

Joint Mode Selection and Resource Allocation for D2D-Enabled NOMA Cellular Networks

Yanpeng Dai, Min Sheng [✉], *Senior Member, IEEE*, Junyu Liu, *Member, IEEE*, Nan Cheng [✉], *Member, IEEE*, Xuemin Shen, *Fellow, IEEE*, and Qinghai Yang [✉]

Abstract—The 5G cellular network employs non-orthogonal multiple access (NOMA) to enhance network connectivity and capacity, and device-to-device (D2D) communications to improve spectrum efficiency. However, the underlay D2D communications may destroy the execution condition for the successive interference cancellation (SIC) decoding of NOMA cellular networks by introducing the extra interference, which degrades the cellular transmission reliability. Thus, we develop the interlay mode as a special D2D mode for NOMA system, which enables the power domain multiplexing of the D2D pair and cellular users to eliminate the strong interference between them by the SIC decoding. When D2D pair conducts the selection between the interlay mode and underlay mode, the SIC decoding constraint should be satisfied at both D2D receiver and NOMA base station. In order to maximize the system sum rate while meeting the SIC decoding constraint, we propose a joint D2D mode selection and resource allocation scheme with interlay mode, which can be formulated as a combinatorial optimization problem. To tackle the combinatorial nature of mode selection and spectrum assignment, we first prove that the original problem can be reformulated as a maximum weight clique problem, and then propose a graph-based algorithm by applying branch-and-bound method to obtain its optimal solution. Finally, simulation results are provided to demonstrate that the interlay mode along with the proposed algorithms can coordinate D2D communications and NOMA cellular network to significantly improve the system sum rate and the D2D access rate.

Index Terms—Channel capacity, device-to-device communications, NOMA, radio spectrum management, 5G mobile communications.

I. INTRODUCTION

THE fifth generation (5G) networks will possess an architecture of heterogeneous networks integrating with

Manuscript received October 15, 2018; revised March 17, 2019; accepted May 5, 2019. Date of publication May 13, 2019; date of current version July 16, 2019. This work was supported in part by the Natural Science Foundation of China under Grants 61725103, 61701363, and 91638202, in part by the Shaanxi Innovation Group under Grant 2017KCT-30-03, in part by the Key Industry Innovation Chain of Shaanxi under Grant 2017ZDCXL-GY-04-04, in part by the China Postdoctoral Science Foundation under Grant 2018T111017, in part by the National S&T Major Project of China under Grants 2017ZX03001010-004 and 2018ZX03001014-004, and in part by the National Key R&D Program of China under Grant 2017YFB1010002. The review of this paper was coordinated by Dr. K. Bian. (*Corresponding author: Min Sheng.*)

Y. Dai, M. Sheng, J. Liu, and Q. Yang are with the State Key Laboratory of Integrated Service Networks, Institute of Information Science, Xidian University, Xi'an 710071, China (e-mail: yp_dai@stu.xidian.edu.cn; msheng@mail.xidian.edu.cn; junyuliu@xidian.edu.cn; qhyang@xidian.edu.cn).

N. Cheng and X. Shen are with the Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON N2L 3G1, Canada (e-mail: n5cheng@uwaterloo.ca; sshen@uwaterloo.ca).

Digital Object Identifier 10.1109/TVT.2019.2916395

multiple advanced communication technologies to satisfy high performance requirements [1]. Device-to-device (D2D) communication is one of the promising technologies to provide the proximity service for mobile users, which is incorporated into cellular networks to enhance the system capacity [2]–[4]. In D2D underlaying cellular networks, a pair of nearby devices reuses cellular resources to communicate without the relaying assistance of base station (BS), which improves spectrum efficiency while saving the radio bandwidth. However, it also introduces the mutual interference between cellular users and D2D pairs. Therefore, the efficient resource allocation and interference management algorithms have been proposed by a large body of literature [5]–[9]. However, all of these works focus on the traditional orthogonal multiple access (OMA) systems. Even though the efficient resource allocation algorithms are adopted, the mutual interference cannot be completely avoided, leading to degraded spectrum resource utilization.

To further enhance the network capacity, advanced multiple access technologies should be utilized for D2D and cellular hybrid networks. Recently, the non-orthogonal multiple access (NOMA) technology has been proposed, which has the potential to support massive connectivity and high capacity [10], [11]. Specifically, the NOMA technology¹ enables the user multiplexing on power domain to realize that several users reuse the same resource to simultaneously transmit the signals. The NOMA receiver applies the successive interference cancellation (SIC) technology to exploit the differences of signal strength to decode the users' signals. Many relative works have been proposed employing NOMA technology under many 5G scenarios [12]–[16]. Pan *et al.* in [17] studied the resource allocation for D2D underlaying NOMA cellular networks, where NOMA technique is used to enhance cellular users' throughputs. In [18], Zhao *et al.* applied NOMA technology into the D2D multicast group underlaying OMA cellular network to alleviate the interference between D2D transmissions. In D2D communications underlaying NOMA cellular networks, due to the non-orthogonal resource competition, the D2D transmission may destroy the execution condition for SIC decoding of NOMA cellular users, which can degrade the transmission reliability of NOMA cellular network. Furthermore, inside the NOMA system, the co-channel interference exists between cellular users, which is determined by the SIC decoding order and can affect the design of the D2D communication scheme. Actually, due to performing the dimension

¹The NOMA technologies mainly includes the code domain NOMA and power domain NOMA. In this paper, we focus on the typical power domain NOMA, and simply refer to it as NOMA in following text.

of power domain, the NOMA system has the potential to provide a new D2D mode to eliminate the mutual interference and improve the efficiency of resource reuse between the cellular user and D2D pair. In this paper, we define the interlay mode as a special D2D communication mode based on the features of NOMA. With the interlay mode, the D2D pair realizes the power-domain multiplexing with the cellular user on the same spectrum. According to the difference of their signals' power levels, the SIC decoding can avoid the mutual interference between the D2D pair and cellular user. Therefore, besides the underlay mode, the interlay mode is an alternative communication mode for D2D in NOMA cellular networks. The coexistence of the interlay mode and underlay mode will be a promising mechanism for D2D in NOMA system to satisfy the diverse requirements of D2D proximity services.

However, the resource management in D2D-enabled NOMA cellular networks with interlay mode is challenging. Specifically, if NOMA system is not properly designed, the co-channel interference may exist between cellular users and the interlay D2D pair, which affects the SIC decoding performance at both D2D receiver and BS. The co-channel interference highly depends on the SIC decoding order which is determined by the users' channel power gains and transmit powers. Therefore, the co-channel interference is more complicated than that in a simple NOMA uplink system where the SIC decoding constraint is only considered among cellular users. To mitigate the co-channel interference and improve the system performance, the mode selection and radio resource allocation should be jointly considered. Firstly, in the interlay mode, the SIC decoder can detect the difference of the signals on power domain to cancel the strong interference, and conversely, the signals in the underlay mode collide on frequency domain such that the interference cannot be avoided. Thus, the mode selection of a D2D pair directly influences the interference between it and other users, which is different from that in the OMA system where the interference is essentially determined by spectrum resource assignment. Secondly, in this hybrid NOMA network, the power allocation takes into account the interference among multiple NOMA cellular users, underlay D2D pairs, and interlay D2D pairs, which is more complicated than that in traditional D2D and cellular hybrid networks. Therefore, it is demanding to design a joint mode selection and power allocation scheme for the hybrid NOMA system to acquire the advantages of NOMA and D2D on spectrum efficiency and network connectivity.

In this paper, we investigate the resource management issue for D2D-enabled NOMA cellular networks where D2D communications are allowed to employ both the interlay mode and underlay mode. We formulate a joint optimization problem of mode selection between interlay mode and underlay mode, subcarrier assignment, and power control. We tackle this problem using a maximum weight clique approach in graph theory, which is a powerful tool for assignment problem and decision problem. As such, a graph model is utilized to depict the interference between cellular users and D2D pairs when D2D pairs use different subcarriers and modes. Moreover, the power control algorithm is designed to guarantee the users' minimum rates. The main contributions of this paper are summarized as follows.

- **Develop interlay mode:** We introduce the interlay mode to support power-domain multiplexing for cellular users and

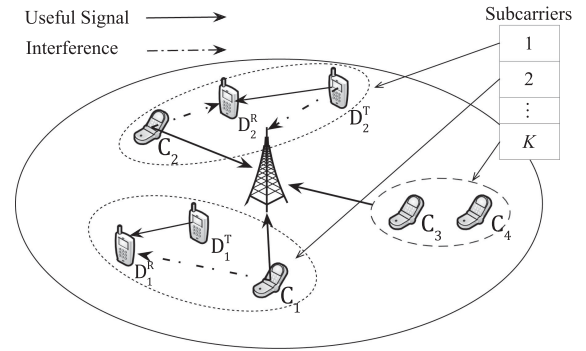


Fig. 1. D2D in NOMA uplink cellular network.

D2D pairs. Accordingly, we propose to enable both the interlay mode and underlay mode to promote the spectrum utilization.

- **Problem reformulation:** We prove that the original optimization problem can be reformulated into a maximum weight clique problem in graph theory, which helps us to tackle the combinatorial and non-convex nature of the original problem. Based on this, the efficient tools in the graph theory and the combinatorial optimization can be exploited to solve our constructed maximum weight clique problem, which increases the efficiency of our algorithm design.
- **Efficient algorithm design:** A graph-based scheme is proposed to maximize the system sum rate. In the proposed scheme, an iterative power allocation algorithm is proposed based on the sequential convex approximation to calculate the vertex weights. Then, a branch-and-bound based algorithm is proposed, which iteratively excludes the vertices representing inefficient allocation to achieve the optimal solution of mode selection and subcarrier assignment. We further propose a heuristic resource management algorithm with low computational complexity.

The remainder of the paper is organized as follows. In Section II, we introduce the network model, D2D communication modes, and problem formulation. The problem reformulation and the proposed joint mode selection and resource allocation scheme are presented in Section III. Then, we present a heuristic resource allocation algorithm in Section IV. Simulation results are given in Section V. Finally, we conclude this paper in Section VI.

II. SYSTEM MODEL AND PROBLEM FORMULATION

In this section, we introduce the network scenario and D2D communication modes, and then formally formulate the mode selection and resource allocation problem.

A. Network Model

As shown in Fig. 1, we consider the D2D communication in a single-cell NOMA uplink network. In this scenario, there are M cellular users, N D2D pairs, and K orthogonal subcarriers, and the sets of them are denoted by $\mathcal{C} = \{C_1, \dots, C_m, \dots, C_M\}$, $\mathcal{D} = \{D_1, \dots, D_n, \dots, D_N\}$, and $\mathcal{K} = \{S_1, \dots, S_k, \dots, S_K\}$, respectively. Furthermore, each D_n consists of a transmitter D_n^T and a receiver D_n^R . We assume that all links experience independent block fading, and the subcarriers follow the frequency-flat

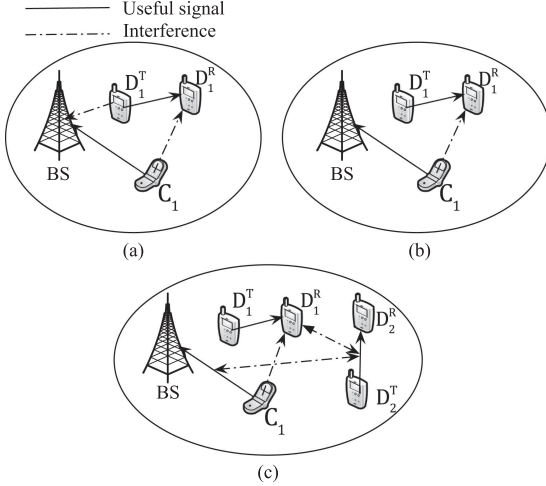


Fig. 2. The underlay mode and interlay mode for D2D in NOMA cellular network.

fading for each link. The channel power gain between cellular users C_m and BS is denoted by $h_{C_m,BS}$, which consists of the path loss, fast fading gain with exponential distribution, and shadow fading with log-normal distribution. Similarly, we denote the channel power gains of C_m - D_n^R , D_m^T - D_n^R , and D_n^T -BS by h_{C_m,D_n^R} , $h_{D_m^T,D_n^R}$, and $h_{D_n^T,BS}$, respectively. In addition, we assume that the BS can acquire the perfect channel state information on the PUCCH and PSCCH.²

NOMA uplink system applies the superposition coding technique at cellular users and adopts SIC technique at the BS. Specifically, multiple cellular users are allowed to transmit their signals to the BS by the same subcarrier and their signals have different power levels. In NOMA uplink system, the users' channel power gains are very different and the user's transmit power is very limited compared with the BS. Thus, the SIC decoder at the BS successively decodes the signals of the users in descending order of the users' channel power gains [20], [21]. Once a user's signal is decoded successfully, it is subtracted from the superposed signal. When one user's signal is being decoded, it treats the residual users' signals as interference. Hence, C_m only receives the co-channel interference from other C_i with $h_{C_m,BS} > h_{C_i,BS}$ in typical NOMA uplink system.

In this work, D2D pairs can operate in the underlay mode or interlay mode. The underlay mode is a traditional mode and has been applied in OMA system. The interlay mode is proposed in this paper and specialized for the NOMA system.

Underlay mode: In underlay mode, D2D pairs directly reuse the subcarriers of cellular users. However, the mutual interference is inevitable between the transmission links assigned with same subcarrier, as shown in Fig. 2a. Therefore, it is required to guarantee that the interference cannot break down the cellular uplink transmission. Similar to [5], [6], [22], we assume each subcarrier is reused by at most one D2D pair in the underlay mode.

However, from Fig. 2a, we observe that if D2D pairs locate in proximity to the BS, their communications would cause severe

mutual interference to cellular uplinks. Such severe interference cannot be relieved in OMA system and degrade the performance of D2D communications.

Interlay mode: The interlay mode is proposed for D2D communication in NOMA system. Different with the underlay mode, the interlay mode enables the power-domain multiplexing between cellular users and D2D pairs on the basis of frequency reuse. In the interlay mode, the D2D pair can join to the NOMA group to be rid of the mutual interference between it and NOMA cellular users. To be specific, if a D2D transmitter's channel power gain to the BS is better than that of cellular users, the BS can exploit SIC to cancel the interference from D2D pairs to cellular users, as is shown in Fig. 2b.

In this paper, we consider the existence of the underlay and interlay mode, as Fig. 2c shows. Specifically, because D_1 works in the interlay mode but D_1^T cannot interfere with C_1 . The cell-edge D_2 works in the underlay mode and reuses the same spectrum with C_1 and D_1 , since it causes weak interference to the BS. It is seen that due to the interference cancellation by the interlay D2D pair, a cell-edge D2D pair could reuse the cellular spectrum in underlay mode. Therefore, the existence of the underlay and interlay mode is adopted to improve the spectrum efficiency.

The signal-to-interference-plus-noise ratio (SINR) of D2D user D_n when using subcarrier S_k in interlay mode is expressed as

$$\xi_{D_n,S_k}^i = \frac{P_{D_n} h_{D_n^T,D_n^R}}{N_0 + I_{D_n,S_k}^{ic} + I_{D_n,S_k}^{id}} \quad (1)$$

$$I_{D_n,S_k}^{id} = \sum_{D_j \in \mathcal{I}_{D_n}^d} s_{D_j,S_k} P_{D_j} h_{D_j^T,D_n^R} + \sum_{D_j \in \mathcal{D}} x_{D_j,S_k} P_{D_j} h_{D_j^T,D_n^R} \quad (2)$$

where $\mathbf{P}_D = \{P_{D_n} \mid D_n \in \mathcal{D}\}$ is the vector of the transmit powers of D2D transmitters and N_0 is the additive white gaussian noise power. $I_{D_n,S_k}^{ic} = \sum_{C_i \in \mathcal{I}_{D_n}^c} y_{C_i,S_k} P_{C_i} h_{C_i,D_n^R}$ is the interference from cellular users, where $\mathbf{P}_C = \{P_{C_m} \mid C_m \in \mathcal{C}\}$ is the vector of the transmit powers of cellular users and $\mathcal{I}_{D_n}^c = \{C_i \in \mathcal{C} \mid h_{C_i,D_n^R} < h_{D_n^T,D_n^R}, D_n \in \mathcal{D}\}$. I_{D_n,S_k}^{id} is the interference from other D2D pairs and consists of two parts. The first part is the interference from interlay D2D pairs, where $\mathcal{I}_{D_n}^d = \{D_j \in \mathcal{D} \mid h_{D_j^T,D_n^R} < h_{D_n^T,D_n^R}, D_n \in \mathcal{D}\}$, and the second part is the interference from underlay D2D pairs. Note that the interlay D2D pair also can use SIC decoding to cancel the interference from other users in the same NOMA group. Let $\mathbf{S} = \{s_{D_n,S_k} \mid D_n \in \mathcal{D}, S_k \in \mathcal{K}\}$ and $\mathbf{X} = \{x_{D_n,S_k} \mid D_n \in \mathcal{D}, S_k \in \mathcal{K}\}$ denote the subcarrier assignment indicators for D2D pairs in the interlay mode and the underlay mode, respectively. $s_{D_n,S_k} = 1$ if D_n is assigned with S_k and works in interlay mode, otherwise $s_{D_n,S_k} = 0$. $x_{D_n,S_k} = 1$ if D_n is assigned with S_k and works in the underlay mode, otherwise $x_{D_n,S_k} = 0$. Similarly, the SINR of D2D pair D_n when using subcarrier S_k in underlay mode is expressed as

$$\xi_{D_n,S_k}^u = \frac{P_{D_n} h_{D_n^T,D_n^R}}{N_0 + I_{D_n,S_k}^{uc} + I_{D_n,S_k}^{ud}} \quad (3)$$

where $I_{D_n,S_k}^{uc} = \sum_{C_i \in \mathcal{C}} y_{C_i,S_k} P_{C_i} h_{C_i,D_n^R}$ and $I_{D_n,S_k}^{ud} = \sum_{D_j \in \mathcal{D}} s_{D_j,S_k} P_{D_j} h_{D_j^T,D_n^R}$ denote the interference from the cellular users and D2D pairs in interlay mode, respectively.

²Cellular users transmit the channel state information via physical uplink control channel (PUCCH). D2D pairs rely on physical sidelink control channel (PSCCH) [19]. In NOMA system, the SIC decoding information should be transmitted to the cellular users and D2D pairs on PUCCH and PSCCH.

The SINR of cellular user C_m when using subcarrier S_k is given by

$$\xi_{C_m, S_k} = \frac{P_{C_m} h_{C_m, BS}}{N_0 + I_{C_m, S_k}^c + I_{C_m, S_k}^d} \quad (4)$$

where $I_{C_m, S_k}^c = \sum_{C_i \in \mathcal{I}_{C_m}^c} y_{C_i, S_k} P_{C_i} h_{C_i, BS}$ is the interference from other cellular users. $\mathcal{I}_{C_m}^c = \{C_i \in \mathcal{C} \mid h_{C_i, BS} < h_{C_m, BS}\}$ is a subset of cellular users which can interfere with C_m . Let $\mathbf{Y} = \{y_{C_m, S_k} \mid C_m \in \mathcal{C}, S_k \in \mathcal{K}\}$ denote the subcarrier assignment indicators for cellular users, where $y_{C_m, S_k} = 1$ if S_k is assigned to C_m , otherwise $y_{C_m, S_k} = 0$. $I_{C_m, S_k}^d = \sum_{D_j \in \mathcal{I}_{C_m}^d} s_{D_j, S_k} p_{D_j} h_{D_j^T, BS} + \sum_{D_j \in \mathcal{D}} x_{D_j, S_k} p_{D_j} h_{D_j^T, BS}$ denotes the interference from D2D pairs using the same subcarrier, where $\mathcal{I}_{C_m}^d = \{D_j \in \mathcal{D} \mid h_{D_j^T, BS} < h_{C_m, BS}, C_m \in \mathcal{C}\}$ is a subset of D2D pairs interfering with C_m if they work in the interlay mode. Since the D2D is a flexible communication technique, this work focuses on the scheduling of D2D communications in the existing cellular network. Under the circumstance, \mathbf{Y} is known. Considering the decoding complexity and energy consumption, we denote d_f as the maximum number of the communication links in a NOMA group [12], [17]. In particular, $d_f = 1$ for OMA system.

Thus, the achievable rates of cellular users and D2D pairs are respectively given by

$$R_{C_m} = \sum_{k=1}^K y_{C_m, S_k} R_{C_m, S_k} = \sum_{k=1}^K y_{C_m, S_k} \log_2(1 + \xi_{C_m, S_k}) \quad (5)$$

$$\begin{aligned} R_{D_n} &= \sum_{k=1}^K (s_{D_n, S_k} R_{D_n, S_k}^i + x_{D_n, S_k} R_{D_n, S_k}^u) \\ &= \sum_{k=1}^K [s_{D_n, S_k} \log_2(1 + \xi_{D_n, S_k}^i) \\ &\quad + x_{D_n, S_k} \log_2(1 + \xi_{D_n, S_k}^u)] \end{aligned} \quad (6)$$

B. Problem Formulation

We formulate the maximization problem of system sum rate, which jointly optimizes the mode selection, subcarrier allocation, and power control. A D2D pair is activated only when the minimum rate requirement is satisfied and other communication links are not interrupted. The formulated optimization problem is as follows.

$$\begin{aligned} \max_{\mathbf{P}_C, \mathbf{P}_D, \mathbf{S}, \mathbf{X}} \quad & \sum_{m=1}^M R_{C_m} + \sum_{n=1}^N R_{D_n} \\ \text{s.t.} \quad & \text{C1: } \sum_{k=1}^K (s_{D_n, S_k} + x_{D_n, S_k}) \leq 1, \forall D_n \in \mathcal{D} \\ & \text{C2: } \sum_{m=1}^M \sum_{n=1}^N (y_{C_m, S_k} + s_{D_n, S_k}) \leq d_f, \forall S_k \in \mathcal{K} \\ & \text{C3: } \sum_{n=1}^N x_{D_n, S_k} \leq 1, \forall S_k \in \mathcal{K} \end{aligned}$$

$$\text{C4: } R_{C_m} \geq R_{\min}^c, \forall C_m \in \mathcal{C}$$

$$\text{C5: } R_{D_n} \geq R_{\min}^d \sum_{k=1}^K (x_{D_n, S_k} + s_{D_n, S_k}), \forall D_n \in \mathcal{D}$$

$$\text{C6: } 0 \leq P_{C_m} \leq P_{\max}^c, \forall C_m \in \mathcal{C}$$

$$\text{C7: } 0 \leq P_{D_n} \leq P_{\max}^d, \forall D_n \in \mathcal{D}$$

$$\text{C8: } s_{D_n, S_k}, x_{D_n, S_k} \in \{0, 1\}, \forall D_n \in \mathcal{D}, \forall S_k \in \mathcal{K} \quad (7)$$

where P_{\max}^c and P_{\max}^d are denoted as the maximum transmit powers for cellular users and D2D pairs, respectively. R_{\min}^c and R_{\min}^d are the minimum rate thresholds of cellular users and D2D pairs, respectively. C1 indicates that each D2D pair selects at most one communication mode. C2 ensures that each subcarrier can be shared with at most d_f cellular users and interlay D2D pairs. C3 ensures that each subcarrier can be reused by at most one underlay D2D pair. C4 and C5 specify the rate requirements of the cellular users and the D2D pairs, respectively. C6 and C7 restrict the maximum transmit power for cellular users and D2D pairs, respectively. C8 are the integer constraints for subcarrier assignment indicators.

III. GRAPH-BASED RESOURCE MANAGEMENT SCHEME

In this section, we first show how the original problem can be reformulated as a maximum weight clique problem with given power allocation. Then, based on the branch-and-bound (BnB) framework [23], we propose an efficient joint mode selection and subcarrier assignment algorithm which can obtain the optimal solution. Afterward, an iterative power control scheme is proposed. In addition, we propose the graph construction method to decrease the graph scale and exploit the special structure of constructed graph to reduce the complexity of the proposed algorithm.

A. Maximum Weight Clique Problem Reformulation

From Section II-B, (7) is a combinatorial optimization problem which includes the integer variables and non-convex functions. Since its computational complexity increases exponentially with the growing number of integer variables, it is very difficult to solve (7) directly. Therefore, we exploit the graph theory to reformulate (7) into a maximum weight clique problem with given feasible power allocation [24].

In order to reformulate (7), we firstly define an undirected graph by $G = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{v_1, \dots, v_i, \dots, v_{|\mathcal{V}|}\}$ is a set of vertices, and $\mathcal{E} = \{e_1, \dots, e_j, \dots, e_{|\mathcal{E}|}\}$ is the set of edges. $|\cdot|$ is the cardinality of a set, i.e., the number of elements in a set. Moreover, let $w(v_i)$ denote the weight of $v_i \in \mathcal{V}$. For a subset of vertices $\mathcal{V}' \in \mathcal{V}$, the weight of \mathcal{V}' is the summation of the weight of vertices in \mathcal{V}' , $w(\mathcal{V}') = \sum_{v_i \in \mathcal{V}'} w(v_i)$. In what follows, we give the definitions of clique and maximum weight clique problem.

Definition 1: In an undirected graph, a *clique* \mathcal{Q} is a subset of vertices in which any pair of vertices is adjacent (connected by an edge).

Definition 2: Maximum weight clique (MWC) problem is to find a clique \mathcal{Q} which has the maximum weight $w(\mathcal{Q})$ among all cliques.

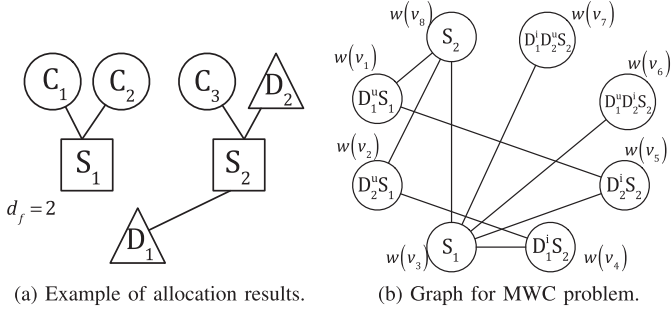


Fig. 3. An example of graph modeling for MWC problem.

Each vertex v_i represents a combination $v_i = (\mathcal{D}_{v_i}^i, \mathcal{D}_{v_i}^u, S_{v_i})$ of D2D pairs and subcarriers, which corresponds to a result of mode selection and subcarrier assignment. $\mathcal{D}_{v_i}^i$ and $\mathcal{D}_{v_i}^u$ are the sets of interlay D2D pairs and underlay D2D pairs using S_{v_i} , respectively. S_{v_i} is the subcarrier included in v_i and can be equivalent to any subcarrier, $S_{v_i} = S_k, S_k \in \mathcal{K}$.

There exists an edge (v_i, v_j) between a pair of vertices, if and only if they do not include the same D2D pairs and same subcarriers. In addition, for each vertex, $w(v_i)$ equals the maximum rate achieved by this combination. It is expressed as

$$w(v_i) = \sum_{D_n \in \mathcal{D}_{v_i}^i} R_{D_n, S_k}^i + \sum_{D_n \in \mathcal{D}_{v_i}^u} R_{D_n, S_k}^u + \sum_{C_m \in \mathcal{C}} y_{C_m, S_k} R_{C_m, S_k} \quad (8)$$

where the rate of each link can be calculated based on given power allocation. Note that if the vertex cannot satisfy C4, C5, C6, and C7 in (7), its weight equals zero. Actually, we can utilize the feature of D2D communication modes to reduce the number of vertices. Hence, we get the following lemma.

Lemma 1: $\forall \mathcal{D}' \subset \mathcal{D}$, if $\sum_{C_m \in \mathcal{C}} y_{C_m, S_k} + |\mathcal{D}'| \leq d_f$, we get

$$w(v_i) \geq w(v_j) \quad (9)$$

where $\mathcal{D}_{v_i}^i \setminus \mathcal{D}' = \emptyset$ and $(\mathcal{D}_{v_j}^i \cup \mathcal{D}_{v_j}^u) \setminus \mathcal{D}' = \emptyset$.

Proof: See the Appendix A. ■

Lemma 1 demonstrates that for a group of D2D pairs, when all of them work in the interlay mode, they can achieve higher rate than that when a part of them work in underlay mode. Thus, there should exist interlay D2D pairs as many as possible on the same subcarrier. Thus, based on Lemma 1, we further provide the following corollary.

Corollary 1: In G , since $w(v_i) \geq w(v_j)$, it is unnecessary to include v_j in MWC. Hence, it does not need to create any v_j for any \mathcal{D}' .

For example, as shown in Fig. 3a, there are three cellular users, two D2D pairs, and two subcarriers, where S_1 is shared by C_1 and C_2 , S_2 is used by C_3 , and $d_f = 2$ for each subcarrier. D_1 and D_2 are to be allocated. The graph model in Fig. 3b includes eight vertices. Specifically, v_1 represents that D_1 is in the underlay mode and assigned with S_1 , which can coexist with v_5 or v_8 . v_6 represents the allocation result shown in Fig. 3a. S_2 is assigned with D_1 and D_2 , which are in the underlay mode and interlay mode, respectively. In addition, v_3 and v_8 indicate there is not any D2D pair using S_1 and S_2 , respectively. In the worst case, the boundary of the number of vertices is $K d_f \frac{N!}{d_f!(N-d_f)!}$. Due

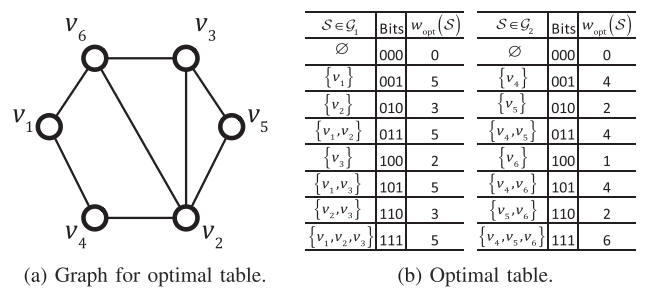


Fig. 4. An example of optimal table.

to the limitation of SIC capability, the value of d_f is generally very limited in practice [10].

Theorem 1: The maximization problem (7) can be equivalently considered as the MWC problem if the feasible power allocation is given.

Proof: See the Appendix B. ■

The MWC problem is a classic problem in graph theory and combinatorial optimization [25]. Although it is also NP-hard, we can exploit the special feature of the constructed graph model to propose a BnB based algorithm.

B. BnB Based Algorithm

The BnB based algorithm consists of preprocessing phase and BnB phase. In preprocessing phase, the vertices are divided into several groups corresponding to different subgraphs. Then, the weights of MWC in different subgraphs are calculated and stored in the optimal table [26]. In BnB phase, the MWC problem is divided into smaller subproblems, in which the upper bound of the weight can be calculated by optimal tables. The optimal upper bound is updated if the subproblem can achieve higher upper bound than current optimal upper bound. Moreover, the unnecessary subproblems are pruned by the optimal tables, which can reduce the search space. Furthermore, we can exploit the gap of vertex weight to further avoid the calculation of unnecessary subproblems.

1) Preprocessing Phase: We firstly use an example to introduce the definition and function of the optimal table. Fig. 4a shows that there are six vertices in the graph \mathcal{G} and their weights $w(v_i) = \{5, 3, 2, 4, 2, 1\}$. The vertices is divided into two groups, $\mathcal{G}_1 = \{v_1, v_2, v_3\}$ and $\mathcal{G}_2 = \{v_4, v_5, v_6\}$. Fig. 4b shows that the optimal table of \mathcal{G} consists of the vertex subsets $S \in \mathcal{G}_1$ and the weight $w_{\text{opt}}(S)$ of the vertex subset's MWC. Each vertex subset is represented by the bit vector with the length of $|S|$. By means of the optimal tables, we can calculate the upper bound U of any vertex subset with $O(1)$. For example, the upper bound of $\mathcal{V}' = \{v_1, v_3, v_5\}$ is calculated as follow.

$$U(\mathcal{V}') = \text{table}[1][101] + \text{table}[2][010] = 7. \quad (10)$$

For BnB phase, the tight upper bound can reduce the complexity and improve the efficiency. The tightness of upper bound is impacted by the number of vertex groups, which is shown in following Lemma 2.

Lemma 2: For $G = (\mathcal{V}, \mathcal{E})$, the vertices are divided into l groups, $\Omega = \{\mathcal{G}_1, \dots, \mathcal{G}_l\}$. We can obtain that $U(\mathcal{V}') \leq l \times w_{\text{opt}}(\mathcal{V}'), \forall \mathcal{V}' \subseteq \mathcal{V}$.

Proof: See the Appendix C. ■

Algorithm 1: Preprocessing Phase.

```

1: Initialization
2: • Set  $\mathcal{U} = \mathcal{V}$ , maximum size for vertex group as  $r$ ;
3: • Set  $t = |\mathcal{V}|$ ,  $s = 0$ ,  $l = 1$  and  $k = 1$ ;
4: • Set  $\mathcal{S}_1 = \emptyset$ ,  $\mathcal{S} = \{\mathcal{S}_1\}$  and  $w_{\text{opt}}(\mathcal{S}_1) = 0$ ;
5: Stage 1: Vertex partitioning
6: while  $\mathcal{U} \neq \emptyset$  do
7:   Set  $s = s + 1$ ,  $\mathcal{I}_s = \emptyset$  and  $\mathcal{U}' = \mathcal{U}$ ;
8:   while  $\mathcal{U}' \neq \emptyset$  and  $|\mathcal{I}_s| < r$  do
9:     Select the vertex  $v_i \in \mathcal{U}'$  with maximum weight;
10:    Set  $v_{\pi(t)} = v_i$ ,  $\mathcal{I}_s = \mathcal{I}_s \cup v_{\pi(t)}$  and  $t = t - 1$ ;
11:     $\mathcal{U}' = \mathcal{U}' \setminus (\{v_i\} \cup \mathcal{N}(v_i))$ ,  $\mathcal{U} = \mathcal{U} \setminus \{v_i\}$ ;
12:   end while
13: end while
14:  $\mathcal{G}_l = \emptyset$ ;
15: for  $i$  from  $s$  to 1 do
16:   if  $|\mathcal{G}_l| + |\mathcal{I}_i| \leq r$  then
17:      $\mathcal{G}_l = \mathcal{G}_l \cup \mathcal{I}_i$ ;
18:   else
19:      $l = l + 1$  and  $\mathcal{G}_l = \mathcal{I}_i$ ;
20:   end if
21: end for
22: Stage 2: Optimal table calculation
23: for all  $\mathcal{G}_l$  do
24:   for all  $u \in \mathcal{G}_l$  do
25:      $\mathcal{S}' = \emptyset$ ;
26:     for  $j$  from 1 to  $|\mathcal{S}|$  do
27:        $k = k + 1$ ,  $\mathcal{S}_k = \mathcal{S}_j \cup u$ ;
28:        $\delta = w(u) + w_{\text{opt}}(\mathcal{S}_j \cap \mathcal{N}(u))$ ;
29:        $w_{\text{opt}}(\mathcal{S}_k) = \max(\delta, w_{\text{opt}}(\mathcal{S}_j))$ ;
30:        $\mathcal{S}' = \mathcal{S}' \cup \mathcal{S}_k$ ;
31:     end for
32:      $\mathcal{S} = \mathcal{S} \cup \mathcal{S}'$ ;
33:   end for
34: end for
35: return  $\mathcal{V}_\pi = \{v_{\pi(1)}, \dots, v_{\pi(|\mathcal{V}|)}\}$  and  $w_{\text{opt}}(\mathcal{S}_k), \forall \mathcal{S}_k \in \mathcal{G}_l, \forall \mathcal{G}_l \in \Omega$ 

```

Lemma 2 reveals that to acquire tight upper bound, the number of groups $|\Omega|$ should be as small as possible [27]. However, since the small $|\Omega|$ brings high computational complexity of optimal table, the appropriate size of each group should be set according to available memory and computational capability of the computer. Thus, we denote the maximum size of each group by r . Moreover, if the vertex group is a independent set, the maximum weight clique is equivalent to the maximum weight vertex [27]. It helps us in getting $w_{\text{opt}}(\mathcal{S})$ and constructing the optimal table easily. Thus, the vertex coloring method is adopted in vertex partitioning.

We design the preprocessing phase algorithm shown in Algorithm 1. The Algorithm 1 first divides the vertices into many small groups and determines the branching order of vertices, and then constructs the optimal tables of the vertex groups. As shown in Algorithm 1, all vertices are partitioned into many independent sets \mathcal{I}_s s and get the vertex sequence. From steps 6 to 13, the vertex v_i with maximum weight in \mathcal{U}' is preferentially selected to join in \mathcal{I}_s and then is given the sequence number $\pi(t)$. Then, the vertex with maximum weight is selected to join in

Algorithm 2: BnB Phase.

```

1: Initialization
2: • Construct  $G = \{\mathcal{V}, \mathcal{E}\}$ ,  $\mathcal{V}_\pi$ ,  $\Omega$  and optimal tables;
3: • Set  $\omega(\cdot) = 0$ ,  $\mathcal{Q}_{\text{max}} = \emptyset$  and  $\bar{\mathcal{V}} = \emptyset$ ;
4: for  $i$  from 1 to  $|\mathcal{V}|$  do
5:   if  $\forall v_j \in \mathcal{V}_{i-1}$ ,  $w(v_{\pi(i)}) < w(v_j)$ ,  $S_{v_i} = S_{v_j}$  and
      $(\mathcal{D}_{v_i}^i \cup \mathcal{D}_{v_i}^u) \setminus (\mathcal{D}_{v_j}^i \cup \mathcal{D}_{v_j}^u) = \emptyset$  then
6:      $\bar{\mathcal{V}} = \bar{\mathcal{V}} \cup \mathcal{V}_i$ ;
7:   else
8:      $\mathcal{V}' = \mathcal{V}_i \setminus \bar{\mathcal{V}}$ ;
9:      $\mathcal{Q} = \emptyset$ ;
10:     $(\mathcal{Q}_{\text{max}}) = \text{RMWCS}(\mathcal{V}', \mathcal{Q}_{\text{max}}, \mathcal{Q})$ ;
11:     $\omega(i) = w(\mathcal{Q}_{\text{max}})$ ;
12:   end if
13: end for
14: return  $\mathcal{Q}_{\text{max}}, w(\mathcal{Q}_{\text{max}})$ 

```

\mathcal{I}_s in the residual vertices except v_i 's neighbor vertices $\mathcal{N}(v_i)$. Until \mathcal{I}_s is maximal or its size achieves r , a new independent set will be created. The principle of vertex partitioning is the vertex coloring method [24], [27]. From 15 to 21, in order to improve the tightness of upper bound, some consecutive independent sets \mathcal{I}_i s are merged into a new vertex group \mathcal{G}_l unless the number of vertices in \mathcal{G}_l is more than r . The stage 2 is to calculate the optimal table of each vertex group. Step 27 shows that the selected vertex u joins in each current vertex subset \mathcal{S}_j to create new subset \mathcal{S}_k . At step 28, the optimal value δ of MWC including u in \mathcal{S}_k is calculated, where $w_{\text{opt}}(\mathcal{S}_j \cap \mathcal{N}(u))$ has been stored in optimal table. Then, if δ is great than or equal to the optimal value $w_{\text{opt}}(\mathcal{S}_j)$, $w_{\text{opt}}(\mathcal{S}_k) = \delta$. Otherwise, we get $w_{\text{opt}}(\mathcal{S}_k) = w_{\text{opt}}(\mathcal{S}_j)$. At steps 30 and 32, \mathcal{S}' is used to store new subsets and contained in \mathcal{S} . The outputs of Algorithm 1 are vertex sequence $\mathcal{V}_\pi = \{v_{\pi(1)}, \dots, v_{\pi(|\mathcal{V}|)}\}$ and the optimal tables of all \mathcal{G}_l , respectively.

2) *BnB Phase:* In BnB phase, we divide G into many subgraphs and each subgraph consisting of $\mathcal{V}_i = \{v_{\pi(1)}, v_{\pi(2)}, \dots, v_{\pi(i)}\}$ corresponds to a subproblem. The MWC of subproblem can be found through recursive approach, which separates the current subproblem into some new subproblems. To reduce recursive procedure, we propose the following theorem to prune unnecessary subproblem.

Theorem 2: In $G = (\mathcal{V}, \mathcal{E})$, there exist v_i and v_j with $(\mathcal{D}_{v_i}^i \cup \mathcal{D}_{v_i}^u) \setminus (\mathcal{D}_{v_j}^i \cup \mathcal{D}_{v_j}^u) = \emptyset$ and $S_{v_i} = S_{v_j}$. if $w(v_i) > w(v_j)$, $v_j \notin \mathcal{Q}_{\text{max}}$.

Proof: First, we assume that v_j is contained in MWC \mathcal{Q}_{max} . Due to $(\mathcal{D}_{v_i}^i \cup \mathcal{D}_{v_i}^u) \setminus (\mathcal{D}_{v_j}^i \cup \mathcal{D}_{v_j}^u) = \emptyset$, v_i and v_j include identical D2D pairs. Hence, both of them cannot coexist in \mathcal{Q}_{max} . Because $w(v_i) > w(v_j)$, \mathcal{Q}_{max} has heavier weight if v_i is added to \mathcal{Q}_{max} instead of v_j , which means that a better MWC is found and contradicts the initial assumption for \mathcal{Q}_{max} . Therefore, v_j cannot be an element in \mathcal{Q}_{max} . ■

Theorem 2 reveals that for two vertices which contain the same subcarrier and D2D pairs but different mode selections, the one with smaller weight must not be in the maximum weight clique. Thus, it can directly prune the subproblems including the vertices with excessively small weight.

The BnB phase is shown in detail in Algorithm 2. In Algorithm 2, we select the vertices in order of \mathcal{V}_π to construct the

Algorithm 3: Recursively MWC solution (RMWCS).

```

1: if  $\mathcal{V}' = \emptyset$  then
2:   if  $w(\mathcal{Q}) > w(\mathcal{Q}_{max})$  then
3:      $\mathcal{Q}_{max} = \mathcal{Q}$ ;
4:   end if
5:   return  $\mathcal{Q}_{max}, \mathcal{Q}$ 
6: end if
7: if  $U(\mathcal{V}') + w(\mathcal{Q}) \geq w(\mathcal{Q}_{max})$  and
    $\omega(m(\mathcal{V}')) + w(\mathcal{Q}) \geq w(\mathcal{Q}_{max})$  then
8:    $\mathcal{Q} = \mathcal{Q} \cup v_m(\mathcal{V}')$ ;
9:    $(\mathcal{Q}_{max}, \mathcal{Q}) = \mathbf{RMWCS}(\mathcal{V}' \cap \mathcal{N}(v_m(\mathcal{V}')),$ 
    $\mathcal{Q}_{max}, \mathcal{Q})$ ;
10:   $\mathcal{Q} = \mathcal{Q} \setminus v_m(\mathcal{V}')$ ;
11:   $(\mathcal{Q}_{max}, \mathcal{Q}) = \mathbf{RMWCS}(\mathcal{V}' \setminus v_m(\mathcal{V}'), \mathcal{Q}_{max}, \mathcal{Q})$ ;
12: end if
13: return  $\mathcal{Q}_{max}, \mathcal{Q}$ 

```

subgraph which consists of \mathcal{V}_i and its optimal solution is stored in $\omega(i)$. $m(\mathcal{V}_i)$ represents the maximum index of vertices in \mathcal{V}_i . At step 5, we employ Theorem 2 to prune the unnecessary problems including the vertices with small weight. These vertices are contained in $\bar{\mathcal{V}}$ and not allowed to join in \mathcal{V}_i . At step 10, we use Algorithm 3 to solve MWC in subgraph. Specifically, from steps 1 to 6 of Algorithm 3, if a subproblem is completely solved and a better solution is found, \mathcal{Q}_{max} will be updated. Step 7 corresponds to the bounding procedure, where the upper bound of the subproblem is obtained by using optimal table or exploiting $\omega(m(\mathcal{V}_i))$. If this upper bound is sufficiently small, this subproblem will be pruned. From steps 8 to 11, the vertex $v_m(\mathcal{V}_i)$ is selected as new branching variable, which generates new leaf nodes of BnB framework. Then, two subproblems are solved for optimal solution recursively, the one subproblem including $v_m(\mathcal{V}_i)$ and the other one excluding $v_m(\mathcal{V}_i)$. Finally, Algorithm 2 outputs optimal MWC \mathcal{Q}_{max} of G and optimal $w(\mathcal{Q}_{max})$. The vertices in \mathcal{Q}_{max} are optimal solution and $w(\mathcal{Q}_{max})$ is the optimal value of (7).

C. Power Allocation

In previous subsection, we have adopted MWC method to solve (7) with given feasible power allocation. In this subsection, we introduce a power allocation scheme for each vertex in $G = (\mathcal{V}, \mathcal{E})$. The power allocation subproblem for the vertex v_i can be expressed as

$$\begin{aligned}
& \max_{\mathbf{P}=[\mathbf{P}_C, \mathbf{P}_D]} w(v_i) \\
& \text{s.t. C9: } R_{C_m, S_k} \geq R_{\min}^c, \forall C_m \in \mathcal{C}_k \\
& \quad \text{C10: } R_{D_n, S_k} \geq R_{\min}^d, \forall D_n \in \mathcal{D}_{v_i}^i \cup \mathcal{D}_{v_i}^u \\
& \quad \text{C6, C7} \tag{11}
\end{aligned}$$

where $\mathcal{C}_{S_k} = \{C_m \mid y_{C_m, S_k} = 1, \forall C_m \in \mathcal{C}\}$. Since the interference exists in the summation of rate expression, the objective function has strong non-convexity. Thus, we exploit sequential convex approximation method to get the near-optimal solution of (11).

The objective function can be rewritten as

$$w(v_i) = f(\mathbf{P}) - h(\mathbf{P}) \tag{12}$$

where $f(\mathbf{P})$ and $h(\mathbf{P})$ are given by

$$\begin{aligned}
f(\mathbf{P}) &= \sum_{C_m \in \mathcal{C}_{S_k}} \log_2 (P_{C_m} h_{C_m, BS} + N_0 + I_{C_m, S_k}^c + I_{C_m, S_k}^d) \\
&+ \sum_{D_n \in \mathcal{D}_{v_i}^i} \log_2 (P_{D_n} h_{D_n^T, D_n^R} + N_0 + I_{D_n, S_k}^{ic} + I_{D_n, S_k}^{id}) \\
&+ \sum_{D_n \in \mathcal{D}_{v_i}^u} \log_2 (P_{D_n} h_{D_n^T, D_n^R} + N_0 + I_{D_n, S_k}^{uc} + I_{D_n, S_k}^{ud}) \tag{13}
\end{aligned}$$

and

$$\begin{aligned}
h(\mathbf{P}) &= \sum_{C_m \in \mathcal{C}_{S_k}} \log_2 (N_0 + I_{C_m, S_k}^c + I_{C_m, S_k}^d) \\
&+ \sum_{D_n \in \mathcal{D}_{v_i}^i} \log_2 (N_0 + I_{D_n, S_k}^{ic} + I_{D_n, S_k}^{id}) \\
&+ \sum_{D_n \in \mathcal{D}_{v_i}^u} \log_2 (N_0 + I_{D_n, S_k}^{uc} + I_{D_n, S_k}^{ud}). \tag{14}
\end{aligned}$$

Since the objective function is the difference between two concave functions, (11) can be rewritten as a standard problem of the D.C. (difference between two convex function) programming [28], [29]. We exploit the first order approximation to transform (12) into a convex function. The gradient of $h(\mathbf{P})$ at the point \mathbf{P} is expressed by

$$\nabla h(\mathbf{P}) = \left(\frac{\partial h}{\partial P_{C_1}}, \dots, \frac{\partial h}{\partial P_{C_m}}, \frac{\partial h}{\partial P_{D_1}}, \dots, \frac{\partial h}{\partial P_{D_n}} \right) \tag{15}$$

where if $U \in \mathcal{C}_{S_k} \cup \mathcal{D}_{v_i}^i$, the partial derivative of h with respect to P_U is presented by

$$\begin{aligned}
\frac{\partial h}{\partial P_U} &= \frac{1}{\ln 2} \sum_{C_m \in \mathcal{J}_U^c} \frac{h_{U, BS}}{N_0 + I_{C_m, S_k}^c + I_{C_m, S_k}^d} \\
&+ \frac{1}{\ln 2} \sum_{D_n \in \mathcal{J}_U^{id}} \frac{h_{U, D_n^R}}{N_0 + I_{D_n, S_k}^{ic} + I_{D_n, S_k}^{id}} \\
&+ \frac{1}{\ln 2} \sum_{D_n \in \mathcal{D}_{v_i}^u} \frac{h_{*, D_n^R}}{N_0 + I_{C_m, S_k}^{uc} + I_{C_m, S_k}^{ud}} \tag{16}
\end{aligned}$$

where $\mathcal{J}_U^c = \{C_m \mid h_{U, BS} \geq h_{C_m, BS}, C_m \in \mathcal{C}_{S_k}\}$ and $\mathcal{J}_U^{id} = \{D_n \mid h_{U, D_n^R} \geq h_{D_n^T, D_n^R}, D_n \in \mathcal{D}_{v_i}^i\}$. If $U \in \mathcal{D}_{v_i}^u$, $\frac{\partial h}{\partial P_U}$ is given by

$$\begin{aligned}
\frac{\partial h}{\partial P_U} &= \frac{1}{\ln 2} \sum_{C_m \in \mathcal{C}_{S_k}} \frac{h_{U, BS}}{N_0 + I_{C_m, S_k}^c + I_{C_m, S_k}^d} \\
&+ \frac{1}{\ln 2} \sum_{D_n \in \mathcal{D}_{v_i}^i} \frac{h_{U, D_n^R}}{N_0 + I_{D_n, S_k}^{ic} + I_{D_n, S_k}^{id}}. \tag{17}
\end{aligned}$$

Accordingly, we approximate $h(\mathbf{P})$ by its first order Taylor expansion $h(\mathbf{P}^k) + \nabla h^T(\mathbf{P})(\mathbf{P} - \mathbf{P}^k)$ at arbitrary \mathbf{P}^k . Therefore, (11) can be transformed to

$$\begin{aligned}
& \max_{\mathbf{P}=[\mathbf{P}_C, \mathbf{P}_D]} f(\mathbf{P}) - h(\mathbf{P}^k) - \nabla h^T(\mathbf{P}^k)(\mathbf{P} - \mathbf{P}^k) \\
& \text{s.t. C6, C7, C9, C10} \tag{18}
\end{aligned}$$

which is a convex problem and can be directly solved by existing algorithms, e.g. interior point method or standard solver, such as CVX [30]. We propose an iteratively based vertex power allocation (IVPA) scheme, which is described in Algorithm 4. Each

Algorithm 4: Iteratively Based Vertex Power Allocation (IVPA) Algorithm.

- 1: **Initialization**
 - 2: • Set \mathbf{P}^0 and calculate $w^0 = f(P^0) - h(P^0)$;
 - 3: • Set $k = 0$, $\varepsilon = 10^{-2}$ and $\varphi = 1000$;
 - 4: **repeat**
 - 5: Solve problem (18) to obtain \mathbf{P}^* ;
 - 6: Calculate $w^k = f(\mathbf{P}^*) - h(\mathbf{P}^*)$;
 - 7: Update $k = k + 1$ and $\mathbf{P}^k = \mathbf{P}^*$;
 - 8: **until** $|w^k - w^{k-1}| \leq \varepsilon$ or $k \geq \varphi$
 - 9: **return** $w(v_i) = w^k$, \mathbf{P}^k
-

Algorithm 5: MWC Based Mode Selection and Resource Allocation (MWC-MSRA) Algorithm.

- 1: **Initialization**
 - 2: • Construct the undirected graph $G = (\mathcal{V}, \mathcal{E})$;
 - 3: Calculate the vertex weight of G by Algorithm 4;
 - 4: Obtain Ω and \mathcal{V}_π from Algorithm 1;
 - 5: Complete mode selection and subcarrier assignment from Algorithm 2;
-

vertex includes at most $d_f + 1$ cellular users and D2D pairs. The maximum number of iterations is φ and the complexity of each iteration is $O((d_f + 1)^{3.5})$ by the interior point method. Thus, the complexity of Algorithm 4 is $O(\varphi(d_f + 1)^{3.5})$.

Finally, we conclude the entire procedure of MWC based mode selection and resource allocation (MWC-MSRA) scheme in Algorithm 5. Note that although BnB method is an efficient approach to solve NP-hard combinatorial problem, the worst-case complexity is same with the exhaustive searching problem. If BnB method cannot prune any subproblem, it will search all possibilities. Thus, it is required to design a low-complexity resource allocation algorithm.

IV. A FAST HEURISTIC RESOURCE ALLOCATION ALGORITHM

In this section, we develop a heuristic resource allocation algorithm with low complexity. The main idea is to activate as many D2D pairs as possible, especially interlay D2D pairs, to improve the system sum rate. The details of the interference-aware heuristic resource allocation (IHRA) algorithm are presented in Algorithm 6.

The first step is to select D2D pairs working in interlay mode and assign them with the subcarriers. During this step, we assume that all users are assigned with unit transmit power and define two utility functions which are the interference utility $\alpha_{n,k}$ and rate utility $\bar{R}_{n,k}^i$, respectively. Specifically, $\alpha_{n,k}$ presents a ratio of the interference when D_n works in underlay mode to the interference when D_n works in interlay mode, and is expressed as

$$\alpha_{n,k} = \frac{I_{n,k}^u}{I_{n,k}^i} \quad (19)$$

$$\begin{aligned} I_{n,k}^u &= \sum_{C_m \in \mathcal{C}_k} (h_{D_n^T, C_m} + h_{C_m, D_n^R}) \\ &+ \sum_{D_j \in \mathcal{D}} s_{D_n, S_k} (h_{D_n^T, D_j^R} + h_{D_j^T, D_n^R}) \end{aligned} \quad (20)$$

Algorithm 6: Interference-Aware Heuristic Resource Allocation (IHRA) Algorithm.

- 1: **Initialization**
 - 2: • Initialize all $s_{D_n, S_k} = 0$, $x_{D_n, S_k} = 0$, $P_{C_m} = 1$, $P_{D_m} = 1$;
 - 3: • Set $\mathcal{D}_0 = \mathcal{D}$, $\mathcal{S}_0 = \{S_k \mid \sum_{C_m \in \mathcal{C}} y_{C_m, S_k} < d_f\}$;
 - 4: **repeat**
 - 5: all $\alpha_{n,k} = 0$;
 - 6: Calculate each $\alpha_{n,k} = I_{n,k}^u / I_{n,k}^i$ and each $\bar{R}_{n,k}^i$, $\forall D_n \in \mathcal{D}_0, \forall S_k \in \mathcal{S}_0$;
 - 7: **if** $\exists (D_n, S_k), \alpha_{n,k} > 1$ **then**
 - 8: Get $\mathcal{P}^* = \{(D_n, S_k) \mid (D_n, S_k) = \arg \max_{n \in \mathcal{D}_0, k \in \mathcal{S}_0} \alpha_{n,k}\}$;
 - 9: **if** $|\mathcal{P}^*| == 1$ **then**
 - 10: Set $s_{D_n, S_k} = 1$, $(D_n, S_k) \in \mathcal{P}^*$ and $\mathcal{D}_0 = \mathcal{D}_0 \setminus D_n$;
 - 11: **else**
 - 12: Select $(D_n, S_k) = \arg \max_{\mathcal{P}^*} \bar{R}_{n,k}^i$;
 - 13: Set $s_{D_n, S_k} = 1$, $\mathcal{D}_0 = \mathcal{D}_0 \setminus D_n$;
 - 14: **end if**
 - 15: **end if**
 - 16: $\mathcal{S}_0 = \{S_k \mid \sum_{C_m \in \mathcal{C}} \sum_{D_n \in \mathcal{D}} y_{C_m, S_k} + s_{D_n, S_k} < d_f\}$;
 - 17: **until** $\mathcal{D}_0 = \emptyset$ or $\mathcal{S}_0 = \emptyset$
 - 18: Use Algorithm 4 to obtain each $w(e_{n,k}), \forall D_n \in \mathcal{D}_0, \forall S_k \in \mathcal{S}$;
 - 19: Construct Matching model and get x_{D_n, S_k} by Kuhn-Munkres Algorithm;
 - 20: **return** $s_{D_n, S_k}, x_{D_n, S_k}, P_{C_m}, P_{D_n}$
-

$$\begin{aligned} I_{n,k}^i &= \sum_{C_m \in \mathcal{C}_k} (h_{D_n^T, C_m} \mathbf{1}_{D_n \in \mathcal{I}_{D_m}^d} + h_{C_m, D_n^R} \mathbf{1}_{C_m \in \mathcal{I}_{D_n}^c}) \\ &+ \sum_{D_j \in \mathcal{D}} s_{D_n, S_k} (h_{D_n^T, D_j^R} \mathbf{1}_{D_n \in \mathcal{I}_{D_j}^d} + h_{D_j^T, D_n^R} \mathbf{1}_{D_j \in \mathcal{I}_{D_n}^d}) \end{aligned} \quad (21)$$

where $\mathbf{1}_f = 1$ if condition f is satisfied, and $\mathbf{1}_f = 0$ otherwise. $\bar{R}_{n,k}^i$ is the sum rate of users when D_n works in interlay mode, which is expressed as $\bar{R}_{n,k}^i = \sum_{C_m \in \mathcal{C}_k} \bar{R}_{C_m} + \sum_{D_i \in \mathcal{D}} \bar{R}_{D_i, S_k}^i + \bar{R}_{D_n, S_k}^i$, where each rate value is calculated by unit transmit power. From Line 7 to 11, we calculate $\alpha_{n,k}$ and $\bar{R}_{n,k}^i$ for every pair of D_n and S_k . Then, two criteria are designed based on these two utilities to determine the allocation order. The second step is to assign the underlay D2D pairs to the subcarriers, which is a maximum-weight matching problem [5], [31]. In this matching model, the set of unassigned D2D pairs \mathcal{D}_0 and the set of subcarriers \mathcal{S} are two disjoint groups of vertices in the bipartite graph. There is an edge $e_{n,k}$ to connect D2D pair D_n with subcarrier S_k , when D_n reuses S_k and cannot break other transmission links on S_k . The weight of $e_{n,k}$ is given by

$$\begin{aligned} w(e_{n,k}) &= R_{D_n, S_k}^u + \sum_{D_i \in \mathcal{D}} s_{D_i, S_k} R_{D_i, S_k}^i \\ &+ \sum_{C_n \in \mathcal{C}} y_{C_m, S_k} R_{C_m, S_k} \end{aligned} \quad (22)$$

where s_{D_i, S_k} and y_{C_m, S_k} is the subcarrier assignment indicators for interlay D2D pairs and cellular users. We can use proposed

TABLE I
SIMULATION PARAMETERS

Parameters	Values
Cell radius	200m
P_{\max}^d	24 dBm
R_{\min}^c, R_{\min}^d	$\log_2(1 + 10)$ bps/Hz
Subcarrier Bandwidth	180 kHz
Noise power, N_0	-174 dBm/Hz \times 180kHz
Path loss model	$37.6 \log_{10}(d[\text{km}]) + 128.1$
Multiple-path fading	Exponential distribution with unit mean
Shadowing	Log-normal distribution with standard deviation of 8 dB

power control algorithm (Algorithm 4) to calculate $w(e_{n,k})$. The maximum-weight matching is a classic problem in graph theory and can be solved by Kuhn-Munkres algorithm [24]. The complexity of first step (lines 4 to 17) is $O(NK)$. The complexity of the calculation of edge weight is $O(NK(d_f + 1)^{3.5})$. The complexity of Kuhn-Munkres algorithm is $O(\max\{K^3, N^3\})$. Thus, the computational complexity of IHRA algorithm equals to $O(\max\{K^3, N^3\} + NK(d_f + 1)^{3.5})$.

V. SIMULATION RESULTS

In this section, we present simulation results to evaluate the performance of the interlay mode and proposed algorithms. We consider a single cell network, where cellular users and D2D pairs are uniformly distributed in cell. The link lengths of D2D pairs are uniformly distributed, where maximum link length is L . The main parameters are summarized in Table I. Limited by our computer's resources, the maximum size for vertex group r is set to be 8. We compare the proposed coexisting scheme of the interlay mode and underlay mode with the conventional underlay scheme [5], [17] and the underlay with SIC scheme [32].³ As [32], the underlay with SIC scheme enables the SIC at both of BS and D2D receivers. Specifically, at each subcarrier, BS can cancel the interference for cellular users which have worse channel gain to BS than corresponding D2D transmitter. The D2D receiver can cancel the interference from the cellular user that has the better channel gain to the D2D receiver than the D2D transmitter and other cellular users using same subcarrier. Moreover, the performance of three benchmark algorithms has also been evaluated for the comparison.

- **Exhaustive search (ES) algorithm:** This algorithm searches all possible results of resource allocation in the constructed graph, which can obtain the optimal solution and is used to investigate the computational efficiency of based MWC resource allocation algorithm.
- **Greedy MWC (GMWC) algorithm:** It is the GGWMIN algorithm derived from [33], [34], which utilizes an efficient utility function to find MWC of constructed graph with low complexity.
- **Random allocation (Random) algorithm:** The algorithm randomly allocates communication modes and subcarriers for all D2D pairs. In order to ensure the fairness, IVPA

³In underlay scheme and underlay with SIC scheme, D2D pairs utilize only the underlay mode. Each D2D pair reuses at most one subcarrier and each subcarrier is reused by at most one D2D pair.

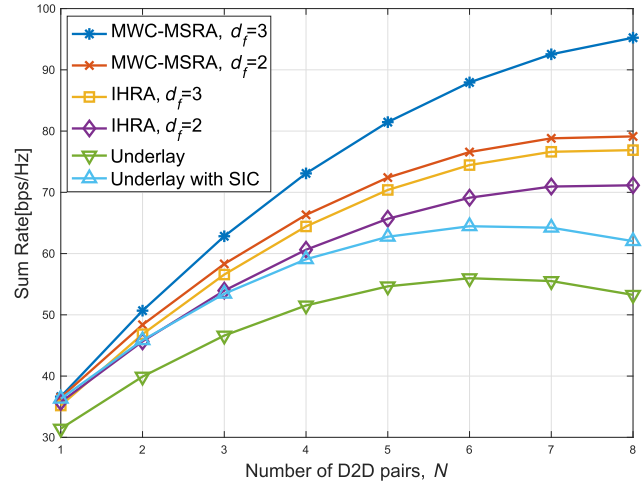


Fig. 5. System sum rate versus the number of D2D pairs ($M = 6$, $K = 4$, $P_{\max}^c = P_{\max}^d = 24$ dBm, $L = 10$ m).

algorithm is adopted to allocate the users transmit power, in which the users sharing the same subcarrier constitute one vertex.

A. Comparison With ES

Table II shows the performance comparison between ES algorithm and MWC-MSRA algorithm with the different number of the subcarriers, cellular users, and D2D pairs. In each cell of the columns of ES and MWC-MSRA, the first value denotes the computation time⁴ used by ES algorithm or MWC-MSRA algorithm, and the second value denotes the corresponding system sum rate. We use MATLAB operated on a computer with 2.2GHz CPU, 8GB memory, and Unix operating system. It can be seen that MWC-MSRA algorithm can achieve the optimal solution and significantly reduce the computation time especially when N and K increase. This is mainly because MWC-MSRA algorithm can exploit the upper bound of subproblem and utilize the gap between the weight of vertices (Theorem 2) to prune the unnecessary subproblems. Furthermore, we can see that for each K , the computation time increases with N and decreases with M . This is because the number of vertices determines the complexity of MWC problems. More D2D pairs lead more allocation results, which increases the number of vertices. However, more cellular users make the number of available resources decline, which reduces the number of vertices.

B. Impact of the Number of D2D Pairs on Network Performance

Fig. 5 shows the sum rate of different schemes versus the number of D2D pairs N . It is observed that the interlay+underlay scheme greatly improves the system sum rate compared with other schemes. This is because the interlay mode can exploit SIC capability to reduce the interference, while the underlay scheme

⁴The computation time does not include the time of graph construction process to emphasize the computational complexity of algorithm execution. This is because both of ES and MWC-MSRA algorithms use the same constructed graph which spends uniform time.

TABLE II
COMPARISON WITH ES ALGORITHM

K	M, N	ES [sec., bps/Hz]	MWC-MSRA [sec., bps/Hz]	Percent of reduction
$K = 2$	$M = 2, N = 2$	0.0257, 33.1604	0.0374, 33.1604	-45%
	$M = 3, N = 2$	0.0235, 24.9812	0.0212, 24.9812	10.7%
	$M = 3, N = 3$	0.0711, 29.2595	0.0252, 29.2595	64.63%
$K = 3$	$M = 4, N = 3$	76.832, 50.0329	0.0685, 50.0329	99.91%
	$M = 5, N = 3$	1.1296, 42.5734	0.0366, 42.5734	96.76%
	$M = 5, N = 4$	555.58, 47.2088	0.0696, 47.2088	99.99%
$K = 4$	$M = 6, N = 3$	1281.3, 54.1118	0.1274, 54.1118	99.99%
	$M = 7, N = 3$	19.743, 55.2671	0.0800, 55.2671	99.59%
	$M = 7, N = 4$	> 13100, -	0.1857, 57.6542	99.99%

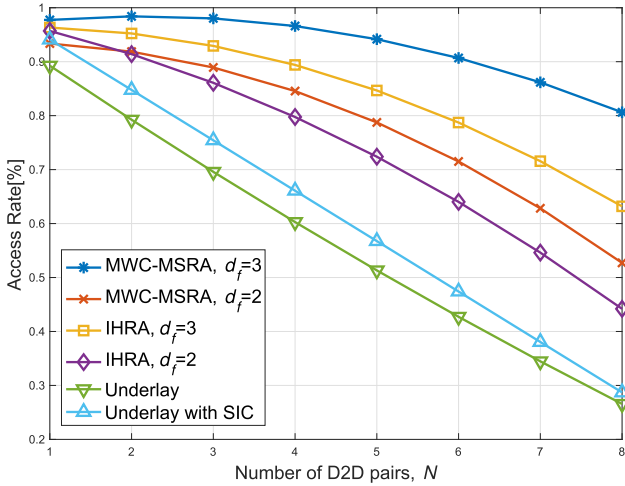


Fig. 6. Access rate of D2D pairs versus the number of D2D pairs ($M = 6$, $K = 4$, $P_{\max}^c = 24$ dBm, $L = 10$ m).

only use power control to restrain the interference. We also observe that because the spectrum resources are deeply reused, the interlay+underlay scheme outperforms the underlay with SIC scheme. Furthermore, we find that when $d_f = 3$, MWC-MSRA and IHRA algorithm can achieve higher sum rate than that when $d_f = 2$. The reason is that when $d_f = 3$, NOMA system can allow more D2D pair working in interlay mode.

Fig. 6 shows the access rate of D2D pairs versus the number of D2D pairs N . We see that with the increase of N , the proposed algorithms achieve higher access rate than other schemes. The reasons are twofold. On the one hand, the coexisting scheme of the interlay mode and underlay mode provides more opportunities for D2D pairs to access to the network. On the other hand, since the interference is efficiently restrained, the rate requirement of each user is more easily satisfied. In addition, it can be seen that MWC-MSRA algorithm achieves higher access rate than other algorithms. This is because taking the advantage of graph model, MWC-MSRA algorithm can more efficiently relieve interference and exploit user diversities to conduct mode selection and subcarrier assignment.

C. Impact of Link Length of D2D Pairs on Network Performance

Fig. 7 shows how the link length of D2D pairs L affects the system sum rate with different algorithms. It can be seen that MWC-MSRA algorithm achieves higher system capacity than other algorithms. This is because it can exploit the graph model for MWC problem to depict the interference relationship and

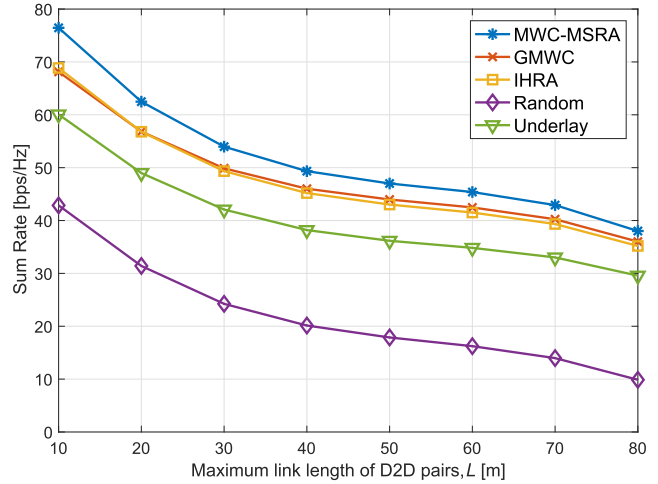


Fig. 7. System sum rate versus link length of D2D pairs ($M = 6$, $K = 4$, $N = 6$, $P_{\max}^c = 24$ dBm, $d_f = 2$).

utilize BnB framework to find optimal solution with low complexity. We also see that the performance of IHRA algorithm is close to that of GMWC algorithm, even though IHRA algorithm does not utilize the graph model. It is because IHRA algorithm can fully release the proximity gain and assign D2D pairs with the appropriate mode causing less interference. From the performance loss of Random algorithm, we can see that although interlay mode can promote network connection and capacity, the efficient resource management algorithm is required.

Fig. 8 shows the cumulative probability functions (CDFs) of the average rate of D2D pairs of different schemes with $L = 10$ m and 50 m. From the results, we see that with the same L , MWC-MSRA algorithm achieves higher average rate of D2D pairs than the underlay scheme. The reason is that MWC-MSRA algorithm can not only leverage interlay mode to eliminate the interference but also utilize graph model and BnB framework to allocate D2D pairs appropriate communication modes and resources. We also see that no matter what schemes are used, the schemes with $L = 10$ m outperform the schemes with $L = 50$ m. It is because that the increasing L degrades link robustness and weakens D2D proximity gain. Thus, we obtain that the proximity gain is the key factor to affect the performance of D2D. Furthermore, with the decrease of L , the gap of D2D performance obviously enlarge. Therefore, the interlay mode can fully exploit D2D proximity feature to improve spectral efficiency.

Fig. 9 plots the CDFs of average rate of cellular users obtained by different schemes with $L = 10$ m and 50 m. It can be seen that MWC-MSRA algorithm achieves higher average rate of cellular

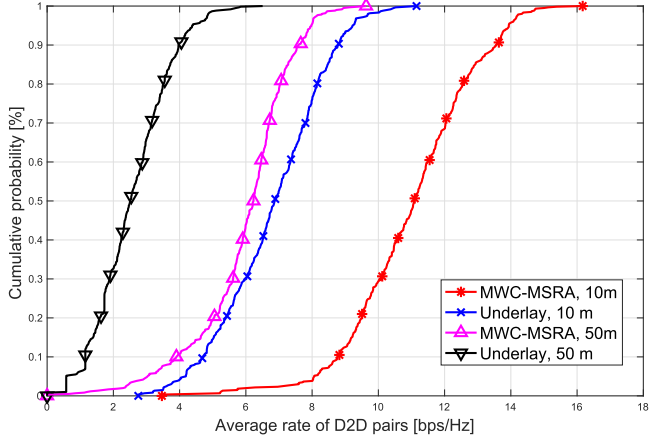


Fig. 8. Cumulative distribution function of average rate of D2D pairs ($M = 6$, $K = 4$, $N = 6$, $P_{\max}^c = 24$ dBm, $d_f = 2$).

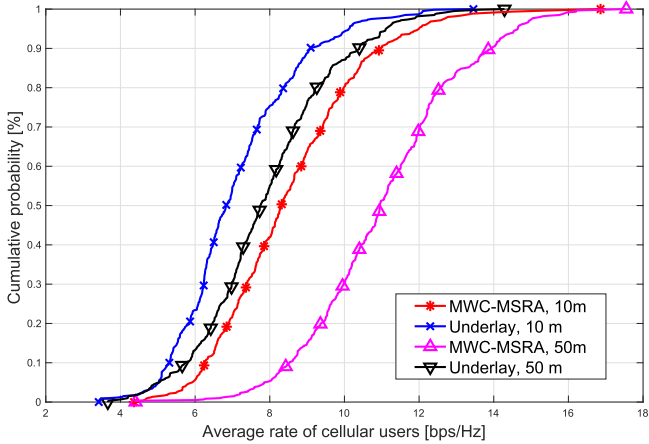


Fig. 9. Cumulative distribution function of average rate of cellular users ($M = 6$, $K = 4$, $N = 6$, $P_{\max}^c = 24$ dBm, $d_f = 2$).

users than underlay scheme. It is because MWC-MSRA algorithm can take the advantage of graph model for MWC problem and interlay mode to guarantee QoS of cellular users. It is worthwhile noticed that with L increasing from 10m to 50m, the average rate achieved by both algorithms is improved. It is because the decline of D2D link length results in that cellular users are allocated with better resources to keep the system throughput. Furthermore, the gap between two curves based on MWC-MSRA algorithm is greater than that based on underlay scheme. Since MWC-MSRA algorithm can find the optimal allocation results, it improves the spectral efficiency of cellular users to guarantee system sum rate when the link quality of D2D pairs is degraded. However, due to applying the maximum matching method, the underlay scheme would activate some D2D pairs with poor link quality to result in network overloading.

D. Impact of Maximum Transmit Power on Network Performance

Fig. 10 shows the system sum rate versus P_{\max}^c of different algorithm. From the results, we see that with the increase P_{\max}^c , the system sum rate obtained by all algorithm gradually increases.

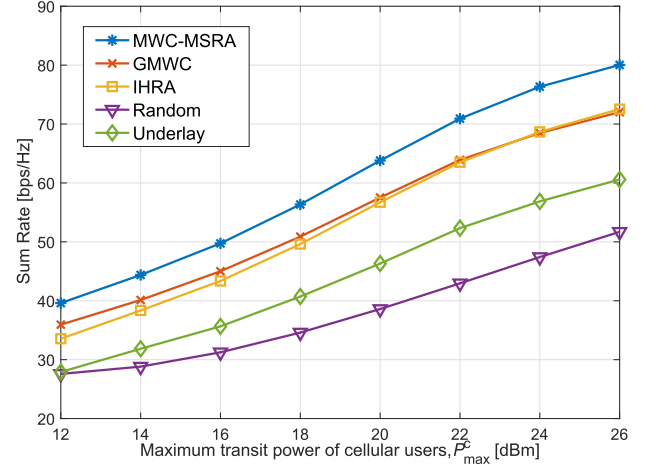


Fig. 10. System sum rate versus maximum transmit power of cellular users ($M = 6$, $K = 4$, $N = 6$, $L = 10$ m, and $d_f = 2$).

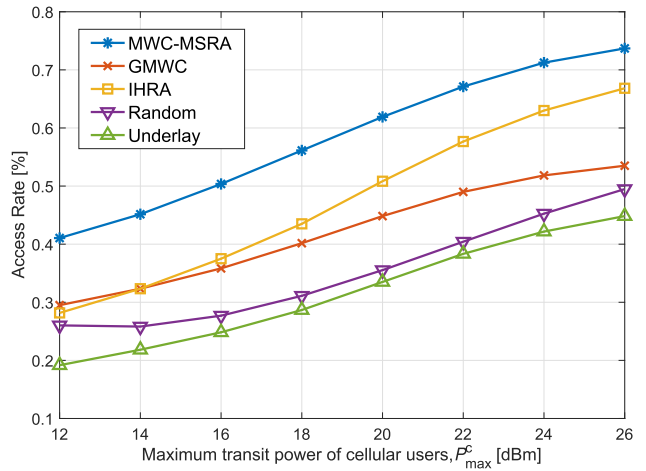


Fig. 11. Access rate of D2D pairs versus maximum transmit power of cellular users ($M = 6$, $K = 4$, $N = 6$, $L = 10$ m, and $d_f = 2$).

The reason is that the higher P_{\max}^c not only makes cellular users more easily achieve R_{\min}^c and but also improves the achievable rate of cellular users. In addition, Fig. 11 shows that higher P_{\max}^c can also release the potential of D2D communications, which enables more D2D pairs to access to the network. It is also seen that MWC-MSRA algorithm can achieve higher sum rate and access rate than other algorithms. This is mainly because it assigns each D2D pair optimal mode and subcarrier, exploiting the constructed graph model to accurately characterize the interference relationship via different allocation results. Furthermore, the IHRA algorithm achieves relatively high sum rate, which is much better than Random scheme. The reason is that IHRA efficiently utilizes the utility functions to select suitable communication mode for D2D pairs, which restrains the interference between cellular users and D2D pairs. On the other hand, the Random algorithm achieves lower sum rate but higher access rate than underlay scheme. This indicates that although interlay mode can activate more D2D pairs, it cannot obtain the potential benefit on spectral efficiency if without effective resource management scheme.

VI. CONCLUSION

In this paper, we have studied the mode selection and resource allocation issue for D2D-enabled NOMA cellular networks. The interlay mode is introduced for D2D communications in NOMA system, which exploits the SIC to cancel the interference between D2D pairs and cellular users. Accordingly, aiming to maximize the system sum rate, we have developed MWC-MSRA algorithm for resource management. MWC-MSRA algorithm is based on graph theory and utilizes BnB method to obtain the optimal solution. To further reduce complexity, IHRA algorithm has been proposed, which can utilize the proximity gain and interference relationship to improve system capacity. Simulation results have shown that the interlay mode can achieve higher D2D access rate than the conventional underlay mode and the proposed algorithms can efficiently improve spectral efficiency.

For the future work, we will extend the considered scenario to the ultra dense network. Due to the large number of BSs, the inter-cell interference should be considered and the effectively distributed scheme should be designed with reasonable complexity.

APPENDIX

A. Proof of Lemma 1

For S_k , if $\sum_{C_m \in \mathcal{C}} y_{C_m, S_k} + |\mathcal{D}'| \leq d_f$, it is valid that there exists v_i representing $(\mathcal{D}', \emptyset, S_k)$ and its weight is calculated by

$$w(v_i) = \sum_{D_n \in \mathcal{D}'} R_{D_n, S_k}^i + \sum_{C_m \in \mathcal{C}} y_{C_m, S_k} R_{C_m, S_k}. \quad (23)$$

For v_j , its combination is expressed as $(\mathcal{D}_{v_j}^i, \mathcal{D}_{v_j}^u, S_k)$ and its weight is calculated by

$$w(v_j) = \sum_{D_n \in \mathcal{D}_{v_j}^i} R_{D_n, S_k}^i + \sum_{D_n \in \mathcal{D}_{v_j}^u} R_{D_n, S_k}^u + \sum_{C_m \in \mathcal{C}} y_{C_m, S_k} \hat{R}_{C_m, S_k} \quad (24)$$

where we denote the rate of cellular users in v_j by \hat{R}_{C_m, S_k} to distinguish the difference rate expressions of cellular users in v_i and v_j . Since the interlay mode can alleviate the interference between D2D pairs and cellular users, $\sum_{C_m \in \mathcal{C}} R_{C_m, S_k} \geq \sum_{C_m \in \mathcal{C}} \hat{R}_{C_m, S_k}$. Besides, the interlay mode can restrain the interference for interlay D2D pairs. The rate achieved by interlay mode is at least the same as that achieved by underlay mode. Therefore, we obtain that $\sum_{D_n \in \mathcal{D}'} R_{D_n, S_k}^i \geq \sum_{D_n \in \mathcal{D}_{v_j}^i} R_{D_n, S_k}^i + \sum_{D_n \in \mathcal{D}_{v_j}^u} R_{D_n, S_k}^u$. Accordingly, we prove that $w(v_i) \geq w(v_j)$.

B. Proof of Theorem 1

If the feasible optimal power allocation has been given, C6 and C7 is satisfied. For v_i , the definition of its weight $w(v_i)$ meets C4, C5, C6, and C7. Meanwhile, $w(v_i)$ represents the sum rate of D2D pairs and cellular users in v_i under given power allocation. Because each vertex represents the unique combination, the vertices are independent for each other. Therefore, the independence and connection between vertices guarantee C1,

C2, C3 and C8 in the resulting maximum weight clique \mathcal{Q}_{\max} . Accordingly, $w(\mathcal{Q}_{\max})$ is just equivalent to the objective function in (7).

C. Proof of Lemma 2

According to the optimal tables, we can obtain the value of $w_{\text{opt}}(\mathcal{V}' \cap \mathcal{G}_n)$. Thus, we prove the inequality directly as follows.

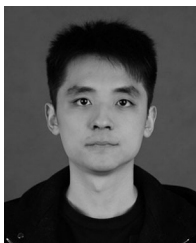
$$U(\mathcal{V}') = \sum_{n=1}^l w_{\text{opt}}(\mathcal{V}' \cap \mathcal{G}_n) \leq \sum_{n=1}^l w_{\text{opt}}(\mathcal{V}') = l \times w_{\text{opt}}(\mathcal{V}'). \quad (25)$$

The inequality has been proved.

REFERENCES

- [1] M. Agiwal, A. Roy, and N. Saxena, "Next generation 5G wireless networks: A comprehensive survey," *IEEE Commun. Surveys Tut.*, vol. 18, no. 3, pp. 1617–1655, Jul.–Sep. 2016.
- [2] A. Asadi, Q. Wang, and V. Mancuso, "A survey on device-to-device communication in cellular networks," *IEEE Commun. Surveys Tut.*, vol. 16, no. 4, pp. 1801–1819, Oct.–Dec. 2014.
- [3] N. Cheng *et al.*, "Performance analysis of vehicular device-to-device underlay communication," *IEEE Trans. Veh. Technol.*, vol. 66, no. 6, pp. 5409–5421, Jun. 2017.
- [4] H. Sun, M. Wildemeersch, M. Sheng, and T. Q. S. Quek, "D2D enhanced heterogeneous cellular networks with dynamic TDD," *IEEE Trans. Wireless Commun.*, vol. 14, no. 8, pp. 4204–4218, Aug. 2015.
- [5] D. Feng, L. Lu, Y. Y. Wu, G. Y. Li, G. Feng, and S. Li, "Device-to-device communications underlying cellular networks," *IEEE Trans. Commun.*, vol. 61, no. 8, pp. 3541–3551, Aug. 2013.
- [6] G. Yu, L. Xu, D. Feng, R. Yin, G. Y. Li, and Y. Jiang, "Joint mode selection and resource allocation for device-to-device communications," *IEEE Trans. Commun.*, vol. 62, no. 11, pp. 3814–3824, Nov. 2014.
- [7] L. Lei, Y. Kuang, N. Cheng, X. Shen, Z. Zhong, and C. Lin, "Delay-optimal dynamic mode selection and resource allocation in device-to-device communications—Part II: Practical algorithm," *IEEE Trans. Veh. Technol.*, vol. 65, no. 5, pp. 3491–3505, May 2016.
- [8] J. Liu, Y. Kawamoto, H. Nishiyama, N. Kato, and N. Kadowaki, "Device-to-device communications achieve efficient load balancing in LTE-advanced networks," *IEEE Wireless Commun.*, vol. 21, no. 2, pp. 57–65, Apr. 2014.
- [9] H. Nishiyama, M. Ito, and N. Kato, "Relay-by-smartphone: Realizing multi-hop device-to-device communications," *IEEE Commun. Mag.*, vol. 52, no. 4, pp. 56–65, Apr. 2014.
- [10] Z. Ding, X. Lei, G. K. Karagiannidis, R. Schober, J. Yuan, and V. K. Bhargava, "A survey on non-orthogonal multiple access for 5G networks: Research challenges and future trends," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 10, pp. 2181–2195, Oct. 2017.
- [11] Z. Ding *et al.*, "Application of non-orthogonal multiple access in LTE and 5G networks," *IEEE Commun. Mag.*, vol. 55, no. 2, pp. 185–191, Feb. 2017.
- [12] L. P. Qian, Y. Wu, H. Zhou, and X. Shen, "Joint uplink base station association and power control for small-cell networks with non-orthogonal multiple access," *IEEE Trans. Wireless Commun.*, vol. 16, no. 9, pp. 5567–5582, Sep. 2017.
- [13] J. Liu, M. Sheng, L. Liu, and J. Li, "Interference management in ultra-dense networks: Challenges and approaches," *IEEE Netw.*, vol. 31, no. 6, pp. 70–77, Nov. 2017.
- [14] D. Zhai, R. Zhang, L. Cai, B. Li, and Y. Jiang, "Energy-efficient user scheduling and power allocation for NOMA based wireless networks with massive IoT devices," *IEEE Internet Things J.*, vol. 5, no. 3, pp. 1857–1868, Jun. 2018.
- [15] L. Lv, Q. Ni, Z. Ding, and J. Chen, "Application of non-orthogonal multiple access in cooperative spectrum-sharing networks over Nakagami- m fading channels," *IEEE Trans. Veh. Technol.*, vol. 66, no. 6, pp. 5506–5511, Jun. 2017.

- [16] B. Di, L. Song, Y. Li, and Z. Han, "V2X meets NOMA: Non-orthogonal multiple access for 5G-enabled vehicular networks," *IEEE Wireless Commun.*, vol. 24, no. 6, pp. 14–21, Dec. 2017.
- [17] Y. Pan, C. Pan, Z. Yang, and M. Chen, "Resource allocation for D2D communications underlying a NOMA-based cellular network," *IEEE Wireless Commun. Lett.*, vol. 7, no. 1, pp. 130–133, Feb. 2018.
- [18] J. Zhao, Y. Liu, K. K. Chai, Y. Chen, and M. ElKashlan, "Joint subchannel and power allocation for NOMA enhanced D2D communications," *IEEE Trans. Commun.*, vol. 65, no. 11, pp. 5081–5094, Nov. 2017.
- [19] *3rd Generation Partnership Project; Technical Specification Group RAN; Study on LTE Device to Device Proximity Services (ProSe)—Radio Aspects (Release 12)*, 3GPP TR 36.843 V1.0.0, Nov. 2013.
- [20] Z. Yang, Z. Ding, P. Fan, and N. Al-Dhahir, "A general power allocation scheme to guarantee quality of service in downlink and uplink NOMA systems," *IEEE Trans. Wireless Commun.*, vol. 15, no. 11, pp. 7244–7257, Nov. 2016.
- [21] B. Xia, J. Wang, K. Xiao, Y. Gao, Y. Yao, and S. Ma, "Outage performance analysis for the advanced SIC receiver in wireless NOMA systems," *IEEE Trans. Veh. Technol.*, vol. 67, no. 7, pp. 6711–6715, Jul. 2018.
- [22] T. D. Hoang, L. B. Le, and T. Le-Ngoc, "Resource allocation for D2D communication underlaid cellular networks using graph-based approach," *IEEE Trans. Wireless Commun.*, vol. 15, no. 10, pp. 7099–7113, Oct. 2016.
- [23] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [24] D. B. West, *Introduction to Graph Theory*, vol. 2. Upper Saddle River, NJ, USA: Prentice-Hall, 2001.
- [25] P. R. J. Östergård, "A new algorithm for the maximum-weight clique problem," *Nordic J. of Comput.*, vol. 8, no. 4, pp. 424–436, Dec. 2001.
- [26] S. Shimizu, K. Yamaguchi, T. Saitoh, and S. Masuda, "Some Improvements on Kumlander's Maximum Weight Clique Extraction Algorithm," in *Proc. Int. Conf. Electr., Comput., Electron. Commun. Eng. (ICECECE)*, 2012, pp. 307–311.
- [27] S. Shimizu, K. Yamaguchi, T. Saitoh, and S. Masuda, "Fast maximum weight clique extraction algorithm: Optimal tables for branch-and-bound," *Discrete Appl. Math.*, vol. 223, pp. 120–134, 2017.
- [28] D. T. Ngo, S. Khakurel, and T. Le-Ngoc, "Joint subchannel assignment and power allocation for OFDMA femtocell networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 1, pp. 342–355, Jan. 2014.
- [29] H. Shan, Y. Zhang, W. Zhuang, A. Huang, and Z. Zhang, "User behavior-aware scheduling based on time-frequency resource conversion," *IEEE Trans. Veh. Technol.*, vol. 66, no. 9, pp. 8429–8444, Sep. 2017.
- [30] M. Grant, S. Boyd, and Y. Ye, "CVX: Matlab software for disciplined convex programming," Apr. 2011. [Online]. Available: <http://cvxr.com/cvx/>
- [31] Y. Gu, W. Saad, M. Bennis, M. Debbah, and Z. Han, "Matching theory for future wireless networks: Fundamentals and applications," *IEEE Commun. Mag.*, vol. 53, no. 5, pp. 52–59, May 2015.
- [32] C. Ma, W. Wu, Y. Cui, and X. Wang, "On the performance of successive interference cancellation in D2D-enabled cellular networks," in *Proc. IEEE Conf. Comput. Commun.*, Apr. 2015, pp. 37–45.
- [33] S. Sakai, M. Togasaki, and K. Yamazaki, "A note on greedy algorithms for the maximum weighted independent set problem," *Discrete Appl. Math.*, vol. 126, no. 2, pp. 313–322, 2003.
- [34] D. Zhai and J. Du, "Spectrum efficient resource management for multi-carrier based NOMA networks: A graph-based method," *IEEE Wireless Commun. Lett.*, vol. 7, no. 3, pp. 388–391, Jun. 2018.



Yanpeng Dai received the B.Eng degree in communication engineering from Shandong Normal University, Jinan, China, in 2014. He is currently working toward the Ph.D degree with the State Key Laboratory of Integrated Service Networks, Xidian University, Xi'an, China. From September 2017 to September 2018, he was a Visiting Student with the University of Waterloo, Waterloo, ON, Canada. His research interests include resource management in heterogeneous cellular networks.



Scholar from Ministry of Education, China, respectively.

Min Sheng (M'03–SM'16) received the M.S. and Ph.D. degrees in communication and information systems from Xidian University, Xi'an, China, in 2000 and 2004, respectively. She is currently a Full Professor and the Director of the State Key Laboratory of Integrated Service Networks, Xidian University. Her general research interests include mobile ad hoc networks, 5G mobile communication systems, and satellite communications networks. She was awarded as a Distinguished Young Researcher from National Natural Science Foundation of China and Changjiang



Junyu Liu (S'15–M'17) received the B.Eng. and Ph.D. degrees in communication and information systems from Xidian University, Xi'an, China, in 2007 and 2016, respectively. He is currently a Postdoctoral Researcher with the State Key Laboratory of Integrated Service Networks, Institute of Information and Science, Xidian University. His research interests include interference management and performance evaluation of wireless heterogeneous networks and ultra-dense wireless networks.



MAC, opportunistic communication for vehicular networks, unmanned aerial vehicles, and artificial intelligence for wireless networks.

Nan Cheng (M'16) received the B.E. and M.S. degrees from Tongji University, and the Ph.D. degree from the University of Waterloo, Waterloo, ON, Canada. He is currently a Joint Professor with the School of Telecommunication, Xidian University, Xi'an, China. He is also a Joint Postdoctoral Fellow with the Department of Electrical and Computer Engineering, University of Toronto, Toronto, ON, Canada, and with the Department of Electrical and Computer Engineering, University of Waterloo. His research interests include performance analysis,



security, social networks, smart grids, and vehicular ad hoc and sensor networks. He is an Engineering Institute of Canada Fellow, a Canadian Academy of Engineering Fellow, a Royal Society of Canada Fellow, and a Distinguished Lecturer of the IEEE Vehicular Technology Society and Communications Society. He received the 2018 James Evans Avant Garde Award from IEEE Vehicular Technology Society and the IEEE Canada R.A. Fessenden Silver Medal Award in 2019. He is the Editor-in-Chief of the IEEE INTERNET OF THINGS JOURNAL. He is the Vice President of publications of the IEEE Communications Society.

Xuemin (Sherman) Shen (M'97–SM'02–F'09) received the B.Sc. degree in electrical engineering from Dalian Maritime University, Dalian, China, in 1982, and the M.Sc. and Ph.D. degrees in electrical engineering from Rutgers University, New Brunswick, NJ, USA, in 1987 and 1990, respectively. He is currently a university Professor and an Associate Chair for graduate studies with the Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON, Canada. His research interests include resource management, wireless network



delivery networks, and LTE-A techniques. He received University President Award while pursuing his Ph.D. degree.

Qinghai Yang received the B.S. degree in communication engineering from the Shandong University of Technology, Zibo, China, in 1998, the M.S. degree in information and communication systems from Xidian University, Xi'an, China, in 2001, and Ph.D. degree in communication engineering from Inha University, Incheon, South Korea, in 2007. From 2007 to 2008, he was a Research Fellow with UWB-ITRC, South Korea. Since 2008, he has been a Full Professor with Xidian University, Xi'an, China. His current research interests include autonomous communication, content