathbf{v}_{c,n}^D) + (\phi_{c,n}^D(m))^\dagger (\mathbf{t}_{c,n}^D(m) - \mathbf{v}_{c,n}^D) \right] \right. \\
& \quad \left. + \frac{c}{2} [\|\mathbf{t}_{c,n}^D(m) - \mathbf{v}_{c,n}^D\|_2^2 + \|\mathbf{t}_{c,n}^D(m) - \mathbf{v}_{c,n}^D\|_2^2] \right\} \\
& + \sum_{m \in \mathcal{M}} \left\{ \psi_{d,m}^U (B^U x_{d,m}^U - y_{d,m}^U) + \psi_{d,m}^D (B^D x_{d,m}^D - y_{d,m}^D) + \frac{c}{2} [(B^U x_{d,m}^U - y_{d,m}^U)^2 + (B^D x_{d,m}^D - y_{d,m}^D)^2] \right\} \\
& + \sum_{m \in \mathcal{M}} \left\{ \psi_{c,m}^U (B^U x_{c,m}^U - y_{c,m}^U) + \psi_{c,m}^D (B^D x_{c,m}^D - y_{c,m}^D) + \frac{c}{2} [(B^U x_{c,m}^U - y_{c,m}^U)^2 + (B^D x_{c,m}^D - y_{c,m}^D)^2] \right\} \\
& + \sum_{m \in \mathcal{M}} \left\{ \theta_{d,m}^U (z_{d,m}^U - d_m - \epsilon) + \theta_{d,m}^D (z_{d,m}^D - d_m - \epsilon) + \theta_{c,m} (z_{c,m} - 1 + d_m - \epsilon) \right. \\
& \quad \left. + \frac{c}{2} [(z_{d,m}^U - d_m - \epsilon)^2 + (z_{d,m}^D - d_m - \epsilon)^2 + (z_{c,m} - 1 + d_m - \epsilon)^2] \right\}
\end{aligned} \tag{22}$$

to find the optimal  $\beta_{d,m}^U$ ; otherwise, we solve the following quartic equation

$$\frac{cB^U x_{d,m}^U + \psi_{d,m}^U + \beta_{d,m}^U}{c} = \frac{(\hat{s}_{d,m}^U)^2}{(\beta_{d,m}^U)^2} \cdot \left[ d_m + \epsilon + \frac{(\hat{s}_{d,m}^U)^2}{c\beta_{d,m}^U} - \frac{\theta_{d,m}^U}{c} \right] \tag{28}$$

to get the optimal  $\beta_{d,m}^U$ . Since both cubic and quartic equations can be solved analytically, the closed-form solution to  $\beta_{d,m}^U$  is achievable.

The problems of  $\{s_{d,m}^D, y_{d,m}^D, z_{d,m}^D\}$  and  $\{s_{c,m}, y_{c,m}, z_{c,m}\}$  can be solved similarly. We omit the details for brevity.

### C. Update $\{\mathbf{T}, \mathbf{d}, \mathbf{x}\}$

It can be easily observed that updating  $\{\mathbf{T}, \mathbf{d}, \mathbf{x}\}$  is completely separable among  $\{\mathbf{T}_d^U(m), x_{d,m}^U\}$ ,  $\{\mathbf{T}_d^D(m), x_{d,m}^D\}$ ,  $\{\mathbf{T}_c^U(m), x_{c,m}^U\}$ ,  $\{\mathbf{T}_c^D(m), x_{c,m}^D\}$  and  $\{d_m\}$ ,  $\forall m \in \mathcal{M}$ . Then, we update the  $5M$  subproblems independently and in parallel.

1) *Update  $\{\mathbf{T}_d^U(m), x_{d,m}^U\}$ ,  $\{\mathbf{T}_d^D(m), x_{d,m}^D\}$ ,  $\{\mathbf{T}_c^U(m), x_{c,m}^U\}$ , and  $\{\mathbf{T}_c^D(m), x_{c,m}^D\}$* : The problem related to  $\{\mathbf{T}_d^U(m), x_{d,m}^U\}$  is expressed as

$$\begin{aligned}
(\text{PN}) : \min_{\{\mathbf{T}_d^U(m), x_{d,m}^U\}} & \frac{c}{2} \left[ \sum_{n \in \mathcal{M}} \|\mathbf{t}_{d,n}^U(m) - \mathbf{v}_n^U + \frac{1}{c} \phi_{d,n}^U(m)\|_2^2 \right. \\
& \quad \left. + (B^U x_{d,m}^U - y_{d,m}^U - \frac{1}{c} \psi_{d,m}^U)^2 \right] \\
\text{s.t.} & \log(w_{d,m}^U) - w_{d,m}^U e_{d,m}^U (\mathbf{T}_d^U(m), \mathbf{u}_{d,m}^U) + 1 \geq x_{d,m}^U.
\end{aligned} \tag{29}$$

It is solved as

$$\begin{cases} \mathbf{t}_{d,m}^U(m) = [2\delta_{d,m}^U w_{d,m}^U (\mathbf{F}_{m,m}^U)^\dagger \mathbf{u}_{d,m}^U (\mathbf{u}_{d,m}^U)^\dagger \mathbf{F}_{m,m}^U + c\mathbf{I}]^{-1} \\ \quad \times [c\mathbf{v}_m^U - \phi_{d,m}^U(m) + 2\delta_{d,m}^U w_{d,m}^U (\mathbf{F}_{m,m}^U)^\dagger \mathbf{u}_{d,m}^U] \\ \mathbf{t}_{d,n}^U(m) = [2\delta_{d,m}^U w_{d,m}^U (\mathbf{F}_{m,n}^U)^\dagger \mathbf{u}_{d,m}^U (\mathbf{u}_{d,m}^U)^\dagger \mathbf{F}_{m,n}^U + c\mathbf{I}]^{-1} \\ \quad \times [c\mathbf{v}_n^U - \phi_{d,n}^U(m)], \quad \forall n \neq m \\ x_{d,m}^U = \frac{cB^U y_{d,m}^U + B^U \psi_{d,m}^U - \delta_{d,m}^U}{c(B^U)^2} \end{cases} \tag{30}$$

where  $\delta_{d,m}^U \geq 0$  is the associated Lagrangian multiplier and should be properly chosen such that the KKT complementarity conditions are satisfied. That is, if (29) holds for  $\delta_{d,m}^U = 0$ , we have  $\delta_{d,m}^U = 0$ ; otherwise, we choose some  $\delta_{d,m}^U > 0$  such that (29) holds for equality, which can be easily done by bisection search.

Similarly, the problems of  $\{\mathbf{T}_d^D(m), x_{d,m}^D\}$ ,  $\{\mathbf{T}_c^U(m), x_{c,m}^U\}$  and  $\{\mathbf{T}_c^D(m), x_{c,m}^D\}$  can be solved.

2) *Update  $d_m$* : The problem of  $d_m$  is

$$\begin{aligned}
(\text{PO}) : \min_{d_m} & \left[ (\theta_{c,m} - \theta_{d,m}^U - \theta_{d,m}^D) d_m + \frac{c}{2} (d_m + \epsilon - z_{d,m}^U)^2 \right. \\
& \quad \left. + \frac{c}{2} [(d_m + \epsilon - z_{d,m}^D)^2 + (d_m - 1 - \epsilon + z_{c,m})^2] \right] \\
\text{s.t.} & 0 \leq d_m \leq 1,
\end{aligned} \tag{31}$$

which is solved as

$$d_m = \left[ \frac{c(z_{d,m}^U + z_{d,m}^D - z_{c,m} + 1 - \epsilon) - (\theta_{c,m} - \theta_{d,m}^U - \theta_{d,m}^D)}{3c} \right]_0^1 \tag{32}$$

where  $[\cdot]_0^1$  denotes the projection onto the range of  $[0, 1]$ .

Since (PH) can be solved distributively, we obtain the fully distributed IA-IWMMSE by updating  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$  based on the ADMM algorithm in Table V.

*Remark 2:* The complexity comparison of the ADMM algorithm and the IP algorithm for (PH) is listed in Table VI. With the problem dimension being  $M(N_u + N_b)$ , the per-iteration complexity of the IP algorithm is  $\mathcal{O}(M^3(N_u + N_b)^3)$ . By contrast, the complexity of the ADMM algorithm is dominated by the update of  $\{\mathbf{T}, \mathbf{d}, \mathbf{x}\}$ . As shown in (30), the per-iteration complexity is  $\mathcal{O}(\max\{M^2 N_u^2, M^2 N_u N_b\})$  if we use the rank-one update rule to compute the matrix inverse. Obviously, the ADMM algorithm is more efficient.

## V. NUMERICAL EXAMPLES

Consider a network consisting of one multi-antenna BS and  $M = 10$  pairs of multi-antenna UEs. The BS is located at the

TABLE VI  
COMPLEXITY COMPARISON OF ADMM ALGORITHM AND IP ALGORITHM

	Per-Iteration Complexity
IP Algorithm	$\mathcal{O}(M^3(N_u + N_b)^3)$
ADMM Algorithm	$\mathcal{O}(\max\{M^2N_u^2, M^2N_uN_b\})$

center of a hexagonal cell with the side length being  $d = 1$  km, and the UEs are randomly deployed in the cell. We assume the Rayleigh channel with zero mean and variance  $L(200/r)^3$ , with  $r$  being the distance between transmitter and receiver, and  $L$  being the shadowing effect satisfying  $10 \log_{10}(L) \sim \mathcal{N}(0, 64)$ . We assume that all TX\_UEs have same transmit power budgets in uplink and downlink, i.e.,  $p_1^U = p_2^U = \dots = p_M^U$ ,  $p_1^D = p_2^D = \dots = p_M^D$ . The background noise spectral density is  $-174$  dBm/Hz, and the noise power depends on the bandwidth values, e.g.,  $B^U$  and  $B^D$ . The parameter  $\epsilon$  is set as  $10^{-6}$  in the following simulations.

We first show the convergence behaviour of the ADMM algorithm for (PH). To this end, we solve (PH) by CVX and then solve the equivalent (PK) by ADMM under the same network configurations. The typical convergence traces (starting from the same initial point) are shown in Fig. 2. It is clear that the ADMM algorithm can converge to the CVX solution, while the convergence behaviour is sensitive to the value of  $c$ . In our simulation, the ADMM algorithm typically converges in about 200 iterations with  $c = 10$ .

Next, we validate the efficiency advantage of the ADMM algorithm against CVX (the IP algorithm) for (PH). In Fig. 3, we compare the CPU running times of the two approaches at different UE pair numbers, which are obtained by running the simulations on a computer with i5-3320M CPU and 8 GB RAM. The results show that the ADMM algorithm has higher efficiency regardless of the UE pair number  $M$ , owing to the simple closed-form solutions in each step. We should mention here that the results in Fig. 3 are not consistent with those in Table VI. This can be explained as follows. In Table VI we actually compare the per-iteration complexities of the two algorithms, while in Fig. 3 we show the total CPU times they need to solve (PH). The mismatch between them is mainly due to the fact that the two algorithms require different numbers of iterations to achieve the optimal solution of (PH).

In Fig. 4, we compare the convergence behaviours of IA-IWMMSE and IA-EWMMSE. In IA-EWMMSE, we update  $\hat{s}$  only after (PG) has been exactly solved (i.e., at the red circles in Fig. 4). Notice that updating  $\hat{s}$  yields apparent performance improvement, especially in the first few DC programming iterations. Generally, IA-EWMMSE converges after 250 WMMSE iterations. By contrast, we update  $\hat{s}$  in each WMMSE iteration in IA-IWMMSE, and thereby the most of performance improvement can be achieved in the first few iterations. As a consequence, IA-IWMMSE is more efficient and converges faster (typically in 100 WMMSE iterations). It can also be observed in Fig. 4 that the performance of IA-IWMMSE is better than that of IA-EWMMSE. This is due to the fact that the two algorithms may converge to different stationary solutions. Since (PF) is non-convex, it is difficult

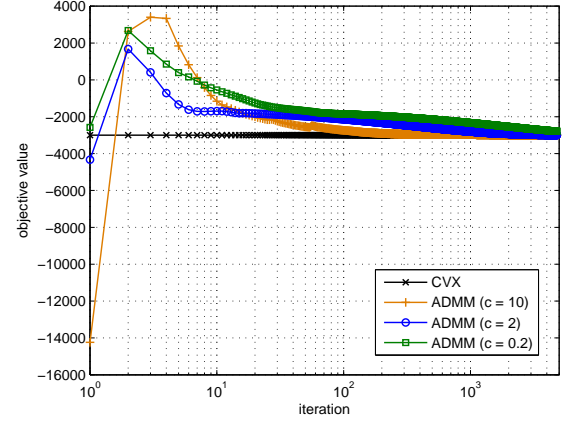


Fig. 2. Typical convergence traces of the ADMM algorithm for (PH), with  $M = 10$ ,  $N_b = 3$ , and  $N_u = 2$ .

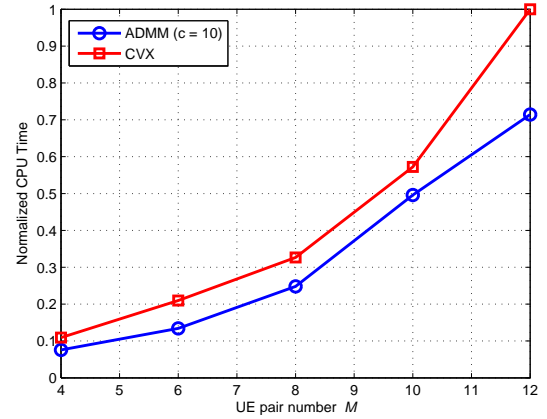


Fig. 3. CPU time comparison at different UE pair numbers  $M$ , with  $N_b = 3$  and  $N_u = 2$ .

to achieve the global optimum. Both IA-EWMMSE and IA-IWMMSE are only guaranteed to converge to a stationary solution (local optimum) of (PF). The two algorithms have the same computational steps, while the order of execution is different (see Tables II and III). In consequence, they have different convergence traces and finally reach different local optimums. However, since the objective function of (PF) depends on the random D2D and cellular channels, there is no theoretical guarantee that one local optimum will be better than the other. This is confirmed in Fig. 5, where we record the sum rates and iteration numbers for convergence of the two algorithms with 50 different channels. Interestingly, although the two algorithms may behave differently in each simulation trial, their average performances are close to each other. The advantage of IA-IWMMSE against IA-EWMMSE mainly lies in the simplified algorithm structure (see Table III) and the accelerated convergence behaviour (see Figs. 4 and 5).

After that, we compare IA-IWMMSE with the following approaches to show its advantages in sum-rate maximization. (1) IA-EWMMSE. (2) All-cellular approach; i.e., all UE pairs communicate through BS. (3) All-D2D approach; i.e., all

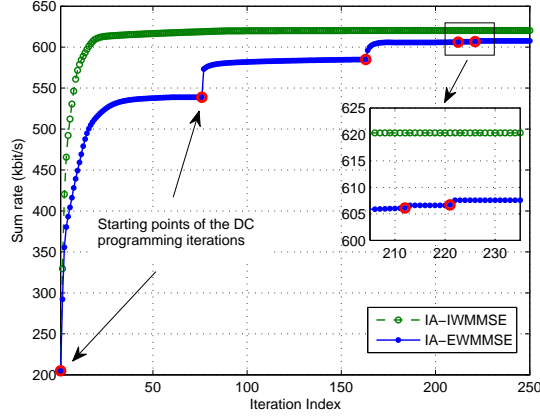


Fig. 4. Typical convergence traces of IA-IWMMSE and IA-EWMMSE, with  $M = 10$ ,  $N_b = 3$ ,  $N_u = 2$ ,  $B^U = B^D = 15$  kHz,  $P_B = 15$  dBm, and  $p_m^D = p_m^U = 0$  dBm,  $\forall m$ .

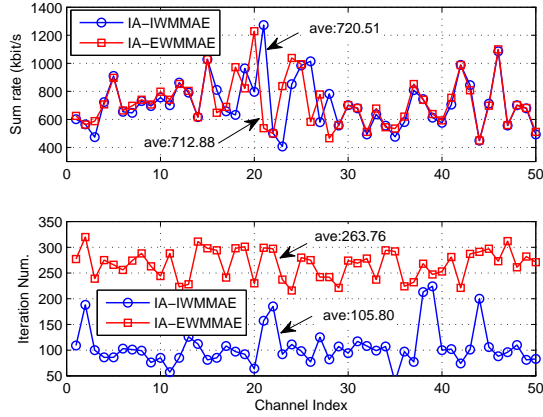


Fig. 5. Comparison of IA-IWMMSE and IA-EWMMSE for 50 different channels with  $M = 10$ ,  $N_b = 3$ ,  $N_u = 2$ ,  $B^U = B^D = 15$  kHz,  $P_B = 15$  dBm,  $p_m^D = p_m^U = 0$  dBm,  $\forall m$ .

UE pairs communicate via the D2D link. (4) Exhaustive search approach, which compares the throughput results of all possible transmission mode settings and select the best one. This approach is computationally intolerable in practice, yet its results can be used to show the gap between our approximate solution and the optimum. (5) The approach in [18], which considers the joint mode selection and transceiver design for sum rate maximization in MIMO D2D networks, and solves it by utilizing the method of SCA. Besides the solution approach, our work differs from [18] mainly in that we employ an exclusive (i.e., either cellular or D2D) mode selection strategy for each UE pair, while [18] allows simultaneous cellular and D2D transmission (i.e., it does not use  $d_m$  when computing the achievable rate of the  $m$ th UE pair, see  $R_m$  in (8)) to further improve the network throughput, which, of course, raises heavy operational burden and signaling cost in practice. To avoid unfair comparison, as simultaneous cellular and D2D transmission occurs for some UE pair in [18], we force it to select the mode contributing more to the throughput.

In Fig. 6, we compare the sum rates of different approaches

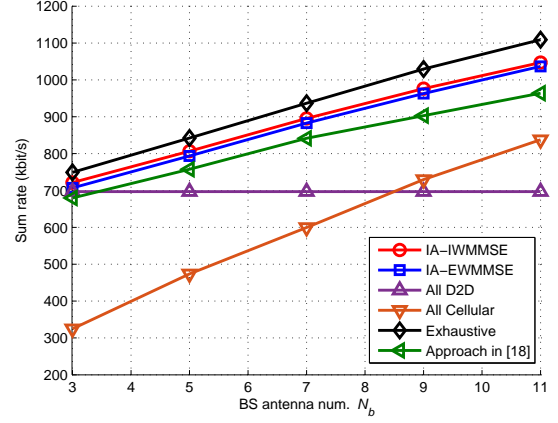


Fig. 6. Sum-rate comparison at different BS antenna numbers  $N_b$ , with  $M = 10$ ,  $N_u = 2$ ,  $B^U = B^D = 15$  kHz,  $P_B = 15$  dBm,  $p_m^U = 0$  dBm,  $p_m^D = 0$  dBm,  $\forall m$ .

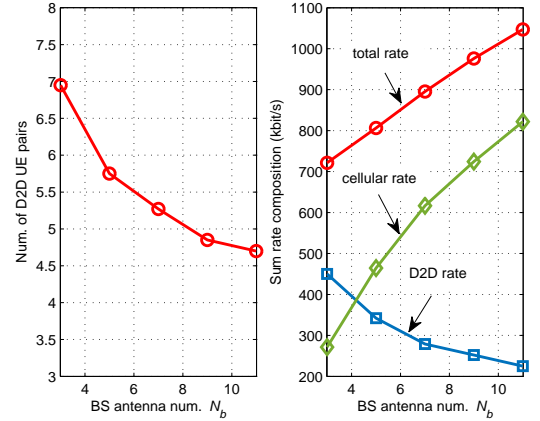


Fig. 7. Some details of IA-IWMMSE at different BS antenna numbers  $N_b$ , with  $M = 10$ ,  $N_u = 2$ ,  $B^U = B^D = 15$  kHz,  $P_B = 15$  dBm,  $p_m^U = 0$  dBm,  $p_m^D = 0$  dBm,  $\forall m$ .

at distinct BS antenna numbers  $N_b$ . Among these approaches, the all-D2D approach works independently of the BS antenna number, while the other approaches achieve higher throughput as  $N_b$  increases, since more space diversity can be exploited. The exhaustive search approach surely achieves the highest sum rate by exploring all the mode selection possibilities. By performing the system-level optimization on transmission mode selection and transceiver design, IA-IWMMSE and IA-EWMMSE achieve the sum rates close to that of the exhaustive search approach and outperform the other approaches. The performance gap between IA-IWMMSE and the exhaustive search approach is mainly due to the relaxation of  $d_m$  and the first-order approximation in solving (PB). Although the approach in [18] also jointly optimizes the transmission mode and the associated transceivers for each UE pair, it is not designed for the exclusive mode selection strategy. Consequently, when required to select only one kind of transmission for each UE pair, it suffers from performance degradation because the mode selection may not be optimal. When all the UE pairs communicate through BS, the network tends to be congested.

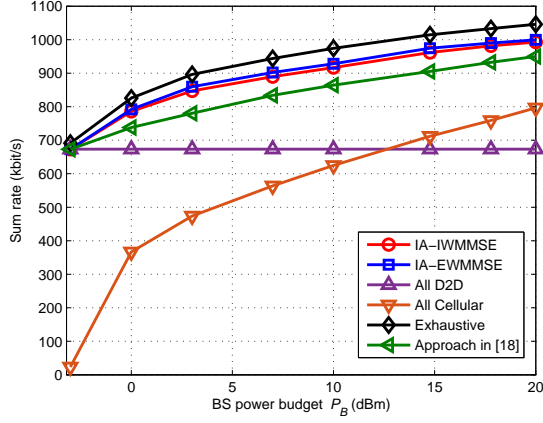


Fig. 8. Sum-rate comparison at different BS power budgets  $P_B$ , with  $M = 10$ ,  $N_b = 10$ ,  $N_u = 2$ ,  $B^U = B^D = 15$  kHz,  $p_m^U = -3$  dBm,  $p_m^D = -3$  dBm,  $\forall m$ .

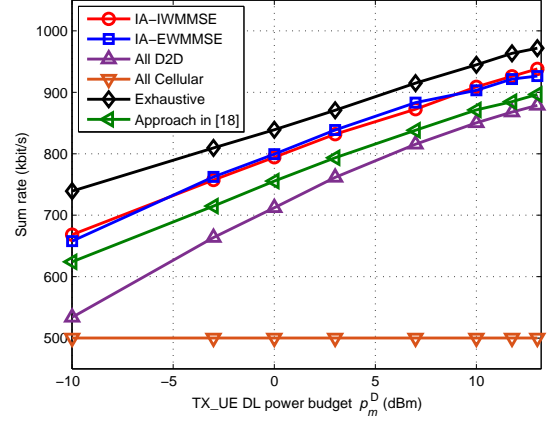


Fig. 10. Sum-rate comparison at different UE downlink power budgets  $p_m^D$ , with  $M = 10$ ,  $N_b = 5$ ,  $N_u = 2$ ,  $B^U = B^D = 15$  kHz,  $P_B = 15$  dBm,  $p_m^U = 0$  dBm,  $\forall m$ .

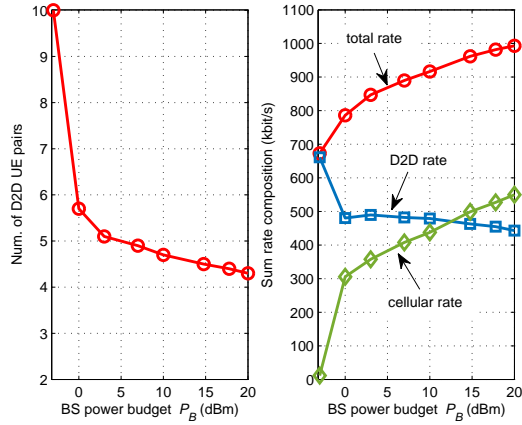


Fig. 9. Some details of IA-IWMMSE at different BS power budgets  $P_B$ , with  $M = 10$ ,  $N_b = 10$ ,  $N_u = 2$ ,  $B^U = B^D = 15$  kHz,  $p_m^U = -3$  dBm,  $p_m^D = -3$  dBm,  $\forall m$ .

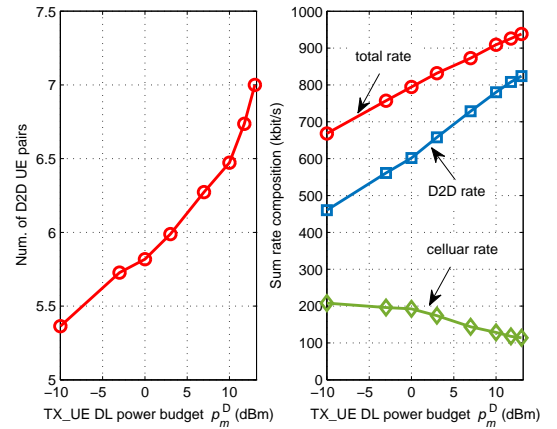


Fig. 11. Some details of IA-IWMMSE at different UE DL power budgets  $p_m^D$ , with  $M = 10$ ,  $N_b = 5$ ,  $N_u = 2$ ,  $B^U = B^D = 15$  kHz,  $P_B = 15$  dBm,  $p_m^U = 0$  dBm,  $\forall m$ .

Therefore, the all-cellular approach has the worst performance in the case of small  $N_b$ . When  $N_b$  increases, the congestion will be alleviated, and hence it gradually outperforms the all-D2D approach. In Fig. 7, we further show some details of IA-IWMMSE. As  $N_b$  increases, the extra space diversity in cellular transmission drives more UE pairs to select the cellular mode. Then, the D2D UE pair number decreases and the traffic through BS increases.

We compare the sum rates of these approaches at different BS power budgets  $P_B$  in Fig. 8, and show some details of IA-IWMMSE in Fig. 9. In the case of very small BS power budget (e.g.,  $P_B = -3$  dBm), IA-IWMMSE suggests that all the UE pairs should select the D2D mode. Therefore, IA-IWMMSE and the all-D2D approach have the same performance. IA-EWMMSE and the approach in [18] achieve similar performance at this point, implying that they also select the D2D mode for all UE pairs. The above approaches outperform the all-cellular approach when  $P_B$  is small. As  $P_B$  increases, the D2D UE pairs are more likely to select the cellular mode for higher rate. As a consequence, the D2D UE pair number

decreases and the sum rate increases. The order of achievable sum rate in this simulation is the same as that in Fig. 6 due to the reason aforementioned. We emphasize that IA-IWMMSE and IA-EWMMSE perform close to the exhaustive approach and outperform the other approaches.

We compare the sum rates of these approaches at different UE downlink power budgets  $p_m^D$  in Fig. 10. In this simulation, the sum rate of the all-cellular approach is independent of  $p_m^D$ , while the sum rates of other approaches increase with  $p_m^D$  since they all involve D2D transmission. The descending order of these approaches in terms of sum rate is the exhaustive search approach, IA-IWMMSE (IA-EWMMSE), the approach in [18], the all-D2D approach, and the all-cellular approach. In Fig. 11, we show more details of IA-IWMMSE. In cellular mode, the rate of one UE pair is determined actually by the smaller one of uplink rate and downlink rate, while in D2D mode the rate is the sum of them. Therefore, increasing the UE downlink power budget improves the rates of D2D UE pairs, thereby encouraging UE pairs to select the D2D mode. In consequence, the D2D UE pair number and the sum rate

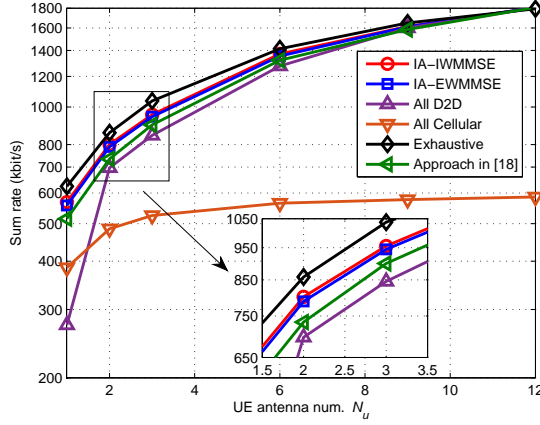


Fig. 12. Sum-rate comparison at different UE antenna number  $N_u$ , with  $M = 10$ ,  $N_b = 5$ ,  $B^U = B^D = 15$  kHz,  $P_B = 15$  dBm,  $p_m^U = 0$  dBm,  $p_m^D = 0$  dBm,  $\forall m$ .

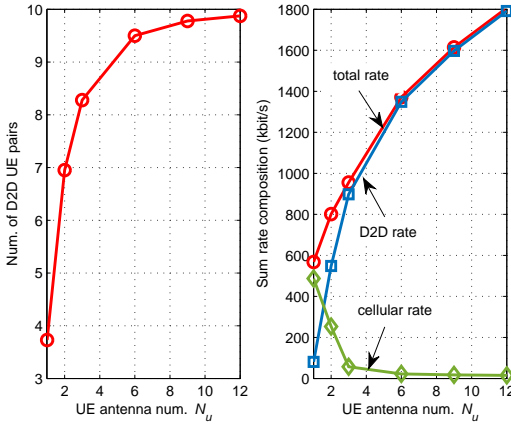


Fig. 13. Some details of IA-IWMMSE at different UE antenna number  $N_u$ , with  $M = 10$ ,  $N_b = 5$ ,  $B^U = B^D = 15$  kHz,  $P_B = 15$  dBm,  $p_m^U = 0$  dBm,  $p_m^D = 0$  dBm,  $\forall m$ .

increase with  $p_m^D$ . The proportion of D2D rate in the sum rate increases with  $p_m^D$  also.

We compare the sum rates of these approaches at different UE antenna numbers  $N_u$  in Fig. 12, and show more details of IA-IWMMSE in Fig. 13. As  $N_u$  increases, all the approaches in comparison achieve higher sum rates. However, the performance of the all-cellular approach grows far more slowly compared with the other approaches, especially in the case of large  $N_u$ . This is reasonable since when  $N_u$  is large, the BS antenna number  $N_b$  becomes the bottleneck of throughput improvement in cellular transmission, thus limiting the performance of the all-cellular approach. By contrast, the D2D transmission benefits from the extra space diversity introduced by increasing  $N_u$ . Then, the other approaches can effectively improve the network throughput by letting more UE pairs to select the D2D mode. Again, among these approaches, IA-IWMMSE achieves the performance close to that of exhaustive search and outperforms the others. It can be observed that both the number of D2D UE pairs and the traffic through D2D link increase with  $N_u$  in IA-IWMMSE. Moreover, in the case of

$N_u = 12$ , the evident advantage of D2D transmission in space diversity drives all the UE pairs to communicate directly, and hence the sum rate curves of all these approaches (except the all-cellular approach) coincide.

## VI. CONCLUSION

In this paper, we consider the joint mode selection and transceiver design problem in underlay D2D multi-user MIMO networks for sum rate maximization. Different from most current studies, we assume that the cellular and D2D transmissions are equally prioritized, and each UE pair can switch between cellular mode and D2D mode freely. We develop a low-complexity algorithm, based on inexact WMMSE, to iteratively compute a stationary solution of the sum-rate maximization problem. To facilitate the algorithm's implementation, we further apply the ADMM framework to decouple the WMMSE problem. Finally, an efficient distributed algorithm is devised, with a simple closed-form update in each step. The efficacy and efficiency of the proposed algorithm have been demonstrated by extensive numerical simulations.

## APPENDIX A PROOF OF PROPOSITION 1

Theoretically, (PA) and (PB) can be solved by the method of exhaustive search; i.e., we compare the throughput results of all possible transmission mode settings and select the best one. Therefore, to establish the equivalence between (PA) and (PB), it is sufficient to show that their resultant problems with the same transmission mode setting are equivalent.

The key step in proof is to show that if UE pair  $m$  works in D2D mode and UE pair  $n$  works in cellular mode, i.e.,  $d_m = 1$  and  $d_n = 0$ , then we have  $\mathbf{v}_{c,m}^D = \mathbf{0}$  and  $\mathbf{v}_{d,n}^D = \mathbf{0}$  at the optimal solution of the resultant problem of (PB). This can be proven by contradiction.

Specifically, assume that UE pair  $m$  works in D2D mode, and  $\{\mathbf{V}^*, \mathbf{U}^*\}$  is the optimal solution of the resultant problem of (PB) with  $\mathbf{v}_{c,m}^D \neq \mathbf{0}$ . Then, we can get another feasible solution of the problem by simply letting  $\mathbf{v}_{c,m}^D = \mathbf{0}$ , which is denoted as  $\{\bar{\mathbf{V}}^*, \mathbf{U}^*\}$ . From the downlink SINR expressions of (6a) and (6b), we conclude that letting  $\mathbf{v}_{c,m}^D = \mathbf{0}$  eliminates the cellular interference from TX\_UE $_m$  to all the other downlink transmissions, and then improves  $\overline{\text{SINR}}_{d,n}^D, \forall n \in \mathcal{M}$ , and  $\overline{\text{SINR}}_{c,n}^D, \forall n \in \mathcal{M}, n \neq m$ . According to (7a) and (7b), this produces increased  $\bar{R}_{d,n}, \forall n \in \mathcal{M}$ , for the  $M$  D2D transmissions, and non-decreased  $\bar{R}_{c,n}, \forall n \in \mathcal{M}, n \neq m$ , for the  $(M-1)$  cellular transmissions. In other words, among the  $2M$  rate terms of  $\{\bar{R}_{d,n}, \bar{R}_{c,n}\}_{n \in \mathcal{M}}$ , only  $\bar{R}_{c,m}$  decreases to zero. Recall that UE pair  $m$  works in D2D mode, i.e.,  $d_m = 1$ , then from (8) we know that  $\bar{R}_{c,m}$  does not count for the throughput. Hence, we claim that  $\{\bar{\mathbf{V}}^*, \mathbf{U}^*\}$  is not the optimal solution since  $\{\mathbf{V}^*, \mathbf{U}^*\}$  yields higher throughput. There comes the contradiction. Similarly, we can show that if UE pair  $n$  works in D2D mode, we have  $\mathbf{v}_{d,n}^D = \mathbf{0}$  at the optimal solution of the resultant problem of (PB).

In conclusion, for the same transmission mode setting, the resultant problems of (PA) and (PB) are equivalent in the sense that they have the same optimal solutions due to the



same SINR and rate expressions. Based on this, we further claim the equivalence between (PA) and (PB).

## APPENDIX B PROOF OF PROPOSITION 2

Using subscript  $k$  to denote the iteration index, we express the updating rules in each iteration of IA-IWMMSE as

$$\{\mathbf{V}^{(k+1)}, \mathbf{d}^{(k+1)}, \mathbf{s}^{(k+1)}\} \in \Omega(\mathbf{U}^{(k)}, \mathbf{W}^{(k)}; \hat{\mathbf{s}}^{(k)}), \quad (33a)$$

$$\mathbf{U}^{(k+1)} = \Gamma(\mathbf{V}^{(k+1)}), \quad (33b)$$

$$\mathbf{W}^{(k+1)} = \Upsilon(\mathbf{V}^{(k+1)}, \mathbf{U}^{(k+1)}), \quad (33c)$$

$$\hat{\mathbf{s}}^{(k+1)} = \mathbf{s}^{(k+1)}. \quad (33d)$$

We next define the following two objective functions

$$Q(\mathbf{d}, \mathbf{s}, \mathbf{V}; \mathbf{U}, \mathbf{W}) \triangleq - \sum_{m \in \mathcal{M}} (s_{d,m}^U)^2 + (s_{d,m}^D)^2 + s_{c,m}^2 \quad (34a)$$

$$\hat{Q}(\mathbf{d}, \mathbf{s}, \mathbf{V}; \mathbf{U}, \mathbf{W}, \hat{\mathbf{s}}) \triangleq \sum_{m \in \mathcal{M}} \begin{bmatrix} (\hat{s}_{d,m}^U)^2 + (\hat{s}_{d,m}^D)^2 + \hat{s}_{c,m}^2 \\ - 2\hat{s}_{d,m}^U s_{d,m}^U - 2\hat{s}_{d,m}^D s_{d,m}^D \\ - 2\hat{s}_{c,m} s_{c,m} \end{bmatrix} \quad (34b)$$

where  $\hat{Q}(\mathbf{d}, \mathbf{s}, \mathbf{V}; \mathbf{U}, \mathbf{W}, \hat{\mathbf{s}})$  is the first-order approximation of  $Q(\mathbf{d}, \mathbf{s}, \mathbf{V}; \mathbf{U}, \mathbf{W})$  at  $\mathbf{s} = \hat{\mathbf{s}}$ ;  $Q$  is the objective of (PF), and  $\hat{Q}$  is a *general form* of the objective of (PG).

Since the WMMSE iterations produce monotonically decreasing objective values (even if only one WMMSE iteration is executed), and the first-order Taylor expansion is a global overestimate of any concave function, we claim that IA-IWMMSE generates monotonically decreasing objective sequence (see (36) and the following discussion). In addition, due to the compactness of the feasible set of  $\{\mathbf{V}, \mathbf{d}\}$  and the continuous mapping among  $\mathbf{V}$ ,  $\mathbf{U}$  and  $\mathbf{W}$ , the objective sequence is lower bounded and thus converges.

We denote  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^*\}$  as a limit point of the iterates generated by IA-IWMMSE, and  $Q^*$  as the corresponding objective value. Denote  $\{\mathbf{V}^{(k_t)}, \mathbf{d}^{(k_t)}, \mathbf{s}^{(k_t)}, \mathbf{U}^{(k_t)}, \mathbf{W}^{(k_t)}, \hat{\mathbf{s}}^{(k_t)}\}_{t=1}^\infty$  as a subsequence converging to this limit point. We show that  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \mathbf{U}^*\}$  is a stationary solution of (PE) by the following two Lemmas.

**Lemma 1:** If  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^*\}$  is a limit point of the iterates generated by IA-IWMMSE, then  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \mathbf{U}^*, \mathbf{W}^*\}$  is a stationary solution of (PG) with  $\hat{\mathbf{s}}$  replaced by  $\hat{\mathbf{s}}^*$ .

*proof:* We establish this by showing that  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \mathbf{U}^*, \mathbf{W}^*\}$  is a KKT point of (PG). To this end, we need to prove that  $\mathbf{U}^* = \Gamma(\mathbf{V}^*)$ ,  $\mathbf{W}^* = \Upsilon(\mathbf{V}^*, \mathbf{U}^*)$ , and  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*\} \in \Omega(\mathbf{U}^*, \mathbf{W}^*; \hat{\mathbf{s}}^*)$ . According to the updating rules in each IA-IWMMSE iteration, i.e., (33), we directly have  $\mathbf{U}^* = \Gamma(\mathbf{V}^*)$ ,  $\mathbf{W}^* = \Upsilon(\mathbf{V}^*, \mathbf{U}^*)$ , and  $\hat{\mathbf{s}}^* = \mathbf{s}^*$ , i.e.,

$$\begin{aligned} & \lim_{t \rightarrow \infty} \{\mathbf{V}^{(k_t)}, \mathbf{d}^{(k_t)}, \mathbf{s}^{(k_t)}, \mathbf{U}^{(k_t)}, \mathbf{W}^{(k_t)}, \hat{\mathbf{s}}^{(k_t)}\} \\ &= \{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^*\} \\ &= \{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \Gamma(\mathbf{V}^*), \Upsilon(\mathbf{V}^*, \mathbf{U}^*), \mathbf{s}^*\}. \end{aligned} \quad (35)$$

We next show that  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*\} \in \Omega(\mathbf{U}^*, \mathbf{W}^*; \hat{\mathbf{s}}^*)$ . Basically, we have

$$\begin{aligned} & \hat{Q}(\mathbf{d}, \mathbf{s}, \mathbf{V}; \mathbf{U}^{(k_t)}, \mathbf{W}^{(k_t)}, \hat{\mathbf{s}}^{(k_t)}) \\ & \geq \hat{Q}(\mathbf{d}^{(k_{t+1})}, \mathbf{s}^{(k_{t+1})}, \mathbf{V}^{(k_{t+1})}; \mathbf{U}^{(k_t)}, \mathbf{W}^{(k_t)}, \hat{\mathbf{s}}^{(k_t)}) \\ & \geq \hat{Q}(\mathbf{d}^{(k_{t+1})}, \mathbf{s}^{(k_{t+1})}, \mathbf{V}^{(k_{t+1})}; \mathbf{U}^{(k_{t+1})}, \mathbf{W}^{(k_{t+1})}, \hat{\mathbf{s}}^{(k_t)}) \\ & \geq Q(\mathbf{d}^{(k_{t+1})}, \mathbf{s}^{(k_{t+1})}, \mathbf{V}^{(k_{t+1})}; \mathbf{U}^{(k_{t+1})}, \mathbf{W}^{(k_{t+1})}) \\ & = \hat{Q}(\mathbf{d}^{(k_{t+1})}, \mathbf{s}^{(k_{t+1})}, \mathbf{V}^{(k_{t+1})}; \mathbf{U}^{(k_{t+1})}, \mathbf{W}^{(k_{t+1})}, \hat{\mathbf{s}}^{(k_{t+1})}). \end{aligned} \quad (36)$$

where the first inequality is due to the fact that  $\{\mathbf{V}^{(k_{t+1})}, \mathbf{d}^{(k_{t+1})}, \mathbf{s}^{(k_{t+1})}\}$  is the global optimal solution of (PH), which minimizes  $\hat{Q}(\cdot)$  with respect to  $\{\mathbf{V}, \mathbf{d}, \mathbf{s}\}$  at  $\mathbf{U} = \mathbf{U}^{(k_t)}$ ,  $\mathbf{W} = \mathbf{W}^{(k_t)}$ , and  $\hat{\mathbf{s}} = \hat{\mathbf{s}}^{(k_t)}$ ; similarly, the second inequality comes from the fact that  $\mathbf{U}^{(k_{t+1})}$  and  $\mathbf{W}^{(k_{t+1})}$  are the optimal solutions of the corresponding minimization problems with respect to  $\mathbf{U}$  and  $\mathbf{W}$ , i.e., (PI) and (PJ), respectively; the third inequality is because of that the first-order Taylor expansion is a global overestimate of a concave function; the last equality can be easily verified in (34). Interestingly, by replacing the first item by  $\hat{Q}(\mathbf{d}^{(k_t)}, \mathbf{s}^{(k_t)}, \mathbf{V}^{(k_t)}; \mathbf{U}^{(k_t)}, \mathbf{W}^{(k_t)}, \hat{\mathbf{s}}^{(k_t)})$ , (36) also indicates that IA-IWMMSE generates monotonically decreasing objective sequence of (PG).

Taking the limits of both sides of (36) with respect to  $t$ , we get<sup>3</sup>

$$\hat{Q}(\mathbf{d}, \mathbf{s}, \mathbf{V}; \mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^*) \geq \hat{Q}(\mathbf{d}^*, \mathbf{s}^*, \mathbf{V}^*; \mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^*) \geq Q^*, \quad (37)$$

which implies  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*\} \in \Omega(\mathbf{U}^*, \mathbf{W}^*; \hat{\mathbf{s}}^*)$ , or equivalently, the KKT conditions of (PH) are satisfied at  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*\}$  with  $\mathbf{U} = \mathbf{U}^*$ ,  $\mathbf{W} = \mathbf{W}^*$  and  $\hat{\mathbf{s}} = \hat{\mathbf{s}}^*$ .

For the sake of conciseness, we define the superscript  $\mathcal{A} \in \{\mathbf{U}, \mathbf{D}\}$  and the subscript  $l \in \{c, d\}$  to simplify the notations. For instance,  $h_{l,m}^{\mathcal{A}}(d_m, \mathbf{s}_m)$  with  $\mathbf{s}_m = \{s_{d,m}^U, s_{d,m}^D, s_{c,m}\}$  may denote the following four functions,

$$\begin{aligned} h_{d,m}^U(d_m, \mathbf{s}_m) &\triangleq \frac{(s_{d,m}^U)^2}{B^U(d_m + \epsilon)}, \quad h_{c,m}^U(d_m, \mathbf{s}_m) \triangleq \frac{s_{c,m}^2}{B^U(1 - d_m + \epsilon)}, \\ h_{d,m}^D(d_m, \mathbf{s}_m) &\triangleq \frac{(s_{d,m}^D)^2}{B^D(d_m + \epsilon)}, \quad h_{c,m}^D(d_m, \mathbf{s}_m) \triangleq \frac{s_{c,m}^2}{B^D(1 - d_m + \epsilon)}. \end{aligned}$$

The KKT conditions of (PH) at  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*\}$  with  $\{\mathbf{U}, \mathbf{W}, \hat{\mathbf{s}}\} = \{\mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^*\}$  are expressed as (38) at the top of next page, where  $\xi_{l,m}^{\mathcal{A}}, \tau_m, \kappa_m, \nu_m^U, \nu_{d,m}^D$  and  $\nu_c^D$  are the associated Lagrangian multipliers.

Due to the facts of  $\mathbf{U}^* = \Gamma(\mathbf{V}^*)$  and  $\mathbf{W}^* = \Upsilon(\mathbf{V}^*, \mathbf{U}^*)$ , we must have

$$\begin{cases} w_{l,m}^{\mathcal{A}} \nabla_{\mathbf{u}_{l,m}^{\mathcal{A}}} e_{l,m}^{\mathcal{A}}(\mathbf{V}^{\mathcal{A},*}, \mathbf{u}_{l,m}^{\mathcal{A},*}) = 0, \\ \quad \forall \mathcal{A} \in \{\mathbf{U}, \mathbf{D}\}, l \in \{c, d\}, m \in \mathcal{M}, \\ \nabla_{w_{l,m}^{\mathcal{A}}} [w_{l,m}^{\mathcal{A}} e_{l,m}^{\mathcal{A}}(\mathbf{V}^{\mathcal{A},*}, \mathbf{u}_{l,m}^{\mathcal{A},*}) - \log(w_{l,m}^{\mathcal{A},*})] = 0, \\ \quad \forall \mathcal{A} \in \{\mathbf{U}, \mathbf{D}\}, l \in \{c, d\}, m \in \mathcal{M}. \end{cases} \quad (39)$$

To complete the proof, we perform a tactful reformulation. In particular, we multiply the equations in (39) by the Lagrangian

<sup>3</sup>We remark that taking the limit inside the objective value  $\hat{Q}(\cdot)$  is legitimate, as the objective function is continuous.

$$\left\{ \begin{array}{l} \sum_{n \in \mathcal{M}} \sum_{l \in \{c,d\}} \xi_{l,n}^U w_{l,n}^{U,*} \nabla_{\mathbf{v}_m^U} e_{l,n}^U(\mathbf{V}^{U,*}, \mathbf{u}_{l,n}^{U,*}) + 2\nu_m^U \mathbf{v}_m^{U,*} = \mathbf{0}, \quad \forall m \in \mathcal{M}, \\ \sum_{n \in \mathcal{M}} \sum_{l \in \{c,d\}} \xi_{l,n}^D w_{l,n}^{D,*} \nabla_{\mathbf{v}_{d,m}^D} e_{l,n}^D(\mathbf{V}^{D,*}, \mathbf{u}_{l,n}^{D,*}) + 2\nu_{d,m}^D \mathbf{v}_{d,m}^{D,*} = \mathbf{0}, \quad \forall m \in \mathcal{M}, \\ \sum_{n \in \mathcal{M}} \sum_{l \in \{c,d\}} \xi_{l,n}^D w_{l,n}^{D,*} \nabla_{\mathbf{v}_{c,m}^D} e_{l,n}^D(\mathbf{V}^{D,*}, \mathbf{u}_{l,n}^{D,*}) + 2\nu_c^D \mathbf{v}_{c,m}^{D,*} = \mathbf{0}, \quad \forall m \in \mathcal{M}, \\ \nabla_{\chi_m} \hat{Q}(\mathbf{d}^*, \mathbf{s}^*, \mathbf{V}^*; \mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^*) + \sum_{\mathcal{A} \in \{U,D\}} \sum_{l \in \{c,d\}} \xi_{l,m}^{\mathcal{A}} \nabla_{\chi_m} h_{l,m}^{\mathcal{A}}(d_m^*, \mathbf{s}_m^*) + \tau_m \nabla_{\chi_m} d_m^* - \kappa_m \nabla_{\chi_m} d_m^* = 0, \\ \forall \chi_m \in \{d_m, s_{d,m}^U, s_{d,m}^D, s_{c,m}\}, \quad m \in \mathcal{M}, \\ 0 \leq \nu_m^U \perp p_m^U - \|\mathbf{v}_m^{U,*}\|_2^2 \geq 0, \quad 0 \leq \nu_{d,m}^D \perp p_m^D - \|\mathbf{v}_{d,m}^{D,*}\|_2^2 \geq 0, \quad \forall m \in \mathcal{M}, \\ 0 \leq \tau_m \perp d_m^* \geq 0, \quad 0 \leq \kappa_m \perp 1 - d_m^* \geq 0, \quad 0 \leq \nu_c^D \perp P_B - \sum_{m \in \mathcal{M}} \|\mathbf{v}_{c,m}^{D,*}\|_2^2 \geq 0, \\ 0 \leq \xi_{l,m}^{\mathcal{A}} \perp \log(w_{l,m}^{\mathcal{A},*}) - w_{l,m}^{\mathcal{A},*} e_{l,m}^{\mathcal{A}}(\mathbf{V}^{\mathcal{A},*}, \mathbf{u}_{l,m}^{\mathcal{A},*}) + 1 - h_{l,m}^{\mathcal{A}}(d_m^*, \mathbf{s}_m^*) \geq 0, \quad \forall \mathcal{A} \in \{U,D\}, \quad l \in \{c,d\}, \quad m \in \mathcal{M}. \end{array} \right. \quad (38)$$

$$\left\{ \begin{array}{l} \xi_{l,m}^{\mathcal{A}} w_{l,m}^{\mathcal{A},*} \nabla_{\mathbf{u}_{l,m}^{\mathcal{A}}} e_{l,m}^{\mathcal{A}} (\mathbf{V}^{\mathcal{A},*}, \mathbf{u}_{l,m}^{\mathcal{A},*}) = \mathbf{0}, \\ \quad \forall \mathcal{A} \in \{\text{U}, \text{D}\}, \quad l \in \{c, d\}, \quad m \in \mathcal{M}, \\ \xi_{l,m}^{\mathcal{A}} \nabla_{w_{l,m}^{\mathcal{A}}} [w_{l,m}^{\mathcal{A},*} e_{l,m}^{\mathcal{A}} (\mathbf{V}^{\mathcal{A},*}, \mathbf{u}_{l,m}^{\mathcal{A},*}) - \log(w_{l,m}^{\mathcal{A},*})] = 0, \\ \quad \forall \mathcal{A} \in \{\text{U}, \text{D}\}, \quad l \in \{c, d\}, \quad m \in \mathcal{M}. \end{array} \right. \quad (40)$$

*proof:* According to the results of Lemma 1, we know that  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \mathbf{U}^*, \mathbf{W}^*\}$  is a stationary solution of (PG), and  $\hat{\mathbf{s}}^* = \mathbf{s}^*$ . Based on the definitions of (34), we obtain

$$\begin{aligned} & \nabla_{\chi_m} Q(\mathbf{d}^*, \mathbf{s}^*, \mathbf{V}^*; \mathbf{U}^*, \mathbf{W}^*) \\ &= \nabla_{\chi_m} \hat{Q}(\mathbf{d}^*, \mathbf{s}^*, \mathbf{V}^*; \mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^* = \mathbf{s}^*), \end{aligned} \quad (41)$$

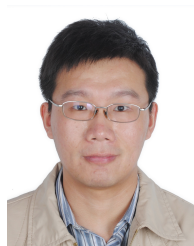
$$\forall \chi_m \in \{d_m, s_{d,m}^U, s_{d,m}^D, s_{c,m}\}, m \in \mathcal{M}.$$

Utilizing the properties of WMMSE reformulation (Theorem {2, 3} in [28] and Proposition 1 in [29]), we conclude that  $\{\mathbf{V}^*, \mathbf{d}^*, \mathbf{s}^*, \mathbf{U}^*, \mathbf{W}^*, \hat{\mathbf{s}}^*\}$  is also a stationary point of (PE).

- [1] J. Lin, Q. Shi, and Q. Li, "Joint device-to-device transmission activation and transceiver design for sum-rate maximization in MIMO interfering channels," in *Proc. IEEE ICASSP*, pp. 3921–3925, Shanghai, Mar. 2016.
- [2] A. Asadi, Q. Wang, and V. Mancuso, "A survey on device-to-device communication in cellular networks," *IEEE Communications Surveys & Tutorials*, vol. 16, no. 4, pp. 1801–1819, Nov. 2014.
- [3] K. Doppler, M. Rinne, C. Wijting, C.B. Ribeiro, and K. Hugl, "Device-to-device communication as an underlay to LTE-advanced networks," *IEEE Communications Magazine*, vol. 47, no. 12, pp. 42–49, Dec. 2009.

- [4] D. Feng, L. Lu, Y.-W. Yi, G.-Y. Li, S. Li, and G. Feng, "Device-to-device communications in cellular networks," *IEEE Communications Magazine*, vol. 52, no. 4, pp. 49–55, Apr. 2014.
- [5] H. Min, W. Seo, J. Lee, S. Park, and D. Hong, "Reliability improvement using receive mode selection in the device-to-device uplink period underlying cellular networks," *IEEE Trans. Wireless Communications*, vol. 10, no. 2, pp. 413–418, Feb. 2011.
- [6] L. Wei, R. Q. Hu, Y. Qian, and G. Wu, "Enable device-to-device communications underlying cellular networks: challenges and research aspects," *IEEE Communications Magazine*, vol. 52, no. 6, pp. 90–96, Jun. 2014.
- [7] G. Fodor, E. Dahlman, G. Mildh, S. Parkvall, N. Reider, G. Miklos, and Z. Turanyi, "Design aspects of network assisted device-to-device communications," *IEEE Communications Magazine*, vol. 50, no. 3, pp. 170–177, Mar. 2012.
- [8] C.-H. Yu, K. Doppler, C.B. Ribeiro, and O. Tirkkonen, "Resource sharing optimization for device-to-device communication underlying cellular networks," *IEEE Trans. Wireless Communications*, vol. 10, no. 8, pp. 2752–2763, Aug. 2011.
- [9] Z. Liu, T. Peng, S. Xiang, and W. Wang, "Mode selection for device-to-device (D2D) communication under LTE-advanced networks," in *Proc. IEEE ICC*, Ottawa, Canada, June 2012, pp. 5563–5567.
- [10] K. Doppler, C.-H. Yu, C. B. Ribeiro, and P. Janis, "Mode selection for device-to-device communication underlying an LTE-advanced network," in *Proc. IEEE WCNC*, Sydney, Australia, April 2010, pp. 1–6.
- [11] L. Song, D. Niyato, Z. Han, and E. Hossain, "Game-theoretic resource allocation methods for device-to-device (D2D) communication," *IEEE Wireless Communications Magazine*, vol. 21, no. 3, pp. 136–144, Jun. 2014.
- [12] R. Zhang, X. Cheng, L. Yang, and B. Jiao, "Interference-aware graph based resource sharing for device-to-device communications underlying cellular networks," in *Proc. IEEE WCNC*, Shanghai, China, April 2013, pp. 140–145.
- [13] L. Su, Y. Ji, P. Wang, and F. Liu, "Resource allocation using particle swarm optimization for D2D communication underlay of cellular networks," in *Proc. IEEE WCNC*, Shanghai, April 2013, pp. 129–133.
- [14] L. Wei, R.Q. Hu, T. He, and Y. Qian, "Device-to-device (D2D) communications underlying MU-MIMO cellular networks," in *Proc. IEEE GLOBECOM*, Atlanta, Dec. 2013, pp. 4902–4907.
- [15] P. Janis, V. Koivunen, C.B. Ribeiro, K. Doppler, and K. Hugl, "Interference-avoiding MIMO schemes for device-to-device radio underlying cellular networks," in *Proc. IEEE PIMRC*, Tokyo, Japan, Sept. 2009, pp. 2385–2389.
- [16] J.C.F. Li, M. Lei, and F. Gao, "Device-to-device (D2D) communication in MU-MIMO cellular networks," in *Proc. IEEE Global Telecommunications Conference (GLOBECOM)*, pp. 3583–3587, Anaheim, CA, Dec. 2012.
- [17] D. Zhu, W. Xu, H. Zhang, C. Zhao, J.C.F. Li, and M. Lei, "Rate maximized transceiver optimization for multi-antenna device-to-device communications," in *Proc. IEEE WCNC*, Shanghai, China, April 2013, pp. 4152–4157.

- [18] A. Tolli, J. Kaleva, and P. Komulainen, "Mode selection and transceiver design for rate maximization in underlay D2D MIMO systems," in *Proc. IEEE ICC*, London, UK, June 2015, pp. 4822–4827.
- [19] N. Reider and G. Fodor, "A distributed power control and mode selection algorithm for D2D communications," *EURASIP Journal on Wireless Communications and Networking*, 2012: 266, pp. 1–25, Aug. 2012.
- [20] E. Naghipour and M. Rasti, "A distributed joint power control and mode selection scheme for D2D communications underlying LET-A networks," in *Proc. IEEE WCNC*, Doha, Qatar, Apr. 2016, pp. 1–6.
- [21] H. Tang and Z. Ding, "Mixed mode transmission and resource allocation for D2D Communication," *IEEE Trans. Wireless Communications*, vol. 15, no. 1, pp. 162–175, Jan. 2016.
- [22] Q. Ye, M. Al-Shalash, C. Caramanis, and J.G. Andrews, "Distributed resource allocation in device-to-device enhanced cellular networks" *IEEE Trans. Communications*, vol. 63, no. 2, pp. 441–454, Feb. 2015.
- [23] D. Wu, Y. Cai, R.Q. Hu, and Y. Qian, "Dynamic distributed resource sharing for mobile MIMO D2D communications," *IEEE Trans. Wireless Communications*, vol. 14, no. 10, pp. 5417–5429, Oct. 2015.
- [24] X. Lin, J.G. Andrews, A. Ghosh, and R. Ratasuk, "An overview on 3GPP device-to-device proximity services," *IEEE Communications Magazine*, vol. 52, no. 4, pp. 40–48, Apr. 2014.
- [25] J. Liu, N. Kato, J. Ma, and N. Kadowaki, "Device-to-device communication in LTE-advanced networks: A survey," *IEEE Communication Surveys & Tutorials*, vol. 17, no. 4, pp.1923–1940, Fourth Quarter 2015.
- [26] P. Mach, Z. Becvar, and T. Vanek, "In-band device-to-device communication in OFDMA cellular networks: A survey and challenges," *IEEE Communication Surveys & Tutorials*, vol. 17, no. 4, pp. 1885–1992, Fourth Quarter 2015.
- [27] Z.-Q. Luo and S. Zhang, "Dynamic spectrum management: complexity and duality," *IEEE Journal of Selected Topics in Signal Processing*, vol. 20, no. 1, pp. 57–73, Feb. 2008.
- [28] Q. Shi, M. Razaviyayn, Z.-Q. Luo, and C. He, "An iteratively weighted MMSE approach to distributed sum-utility maximization for a MIMO interfering broadcast channel," *IEEE Trans. Signal Processing*, vol. 59, no. 9, pp. 4331–4340, Sep. 2011.
- [29] M. Razaviyayn, M. Hong, and Z.-Q. Luo, "Linear transceiver design for a MIMO interfering broadcast channel achieving max-min fairness," *Signal Processing*, vol. 93, no. 12, pp. 3327–3340, 2013.
- [30] M. Fukushima, "Application of alternating direction method of multipliers to separable convex programming problems," *Computational Optimization and Applications*, vol. 1, no. 1, pp. 93–111, Oct. 1992.
- [31] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge, UK: Cambridge University Press, 2004.
- [32] R. Horst and N.V. Thoai, "DC programming: overview," *J. of Optim. Theory and Appl.*, vol. 103, no. 1, pp. 1–43, 1999.
- [33] M. Grant, S. Boyd, "The CVX users' guide," *CVX Research Inc.*, Oct. 24, 2014 [on-line]. <http://cvxr.com/cvx/>.

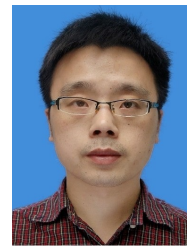


**Jingran Lin** received the B.S. degree in Computer Communication from University of Electronic Science and Technology of China (UESTC), Chengdu, China, in 2001, and the M.S. and Ph.D. degrees in Signal and Information Processing from UESTC in 2005 and 2007, respectively.

After his graduation in July 2007, he joined the School of Information and Communication Engineering, UESTC, where he is currently an Associate Professor. From January 2012 to January 2013, he was a Visiting Scholar with the University of Minnesota (Twin Cities), Minneapolis, MN, USA. His research interests include the design and analysis of efficient optimization algorithms for the signal processing problems arising from modern communication systems.



Dr. Shi is an Associate Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING. He was awarded Golden Medal at the 46th International Exhibition of Inventions of Geneva in 2018, and also was the recipient of the First Prize of Science and Technology Award from China Institute of Communications in 2017, the National Excellent Doctoral Dissertation Nomination Award in 2013, the Shanghai Excellent Doctorial Dissertation Award in 2012, and the Best Paper Award from the IEEE PIMRC'09 conference.



**Qingjiang Shi** received his Ph.D. degree in electronic engineering from Shanghai Jiao Tong University, Shanghai, China, in 2011. From September 2009 to September 2010, he visited Prof. Z.-Q. (Tom) Luo's research group at the University of Minnesota, Twin Cities. In 2011, he worked as a Research Scientist at Bell Labs China. From 2012, he was with the School of Information and Science Technology at Zhejiang Sci-Tech University. From Feb. 2016 to Mar. 2017, he worked as a research fellow at Iowa State University, USA. From 2018, he is currently a professor with the School of Software Engineering at Tongji University. His interests lie in algorithm design and analysis with applications in machine learning, signal processing and wireless networks. So far he has published more than 40 IEEE journals (five papers were ESI highly cited papers and one was nominated as the best paper award of IEEE Signal Processing Society in 2016) and filed more than 20 national patents.

Dr. Shi is an Associate Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING. He was awarded Golden Medal at the 46th International Exhibition of Inventions of Geneva in 2018, and also was the recipient of the First Prize of Science and Technology Award from China Institute of Communications in 2017, the National Excellent Doctoral Dissertation Nomination Award in 2013, the Shanghai Excellent Doctorial Dissertation Award in 2012, and the Best Paper Award from the IEEE PIMRC'09 conference.

**Qiang Li** received the B.Eng. and M.Phil. degrees in Communication and Information Engineering from University of Electronic Science and Technology of China (UESTC), Chengdu, China, and the Ph.D. degree in Electronic Engineering from the Chinese University of Hong Kong (CUHK), Hong Kong, in 2005, 2008, and 2012, respectively. From August 2011 to January 2012, he was a Visiting Scholar with the University of Minnesota, Minneapolis, MN, USA. From February 2012 to October 2013, he was a Research Associate with the Department of Electronic Engineering and the Department of Systems Engineering and Engineering Management, CUHK. Since November 2013, he has been with the School of Information and Communication Engineering, UESTC, where he is currently an Associate Professor. His research interests include convex optimization and its application in signal processing with an emphasis on the physical-layer security and full-duplex communications.

He received the First Prize Paper Award in the IEEE Signal Processing Society Postgraduate Forum Hong Kong Chapter in 2010, a Best Paper Award of IEEE PIMRC 2016, and the Best Paper Award of the IEEE Signal Processing Letters 2016.



**Dongmei Zhao** received the B.S. degree in Communication Engineering from Chongqing University, Chongqing, China, in 2015, and received the M.S. degree in Electronic and Communication Engineering from the University of Electronic Science and Technology of China (UESTC), Chengdu, China, in 2018.

She is currently working as a Software Engineer in ZTE Corporation. Her research interests include cooperative communication and communication protocol.