

Optimization-Based Resource Allocation and Transmission Scheduling for Wireless Networks

by

Shiwei Huang

A Thesis submitted to the Faculty of Graduate Studies of
The University of Manitoba
in partial fulfillment of the requirements of the degree of

DOCTOR OF PHILOSOPHY

Department of Electrical and Computer Engineering
University of Manitoba
Winnipeg

© Copyright 2018 by Shiwei Huang

Supervisor: **Prof. Jun Cai**

Abstract

Future wireless communication networks are expected to be more energy-efficient and to provide higher throughput, in order to satisfy the demands for the increasing number of mobile users. Resource allocation and transmission scheduling play more and more important roles in improving the performance of wireless networks, in terms of energy saving, throughput, delay, etc. In this thesis, we consider three networks with different characteristics and objectives, i.e., wireless relay networks for distant transmissions, dense multi-user coexisting networks, and device-to-device (D2D) assisted mobile edge computing systems for compute-intensive mobile applications. We aim to investigate the key resource allocation and/or transmission scheduling issues in these networks. In particular, i) a transmit power allocation scheme with reduced overheads for amplify-and-forward relay networks is proposed to reduce energy consumption, based on the two-stage stochastic programming method, ii) an analysis framework for buffer-aided decode-and-forward relay networks under time-correlated fading channels is developed and an improved link scheduling/selection policy is presented, through the analyses to two quasi-birth-death Markov chains, iii) an interference-avoidance scheduling scheme for dense multi-user coexisting networks with heterogeneous priorities and demands is presented to increase the number of admitted users, on the basis of the column generation method, and iv) a joint optimization of admission control, link scheduling, and resource management for D2D-assisted mobile edge computing is carried out, according to the branch-and-price method. Simulations are performed

to verify the effectiveness of the proposed schemes where the performance of networks is shown to be improved significantly.

Keywords: Wireless relay networks, multi-user coexistence, mobile edge computing, resource allocation, power control, link scheduling, two-stage stochastic programming, quasi-birth-death Markov chains, column generation, branch-and-price.

Acknowledgments

First and foremost, I would like to express my deepest gratitude to my advisor Prof. Jun Cai for his continuous support during my Ph.D. study. Prof. Cai taught me focusing on the key points of issues in the research and he helped me revise papers very patiently and carefully. Besides, I would like to thank my examining committee members: Prof. Pradeepa Yahampath, Prof. Yunhua Luo, and Prof. Xavier Fernando for their valuable feedback and constructive suggestions. In addition, I would also like to express my sincere thanks to my previous supervisor Prof. Hongbin Chen for his moral encouragement and invaluable support. All of these are very helpful in finishing my research work and the writing of this thesis.

I thank all my colleagues in the Network Intelligence and Innovation (NI2) Lab for their insightful comments and suggestions during our discussions. My sincere thanks go also to the teaching, administrative and technical staff of the department of electrical and computer engineering, the University of Manitoba, which provides a peaceful and productive working environment.

Last but not the least, I must express my profound gratitude to my wife for her continuous support, constant encouragement and genuine dedications, and to my parents, sister, and brother who are always behind me. Without them, this thesis would not have been complete. I also appreciate the continued encouragement that my dear friends gave me, and meanwhile, I am happy and excited for the arrival of my newborn baby. I am truly lucky to have them in my life.

I acknowledge the University of Manitoba and the Government of Manitoba for awarding me a Manitoba Graduate Scholarship (MGS) offering financial support.

This thesis is dedicated to my family.

Contents

Abstract	i
Acknowledgments	iii
Dedication	iv
Table of Contents	vii
List of Figures	viii
List of Tables	x
List of Abbreviations	xi
1 Introduction	1
1.1 Wireless Relay Networks	2
1.1.1 Power Allocation in Amplify-and-Forward Relaying	3
1.1.2 Enhanced Decode-and-Forward Relaying	5
1.2 Dense Multi-User Coexisting Networks	8
1.2.1 Medium Access Control	8
1.2.2 Priority-Aware Interference-Avoidance Scheduling	10
1.3 Mobile Edge Computing Systems	14
1.3.1 Mobile Edge Computing	14
1.3.2 D2D-Assisted Mobile Edge Computing	15
1.4 Contributions	18
1.5 Organization of the Thesis	22
2 Transmit power allocation for amplify-and-forward relay networks with reduced overheads	23
2.1 Related Works	24
2.2 Network model and Problem formulation	28
2.2.1 Strategy I	29
2.2.2 Strategy II	32
2.2.3 Strategy III	32
2.3 Optimal power control	35
2.3.1 Strategy I	35
2.3.2 Strategy II	39
2.3.3 Strategy III	42

2.4	Simulation Results	45
2.4.1	Effectiveness of the Proposed Algorithms	46
2.4.2	Comparisons among Three Strategies in Terms of Total Power Consumption	49
2.4.3	Distinctiveness of Strategy III and its Applicable Area	51
3	An Analysis Framework for Buffer-Aided Relaying under Time-Correlated Fading Channels	55
3.1	System Model	57
3.1.1	Correlated Fading Channel Model	58
3.1.2	Link Scheduling Policy	60
3.1.3	Queue at The Relay	63
3.2	Queueing Behavior under Policy I: Infinite Buffer Size	64
3.3	Queueing Behavior under Policy II: Finite Buffer Size	71
3.4	Numerical and Simulation Results	75
3.5	Guidelines on Performance Improvement under Correlated Fading	80
3.5.1	Performance Degradation due to Channel Correlations	81
3.5.2	An Improved Policy	82
3.5.3	Unknown Channel Outage State Information	84
4	Interference-Avoidance Scheduling for Dense Multi-User Coexisting Networks with Heterogeneous Priorities and Demands	87
4.1	System Model	88
4.1.1	Network Model	88
4.1.2	Interference Model	89
4.2	Problem Formulation and Solution Framework	90
4.2.1	Admission Control	90
4.2.2	Throughput Maximization of Admitted Users	94
4.3	Column Generation-Based Solution to (P2)	95
4.3.1	Column Generation Preliminaries	96
4.3.2	Restricted Master Problem (RMP)	97
4.3.3	Pricing Problem (PP)	99
4.3.4	Computational Cost Reduction	102
4.4	Column Generation Solution to (P3)	103
4.5	Discussions	105
4.5.1	Algorithm Summary and Computational Complexity	105
4.5.2	Extension to the SINR Model	108
4.6	Numerical Results	110
5	Joint Admission Control and Resource Management for D2D-Assisted Mobile Edge Computing	123
5.1	Related Works	124

5.1.1	Mobile Cloud Computing	124
5.1.2	D2D Communication	127
5.2	System Model and Problem Formulation	128
5.2.1	Computation Model	130
5.2.2	Communication Model	131
5.2.3	Problem Formulation	135
5.3	Extended Formulation and Solution Framework	137
5.3.1	Extended Formulation	137
5.3.2	Solution Framework	140
5.4	Column Generation Solution	144
5.4.1	Solution to Relaxed IP1	146
5.4.2	Solution to Relaxed IP2	150
5.4.3	Computational Complexity Discussions	152
5.4.4	Greedy Algorithm	154
5.5	Simulation Results	157
6	Conclusions and Future Work	164
6.1	Conclusions	164
6.2	Future Work	167
	Bibliography	170
A	Appendixes of Chapter 3	183
A.1	Transition Sub-matrices under Policy I	183
A.2	Transition Sub-matrices under Policy II	184
A.3	Proofs of (3.18) and (3.19)	185
A.4	Proof of Theorem 3.1	185
B	Appendixes of Chapter 4	188
B.1	Proof of Theorem 4.1	188
B.2	Proof of Theorem 4.2	189
C	Appendixes of Chapter 5	190
C.1	Proof of Theorem 5.1	190
C.2	Proof of Theorem 5.2	191
	List of Publications	192

List of Figures

1.1	Conventional Relaying.	3
2.1	Power consumptions under three strategies when $\bar{\gamma}_r = 30$ dB: (a) Strategy I; (b) Strategy II; (c) Strategy III. The results obtained by the proposed method are denoted by solid or dashed lines, and the enumeration method by circles or squares.	47
2.2	Power consumptions under three strategies when $\bar{\gamma}_s = 30$ dB: (a) Strategy I; (b) Strategy II; (c) Strategy III. The results obtained by the proposed method are denoted by solid or dashed lines, and the enumeration method by circles or squares.	48
2.3	Power consumption versus $\bar{\gamma}_s$ under three strategies.	49
2.4	Power consumption versus $\bar{\gamma}_r$ under three strategies.	50
2.5	Total power consumption versus the distance ratio d_1/d	53
3.1	System model.	57
3.2	State transition diagram of the aggregate Markov chain $Y(i) = (\beta(i), Q(i))$ under Policy I.	66
3.3	State transition diagram of the aggregate Markov chain $Y(i) = (\beta(i), Q(i))$ under Policy II.	72
3.4	Effects of P_C in Policy I. Fading margins $F_0 = F_1 = 10$ dB.	76
3.5	Effects of buffer size L in Policy II. $F_0 = F_1 = 10$ dB.	77
3.6	Average throughput versus fading margin $F = F_0 = F_1$ under Policy I.	78
3.7	Average throughput versus fading margin $F = F_0 = F_1$ under Policy II.	79
3.8	Average throughput versus average end-to-end delay for both Policies I and II. $F_0 = F_1 = 10$ dB.	80
3.9	Average throughput versus average end-to-end delay under Policy I and the improved policy. $F_0 = F_1 = 10$ dB.	84
3.10	Comparisons among the blind-COSI policy, the outdated-COSI policy and Policy I with known COSI under the correlated channels. $F_0 = F_1 = 10$ dB.	86
4.1	Multi-user coexisting wireless network.	88

4.2	Flow chart of column generation method.	97
4.3	Convergence performance of the column generation method in solving (P2) when $N = 100$ and $k = k^*$	112
4.4	Convergence performance of the column generation method in solving (P3) when $N = 100$	112
4.5	Total average throughput of admitted users under different algorithms.	113
4.6	Comparisons of the proposed sequential search and the binary search when $K = 10$	116
4.7	Comparisons of the proposed sequential search and the binary search when $K = 5$	117
4.8	Average number of admitted users versus total number of users.	118
4.9	Average number of admitted users versus the coexisting area size.	120
4.10	Average number of admitted users versus the traffic demand.	121
4.11	SINR comparisons of the proposed algorithm and the coloring algorithm.	122
5.1	D2D-assisted mobile edge computing.	129
5.2	Flow chart of the solution framework	142
5.3	Average number of admitted requesters (ANAR) versus the number of channels when $F_{ES} = 10$ GHz, $N_v = N_r = 10$	158
5.4	Average number of admitted requesters (ANAR) versus the computing capacity of the edge server (ES) when $ \mathcal{M} = 5$, $N_v = N_r = 10$	159
5.5	Average number of admitted requesters vs. the number of channels under homogeneous channel conditions when $N_v = N_r = 10$	161
5.6	Average number of admitted requesters vs. the number of channels under heterogeneous channel conditions when $N_v = N_r = 10$	161
5.7	Average number of admitted requesters (ANAR) versus the number of requesters (N_r) under heterogeneous channel conditions when $N_v = N_r$ and $M = 5$	162

List of Tables

2.1	Comparisons of total power consumptions when exchanging the values of $\bar{\gamma}_s$ and $\bar{\gamma}_r$	52
4.2	Probabilities of correctly finding the critical priority level k^*	114
4.3	Average computation time.	115
5.1	Main simulation parameters.	157
5.2	Average computation time of optimal and greedy algorithms.	163

List of Abbreviations

ACK	Acknowledgement
AF	Amplify-and-forward
BBA	Branch and bound algorithm
BS	Base station
BSI	Buffer state information
CDMA	Code division multiple access
COSI	Channel outage state information
CPU	Central processing unit
CSMA/CA	Carrier sense multiple access / collision avoidance
CSI	Channel state information
CTS	Clear-to-send
D2D	Device-to-device
DF	Decode-and-forward
ES	Edge server
FCG	Feasible candidate group
FCLS	Feasible candidate link subset
FDMA	Frequency division multiple access
GIA	Greedy initialization algorithm
i.i.d.	Independent and identically distributed
IP	Integer programming

LTE-A	Long term evolution advanced
MAC	Medium access control
MDP	Markov decision process
MEC	Mobile edge computing
MIMO	Multiple-input-and-multiple-output
MP	Master problem
MU	Mobile user
NACK	Negative acknowledgement
OFDMA	Orthogonal frequency-division multiple access
PP	Pricing problem
QBD	Quasi-birth-death
QoS	Quality of service
RAM	Random-access memory
RMP	Restricted master problem
RTS	Request-to-send
SINR	Signal-to-interference-plus-noise ratio
SNR	Signal-to-noise ratio
STDMA	Spatial reuse time division multiple access
TDMA	Time division multiple access
WBAN	Wireless body area network
WGL	Weighted greedy algorithm based on linear programming relaxation
WISP	Weighted independent set problem

Chapter 1

Introduction

Recently, some advanced wireless communication systems are being developed to meet the demands for wide coverage and high traffic capacities, as well as to satisfy the requirements of many emerging compute-intensive applications. Specifically, wireless relay networks provide reliable long-distance connections for cell-edge users or implement remote wireless backhauling. Dense multiple-user coexisting networks implement non-interfering communications among multiple co-located users with improved spectrum utilization efficiency of wireless networks. Mobile edge computing systems can facilitate the implementations of compute-intensive and delay-sensitive mobile applications such as augmented reality, face recognition, etc. In this thesis, we focus our studies on resource allocation and transmission scheduling issues in these systems, by applying optimization and queueing theory. In the following subsections, we will introduce the background and characteristics of these communication systems, and provide our research motivations.

1.1 Wireless Relay Networks

Wireless networks are experiencing the evolution from voice-centric services to data-centric services, making the distant/cell-edge users become an important system bottleneck in terms of throughput performance. This problem will be even worse for systems operating at high frequencies (e.g., millimeter-waves [1]) that result in larger path loss than low frequencies. The relay communication technique is a promising solution to solve this issue, by adding an extra node between distant transmitter and receiver so as to reduce the one-hop transmission distance and improve throughput.

The seminal works on relay communications can be traced back to around 1960s when Meulen first introduced a three-terminal relay channel in [2] [3], and then in 1976, Cover et al. analyzed the capacity of the three-terminal relay channel in [4]. After that, little work on relay communications was done until about ten years ago. In around 2003, Sendonaris and Laneman et al. discussed further the benefits of user cooperation (called cooperative diversity) from the information-theoretic point of view in [5]–[8]. Since then relay communication techniques attract a lot of attention because they demonstrate great potential to improve the performance of wireless networks in terms of system throughput, coverage, and robustness to channel variations. Nowadays, relay communication techniques have already been included in some wireless communication standards such as Long Term Evolution-Advanced (LTE-A) [9] [10], IEEE 802.16j [11] and IEEE 802.16m [12].

A classic two-hop three-node relay network consists of a source node, a relay node and a destination node, where the relay assists in transmitting the data from the source to the destination, as shown in Fig. 1.1. Two main types are amplify-and-

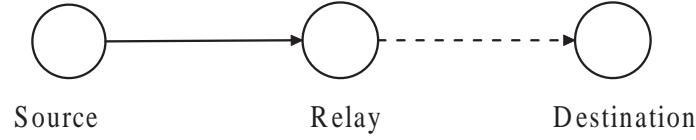


Figure 1.1: Conventional Relaying.

forward (AF) and decode-and-forward (DF) relaying. Under AF relaying, the relay first receives the signal from the source, and then amplifies and forwards the received signal to the destination. AF relaying is favored for its simple implementation since the relay does not need to carry out complex signal processing. Under DF relaying, the relay first receives and decodes the signal from the source, and then forwards the recovered signal to the destination.

1.1.1 Power Allocation in Amplify-and-Forward Relaying

Power allocation/control plays an important role in implementing amplify-and-forward relaying. This is because the received signal of the relay is simply amplified and forwarded to the destination, and it is not processed and refined by the relay unlike decode-and-forward relaying. Due to imperfect electronic components, communication receivers are normally affected by thermal noises. In the relaying process, the noise introduced in the receiver of the relay is also amplified and transferred to the destination. This means that the receiver of the destination will be destructed by two noises: the noise introduced by its own components and the amplified noise from the relay. The noise accumulation degrades the relaying performance. Thus, it is necessary to optimize the transmit power use and limit the noise.

In the literature [13]– [24], extensive work on optimizing power allocation focused

on two categories of problems: 1) rate/signal-to-noise ratio (SNR) maximization or outage probability minimization under a maximum power constraint [19]– [24]; and 2) transmit power minimization for saving energy under a minimum rate/SNR requirement or a maximum allowable outage probability constraint [13]– [18]. The performance of power allocation determines whether cooperative benefits can be obtained and how much cooperative gain can be achieved, while itself greatly depending on the availability of channel state information (CSI). In general, an optimal power allocation algorithm requires complete instantaneous CSI. However, in practice, the acquisition of complete instantaneous CSI will introduce a considerable amount of control channel overheads. For example, in a wireless relay network consisting of a source, a relay and a destination, the CSI on the source-relay link is ordinarily estimated by the relay and the CSI on the relay-destination link is ordinarily estimated by the destination. All of these estimates are required to be fed back to a control center for executing the power allocation algorithm. In addition, after performing the optimization algorithm, the control center needs to report the optimal power values to the transmitters (i.e., source and relay).

In Chapter 2, we will present a power allocation scheme with reduced overheads for AF relay networks while avoiding a significant performance loss caused by the lack of complete instantaneous CSI.

1.1.2 Enhanced Decode-and-Forward Relaying

Conventional Half-duplex Decode-and-Forward Relaying

Normally, relaying systems adopt a half-duplex pattern, i.e., the relay receives and forwards signals on two different time intervals to avoid self-interference. This is because it is quite demanding in practice to implement the full-duplex pattern where the relay simultaneously receives and forwards the signals on the same frequency band. Under the full-duplex relaying, there is a large imbalance between the transmitting and the receiving power, where the strong transmitting power saturates the radio frequency (RF) receiver chain (especially the analog-to-digital converter having a small dynamic range). Such effect of the relay's transmitting signals on its own receiver chain is called self-interference. Although the implementation of full duplex relays may be possible by self-interference cancellation techniques in both analog and digital domains [25] [26], much more complex processing is required in both hardware and software. In addition, it is inefficient and impractical to use two different frequency bands for transmitting and receiving, because of the scarcity of radio spectrum resources available for wireless communications. Hence, half-duplex relays are preferred in practice due to their easier implementations than full-duplex relays.

The conventional half-duplex DF relaying works as follows [7] [8]: Time is divided into frames with an equal length and each frame is further split into two equal-duration slots. In any frame, the source first transmits a signal to the relay which decodes this signal in the first slot, and then the relay immediately forwards the decoded signal to the destination in the second slot. In this way, the source transmits signals in only half time, and thus the half-duplex relaying suffers from a loss of $1/2$ factor in the

average throughput, which can be represented as

$$C = \frac{1}{2}E[\min\{\log(1 + \gamma_{sr}), \log(1 + \gamma_{rd})\}], \quad (1.1)$$

where γ_{sr} and γ_{rd} denote the instantaneous signal-to-noise ratios (SNRs) over the source-relay and the relay-destination links, respectively, which may vary from one frame to another. $E[\cdot]$ denotes the expectation. In addition, from the above formula, we can see that the performance of conventional DF relaying is dominated by the worse link between the source-relay and the relay-destination ones.

Buffer-Aided Decode-and-Forward Relaying

Conventional DF relaying [5]–[8], [13] [27] adopts a predetermined transmission schedule where the transmission of one packet is finished in two successive slots, i.e., the source transmits a packet to the relay in the first slot and the relay immediately forwards the received packet to the intended destination in the second slot. However, such a fixed scheduling strategy greatly limits the throughput of relaying systems, and this issue may be overcome by introducing data buffers at the relays to allow the transmissions of two distinct packets in two successive slots. In [28] [29], using the buffering abilities of the relays, the authors proposed a max-max relay selection scheme for multi-relay systems where one packet is transmitted over the best source-relay link in the first slot and possibly another different packet over the best relay-destination link in the second slot. However, the max-max relay selection is still limited by the fixed two-slot link scheduling structure where the first slot is always allocated for the source-relay links and the second slot for the relay-destination links. This limitation is relaxed in [30] where the authors proposed a max-link selection scheme which selects

the best link among all available source-relay and relay-destination links in any time slot. Such an adaptive link selection allows relaying systems to further exploit the spatial diversity of wireless channels so as to result in improved system throughput.

Delay Problems in Buffer-Aided DF Relaying

In the buffer-aided relaying system, the relay may not immediately forward the received packet to the destination, and thus the buffering at the relay inevitably results in extra delay, finally affecting the experience of users (for example, online videos may not run smoothly). In particular, under time-correlated fading channels, the buffer-aided relaying may lead to degraded delay performance.

In practice, fading channels are commonly time-correlated especially for low-speed mobile nodes [31]– [35]. For example, in [31] [32], the authors pointed out that, for “slow” fading with $f_D T < 0.1$ (f_D is the Doppler frequency and T is the time slot length), the time-correlation of channel fading between adjacent slots cannot be neglected. In addition, buffer-aided relaying prefers more to “slow” fading scenarios since adaptive link scheduling/selection requires channel fading remaining invariant for at least one slot. Moreover, buffer-aided relaying may also request a relatively long slot length in order to compensate the time overhead for potentially frequent switching between reception and transmission modes. In fact, unlike the conventional relaying without buffers, correlated fading may have great effects on the delay performance of buffer-aided relaying. For example, the correlated fading may cause the source-relay link and the relay-destination link to keep staying in “good” and “bad” states, respectively, for a long time (multiple successive slots). Then, the link scheduling scheme

based on only link states may always choose the source-relay link, which may cause packets backlogged in the relay's buffer and lead to potential buffer overflow. These effects may, in turn, result in a long transmission delay and throughput degradation. On the contrary, if the source-relay link is continuously in "bad" state while the relay-destination link is in "good" state, the relay's buffer may be emptied and suffer from buffer underflow.

However, to our best knowledge, almost all existing works about buffered-aided relaying in the literature assumed independent and identically distributed (i.i.d.) fading channels due to their simplicity and tractability. Only in [36], time-correlated fading channels were briefly mentioned without detailed performance analyses. Thus, it is important and necessary to investigate the performance of buffer-aided relaying under time-correlated fading channels. We will investigate this issue in details in Chapter 3.

1.2 Dense Multi-User Coexisting Networks

1.2.1 Medium Access Control

In traditional wireless ad hoc networks, medium access control (MAC) protocols normally adopt random contention access methods, where each user contends for the channel when it has packet arrivals for transmission. One of frequently adopted protocols is carrier sense multiple access/collision avoidance (CSMA/CA) [37], where each user detects if the channel is occupied by other users before transmission. It transmits only if the channel is idle. The user will back off a random period to avoid interference if it finds the channel is busy. However, carrier sense is only carried out at

the transmitter whereas packet collision actually occurs at the receiver. This may lead to the hidden terminal problem where the transmitter does not detect interference while the receiver is actually in the interference range. One way to avoid this is handshaking: the transmitter first sends a short request-to-send (RTS) signal to the receiver if it detects an idle channel, and then the receiver feeds back a clear-to-send (CTS) signal to the transmitter if it also detects the idle channel. The transmitter sends a message only if it successfully receives the CTS signal. If the receiver finds the channel is busy, it will not feed back the CTS signal.

Although much work has been done to improve random access, it leads to a very low channel utilization rate due to its inherent random collisions, which makes it inapplicable to user-dense coexisting networks. The analyses in [38] [39] showed that adopting the IEEE 802.15.4-based random access protocol for multi-user coexistence, leads to a large loss of throughput and very low channel utilization when there are a large number of users. Also, it is difficult for random access to guarantee delay performance due to the inherent random collisions. Moreover, RTS/CTS signals may also be lost, which further exacerbates delay performance.

In traditional cellular networks, MAC protocols are usually based on orthogonal multiple access patterns. The most frequently used protocols include time division multiple access (TDMA), frequency division multiple access (FDMA), and code division multiple access (CDMA). In these protocols, each user has exclusive resource blocks (time intervals, frequency bands, or codewords) for transmissions. The same frequency band can only be reused in non-adjacent cells for capacity enhancement. Such frequency reuse is only effective for macro-cell networks where the locations

and coverage of macro base stations (BSs) are fixed and well planned in advance. It is inappropriate for many emerging wireless networks, such as unplanned small cell networks with random positions and scalable coverage, device-to-device (D2D) communication underlaying cellular networks where D2D links reuse spectrum with cellular links, wireless body area networks (WBANs) where a lot of WBANs may coexist and share spectrum, etc. A common feature of these networks is that a large number of users may randomly gather in a small area, forming dense multiple user co-existing networks. To satisfy these users' traffic demands, advanced multi-user access protocols are required.

To address the dense multi-user coexisting issue, one potential solution is to apply spatial reuse TDMA (STDMA), where users in close proximity are scheduled in different time intervals, but users geographically separated are allowed to reuse a time interval. Geographical separation significantly decreases the interference among users, and spectrum sharing among these users will greatly boost the network capacity. To implement STDMA, we have to answer a series of questions: which users should be grouped together to share a time interval and how much time should be assigned to each user group. User scheduling with interference awareness and time assignment are the key issues to improve the spectral utilization of STDMA, which requires a fast and effective optimization algorithm.

1.2.2 Priority-Aware Interference-Avoidance Scheduling

As mentioned before, spectrum sharing that allows multiple users to concurrently transmit has been considered as an important way to enhance the spectral utilization

rates and capacities of wireless networks [40]–[42]. A key issue is to properly schedule the transmissions of coexisting users to avoid interference while maintaining high concurrency. The interference-avoidance scheduling for multi-user coexisting networks has been attracting a lot of interest in different areas of wireless communication environments, such as wireless sensor networks [43] [44], cognitive radio networks [45] [46], wireless mesh networks [47] [48], wireless multi-hop networks [49], wireless body area networks [50] [51], etc.

In addition to spectrum sharing, priority-aware scheduling is also an urgent issue for multi-user coexisting wireless networks. One typical example is the wireless body area network (WBAN), which consists of wireless sensors implanted in and/or placed on patients' bodies to collect physiological signals. Extensive deployment of sensors in WBANs (the IEEE 802.15.6 standard [52] requires that a WBAN can admit up to 64 sensor nodes) may introduce high traffic so that system overload may become inevitable [53], especially in patient gathering areas like emergency wards in a hospital or dining halls in an assisted living home. In addition, different patients may have heterogeneous priorities with respect to their health conditions and the criticality of diseases [54] [55]. For example, intensive care/monitoring patients such as those suffering from acute diseases or heart diseases should be given high priority while general care patients can be given lower priority. Priority awareness is an important feature in WBANs and has gained a lot of attention. In the IEEE 802.15.6 standard [52], user priorities have been defined. In [56] [57], the authors considered the seriousness of health parameters and proposed a priority-aware adaptive medium access control (MAC) mechanism for improving the reliability of critical nodes in WBANs. In [58],

a priority-based data rate tuning mechanism was developed for WBANs to improve the quality of service of sensors with critical physiological data.

In the literature, much work has been done on multi-user coexistence interference-avoidance scheduling. Most of the studies are based on random access/scheduling, graph coloring or optimization. Carrier sense multiple access/collision avoidance (CSMA/CA) [37] was widely used for random access where a user could back off for a random period to avoid interference when it detects a collision. Another random access method is the slotted ALOHA [59]. Different from CSMA/CA, the slotted ALOHA does not carry out carrier sensing, but uses a control parameter, called contention probability, to manage the access of a user. The priority-aware random access can be achieved by using different backoff periods or contention probabilities [59]. In addition, the superframe-level frequency-hopping or time-hopping scheme was also proposed to coordinate interference [60]. However, it has been shown in the previous subsection that, random access leads to very low channel utilization, which makes it inapplicable to user-dense coexisting networks.

In the graph coloring method, a network is modeled as a graph $G = (V, E)$ where vertices V denote nodes and edges E represent conflicts between mutually interfering nodes. Each time/channel resource unit is represented by a color. Then, the coexistence scheduling was realized by vertex coloring [50] or edge coloring [61], where interference is avoided by assigning different colors to adjacent vertices (vertex coloring) or edges incident on the same node (edge coloring). Although complete coloring could find the minimum number of colors, it is known to be an NP-complete problem [61] and has an exponentially increasing complexity. Thus, a lot of work has been

done focusing on time-complexity reduction. For example, a quick coloring algorithm was proposed for mobile ad hoc networks (MANETs) [62] and a random incomplete coloring was used to balance time complexity and scheduling performance [50]. A drawback of graph coloring is that it requires a fixed-size resource unit for each color. Thus, this method cannot be well applied to the networks with heterogeneous/variable traffic demands.

The coexistence scheduling can also be mathematically modeled as an optimization problem, where the objective is to effectively allocate time/channel resources to all feasible candidate groups (FCGs) [44]- [47], [49] [63] [64]. Here, an FCG refers to a group of coexisting users that can concurrently operate on the same channel without interference. In this method, the resource unit size is not necessarily fixed and co-existing users are allowed to have heterogeneous demands. In addition, this method enables the joint optimization of scheduling, routing, and power control. For example, joint routing and link scheduling were studied for multi-hop backhaul networks [49] and software-defined networks [65]. The authors in [63] and [44] discussed the issue of joint power control and scheduling in wireless ad hoc networks and sensor networks, respectively. In these works, it has been shown that the optimization framework has a great potential to further improve the system capacity of multiple user coexisting networks. However, in the literature, heterogeneous priorities among coexisting users are seldom taken into account in the optimization model.

Challenges in Priority-Aware Scheduling Optimization

- The consideration of priorities significantly increases the difficulty in scheduling design, since the number of user priorities that a system can accommodate

is uncertain and unknown, especially when the network becomes so crowded that it has to drop some low priority users. This requires a fast and effective searching algorithm for admission control, in order to find the maximum number of priority levels that the system can support so that all users in these high priority levels can be admitted. This is different from existing works that mainly focusing on two types of scheduling problems: i) time length minimization to satisfy users' traffic demands [64] [66]; and ii) throughput maximization [67].

- In the optimization model, it is very challenging to find all potential FCGs, the number of which grows exponentially with the number of coexisting users. In the worst case, we need to search all 2^N possible groups for an N -user coexisting network, which is definitely a time-consuming task if N is large. Furthermore, the huge number of FCGs makes the resource allocation problem involve too many variables to be solved efficiently

In Chapter 4, we will address the interference-avoidance scheduling issue of multi-user coexisting networks by considering heterogeneity in terms of user priorities and traffic demands.

1.3 Mobile Edge Computing Systems

1.3.1 Mobile Edge Computing

Nowadays, smart mobile devices (e.g., smartphones and tablets) have played more and more roles in people's daily lives for business, learning, and entertainment [68]. Many emerging mobile applications are gaining increasing attention, such as face

recognition, natural language processing, interactive gaming, and augmented reality [69]. These applications are ordinarily hungry in the resource, requiring executing huge computations and consuming a large amount of energy. However, mobile devices normally have limited processing abilities and battery capacities due to their small physical sizes. This causes a strong conflict between resource-hungry mobile applications and resource-limited mobile devices.

To address such a conflict, mobile edge computing (MEC) [69]- [71] is developed as a promising solution. It provides cloud computing services at the edge of cellular networks (e.g., cloud computing servers can be installed in/near the base stations of cellular networks), and mobile users can offload their computation tasks to cloud servers via cellular connections. MEC is beneficial for real-time interactive applications since it requires only one-hop low-latency connection for offloading computation data from mobile users to cloud servers. This is different from remote public clouds, such as Amazon EC2 and Windows Azure, where the data exchanges between users and remote cloud servers usually go through multi-hop wide area networks, possibly leading to a long unreliable latency. In addition, compared to cloudlet computing via one-hop WiFi access, MEC provides a larger coverage via cellular connections, satisfying ubiquitous service demands. Thus, MEC has been envisioned as an effective complement to remote public clouds and cloudlets.

1.3.2 D2D-Assisted Mobile Edge Computing

Although MEC has great potentials to augment the computing capabilities of mobile devices for resource-intensive applications, it encounters many significant chal-

lenges. One major challenge is the tremendously increased mobile data traffic in cellular networks. It has been estimated by Cisco [72] that global mobile data traffic grew 63 percent in 2016, reaching 7.2 exabytes (7.2 billion gigabytes) per month, and is expected to continue growing nearly 7 folds between 2016 and 2021. Unfortunately, the transmission capacities of cellular networks are rather limited, leading to a big shift of mobile traffic (60 percent in 2016 [72]) from mobile networks to fixed networks. Thus, the bottleneck of cellular network capacities is going to severely block the implementation of MEC, especially when enormous mobile users swarm to networks. Another major challenge results from the fact that MEC servers are usually not built as powerful as public clouds, due to the limitations on the site space for installing base stations, initial investments, and subsequent operational costs. Thus, the computing resources at MEC servers are quite limited instead of infinity as commonly assumed for public clouds and may not be able to support enormous mobile computation tasks.

To alleviate the stress on both MEC servers and cellular networks, the MEC system can be enhanced by introducing device-to-device (D2D) communication where adjacent users are allowed to directly communicate without the need of going through the BS. In such a system, part of mobile computation tasks can be offloaded to vacant mobile devices via D2D communication links. The benefits are two-fold: i) cellular traffic is reduced by offloading traffic from cellular links to D2D links; and ii) the computing capability of the MEC server is expanded by leveraging vacant mobile devices. Moreover, the wide coverage of cellular networks can effectively alleviate the mobility issue that a mobile user may move out of the communication range of its

serving node (e.g., a vacant device), and the computing result cannot be transmitted back to the mobile user. Different from cloudlet computing where the mobility issue is severe [73], in the mobile edge computing based on cellular networks, the computing results of serving nodes can be fed back to mobile users via cellular base stations. Since the amount of the computation result data is normally much smaller than that of the computation input data for most mobile applications like face recognition and virus scanning [69] [74], the feedback of computation results won't impose too much traffic load on cellular networks.

Resource Management and Challenges

Due to the scarcity of spectrum resource, D2D connections may have to share the same spectrum with cellular connections. This may cause severe interference among D2D and cellular connections. In order to capitalize on the benefits introduced by D2D communication, we have to effectively control the interference. Otherwise, the potential benefits will disappear with the increased interference. Specifically, we have to address the following issues: i) which computation requesters should be admitted (i.e., admission control); ii) which vacant mobile device or the server itself should be selected to serve each admitted requester (i.e., link scheduling/user association); iii) which channel should be chosen for each scheduled link; and iv) how much transmit power should be set on each link for offloading computation input data. To answer all these questions, a joint admission control, link scheduling and resource management optimization problem has to be solved. However, the combination property from user association and channel assignment makes our problem fall into the scope of combinatorial optimization, which is extremely hard to solve.

In Chapter 5, we will investigate in details the resource management issue in D2D-assisted MEC systems.

1.4 Contributions

The main objective of this research is to design effective and efficient resource management and transmission scheduling schemes for wireless communication networks. These designs are based on the distinct characteristics of networks. Specifically, three networks are considered, i.e., wireless relay networks, dense multiple user coexisting networks, and mobile edge computing systems. The main contributions are summarized as follows.

1) Transmit power allocation for amplify-and-forward relay networks with reduced overheads (this work has been published in the IEEE Transactions on Vehicular Technology).

- We propose a new hybrid/partial channel state information (CSI) based power allocation strategy for amplify-and-forward relay networks where the source power and relay power are determined based on statistical and instantaneous CSI, respectively. The objective is to minimize the total power consumption while satisfying a minimum rate constraint or a minimum non-violation probability constraint.
- This new strategy can reduce the control channel overheads by at least 50% compared with the traditional strategy based on complete instantaneous CSI while avoiding a large performance loss as in the statistical CSI strategy. Besides, it

is shown that the proposed partial CSI strategy can obtain near-optimal performance as the complete instantaneous CSI strategy when the relay is located close to the source.

- The two-stage stochastic programming method is applied to formulate the considered optimization problem. Moreover, an N -section method is proposed to solve the power allocation problems for both statistical and partial CSI strategies, which comes close to optimal solutions in simulations.
- 2) An analysis framework for buffer-aided decode-and-forward relaying under time-correlated fading channels (this work has been published in the IEEE Transactions on Vehicular Technology).
- A framework is formulated to analyze the performance of buffer-aided relaying under time-correlated fading channels in terms of the queueing behaviors of packets in the relay buffer. The average throughput, outage probability and average end-to-end delay are provided. Two delay-controllable link scheduling/selection policies with respect to infinite and finite buffers are considered in the formulated framework.
 - The results show that compared with independently and identically distributed (i.i.d.) fading, correlated fading causes greater throughput degradation for transmissions with stringent delay requirements in low fading margins. The fading margin means the maximum fading attenuation that the system can tolerate for successfully transmitting one packet. In particular, a throughput loss of about 16% under an infinite buffer (28% under a finite buffer) is observed for

the requirement of an average delay of 20 slots and a fading margin of 5 dB. This means that the link scheduling design based on i.i.d. fading is not always fit for correlated fading. According to these observations, some insights on performance degradation and guidelines on redesigning improved link scheduling policies under correlated fading are provided.

- In addition, the results show that the link scheduling policy with an infinite buffer can obtain higher average throughput than that with a finite buffer under the same average delay, which implies the superiority and flexibility of a large buffer size in link scheduling.

3) Interference-avoidance scheduling for dense multi-user coexisting networks with heterogeneous priorities and demands (this work has been accepted for publication in the IEEE Transactions on Wireless Communications).

- Both admission control and throughput maximization issues for dense multi-user coexisting networks are investigated, where the objective is to first allow the system admit as many high priority users as possible and then maximize the throughput of admitted users.
- A sequential solution framework based on priority constraints is presented to find the maximum number of high priority levels that the system can support so that all users in these high priority levels can be admitted. At each step, a large-scale linear subproblem requires being solved. After that, the throughput of admitted users is maximized by formulating another large-scale linear programming subproblem.

- To solve these large-scale linear programming subproblems, the column generation method is introduced. In this method, by capitalizing on the special structure of the sequential solution framework, a greedy initialization algorithm is proposed to warmly start the column generation process. In addition, both upper and lower bounds on the optimal objective function of each subproblem are derived. By applying these bounds to the sequential solution framework, it is not necessary to optimally solve every problem, which significantly reduces the computation time.
 - Simulation results verify that the proposed algorithm can effectively and efficiently address the interference-avoidance scheduling issue for dense multi-user coexisting networks with heterogeneous user priorities and traffic demands.
- 4) Joint admission control, link scheduling, and resource management for D2D-assisted mobile edge computing (this work has been submitted to the IEEE Transactions on Mobile Computing).
- A joint admission control, link scheduling and resource management issue for D2D-assisted mobile edge computing (MEC) systems is studied, which aims to determine i) which requests should be admitted under limited radio and computing resources (i.e., admission control), ii) which vacant mobile device or the MEC server should be selected to serve each admitted requester (i.e., link scheduling/user association), iii) which channel should be assigned to each scheduled link, and iv) how much transmit power should be used on each link for offloading computation input data. The objective is to maximize the number of admitted requests.

- An optimal branch-and-price based algorithm is developed to solve the joint optimization problem of admission control, link scheduling, channel assignment, and power control. Although the computational complexity of the optimal solution is high, it can be used as a performance benchmark. In addition, a low-complexity suboptimal algorithm is proposed for the practical implementation purpose.
- Simulations show that the proposed D2D-assisted MEC scheme can significantly increase the number of admitted requesters under limited radio and cloud computing resources compared to counterparts.

1.5 Organization of the Thesis

The rest of the thesis is organized as follows. In Chapter 2, amplify-and-forward relaying is considered where the transmit power allocation with reduced overheads is investigated. In Chapter 3, buffer-aided decode-and-forward relaying is considered where an analysis framework under time-correlated fading channels is presented. The priority-aware interference-avoidance scheduling for dense multi-user coexisting networks is discussed in Chapter 4. In Chapter 5, the joint admission control, link scheduling, and resource management issue for D2D-assisted mobile edge computing systems is studied, followed by conclusions and future work in Chapter 6.

Chapter 2

Transmit power allocation for amplify-and-forward relay networks with reduced overheads

In this chapter, we consider a typical three-node amplify-and-forward relay network consisting of a source, a relay, and a destination, as shown in Fig. 1.1. This system model can also be considered as a part of selection relay networks where only a “best” relay node is selected to forward the signal from the source. We aim to minimize the total transmit power with a consideration of control channel overheads under a minimum rate constraint (equivalent to a signal-to-noise ratio constraint) or a minimum non-violation probability constraint (equivalent to a maximum allowable outage probability constraint). For the sake of completeness and comparison, we first define and analyze two schemes based on traditional routines, called Strategy I and Strategy II. In Strategy I, the transmit powers of both the source and the relay are

adjusted based on complete instantaneous CSI, while in Strategy II, the powers of both the source and the relay are based on statistical CSI. After that, we propose a new hybrid/partial strategy, called Strategy III, where the source power is based on statistical CSI while the relay power is determined based on instantaneous CSI. We then formulate and solve Strategy III by the *two-stage stochastic programming method* [75]. The first-stage optimization problem aims to find an optimal source power to minimize the power sum of the source and the relay based on statistical CSI. The optimal relay power is derived from the second-stage optimization problem, where the relay power can be treated as a “penalty” due to the inadequate source power for satisfying a minimum transmission rate. Namely, if the source power is smaller, then the relay requires a larger transmit power to compensate for the shortage of the source power so that the minimum rate can be met.

2.1 Related Works

In order to reduce the overheads for reporting CSI, the statistical CSI (such as the mean or distribution of CSI) instead of instantaneous CSI is used for optimizing power allocation. In [76] and [77], the authors proposed a power allocation scheme which requires only knowledge of average channel gain for decode-and-forward relay networks, aiming to minimize the outage probability or the upper bound of outage probability. In [78], with only statistical CSI (fading distribution and path loss information) at transmitters, a power allocation scheme was derived for minimizing the approximations of outage probabilities in the high SNR regime under decode-and-forward, amplify-and-forward, and distributed space-time coded relaying protocols.

Focusing on multi-hop relay networks, the authors in [79] discussed an optimal power allocation strategy based on statistical CSI for both decode-and-forward and amplify-and-forward relay networks, with the objective to minimize the outage probability. In [80], with only knowledge of statistical CSI, an optimal power allocation algorithm was proposed to jointly minimize both upper and lower bounds of outage probability for amplify-and-forward parallel multi-relay networks. In addition, for reducing the CSI feedback overheads, the authors in [81] proposed a method to directly estimate at the destination the overall source-relay-destination channel instead of estimating the source-relay and relay-destination channels separately, so that the CSI on the source-relay link does not need to be transmitted to the destination. However, in such a way, the overheads for reporting the power allocation results from the destination to the source and the relay are still large.

Due to the lack of instantaneous CSI, the power allocation strategy based on statistical CSI will inevitably result in performance degradation. To prevent a significant performance loss, some researchers turned to the power allocation algorithms based on partial CSI where each node has only knowledge of local instantaneous CSI on links connected to itself and/or statistical CSI of other links. In [82], an adaptive power control algorithm for adjusting the transmit power of the selected relay node was proposed for decode-and-forward selection relay networks where the selected relay node requires only the instantaneous CSI on the itself-destination link, while the source node is allocated a fixed transmit power. In [83], the authors discussed a distributed power allocation strategy based on local CSI for decode-and-forward parallel relay networks where the source requires the instantaneous CSI of all source-relay

links and the source-destination link, and the statistical CSI of all relay-destination links, while each relay node requires the instantaneous CSI of source-itself and itself-destination links. In [84], for optimizing the transmit powers of relay nodes in amplify-and-forward parallel relay networks, the authors derived power allocation algorithms under two assumptions on partial CSI: i) Each relay has the instantaneous CSI of all source-relay links; and ii) Each relay has only the local instantaneous CSI of the source-itself link and has no knowledge of links from the source to other relays. Under both assumptions, only the statistical CSI of source-destination and relay-destination links is available at each relay. Following the similar assumption as in [84], the authors in [85] proposed a power allocation algorithm for both parallel and selection amplify-and-forward relay networks. In [84], the average signal-to-noise ratio (SNR) maximization problem under a maximum power budget was addressed while in [85] the total power minimization problem under an average SNR constraint was studied. However, as pointed out in [85], the outage probability is still much higher under an average SNR constraint, which is not good enough to guarantee users' quality of experience.

Although some power allocation strategies based on partial CSI have been proposed in literature [82]–[85], most of these works have not fully and systematically answered the following questions: 1) How much overheads of control channels can be reduced by partial CSI strategies when compared to complete instantaneous CSI strategies; 2) How much performance degradation partial CSI strategies may lead to; 3) In which cases the performance of partial CSI strategies can approach to the optimal one based on complete instantaneous CSI. In this chapter, we are going to

answer these questions via analytical or numerical results.

The main contributions of this chapter are summarized as follows:

(1) We propose a new partial CSI based power allocation strategy (Strategy III) for amplify-and-forward relay networks where the source and the relay powers are determined based on statistical and instantaneous CSI, respectively. This new strategy can reduce the control channel overheads by 50% compared with the complete instantaneous CSI strategy (Strategy I) while avoiding a large performance loss as in the statistical CSI strategy (Strategy II). Besides, it is shown that the proposed partial CSI strategy can obtain near-optimal performance as the complete instantaneous CSI strategy when the relay is located close to the source.

(2) The two-stage stochastic programming method is applied to formulate the optimization problem of Strategy III. Moreover, an N -section method is proposed to solve power allocation problems for both statistical and partial CSI strategies (Strategies II and III), which can lead to almost-optimal solutions.

(3) Different from [84] and [85] where the average SNR was treated as a performance metric for simplicity, the non-violation probability constraint (equivalent to outage probability constraint) is explicitly applied in this chapter.

(4) Different from [78] and [80] which used the approximations or bounds of outage probabilities for parallel relay networks, the exact expressions of non-violation probabilities are applied in this chapter. In addition, different from [78]–[80] where outage probability minimization problems were studied, the total power minimization problem under a non-violation probability constraint is investigated in this chapter.

(5) Different from [86]–[88] where only the relay power allocation was optimized,

in this chapter, both the source and the relay powers are simultaneously adjusted by the two-stage stochastic programming method. Although the source power is also fixed in the proposed strategy, its value is chosen properly. In addition, different from [86]- [88] which mainly focused on SNR or capacity maximization, we aim to minimize the total transmit power.

2.2 Network model and Problem formulation

Consider a typical two-hop relay network consisting of a source node, a half-duplex amplify-and-forward relay node, and a destination node, as shown in Fig. 1.1. We consider the case where there is no direct link between the source and the destination, which may happen when the source and the destination are far away from each other or there exists a barrier between them. Time is divided into frames with equal lengths and each frame is further split into two equal-duration slots. In the first slot, the relay receives signals from the source. In the second slot, the relay amplifies its received signals during the first slot and transfers the amplified signals to the destination.

In the following, we will introduce three power control strategies with different levels of CSI and formulate an optimization problem for each strategy with an objective to minimize the total power consumption of both the source and the relay while ensuring a minimum transmission rate or non-violation probability. Besides, we will compare these strategies in terms of control channel overheads resulting from the exchanges of CSI and optimization results.

2.2.1 Strategy I

In Strategy I, one node is selected from the source, relay or destination as the central node with tasks of i) collecting the complete instantaneous CSI of all links, ii) determining the optimal transmit powers of both the source and the relay according to the collected CSI, and iii) distributing the optimal solutions (i.e., transmit power values) to the source and the relay. Obviously, the collection of instantaneous CSI and the distribution of optimal power values contribute to overheads on control channels.

Let P_s and P_r denote the transmit powers of source and relay, respectively. $\gamma_s = \frac{|h_{sr}|^2}{\sigma_r^2}$ ($\gamma_r = \frac{|h_{rd}|^2}{\sigma_d^2}$) denote the ratio of channel gain to noise power over the link from source to relay (from relay to destination), where h_{sr} (h_{rd}) is the channel coefficient and σ_r^2 (σ_d^2) is the noise power received at the relay (destination). Both channels are block faded so that γ_s and γ_r are constant during each frame, but may be varying from one frame to another (i.e., γ_s and γ_r are random variables). For Rayleigh fading, γ_s (γ_r) follows an exponential distribution with a parameter $\frac{1}{\bar{\gamma}_s}$ ($\frac{1}{\bar{\gamma}_r}$) where $\bar{\gamma}_s$ ($\bar{\gamma}_r$) is the expected value of γ_s (γ_r). The probability density functions of γ_s and γ_r can be represented, respectively, as

$$f_{\gamma_s}(\gamma_s) = \begin{cases} \frac{1}{\bar{\gamma}_s} e^{-\frac{\gamma_s}{\bar{\gamma}_s}}, & \gamma_s \geq 0 \\ 0, & \gamma_s < 0 \end{cases} \quad (2.1)$$

$$f_{\gamma_r}(\gamma_r) = \begin{cases} \frac{1}{\bar{\gamma}_r} e^{-\frac{\gamma_r}{\bar{\gamma}_r}}, & \gamma_r \geq 0 \\ 0, & \gamma_r < 0 \end{cases} \quad (2.2)$$

The optimization problem with the objective to minimize the power sum of the source and the relay while satisfying the requirement of a minimum transmission rate can

be formulated as

$$(P1) \quad \min_{P_s, P_r} \quad P_s + P_r \quad (2.3a)$$

$$\text{s.t.} \quad r(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r) \geq r_{min}, \quad (2.3b)$$

$$0 \leq P_s \leq P_s^{max}, \quad 0 \leq P_r \leq P_r^{max}, \quad (2.3c)$$

where P_s^{max} (P_r^{max}) is the maximum power budget of source (relay). r_{min} is the required minimum rate and $r(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r)$ denotes the transmission rate, which can be calculated based on Shannon's capacity equation as

$$r(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r) = (B/2) \log_2 (1 + \gamma(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r)), \quad (2.4)$$

where $\gamma(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r)$ is the received SNR at the destination, $\tilde{\gamma}_s$ ($\tilde{\gamma}_r$) denotes one realization of random variable γ_s (γ_r), i.e., instantaneous CSI, and B is the channel bandwidth. We consider the additive white Gaussian noises with zero means and assume that the transmitted signal is independent from the noises. Then the received SNR at the destination can be calculated as [8] [80]

$$\gamma(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r) = \frac{P_s \tilde{\gamma}_s P_r \tilde{\gamma}_r}{P_s \tilde{\gamma}_s + P_r \tilde{\gamma}_r + 1}. \quad (2.5)$$

Because of Constraint (2.3b), the rate violation can only happen when P_s and P_r have reached their maximum values. Thus, the probability without violating the minimum rate requirement can be represented by

$$p_{no} = Pr\{r(P_s^{max}, P_r^{max}, \gamma_s, \gamma_r) \geq r_{min}\}. \quad (2.6)$$

In Strategy I, the instantaneous channel states $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$ are first measured at the relay and the destination, respectively, and are then reported to the central node for

optimizing power allocation. Next, we briefly discuss overheads on control channels based on different selections of the central node.

1) The source is selected as the central node. In this case, the channel state $\tilde{\gamma}_s$ needs to be transmitted from the relay to the source, and $\tilde{\gamma}_r$ from the destination to the source (note that this transmission may need two hops, i.e., from the destination to the relay and from the relay to the source when the direct link between the source and destination is not available). After solving the optimization problem, the source keeps its own optimal power value P_s^* and feeds back P_r^* to relay. Let K denote the units of overheads used for transmitting one value ($\tilde{\gamma}_s$, $\tilde{\gamma}_r$, P_s^* , or P_r^*) over one hop. Then the total overheads are $3K$ units ($4K$ units if the direct link is not available).

2) The relay is selected as the central node. The relay has the channel state $\tilde{\gamma}_s$ measured by itself and receives $\tilde{\gamma}_r$ from the destination. After computing the optimal power values P_s^* and P_r^* , the relay keeps its own power value P_r^* and transmits P_s^* to the source. The total overheads are $2K$ units.

3) The destination is selected as the central node. The destination measures the channel state $\tilde{\gamma}_r$ by itself and receives $\tilde{\gamma}_s$ from the relay. After calculations, the destination transmits P_r^* to the relay through one hop and transmits P_s^* to the source through one or two hops depending on whether the direct link is available between the source and the destination. The total overheads are $3K$ units ($4K$ units if the direct link is not available).

2.2.2 Strategy II

In Strategy II, we utilize only means $\bar{\gamma}_s$ and $\bar{\gamma}_r$ instead of instantaneous CSI $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$ for power control. Note that, at the end of the first slot in any frame i , the relay can estimate the average power of its received signal during this time slot, denoted as $|y(i)|^2 = E(|y(t)|^2)$, $iT \leq t \leq iT + \frac{1}{2}T$ where T is the frame length. Given $|y(i)|^2$, the relay with adaptive power control can adjust its amplification gain as $\alpha(i) = P_r/|y(i)|^2$ [80] [20]. Thus, there is no need for the instantaneous CSI on the source-relay link to adjust the relay amplification gain. Note that, in this strategy, since power control based on statistic CSI is performed only once instead of in each frame as in Strategy I, the overheads used for reporting CSI and optimization results can be significantly reduced.

We can formulate a similar optimization problem which minimizes the total power consumption of the source and the relay while ensuring the probability without violating a minimum rate above a threshold ϵ as

$$(P2) \min_{P_s, P_r} P_s + P_r \quad (2.7a)$$

$$\text{s.t. } p_{no}(P_s, P_r) = Pr\{r(P_s, P_r, \gamma_s, \gamma_r) \geq r_{min}\} \geq \epsilon, \quad (2.7b)$$

$$0 \leq P_s \leq P_s^{max}, \quad 0 \leq P_r \leq P_r^{max}. \quad (2.7c)$$

2.2.3 Strategy III

In Strategy I, overheads on control channels arise from reporting instantaneous CSI ($\tilde{\gamma}_s$ and/or $\tilde{\gamma}_r$) and optimization results (P_s^* and/or P_r^*). Such overheads happen in every frame. On the contrary, overheads are reduced to a minimum level in Strategy

II since all information exchanges only once. However, utilizing only statistical CSI in Strategy II may result in significant performance degradation, as shown in Section 2.4. In order to balance overhead reduction and performance degradation, we propose a new strategy (called Strategy III). In Strategy III, the relay is selected as the central node. The optimal transmit power of the source, P_s^* , is calculated based on statistical CSI $\bar{\gamma}_s$ and $\bar{\gamma}_r$ so that P_s^* requires to be fed back only once from the central node to the source. In each frame, the relay tunes its own transmit power P_r to guarantee a minimum transmission rate requirement based on the fixed P_s^* and instantaneous $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$. Obviously, in this strategy, only $\tilde{\gamma}_r$ needs to be transmitted from the destination to the relay in each frame, and the total overheads are K units.

Note that in each frame, collecting the full instantaneous CSI at the relay requires only one-hop feedback overheads (K units) for acquiring $\tilde{\gamma}_r$ (note that $\tilde{\gamma}_s$ is estimated at the relay itself and thus the feedback of $\tilde{\gamma}_s$ is not required). However, if the full instantaneous CSI is collected at the source, both $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$ need to be fed back to the source. The total feedback overheads are at least $2K$ units. If there is no direct link between the source and the destination, transmitting $\tilde{\gamma}_r$ from the destination to the source requires two hops and the total feedback overheads will become $3K$ units. Therefore, collecting the full instantaneous CSI at the source always leads to more feedback overheads than doing that at the relay. In addition, the overheads resulting from collecting the full instantaneous CSI at the source are no less than that of Strategy I with the relay serving as the central node (total feedback overheads are $2K$ units). In summary, the strategy that the source collects and uses the full CSI is not recommended.

In Strategy III, the source power is determined based on statistical CSI before instantaneous CSI is known, and the relay power can be treated as a “penalty” due to the insufficient source power to satisfy the rate requirement after instantaneous CSI is known. Thus, the power allocation problem can be formulated as a two-stage stochastic programming problem [75]. The optimization problem in the first stage is

$$(P3-1) \quad \min_{P_s} \quad P_s + E_{\gamma_s, \gamma_r} [Q(P_s, \gamma_s, \gamma_r)] \quad (2.8a)$$

$$\text{s.t.} \quad p_{no}(P_s) = Pr(r(P_s, P_r^{max}, \gamma_s, \gamma_r) \geq r_{min}) \geq \epsilon, \quad (2.8b)$$

$$0 \leq P_s \leq P_s^{max}, \quad (2.8c)$$

where constraint (2.8b) ensures that the probability without violating the minimum rate is no less than a threshold ϵ , and P_r^{max} in (2.8b) implies that the relay tries its best to satisfy the transmission rate. In the objective function, $\bar{P}_r^*(P_s) = E_{\gamma_s, \gamma_r} [Q(P_s, \gamma_s, \gamma_r)]$ denotes the optimal average transmit power of the relay for a given P_s , and $Q(P_s, \tilde{\gamma}_s, \tilde{\gamma}_r)$ is the optimal instantaneous transmit power of the relay given P_s , $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$. $Q(P_s, \tilde{\gamma}_s, \tilde{\gamma}_r)$ is the outcome of the second-stage optimization problem, which can be formulated as

$$(P3-2) \quad Q(P_s, \tilde{\gamma}_s, \tilde{\gamma}_r) = \min_{P_r} P_r \quad (2.9a)$$

$$\text{s.t.} \quad r(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r) \geq r_{min}, \quad \forall \tilde{\gamma}_s, \tilde{\gamma}_r \quad (2.9b)$$

$$0 \leq P_r \leq P_r^{max}. \quad (2.9c)$$

When channels are in deep fading, i.e., $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$ are small, the second-stage optimization problem may not have a feasible solution due to the limited power budget. When this case happens, the relay will stop transmitting, i.e., $P_r = 0$, since it cannot satisfy the minimum rate requirement.

2.3 Optimal power control

In this section, we will derive the optimal transmit power by solving formulated optimization problems for the three strategies.

2.3.1 Strategy I

The optimization problem (P1) can be solved by following a similar way as in [85]. From (2.3b), (2.4) and (2.5), constraint (2.3b) is equivalent to

$$\gamma(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r) = \frac{P_s \tilde{\gamma}_s P_r \tilde{\gamma}_r}{P_s \tilde{\gamma}_s + P_r \tilde{\gamma}_r + 1} \geq \gamma_{th} \quad (2.10)$$

where $\gamma_{th} = 2^{(2r_{min}/B)} - 1$ denotes the SNR threshold to guarantee the minimum transmission rate r_{min} . Since P_s , P_r , $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$ take nonnegative values, it is easy to verify that inequality (2.10) is equivalent to

$$P_r \geq \frac{\gamma_{th}(P_s \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s \tilde{\gamma}_s - \gamma_{th})} \text{ and } P_s \tilde{\gamma}_s - \gamma_{th} > 0 \text{ and } \tilde{\gamma}_r > 0 \quad (2.11)$$

Obviously, (P1) reaches optimality when (2.11) is satisfied at equality. Hence, (P1) could be simplified to include only one variable as

$$\min_{P_s} \quad P_s + \frac{\gamma_{th}(P_s \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s \tilde{\gamma}_s - \gamma_{th})} \quad (2.12a)$$

$$\text{s.t.} \quad P_s \tilde{\gamma}_s - \gamma_{th} > 0, \quad 0 \leq P_s \leq P_s^{max}, \quad (2.12b)$$

$$0 \leq P_r = \frac{\gamma_{th}(P_s \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s \tilde{\gamma}_s - \gamma_{th})} \leq P_r^{max}, \quad (2.12c)$$

$$\tilde{\gamma}_s > 0, \quad \tilde{\gamma}_r > 0. \quad (2.12d)$$

From (2.12b), the range of P_s can be determined as $P_s \in (\frac{\gamma_{th}}{\tilde{\gamma}_s}, P_s^{max}]$. Since $P_s \tilde{\gamma}_s - \gamma_{th} > 0$, constraint (2.12c) is equivalent to

$$P_s \geq \frac{\gamma_{th}(P_r^{max} \tilde{\gamma}_r + 1)}{\tilde{\gamma}_s(P_r^{max} \tilde{\gamma}_r - \gamma_{th})} \text{ and } P_r^{max} \tilde{\gamma}_r - \gamma_{th} > 0. \quad (2.13)$$

Define

$$P_s^{min} = \frac{\gamma_{th}(P_r^{max} \tilde{\gamma}_r + 1)}{\tilde{\gamma}_s(P_r^{max} \tilde{\gamma}_r - \gamma_{th})}. \quad (2.14)$$

We have $P_s^{min} = \frac{\gamma_{th}(P_r^{max} \tilde{\gamma}_r + 1)}{\tilde{\gamma}_s(P_r^{max} \tilde{\gamma}_r - \gamma_{th})} > \frac{\gamma_{th}}{\tilde{\gamma}_s}$. Hence, the range of P_s can be updated to be $P_s \in [P_s^{min}, P_s^{max}]$ when $P_r^{max} \tilde{\gamma}_r - \gamma_{th} > 0$.

Let

$$f(P_s) = P_s + \frac{\gamma_{th}(P_s \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s \tilde{\gamma}_s - \gamma_{th})}. \quad (2.15)$$

Then, the first derivative of $f(P_s)$ with respect to P_s is

$$f'(P_s) = 1 - \frac{\tilde{\gamma}_s \gamma_{th} (\gamma_{th} + 1)}{\tilde{\gamma}_r (P_s \tilde{\gamma}_s - \gamma_{th})^2}. \quad (2.16)$$

It is easy to verify that $f'(P_s)$ is an increasing function of P_s when $P_s \tilde{\gamma}_s - \gamma_{th} > 0$ (i.e., $P_s > \frac{\gamma_{th}}{\tilde{\gamma}_s}$). Let $f'(P_s) = 0$. We can obtain two roots

$$P'_s = \frac{\gamma_{th} + \sqrt{\frac{\tilde{\gamma}_s \gamma_{th} (\gamma_{th} + 1)}{\tilde{\gamma}_r}}}{\tilde{\gamma}_s}, \quad \tilde{\gamma}_s \neq 0, \quad (2.17)$$

$$P''_s = \frac{\gamma_{th} - \sqrt{\frac{\tilde{\gamma}_s \gamma_{th} (\gamma_{th} + 1)}{\tilde{\gamma}_r}}}{\tilde{\gamma}_s}, \quad \tilde{\gamma}_s \neq 0. \quad (2.18)$$

Obviously, $P''_s < \frac{\gamma_{th}}{\tilde{\gamma}_s} < P'_s$. Besides, since $\frac{\gamma_{th}}{\tilde{\gamma}_s} < P_s^{min}$ from (2.14), we have $P''_s < \frac{\gamma_{th}}{\tilde{\gamma}_s} < P_s^{min}$ and $P'_s \notin [P_s^{min}, P_s^{max}]$. For P'_s , there are three cases needing to be considered:

1) Case 1: $P'_s < P_s^{min}$. Since $\frac{\gamma_{th}}{\tilde{\gamma}_s} < P'_s$ and $f'(P_s)$ is increasing with P_s when $P_s > \frac{\gamma_{th}}{\tilde{\gamma}_s}$, $f'(P_s)$ is increasing with P_s when $P_s > P'_s$. In addition, since $f'(P'_s) = 0$, we

have $f'(P_s) > 0$ for $P_s \in [P_s^{min}, P_s^{max}]$. Therefore, $f(P_s)$ is an increasing function for $P_s \in [P_s^{min}, P_s^{max}]$. Thus, the minimum power-sum of $f(P_s)$ is achieved at $P_s = P_s^{min}$, i.e., the optimal source power is $P_s^* = P_s^{min}$.

2) Case 2: $P_s^{min} \leq P'_s \leq P_s^{max}$. Since $\frac{\gamma_{th}}{\tilde{\gamma}_s} < P_s^{min}$ and $f'(P_s)$ is increasing with P_s when $P_s > \frac{\gamma_{th}}{\tilde{\gamma}_s}$, $f'(P_s)$ is also increasing with P_s when $P_s > P_s^{min} > \frac{\gamma_{th}}{\tilde{\gamma}_s}$. Note that since $f'(P'_s) = 0$, we have $f'(P_s) < 0$ for $P_s \in [P_s^{min}, P'_s)$ and $f'(P_s) > 0$ for $P_s \in (P'_s, P_s^{max}]$. Hence, $f(P_s)$ is a decreasing function for $P_s \in [P_s^{min}, P'_s)$ and an increasing function for $P_s \in (P'_s, P_s^{max}]$. Thus, the optimal source power is achieved at $P_s^* = P'_s$.

3) Case 3: $P'_s > P_s^{max}$. Similarly, it can be verified that $f'(P_s)$ is an increasing function when $P_s > P_s^{min}$. Also, since $f'(P'_s) = 0$, we have $f'(P_s) < 0$ for $P_s \in [P_s^{min}, P_s^{max}]$. Thus, $f(P_s)$ is a decreasing function for $P_s \in [P_s^{min}, P_s^{max}]$ and the optimal source power should be $P_s^* = P_s^{max}$.

Given P_s^* , the optimal transmit power of the relay can be calculated as

$$P_r^* = \frac{\gamma_{th}(P_s^* \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s^* \tilde{\gamma}_s - \gamma_{th})}. \quad (2.19)$$

We summarize the solution procedure in Algorithm 1.

Assume that random variables γ_s and γ_r are independent of each other. The probability that the minimum rate r_{min} is not violated can be calculated as:

$$\begin{aligned} & p_{no}(P_s^{max}, P_r^{max}) \\ &= Pr\{r(P_s^{max}, P_r^{max}, \gamma_s, \gamma_r) > r_{min}\} \\ &= Pr\{\gamma(P_s^{max}, P_r^{max}, \gamma_s, \gamma_r) > \gamma_{th}\} \\ &= Pr\left\{\frac{P_s^{max} \gamma_s P_r^{max} \gamma_r}{P_s^{max} \gamma_s + P_r^{max} \gamma_r + 1} > \gamma_{th}\right\} \end{aligned}$$

Algorithm 1: Solution to Problem (P1)

Input: $\tilde{\gamma}_s, \tilde{\gamma}_r, P_s^{max}, P_r^{max}, r_{min}$
Output: P_s^*, P_r^*

- 1 Calculate P_s^{min} by (2.14) and calculate $\gamma_{th} = 2^{\frac{2r_{min}}{B}} - 1$;
- 2 **if** $(\tilde{\gamma}_s, \tilde{\gamma}_r > 0), (P_s^{max} \geq P_s^{min})$ **and** $(P_r^{max} \tilde{\gamma}_r > \gamma_{th})$ **then**
- 3 Calculate P'_s by (2.17);
- 4 **if** $P'_s < P_s^{min}$ **then**
- 5 $P_s^* = P_s^{min}$;
- 6 **else if** $P_s^{min} \leq P'_s \leq P_s^{max}$ **then**
- 7 $P_s^* = P'_s$;
- 8 **else**
- 9 $P_s^* = P_s^{max}$
- 10 Calculate P_r^* by (2.19);
- 11 **else**
- 12 $P_s^* = P_r^* = 0$ and there is no feasible solution.

$$\begin{aligned}
 &= \iint_D f_{\gamma_s}(\gamma_s) f_{\gamma_r}(\gamma_r) d\gamma_s d\gamma_r \\
 &= \int_{\frac{\gamma_{th}}{P_s^{max}}}^{\infty} \int_{g(P_s^{max}, P_r^{max}, \gamma_s)}^{\infty} f_{\gamma_s}(\gamma_s) f_{\gamma_r}(\gamma_r) d\gamma_s d\gamma_r \\
 &= \int_{\frac{\gamma_{th}}{P_s^{max}}}^{\infty} f_{\gamma_s}(\gamma_s) \left[\int_{g(P_s^{max}, P_r^{max}, \gamma_s)}^{\infty} f_{\gamma_r}(\gamma_r) d\gamma_r \right] d\gamma_s \\
 &= \int_{\frac{\gamma_{th}}{P_s^{max}}}^{\infty} \frac{1}{\bar{\gamma}_s} \exp\left(-\frac{1}{\bar{\gamma}_s} \gamma_s\right) \exp\left(-\frac{1}{\bar{\gamma}_r} g(P_s^{max}, P_r^{max}, \gamma_s)\right) d\gamma_s, \tag{2.20}
 \end{aligned}$$

where

$$\begin{aligned}
 D &= \left\{ (\gamma_s, \gamma_r) : \frac{P_s^{max} \gamma_s P_r^{max} \gamma_r}{P_s^{max} \gamma_s + P_r^{max} \gamma_r + 1} > \gamma_{th} \right\} \\
 &= \left\{ (\gamma_s, \gamma_r) : \gamma_r \geq \frac{\gamma_{th} (P_s^{max} \gamma_s + 1)}{P_r^{max} (P_s^{max} \gamma_s - \gamma_{th})}, P_s^{max} \gamma_s > \gamma_{th} \right\} \\
 &= \left\{ (\gamma_s, \gamma_r) : \gamma_r \geq g(P_s^{max}, P_r^{max}, \gamma_s), \gamma_s > \frac{\gamma_{th}}{P_s^{max}} \right\}, \tag{2.21}
 \end{aligned}$$

$$g(P_s^{max}, P_r^{max}, \gamma_s) = \frac{\gamma_{th} (P_s^{max} \gamma_s + 1)}{P_r^{max} (P_s^{max} \gamma_s - \gamma_{th})}. \tag{2.22}$$

For expression simplicity, we can rewrite (2.20) as

$$p_{no}(P_s^{max}, P_r^{max}) = G(P_s^{max}, P_r^{max}), \quad (2.23)$$

where

$$G(X, Y) = \int_{\frac{\gamma_{th}}{X}}^{\infty} \frac{1}{\bar{\gamma}_s} \exp(-\frac{1}{\bar{\gamma}_s} \gamma_s) \exp(-\frac{1}{\bar{\gamma}_r} g(X, Y, \gamma_s)) d\gamma_s, \quad (2.24)$$

$$g(X, Y, u) = \frac{\gamma_{th}(X \cdot u + 1)}{Y(X \cdot u - \gamma_{th})}. \quad (2.25)$$

2.3.2 Strategy II

The probability with which the minimum rate is not violated in Strategy II can be written as

$$p_{no}(P_s, P_r) = Pr\{\gamma(P_s, P_r, \gamma_s, \gamma_r) > \gamma_{th}\} = G(P_s, P_r), \quad (2.26)$$

where $\gamma(P_s, P_r, \gamma_s, \gamma_r) = \frac{P_s \gamma_s P_r \gamma_r}{P_s \gamma_s + P_r \gamma_r + 1}$ and $G(X, Y)$ is defined in (2.24). It is easy to prove that $p_{no}(P_s, P_r)$ is an increasing function of either P_s or P_r , since $\gamma(P_s, P_r, \gamma_s, \gamma_r)$ is an increasing function of either P_s or P_r . Therefore, by letting $P_r = P_r^{max}$ and solving equation $p_{no}(P_s, P_r^{max}) = G(P_s, P_r^{max}) = \epsilon$, we can obtain the minimum source power P_s^{min} that satisfies the non-violation probability requirement. Since $p_{no}(P_s, P_r^{max})$ is an increasing function of P_s , the equation $p_{no}(P_s, P_r^{max}) = \epsilon$ can be easily solved by using the bisection method. Then, the feasible region of Problem (P2) becomes $P_s \in [P_s^{min}, P_s^{max}]$.

For any $\hat{P}_s \in [P_s^{min}, P_s^{max}]$, we can use the bisection method again to find a minimum relay power \hat{P}_r^{min} by solving equation $p_{no}(\hat{P}_s, P_r) = \epsilon$ with P_r being a variable. If P_s is a discrete and countable variable, we can enumerate all values of \hat{P}_s

in $[P_s^{min}, P_s^{max}]$ to figure out the exact optimal solution of Problem (P2). However, since P_s is continuous in most cases, enumerating all values of $\hat{P}_s \in [P_s^{min}, P_s^{max}]$ becomes infeasible.¹

Actually, \hat{P}_r^{min} is a function of \hat{P}_s and can be denoted as $\hat{P}_r^{min}(\hat{P}_s)$. Thus, the optimal solution is the value of \hat{P}_s which minimizes the total power $P_{total}(\hat{P}_s) = \hat{P}_s + \hat{P}_r^{min}(\hat{P}_s)$ in the feasible region $\hat{P}_s \in [P_s^{min}, P_s^{max}]$. However, it is not easy to find the global minimum since there is no closed-form expression and the properties (such as monotonicity, convexity, etc.) of the function $P_{total}(\hat{P}_s)$ are unknown. Therefore, in order to reduce computational complexity, instead of seeking the global minimum of $P_{total}(\hat{P}_s)$, we propose a new method, called N -section method, to find a local minimum.

The N -section method ($N \geq 3$) can be described as follows. *Initialization*: let $A = P_s^{min}$ and $B = P_s^{max}$. *Step 1*: Partition the interval $[A, B]$ equally into N sections. Note that, there are $N + 1$ segmentation points \hat{P}_s (including two end points). *Step 2*: Calculate the total power values $P_{total}(\hat{P}_s) = \hat{P}_s + \hat{P}_r^{min}(\hat{P}_s)$ at $N + 1$ segmentation points. *Step 3*: Select the point C which provides the minimum value of $P_{total}(\hat{P}_s)$ among $N + 1$ points. If $C = A$, let B be the nearest point after C . If $C = B$, then let A be the nearest point before C . Otherwise, let A be the nearest point before C and B be the nearest point after C . Repeat Steps 1 – 3 till $B - A \leq \delta$, where δ denotes the termination threshold.

The N -section method attempts to find a local minimum point by gradually nar-

¹Note that, since independent variables (P_s and P_r) and integral variable γ_s in the integrand of $p_{no}(P_s, P_r)$ cannot be separated (here, separation means that the integrand can be represented in the form of $f_1(P_s, P_r) * f_2(\gamma_s)$), the derivative of $p_{no}(P_s, P_r)$ with respect to P_s or P_r cannot be easily obtained to further analyze the convexity and other properties of function $p_{no}(P_s, P_r)$. From this point, it is also difficult to attain the exact optimal solution of (P2).

Algorithm 2: Solution to Problem (P2)

Input: $\bar{\gamma}_s, \bar{\gamma}_r, P_s^{max}, P_r^{max}, r_{min}, \epsilon$
Output: P_s^*, P_r^*

- 1 Solve $p_{no}(P_s, P_r^{max}) = \epsilon$ by bisection method for the root $P_s = P_s^{min}$;
- 2 $A = P_s^{min}$; $B = P_s^{max}$;
- 3 **while** $B - A > \delta$ **do**
- 4 $\Delta = (B - A)/N$ where $N \geq 3$ and Δ is the interval between two segmentation points;
- 5 **for** $i = 1, 2, \dots, N + 1$ **do**
- 6 $\hat{P}_{s,i} = A + \Delta \cdot (i - 1)$;
- 7 Solve $p_{no}(\hat{P}_{s,i}, P_r) = \epsilon$ by bisection method for the root $P_r = \hat{P}_{r,i}^{min}$;
- 8 $P_{total,i} = \hat{P}_{s,i} + \hat{P}_{r,i}^{min}$;
- 9 $i^* = \arg \min_i \{P_{total,i}, i = 1, 2, \dots, N + 1\}$; $C = \hat{P}_{s,i^*}$;
- 10 **if** $C = A$ **then**
- 11 $B = \hat{P}_{s,i^*+1}$;
- 12 **else if** $C = B$ **then**
- 13 $A = \hat{P}_{s,i^*-1}$;
- 14 **else**
- 15 $A = \hat{P}_{s,i^*-1}$; $B = \hat{P}_{s,i^*+1}$;
- 16 $P_s^* = \hat{P}_{s,i^*} = C$; $P_r^* = \hat{P}_{r,i^*}^{min}$.

rowing the searching range. Let n denote the number of iterations required to meet the allowable error δ . We have $(B - A)(\frac{2}{N})^n \leq \delta$, from which we can derive the required number of iterations as $n \geq \ln(\frac{B-A}{\delta}) / \ln(\frac{N}{2})$. Let ρ denote the total number of operations required for getting a value of $P_{total}(\hat{P}_s)$ when \hat{P}_s is given. Then, each iteration needs $(N + 1)\rho$ operations for computing $P_{total}(\hat{P}_s)$ plus N comparison operations. Therefore, the entire computational complexity is $n[(N + 1)\rho + Nc]$, where c denotes the complexity for a comparison operation.

The solution procedure to Problem (P2) is summarized in Algorithm 2.

2.3.3 Strategy III

(a) Solution for $E_{\gamma_s, \gamma_r}[Q(P_s, \gamma_s, \gamma_r)]$ in Problem (P3-1)

Let us first solve the second-stage optimization problem (P3-2) to get $Q(P_s, \tilde{\gamma}_s, \tilde{\gamma}_r)$.

Similar to (2.10) and (2.11), the first constraint (2.9b) in (P3-2) is equivalent to

$$\gamma(P_s, P_r, \tilde{\gamma}_s, \tilde{\gamma}_r) = \frac{P_s \tilde{\gamma}_s P_r \tilde{\gamma}_r}{P_s \tilde{\gamma}_s + P_r \tilde{\gamma}_r + 1} \geq \gamma_{th}, \quad (2.27)$$

which can be further rewritten as

$$P_r \geq \frac{\gamma_{th}(P_s \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s \tilde{\gamma}_s - \gamma_{th})} \text{ and } P_s \tilde{\gamma}_s - \gamma_{th} > 0 \text{ and } \tilde{\gamma}_r > 0. \quad (2.28)$$

Obviously, (P3-2) achieves optimality when (2.28) is satisfied at equality, i.e.,

$$Q(P_s, \tilde{\gamma}_s, \tilde{\gamma}_r) = P_r^* = \frac{\gamma_{th}(P_s \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s \tilde{\gamma}_s - \gamma_{th})}. \quad (2.29)$$

Note that the above equality holds only if at least one feasible solution exists in (P3-2). For clarity, the conditions for which problem (P3-2) has feasible solutions are listed as follows:

$$P_s \tilde{\gamma}_s - \gamma_{th} > 0 \text{ and } \tilde{\gamma}_r > 0, \quad (2.30)$$

$$0 \leq P_r^* = \frac{\gamma_{th}(P_s \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s \tilde{\gamma}_s - \gamma_{th})} \leq P_r^{max}, \quad (2.31)$$

where (2.30) and (2.31) are derived from (2.28) and (2.9c), respectively.

Combining (2.30) and (2.31), we can simplify the conditions as

$$\tilde{\gamma}_s > \gamma_{th}/P_s > 0 \text{ and } P_s > 0, \quad (2.32)$$

$$\tilde{\gamma}_r \geq g(P_s, P_r^{max}, \tilde{\gamma}_s), \quad (2.33)$$

where $g(X, Y, u)$ is defined in (2.25).

In summary, the optimal objective function of (P3-2) $Q(P_s, \tilde{\gamma}_s, \tilde{\gamma}_r) = P_r^*$ is given by (2.29) when conditions (2.32) and (2.33) are satisfied. Otherwise, $Q(P_s, \tilde{\gamma}_s, \tilde{\gamma}_r) = P_r^* = 0$, i.e., the relay does not forward data from the source. Thus, the average power of the relay can be calculated as

$$\begin{aligned}
 \bar{P}_r^*(P_s) &= E_{\gamma_s, \gamma_r}[Q(P_s, \gamma_s, \gamma_r)] \\
 &= \int_{\frac{\gamma_{th}}{P_s}}^{\infty} \int_{g(P_s, P_r^{max}, \gamma_s)}^{\infty} Q(P_s, \gamma_s, \gamma_r) f_{\gamma_s}(\gamma_s) f_{\gamma_r}(\gamma_r) d\gamma_s d\gamma_r \\
 &= \int_{\frac{\gamma_{th}}{P_s}}^{\infty} f_{\gamma_s}(\gamma_s) \left[\int_{g(P_s, P_r^{max}, \gamma_s)}^{\infty} Q(P_s, \gamma_s, \gamma_r) f_{\gamma_r}(\gamma_r) d\gamma_r \right] d\gamma_s \\
 &= \frac{\gamma_{th}}{\bar{\gamma}_s \bar{\gamma}_r} \int_{\frac{\gamma_{th}}{P_s}}^{\infty} \frac{P_s \gamma_s + 1}{P_s \gamma_s - \gamma_{th}} \exp\left(-\frac{1}{\bar{\gamma}_s} \gamma_s\right) \left[\int_{g(P_s, P_r^{max}, \gamma_s)}^{\infty} \frac{1}{\bar{\gamma}_r} \exp\left(-\frac{1}{\bar{\gamma}_r} \gamma_r\right) d\gamma_r \right] d\gamma_s. \quad (2.34)
 \end{aligned}$$

Let $t = \frac{1}{\bar{\gamma}_r} \gamma_r$. We have

$$\begin{aligned}
 &\int_{g(P_s, P_r^{max}, \gamma_s)}^{\infty} \frac{1}{\bar{\gamma}_r} \exp\left(-\frac{1}{\bar{\gamma}_r} \gamma_r\right) d\gamma_r \\
 &= \int_{g(P_s, P_r^{max}, \gamma_s)}^{\infty} \frac{\frac{1}{\bar{\gamma}_r}}{\frac{1}{\bar{\gamma}_r} \gamma_r} \exp\left(-\frac{1}{\bar{\gamma}_r} \gamma_r\right) d\gamma_r \\
 &= \int_{\frac{1}{\bar{\gamma}_r} g(P_s, P_r^{max}, \gamma_s)}^{\infty} \frac{e^{-t}}{t} dt = E_1\left(\frac{1}{\bar{\gamma}_r} g(P_s, P_r^{max}, \gamma_s)\right), \quad (2.35)
 \end{aligned}$$

where $E_1(x) = \int_x^{\infty} \frac{e^{-t}}{t} dt$ is the exponential integral function. Substituting (2.35) into (2.34), we have

$$\begin{aligned}
 \bar{P}_r^*(P_s) &= E_{\gamma_s, \gamma_r}[Q(P_s, \gamma_s, \gamma_r)] \\
 &= \frac{\gamma_{th}}{\bar{\gamma}_s \bar{\gamma}_r} \int_{\frac{\gamma_{th}}{P_s}}^{\infty} \frac{P_s \gamma_s + 1}{P_s \gamma_s - \gamma_{th}} \exp\left(-\frac{\gamma_s}{\bar{\gamma}_s}\right) E_1\left(\frac{g(P_s, P_r^{max}, \gamma_s)}{\bar{\gamma}_r}\right) d\gamma_s. \quad (2.36)
 \end{aligned}$$

(b) Determine the range of source power P_s in (P3-1)

The probability with which the minimum rate is not violated in Strategy III can be written as

$$p_{no}(P_s, P_r^{max}) = Pr\{\gamma(P_s, P_r^{max}, \gamma_s, \gamma_r) > \gamma_{th}\} = G(P_s, P_r^{max}), \quad (2.37)$$

where $\gamma(P_s, P_r^{max}, \gamma_s, \gamma_r) = \frac{P_s \gamma_s P_r^{max} \gamma_r}{P_s \gamma_s + P_r^{max} \gamma_r + 1}$ and $G(X, Y)$ are defined in (2.24).

Similar to Strategy II, it can be verified that $p_{no}(P_s, P_r^{max})$ is an increasing function of P_s . Thus, the minimum source power P_s^{min} can be derived from the equation $p_{no}(P_s, P_r^{max}) = \epsilon$, which can be easily solved by using the bisection method. Hence, the feasible region of Problem (P3-1) is $P_s \in [P_s^{min}, P_s^{max}]$.

(c) Solve Problem (P3-1) and find the optimal source power P_s^*

Similar to Strategy II, the N -section method can also be applied to Strategy III to find a local minimum point of function $P_{total}(\hat{P}_s) = \hat{P}_s + \bar{P}_r^*(\hat{P}_s)$. The only difference is in Step 2, where $\bar{P}_r^*(\hat{P}_s)$ can be calculated directly from (2.36) for Strategy III, while for Strategy II, the calculation of \hat{P}_r^{min} requires the bisection method to solve equation $p_{no}(\hat{P}_s, P_r) = \epsilon$. Strategy II needs a two-level nested loop (one is for the N -section method and the other is for the bisection method) while Strategy III needs only a single-level loop for the N -section method. This implies that Strategy III has lower computational complexity compared with Strategy II.

(d) The optimal relay power P_r^* :

After the optimal source power P_s^* is derived, if conditions (2.32) and (2.33) hold, then the optimal relay power P_r^* is given by (2.29) by letting $P_s = P_s^*$, i.e., $P_r^* = \frac{\gamma_{th}(P_s^* \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s^* \tilde{\gamma}_s - \gamma_{th})}$. Otherwise, $P_r^* = 0$. Note that, in Strategy III, P_s^* remains the same while P_r^* may vary from one frame to another according to the channel fading realizations $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$.

For clarity, we summarize the power allocation procedure for Strategy III in Algorithm 3.

Algorithm 3: Power allocation procedure for Strategy III

Input: $\tilde{\gamma}_s, \tilde{\gamma}_r, \tilde{\gamma}_s, \tilde{\gamma}_r, P_s^{max}, P_r^{max}, r_{min}, \epsilon$
Output: P_s^*, P_r^*

- 1 Solve $p_{no}(P_s, P_r^{max}) = \epsilon$ by bisection method for the root $P_s = P_s^{min}$;
- 2 $A = P_s^{min}$; $B = P_s^{max}$;
- 3 **while** $B - A > \delta$ **do**
- 4 $\Delta = (B - A)/N$ where $N \geq 3$ and Δ is the interval between two segmentation points;
- 5 **for** $i = 1, 2, \dots, N + 1$ **do**
- 6 $\hat{P}_{s,i} = A + \Delta(i - 1)$;
- 7 Calculate $\bar{P}_{r,i}^* = \bar{P}_r^*(P_s)|_{P_s=\hat{P}_{s,i}}$ by (2.36);
- 8 $P_{total,i} = \hat{P}_{s,i} + \bar{P}_{r,i}^*$;
- 9 $i^* = \underset{i}{\operatorname{argmin}} \{P_{total,i}, i = 1, 2, \dots, N + 1\}$; $C = \hat{P}_{s,i^*}$;
- 10 **if** $C = A$ **then**
- 11 $B = \hat{P}_{s,i^*+1}$;
- 12 **else if** $C = B$ **then**
- 13 $A = \hat{P}_{s,i^*-1}$;
- 14 **else**
- 15 $A = \hat{P}_{s,i^*-1}$; $B = \hat{P}_{s,i^*+1}$;
- 16 $P_s^* = \hat{P}_{s,i^*} = C$;
- 17 **if** $\tilde{\gamma}_s > \frac{\gamma_{th}}{P_s^*} > 0$, $P_s^* > 0$ and $\tilde{\gamma}_r \geq g(P_s^*, P_r^{max}, \tilde{\gamma}_s)$ **then**
- 18 $P_r^* = \frac{\gamma_{th}(P_s^* \tilde{\gamma}_s + 1)}{\tilde{\gamma}_r(P_s^* \tilde{\gamma}_s - \gamma_{th})}$;
- 19 **else**
- 20 $P_r^* = 0$;

2.4 Simulation Results

In this section, simulation results are demonstrated to verify the effectiveness of the proposed algorithms, and to compare three strategies in terms of the total power consumption. In simulations, the SNR threshold is set to $\gamma_{th} = 10$ dB. The power budgets of the source and the relay are $P_s^{max} = P_r^{max} = 33$ dBm, and the non-violation probability threshold is $\epsilon = 0.95$.

For Strategy I, the total power consumption is obtained by averaging the sum $P_s^* + P_r^*$ over 10^6 Monte Carlo simulations. For Strategy II, the total power consumption

is $P_{total}^* = P_s^* + P_r^*$, and for Strategy III, $P_{total}^* = P_s^* + \bar{P}_r^*$. Note that, in Strategy I, since P_s^* and P_r^* change slot by slot according to instantaneous channel realizations $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$, we use their mean values for comparisons. Similarly, the mean value of P_r^* is used for comparisons in Strategy III.

2.4.1 Effectiveness of the Proposed Algorithms

For evaluation purpose, we introduce an enumeration method as the comparison benchmark. In the enumeration method, we discretize values of P_s and/or P_r within $[0, P_s^{max}]$ and/or $[0, P_r^{max}]$, respectively, with a resolution of 10^{-6} . For Strategy I, we check if constraint (2.3b) is satisfied by calculating (2.4) for each pair of (P_s, P_r) with $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$ in each simulation. For Strategy II, constraint (2.7b) is checked by calculating (2.26) for each pair of (P_s, P_r) with $\bar{\gamma}_s$ and $\bar{\gamma}_r$. In Strategies I and II, the optimal solution is the pair of (P_s, P_r) that is feasible and results in the minimum total power $P_s + P_r$. For Strategy I, 10^6 Monte Carlo simulations are conducted for each channel condition $(\bar{\gamma}_s, \bar{\gamma}_r)$ to get \bar{P}_s and \bar{P}_r . For Strategy III, we check if constraint (2.8b) is satisfied by calculating (2.37) for each pair of (P_s, P_r^{max}) with $\bar{\gamma}_s$ and $\bar{\gamma}_r$. For each feasible P_s , P_r^* is derived by (2.29) when conditions (2.32) and (2.33) are satisfied. Otherwise, $P_r^* = 0$. Similarly, 10^6 Monte Carlo simulations are conducted for each feasible P_s to get \bar{P}_r^* . The optimal source power P_s^* is the one resulting in the minimum total power $P_s + \bar{P}_r^*$.

Fig. 2.1 shows the transmit powers of the source and the relay versus the average channel gain-to-noise ratio on the link from the source to the relay $\bar{\gamma}_s$ while the average channel gain-to-noise ratio from the relay to the destination is set to $\bar{\gamma}_r =$

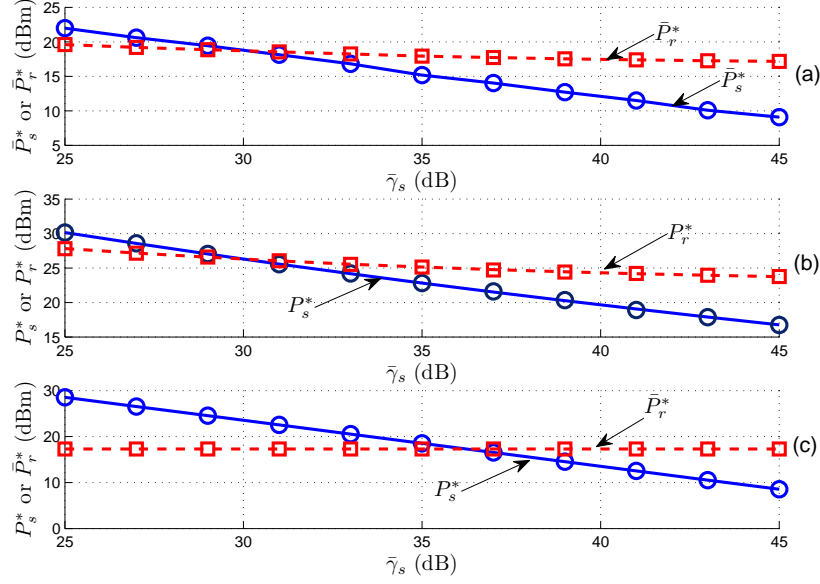


Figure 2.1: Power consumptions under three strategies when $\bar{\gamma}_r = 30$ dB: (a) Strategy I; (b) Strategy II; (c) Strategy III. The results obtained by the proposed method are denoted by solid or dashed lines, and the enumeration method by circles or squares.

30 dB. Note that, in Fig. 2.1, the low channel gain region (i.e., $\bar{\gamma}_s < 25$ dB) is not shown because they cannot satisfy the target non-violation probability ($\epsilon = 0.95$) in our simulation settings even using the maximum power $P_s^{max} = P_r^{max} = 33$ dBm. From this figure, we can see that the simulation results by the proposed methods are the same as those by the enumeration method, which verify the effectiveness of the proposed methods. Although the N -section method applying for Strategy II (Algorithm 2) and Strategy III (Algorithm 3) is sub-optimal, it leads to numerically optimal solutions in our simulations. This is because the N -section method can converge to the optimal solution if the function $P_{total}(\hat{P}_s)$ has only one local minimum point or is monotonic in $[P_s^{min}, P_s^{max}]$. Although changing simulation settings showed

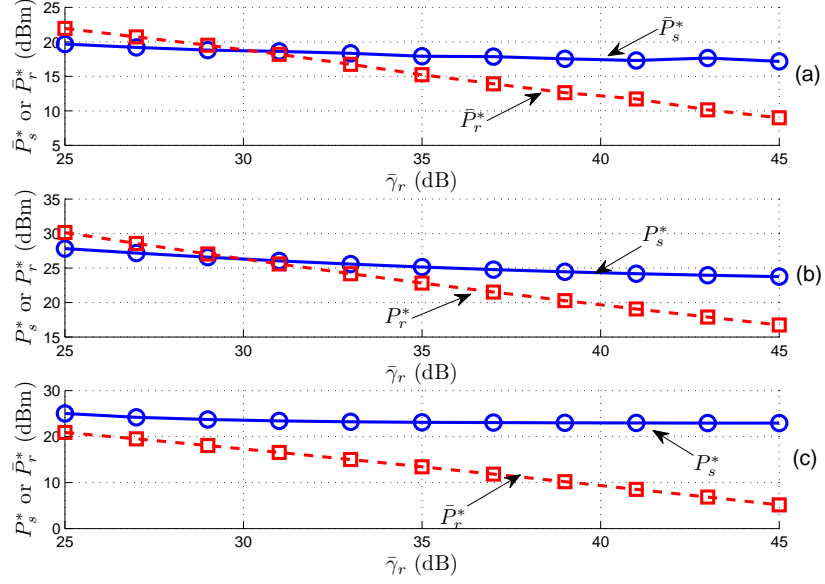


Figure 2.2: Power consumptions under three strategies when $\bar{\gamma}_s = 30$ dB: (a) Strategy I; (b) Strategy II; (c) Strategy III. The results obtained by the proposed method are denoted by solid or dashed lines, and the enumeration method by circles or squares.

the same observation, we cannot theoretically prove the existence of these properties due to the very complicated expressions.

We repeat similar simulations by fixing $\bar{\gamma}_s = 30$ dB and varying $\bar{\gamma}_r$, as shown in Fig. 2.2. The same observations can be obtained. From both figures, we can observe an interesting phenomenon. The relay power consumption (\bar{P}_r^* or P_r^*) does not change much with $\bar{\gamma}_s$ when $\bar{\gamma}_r$ is fixed (see Fig. 2.1), while the source power consumption (\bar{P}_s^* or P_s^*) remains almost invariant when $\bar{\gamma}_s$ is fixed (see Fig. 2.2). This means that the relay power consumption is mainly determined by the CSI on the link from the relay to the destination γ_r , and the source power consumption is mainly determined by the CSI on the link from the source to the relay γ_s . This can be explained from

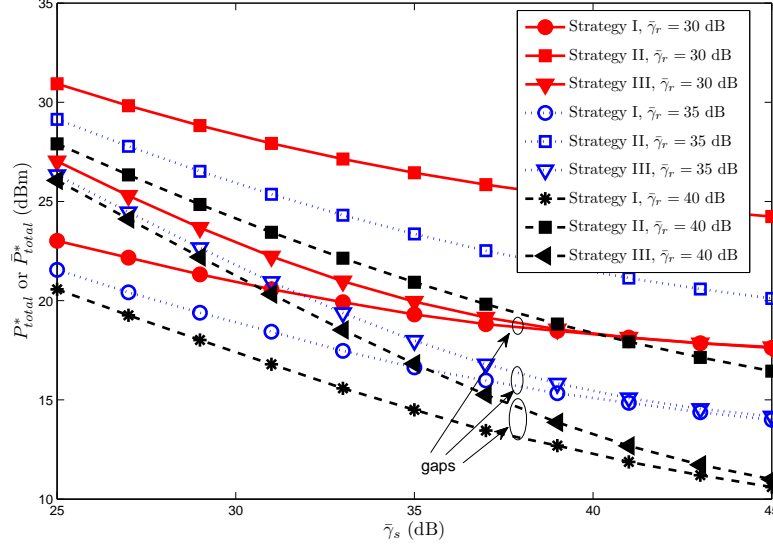


Figure 2.3: Power consumption versus $\bar{\gamma}_s$ under three strategies.

the expression of SNR, which is calculated by $\gamma = \frac{P_s \gamma_s P_r \gamma_r}{P_s \gamma_s + P_r \gamma_r + 1}$. The product of P_s and γ_s as a whole affects the value of SNR, and so does the product of P_r and γ_r . It is easy to show that $P_s \gamma_s > \gamma$, $P_r \gamma_r > \gamma$. Hence, in order to satisfy the SNR threshold requirement, i.e., $\gamma > \gamma_{th}$, it must require $P_s \gamma_s > \gamma_{th}$ and $P_r \gamma_r > \gamma_{th}$.

2.4.2 Comparisons among Three Strategies in Terms of Total Power Consumption

Figs. 2.3 and 2.4 show comparisons of three strategies in terms of total power consumption when varying γ_s and γ_r , respectively. From these two figures, we can see that the total power consumption of Strategy I is always lowest, Strategy II is always highest, and Strategy III is always between the first two strategies. This is because Strategy I utilizes instantaneous CSI $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$ to adapt P_s and P_r , while

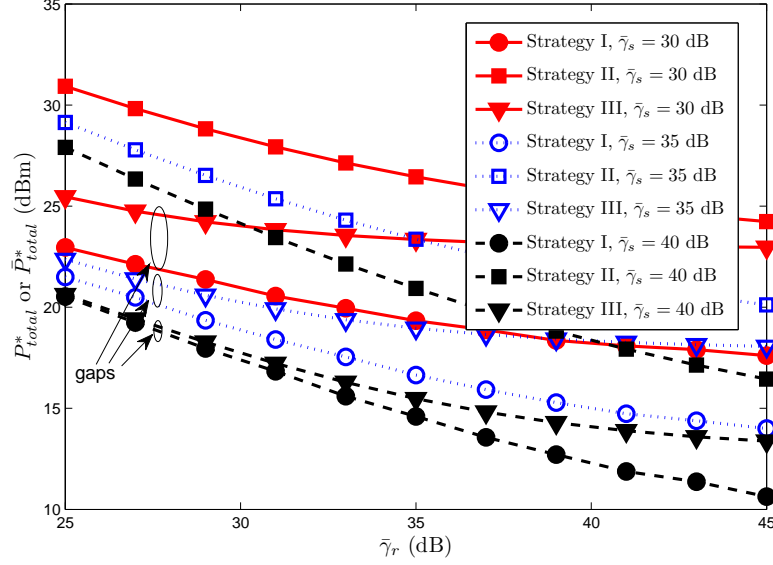


Figure 2.4: Power consumption versus $\bar{\gamma}_r$ under three strategies.

Strategy II only utilizes statistical CSI $\bar{\gamma}_s$ and $\bar{\gamma}_r$. For Strategy III, P_s is adjusted according to $\bar{\gamma}_s$ and $\bar{\gamma}_r$, but P_r is adjusted based on $\tilde{\gamma}_s$ and $\tilde{\gamma}_r$. Therefore, Strategy III well balances the tradeoff between Strategy I and Strategy II. The gap of power consumptions between Strategy II and Strategy I remains almost stable at around 8 dB (i.e., 6.3 times) independently of the difference between $\bar{\gamma}_s$ and $\bar{\gamma}_r$. However, it is not the case for Strategy III. The gap of power consumptions between Strategy III and Strategy I decreases with $\bar{\gamma}_s$ when giving a fixed $\bar{\gamma}_r$, and increases with $\bar{\gamma}_r$ when $\bar{\gamma}_s$ is fixed. In other words, the gap between Strategy III and Strategy I is decreasing with $\bar{\gamma}_s - \bar{\gamma}_r$. In addition, when $\bar{\gamma}_s - \bar{\gamma}_r$ is large enough (around 7 dB above), Strategy III almost achieves the same performance as Strategy I. The main reason is as follows. From Section 2.4.1, we have known that the product $P_s \gamma_s$ ($P_r \gamma_r$) as a whole affects the value of SNR and thus a larger $\bar{\gamma}_s$ ($\bar{\gamma}_r$) results in a smaller P_s (P_r). Therefore,

when $\bar{\gamma}_s$ is much larger than $\bar{\gamma}_r$, P_s is much smaller than P_r so that P_r dominates the total power consumption. In addition, since the relay powers in both Strategy III and Strategy I are based on complete instantaneous CSI, it is expected that the derived values for P_r approach each other in two strategies. On the contrary, when $\bar{\gamma}_s$ is much smaller than $\bar{\gamma}_r$, P_s becomes the major contributor to the total power consumption. Since P_s in Strategy III is based on statistical CSI instead of instantaneous CSI as in Strategy I, an obvious performance gap between two strategies exists.

2.4.3 Distinctiveness of Strategy III and its Applicable Area

Table 2.1 compares the total power consumptions when we exchange the values of $\bar{\gamma}_s$ and $\bar{\gamma}_r$. The numbers on the left side of slashes represent the total power consumption (dBm) when $(\bar{\gamma}_s, \bar{\gamma}_r)$ take values of the left-most column, while the numbers on the right side of slashes correspond to that when $(\bar{\gamma}_s, \bar{\gamma}_r)$ take values of the right-most column. For example, in this table, "26.4/28.9" means that the total power consumption is 26.4 dBm and 28.9 dBm when $(\bar{\gamma}_s = 30 \text{ dB}, \bar{\gamma}_r = 25 \text{ dB})$ and $(\bar{\gamma}_s = 25 \text{ dB}, \bar{\gamma}_r = 30 \text{ dB})$, respectively.

From Table 2.1, we can see that, for Strategy I and Strategy II, the total power consumption remains the same after exchanging the values of $\bar{\gamma}_s$ and $\bar{\gamma}_r$. This is because both P_s^* and P_r^* are treated equally in these two strategies. Moreover, from the expression of SNR $\gamma = \frac{P_s \gamma_s P_r \gamma_r}{P_s \gamma_s + P_r \gamma_r + 1}$, the products of $P_s \gamma_s$ and $P_r \gamma_r$ play the same role for SNR. Hence, if we exchange the values of $\bar{\gamma}_s$ and $\bar{\gamma}_r$, we only need to exchange the values of P_s^* and P_r^* to achieve optimality so that the sum of power consumptions $P_{total}^* = P_s^* + P_r^*$ keeps the same.

Table 2.1: Comparisons of total power consumptions when exchanging the values of $\bar{\gamma}_s$ and $\bar{\gamma}_r$

$(\bar{\gamma}_s, \bar{\gamma}_r)(\text{dB})$	Strategy I	Strategy II	Strategy III	$(\bar{\gamma}_s, \bar{\gamma}_r)(\text{dB})$
(30, 25)	24.0 / 24.0	32.1 / 32.1	26.4 / 28.9	(25, 30)
(30, 35)	19.9 / 19.9	27.1 / 27.1	23.5 / 21.0	(35, 30)
(30, 40)	18.7 / 18.7	25.6 / 25.6	23.2 / 18.8	(40, 30)
(35, 25)	22.7 / 22.7	30.6 / 30.6	23.5 / 28.3	(25, 35)
(35, 40)	15.6 / 15.6	22.1 / 22.1	18.5 / 16.3	(40, 35)

Notes: The numbers on the left side of slashes represent the total power consumption (dBm) when $(\bar{\gamma}_s, \bar{\gamma}_r)$ take values of the left-most column, while the numbers on the right side of slashes correspond to that when $(\bar{\gamma}_s, \bar{\gamma}_r)$ take values of the right-most column.

However, this is not true for Strategy III due to the different treatments for P_s^* and P_r^* . From Table 2.1, we can see that the total power consumption will become smaller if $\bar{\gamma}_s > \bar{\gamma}_r$ after exchanging, and vice versa. Namely, a larger $\bar{\gamma}_s$ leads to better performance. This is consistent with the conclusion in Section 2.4.2 that the gap between Strategy III and Strategy I is decreasing with $\bar{\gamma}_s - \bar{\gamma}_r$. In Strategy III, P_s^* has less flexibility than P_r^* since P_s^* is based on statistical CSI and remains invariant while P_r^* is based on instantaneous CSI and changes frame by frame. Thus, we should give priority to P_s^* instead of P_r^* in order to achieve a smaller average total power consumption. In other words, if applying Strategy III for power control, we should make the relay station close to the source instead of the destination when deploying a relay station or selecting a relay node so that $\bar{\gamma}_s$ becomes larger and P_s^* becomes smaller.

The above observation can be further verified in Fig. 2.5. In this figure, we

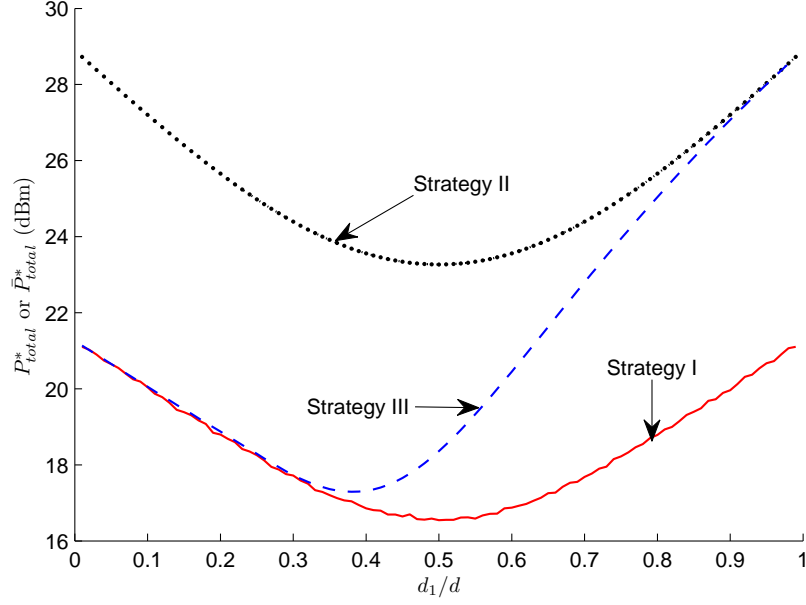


Figure 2.5: Total power consumption versus the distance ratio d_1/d .

simulate a scenario where the relay is placed on the line between the source and the destination and moves from one end to the other. The distance between the source and the destination is set to $d = 1000$ m. The distances from the relay to the source and the destination are denoted as d_{sr} and d_{rd} , respectively, where $d_{rd} = d - d_{sr}$. The average channel gains are set to $E(|h_{sr}|^2) = d_{sr}^{-\alpha}$ and $E(|h_{rd}|^2) = d_{rd}^{-\alpha}$ where $\alpha = 4$ is the path loss exponent. The noise power is $\sigma_r^2 = \sigma_d^2 = N_0 B$ where the noise power spectral density is set to $N_0 = -174$ dBm/Hz and the channel bandwidth is set to $B = 1$ MHz. The average channel gain-to-noise ratios are $\bar{\gamma}_s = E(|h_{sr}|^2)/\sigma_r^2$ and $\bar{\gamma}_r = E(|h_{rd}|^2)/\sigma_d^2$. Fig. 2.5 shows the total power consumptions versus the distance ratio d_1/d under Strategies I, II and III where $d_1 = d_{sr}$. It can be clearly seen that the total power consumption of Strategy III approaches that of Strategy I when the relay is close to the source (i.e., $d_1/d < 0.5$). This is because the relay power P_r dominates

the total power consumption in this case due to small $\bar{\gamma}_r$ while P_r is treated equally in both Strategies I and III (i.e., based on complete instantaneous CSI). When the relay is close to the destination (i.e., d_1/d is close to 1), the total power consumption of Strategy III approaches that of Strategy II since P_s dominates the total power consumption in this case and P_s in both Strategies II and III is based on statistic CSI. Therefore, it is also recommended to deploy the relay station close to the source rather than the destination when applying Strategy III.

Chapter 3

An Analysis Framework for Buffer-Aided Relaying under Time-Correlated Fading Channels

In this chapter, we will consider a typical three-node buffer-aided system [89] [90], as shown in Fig. 3.1, and analyze its performance under time-correlated fading channels in terms of average throughput, outage probability, and end-to-end delay. The key element of performance analysis is to examine the queueing behavior of packets in the relay's buffer (i.e., the distribution of buffer occupancy $Q(i)$). However, unlike independent and identically distributed (i.i.d.) fading channels, correlated fading brings challenges in performance analysis since the transition probabilities of buffer occupancy states become time-variant. To avoid this difficulty, we first define an aggregate chain $Y(i) = (\beta(i), Q(i))$ which integrates both channel state $\beta(i)$ and buffer state $Q(i)$. We then analyze the stationary distribution of the aggregate

chain $Y(i)$, and finally extract the stationary distribution of $Q(i)$ from $Y(i)$. Two delay-controllable link scheduling/selection policies are considered for infinite and finite buffers, respectively. For both policies, the aggregate chain $Y(i)$ is formulated as quasi-birth-death (QBD) chains with infinite and finite lengths. The traditional matrix-geometric method [91] [92] and a modified method based on [93] are used for solving the stationary distribution of $Y(i)$ under two policies, respectively.

In the literature, buffer-aided relaying has attracted a lot of interest from researchers. Zlatanov and Schober investigated a three-node buffer-aided relaying system which consists of a source, a relay with a buffer and a destination, for variable-rate [89], fixed-rate and mixed-rate transmissions [90]. Later, Jamali *et al.* extended these works to bidirectional buffer-aided relaying systems in [94]–[96]. Besides, Zafar *et al.* discussed the scenarios that two source-destination pairs share a single relay with a buffer in [97] and a source broadcasts to multiple destinations via a relay in [98]. A network where multiple sources transmit data to a common destination via multiple relays with buffers was discussed in [99]. In addition, multi-relay transmissions [28]–[30] and cognitive relay transmissions [100] [101] assisted by buffers were also investigated. To our best knowledge, almost all existing works about buffered-aided relaying in the literature assumed independent and identically distributed (i.i.d.) fading channels due to their simplicity and tractability. Only in [36], time-correlated fading channels are briefly mentioned without detailed performance analyses.

Notations: Throughout this chapter, scalars are denoted by lowercase letters, vectors by lowercase and bold letters, and matrices by uppercase and bold letters. \mathbf{I} stands for an identity matrix with an appropriate dimension, $\mathbf{0}$ represents an all-zero

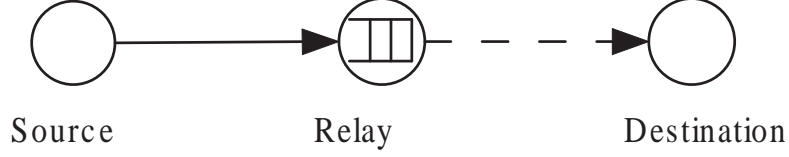


Figure 3.1: System model.

matrix, and $\mathbf{1}$ denotes an all-one column vector. $(\cdot)^{-1}$, $(\cdot)^T$ and $\|\cdot\|_\infty$ denote the inverse, transpose and infinity norm of a matrix, respectively.

3.1 System Model

The considered buffer-aided relay network consists of a source (S), a half-duplex decode-and-forward relay (R) with a buffer, and a destination (D), as shown in Fig. 3.1. We consider the scenario that there is no direct link between the source and the destination, which may happen when the source is far away from the destination or there exists a barrier between them. In this relay system, the source first sends data packets to the relay. Then, the relay decodes the received packets, stores them in its buffer, and eventually forwards them to the destination. A time-slot structure is considered where each slot has an equal length. Based on the half-duplex relaying protocol, the relay can only receive a packet from the source or send a packet to the destination in a given slot, i.e, the relay simultaneously receiving and sending packets is not allowed. Whether the relay receives or sends a packet depends on the designed link scheduling policy, which will be discussed later. Throughout this chapter, we assume that there is always data available for transmissions at the source [89] [90].

3.1.1 Correlated Fading Channel Model

Consider a slow fading channel where the instantaneous signal-to-noise ratio (SNR) remains the same in a slot. Let $h_j(i)$ denote the instantaneous channel coefficient during slot i ($i = 1, 2, \dots$) over link j ($j = 0, 1$) where $j = 0$ and $j = 1$ stand for the source-relay and relay-destination links, respectively. P_j and σ_j^2 represent the transmit power at the transmitter and the noise power at the receiver over link j , respectively. Then, the instantaneous received SNR equals $\gamma_j(i) = \frac{P_j}{\sigma_j^2} E(|h_j(i)|^2) |\alpha_j(i)|^2$ where $|\alpha_j(i)|^2 = |h_j(i)|^2 / E(|h_j(i)|^2)$ is the normalized channel gain with a unit mean. Here, $E(\cdot)$ denotes the expectation operation. A widely accepted method for determining a packet's success/failure is to compare the instantaneous SNR $\gamma_j(i)$ with a threshold γ_j^{th} [90]. If the received SNR is above the threshold, the packet can be successfully decoded. Otherwise, the packet is lost. Then, we can use a binary variable $\beta_j(i)$ to represent the packet's success/failure as

$$\beta_j(i) = \begin{cases} 1 & \text{if } \gamma_j(i) > \gamma_j^{th} \text{ or } |\alpha_j(i)|^2 > 1/F_j, \\ 0 & \text{if } \gamma_j(i) \leq \gamma_j^{th} \text{ or } |\alpha_j(i)|^2 \leq 1/F_j. \end{cases} \quad (3.1)$$

In (3.1), $F_j = \frac{P_j}{\sigma_j^2} E(|h_j(i)|^2) \frac{1}{\gamma_j^{th}}$ is called the fading margin, i.e., the maximum fading attenuation that the system allows for transmitting one packet successfully. The packet success/failure process $\beta_j(i)$ can be modeled as a first-order Markov chain where the transition probability matrix equals [31]–[33]

$$\mathbf{M}_j = \begin{bmatrix} p_j & 1 - p_j \\ 1 - q_j & q_j \end{bmatrix}, \quad (3.2)$$

where $p_j = \text{Pr}\{\beta_j(i) = 1 | \beta_j(i-1) = 1\}$ and $1 - q_j = \text{Pr}\{\beta_j(i) = 1 | \beta_j(i-1) = 0\}$ denote the probabilities that the packet in slot i is decoded successfully, given that the

packet in slot $i - 1$ was decoded successfully or unsuccessfully, respectively. Note that when $p_j = 1 - q_j$, the packet's success/failure process is reduced to be independent and identically distributed (i.i.d.), i.e., $P_r\{\beta_j(i) = 1\} = p_j = 1 - q_j$ regardless of $\beta_j(i - 1) = 1$ or 0. Given the transition probability matrix \mathbf{M}_j , the average packet error rate over link j is given by [31]

$$P_{e,j} = \frac{1 - p_j}{2 - p_j - q_j}. \quad (3.3)$$

For Rayleigh fading, we have [31]

$$P_{e,j} = 1 - e^{-1/F_j}, \quad (3.4)$$

$$q_j = 1 - \frac{\mathcal{Q}(\theta_j, \rho_j \theta_j) - \mathcal{Q}(\rho_j \theta_j, \theta_j)}{e^{1/F_j} - 1}, \quad (3.5)$$

where

$$\theta_j = \sqrt{\frac{2/F_j}{1 - \rho_j^2}}, \quad (3.6)$$

$\rho_j = J_0(2\pi f_{D,j}T)$ is the correlation coefficient between two successive channel gains $\alpha_j(i - 1)$ and $\alpha_j(i)$, $f_{D,j}$ is the Doppler frequency, T is the length of a time slot, $f_{D,j}T$ is the normalized Doppler frequency, and $J_0(\cdot)$ is the Bessel function of the first kind and zeroth order. $\mathcal{Q}(\cdot, \cdot)$ is the Marcum \mathcal{Q} function defined by

$$\mathcal{Q}(x, y) = \int_y^\infty e^{-\frac{x^2+w^2}{2}} I_0(xw) w dw, \quad (3.7)$$

where $I_0(\cdot)$ is the modified Bessel function of the first kind and zeroth order. Therefore, given $f_{D,j}T$ and F_j , the Markov parameters q_j can be uniquely determined by (3.5). In addition, p_j can be derived by using (3.3) as $p_j = \frac{1 - P_{e,j}(2 - q_j)}{1 - P_{e,j}}$.

Since the time correlation of channel fading depends only on the normalized Doppler frequency of $f_{D,j}T$, the fading channel models with different degrees of correlation can be established by choosing different values of $f_{D,j}T$. When $f_{D,j}T$ is small (< 0.1), the channel fading process is highly correlated (“slow” fading), while for a large value of $f_{D,j}T$ (> 0.2), the channel fading process becomes almost independent (“fast” fading) [31] [32]. In many scenarios, the channel can be considered to be “slow” fading. For example, for a carrier frequency of 2.1 GHz and a mobile speed of 20 km/h, the “slow” fading condition of $f_{D,j}T < 0.1$ is always satisfied when $T < 2.57$ ms ($f_{D,j} = v/\lambda$ where v is the mobile speed and λ is the carrier wavelength). In the “slow” fading case, the dependence between transmission successes/failures of consecutive packets cannot be neglected.

3.1.2 Link Scheduling Policy

In this chapter, we consider link scheduling policies which are able to limit transmission delay. In buffer-aided relaying, the end-to-end transmission delay of a packet from the source to the destination mainly results from the queueing time in the relay buffer. Therefore, in order to limit transmission delay, two methods can be applied [89]: i) starving the buffer (i.e., reducing the arrival rate and increasing the departure rate) if the buffer size is infinite, and ii) limiting the buffer size. Based on these two methods, we can introduce the following delay-controllable policies.

Policy I: Infinite buffer size. Let $d_i \in \{0, 1\}$ denote a link scheduling/selection indicator in slot i : $d_i = 0$ means the source-relay link is scheduled/selected for transmissions (i.e., the source transmits and the relay receives) and $d_i = 1$ means the

relay-destination link is selected (i.e., the relay transmits and the destination receives). $Q(i)$ denotes the number of packets in the relay's buffer at the beginning of slot i . Then, this policy can be described as [90]: If $Q(i) = 0$, then $d_i = 0$; Otherwise, d_i is given as

$$d_i = \begin{cases} \varepsilon, & \beta_0(i) = 0 \text{ and } \beta_1(i) = 0, \\ 1, & \beta_0(i) = 0 \text{ and } \beta_1(i) = 1, \\ 0, & \beta_0(i) = 1 \text{ and } \beta_1(i) = 0, \\ \mathcal{C}, & \beta_0(i) = 1 \text{ and } \beta_1(i) = 1, \end{cases} \quad (3.8)$$

where $\mathcal{C} \in \{0, 1\}$ denotes the possible outcomes of coin flipping, i.e., \mathcal{C} is a Bernoulli distributed random variable. ε can take any value in $\{0, 1\}$ since both the source-relay and relay-destination links are in outage when $\beta_0(i) = 0$ and $\beta_1(i) = 0$ so that both of them will remain silent regardless of $d_i = 0$ or 1. For simplicity of notations, let $\Pr\{\mathcal{C} = 1\} = P_C$ and $\Pr\{\mathcal{C} = 0\} = 1 - P_C$. Obviously, the delay can be limited by adjusting the value of P_C . The minimum delay can be achieved by setting $P_C = 1$ since the buffer has the maximum departure rate and the minimum arrival rate. Similarly, the maximum delay appears at $P_C = 0$.

Policy II: Finite buffer size. In this policy, if $Q(i) = 0$, then $d_i = 0$; if $Q(i) = L$ (i.e., the buffer is full) where L denotes the buffer size, then $d_i = 1$; otherwise, d_i is given by (3.8). In this case, delay is controlled by limiting the buffer size L . P_C can be fixed at a feasible value. One possible selection of P_C is to balance the arrival rate and the departure rate. As proved in our previous work [102], such balancing can effectively avoid potential buffer overflow and underflow so as to improve system throughput. Let $E[(1 - d_i)\beta_0(i)] = E[d_i\beta_1(i)]$, where the left (right) side of the

equality stands for the arrival (departure) rate excluding the cases of empty and full buffers. We have

$$P_C = \begin{cases} \frac{1-2P_{e,0}+P_{e,0}P_{e,1}}{2(1-P_{e,0})(1-P_{e,1})}, & P_{e,0} \leq \frac{1}{2-P_{e,1}} \text{ and } P_{e,1} \leq \frac{1}{2-P_{e,0}} \text{ (Case I)}, \\ 0, & P_{e,0} > \frac{1}{2-P_{e,1}} \text{ (Case II)}, \\ 1, & P_{e,1} > \frac{1}{2-P_{e,0}} \text{ (Case III)}. \end{cases} \quad (3.9)$$

Note that Cases I, II and III in (3.9) are mutually exclusive, i.e., for any combination of $P_{e,0}$ and $P_{e,1}$, only one of these cases applies (the proof is similar to Appendix B in [90]). For Case I, the arrival and departure rates are exactly equal while for Cases II and III, P_C is chosen such that two rates are balanced as much as possible.

Overheads: In order to perform link scheduling in a given slot i , a central node selected from the source, the relay or the destination requires knowledge of outage states of both links (i.e., $\beta_0(i)$ and $\beta_1(i)$) as well as the buffer underflow/overflow state (i.e., if the buffer is empty/full). The outage state of the source-relay link, $\beta_0(i)$, can be determined based on γ_0^{th} and $\gamma_0(i)$ at the relay, and may be fed back to the central node using one bit of feedback overhead. Similarly, $\beta_1(i)$ can be determined based on γ_1^{th} and $\gamma_1(i)$ at the destination. The instantaneous SNR $\gamma_j(i)$, $j = 0, 1$, is calculated by $\gamma_j(i) = \frac{P_j}{\sigma_j^2} |h_j(i)|^2$ where the instantaneous channel state information (CSI) $h_j(i)$ is estimated at the receiver over a link by using pilot symbols. After executing link scheduling, the central node needs to broadcast the result to the source and the relay. Note that, only channel outage state information (COSI) $\beta_j(i)$, rather than complete CSI $h_j(i)$, is required at the central node. The transmission rate is fixed. This is because adaptively variable rate transmissions may require feeding back complete CSI (instead of outage states only) to the central node for link scheduling as well

as to the transmitter associated with the selected link for rate adaptation. These processes induce much more overheads.

In practice, the central node selection usually depends on applications. For example, in an ad-hoc relay network, the relay serving as the central node may have more advantages since it requires less feedback overheads (the feedback of the relay's buffer state information is not required) while in the downlink of a cellular network, it may be better to select the source (i.e., the base station) as the central node since it can offer more processing power and thus afford the computational complexity for link scheduling and other resource allocation algorithms [89].

3.1.3 Queue at The Relay

The queue length $Q(i)$ may change from one slot to another depending on the applied link scheduling policy, the outage states of both the source-relay and relay-destination links, and the state of the buffer itself. Specifically, the arrival and departure processes of the queue can be described as follows.

Arrival process: In a given slot i , if i) the source-relay link is selected (i.e., $d_i = 0$), ii) the source-relay link is not in outage (i.e., $\beta_0(i) = 1$), and iii) the relay buffer is not full, then the source transmits one packet to the queue, i.e., $Q(i+1) = Q(i) + 1$. Otherwise, there is no packet arrival.

Departure process: Similarly, in a given slot i , if i) the relay-destination link is selected (i.e., $d_i = 1$), ii) the relay-destination link is not in outage (i.e., $\beta_1(i) = 1$), and iii) the relay buffer is nonempty, then one packet departs from the queue, i.e., $Q(i+1) = Q(i) - 1$. Otherwise, there is no packet departure.

Note that, the case that there is one packet arrival and one packet departure at the same time is not feasible since only one of the source-relay and relay-destination links can be selected for transmissions in a given slot subject to the half-duplex constraint of the relay.

3.2 Queueing Behavior under Policy I: Infinite Buffer Size

In this section, we will examine the queueing behavior of packets in the relay buffer when applying Policy I with an infinite buffer size. Under time-correlated fading channels, the queueing behaviour is closely related to the outage processes of both the source-relay and relay-destination links (i.e., $\beta_0(i)$ and $\beta_1(i)$), i.e, the evolution of buffer state $Q(i)$ closely depends on the changing processes of link outage states. Unlike i.i.d. fading where the state transition probability of $Q(i)$ is independent from time slot index i (i.e., time-invariant) and thus the stationary distribution of $Q(i)$ can be easily analyzed [90], for time-correlated fading, such transition probability becomes time-variant. For example, according to the link scheduling policy (3.8) and the arrival/departure process of the queue, we have the state transition probability $\Pr\{Q(i+1) = q+1|Q(i) = q\} = \Pr\{\beta_0(i) = 1, \beta_1(i) = 0\} + (1 - P_C)\Pr\{\beta_0(i) = 1, \beta_1(i) = 1\}$ for $q > 0$. For i.i.d. fading, both $\Pr\{\beta_0(i) = 1, \beta_1(i) = 0\}$ and $\Pr\{\beta_0(i) = 1, \beta_1(i) = 1\}$ are not associated with time index i which makes $\Pr\{Q(i+1) = q+1|Q(i) = q\}$ independent of time index i . For correlated fading, however, both $\Pr\{\beta_0(i) = 1, \beta_1(i) = 0\}$ and $\Pr\{\beta_0(i) = 1, \beta_1(i) = 1\}$

depend on $\beta_0(i-1)$ and $\beta_1(i-1)$ (i.e., associated with time index i). In order to overcome this issue, we introduce an aggregate process $Y(i) = (\beta_0(i), \beta_1(i), Q(i))$ by combining the channel fading process and the buffer occupancy process. The state transition probability of $Y(i)$ can be represented as

$$\begin{aligned} & \Pr\{Y(i+1) = (a', b', q') | Y(i) = (a, b, q)\} \\ &= \Pr\{\beta_0(i+1) = a' | \beta_0(i) = a\} \times \Pr\{\beta_1(i+1) = b' | \beta_1(i) = b\} \\ & \times \Pr\{Q(i+1) = q' | \beta_0(i) = a, \beta_1(i) = b, Q(i) = q\}, \end{aligned} \quad (3.10)$$

where $a, a', b, b' \in \{0, 1\}$ and $q, q' \in \{0, 1, 2, \dots\}$. According to (3.2) and (3.8), it is obvious that all terms in the right hand side of (3.10) are not associated with time index i . Thus, the state transition probabilities of $Y(i)$ are time invariant. Next, we will first analyze the stationary distribution of $Y(i)$ and then extract the stationary distribution of $Q(i)$ from $Y(i)$.

$Y(i)$ is a discrete-time Markov chain (DTMC). For a compact representation, we introduce a new variable $\beta(i)$ as

$$\beta(i) = \begin{cases} 0, & \beta_0(i) = 0 \text{ and } \beta_1(i) = 0, \\ 1, & \beta_0(i) = 0 \text{ and } \beta_1(i) = 1, \\ 2, & \beta_0(i) = 1 \text{ and } \beta_1(i) = 0, \\ 3, & \beta_0(i) = 1 \text{ and } \beta_1(i) = 1. \end{cases} \quad (3.11)$$

Then, the aggregate Markov chain $Y(i)$ can be rewritten as $Y(i) = (\beta(i), Q(i))$ and its state space can be represented as $\{(0, 0), (1, 0), (2, 0), (3, 0), (0, 1), (1, 1), (2, 1), (3, 1), \dots\}$ where the states $\{(0, q), (1, q), (2, q), (3, q)\}$ constitute the q -th level of the Markov chain $Y(i)$, $q \in \{0, 1, \dots\}$.

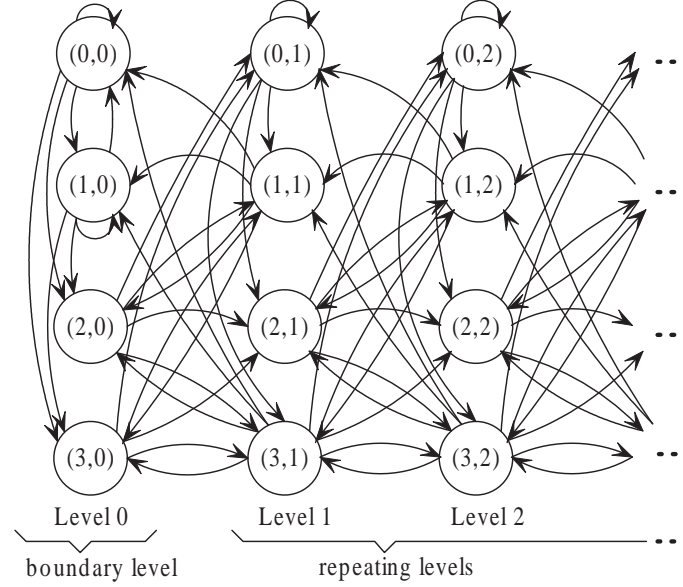


Figure 3.2: State transition diagram of the aggregate Markov chain $Y(i) = (\beta(i), Q(i))$ under Policy I.

The transition probability matrix of the aggregate Markov chain $Y(i) = (\beta(i), Q(i))$ under Policy I can be represented as a semi-infinite matrix

$$\mathbf{T} = \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_0 & & & \\ \mathbf{A}_2 & \mathbf{A}_1 & \mathbf{A}_0 & & \\ & \mathbf{A}_2 & \mathbf{A}_1 & \mathbf{A}_0 & \\ & & \ddots & \ddots & \ddots \end{bmatrix}, \quad (3.12)$$

where \mathbf{B}_1 stands for the transition sub-matrix within level 0 and \mathbf{B}_0 represents the transition sub-matrix from level 0 to level 1. Similarly, \mathbf{A}_2 , \mathbf{A}_1 , and \mathbf{A}_0 stand for the transition sub-matrices from level q to $q - 1$, within q , and from level q to $q + 1$ for $q > 0$, respectively. All of these sub-matrices \mathbf{B}_1 , \mathbf{B}_0 , \mathbf{A}_2 , \mathbf{A}_1 , \mathbf{A}_0 are 4×4 real matrices, which are derived in Appendix A.1. The possible state transitions of the aggregate chain $Y(i)$ under Policy I are illustrated in Fig. 3.2. From this figure, we

can see that the chain $Y(i)$ is a discrete-time quasi-birth-death (QBD) chain where the state transitions occur only at the same level or between adjacent levels.

Stability condition: For stability, the drift of the QBD chain $Y(i)$ to higher levels must be smaller than the drift to lower levels so that $Y(i)$ is recurrent. Specifically, the stability condition is given as [91] [92]

$$\pi_{\mathbf{A}} \mathbf{A}_0 \mathbf{1} < \pi_{\mathbf{A}} \mathbf{A}_2 \mathbf{1}, \quad (3.13)$$

where $\pi_{\mathbf{A}}$ is the stationary probability vector of the generator matrix \mathbf{A} ($\mathbf{A} = \mathbf{A}_0 + \mathbf{A}_1 + \mathbf{A}_2$). $\pi_{\mathbf{A}}$ is a 1×4 row vector and $\mathbf{1}$ is a 4×1 column vector with all elements equal to 1. In fact, \mathbf{A} is the state transition matrix of the channel outage process $\beta(i)$ and $\pi_{\mathbf{A}}$ is the stationary distribution of $\beta(i)$. Therefore, the stability condition (3.13) is equivalent to $E[(1 - d_i)\beta_0(i)] < E[d_i\beta_1(i)]$ (d_i is defined in (3.8)). This means that the stability condition requires the arrival rate should be less than the departure rate excluding the case of the buffer being empty.

Stationary distribution: $Y(i)$ is a typical QBD chain and thus its stationary probability distribution can be solved by the standard Matrix-geometric method [91] [92]. Specifically, we denote the stationary distribution of $Y(i)$ as a semi-infinite vector $\pi = \{\pi_0, \pi_1, \dots\}$ where $\pi_q = \{\pi(4q + 0), \pi(4q + 1), \pi(4q + 2), \pi(4q + 3)\}$ represents the stationary probability vector of the q -th level of $Y(i)$ and $\pi(4q + a)$ stands for the steady probability that the Markov chain $Y(i) = (\beta(i), Q(i))$ is in the state of $\beta(i) = a$ and $Q(i) = q$ ($a \in \{0, 1, 2, 3\}$, $q \in \{0, 1, \dots\}$). The stationary probability vectors of levels 0 and 1, π_0 and π_1 respectively, can be obtained by

solving the following equations

$$\begin{cases} (\boldsymbol{\pi}_0, \boldsymbol{\pi}_1) = (\boldsymbol{\pi}_0, \boldsymbol{\pi}_1) \mathbf{B}[\mathbf{R}], & (3.14a) \\ \boldsymbol{\pi}_0 \mathbf{1} + \boldsymbol{\pi}_1 (\mathbf{I} - \mathbf{R})^{-1} \mathbf{1} = 1, & (3.14b) \end{cases} \quad (3.14)$$

where

$$\mathbf{B}[\mathbf{R}] = \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_0 \\ \mathbf{A}_2 & \mathbf{A}_1 + \mathbf{R} \mathbf{A}_2 \end{bmatrix}. \quad (3.15)$$

Here, \mathbf{R} is the minimal non-negative solution to the following equation

$$\mathbf{R} = \mathbf{A}_0 + \mathbf{R} \mathbf{A}_1 + \mathbf{R}^2 \mathbf{A}_2. \quad (3.16)$$

Equation (3.16) can be solved by an iterative method where the iterative formula is $\mathbf{R}(k+1) = \mathbf{A}_0 + \mathbf{R}(k) \mathbf{A}_1 + \mathbf{R}^2(k) \mathbf{A}_2$ or $\mathbf{R}(k+1) = (\mathbf{A}_0 + \mathbf{R}^2(k) \mathbf{A}_2)(\mathbf{I} - \mathbf{A}_1)^{-1}$, the initial value can be set as $\mathbf{R}(0) = \mathbf{0}$, and the iteration is repeated until $\|\mathbf{R}(n+1) - \mathbf{R}(n)\|_\infty < \epsilon$ ($\|\cdot\|_\infty$ denotes the infinity norm of a matrix and ϵ is the predefined tolerance error). Note that the matrix equation (3.16) can also be solved by other advanced methods, such as the logarithmic reduction (LR) method [103] and the cyclic reduction (CR) method [104] [105]. The stationary probability vectors of higher levels (≥ 2) of $Y(i)$ can be derived recursively as

$$\boldsymbol{\pi}_q = \boldsymbol{\pi}_{q-1} \mathbf{R} = \boldsymbol{\pi}_1 \mathbf{R}^{q-1}, \quad q \geq 2. \quad (3.17)$$

Remark: If the appropriate inverse matrices exist, we can write $\boldsymbol{\pi}_0$ and $\boldsymbol{\pi}_1$ in a closed form based on (3.14), as

$$(\boldsymbol{\pi}_0, \boldsymbol{\pi}_1) = \mathbf{1}^T (\mathbf{I} - \mathbf{B}[\mathbf{R}] + \mathbf{U})^{-1}, \quad (3.18)$$

or

$$\begin{cases} \boldsymbol{\pi}_1 = \boldsymbol{\beta}(\mathbf{A}_2(\mathbf{I} - \mathbf{B}_1)^{-1} + (\mathbf{I} - \mathbf{R})^{-1})^{-1}, & (3.19a) \\ \boldsymbol{\pi}_0 = \boldsymbol{\pi}_1 \mathbf{A}_2(\mathbf{I} - \mathbf{B}_1)^{-1}. & (3.19b) \end{cases} \quad (3.19)$$

In (3.18), $\mathbf{U} = [\mathbf{u}, \mathbf{u}, \dots, \mathbf{u}]$ is a 8×8 matrix with each column equal to \mathbf{u} ,

$$\mathbf{u} = \begin{bmatrix} \mathbf{1} \\ (\mathbf{I} - \mathbf{R})^{-1} \mathbf{1} \end{bmatrix}. \quad (3.20)$$

In (3.19), $\boldsymbol{\beta}$ is the stationary probability vector of channel outage states, i.e.,

$$\begin{aligned} \boldsymbol{\beta} &= [\Pr\{\beta(i) = 0\}, \Pr\{\beta(i) = 1\}, \Pr\{\beta(i) = 2\}, \Pr\{\beta(i) = 3\}] \\ &= [P_{e,0}P_{e,1}, P_{e,0}(1 - P_{e,1}), (1 - P_{e,0})P_{e,1}, (1 - P_{e,0})(1 - P_{e,1})]. \end{aligned} \quad (3.21)$$

Otherwise, if the associated inverse matrices do not exist, $\boldsymbol{\pi}_0$ and $\boldsymbol{\pi}_1$ can be obtained by using Gaussian elimination method [92] to solve the equation set (3.14).

The proofs of (3.18) and (3.19) are shown in Appendix A.3.

Average end-to-end delay: The end-to-end delay of a packet transmitted from the source to the destination includes three parts ¹: i) transmission delay from the source to the relay, denoted as τ_1 ; ii) waiting delay in the relay queue, denoted as τ_2 ; and iii) transmission delay from the relay to the destination, denoted as τ_3 . Since the link outage state is known in advance by channel estimates and a packet is transmitted only if the link is not in outage, it requires only one slot to transmit a packet from one node to another, i.e., $\tau_1 = \tau_3 = 1$ slot. The sum of τ_2 and τ_3 comprises the total system delay of a packet in the relay queue (i.e., the sum of waiting and service

¹Since we assume the source always has data available for transmissions, the waiting delay at the source is not included as in [89] [90]. Thus, we mainly focus on the additional delay caused by the relay buffer. When there are constant or random data arrivals at the source, the waiting delay at the source should be considered.

times), denoted as $\tau_s = \tau_2 + \tau_3$. Therefore, the average end-to-end delay of packets, $\bar{\tau}$, can be represented as

$$\bar{\tau} = \bar{\tau}_1 + \bar{\tau}_s, \quad (3.22)$$

where $\bar{\tau}_1 = 1$ and the average system delay $\bar{\tau}_s$, according to Little's law, is given as

$$\bar{\tau}_s = E(Q(i))/\bar{r}_A, \quad (3.23)$$

where $E(Q(i))$ is the average queue length and \bar{r}_A is the average arrival rate of the queue. The average queue length can be represented as

$$E(Q(i)) = \sum_{q=0}^{\infty} q \cdot \Pr\{Q(i) = q\}, \quad (3.24)$$

$$\Pr\{Q(i) = q\} = \boldsymbol{\pi}_q \mathbf{1}. \quad (3.25)$$

According to Policy I, the average arrival rate of the queue, \bar{r}_A (packets/slot), can be calculated by

$$\bar{r}_A = \boldsymbol{\pi}_0 \mathbf{c}_1 + \left(\sum_{q=1}^{\infty} \boldsymbol{\pi}_q \right) \mathbf{c}_2 = \boldsymbol{\pi}_0 \mathbf{c}_1 + \boldsymbol{\pi}_1 (\mathbf{I} - \mathbf{R})^{-1} \mathbf{c}_2, \quad (3.26)$$

where $\mathbf{c}_1 = [0, 0, 1, 1]^T$ and $\mathbf{c}_2 = [0, 0, 1, 1 - P_C]^T$.

Average system throughput: The throughput of the considered system is defined as the number of packets arriving at the destination from the source per slot, in a unit of packets/slot. Thus, the average throughput, \bar{r} , is equal to the average departure rate of the queue, \bar{r}_D , as

$$\bar{r} = \bar{r}_D = \left(\sum_{q=1}^{\infty} \boldsymbol{\pi}_q \right) \mathbf{c}_3 = \boldsymbol{\pi}_1 (\mathbf{I} - \mathbf{R})^{-1} \mathbf{c}_3, \quad (3.27)$$

where $\mathbf{c}_3 = [0, 1, 0, P_C]^T$.

System outage probability: The outage probability of a system is defined as the frequency that the link quality cannot afford the successful transmission of a packet. Note that in our system, the link outage does not lead to packet loss since the outage states of both the source-relay and relay-destination links are known in advance by channel estimates and the system remains silent if both links are in outage. However, the reduction in system throughput is inevitable due to the presence of link outages. Thus, the system outage probability can also be interpreted as the fraction of throughput loss caused by link outages as [90]

$$P_{out} = 1 - \frac{\bar{r}}{r_{max}}, \quad (3.28)$$

where r_{max} represents the maximum throughput in the absence of link outages. In our system, $r_{max} = 0.5$ since a packet requires two slots to transmit from the source to the destination when both links are not in outage.

3.3 Queueing Behavior under Policy II: Finite Buffer Size

In this section, we will study Policy II with a finite buffer size and focus on the queueing behaviour of packets in the relay buffer. Similar to the analyses of Policy I, we investigate the aggregate Markov chain $Y(i) = (\beta(i), Q(i))$.

Under Policy II, the state transitions of the chain $Y(i)$ are illustrated in Fig. 3.3

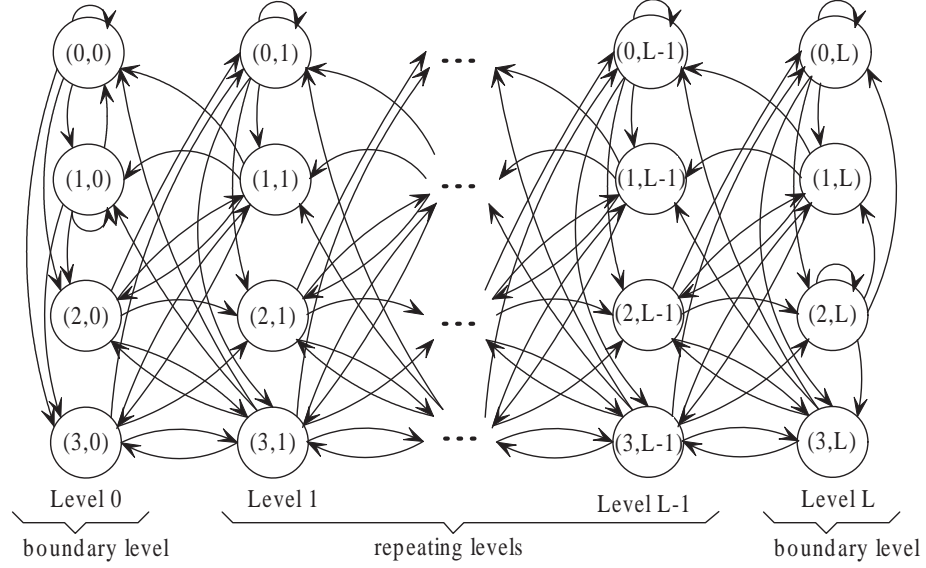


Figure 3.3: State transition diagram of the aggregate Markov chain $Y(i) = (\beta(i), Q(i))$ under Policy II.

and the transition probability matrix can be described as

$$\mathbf{T} = \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_0 & & & \\ \mathbf{A}_2 & \mathbf{A}_1 & \mathbf{A}_0 & & \\ & \mathbf{A}_2 & \mathbf{A}_1 & \mathbf{A}_0 & \\ & & \ddots & \ddots & \ddots \\ & & & \mathbf{A}_2 & \mathbf{A}_1 & \mathbf{A}_0 \\ & & & & \mathbf{C}_2 & \mathbf{C}_1 \end{bmatrix}, \quad (3.29)$$

where \mathbf{B}_1 , \mathbf{B}_0 , \mathbf{A}_2 , \mathbf{A}_1 , \mathbf{A}_0 are given by (A.1). \mathbf{C}_1 and \mathbf{C}_2 stand for the probabilities of transitions within level L and from level L to level $L - 1$, respectively, which are derived in Appendix A.2.

Stationary distribution: From Fig. 3.3 and the transition probability matrix (3.29), we can observe that the aggregate chain $Y(i)$ under Policy II is also a QBD

chain. However, different from Policy I with an infinite buffer, since Policy II is limited by the buffer size, $Y(i)$ under Policy II has two boundary levels (level 0 and level L). In this case, the stationary distribution of $Y(i)$ cannot be directly calculated by the standard matrix-geometric method. To address the problem of two boundary levels, we modify the matrix-geometric method in [93]. Note that although there are already some solutions on a general finite QBD chain [93], [106]–[109], applying existing solutions for the specific problem as formulated in this chapter needs some necessary modifications in boundaries. Our solution is based on the framework of the matrix-geometric method as proposed in [93] so that it has a logarithmic time complexity on the number of levels ($\mathcal{O}(\log_2 L)$). It is better than other solutions with a linear time complexity of $\mathcal{O}(L)$ in [106]–[109], especially in the case under our consideration where there are a small number of phases (i.e., the order of transition sub-matrices \mathbf{A}_0 , \mathbf{A}_1 , and \mathbf{A}_2), but a large number of levels. The modified matrix-geometric method is described in the following theorem.

Theorem 3.1. *Define rate matrices \mathbf{R}_1 and \mathbf{R}_2 as the minimal nonnegative solutions to the following matrix equations*

$$\mathbf{R}_1 = \mathbf{A}_0 + \mathbf{R}_1 \mathbf{A}_1 + \mathbf{R}_1^2 \mathbf{A}_2, \quad (3.30)$$

$$\mathbf{R}_2 = \mathbf{A}_2 + \mathbf{R}_2 \mathbf{A}_1 + \mathbf{R}_2^2 \mathbf{A}_0, \quad (3.31)$$

respectively. Then, the stationary probability vector of level q ($1 \leq q \leq L - 1$) can be given by the following recursive structure

$$\boldsymbol{\pi}_q = \mathbf{v}_1 \mathbf{R}_1^{q-1} + \mathbf{v}_2 \mathbf{R}_2^{L-q-1}, \quad (3.32)$$

where \mathbf{v}_1 and \mathbf{v}_2 are two constant vectors. \mathbf{v}_1 , \mathbf{v}_2 , $\boldsymbol{\pi}_0$ and $\boldsymbol{\pi}_L$ together can be obtained by solving the following equations

$$\begin{cases} [\boldsymbol{\pi}_0, \mathbf{v}_1, \mathbf{v}_2, \boldsymbol{\pi}_L] \mathbf{B}[\mathbf{R}_1, \mathbf{R}_2] = [\boldsymbol{\pi}_0, \mathbf{v}_1, \mathbf{v}_2, \boldsymbol{\pi}_L], & (3.33a) \\ (\boldsymbol{\pi}_0 + \mathbf{v}_1 \mathbf{S}_1 + \mathbf{v}_2 \mathbf{S}_2 + \boldsymbol{\pi}_L) \mathbf{1} = 1, & (3.33b) \end{cases} \quad (3.33)$$

where

$$\mathbf{S}_1 = \sum_{q=0}^{L-2} \mathbf{R}_1^q, \quad \mathbf{S}_2 = \sum_{q=0}^{L-2} \mathbf{R}_2^q, \quad (3.34)$$

$$\mathbf{B}[\mathbf{R}_1, \mathbf{R}_2] = \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_0 & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_2 & \mathbf{J}_1 & \mathbf{J}_3 & \mathbf{R}_1^{L-2} \mathbf{A}_0 \\ \mathbf{R}_2^{L-2} \mathbf{A}_2 & \mathbf{J}_2 & \mathbf{J}_4 & \mathbf{A}_0 \\ \mathbf{0} & \mathbf{0} & \mathbf{C}_2 & \mathbf{C}_1 \end{bmatrix}. \quad (3.35)$$

In (3.35), $\mathbf{J}_1 = \mathbf{A}_1 + \mathbf{R}_1 \mathbf{A}_2$, $\mathbf{J}_2 = \mathbf{R}_2^{L-3}(\mathbf{R}_2 \mathbf{A}_1 + \mathbf{A}_2 - \mathbf{R}_2)$, $\mathbf{J}_3 = \mathbf{R}_1^{L-3}(\mathbf{A}_0 + \mathbf{R}_1 \mathbf{A}_1 - \mathbf{R}_1)$, and $\mathbf{J}_4 = \mathbf{R}_2 \mathbf{A}_0 + \mathbf{A}_1$. Note that $L \geq 3$.

Proof. Please see Appendix A.4. □

Remark: The condition of $L \geq 3$ should be satisfied to apply the modified matrix-geometric method. Otherwise, if $L < 3$, the stationary distribution of $Y(i)$ can be obtained by directly solving the equations

$$\begin{cases} \boldsymbol{\pi} \mathbf{T} = \boldsymbol{\pi}, & (3.36a) \\ \boldsymbol{\pi} \mathbf{1} = 1, & (3.36b) \end{cases} \quad (3.36)$$

where (3.36a) is the balance equation associated with the chain $Y(i)$ and (3.36b) is the normalized condition. Note that, in fact, one can always solve (3.36) for the

stationary distribution of $Y(i)$ regardless of how large L is. However, when L is large, directly solving (3.36) will cause a much higher computational complexity. The modified matrix-geometric method utilizes the repetitive structure in the transition matrix of $Y(i)$ and has a relatively low complexity when L is large.

Average end-to-end delay: As in Policy I, the average end-to-end delay of packets in Policy II can also be expressed as $\bar{\tau} = \bar{\tau}_1 + \bar{\tau}_s$ where $\bar{\tau}_1 = 1$ and $\bar{\tau}_s = E(Q(i))/\bar{r}_A$. $E(Q(i))$ can be calculated as $E(Q(i)) = \sum_{q=0}^L q \cdot \Pr\{Q(i) = q\}$ where $\Pr\{Q(i) = q\} = \pi_q \mathbf{1}$, and the average arrival rate is $\bar{r}_A = \pi_0 \mathbf{c}_1 + (\sum_{q=1}^{L-1} \pi_q) \mathbf{c}_2 = \pi_0 \mathbf{c}_1 + (\mathbf{v}_1 \mathbf{S}_1 + \mathbf{v}_2 \mathbf{S}_2) \mathbf{c}_2$.

Average system throughput: The average system throughput, \bar{r} , equals the average departure rate of the queue, \bar{r}_D , given as $\bar{r} = \bar{r}_D = (\sum_{q=1}^{L-1} \pi_q) \mathbf{c}_3 + \pi_L \mathbf{c}_4 = (\mathbf{v}_1 \mathbf{S}_1 + \mathbf{v}_2 \mathbf{S}_2) \mathbf{c}_3 + \pi_L \mathbf{c}_4$ where $\mathbf{c}_4 = [0, 1, 0, 1]^T$.

System outage probability: As in Policy I, the system outage probability equals $P_{out} = 1 - \bar{r}/r_{max}$.

3.4 Numerical and Simulation Results

In this section, the performance of buffer-aided relaying on time-correlated fading channels will be evaluated. In the following evaluations, if not particularly indicated, we consider a case of correlated fading channels with a normalized Doppler bandwidth of 0.02 (i.e., $f_{D,0}T = f_{D,1}T = 0.02$) which corresponds to a pedestrian Doppler frequency of 1.34 Hz and a slot length of 16 ms [31] [32]. For comparison, the case of i.i.d. fading with the same packet error rate will also be considered. Note that all of the following analytical results have been verified by 10^6 Monte Carlo simulations.

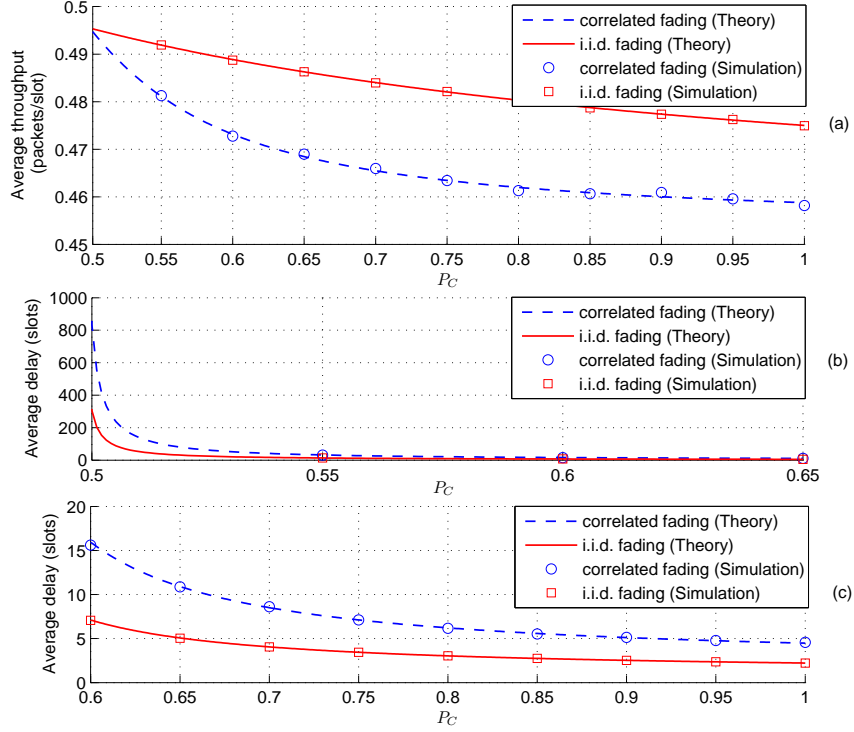
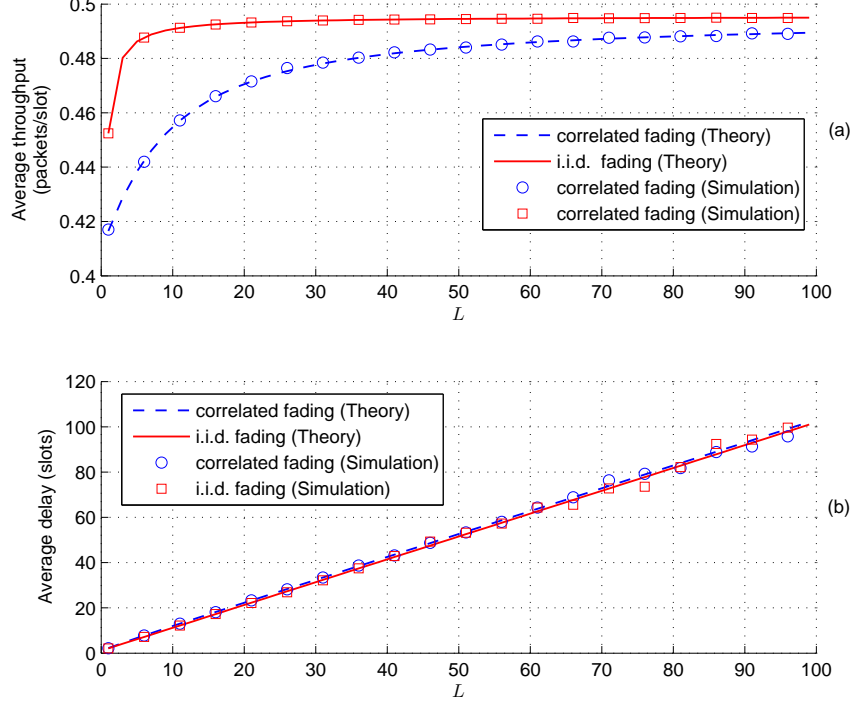


Figure 3.4: Effects of P_C in Policy I. Fading margins $F_0 = F_1 = 10$ dB.

Some simulation results may not be plotted in the figures for clarity.

1) *Effects of P_C in Policy I.* Fig. 3.4 illustrates the effects of P_C on average throughput and end-to-end delay in Policy I with an infinite buffer size, where the sub-figure (c) is a zoomed version of the sub-figure (b) on $P_C \in [0.6, 1]$. Note that the throughput and delay curves corresponding to $P_C < 0.5$ are not plotted since the queue is unstable in our simulation settings. $P_C > 0.5$ corresponds to the stable region and $P_C = 0.5$ is the critical point at which the queue is on the edge of unstable (absorbing) and stable (non-absorbing) states. From this figure, we can see that the analytical results match with the simulation ones very well, which verifies the correctness of our analysis framework. Delay is decreasing with P_C in the stable


 Figure 3.5: Effects of buffer size L in Policy II. $F_0 = F_1 = 10$ dB.

region for both cases of correlated and i.i.d. fading, and thus delay can be limited by adjusting P_C as expected. The minimal delay is achieved at $P_C = 1$ and the infinite delay but the maximal throughput are observed at the critical point ($P_C = 0.5$). Moreover, at the critical point, the throughput under correlated fading is the same as that under i.i.d. fading. This means that for transmissions without delay constraints, the queue may operate at the critical point in order to maximize throughput and correlated fading does not affect the throughput performance of buffer-aided relaying. In addition, for the same value of P_C and $F_0 = F_1 = 10$ dB, correlated fading causes a small loss of throughput (about 4% at $P_C = 1$) and a larger delay compared with i.i.d. fading.

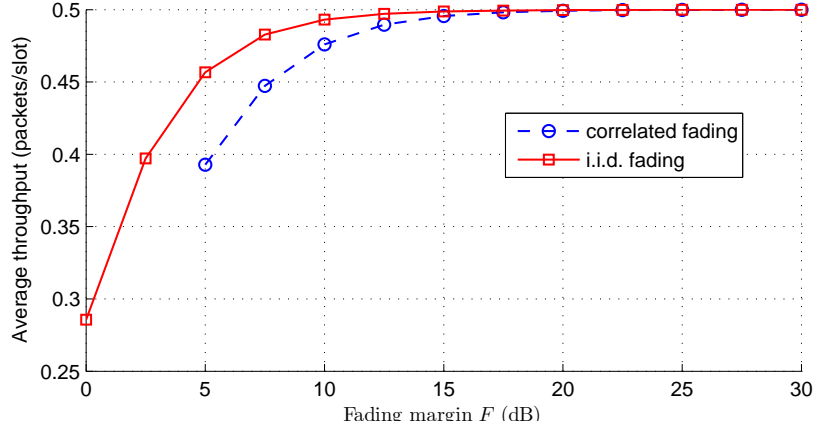


Figure 3.6: Average throughput versus fading margin $F = F_0 = F_1$ under Policy I.

2) *Effects of buffer size L in Policy II.* The average throughput and the average end-to-end delay versus the buffer size L for Policy II are shown in Fig. 3.5. From this figure, we can also observe that the analytical results are almost consistent with simulations. Since both throughput and delay are increasing with L , there exists a tradeoff between increasing throughput and reducing delay by adjusting the buffer size L . In addition, for a small L and $F_0 = F_1 = 10$ dB, correlated fading leads to a large throughput degradation (about 10% at $L = 5$) compared with i.i.d. fading, but the delays on correlated and i.i.d. fading are almost the same.

3) *Effects of fading margins and stringent delay constraints:* The average throughput versus the fading margin are shown in Fig. 3.6 for Policy I with an average delay of $\bar{\tau} = 20$ slots. The curves are obtained by adjusting P_C to satisfy the delay constraint. The curve on correlated fading does not extend to low fading margins ($F < 5$ dB) because low fading margins cannot meet the delay requirement of $\bar{\tau} = 20$ slots even $P_C = 1$ in our settings. From the figure, we can observe that, for the stringent delay requirement of $\bar{\tau} = 20$ slots, correlated fading leads to a non-negligible loss on

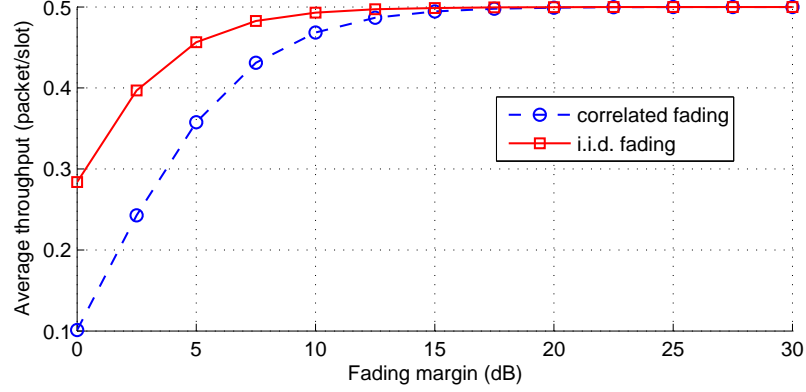


Figure 3.7: Average throughput versus fading margin $F = F_0 = F_1$ under Policy II.

throughput in low fading margins compared with i.i.d. fading. Specifically, such a throughput loss reaches 16% when the fading margin is $F = 5$ dB. The throughput under correlated fading can only approach that under i.i.d. fading in high fading margins. Unfortunately, in practice, high fading margins may not be feasible since a higher fading margin $F_j = \frac{P_j}{\sigma_j^2} E(|h_j(i)|^2) \frac{1}{\gamma_j^{th}}$ means higher power consumption P_j and/or a lower SNR threshold γ_j^{th} (lower rate).

We repeat a similar simulation for Policy II, as shown in Fig. 3.7. Similar observations can be obtained. In particular, a throughput loss of about 28% can be observed when $F = 5$ dB.

4) *Comparisons between Policies I and II.* In Fig. 3.8, we illustrate the average throughput versus the average end-to-end delay for both Policies I and II. From the figure, we can observe that Policy I has higher throughput than Policy II under the same average delay for both correlated and i.i.d. fading. It is because Policy II with a finite buffer size may experience buffer overflow which limits the flexibility in link scheduling. The throughput is increasing with the affordable delay for both

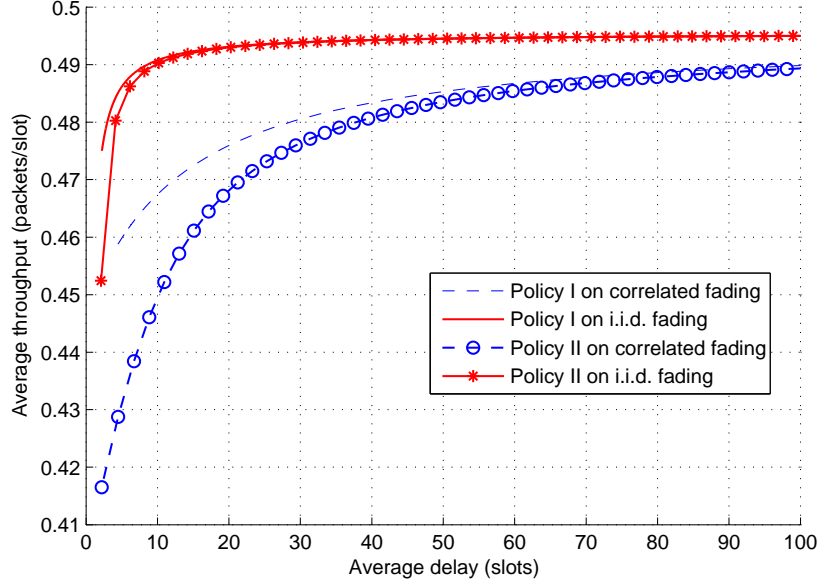


Figure 3.8: Average throughput versus average end-to-end delay for both Policies I and II. $F_0 = F_1 = 10$ dB.

policies. When the affordable delay is sufficiently large, Policy II has almost the same throughput as Policy I since the buffer size L in Policy II is allowed to increase as the affordable delay becomes large, and Policies I and II are equivalent when P_C is at the critical point in Policy I and $L \rightarrow \infty$ in Policy II.

3.5 Guidelines on Performance Improvement under Correlated Fading

In this section, we will provide more discussions on the effects of correlated channels on the performance of buffer-aided relaying, by briefly showing why the time correlation in channel fading may degrade the performance tradeoff between throughput

and delay. Then, we will present some guidelines on how to improve the performance under correlated channels, as well as how to design an effective policy without channel outage state information (COSI). Some simple examples of improved policies are provided.

3.5.1 Performance Degradation due to Channel Correlations

For correlated fading, the transition probability of a link within the same state (i.e., from “good” to “good”, $\beta_j(i-1) = 1 \rightarrow \beta_j(i) = 1$, or from “bad” to “bad”, $\beta_j(i-1) = 0 \rightarrow \beta_j(i) = 0$) is relatively high, while the transition probability between different states is relatively low. Such characteristics tend to make each link keep staying on the same state for a long time (i.e., multiple successive slots). Let’s consider two scenarios. In Scenario I where $\beta_0(i) = 1$ and $\beta_1(i) = 0$ occur in several successive slots, only the source-relay link can be selected for transmissions and there will be multiple packet arrivals in succession. On the contrary, in Scenario II where $\beta_0(i) = 0$ and $\beta_1(i) = 1$ happen in several successive slots, the system can only provide successive services. The probability of Scenario I (Scenario II) lasting k successive slots given the scenario has already occurred can be calculated as $P_k^I = (p_0q_1)^{k-1}(1-p_0q_1)$ ($P_k^{II} = (q_0p_1)^{k-1}(1-q_0p_1)$). For a large k , such probabilities of successive occurrences under correlated fading become much higher than those under i.i.d. fading. For example, under the correlated fading with $f_{D,0}T = f_{D,1}T = 0.02$ and $F_0 = F_1 = 10$ dB, we have $p_0 = p_1 = 0.98$, $q_0 = q_1 = 0.85$, $P_{e,0} = P_{e,1} = 0.095$, and $P_k^I = P_k^{II} = (0.83)^{k-1} \times 0.17$. While, under the i.i.d. fading with the same packet error rates, we have $P_k^I = P_k^{II} = [(1-P_{e,0})P_{e,1}]^{k-1}[1-(1-P_{e,0})P_{e,1}] = 0.086^{k-1} \times 0.914$.

This observation indicates that buffer-aided relaying has the characteristics of batch arrivals and batch services under correlated fading, instead of scattered arrivals and scattered services under i.i.d. fading. Compared to scattered arrivals/services, batch arrivals/services lead to a larger delay, which can be explained by the following two simple examples. An extreme example with scattered arrivals/services is that there is an arrival followed by a service in every other slot so that the end-to-end delay is 2 slots. The other example with batch arrivals/services is that there are two successive arrivals followed by two successive services in every four slots, where each packet requires a one-slot wait in the relay buffer and the end-to-end delay becomes 3 slots. Moreover, batch arrivals/services are prone to buffer overflow/underflow, which reduces the flexibility in link scheduling and eventually decreases system throughput.

3.5.2 An Improved Policy

Based on the previous discussions, we can find that the characteristics of batch arrivals/services under correlated fading channels make the system become more sensitive to the buffer occupancy state. Thus, a potential improvement on link scheduling is to take the buffer state information (BSI) into consideration. In this subsection, as an example, we are going to show a simple improvement based on Policy I.

For Policy I with an infinite buffer, based on (3.8), there are two scenarios that may result in throughput loss: i) both links are in outage ($\beta_0 = 0$ and $\beta_1 = 0$); and ii) $\beta_0 = 0$ and $\beta_1 = 1$ while the buffer is empty (i.e., buffer underflow). In general, the occurring frequency of scenario i) cannot be controlled when the transmit power and rate are fixed. However, the occurring frequency of scenario ii) could be

managed by modifying the link scheduling policy. For example, in order to reduce the occurrences of scenario ii), if the stock in the buffer $Q(i)$ is small, one can increase the probability of the source-relay link being selected when both links are in “good” state (i.e., $\beta_0(i) = \beta_1(i) = 1$) and thus increase the arrival rate to avoid buffer underflow. On the other hand, in order to reduce the queueing delay, the relay-destination link should be given higher priority over the source-relay link so as to maximize the service rate when both links are in “good” state. Therefore, a tradeoff between increasing throughput and decreasing delay exists. The designed link scheduling policy should remain a certain stock in the relay buffer while giving priority to the relay-destination link if possible. Motivated by this, the improved policy redefines the probability P_C in (3.8) as a function of the buffer state $Q(i)$, i.e.,

$$P_C = \begin{cases} 0, & Q(i) = 0, \\ \frac{1}{M-Q(i)}, & 1 \leq Q(i) \leq M-1, \\ 1, & Q(i) \geq M, \end{cases} \quad (3.37)$$

where the value of M can be adjusted to balance throughput and delay.

Fig. 3.9 compares Policy I and the improved policy. The results are based on the same setting as shown in Section 3.4 and obtained by averaging over 10^6 Monte Carlo simulations. From this figure, we can see that the improved policy can always achieve a better tradeoff between throughput and delay compared with Policy I. Moreover, as the tolerable delay increases, the improved policy can narrow the performance gap between i.i.d. and correlated fading channels in a faster way than Policy I.

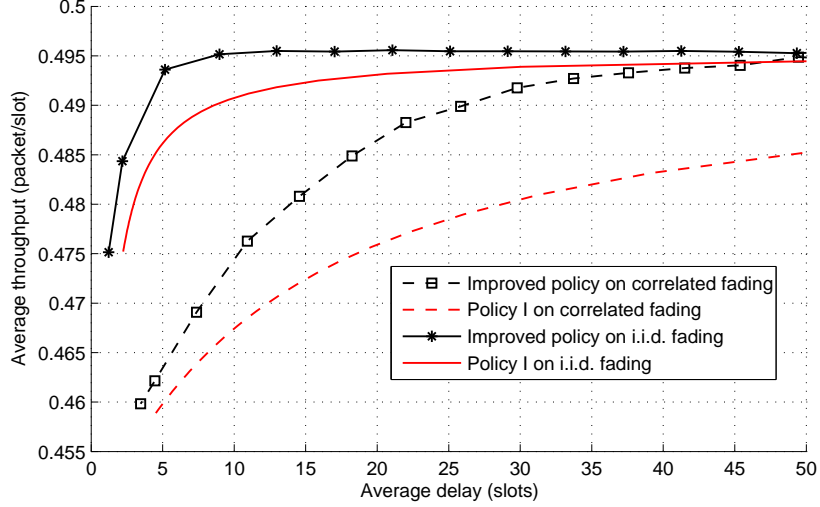


Figure 3.9: Average throughput versus average end-to-end delay under Policy I and the improved policy. $F_0 = F_1 = 10$ dB.

3.5.3 Unknown Channel Outage State Information

In practice, channel outage state information (COSI) may not be available at the central node due to the consideration of potential overheads for channel estimates and COSI exchanges. In such a case, the link scheduling policy can be designed as a random strategy. For example, for an infinite buffer, a simple random policy without COSI, called the blind-COSI policy, can be designed as

$$d_i = \begin{cases} 0, & Q(i) = 0, \\ \mathcal{C}, & Q(i) > 0. \end{cases} \quad (3.38)$$

The blind-COSI policy is a modified version of Policy I where the probability $P_C = \Pr\{\mathcal{C} = 1\}$ is independent of COSI, and tradeoffs between throughput and delay can also be achieved by adjusting the value of P_C . However, such a policy may lead to a great throughput loss and a large delay since the lack of COSI may cause

inappropriate link selections (for example, the source-relay link may also be selected when it is in outage). In order to alleviate such loss, we can use past/outdated COSI in designing link scheduling policies under correlated fading channels since the current COSI may be predicted based on the past COSI. Specifically, the central node can obtain the outdated COSI by listening to the acknowledgement (ACK)/negative-acknowledgement (NACK) signals from the receiver on the selected link at the end of the previous slot. For explanation purpose, we propose a new policy based on the COSI in the previous slot, called the outdated-COSI policy, by redefining the probability P_C in the blind-COSI policy as follows

$$P_C = \begin{cases} c(1 - p_0), & d(i - 1) = 0 \text{ and } \beta_0(i - 1) = 1, \\ q_0, & d(i - 1) = 0 \text{ and } \beta_0(i - 1) = 0, \\ p_1, & d(i - 1) = 1 \text{ and } \beta_1(i - 1) = 1, \\ 1 - q_1, & d(i - 1) = 1 \text{ and } \beta_1(i - 1) = 0, \end{cases} \quad (3.39)$$

where $1 - p_0 = \Pr\{\beta_0(i) = 0 | \beta_0(i - 1) = 1\}$, $q_0 = \Pr\{\beta_0(i) = 0 | \beta_0(i - 1) = 0\}$, $p_1 = \Pr\{\beta_1(i) = 1 | \beta_1(i - 1) = 1\}$, and $1 - q_1 = \Pr\{\beta_1(i) = 1 | \beta_1(i - 1) = 0\}$. The constant c is used to control the frequency of selecting the relay-destination link so as to balance throughput and delay. The range of c is $0 \leq c \leq 1/(1 - p_0)$.

We compare three link scheduling policies: the blind-COSI policy, the outdated-COSI policy, and Policy I with known COSI, through simulations and demonstrate the results in Fig. 3.10. Note that the same setting as in Section 3.4 is used in simulations. From this figure, we can observe that the outdated-COSI policy can significantly improve the system performance compared to the blind-COSI policy.

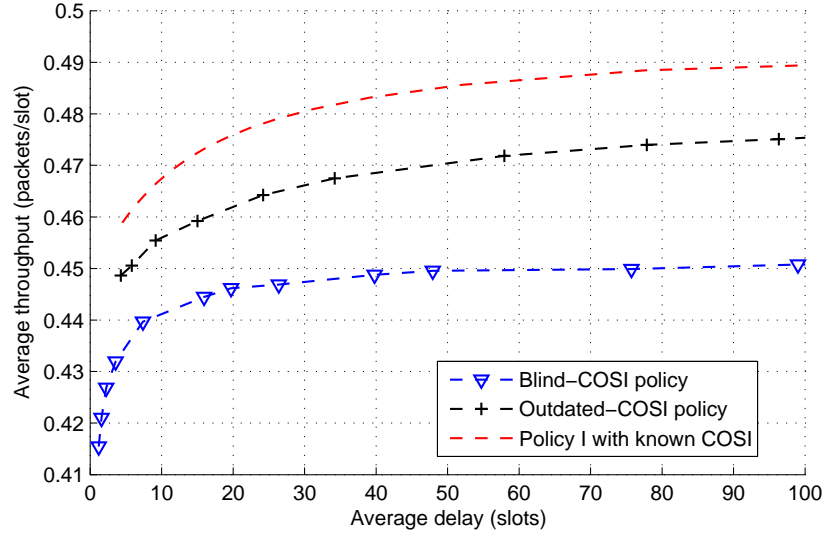


Figure 3.10: Comparisons among the blind-COSI policy, the outdated-COSI policy and Policy I with known COSI under the correlated channels. $F_0 = F_1 = 10$ dB.

Chapter 4

Interference-Avoidance Scheduling for Dense Multi-User Coexisting Networks with Heterogeneous Priorities and Demands

In this chapter, we address the interference-avoidance scheduling issue for dense multi-user coexisting networks, by considering heterogeneity in terms of user priorities and traffic demands. For the explanation purpose, wireless body area networks (WBANs) are used as an example of multi-user coexisting networks. The proposed scheme can also be applied to other multi-user coexisting networks with heterogeneous priorities and demands. Considering limited resources, the priority-aware admission control issue is first investigated, where the objective is to accommodate as many high priority users as possible. To achieve this, we determine the maximum number

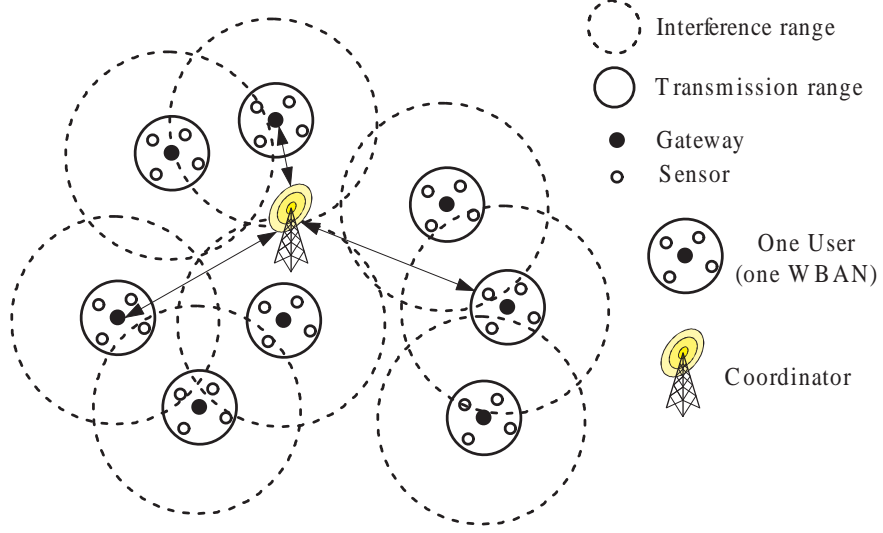


Figure 4.1: Multi-user coexisting wireless network.

of high priority levels that the system can support, where all users in these high priority levels will be admitted. We also discuss how to further admit a subset of users in low priorities. After that, the throughput of admitted users is optimized.

4.1 System Model

In this section, we present the system model under consideration and formulate the problem of priority-aware interference-avoidance scheduling for multi-user coexisting networks.

4.1.1 Network Model

Consider a system including a coordinator and N users as shown in Fig. 4.1. Each user is covered by a WBAN which further consists of a gateway and multiple

sensors. The gateway is responsible for collecting and processing health-related data from sensors while the coordinator scheduling the transmissions of N WBANs so that they do not interfere with each other. The coordinator may be an access point, a base station, or a gateway from users. Here, we consider only intra-WBAN transmissions from sensors to gateways and focus on the scheduling among multiple WBANs, while for beyond-WBAN transmissions from gateways to the coordinator, interested readers can refer to [110] [111]. Different WBANs may have heterogeneous transmission priorities based on patients' health conditions. The priority of each WBAN is selected from a set of priority levels $\mathcal{K} = \{1, \dots, K\}$ where 1 represents the highest priority level and K the lowest one. Let $\mathcal{N}_k \subset \mathcal{N}$ be the subset of users in priority level k , $k \in \mathcal{K}$, where \mathcal{N} is the set of all users. For simplicity, we consider only one channel. Users access the channel by using the spatial-reuse time division multiple access (STDMA) where users geographically separated are allowed to use the channel at the same time. We consider a time frame structure with an equal frame length T .

4.1.2 Interference Model

Define $c_{i,j}$ as a binary interference indicator: $c_{i,j} = 0$ if there is no interference between users i and j , and $c_{i,j} = 1$ otherwise. For example, the distance-based interference model [50] defines

$$c_{i,j} = \begin{cases} 1; & \text{if } d_{i,j} \leq d^{th} \\ 0; & \text{if } d_{i,j} > d^{th} \end{cases} \quad (4.1)$$

where $d_{i,j}$ is the distance between two distinct users i and j , and d^{th} is a threshold. Note that $c_{i,i} = 0$. A feasible candidate group (FCG) is defined as a group of users that

can simultaneously transmit data without interference. The set of all FCGs selected from \mathcal{N} (\mathcal{N}_k) is denoted as \mathcal{F} (\mathcal{F}_k). The m -th FCG in \mathcal{F} (\mathcal{F}_k) can be represented by a $|\mathcal{N}|$ -length ($|\mathcal{N}_k|$ -length) binary column vector \mathbf{f}_m with the i -th element as

$$f_{i,m} = \begin{cases} 1; & \text{if user } i \text{ belongs to the } m\text{-th FCG} \\ 0. & \text{otherwise.} \end{cases} \quad (4.2)$$

4.2 Problem Formulation and Solution Framework

In this section, we will present and formulate mathematically our problems of interest, and discuss the associated solution framework.

4.2.1 Admission Control

In dense user areas, admission control with priority guarantee plays an important role since normally limited resources cannot satisfy all users' transmission requirements. Thus, considering the heterogeneous priorities and traffic demands among users, our first problem of interest is to design a scheduling scheme so that the system can accommodate as many high priority users as possible, while the minimum traffic demands of admitted users are satisfied. Low priority users are only allowed to access the system if all higher priority users have been fully satisfied in terms of their minimum traffic demands. Here, the minimum traffic demand, u_i , refers to the minimum amount of data required to be transmitted in a frame by user i , which is equivalent to a minimum transmission time demand $\tau_i = u_i/r_i$ ($\tau_i < T$), where r_i is the transmission rate of user i . Denote by x_i a binary indicator: $x_i = 1$ represents user i accesses the system; otherwise, $x_i = 0$. Given that all FCGs are known, the

admission control problem can be formulated as

$$(\mathbf{P0}) : \max_{t_m \geq 0, x_i \in \{0,1\}} \sum_{i \in \mathcal{N}} x_i \quad (4.3a)$$

$$\text{s.t.} \quad \sum_{m=1}^{|\mathcal{F}|} t_m \leq T, \quad (4.3b)$$

$$\left(\sum_{m=1}^{|\mathcal{F}|} f_{i,m} t_m - \tau_i \right) x_i \geq 0, \quad \forall i \in \mathcal{N}, \quad (4.3c)$$

$$x_j \leq x_i, \quad \forall i \in \mathcal{N}_k, j \in \mathcal{N}_{k+1}, k \in \mathcal{K} \setminus \{K\}, \quad (4.3d)$$

$$x_j = x_i, \quad \forall i, j \in \mathcal{N}_k, k \in \mathcal{K} \quad (4.3e)$$

where t_m represents the time assignment to the m -th FCG in a frame. Constraint (4.3b) ensures that the total time assignment has to be no more than the frame length. Constraint (4.3c) guarantees that the minimum traffic demand of any admitted user is satisfied. In other words, any user i can access the system only if its minimum demand can be satisfied, i.e., $x_i = 1$ only if $\sum_{m=1}^{|\mathcal{F}|} f_{i,m} t_m \geq \tau_i$. Constraints (4.3d) and (4.3e) represent priority constraints. Constraint (4.3d) forces a low priority user $j \in \mathcal{N}_{k+1}$ remaining silent if any higher priority user $i \in \mathcal{N}_k$ cannot access the system. Constraint (4.3e) ensures that users in the same priority have the same access opportunity.

Unfortunately, directly solving (P0) is a very challenging task because of the following reasons. First, there are a total of $2^{|\mathcal{N}|}$ possible transmission groups in \mathcal{N} and such number increases exponentially with the number of users. Thus, in order to find all FCGs (i.e., coefficients $f_{i,m}$ in (P0)), all $2^{|\mathcal{N}|}$ possible groups need to be examined in the worst case. The time complexity of this examination procedure increases exponentially with the number of users. Even if we could find all FCGs, the number of them, $|\mathcal{F}|$, may also be huge and increase significantly with the number of

users. This will lead to a substantial number of variables since each FCG \mathbf{f}_m requires a variable t_m to denote its time assignment. In addition, constraint (4.3c) makes (P0) become a mixed nonlinear program. It is well known that solving a nonlinear program with a large number of variables is a very difficult task. All these motivate us to seek a novel solution.

Based on priority constraints (4.3d) and (4.3e), it is obvious that (P0) is equivalent to find a critical priority level, denoted as k^* , so that the minimum demands of all users in the first k^* high priority levels $\mathcal{N}_1 \cup \dots \cup \mathcal{N}_{k^*}$ can be satisfied while all requests from the first $k^* + 1$ high priority levels cannot be met any more. Thus, instead of directly solving (P0), we can examine priority levels one by one starting from the highest. Specifically, we first check if the minimum demands of all users in the highest priority level \mathcal{N}_1 can be met. If the check fails, then the examination process stops and outputs $k^* = 0$. Otherwise, we proceed to check if the demands of all users in the first two high priority levels $\mathcal{N}_1 \cup \mathcal{N}_2$ can be met or not. The failed check stops the examination process and outputs $k^* = 1$. Otherwise, we continue to check the first three high priority levels $\mathcal{N}_1 \cup \mathcal{N}_2 \cup \mathcal{N}_3$. This process will be repeated until adding one more priority level makes the system become overloaded or the lowest priority level has been considered. Obviously, at each examination k , we need to verify if the following problem has feasible solutions:

$$(\mathbf{P1}) : \max_{t_m \geq 0, x_i \in \{0,1\}} \sum_{i \in \mathcal{N}} x_i \quad (4.4a)$$

$$\text{s.t.} \quad \sum_{m=1}^{|\hat{\mathcal{F}}_k|} t_m \leq T, \quad (4.4b)$$

$$\sum_{m=1}^{|\hat{\mathcal{F}}_k|} f_{i,m} t_m \geq \tau_i, \forall i \in \hat{\mathcal{N}}_k, \quad (4.4c)$$

where $\hat{\mathcal{N}}_k = \mathcal{N}_1 \cup \dots \cup \mathcal{N}_k$ denotes the set of all users in the first k high priority levels

and $\hat{\mathcal{F}}_k$ denotes the set of all FCGs from $\hat{\mathcal{N}}_k$. Note that (P1) is derived from (P0) by replacing demand constraint (4.3c) as (4.4c) and removing priority constraints (4.3d) (4.3e) in (P0). It is because in (P1) priority levels are checked one by one starting from the highest so that the priority order has already been guaranteed. If (P1) has feasible solutions, it means there exist suitable time assignments so that the demands from the first k priority levels can be satisfied. In order to determine the existence of a feasible time allocation t_m satisfying both constraints (4.4b) and (4.4c), we formulate a time minimization problem as

$$(\mathbf{P2}) : \min_{t_m \geq 0} z = \sum_{m=1}^{|\hat{\mathcal{F}}_k|} t_m \quad (4.5a)$$

$$\text{s.t. } \sum_{m=1}^{|\hat{\mathcal{F}}_k|} f_{i,m} t_m \geq \tau_i, \forall i \in \hat{\mathcal{N}}_k, \quad (4.5b)$$

Let z_{P2} and z_{P2}^* denote the feasible and the optimal values of the objective function of (P2), respectively. We have $z_{P2}^* \leq z_{P2}$. Obviously, $z_{P2}^* \leq T$ is a sufficient and necessary condition for feasibility of (P1), and $z_{P2} \leq T$ is a sufficient condition. In some cases, applying the sufficient condition can simplify the decision process significantly by avoiding the calculation of optimal solution z_{P2}^* . This idea will be applied in our solution design in Subsection 4.3.2. By sequentially solving (P2) with $k = 1, 2, \dots$, we can eventually find the critical priority k^* .

Note that we may further admit a subset of users in the priority levels $\{k^* + 1, \dots, K\}$ to allow more admitted users. However, it is difficult to find an optimal subset of users from priority levels $\{k^* + 1, \dots, K\}$ such that the total number of admitted users reaches its maximum. This is because the total number of user subsets may become too large, which can be calculated as $\binom{n}{1} + \dots + \binom{n}{n}$ where $n = |\mathcal{N} \setminus \hat{\mathcal{N}}_{k^*}|$. Obviously, this number grows significantly with n so that it is too time-consuming

to search the optimal one from all user subsets. In addition, since users may have heterogeneous transmission demands and locations, it is hard to determine an optimal search order for including users one by one into the system. To reduce the search cost, we design a greedy search algorithm. First, for each user i , we define an interference index as

$$\rho_i = \frac{\tau_i}{\sum_{j \in \mathcal{N}} \tau_j} \cdot \sum_{j \in \mathcal{N}} \tau_j c_{i,j}. \quad (4.6)$$

Note that a smaller interference index ρ_i means that user i will potentially cause less interference to other users. Then, the greedy search algorithm works as follows: The users in $\mathcal{N} \setminus \hat{\mathcal{N}}_{k^*}$ are checked one by one in the increasing order of their interference indices. At each checking, we examine if the current user i can join a subset of existing FCGs, denoted as \mathcal{F}_s , such that its minimum traffic demand can be met (i.e., $\sum_{m=1}^{|\mathcal{F}_s|} t_m \geq \tau_i$). If the checking succeeds, the user is admitted into the system and added to the associated FCGs. Otherwise, it is rejected. The checking process is repeated until all users in $\mathcal{N} \setminus \hat{\mathcal{N}}_{k^*}$ have been considered.

4.2.2 Throughput Maximization of Admitted Users

After deriving the critical priority k^* , we have determined the set of admitted users $\hat{\mathcal{N}}_{k^*}$. To make full use of limited resources, we further formulate a throughput (sum rate) maximization problem for admitted users as follows

$$(\mathbf{P3}) : \max_{t_m \geq 0} \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} \sum_{m=1}^{|\hat{\mathcal{F}}_{k^*}|} r_i f_{i,m} t_m \quad (4.7a)$$

$$\text{s.t.} \quad \sum_{m=1}^{|\hat{\mathcal{F}}_{k^*}|} t_m \leq T, \quad (4.7b)$$

$$\sum_{m=1}^{|\hat{\mathcal{F}}_{k^*}|} f_{i,m} t_m \geq \tau_i, \forall i \in \hat{\mathcal{N}}_{k^*}, \quad (4.7c)$$

Algorithm 4: Sequential Solution Framework

- 1 Step i): Find the critical priority level k^*
 - 2 Initialize $k^* = K$;
 - 3 **for** $k = 1, 2, \dots, K$ **do**
 - 4 Solve (P2) for z_{P2}^* ;
 - 5 **if** $z_{P2}^* > T$ **then**
 - 6 $k^* = k - 1$; break;
 - 7 Step ii): Maximize throughput of admitted users by solving (P3).
-

where constraint (4.7c) guarantees that the minimum demands of users in the first k^* high priority levels can be met. Note that, different from (P1), (P3) has been guaranteed to be feasible.

Obviously, the overall procedure to solve both admission control and throughput maximization represents a sequential solution framework, which has been summarized in Algorithm 4. Note that we mainly focus on a quasi-static scenario [55] [58] where the order of user priorities keeps invariant during a scheduling period. For the dynamic scenario, if a new high-priority user joins the system, we may choose either to rerun the admission control algorithm starting from the new user's priority level or to defer the new user's transmission until the next scheduling period. In the following sections, the details on solving (P2) and (P3) will be presented.

4.3 Column Generation-Based Solution to (P2)

In this section, we present how to efficiently solve (P2). It can be observed that although (P2) is a linear program, directly solving it still leads to an exponentially increased complexity. It is because the number of variables increases exponentially with the number of users $|\hat{\mathcal{N}}_k|$. Since there are a total of $2^{|\hat{\mathcal{N}}_k|}$ potential FCGs in

(P2). Fortunately, although the number of variables is huge, we can find an optimal solution where most of the variables equal zero because there are at most q non-zero basic variables in a linear program with q constraints. This implies that only a small number of FCGs have non-zero time assignments, which actually contribute to the objective function. Motivated by this property, we introduce the column generation method [112].

4.3.1 Column Generation Preliminaries

In general, the column generation method decomposes a master problem (MP) into two subproblems: a restricted master problem (RMP) and a pricing problem (PP). We then find a subset of columns (i.e., FCGs in this chapter) satisfying all constraints of MP as the initialization to RMP. Starting from this initial subset of columns, we solve RMP to obtain optimal dual variables (called pricing factors). These pricing factors are then passed to PP. After solving PP, we can obtain a new column (i.e., a new FCG in this chapter) with the most negatively (positively) reduced cost for a minimization (maximization) MP. If the most negatively (positively) reduced cost is less (greater) than zero, this new column will be included in the initial subset of RMP for the next round of iteration $n + 1$. Otherwise, the optimal solution to RMP at the current iteration n is already the optimal one to MP and the iteration process stops. Thus, the column generation method is an iterative algorithm, which alternatively solves RMP and PP so that the solution to RMP eventually converges to that of MP [112]. The flow chart of the solution procedure is shown in Fig. 4.2. In the following subsections, we will present an enhanced column generation method

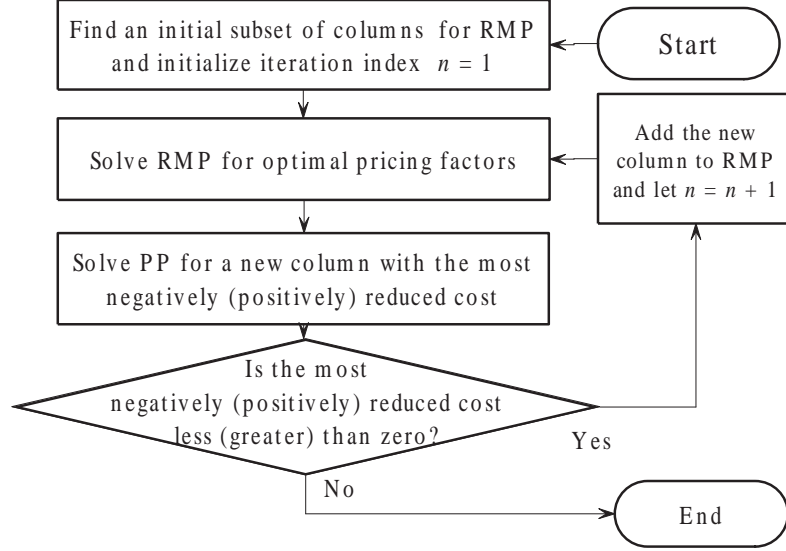


Figure 4.2: Flow chart of column generation method.

by integrating specific properties of (P2).

4.3.2 Restricted Master Problem (RMP)

Let $\hat{\mathcal{F}}'_k \subset \hat{\mathcal{F}}_k$ be a subset of FCGs from \mathcal{N}_k . Then, the RMP with respect to (w.r.t.) (P2) can be formulated as

$$(\mathbf{RMP-P2}) : \min_{t_m \geq 0} z = \sum_{m=1}^{|\hat{\mathcal{F}}'_k|} t_m \quad (4.8a)$$

$$\text{s.t.} \quad \sum_{m=1}^{|\hat{\mathcal{F}}'_k|} f_{i,m} t_m \geq \tau_i, \forall i \in \hat{\mathcal{N}}_k. \quad (4.8b)$$

An initial $\hat{\mathcal{F}}'_k$ is required to make (RMP-P2) feasible. Obviously, any set of FCGs which contains all users in $\hat{\mathcal{N}}_k$ is feasible since constraint (4.8b) can be satisfied by setting a sufficiently large time assignment t_m (e.g., $t_m = \max_{i \in \hat{\mathcal{N}}_k} (\tau_i), \forall m = 1, \dots, |\hat{\mathcal{F}}'_k|$). Thus, a most straightforward selection of $\hat{\mathcal{F}}'_k$ is a set of $|\hat{\mathcal{N}}_k|$ groups, each of which includes only one distinct user. However, such selection may not be

efficient. Intuitively, a good initial set should include as few groups as possible so that (RMP-P2) consists of fewer variables. Moreover, fewer groups implies that each group contains more users transmitting simultaneously, which in turn leads to a smaller total time assignment (i.e., a smaller initial objective function of (RMP-P2)) and thus reduces the number of iterations in the column generation method. Motivated by these, we propose a greedy initialization algorithm as follows to seek a good initial subset of FCGs.

A greedy initialization algorithm (GIA): The algorithm starts from an empty set and then considers all users in $\hat{\mathcal{N}}_k$ one by one. At each step, only one user (say user j) is considered. User j is added into one existing FCG if its introduction can keep the group feasible; otherwise, a new group is created with user j only. Since we only need to find a feasible set, the user selection in each step can be done randomly. In fact, by considering the overall solution framework as shown in Algorithm 4, the GIA can be further simplified. Consider the k -th examination in Algorithm 4. The user set $\hat{\mathcal{N}}_k$ can be divided into two parts: i) existing users at the previous examination $k - 1$, defined as $\hat{\mathcal{N}}_{k-1} = \{1, 2, \dots, |\hat{\mathcal{N}}_{k-1}|\}$; and ii) newly added users at the current examination k , defined as $\mathcal{N}_k = \{|\hat{\mathcal{N}}_{k-1}| + 1, |\hat{\mathcal{N}}_{k-1}| + 2, \dots, |\hat{\mathcal{N}}_k|\}$. In the solution to (RMP-P2) at examination $k - 1$, the FCGs having non-zero time assignments forms a set, denoted as $\hat{\mathcal{F}}_{k-1}^*$ (note that $\hat{\mathcal{F}}_0^* = \emptyset$). Then, all existing users in $\hat{\mathcal{N}}_{k-1}$ have already been well scheduled in $\hat{\mathcal{F}}_{k-1}^*$. Therefore, we can start from the set $\hat{\mathcal{F}}_{k-1}^*$ to run the GIA rather than an empty set. The details of this algorithm are summarized in Algorithm 5.

Given the initial $\hat{\mathcal{F}}'_k$, we can use the simplex method to solve (RMP-P2) to obtain

Algorithm 5: Greedy initialization algorithm.

```

1 Initialize:  $\hat{\mathcal{F}}'_k = \hat{\mathcal{F}}^*_{k-1}$ ;
2 for  $j = |\hat{\mathcal{N}}_{k-1}| + 1, \dots, |\hat{\mathcal{N}}_k|$  do
3     if  $j == 1$  then
4         Create a new group with only user  $j$  and add this group into  $\hat{\mathcal{F}}'_k$ ;
5     else
6          $flag1 = 0$ ;
7         for  $m = 1, 2, \dots, |\hat{\mathcal{F}}'_k|$  do           // Decide if user  $j$  can join one of
            existing groups in  $\hat{\mathcal{F}}'_k$ 
8              $flag2 = 1$ ;
9             for  $i = 1, 2, \dots, j - 1$  do       // Decide if user  $j$  can be added into
                group  $m$ 
10                if (User  $i$  is in group  $m$ )  $\&\&$  ( $c_{i,j} == 1$ ) then
11                     $flag2 = 0$ ; break;
12                if  $flag2 == 1$  then               // User  $j$  can join group  $m$ 
13                    Add user  $j$  into group  $m$ ;  $flag1 = 1$ ; break;
14        if  $flag1 == 0$  then                       // User  $j$  cannot join existing groups
15            Create a new group with only user  $j$  and add this group into  $\hat{\mathcal{F}}'_k$ ;
    
```

a primal optimal solution t_m^* , $m = 1, 2, \dots, |\hat{\mathcal{F}}'_k|$, and a dual optimal solution v_i^* , $i = 1, 2, \dots, |\hat{\mathcal{N}}_k|$.

4.3.3 Pricing Problem (PP)

The pricing problem with respect to (P2) is formulated as

$$(\mathbf{PP-P2}) : \min_{f_{i,m} \in \{0,1\}} \Delta_m = 1 - \sum_{i=1}^{|\hat{\mathcal{N}}_k|} v_i^* f_{i,m} \quad (4.9a)$$

$$\text{s.t. } f_{a,m} + f_{b,m} \leq 1, \forall a, b \in \hat{\mathcal{N}}_k \text{ with } c_{a,b} = 1, a \leq b, \quad (4.9b)$$

where constraint (4.9b) means that any two users interfering with each other cannot transmit simultaneously. The objective function of (PP-P2), Δ_m , represents the reduced cost caused by column \mathbf{f}_m , i.e., the amount that the objective function of

(RMP-P2) increases if \mathbf{f}_m is added to (RMP-P2) and the corresponding variable t_m increases by one unit. Note that $\Delta_m < 0$ implies a negative increase. The column with the most negatively reduced cost $\Delta^* < 0$ (Δ^* is the optimal objective function of (PP-P2)) can decrease the objective function of (RMP-P2) most. Therefore, this column should be passed to (RMP-P2) (i.e., the associated FCG should be added to $\hat{\mathcal{F}}'_k$) for the next round of iteration $n+1$. If $\Delta^* \geq 0$, the optimal solution to (RMP-P2) at the current iteration n is already the optimal one to (P2) and the iteration stops.

However, it has to be mentioned that (PP-P2) is an NP-hard problem (i.e., 0-1 integer program) and there are no known polynomial-time algorithms to find an optimal solution. Although existing integer linear programming algorithms (e.g., branch and bound algorithm (BBA)) can solve it optimally because of the relatively small number of variables, it is still necessary to reduce the dependence on its optimality so as to decrease the computational cost. It can be observed that (PP-P2) is equivalent to a maximum weighted independent set problem (WISP) in the graph theory. Thus, we define a graph G_k , where each user is treated as a vertex with weight v_i^* and there is an edge between any two interfered users. With G_k , our objective is translated to find a set of non-adjacent vertexes with the maximum weight sum (i.e., maximum weighted independent set). The newly formulated optimization problem is as follows

$$\text{(PP-P2 Equivalence)} : \max_{f_{i,m} \in \{0,1\}} \sum_{i=1}^{|\hat{\mathcal{N}}_k|} v_i^* f_{i,m} \text{ s.t. (4.9b).}$$

As shown in [113], (PP-P2 Equivalence) can be solved by a weighted greedy algorithm based on linear programming relaxation (WGL), which has an approximation ratio of $(\bar{d}_w(G_k) + 1)/2$, where $\bar{d}_w(G_k)$ indicates the average weighted degree of graph G .

Let $\alpha(G_k)$ denote the weight of a maximum weighted independent set on G_k and $A(G_k)$ represent the weight of the independent set obtained by WGL on G_k . From the approximation ratio, we have

$$\alpha(G_k)/A(G_k) \leq (\bar{d}_w(G_k) + 1)/2. \quad (4.11)$$

Since $\Delta^* = 1 - \alpha(G_k) \geq 1 - A(G_k) \cdot (\bar{d}_w(G_k) + 1)/2$, a lower bound of Δ^* can be represented as $\Delta^l = 1 - A(G_k) \cdot (\bar{d}_w(G_k) + 1)/2$. Also, based on $\alpha(G_k) \geq A(G_k)$, we can easily derive an upper bound of Δ^* as $\Delta^u = 1 - A(G_k)$.

The column generation process is stopped whenever $\Delta^* \geq 0$. However, if we run the approximation algorithm WGL for (PP-P2), Δ^* is unknown. In this case, the termination rule of column generation is modified as

- Scenario I: $\Delta^l < \Delta^u < 0$. In this scenario, $\Delta^* < 0$ and the optimality of (P2) has not been reached. The obtained column should be passed to (RMP-P2).
- Scenario II: $0 \leq \Delta^l < \Delta^u$. In this scenario, we have $\Delta^* > 0$ and the optimal solution of (P2) is already reached. The column generation process stops.
- Scenario III: $\Delta^l < 0$ and $\Delta^u \geq 0$. In this scenario, we cannot determine the sign of Δ^* , and the column obtained by WGL may become meaningless. If this scenario happens, we can either figure out the optimal solution to (PP-P2) by applying some optimal solvers or simply stop the column generation process and output an approximated solution.

4.3.4 Computational Cost Reduction

In this subsection, we further reduce the computational cost, by utilizing the sequential solution framework of Algorithm 4 and deriving bounds on the objective function of (P2). Recall that, in the sequential solution framework, we only need to decide if the objective function of (P2) satisfies condition $z_{P2}^* \leq T$ or $z_{P2}^* > T$, rather than computing its actual value. Thus, it is not necessary to solve (P2) optimally if the lower bound z_{P2}^l and/or the upper bound z_{P2}^u of z_{P2}^* can be calculated and they satisfy bound conditions $z_{P2}^l > T$ and/or $z_{P2}^u \leq T$. This also means pricing subproblem (PP-P2) is not necessarily solved optimally as long as the approximation solution can find some columns so that bound conditions are met. Following this way, the column generation process can be stopped in advance whenever bound conditions are met, regardless if it has reached optimality of (P2). Thus, the computational complexity can be reduced significantly.

Theorem 4.1. *Let z_{RMP-P2}^* represent the optimal objective function of (RMP-P2) at any iteration n . Then, the optimal objective function of (P2), z_{P2}^* , has an upper bound z_{RMP-P2}^* and a lower bound $\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i \tilde{v}_i$ where $\tilde{v}_i = \Delta^* + v_i^*$, i.e., $\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i \tilde{v}_i \leq z_{P2}^* \leq z_{RMP-P2}^*$.*

Proof. Please see Appendix B.1. □

Corollary 4.1. *Under the approximated solution $A(G_k)$ to (PP-P2 Equivalence) by WGL, z_{P2}^* has an upper bound z_{RMP-P2}^* and a lower bound $\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i (\Delta^l + v_i^*)$, where $\Delta^l = 1 - A(G_k) \cdot (\bar{d}_w(G_k) + 1)/2$ is a lower bound of Δ^* .*

Proof. i) The proof of the upper bound is the same as Theorem 4.1 (refer to Ap-

pendix B.1).

ii) The proof of the lower bound: Based on Theorem 4.1 as well as $\Delta^* \geq \Delta^l$, we can easily derive $z_{P2}^* \geq \sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i(\Delta^* + v_i^*) \geq \sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i(\Delta^l + v_i^*)$. \square

4.4 Column Generation Solution to (P3)

In this section, the solution details of (P3) is presented. Similar to (P2), the RMP w.r.t. (P3) can be formulated as

$$(\text{RMP-P3}) : \max_{t_m \geq 0} \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} \sum_{m=1}^{|\hat{\mathcal{F}}'_{k^*}|} r_i f_{i,m} t_m \quad (4.12a)$$

$$\text{s.t.} \quad \sum_{m=1}^{|\hat{\mathcal{F}}'_{k^*}|} t_m \leq T, \quad (4.12b)$$

$$\sum_{m=1}^{|\hat{\mathcal{F}}'_{k^*}|} f_{i,m} t_m \geq \tau_i, \forall i \in \hat{\mathcal{N}}_{k^*}, \quad (4.12c)$$

where $\hat{\mathcal{F}}'_{k^*} \subset \hat{\mathcal{F}}_{k^*}$. Since the optimal objective function of (P2) at $k = k^*$, z_{P2}^* , is less than the frame length T , the associated set of FCGs can be used as an initial $\hat{\mathcal{F}}'_{k^*}$ for (RMP-P3). Since (RMP-P3) is a linear program, it can be solved by the simplex method. Let $(t_m^*, m = 1, \dots, |\hat{\mathcal{F}}'_{k^*}|)$ and $(\varphi^*, w_i^*: i = 1, \dots, |\hat{\mathcal{N}}_{k^*}|)$ denote the primal and the dual optimal solutions to (RMP-P3), respectively, where φ^* corresponds to the first constraint and w_i^* the $(i + 1)$ -th constraint.

The pricing problem w.r.t. (P3) is

$$(\text{PP-P3}) : \max_{f_{i,m} \in \{0,1\}} \delta_m = \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} (r_i + w_i^*) f_{i,m} - \varphi^* \quad (4.13a)$$

$$\text{s.t.} \quad f_{a,m} + f_{b,m} \leq 1, \forall a, b \in \hat{\mathcal{N}}_{k^*}, a \leq b, c_{a,b} = 1. \quad (4.13b)$$

Note that different from (P2) which is a minimization problem, (P3) is a maximization problem. Thus, if the optimal objective function of (PP-P3) δ^* is larger than zero,

the derived FCG should be included in the initial $\hat{\mathcal{F}}'_{k^*}$ for the next round of iteration. Otherwise, the iteration process has already converged to the optimal solution to (P3) and the iteration terminates. Similar to (PP-P2), (PP-P3) can also be solved by WGL on a new graph G_{k^*} with the weight $r_i + w_i^*$ of vertex/user i . Accordingly, the lower and upper bounds of δ^* can be derived as $\delta^l = A(G_{k^*}) - \varphi^*$ and $\delta^u = A(G_{k^*})(\bar{d}_w(G_{k^*}) + 1)/2 - \varphi^*$, respectively, where $A(G_{k^*})$ denotes the weight of the independent set obtained by WGL on G_{k^*} and $\bar{d}_w(G_{k^*})$ is the average weighted degree on G_{k^*} .

At any iteration n of column generation, let z_{RMP-P3}^* and z_{P3}^* represent the optimal objective functions of (RMP-P3) and (P3), respectively. Obviously, z_{RMP-P3}^* is a lower bound of z_{P3}^* , i.e., $z_{RMP-P3}^* \leq z_{P3}^*$, and the gap $z_{P3}^* - z_{RMP-P3}^*$ represents how far the current solution z_{RMP-P3}^* is away from the optimal one. Thus, the iteration process can be stopped if $z_{P3}^* - z_{RMP-P3}^* \leq \epsilon$ where ϵ is an error tolerance. Unfortunately, z_{P3}^* is not known *a priori*. To address this issue, we replace z_{P3}^* by its upper bound, denoted as z_{P3}^u . Obviously, if $z_{P3}^u - z_{RMP-P3}^* < \epsilon$, then $z_{P3}^* - z_{RMP-P3}^* \leq \epsilon$. A possible upper bound can be derived based on the following theorem.

Theorem 4.2. *At any iteration n of the column generation method, an upper bound of the optimal objective function of (P3) is $z_{P3}^u = T\tilde{\varphi} - \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} \tau_i w_i^*$ where $\tilde{\varphi} = \delta^* + \varphi^*$.*

Proof. Please see Appendix B.2. □

Corollary 4.2. *With the approximated solution to (PP-P3) by WGL, an upper bound of z_{P3}^* can be represented as $z_{P3}^u = T(\delta^u + \varphi^*) - \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} \tau_i w_i^*$ where $\delta^u = A(G_{k^*})(\bar{d}_w(G_{k^*}) + 1)/2 - \varphi^*$.*

Proof. Based on Theorem 4.2 and $\delta^* \leq \delta^u$, it is easy to derive $z_{P_3}^* \leq T(\delta^* + \varphi^*) - \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} \tau_i w_i^* \leq (\delta^u + \varphi^*) - \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} \tau_i w_i^*$. \square

4.5 Discussions

In this section, we will summarize the proposed algorithm, discuss its computational complexity, and extend the binary protocol model to the signal-to-interference-plus-noise ratio model.

4.5.1 Algorithm Summary and Computational Complexity

Till now, we have presented the solution framework for admission control and throughput maximization in Section 4.2, and discussed the column generation based solution details to (P2) and (P3) in Sections 4.3 and 4.4. Also, we have derived both upper and lower bounds to accelerate the solution procedure. In this subsection, we summarize the proposed algorithm as shown in Algorithm 6, and discuss its computational complexity.

Note that pricing problem (PP-P2) can be solved by either the optimal BBA or the approximated WGL. In fact, we can combine both of them for solving (PP-P2). Specifically, we first adopt WGL with low complexity to generate a new column and if this new column can provide an improved solution to (RMP-P2), then we pass it to (RMP-P2) and continues the next iteration. Otherwise, if this new column cannot provide an improved solution, then we use the optimal BBA to find another new column. We call this algorithm as the hybrid WGL-BBA, which is used for the solution procedure of (P2) in Algorithm 6 for the explanation purpose.

Algorithm 6: Accelerated Algorithm

```

1 Step i): Find the critical priority level  $k^*$ 
2 Initialize  $k^* = K$ ;
3 for  $k = 1, 2, \dots, K$  do
4     Initialize (RMP-P2) by Algorithm 5;
5     while true do
6         Solve (RMP-P2) by the simplex method and update  $z_{P2}^u = z_{RMP-P2}^*$ ;
7         if  $z_{P2}^u \leq T$  then
8              $\quad$  continue;
9         Solve (PP-P2) by WGL and update  $z_{P2}^l$  (see Corollary 4.1),  $\Delta^u$  and  $\Delta^l$  (see
            Section 4.3.3);
10        if  $z_{P2}^l > T$  then
11             $\quad$  break;
12        if  $\Delta^l < \Delta^u < 0$  then
13             $\quad$  Add the new column generated by WGL to (RMP-P2);
14        else if  $0 \leq \Delta^l < \Delta^u$  then
15             $\quad$  Reach optimality of (P2) and let  $z_{P2}^* = z_{RMP-P2}^*$ ; break;
16        else if  $\Delta^l < 0$  and  $\Delta^u \geq 0$  then
17             $\quad$  Solve (PP-P2) by any existing optimal solver and update  $z_{P2}^l$  (see
                Theorem 4.1) and  $\Delta^*$ ;
18             $\quad$  if  $z_{P2}^l > T$  then
19                 $\quad$  break;
20             $\quad$  if  $\Delta^* < 0$  then
21                 $\quad$  Add the newly generated column to (RMP-P2);
22             $\quad$  else
23                 $\quad$  Reach optimality of (P2) and let  $z_{P2}^* = z_{RMP-P2}^*$ ; break;
24        if  $z_{P2}^l > T$  or  $z_{P2}^* > T$  then
25             $\quad$   $k^* = k - 1$ ; break;
26 Step ii): Solve (P3) similarly (details omitted here).
```

In Algorithm 6, the large-scale problems (P2) and (P3) with exponentially increased numbers of variables are solved by iteratively working on the decomposed subproblems with much smaller sizes so that the computational complexities mainly rely on the number of iterations in the course of column generation. As shown in [112], the column generation method inherits finiteness and correctness from the simplex-

type method in the number of column generations. Thus, for explanations, we borrow the theoretical results of the simplex method on computational complexity. It can be shown [114] that the average number of iterations (column generations in solving either (P2) or (P3)) is bounded by $O([\min(n_c, n_v)]^2)$ where n_c and n_v denote the numbers of constraints and unknown variables, respectively. For Algorithm 6, one can see that in the worst case, the average number of column generations is bounded by $O(KN^2 + (N + 1)^2) = O(KN^2)$ in the whole solution process since (P2) needs to be solved at most K times and there are total of $|\hat{\mathcal{N}}_k| \leq N$ and $|\hat{\mathcal{N}}_{k^*}| + 1 \leq N + 1$ constraints in (P2) and (P3), respectively. As mentioned earlier, it is not necessary to solve (P2) optimally in many cases by applying corresponding bounds, and thus the actual number of column generations can be largely reduced.

At each column generation for solving (P2), problems (RMP-P2) and (PP-P2) are required to be solved. (RMP-P2) is a linear program with at most $|\hat{\mathcal{N}}_k| \leq N$ constraints and the associated computational complexity is $O(N^2)$ in the worst case. Applying the WGL to (PP-P2) requires solving its relaxed linear program and thus it has a computational complexity of $O(N^2)$. Hence, the computational complexity of each column generation for solving (P2) is $O(N^2 + N^2) = O(N^2)$. Similarly, we can obtain that the computational complexity of each column generation for solving (P3) is also $O(N^2)$. In summary, the overall computational complexity of the accelerated algorithm for admission control and throughput maximization is $O(KN^2 \cdot N^2) = O(KN^4)$ in the worst case, which is polynomial with respect to the numbers of priority levels and users.

Note that we may also use the binary search algorithm for admission control.

With the binary search, the set of priority levels is divided by 2 and the medium priority level is checked at each time. If the check succeeds, i.e., the system can accommodate users in the first half priorities, then only the last half priorities need to be further checked. Otherwise, only the first half priorities need to be checked. It means the searching space is reduced by half at each time. This process is repeated until the searching space becomes empty. In this way, the computational complexity of searching the critical priority level can be reduced to $O(\log_2(K)N^4)$. However, in our case, the searching process is actually closely related to the follow-up implementation of column generation, and the overall computational complexity is dominated by the column generation procedure. With the sequential search framework, we can achieve a warm-start of the column generation method via the proposed greedy initialization algorithm. Specifically, the solution at the priority level k can be used to efficiently initialize the solution process of the next priority level $k + 1$, i.e., the solution at level k provides an effective initial subset of columns for level $k + 1$. By doing this, the number of column generations can be greatly reduced, so is the overall computational overhead. However, such benefit cannot be available for the binary search framework because priority levels are explored in a non-deterministic and non-sequential order. In the next section, we will compare these two algorithms through simulations.

4.5.2 Extension to the SINR Model

Under the signal-to-interference-plus-noise ratio (SINR) interference model [66] [67], a feasible candidate group (FCG) \mathbf{f}_m in \mathcal{N} requires to satisfy the following

conditions

$$\gamma_i = \frac{f_{i,m} P_i G_{i,i}}{\sigma^2 + \sum_{i' \in \mathcal{N} \setminus \{i\}} f_{i',m} P_{i'} G_{i',i}} \geq f_{i,m} \gamma_i^{th}, \quad \forall i \in \mathcal{N}, \quad (4.14a)$$

$$0 \leq P_i \leq P_i^{max}, \quad \forall i \in \mathcal{N}, \quad (4.14b)$$

where (4.14a) and (4.14b) are the SINR and power constraints, respectively. The SINR constraint means that if user i is active in the m -th FCG (i.e., $f_{i,m} = 1$), its received SINR γ_i should be no less than the associated SINR threshold γ_i^{th} . P_i represents the transmit power of user i , which should be non-negative and no greater than its power budget P_i^{max} as shown in (4.14b). $G_{i',i}$ denotes the inter-WBAN channel gain from user i' to user i , and $G_{i,i}$ represents the intra-WBAN channel gain from the sensor to the gateway of user i .

To integrate the SINR interference model into the column generation solution framework, only the pricing problems for generating new feasible candidate groups (FCGs) need to be modified accordingly. The new pricing problem with respect to (P2) can be formulated as

$$\begin{aligned} \text{(New PP-P2)} : \quad & \min_{f_{i,m}, P_i} \Delta_m = 1 - \sum_{i=1}^{|\hat{\mathcal{N}}_k|} v_i^* f_{i,m} \\ \text{s.t.} \quad & \frac{f_{i,m} P_i G_{i,i}}{\sigma^2 + \sum_{i' \in \mathcal{N} \setminus \{i\}} f_{i',m} P_{i'} G_{i',i}} \geq f_{i,m} \gamma_i^{th}, \quad \forall i \in \hat{\mathcal{N}}_k, \end{aligned} \quad (4.15a)$$

$$0 \leq P_i \leq P_i^{max}, \quad \forall i \in \hat{\mathcal{N}}_k, \quad (4.15b)$$

$$f_{i,m} \in \{0, 1\}, \quad (4.15c)$$

where (4.15a) and (4.15b) are the SINR and power constraints with respect to user i in $\hat{\mathcal{N}}_k$. In order to solve (New PP-P2), we first transform the SINR and power

constraints (4.15a) (4.15b) as

$$\frac{P_i G_{i,i}}{\sigma^2 + \sum_{i' \in \mathcal{N} \setminus \{i\}} P_{i'} G_{i',i}} \geq f_{i,m} \gamma_i^{th}, \quad \forall i \in \hat{\mathcal{N}}_k, \quad (4.16a)$$

$$0 \leq P_i \leq f_{i,m} P_i^{max}, \quad \forall i \in \hat{\mathcal{N}}_k. \quad (4.16b)$$

Then, we can introduce an auxiliary variable $Y_{i,i'} = f_{i,m} P_{i'}$ so that (4.16a) can be linearized as

$$P_i G_{i,i} - f_{i,m} \gamma_i^{th} \sigma^2 - \sum_{i' \in \mathcal{N} \setminus \{i\}} Y_{i,i'} G_{i',i} \geq 0, \quad \forall i \in \hat{\mathcal{N}}_k, \quad (4.17a)$$

$$0 \leq Y_{i,i'} \leq f_{i,m} V, \quad \forall i \in \hat{\mathcal{N}}_k, i' \in \mathcal{N} \setminus \{i\}, \quad (4.17b)$$

$$P_{i'} + (f_{i,m} - 1)V \leq Y_{i,i'} \leq P_{i'}, \quad \forall i \in \hat{\mathcal{N}}_k, i' \in \mathcal{N} \setminus \{i\}, \quad (4.17c)$$

where V is a sufficiently large number satisfying $V \geq \max\{P_i^{max}, \forall i \in \hat{\mathcal{N}}_k\}$. Hence, after transformation, (PP-P2) becomes a mixed integer linear program with $2|\hat{\mathcal{N}}_k| + |\hat{\mathcal{N}}_k|(|\hat{\mathcal{N}}_k| - 1)$ variables and $2|\hat{\mathcal{N}}_k| + 2|\hat{\mathcal{N}}_k|(|\hat{\mathcal{N}}_k| - 1)$ constraints, which can be solved by any existing linear programming solver. Following the similar way, the pricing problem with respect to (P3) can also be modified and solved, which is omitted here for brevity.

4.6 Numerical Results

In this section, we will evaluate the performance of the proposed algorithm via simulations. The following simulation settings are used if not particularly specified. We consider a network where users are uniformly distributed in a $10 \text{ m} \times 10 \text{ m}$ area [50]. The number of priorities is set as $K = 5$ and each user randomly chooses a priority from \mathcal{K} with equal probabilities. The interference range between users is set

to $d^{th} = 3$ m and the frame length is set to $T = 1$ s. The minimum traffic demand of each user is randomly selected in $0 \sim 50,000$ bits and the transmission rate is set to $r = 250$ kbps. Note that parameter selection in our simulations is used for the purpose of explanations only. Similar observations can be obtained if other parameter values are chosen.

The convergence of the column generation method in solving (P2) is shown in Fig. 4.3, where the optimal branch and bound (BBA) is used for pricing problems. In two sub-figures, we plot the upper and lower bounds of the optimal objective function of (P2), z_{P2}^* , and the most negatively reduced cost Δ^* . From the figure, we can observe that the upper and lower bounds are merged to each other over iterations, i.e., both of them converge to the optimal value. The upper bound z_{RMP-P2}^* is decreasing (i.e., improving) with iterations because the column generation method adds a new column into (RMP-P2) at each iteration so that the optimization space of (RMP-P2) is expanded with iterations. Although the lower bound may not be necessarily increasing with iterations, its trend does. This is because the lower bound $\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i(\Delta^* + v_i^*)$ is proportional to the most negatively reduced cost Δ^* which does not necessarily increase with iterations. As the iteration goes, the most negatively reduced cost Δ^* approaches to zero and the lower bound converges to the optimal objective function of the dual problem of (P2), $\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i v_i^*$, which equals the optimal objective function of (P2), z_{P2}^* , based on the *Strong Duality Theorem*. We repeat simulations for (P3), as shown in Fig. 4.4, and similar observations can be obtained.

As mentioned earlier, pricing problems (PP-P2) and (PP-P3) are NP-hard and

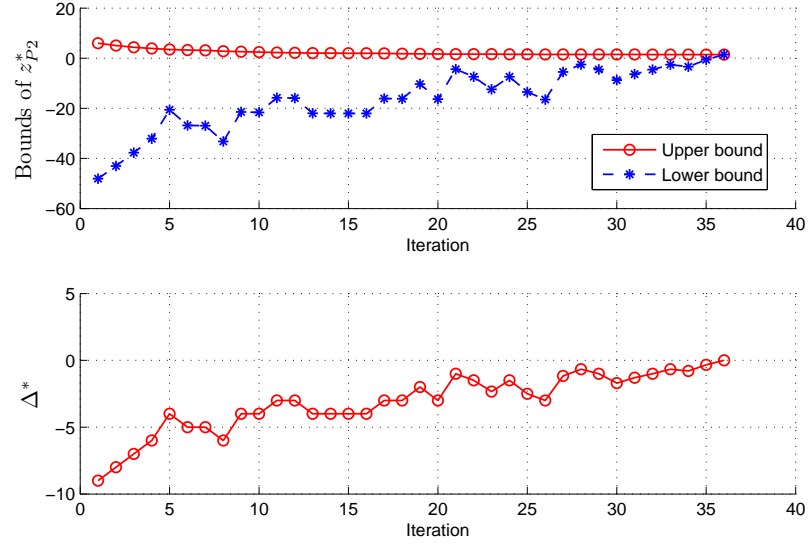


Figure 4.3: Convergence performance of the column generation method in solving (P2) when $N = 100$ and $k = k^*$.

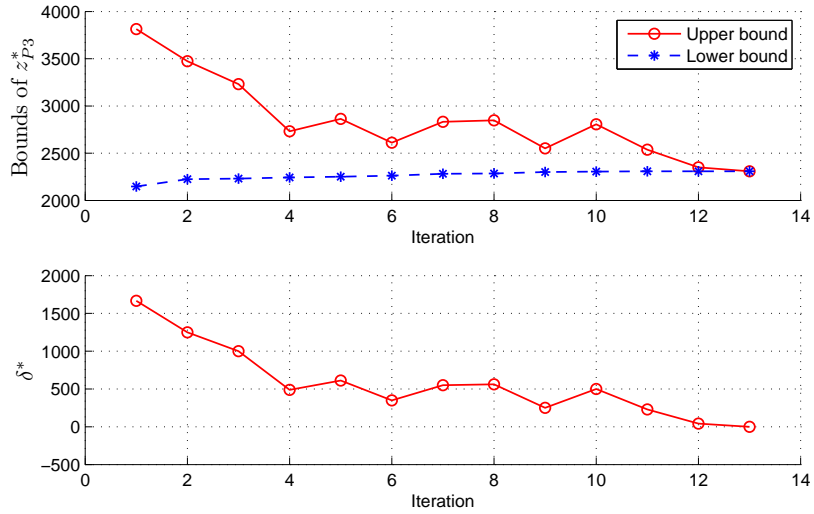


Figure 4.4: Convergence performance of the column generation method in solving (P3) when $N = 100$.

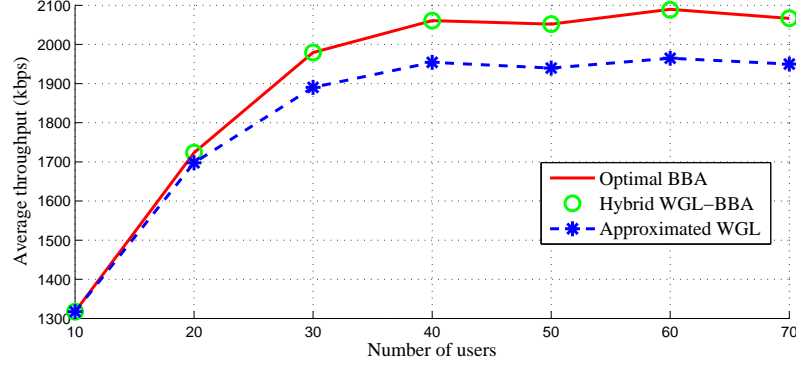


Figure 4.5: Total average throughput of admitted users under different algorithms.

they may have a great impact on the computational complexity of the proposed algorithm. Thus, for comparisons, we consider three algorithms separately for pricing problems: i) optimal branch and bound algorithm (BBA); ii) approximated weighted greedy algorithm based on linear programming relaxation (WGL); and iii) hybrid WGL-BBA algorithm, i.e., first apply WGL and then use BBA only if necessary (as shown in Section 4.5). Note that when applying these algorithms, different bounds are required accordingly, e.g. bounds of $z_{P_2}^*$ are given by Theorem 4.1 for BBA and they are provided by Corollary 4.1 for WGL.

Fig. 4.5 plots the total average throughput of admitted users under the proposed accelerated algorithm, where we only admit users in the first k^* high priorities. The results are averaged over 100 Monte Carlo simulations. From this figure, we can see that under the hybrid WGL-BBA algorithm and the pure WGL for pricing sub-problems, the overall accelerated algorithm can still obtain optimal and near-optimal throughput performance, respectively. In particular, the average throughput obtained by the hybrid WGL-BBA matches very well with that of the optimal BBA, while the average throughput obtained by the WGL is very close to the optimum in the low

Table 4.2: Probabilities of correctly finding the critical priority level k^* .

N	10	20	30	40	50	60	70
WGL-BBA	1.00	1.00	1.00	1.00	1.00	1.00	1.00
WGL	1.00	1.00	1.00	1.00	0.98	0.99	0.99

user density region ($N \leq 20$). Even at the extremely high density $N = 100$, the WGL can still result in 94% optimality. In addition, we can observe that in the low user density region ($N \leq 40$), the total system throughput improves with the number of users. It is because the spatial-reuse utilization increases with the user density. When the system becomes crowded ($N > 40$), however, the total system throughput fluctuates slightly as the number of users keeps increasing. This results from the priority constraint since only users in first k^* high priorities are admitted. Specifically, the number of priority levels that the system can support k^* may decrease when the total number of users increases and the system gets saturated, and the associated number of users that can be accommodated in a frame may decrease slightly as well, giving rise to a small reduction in throughput. Table 4.2 shows the probabilities of correctly searching the critical priority level k^* by the hybrid WGL-BBA and the pure WGL individually. These results are obtained via comparing with the optimal BBA. As we expected, the hybrid WGL-BBA can obtain the same critical priority k^* as the optimal BBA while WGL finds k^* with probabilities very close to one. These results demonstrate the effectiveness of the accelerated algorithm combined with hybrid WGL-BBA or WGL, which can provide optimal and near-optimal performance, respectively.

Table 4.3 shows the average computation time of the sequential solution process

Table 4.3: Average computation time.

N	10	20	30	40	50	60	70
BBA with GIA and bounds (s)	0.21	0.37	1.54	13.7	80.8	239	563
WGL-BBA with GIA and bounds (s)	0.24	0.59	1.94	10.5	54.2	141	211
WGL with GIA and bounds (s)	0.21	0.41	0.70	1.24	1.84	2.35	2.74
WGL without GIA and bounds (s)	0.67	1.89	3.19	5.65	7.21	9.73	13.34

under different scenarios. Three different algorithms (BBA, WGL, and hybrid WGL-BBA) are applied for pricing problems separately. In addition, to show the effects of the greedy initialization algorithm (GIA) and the derived bounds, we further compare two methods: i) One applies GIA and bounds to the sequential solution framework and ii) the other one does not. The computation time is obtained on a computer with a single-core central processing unit (CPU) of 3.6 GHz (Intel Core i5) and a random-access memory (RAM) of 4 GB. From the first three rows of the table, we can see that when the system has a small number of users ($N \leq 30$), the computation time is comparable, and the optimal BBA even outperforms other two algorithms. This is because the BBA can quickly find the optimal solutions to pricing problems when there are only a few users. However, when more users join the network, both the hybrid WGL-BBA and WGL show great advantages in computation time. The computation time using the optimal BBA increases substantially with the number of users, while the computation time using WGL increases only slightly. From the last two rows of the table, we can see that by applying the GIA and bounds, the average computation time is significantly reduced to about 25% of that based on the method without using GIA and bounds. Moreover, it is noted that applying the GIA and bounds to P2 does not affect the result of admission control, which means both

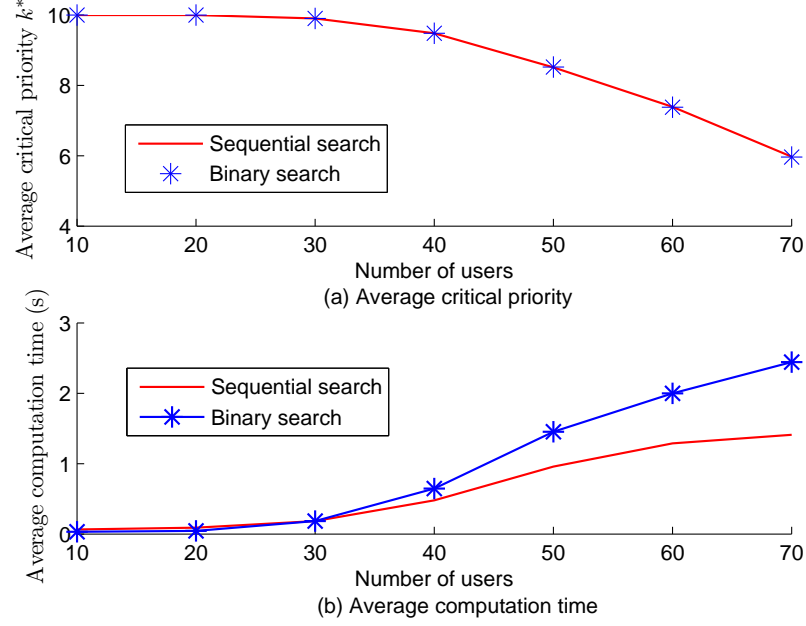


Figure 4.6: Comparisons of the proposed sequential search and the binary search when $K = 10$.

methods achieve the same performance in terms of the number of admitted users. The impact of applying GIA and bounds on throughput can be controlled to be very small (the gap between lower and upper bounds in simulations is set as $\epsilon = 0.01$) and thus both methods have very close throughput performance. By considering the results in Fig. 4.5, we can conclude that the accelerated algorithm combined with WGL, GIA, and bounds can more effectively and efficiently solve the scheduling problem for multi-user coexistence.

Fig. 4.6 compares the sequential and the binary search algorithms for admission control when the number of priority levels are set as $K = 10$. Fig. 4.6(a) shows the average critical priority level k^* versus the number of users, while Fig. 4.6(b) illustrates the average computation time versus the number of users. From Fig.

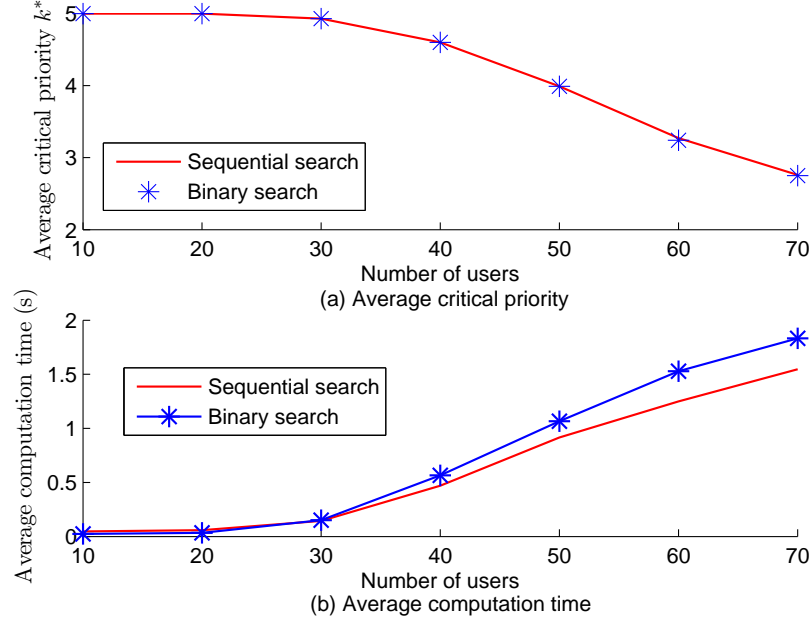


Figure 4.7: Comparisons of the proposed sequential search and the binary search when $K = 5$.

4.6(a), we can see that both sequential and binary search algorithms obtain the same critical priority level k^* . From Fig. 4.6(b), we can observe that the computation times of both the sequential and binary search algorithms are comparable when there are only a few users. However, when the number of users becomes large, the proposed sequential search leads to less computation time compared to the binary search. This is because the potential interference among users increases with the number of users, and in order to reduce the computation time, the column generation method requires an effective initialization algorithm to construct the initial subset of columns. The proposed sequential search effectively takes advantage of the greedy initialization algorithm, where the users in the first k priorities are effectively scheduled after checking priority k and this scheduling is passed to initialize the checking process of

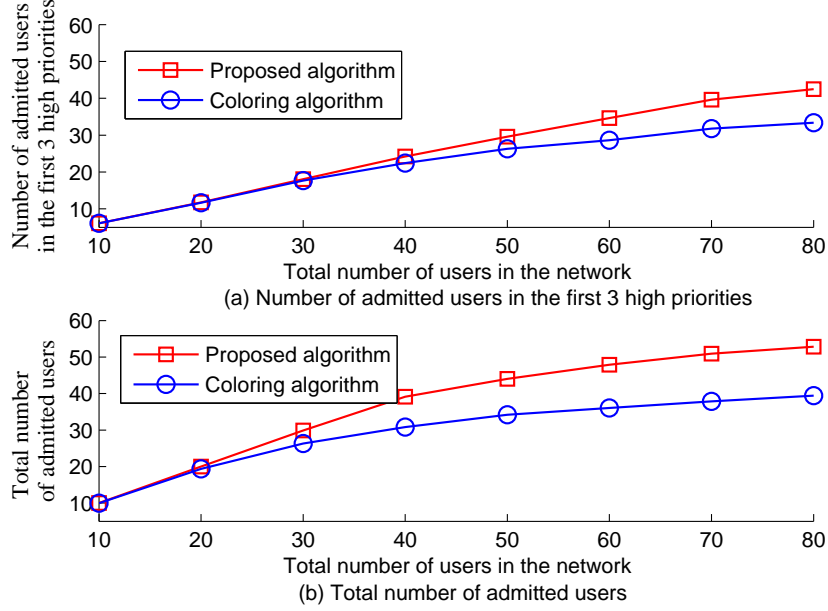


Figure 4.8: Average number of admitted users versus total number of users.

the next priority $k + 1$. However, the binary search does not effectively capitalize on this point. Similar observations can be seen when there are 5 priority levels, i.e., $K = 5$, as shown in Fig. 4.7. Hence, we can conclude that the proposed sequential search combined with the developed greedy initialization algorithm is more efficient than the binary search when there are a large number of users.

In the following, the coloring algorithm [50] is introduced as a benchmark for comparisons, where a frame T is divided into slots with equal lengths and each slot is represented by a color (the slot length τ is set as the maximum value of users' minimum time demands, i.e., $\tau = \max\{\tau_i, i \in \mathcal{N}\}$). Users are colored one by one in the decreasing order of priorities. Note that, in the following simulations, the system is allowed to not only admit users in the first k^* priorities but also include users in lower priorities $\{k^* + 1, \dots, K\}$.

Fig. 4.8 shows the effects of the total number of users in the network under both the proposed and the coloring algorithms. Fig. 4.8(a) and Fig. 4.8(b) illustrate the average number of admitted users in the first 3 high priorities and the total average number of admitted users, respectively. From this figure, we can see that when the total number of users is small, both the proposed algorithm and the coloring algorithm have the same number of admitted users. This is because all users' demands can be easily satisfied in the low user density scenario. However, when the total number of users becomes large, the proposed algorithm not only admits much more high priority users but also includes much more users in total, compared to the coloring algorithm. In particular, when the number of users is 80, the performance gain of the proposed algorithm over the coloring algorithm is about 1.35, in terms of the total number of admitted users. That is, the proposed algorithm increases by about 35% the total number of admitted users. This implies that the proposed algorithm has great advantages for handling the priority-aware admission control.

Fig. 4.9 shows the effects of the coexisting area size when the number of users is fixed at $N = 50$. Fig. 4.9(a) and Fig. 4.9(b) illustrate the average number of admitted users in the first 3 high priorities and the total average number of admitted users, respectively. From this figure, we can see that the average number of admitted users increases with the coexisting area for both the proposed and coloring algorithms. This is because the interference among users decreases when the coexisting area becomes larger. When the coexisting area is large enough ($\geq 20 \times 20 \text{ m}^2$), both the proposed and coloring algorithms can admit all users. However, when the coexisting area is small, the proposed algorithm can admit much more users (either high priority users

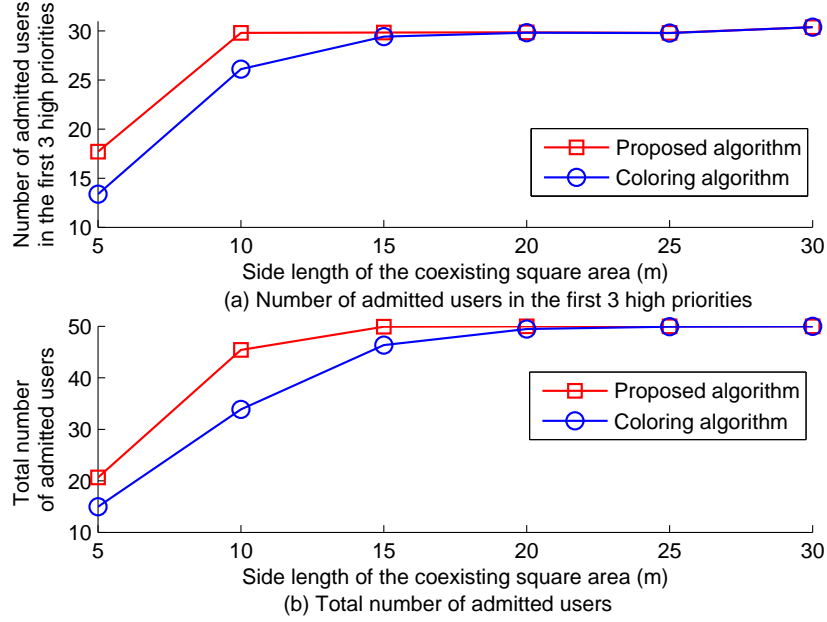


Figure 4.9: Average number of admitted users versus the coexisting area size.

or users in total) than the coloring algorithm. It means that the proposed algorithm is superior to the coloring algorithm in the high user density area, which is consistent with the observation results in Fig. 4.8.

Fig. 4.10 illustrates the effects of the traffic demand D (the minimum traffic demand of each user is randomly chosen from $[0, D]$) when the number of users is fixed at $N = 50$ and the coexisting area is set as $10 \times 10 \text{ m}^2$. Similar to Figs. 4.8 and 4.9, Fig. 4.10(a) and Fig. 4.10(b) illustrate the average number of admitted users in the first 3 high priorities and the total average number of admitted users, respectively. From this figure, we can see that the number of admitted users decreases with the traffic demand for both the proposed and coloring algorithms. When the traffic demand is small ($D < 20 \text{ kbits}$), both the proposed and coloring algorithms can admit all users. However, when the traffic demand becomes large, the proposed

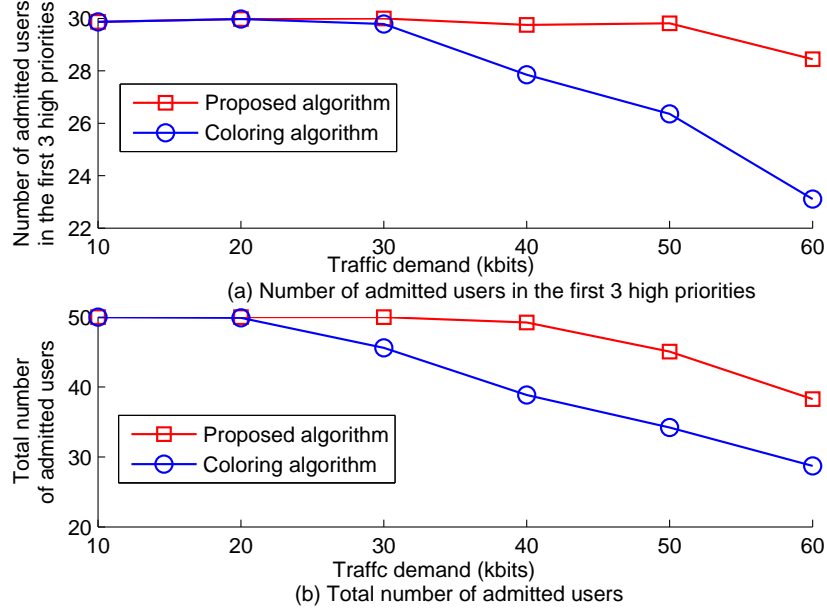


Figure 4.10: Average number of admitted users versus the traffic demand.

algorithm can accommodate more users (either high priority users or total users) than the coloring algorithm. This exhibits the superiority of the proposed algorithm in high traffic demand scenarios. To sum up, we can conclude that the proposed algorithm outperforms the coloring algorithm especially under the high user density region and the high traffic demand scenario. Similar observation results can be seen in terms of throughput, which are omitted here for brevity.

Fig. 4.11 compares the proposed algorithm and the coloring algorithm [50] in terms of the average received SINR. For fairness, the comparisons are done under the same admitted users and the same channel use time. To achieve this, the proposed algorithm is first applied to determine the set of admitted users \mathcal{N}_a and the channel use time $t = \sum_{m \in \mathcal{M}} t_m$. Then, in order to fit the coloring algorithm, the total channel use time t is divided into slots with equal lengths and each slot is represented by a color.

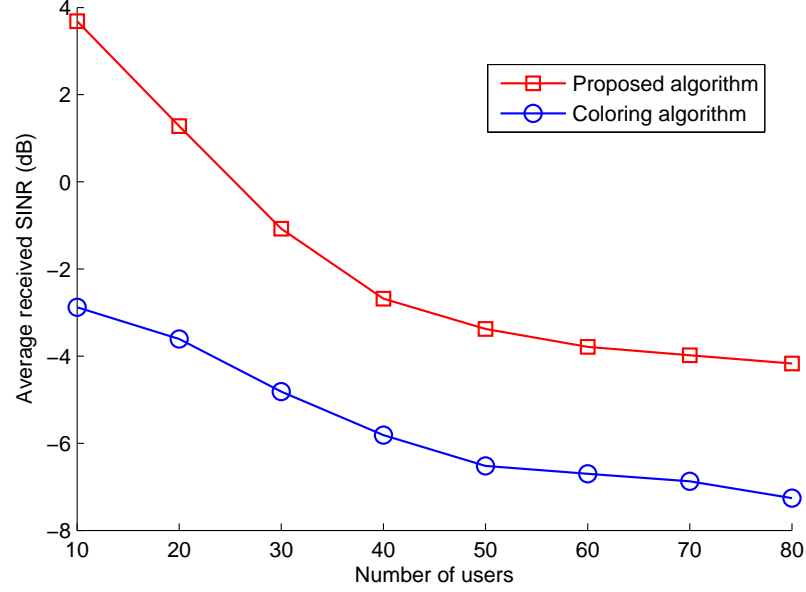


Figure 4.11: SINR comparisons of the proposed algorithm and the coloring algorithm.

We allocate users in \mathcal{N}_a slots by the coloring algorithm. If the coloring algorithm is not able to admit all users in \mathcal{N}_a , then the rest of uncolored users are randomly allocated slots in order to maintain the same number of admitted users. To calculate the SINR, we follow the channel model recommended by the IEEE 802.15.6 standard for WBANs [115] and the transmit power is set as -10 dBm. From this figure, we can see that the proposed algorithm always outperforms the coloring algorithm in terms of the average received SINR. This is because the proposed algorithm effectively deals with heterogeneous traffic demands while the coloring algorithm does not. In particular, when the number of users is 10, the proposed algorithm achieves an SINR gain of about 6 dB compared to the coloring algorithm. It implies that the proposed algorithm can more effectively control the interference among coexisting users.

Chapter 5

Joint Admission Control and Resource Management for D2D-Assisted Mobile Edge Computing

In this chapter, we will investigate the issues of admission control and resource management for mobile edge computing (MEC) systems with the assistance of device-to-device communication. Specifically, we will discuss the following issues: i) which computation requesters should be admitted (i.e., admission control); ii) which vacant mobile device or the server itself should be selected to serve each admitted requester (i.e., link scheduling/user association); iii) which channel should be chosen for each scheduled link; and iv) how much transmit power should be used on each link for offloading computation input data. To answer these questions, we mathematically

formulate an optimization problem where the objective is to maximize the number of admitted requesters and the limited radio and computing resource constraints are considered.

5.1 Related Works

5.1.1 Mobile Cloud Computing

Mobile cloud computing has appeared for many years and the related works can be roughly categorized into two branches: i) computation partitioning and offloading at mobile terminals, and ii) resource management and admission control at cloud servers.

Some applications at mobile devices may involve a large number of computations which require being partitioned into multiple tasks/components. Each computation task can be either run on the mobile device itself or offloaded to a cloud server for executions. To make mobile users achieve the better quality of experience, it is necessary to carefully partition mobile applications and choose appropriate computation tasks for offloading. This needs to be done based on the characteristics of mobile applications, the computing capabilities of cloud servers, and the capacities and loads of radio access networks. Yang et al. in [116] investigated the partitioning of multi-user computations aiming to minimize the average completion time. Zhang et al. in [73] developed an optimal offloading algorithm via the Markov decision process (MDP) where the intermittent connectivity between mobile users and cloudlets was considered. Chen et al. in [69] designed a distributed multi-user computation offloading

scheme based on the game theory for mobile edge cloud computing.

Resource management plays an important role in improving the performance of mobile cloud computing systems, which has been attracting a lot of attention from researchers. For examples, in [117], Kosta et al. designed a few types of virtual machines to provide on-demand resource allocations based on the mobile application requirements and the cloud's computing capacity. In [118], Liang et al. considered a geographically distributed mobile cloud system and investigated cloud resource management among multiple cloud domains. However, these works studied only the computing resource without taking the radio resource into consideration.

To further enhance the quality of service of mobile cloud computing systems, joint computing and radio resource allocations were studied from different application settings. Sardellitti et al. in [119] jointly optimized the transmit precoding matrices of mobile users and the computing capacity assignment, so as to minimize the total energy consumption of users while guaranteeing their latency constraints. Munoz et al. in [120] studied both the transmit precoder of a mobile user and the computation load distribution between the mobile user and the cloud provider in a femtocell cloud computing system. Zhang et al. in [74] considered joint computation offloading and channel assignment for the mobile edge computing system consisting of heterogeneous cells, with the objective to minimize the total energy consumption. Kaewpuang et al. in [121] investigated the sharing of both wireless bandwidth and cloud computing resources among multiple mobile service providers and their goal was to find the optimal number of applications supported to allow service providers achieve their maximum revenues. Liu et al. in [122] studied the joint allocations of wireless bandwidths and

virtual machines in a mobile cloudlet system, so as to improve the system QoS in terms of the number of admitted applications and the service latency.

Under limited computing and transmission capacities, admission control is also the key for performance enhancement, especially when enormous mobile users request cloud resources at the same time. The cloud server needs to decide which users can access the system and which users should be rejected, after receiving the computation offloading requests from mobile users. Normally, admission control is closely tied to the resource management of cloud servers. In [122], Liu et al. studied both the admission control of computation requests and the resource allocation for a cloudlet system, where a semi-Markov decision process was applied to solve the considered problem. In [123], Almeida et al. investigated a joint issue of admission control and computing capacity allocation on virtual machines, aiming to improve the system quality of service (QoS).

However, these works did not address user association in mobile cloud computing systems, i.e., how to pair computation requesters to computing service nodes. With the assistance of D2D communication, each computation requester may have multiple potential service nodes (the cloud server or vacant mobile devices). Therefore, it is necessary to take user association into consideration, and a joint resource management and admission control scheme is required for achieving effective and efficient mobile cloud computing.

5.1.2 D2D Communication

With the rapid growth of mobile traffic, D2D communication has appeared as a promising technique to expand the capacities of cellular networks. In [124], Xu et al. discussed the assignment of channels to multiple D2D communication pairs where a reverse iterative combinatorial spectrum auction was proposed to improve the system sum rate. In [125], Gu et al. presented both optimal and heuristic algorithms for carrier assignment with proportional fairness in D2D-based cellular networks. However, these works considered only wireless bandwidth allocation with fixed transmit power. In [126], Wang et al. studied the allocation of both radio and power resources to D2D communication pairs by applying an iterative combinatorial auction, with the objective to extend the battery lives of users. In [127], Zhang et al. jointly optimized link scheduling and power allocation for the system throughput maximization of D2D-assisted wireless caching networks. In [128], Cheng et al. introduced an optimal power control algorithm for D2D underlaying cellular networks, where their objective was to maximize the system throughput while making provision for statistical delay constraints. Besides, the interference coordination among D2D and cellular users was investigated in [129] by the joint spectrum and power allocation as well as the price control. The joint issue of mode selection, channel allocation and power control for D2D underlaying OFDMA-based wireless networks was analyzed in [130], where the goal was to minimize the total power consumption while guaranteeing individual user rates.

In summary, mobile cloud computing with the assistance of D2D communication was seldom investigated. In [131], Vallati et al. presented a mobile edge computing

(MEC) enabled smart home system architecture, where D2D communication was leveraged to enhance system performance. In [132], Jo et al. introduced a hierarchical cloud computing architecture where a dynamic cloud formed by mobile devices was added to a traditional static cloud by using D2D communication. In [133], Li et al. analyzed two schemes of access to the mobile cloud resource: optimal and periodic access schemes, where the intermittent connectivity of D2D links was taken into account. However, these works have not addressed the resource management for D2D-assisted mobile cloud computing systems. Different from existing works, this chapter focuses on MEC in D2D underlaying cellular networks and addresses the following issues:

- The management of multi-dimensional resources (power, spectrum and cloud computing resources) and the admission control of computation requests in D2D-assisted MEC systems are investigated.
- The offloading of cellular data traffic (cloud computing tasks) to D2D communication links (vacant mobile devices) is studied.
- The joint optimization problem of power allocation, channel assignment, and link scheduling is formulated and solved, where the objective is to maximize the cloud server's utility under limited resources.

5.2 System Model and Problem Formulation

In this section, we describe the system model of mobile edge computing with device-to-device (D2D) communication underlaying cellular networks. Consider a

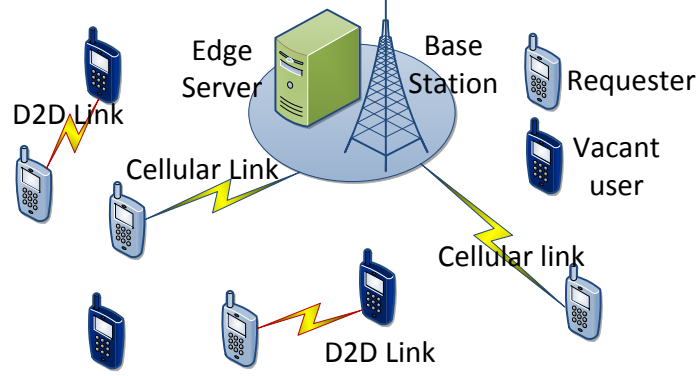


Figure 5.1: D2D-assisted mobile edge computing.

system which consists of a mobile edge computing server (simply called edge server (ES)) installed in a base station (BS), and a set of mobile users (MUs) \mathcal{N} , as shown in Fig. 5.1. The edge server is deployed to provide cloud computing services. The mobile users are divided into two groups: the users with computation requests, \mathcal{N}_r with cardinality of $|\mathcal{N}_r| = N_r$, and the rest of users, \mathcal{N}_v with cardinality of $|\mathcal{N}_v| = N_v$. Obviously, $\mathcal{N} = \mathcal{N}_r \cup \mathcal{N}_v$ with cardinality of $|\mathcal{N}| = N_r + N_v = N$. Each user in \mathcal{N}_r can obtain the computing service by offloading its computation task either to the ES via the cellular link or to one of nearby vacant users via the D2D link.

Following many previous works on mobile edge computing [69] [74] and D2D communication [126] [127], we consider a quasi-static scenario where the set of mobile users \mathcal{N} remains invariant during a computation offloading period but may change from one period to another. Each user in \mathcal{N}_r has a computation task requiring a computing service from the ES or one nearby vacant user. The admission control of computation requests and the scheduling of computing resources are managed by the ES as a central entity, i.e., the ES is responsible for deciding which computation requests should be admitted into the system and determining where each admitted

request should be sent to (the ES or one nearby vacant user). Next, we will discuss the computation and communication models separately.

5.2.1 Computation Model

We characterize the computation task of mobile user $i \in \mathcal{N}_r$ by a tuple $\mathcal{J}_i = \{b_i, d_i, T_i\}$, where b_i denotes the size of computation input data bits (such as program codes, input files, input parameters, etc.), d_i is the total amount of CPU cycles that the computation task \mathcal{J}_i requires, and T_i is the deadline requirement of task \mathcal{J}_i . Each user $i \in \mathcal{N}_r$ reports its computation task information \mathcal{J}_i to the ES at the beginning of each computation offloading period.

The computing speed (i.e., CPU cycles per second) for task \mathcal{J}_i of user i is ordinarily pre-determined based on the contract signed by user i with the cloud service provider [69] [74], which is denoted by f_i . The computation capacities (i.e., maximum CPU cycles per second) of the ES and the vacant user $j \in \mathcal{N}_v$ are denoted by F_{ES} and F_j , respectively. Then the computation execution time of task \mathcal{J}_i is

$$t_i^{exe} = d_i/f_i, \forall i \in \mathcal{N}_r. \quad (5.1)$$

The energy consumption of serving node j for accomplishing task \mathcal{J}_i is

$$e_j^{J_i} = \theta_j d_i, \forall j \in \mathcal{N}_v \cup \{ES\}, i \in \mathcal{N}_r, \quad (5.2)$$

where θ_j is the energy consumption per CPU cycle consumed by serving node j . Let c_j^e be the cost per unit energy of node j . Then, the total energy cost of node j for computing task \mathcal{J}_i is

$$C_j^{J_i} = c_j^e e_j^{J_i} = c_j^e \theta_j d_i, \forall j \in \mathcal{N}_v \cup \{ES\}, i \in \mathcal{N}_r. \quad (5.3)$$

If the ES delegates the computation task \mathcal{J}_i to one of the vacant users, it needs to pay the delegated user in order to encourage users' participation. Denote by $p_{ES}^{f_i}$ the price per CPU cycle that the ES pays to the delegated user when the required service speed is f_i . Such payment can be determined based on service contracts and can be in form of monetary remuneration or free data access. Then, the total price paid by the ES for delegating computation task \mathcal{J}_i is

$$\zeta_{ES}^{J_i} = p_{ES}^{f_i} \cdot d_i. \quad (5.4)$$

Thus, we can represent the set of vacant candidate users that i) can satisfy the computing speed requirement of task \mathcal{J}_i (i.e., $F_j \geq f_i$) and ii) are ensured to benefit from executing task \mathcal{J}_i (i.e., $\zeta_{ES}^{J_i} > C_j^{J_i}$) as

$$\mathcal{N}_i = \{j : j \in \mathcal{N}_v, F_j \geq f_i, \zeta_{ES}^{J_i} > C_j^{J_i}\}. \quad (5.5)$$

5.2.2 Communication Model

We consider the case that D2D links reuse the uplink channels occupied by cellular links because D2D links won't suffer serious interference from the high-power BS [126]. The set of M orthogonal uplink cellular channels is denoted as $\mathcal{M} = \{1, 2, \dots, M\}$. Without loss of generality, multiple D2D links are allowed to share the same channel with a cellular link. Considering the hardware limitation of mobile devices, we follow the convention in the literature [126] [127] that there is at most one outgoing link at each user in \mathcal{N}_r and one incoming link at each user in \mathcal{N}_v , where each link operates on a single channel. Let $x_{i,j}^m \in \{0, 1\}$ denote a binary indicator of joint requester-server pairing and channel allocation: $x_{i,j}^m = 1$ if the computation task of requester $i \in \mathcal{N}_r$ is

sent to the serving node $j \in \mathcal{N}_i \cup \{ES\}$ via the channel $m \in \mathcal{M}$; otherwise, $x_{i,j}^m = 0$.

Then, we have

$$\sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{m \in \mathcal{M}} x_{i,j}^m \leq 1, \quad \forall i \in \mathcal{N}_r, \quad (5.6)$$

$$\sum_{i \in \mathcal{N}_r} \sum_{m \in \mathcal{M}} x_{i,j}^m \leq 1, \quad \forall j \in \mathcal{N}_v, \quad (5.7)$$

where inequalities (5.6) and (5.7) impose the limitations of at most one outgoing link at each request user $i \in \mathcal{N}_r$ and at most one incoming link at each vacant user $j \in \mathcal{N}_v$, respectively.

Different from short-distance D2D transmissions between proximity users, cellular transmissions normally require high transmit power from the mobile users to the distant BS, which may cause large interference to D2D links sharing the same channel. Thus, to efficiently control the interference from cellular links, we restrict each channel to contain at most one cellular link [126] [124]. This constraint can be represented as

$$\sum_{i \in \mathcal{N}_r} x_{i,ES}^m \leq 1, \quad \forall m \in \mathcal{M}. \quad (5.8)$$

Denote by (i, j) the link from transmitter i to receiver j , and the set of all links is represented as $\mathcal{L} = \{(i, j) : i \in \mathcal{N}_r, j \in \mathcal{N}_i \cup ES\}$ with cardinality of $|\mathcal{L}| = L$. Let $P_{i,j}^m$ be the transmit power of transmitter i assigning to link (i, j) over channel m . Then, the transmit energy of requester i for offloading its computation input data to serving node j over channel m can be calculated by

$$e_{i,j}^{off,m} = P_{i,j}^m \cdot (T_i - t_i^{exe}), \quad \forall i \in \mathcal{N}_r, j \in \mathcal{N}_i \cup \{ES\}, m \in \mathcal{M}, \quad (5.9)$$

where $T_i - t_i^{exe}$ is the transmission time available for requester i to offload its data. Because of the limited battery capacities of mobile devices, the offloading energy

consumption of requester i requires to satisfy

$$\sum_{m \in \mathcal{M}} \sum_{j \in \mathcal{N}_i \cup \{ES\}} x_{i,j}^m e_{i,j}^{off,m} \leq e_i^{th}, \quad \forall i \in \mathcal{N}_r, \quad (5.10)$$

where e_i^{th} denotes the maximum energy threshold that requester i allows. This constraint is equivalent to

$$\sum_{m \in \mathcal{M}} \sum_{j \in \mathcal{N}_i \cup \{ES\}} x_{i,j}^m P_{i,j}^m \leq P_i^{th}, \quad \forall i \in \mathcal{N}_r, \quad (5.11)$$

where $P_i^{th} = \frac{e_i^{th}}{(T_i - t_i^{exe})}$ is the power threshold. Besides, since each transmitter $i \in \mathcal{N}_r$ has a total transmit power budget P_i^{max} , we must have

$$\sum_{m \in \mathcal{M}} \sum_{j \in \mathcal{N}_i \cup \{ES\}} x_{i,j}^m P_{i,j}^m \leq P_i^{max}, \quad \forall i \in \mathcal{N}_r. \quad (5.12)$$

Then, power threshold constraint (5.11) and power budget constraint (5.12) can be combined into a single constraint as

$$\sum_{m \in \mathcal{M}} \sum_{j \in \mathcal{N}_i \cup \{ES\}} x_{i,j}^m P_{i,j}^m \leq \min\{P_i^{th}, P_i^{max}\}, \quad \forall i \in \mathcal{N}_r. \quad (5.13)$$

Since each transmitter i is restricted on at most one channel (5.6), power constraint (5.13) can be rewritten as

$$\sum_{j \in \mathcal{N}_i \cup \{ES\}} x_{i,j}^m P_{i,j}^m \leq \min\{P_i^{th}, P_i^{max}\}, \quad \forall i \in \mathcal{N}_r, m \in \mathcal{M}. \quad (5.14)$$

The received signal-to-interference-plus-noise ratio (SINR) of receiver j on link (i, j) and channel m can be represented as

$$\gamma_{i,j}^m = \frac{x_{i,j}^m P_{i,j}^m G_{i,j}^m}{\sigma^2 + \sum_{(i',j') \in \mathcal{L} \setminus (i,j)} x_{i',j'}^m P_{i',j'}^m G_{i',j'}^m}, \quad \forall (i, j) \in \mathcal{L}, m \in \mathcal{M}, \quad (5.15)$$

where $G_{i,j}^m$ is the power gain of link (i, j) on channel m , and σ^2 denotes the background noise power. Then, according to the Shannon capacity formula, the transmission rate

of transmitter i on link (i, j) and over channel m is

$$R_{i,j}^m = W_m \log_2(1 + \gamma_{i,j}^m), \quad \forall (i, j) \in \mathcal{L}, \quad (5.16)$$

where W_m is the bandwidth of channel m . Recall that each requester $i \in \mathcal{N}_r$ is restricted to operate on at most one link and it uses at most one channel. Thus, the transmission rate of requester i needs to satisfy the following condition in order to finish offloading input data b_i before the available transmission time $T_i - t_i^{exe}$

$$R_{i,j}^m \geq b_i / (T_i - t_i^{exe}) \cdot x_{i,j}^m, \quad \forall (i, j) \in \mathcal{L}, m \in \mathcal{M}. \quad (5.17)$$

Note that the time overhead that the serving node requires to feed back the computation result to the requester is neglected. It is because the data size of computation result is normally much smaller than the input data size for many mobile applications (e.g., face recognition, virus scanning, etc.) [69] [74]. In fact, such small-size data of computation result can be sent back to the requester reliably via the control channel at a constant rate, so that the time needed can be regarded as a small constant [119], which does not affect our analyses. Equation (5.17) can be further represented as the SINR constraint

$$\gamma_{i,j}^m \geq \gamma_{i,j}^{th,m} \cdot x_{i,j}^m, \quad \forall (i, j) \in \mathcal{L}, m \in \mathcal{M}, \quad (5.18)$$

where $\gamma_{i,j}^{th,m}$ is the SINR threshold when link (i, j) operates on channel m and equals

$$\gamma_{i,j}^{th,m} = 2^{b_i / (W_m(T_i - t_i^{exe}))} - 1, \quad \forall (i, j) \in \mathcal{L}, m \in \mathcal{M}. \quad (5.19)$$

In addition, the scheduling of computation resource has to satisfy the ES's total computation capacity requirement

$$\sum_{i \in \mathcal{N}_r} \sum_{m \in \mathcal{M}} f_i \cdot x_{i,ES}^m \leq F_{ES}. \quad (5.20)$$

5.2.3 Problem Formulation

We first define the utility function of the ES as

$$U_{ES} = \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{m \in \mathcal{M}} a_{i,j} \cdot x_{i,j}^m, \quad (5.21)$$

where $a_{i,j}$ denotes the benefit that the ES can obtain when it schedules a serving node $j \in \mathcal{N}_i \cup \{ES\}$ to provide computing service for requester i . For example, if we let $a_{i,j} = 1$, then the utility function U_{ES} represents the total number of requesters that the ES allows to access the system. Another example is to define $a_{i,j}$ as

$$a_{i,j} = \begin{cases} p_{MU}^{f_i} d_i - C_{ES}^{J_i}; & j = ES \\ p_{MU}^{f_i} d_i - \zeta_{ES}^{J_i}; & j \in \mathcal{N}_i \end{cases} \quad (5.22)$$

where $p_{MU}^{f_i}$ denotes the price per CPU cycle that the ES charges from requester i who requires a computing speed f_i , $C_{ES}^{J_i}$ is the energy cost of the ES for computation task \mathcal{J}_i (5.3), and $\zeta_{ES}^{J_i}$ is the delegation cost of the ES (5.4). Then, the utility function implies the total net revenue of the ES.

The objective is to maximize the utility function of the ES, by jointly optimizing the admission control of computation requests, the scheduling of computing resources, channel assignment and power control, under a series of computation and communication resource constraints. This problem can be mathematically formulated as

$$\begin{aligned} [\text{P}] \quad & \max_{x_{i,j}^m, P_{i,j}^m} U_{ES} = \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{m \in \mathcal{M}} a_{i,j} x_{i,j}^m \\ \text{s.t.} \quad & (5.6) - (5.8), (5.14), (5.15), (5.18), (5.20) \\ & x_{i,j}^m \in \{0, 1\}, P_{i,j}^m \geq 0, \forall i \in \mathcal{N}_r, j \in \mathcal{N}_i, m \in \mathcal{M}. \end{aligned}$$

Solving the problem P is extremely challenging due to the following reasons: i) there are abundant combinations resulting from the requester-server association and

channel assignment, which make the formulated problem fall in the area of combinatorial optimization [126]; ii) the co-channel interference between cellular and D2D links is complicated to handle; and iii) even after removing those complicated constraints (5.14), (5.15), (5.18), (5.20), the relaxed problem turns out to be the maximum weighted independent set problem, which is well-known as NP-hardness [134].

In addition, although the assumption of homogenous channel conditions (i.e., $W_m = W, G_{i,j}^m = G_{i,j}, \forall m \in \mathcal{M}$) seems to be able to simplify the solution, it is unexpectedly not true when we apply the branch-and-bound method for integral optimal solutions. This is due to the inherent assignment symmetry in the homogeneous case, i.e., exchanging the link assignments on any two channels does not change the objective function value but only leads to different variable values. The symmetry implies that there are likely a large number of alternative optimal solutions scattering throughout the branch-and-bound tree, and these optimal solutions differ only by the indexes of channels. Thus, during exploring the branch-and-bound tree for determining optimal solutions, pruning branches by bounds becomes nearly useless and a mass of branches have to be explored before reaching optimality.

To address all these challenges, in the follows, we will propose a new solution framework for both homogeneous and heterogeneous channel conditions.

5.3 Extended Formulation and Solution Framework

In this section, we first reformulate the original problem P in an extended form, and then propose a new solution framework based on branch-and-price.

5.3.1 Extended Formulation

In this subsection, we will discuss the extended formulations of P under heterogeneous and homogenous channel conditions separately.

Heterogeneous Channel Conditions

It can be observed that problem P is a set partitioning problem where the task is to partition the set of links \mathcal{L} into M feasible link subsets that will be assigned to M channels, and the objective is to find a maximum-utility partitioning. Define a feasible candidate link subset (FCLS) on any channel m as a group of links that can share channel m to accomplish their own tasks of data transmission. Then, problem P is equivalent to find an optimal FCLS for each channel from all potential FCLSs.

Any FCLS k on channel m can be characterized by a vector $\mathbf{y}_k^m = [y_{i,j,k}^m : \forall (i, j) \in \mathcal{L}]$, where $y_{i,j,k}^m$ is a binary indicator: $y_{i,j,k}^m = 1$ means that link (i, j) belongs to this FCLS; $y_{i,j,k}^m = 0$, otherwise. Because of its feasibility, the vector \mathbf{y}_k^m satisfies both the power limitation and the SINR constraint as

$$\sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k}^m P_{i,j}^m \leq \min\{P_i^{th}, P_i^{max}\}, \quad \forall i \in \mathcal{N}_r, m \in \mathcal{M}, \quad (5.23)$$

$$\frac{y_{i,j,k}^m P_{i,j}^m G_{i,j}^m}{\sigma^2 + \sum_{(i',j') \in \mathcal{L} \setminus (i,j)} y_{i',j',k}^m P_{i',j'}^m G_{i',j'}^m} \geq y_{i,j,k}^m \gamma_{i,j}^{th,m}, \quad \forall (i, j) \in \mathcal{L}, m \in \mathcal{M}, \quad (5.24)$$

where (5.23) and (5.24) are based on (5.14) and (5.18), respectively, by replacing $x_{i,j}^m$ as $y_{i,j,k}^m$ to indicate that the k -th FCLS operates on channel m .

We represent the set of all FCLSs on channel m as $\mathcal{Y}_m = \{\mathbf{y}_1^m, \dots, \mathbf{y}_k^m, \dots, \mathbf{y}_{K_m}^m\}$ and denote the associated set of indexes as $\mathcal{K}_m = \{1, 2, \dots, K_m\}$, where K_m is the number of FCLSs on channels m . Then, we have a relation equality

$$x_{i,j}^m = \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m, \quad \forall (i, j) \in \mathcal{L}, m \in \mathcal{M}, \quad (5.25)$$

where λ_k^m is a binary variable: $\lambda_k^m = 1$ indicates channel m chooses the k -th FCLS; otherwise, $\lambda_k^m = 0$. Applying the relation equality (5.25) to P leads to the following extended formulation

$$[\text{EP1}] \quad \max_{\lambda_k^m} U_{ES} = \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}_m} a_{i,j} y_{i,j,k}^m \lambda_k^m$$

$$\text{s.t.} \quad \sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m \leq 1, \quad \forall i \in \mathcal{N}_r, \quad (5.26a)$$

$$\sum_{i \in \mathcal{N}_r} \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m \leq 1, \quad \forall j \in \mathcal{N}_v, \quad (5.26b)$$

$$\sum_{i \in \mathcal{N}_r} \sum_{k \in \mathcal{K}_m} y_{i,ES,k}^m \lambda_k^m \leq 1, \quad \forall m \in \mathcal{M}, \quad (5.26c)$$

$$\sum_{i \in \mathcal{N}_r} \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}_m} f_i \cdot y_{i,ES,k}^m \lambda_k^m \leq F_{ES}, \quad (5.26d)$$

$$\sum_{k \in \mathcal{K}_m} \lambda_k^m \leq 1, \quad \forall m \in \mathcal{M}, \quad (5.26e)$$

$$\lambda_k^m \in \{0, 1\}, \quad \forall m \in \mathcal{M}, k \in \mathcal{K}_m. \quad (5.26f)$$

Note that we assume that all FCLSs $\mathcal{Y}_m, \forall m \in \mathcal{M}$, have already been enumerated in advance based on constraints (5.23) and (5.24). The detailed procedure of obtaining FCLSs will be discussed in the next section. Constraints (5.26a), (5.26b) and (5.26c) of EP1 are based on (5.6), (5.7) and (5.8) of P, respectively. Constraint (5.26d)

represents the computing capacity limitation which is based on (5.20). Constraint (5.26e) implies that each channel m serves at most one FCLS. Note that the power limitation (5.14) and the SINR constraint (5.18) of P have been transferred to the restrictions on FCLSs.

The extended formulation is beneficial because the integer programming relaxation of EP1 provides a tighter bound than directly relaxing P. The explanation is as follows. Based on the relation equality (5.25), we can see that in the relaxation of EP1, any solution $x_{i,j}^m$ is a convex combination of $y_{i,j,k}^m, \forall k \in \mathcal{K}_m$, whereas all non-convex combination solutions are excluded. However, when relaxing P, $x_{i,j}^m \in [0, 1]$ includes both convex and non-convex combinations of $y_{i,j,k}^m, \forall k \in \mathcal{K}_m$. This means the relaxation of EP1 narrows the solution space compared to the relaxation of P, leading to an improved bound.

Homogeneous Channel Conditions

Under homogeneous channel conditions, a same set of FCLSs is expected on all channels (i.e., $\mathcal{Y} = \mathcal{Y}_m, \forall m \in \mathcal{M}$), and these FCLSs share a same index set (i.e., $\mathcal{K} = \mathcal{K}_m, \forall m \in \mathcal{M}$). Let $y_{i,j,k} = y_{i,j,k}^m, \forall m \in \mathcal{M}$ and introduce an aggregated variable

$$\lambda_k = \sum_{m \in \mathcal{M}} \lambda_k^m. \quad (5.27)$$

Then under homogeneous channel conditions, EP1 can be rewritten as

$$\begin{aligned} [\text{EP2}] \quad & \max_{\lambda_k} U_{ES} = \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{k \in \mathcal{K}} a_{i,j} y_{i,j,k} \lambda_k \\ \text{s.t.} \quad & \sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{k \in \mathcal{K}} y_{i,j,k} \lambda_k \leq 1, \quad \forall i \in \mathcal{N}_r, \end{aligned} \quad (5.28a)$$

$$\sum_{i \in \mathcal{N}_r} \sum_{k \in \mathcal{K}} y_{i,j,k} \lambda_k \leq 1, \quad \forall j \in \mathcal{N}_v, \quad (5.28b)$$

$$\sum_{i \in \mathcal{N}_r} \sum_{k \in \mathcal{K}} f_i \cdot y_{i,ES,k} \lambda_k \leq F_{ES}, \quad (5.28c)$$

$$\sum_{k \in \mathcal{K}} \lambda_k \leq M, \quad (5.28d)$$

$$\lambda_k \in \{0, 1\}, \quad \forall k \in \mathcal{K}, \quad (5.28e)$$

where constraints (5.28a), (5.28b) and (5.28c) come from (5.26a), (5.26b) and (5.26d), respectively. Constraint (5.28d) results from (5.26e) where we have $\sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}} \lambda_k^m \leq M$. By exchanging the summation order, we get $\sum_{k \in \mathcal{K}} \sum_{m \in \mathcal{M}} \lambda_k^m = \sum_{k \in \mathcal{K}} \lambda_k \leq M$. Note that constraint (5.26c) in EP1 can be transferred to a condition on any FCLS $k \in \mathcal{K}$, i.e.,

$$\sum_{i \in \mathcal{N}_r} y_{i,ES,k} \leq 1. \quad (5.29)$$

The condition (5.29) implies the limitation of at most one cellular link on each channel.

Note that by the variable aggregation (5.27), the extended formulation EP2 eliminates the assignment symmetry of P under homogeneous channel conditions (i.e., the channel index m has been removed), so that the number of optimal candidate solutions can be significantly reduced. Thus, this reformulation is much more suitable for applying pruning rules on the branch-and-bound tree, so that the integral optimal solution can be found in a much faster way.

5.3.2 Solution Framework

In this subsection, we propose a solution framework for deriving the integral optimal solution to the extended formulation problems EP1/EP2, following the idea of branch-and-price [135]–[137]. The branch-and-price integrates the column generation method into the branch-and-bound framework, where the column generation

method is used to solve the relaxed integer programming (IP) at each node of the branch-and-bound tree. Like branch-and-bound, the set of candidate solutions forms a rooted tree where the root node has the full solution set whereas each child node includes only a subset of candidate solutions. Each node corresponds to a relaxed IP with the same objective function as EP1/EP2 but with a distinct set of constraints. The flow chart of the solution framework is shown in Fig 5.2, which mainly consists of two components: branching and pruning rules. The branching rule determines how to form the branches of the tree, while the pruning rule aims to bypass the branches that cannot produce the optimal solutions. These two rules are discussed as follows in details.

Branching occurs when the relaxed IP at any node has a fractional solution. If this situation happens, this node, as a parent node, can be divided into two child nodes (sub-problems) by adding two extra constraints. By doing this, the solution space of the parent node is split into two subsets. The branching process starts with the relaxed EP1/EP2 as the root node of the whole tree and proceeds until the associated branch is pruned. Note that the added constraint at each child node is required to be compatible with the column generation method for solving the relaxed IP. In the following, we discuss branching constraints under heterogeneous and homogeneous channel conditions separately.

1) *Heterogeneous Channel Conditions:* In this case, we can choose to branch on the original variable $x_{i,j}^m$. Without loss of generality, assume that $x_{i,j}^m = \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m$ at a parent node with a fractional value of α ($0 < \alpha < 1$). Then this parent node is

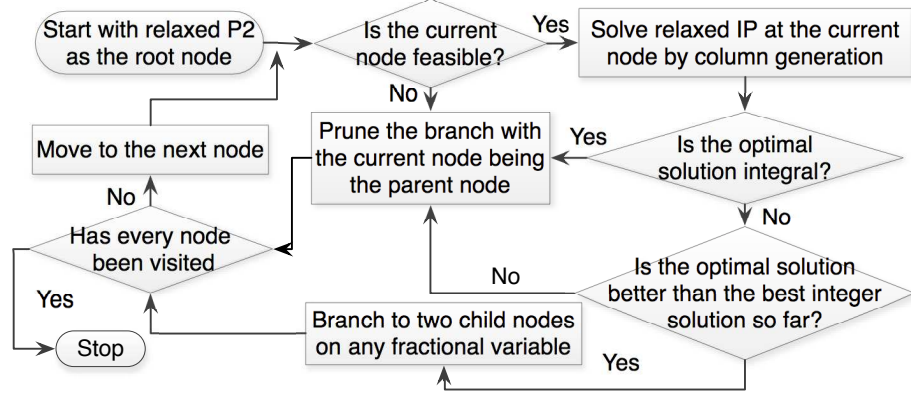


Figure 5.2: Flow chart of the solution framework

branched into two child nodes by adding the following two constraints respectively

$$x_{i,j}^m = \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m = 0, \quad (5.30a)$$

$$x_{i,j}^m = \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m = 1. \quad (5.30b)$$

2) *Homogeneous Channel Conditions:* In this case, branching on the original variable $x_{i,j}^m$ will reintroduce the channel assignment symmetry since the channel index m is re-added and the aggregated variable (5.27) is expanded. Observing that if a parent node has a fractional component in the solution vector $\lambda = [\lambda_1, \dots, \lambda_{|\mathcal{K}|}]$, there must be a variable index subset $\hat{\mathcal{K}}_q \subseteq \mathcal{K}$ such that $\sum_{k \in \hat{\mathcal{K}}_q} \lambda_k = \beta_q$ is fractional. Thus, the branching constraints at two child nodes can take the following forms respectively

$$\sum_{k \in \hat{\mathcal{K}}_q} \lambda_k \leq \lfloor \beta_q \rfloor, \quad (5.31a)$$

$$\sum_{k \in \hat{\mathcal{K}}_q} \lambda_k \geq \lceil \beta_q \rceil. \quad (5.31b)$$

The index subset $\hat{\mathcal{K}}_q$ can be defined by restricting part of links, i.e.,

$$\hat{\mathcal{K}}_q = \{k : y_{i,j,k} = \rho_{i,j}^q, \forall (i,j) \in \hat{\mathcal{L}}_q\}, \quad (5.32)$$

where $\rho_{i,j}^q \in \{0, 1\}$ denotes the state of link (i, j) and $\hat{\mathcal{L}}_q \subseteq \mathcal{L}$ is a link subset. The index subset $\hat{\mathcal{K}}_q$ is determined exclusively by the link subset $\hat{\mathcal{L}}_q$ and the associated state vector $\boldsymbol{\rho}_{\hat{\mathcal{L}}_q} = [\rho_{i,j}^q, (i, j) \in \hat{\mathcal{L}}_q]$.

A branch can be pruned safely if we can determine that this branch does not contain the optimal solution to EP1/EP2. To this end, we can work on the parent node of this branch, which includes the candidate solution subset of the whole branch and the associated optimal solution is the best one on this branch. Specifically, a branch can be excluded safely if the parent node of this branch satisfies one of the following three conditions: i) the relaxed IP at the parent node is infeasible, which means that there is no feasible solution in the whole branch; ii) the optimal solution at the parent node is already integral, which implies that the integral optimal solution of this branch is already found; and iii) either the optimal objective function of the relaxed IP at the parent node or its any upper bound is less than the objective function of EP1/EP2 under the best integer solution found so far, i.e., this branch is impossible to contain the integral optimal solution. The best integer solution is a lower bound on the objective function of EP1/EP2, which is determined from all feasible integer solutions found so far. A feasible integer solution may be obtained by rounding the fractional solution at any node to the nearest integer.

In the next section, we will apply the column generation method to solve the relaxed IP at each node in an iterative way, and at the same time derive the upper bound on the objective function at each iteration. The upper bound will help us avoid searching optimality at each node, so that the computational cost can be significantly reduced.

5.4 Column Generation Solution

In this section, we will discuss the solution procedure of the relaxed IP problem at each node of the branch-and-price tree. The relaxed IP at each node results from EP1/EP2 by adding branching constraints to narrow the solution space. In the follows, we consider any node u under heterogeneous and homogeneous channel conditions separately.

1) *Heterogeneous Channel Conditions:* In this case, the additional branching constraints at any node u can be represented as

$$x_{i,j}^m = \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m = 0, \quad \forall (i, j) \in \underline{\mathcal{L}}_m^u, m \in \underline{\mathcal{M}}^u, \quad (5.33a)$$

$$x_{i,j}^m = \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m = 1, \quad \forall (i, j) \in \overline{\mathcal{L}}_m^u, m \in \overline{\mathcal{M}}^u, \quad (5.33b)$$

where $\overline{\mathcal{L}}_m^u$ ($\underline{\mathcal{L}}_m^u$) denotes the set of links that are forced to operate (not to operate) on channel m at node u , and $\overline{\mathcal{M}}^u$ ($\underline{\mathcal{M}}^u$) is the associated channel set. Then, the relaxed IP at node u can be obtained by relaxing EP1 and adding the above branching constraints as

$$\begin{aligned} \text{[Relaxed IP1]} \quad \max_{\lambda_k^m \geq 0} U_{ES} &= \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}_m} a_{i,j} y_{i,j,k}^m \lambda_k^m \\ \text{s.t. } &(5.26a) - (5.26e), (5.33a), (5.33b). \end{aligned}$$

2) *Homogeneous Channel Conditions:* In this case, the branching constraints at node u become

$$\sum_{k \in \hat{\mathcal{K}}_q} \lambda_k \leq \lfloor \beta_q \rfloor, \quad \forall q \in \underline{\mathcal{Q}}^u, \quad (5.34a)$$

$$\sum_{k \in \hat{\mathcal{K}}_q} \lambda_k \geq \lceil \beta_q \rceil, \quad \forall q \in \overline{\mathcal{Q}}^u, \quad (5.34b)$$

where \underline{Q}^u (\overline{Q}^u) is the set of indexes associated with less-than-equal-to (greater-than-equal-to) branching constraints at node u . Then, the relaxed IP at node u becomes

$$\begin{aligned} \text{[Relaxed IP2]} \quad & \max_{\lambda_k \geq 0} U_{ES} = \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} \sum_{k \in \mathcal{K}} a_{i,j} y_{i,j,k} \lambda_k \\ \text{s.t.} \quad & (5.28a) - (5.28d), (5.34a), (5.34b). \end{aligned}$$

Note that constraint (5.28a) implicitly represents the relaxation of integer variables λ_k , i.e., $\lambda_k \leq 1, \forall k \in \mathcal{K}$.

Unfortunately, under both heterogeneous and homogeneous channel conditions, directly solving the relaxed IP is a challenging task due to the following reasons. First, we have to find all FCLSs before solving the relaxed IP. However, the searching procedure is time-consuming since the number of all possible link subsets (i.e., $2^{|\mathcal{L}|}$) increases exponentially with the number of links. Thus, in order to find all FCLSs, we may have to examine all possible link subsets in the worst case. Second, even if we could enumerate all FCLSs, the induced huge number of variables hinder us from directly solving the relaxed IP.

In order to circumvent this difficulty, we introduce the column generation method to solve the relaxed IP where FCLSs are generated on-the-fly instead of in advance. More importantly, only a small portion of FCLSs are required in column generation because we can find an optimal solution with at most η non-zero basic variables in a linear program with η constraints while most of other variables are zero.

A brief introduction of the column generation method [112] can be seen in Section 4.3.1. It decomposes a master problem (MP) (i.e., Relaxed IP1/IP2 in this chapter) into a restricted master problem (RMP) and a pricing problem (PP). The RMP is

first initialized by a subset of columns (i.e., FCLSs in this chapter) and then solved optimally for a dual solution (i.e., a set of pricing factors). The pricing factors are passed to the PP for generating a new column. If the new column has a positively increased profit, then it is added into the initial column subset of the RMP for solution improvement, and the next round of iteration begins. Otherwise, the optimal solution to the RMP has already converged to that of the MP and the iteration process terminates. In the following subsections, we show the solution procedures in details for Relaxed IP1 and IP2.

5.4.1 Solution to Relaxed IP1

The RMP can be formulated as

$$\begin{aligned}
 \text{[RMP1]} \quad & \max_{\lambda_k^m \geq 0} U_{ES} = \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}'_m} \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k}^m \lambda_k^m \\
 \text{s.t.} \quad & \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}'_m} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k}^m \lambda_k^m \leq 1, \forall i \in \mathcal{N}_r, \tag{5.35a}
 \end{aligned}$$

$$\sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}'_m} \sum_{i \in \mathcal{N}_r} y_{i,j,k}^m \lambda_k^m \leq 1, \forall j \in \mathcal{N}_v, \tag{5.35b}$$

$$\sum_{k \in \mathcal{K}'_m} \sum_{i \in \mathcal{N}_r} y_{i,ES,k}^m \lambda_k^m \leq 1, \forall m \in \mathcal{M}, \tag{5.35c}$$

$$\sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}'_m} \sum_{i \in \mathcal{N}_r} f_i \cdot y_{i,ES,k}^m \lambda_k^m \leq F_{ES}, \tag{5.35d}$$

$$\sum_{k \in \mathcal{K}'_m} \lambda_k^m = 1, \forall m \in \mathcal{M}, \tag{5.35e}$$

where $\mathcal{K}'_m \subset \mathcal{K}_m$. Note that the branching constraints (5.33a) (5.33b) in Relaxed IP1 are included in both the RMP1 and the subsequent pricing problem, which will be discussed later. On the branch-and-price tree, the root node can start from an empty initial \mathcal{K}'_m , and the child node can inherit the existing FCLSs from its parent node

while removing those FCLSs that do not satisfy the newly added branching constraint. Given \mathcal{K}'_m , RMP1 can be solved by the simplex method for a primal optimal solution (λ_k^{m*}) and a dual optimal solution $(\varphi_i^{r*}, \varphi_j^{v*}, \psi_m^*, \psi_{ES}^*, \mu_m^*)$ that are associated with constraints (5.35a), (5.35b), (5.35c), (5.35d), and (5.35e), respectively).

The pricing problem (PP) can be formulated as

$$\begin{aligned}
 \text{[PP1]} \quad \max_{P_{i,j}^m \geq 0, y_{i,j,k}^m} \quad & \Delta_k^m = \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k}^m - \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k}^m \varphi_i^{r*} \\
 & - \sum_{j \in \mathcal{N}_v} \sum_{i \in \mathcal{N}_r} y_{i,j,k}^m \varphi_j^{v*} - \sum_{i \in \mathcal{N}_r} y_{i,ES,k}^m \psi_m^* - \sum_{i \in \mathcal{N}_r} f_i \cdot y_{i,ES,k}^m \psi_{ES}^* - \mu_m^* \\
 \text{s.t.} \quad & y_{i,j,k}^m = 0, \quad \forall (i,j) \in \underline{\mathcal{L}}_m^u, m \in \underline{\mathcal{M}}^u, \tag{5.36a}
 \end{aligned}$$

$$y_{i,j,k}^m = 1, \quad \forall (i,j) \in \overline{\mathcal{L}}_m^u, m \in \overline{\mathcal{M}}^u, \tag{5.36b}$$

$$(5.23), (5.24), y_{i,j,k}^m \in \{0, 1\},$$

where (5.36a) and (5.36b) are based on the branching constraints (5.33a) and (5.33b), respectively. (5.23) and (5.24) are the power and SINR constraints, respectively. The objective function of PP1, Δ_k^m , implies the increased profit from the k -th FCLS ($k \in \mathcal{K}_m$), i.e., the amount that the objective function of RMP1 increases when the k -th FCLS is included into RMP1 and the associated variable λ_k^m increases by one unit. Denote by Δ_k^{m*} the most positively increased profit, i.e., the optimal objective function of PP1. Then, if $\Delta_k^{m*} \leq 0, \forall m \in \mathcal{M}$, the objective function of RMP1 cannot be improved anymore via adding new FCLSs and the current solution to RMP1 is already the optimal one to Relaxed IP1. Otherwise, it needs to seek a new FCLS with $\Delta_k^m > 0$ by solving PP1.

Branching constraints and the feasibility of Relaxed IP1: In the column generation formulation, branching constraints (5.33a) (5.33b) of Relaxed IP1

are guaranteed by constraint (5.35e) of RMP1 and constraints (5.36a) (5.36b) of PP1 together, i.e.,

$$x_{i,j}^m = \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m = \sum_{k \in \mathcal{K}_m} 0 \cdot \lambda_k^m = 0, \quad \forall (i,j) \in \underline{\mathcal{L}}_m^u, m \in \underline{\mathcal{M}}^u, \quad (5.37a)$$

$$x_{i,j}^m = \sum_{k \in \mathcal{K}_m} y_{i,j,k}^m \lambda_k^m = \sum_{k \in \mathcal{K}_m} 1 \cdot \lambda_k^m = 1, \quad \forall (i,j) \in \overline{\mathcal{L}}_m^u, m \in \overline{\mathcal{M}}^u. \quad (5.37b)$$

Due to these additional branching constraints at node u , Relaxed IP1 at this node may become infeasible. The Relaxed IP1 is feasible if and only if the set of active links determined by the branching constraints, $\overline{\mathcal{L}}_m^u$, satisfies conditions (5.23), (5.24), (5.35a)-(5.35d). Specifically, the feasibility of Relaxed IP1 can be determined by the following two conditions: i) Given $\overline{\mathcal{L}}_m^u, m \in \mathcal{M}$, an appropriate power allocation scheme should be found such that both the power and SINR constraints (5.23) (5.24) are satisfied. This can be easily determined by solving a set of linear SINR equations in $\overline{\mathcal{L}}_m^u$ (i.e., $P_{i,j} G_{i,j}^m / \gamma_{i,j}^{th,m} - \sum_{(i',j') \in \overline{\mathcal{L}}_m^u \setminus (i,j)} P_{i',j'} G_{i',j}^m = \sigma^2, \forall (i,j) \in \overline{\mathcal{L}}_m^u$), and the solution to these equations, $P_{i,j}$, should satisfy power constraint $0 \leq P_{i,j} \leq P_i^{max}, \forall (i,j) \in \overline{\mathcal{L}}_m^u$. ii) Denote the FCLS containing only the active link set $\overline{\mathcal{L}}_m^u$ as $\mathbf{y}_{\tilde{k}_m}^m = [y_{i,j,\tilde{k}_m}^m = 1 : (i,j) \in \overline{\mathcal{L}}_m^u; y_{i,j,\tilde{k}_m}^m = 0 : (i,j) \in \mathcal{L} \setminus \overline{\mathcal{L}}_m^u], \forall m \in \mathcal{M}$, and then these FCLSs $\{\mathbf{y}_{\tilde{k}_m}^m, m \in \mathcal{M}\}$ should satisfy constraints (5.35a)-(5.35d) by letting $\mathcal{K}_m = \{\tilde{k}_m\}$ and $\lambda_{\tilde{k}_m}^m = 1, \forall m \in \mathcal{M}$. If these conditions are not satisfied, Relaxed IP1 is infeasible and the associated branch can be pruned from the branch-and-price tree.

Solution to PP1: To facilitate the solution to PP1, we first transform power constraint (5.23) and SINR constraint (5.24) as

$$\sum_{j \in \mathcal{N}_i \cup \mathcal{E}\mathcal{S}} P_{i,j}^m \leq \min\{P_i^{th}, P_i^{max}\}, \quad \forall i \in \mathcal{N}_r, m \in \mathcal{M}, \quad (5.38)$$

$$\frac{P_{i,j}^m G_{i,j}^m}{\sigma^2 + \sum_{(i',j') \in \mathcal{L} \setminus (i,j)} P_{i',j'}^m G_{i',j}^m} \geq y_{i,j,k}^m \gamma_{i,j}^{th,m}, \quad \forall (i,j) \in \mathcal{L}, m \in \mathcal{M}, \quad (5.39)$$

$$P_{i,j}^m \leq y_{i,j,k}^m \cdot \min\{P_i^{th}, P_i^{max}\}, \quad \forall (i, j) \in \mathcal{L}, m \in \mathcal{M}. \quad (5.40)$$

Note that in the case that there is at most one outgoing link at each requester, constraint (5.38) is redundant because of constraint (5.40) and thus it can be removed.

Then, by introducing an auxiliary variable $Y_{i,j,i',j'}^m = y_{i,j,k}^m \cdot P_{i',j'}$, constraint (5.39) can be linearized as

$$P_{i,j} G_{i,j}^m - y_{i,j,k}^m \gamma_{i,j}^{th,m} \sigma^2 - \sum_{(i',j') \in \mathcal{L} \setminus (i,j)} Y_{i,j,i',j'}^m \gamma_{i,j}^{th,m} G_{i',j}^m \geq 0, \quad \forall (i, j) \in \mathcal{L}, m \in \mathcal{M}, \quad (5.41)$$

$$0 \leq Y_{i,j,i',j'}^m \leq y_{i,j,k}^m \cdot V, \quad \forall (i, j), (i', j') \in \mathcal{L}, m \in \mathcal{M}, \quad (5.42)$$

$$P_{i',j'} + (y_{i,j,k}^m - 1)V \leq Y_{i,j,i',j'}^m \leq P_{i',j'}, \quad \forall (i, j), (i', j') \in \mathcal{L}, m \in \mathcal{M}, \quad (5.43)$$

where V is a sufficiently large number such that $V > \max\{P_i^{max}, i \in \mathcal{N}_r\}$. Thus, PP1 becomes a small-scale mixed integer linear program with at most $2|\mathcal{L}| + |\mathcal{L}|^2$ variables, so that it can be solved by existing solvers (e.g., the function “intlinprog” in Matlab).

In order to further reduce the computational complexity, the column generation process can be terminated before it converges to the optimal solution to Relaxed IP1, by deriving the upper bound of the objective function. It is because if the upper bound is less than the objective function under the best integer solution found so far in the branch-and-price tree, then the branch with the parent node u can be confirmed not including the optimal solution and thus can be pruned safely from the branch-and-price tree. At each iteration of the column generation solution process, an upper bound can be derived based on the following theorem.

Theorem 5.1. *In each column generation, an upper bound on the optimal objective function of Relaxed IP1 can be represented as $\bar{U}_{IP1} = \sum_{i \in \mathcal{N}_r} \varphi_i^{r*} + \sum_{j \in \mathcal{N}_v} \varphi_j^{v*} + \sum_{m \in \mathcal{M}} \psi_m^* + F_{ES} \psi_{ES}^* + \sum_{m \in \mathcal{M}} (\mu_m^* + \Delta_k^{m*})$.*

Proof. Please see Appendix C.1. □

5.4.2 Solution to Relaxed IP2

The RMP with respect to Relaxed IP2 can be formulated as

$$\begin{aligned} [\text{RMP2}] \quad & \max_{\lambda_k \geq 0} U_{ES} = \sum_{k \in \mathcal{K}'} \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k} \lambda_k \\ \text{s.t.} \quad & \sum_{k \in \mathcal{K}'} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k} \lambda_k \leq 1, \quad \forall i \in \mathcal{N}_r, \end{aligned} \quad (5.44a)$$

$$\sum_{k \in \mathcal{K}'} \sum_{i \in \mathcal{N}_r} y_{i,j,k} \lambda_k \leq 1, \quad \forall j \in \mathcal{N}_v, \quad (5.44b)$$

$$\sum_{k \in \mathcal{K}'} \sum_{i \in \mathcal{N}_r} f_i \cdot y_{i,ES,k} \lambda_k \leq F_{ES}, \quad (5.44c)$$

$$\sum_{k \in \mathcal{K}'} \lambda_k \leq M, \quad (5.44d)$$

$$\sum_{k \in \hat{\mathcal{K}}'_q} \lambda_k \leq \lfloor \beta_q \rfloor, \quad \forall q \in \underline{\mathcal{Q}}^u, \quad (5.44e)$$

$$\sum_{k \in \hat{\mathcal{K}}'_q} \lambda_k \geq \lceil \beta_q \rceil, \quad \forall q \in \overline{\mathcal{Q}}^u, \quad (5.44f)$$

where $\hat{\mathcal{K}}'_q \subseteq \hat{\mathcal{K}}_q$ and $\hat{\mathcal{K}}'_q \subseteq \mathcal{K}' \subseteq \mathcal{K}$. To obtain the initial subset \mathcal{K}' , the child node u can inherit the existing FCLSs from its parent node. Note that the root node can start from an empty initial subset $\mathcal{K}' = \emptyset$.

Let $\varphi_i^{r*}, \varphi_j^{v*}, \psi_{ES}^*, \mu^*, \underline{\mu}_q^*, \overline{\mu}_q^*$ denote the dual optimal variables associated with constraints (5.44a), (5.44b), (5.44c), (5.44d), (5.44e), and (5.44f), respectively. Then, the pricing problem (PP) with respect to Relaxed IP2 can be formulated as

$$\begin{aligned} [\text{PP2}] \quad & \max_{P_{i,j} \geq 0, y_{i,j,k}, I_{q,k}} \Delta_k = \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k} - \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k} \varphi_i^{r*} \\ & - \sum_{\forall j \in \mathcal{N}_v} \sum_{i \in \mathcal{N}_r} y_{i,j,k} \varphi_j^{v*} - \sum_{i \in \mathcal{N}_r} f_i \cdot y_{i,ES,k} \psi_{ES}^* - \mu^* - \sum_{q \in \underline{\mathcal{Q}}^u} I_{q,k} \cdot \underline{\mu}_q^* + \sum_{q \in \overline{\mathcal{Q}}^u} I_{q,k} \cdot \overline{\mu}_q^* \end{aligned}$$

$$\text{s.t. } I_{q,k} = \begin{cases} 1, & \text{if } k \in \hat{\mathcal{K}}_q \\ 0, & \text{otherwise} \end{cases} \quad (5.45a)$$

$$(5.23), (5.24), (5.29), y_{i,j,k} \in \{0, 1\}.$$

Branching constraints and feasibility of Relaxed IP2: In the above column generation formulation, branching constraints (5.34a) (5.34b) of Relaxed IP2 are represented by constraints (5.44e) (5.44f) of RMP2. However, since RMP2 contains only partial FCLSs $\mathcal{K}' \subseteq \mathcal{K}$, it may be insufficient to decide the feasibility of Relaxed IP2. Fortunately, by carefully selecting some special FCLSs added into RMP2, we can justify the feasibility of Relaxed IP2 successfully. Observing that the “greater-than-equal-to” constraint (5.44f) is the key to prevent RMP2 from having feasible solutions, we can choose $|\overline{\mathcal{Q}}^u|$ special FCLSs where the q -th FCLS contains only those active links enforced by the index subset $\hat{\mathcal{K}}_q$ (i.e, the link set $\{(i, j) | \rho_{i,j}^q = 1\}$). It is obvious that these special FCLSs represents the minimum number of active links required by node u . Thus, if the RMP2 is still infeasible after adding these special FCLSs, then Relaxed IP2 must be infeasible, and the branch with u as the parent node can be pruned from the branch-and-price tree. Otherwise, Relaxed IP2 is feasible.

Solution to PP2: Based on the definition of index subset $\hat{\mathcal{K}}_q$ (5.32), constraint (5.45a) is equivalent to

$$I_{q,k} = \begin{cases} 1, & \text{if } \delta_{q,k} = 0 \\ 0, & \text{otherwise} \end{cases} \quad (5.46)$$

where $\delta_{q,k}$ is defined as

$$\delta_{q,k} = \sum_{(i,j) \in \hat{\mathcal{L}}_q, \rho_{i,j}^q=1} (\rho_{i,j}^q - y_{i,j,k}) + \sum_{(i,j) \in \hat{\mathcal{L}}_q, \rho_{i,j}^q=0} (y_{i,j,k} - \rho_{i,j}^q). \quad (5.47)$$

Note that $\delta_{q,k} \geq 0$. $\delta_{q,k} = 0$ implies that the current FCLS matches the definition of subset $\hat{\mathcal{K}}_q$, and otherwise, it does not. The conditional constraint (5.46) can be rewritten as the following linear constraints

$$-V_1 \cdot V_2 \cdot \delta_{q,k} \leq (I_{q,k} - 1) \cdot V_1 + \delta_{q,k} \leq 0, \quad (5.48a)$$

$$I_{q,k} \in \{0, 1\}, \quad (5.48b)$$

where V_1 and V_2 should satisfy conditions $V_1 \geq \delta_{q,k}$ and $V_2 \geq 1/\delta_{q,k}$, $\forall \delta_{q,k} \neq 0$. Since $0 \leq \delta_{q,k} \leq |\mathcal{L}|$, we can take $V_1 = |\mathcal{L}|$ and $V_2 = 1$. In addition, like the previous subsection, the power and SINR constraints (5.23)–(5.24) can be linearized as (5.38)–(5.43). Thus, PP2 becomes a small-scale mixed integer linear program problem and the existing solvers can be applied to solve it.

Similar to Relaxed IP1, in order to reduce the computational complexity, we derive the upper bound on the objective function of Relaxed IP2 in the following theorem.

Theorem 5.2. *In each column generation, an upper bound on the optimal objective function of Relaxed IP2 can be derived as $\bar{U}_{IP2} = \sum_{i \in \mathcal{N}_r} \varphi_i^{r*} + \sum_{j \in \mathcal{N}_v} \varphi_j^{v*} + F_{ES} \psi_{ES}^* + M(\mu^* + \Delta_k^*) + \sum_{q \in \underline{\mathcal{Q}}^u} \lfloor \beta_q \rfloor \underline{\mu}_q^* - \sum_{q \in \bar{\mathcal{Q}}^u} \lceil \beta_q \rceil \bar{\mu}_q^*$.*

Proof. Please see Appendix C.2. □

5.4.3 Computational Complexity Discussions

Till now, we have presented the overall solution procedure to the original problem P, which is first reformulated in an extended form (i.e., EP1/EP2) and then solved by the branch-and-price based solution framework. For the subproblem at each node of the branch-and-price tree, the column generation based algorithm is proposed to

solve it. The overall computational complexity mainly depends on the number of nodes that require to be explored on the branch-and-price tree and the number of column generations at each node.

Firstly, to reduce the number of nodes to explore, branching pruning by bounds is introduced to avoid enumerating all nodes. The bound is tightened by reformulating our original problem in an extended form. Furthermore, to avoid the inefficiency of branch pruning under homogeneous channel conditions with channel assignment symmetry, variable aggregation is applied to reformulate our problem. However, it has to mention that although many efforts have been made to reduce the number of nodes for exploring, it cannot be guaranteed to grow polynomially with the size of problem (i.e., the number of mobile users) [135]–[137].

Secondly, to efficiently solve each subproblem, the column generation method is introduced to address the exponential number of FCLSs/variables, where its computational complexity mainly relies on the number of column generations before reaching optimality. As shown in [112], column generation is computationally efficient in practice for linear programming with extensive implicit variables, and inherits finiteness and correctness from the simplex-type method in the number of column generations. For explanations, we borrow the theoretical results of the simplex method on computational complexity. It can be shown [114] that the average number of iterations (column generations) is bounded by $O([\min(n_c, n_v)]^2)$, where n_c and n_v denote the numbers of constraints and unknown variables, respectively. Thus, for Relaxed IP1 with at most $N + 2M + 1 + ML$ constraints, the average number of column generations is bounded by $O((N + 2M + 1 + ML)^2) = O(N^2 + M^2L^2 + NML)$.

While for Relaxed IP2, the average number of column generations is bounded by $O((N + 2 + 2L)^2) = O(N^2 + L^2 + NL)$, which is resulted from the facts that there are $N + 2 + |\underline{\mathcal{Q}}^u| + |\overline{\mathcal{Q}}^u|$ constraints in all and the cardinality of $\underline{\mathcal{Q}}^u$ ($\overline{\mathcal{Q}}^u$) is at most L on average (this is because in practice we seldom need to limit more than one link in determining $\hat{\mathcal{K}}_q, q \in \underline{\mathcal{Q}}^u(\overline{\mathcal{Q}}^u)$ [135] [136]).

5.4.4 Greedy Algorithm

To address the high computational complexity of the proposed optimal algorithm, in this subsection, we propose a greedy algorithm for obtaining a sub-optimal solution with low complexity. Note that the greedy algorithm can also be used to initialize the proposed branch-and-price based algorithm.

The basic idea of the greedy algorithm is to explore channels one by one. For a given channel, computation requesters are selected one after the other to operate on the channel. Specifically, denote by \mathcal{N}_r^{rest} the rest of requesters that have not been admitted in the system, which is initialized as $\mathcal{N}_r^{rest} = \mathcal{N}_r$. For each channel, we keep tracking two sets: \mathcal{N}_r^{un} represents the set of requesters that have not been considered and \mathcal{N}_r^{in} represents the set of requesters that have been allowed to operate on the channel. They are initialized as $\mathcal{N}_r^{un} = \mathcal{N}_r^{rest}$ and $\mathcal{N}_r^{in} = \emptyset$. At each time, the requester with the maximum number of potential D2D links (which implies the maximum number of potential vacant users) in \mathcal{N}_r^{un} is first chosen to check if it can operate on the channel with those existing requesters in \mathcal{N}_r^{in} simultaneously. If the check is successful, then this requester is moved from \mathcal{N}_r^{un} to \mathcal{N}_r^{in} . Otherwise, this requester is excluded from \mathcal{N}_r^{un} but not included in \mathcal{N}_r^{in} . This checking process is

Algorithm 7: Greedy Algorithm

```

1 Initialize:  $\mathcal{N}_r^{rest} = \mathcal{N}_r$  and  $F_{ES}^{rest} = F_{ES}$ ;
2 Sort channels  $\mathcal{M}$  in the decreasing order of bandwidths and assume
    $W_1 \geq \dots W_m \geq \dots \geq W_M$ ;
3 for Channel  $m = 1, \dots, M$  do
4   Let  $\mathcal{N}_r^{un} = \mathcal{N}_r^{rest}$ ,  $\mathcal{N}_r^{in} = \mathcal{N}_v^{in} = \mathcal{L}^{in} = \emptyset$ ;
5   while  $\mathcal{N}_r^{un}$  is nonempty do
6     Choose from  $\mathcal{N}_r^{un}$  the requester with the maximum number of vacant
       candidate users, denoted by  $i^{max} = \arg \max_{i \in \mathcal{N}_r^{un}} |\mathcal{N}_i|$ ;
7     Sort vacant candidate users  $\mathcal{N}_{i^{max}}$  in the decreasing order of channel qualities
       on links from requester  $i^{max}$  to them over channel  $m$ ;
8     Let  $\mathcal{N}_v^{ca} = \mathcal{N}_{i^{max}} \cup \{ES\}$  be the set of vacant candidate nodes by adding ES
       to the end of sorted  $\mathcal{N}_{i^{max}}$ ;
9     while  $k = 1, \dots, |\mathcal{N}_v^{ca}|$  do
10      Choose from  $\mathcal{N}_v^{ca}$  the  $k$ -th element, denoted as vacant node  $j$ ;
11      Check if requester  $i^{max}$  can work on channel  $m$  with existing requesters
         $\mathcal{N}_r^{in}$  simultaneously, for offloading its computing task to vacant node  $j$ ,
        i.e., examine if links  $\mathcal{L}^{in} \cup \{(i^{max}, j)\}$  can operate on channel  $m$ 
        simultaneously to satisfy both SINR and power constraints of requesters;
12      if  $j = ES$  then
13        If  $\mathcal{N}_v^{in}$  contains  $ES$ , there is already a cellular link on channel  $m$  and
        the check fails;
14        If  $F_{ES}^{rest} < f_{i^{max}}$ , the rest computing capacity of ES cannot satisfy the
        computing speed requirement of requester  $i^{max}$  and the check fails;
15      if the check succeeds then
16        Let  $x_{i^{max},j}^m = 1$ ;
17        Update  $\mathcal{N}_r^{in} = \mathcal{N}_r^{in} \cup \{i^{max}\}$ ,  $\mathcal{N}_v^{in} = \mathcal{N}_v^{in} \cup \{j\}$ ,
         $\mathcal{L}^{in} = \mathcal{L}^{in} \cup \{(i^{max}, j)\}$ ;
18        Update vacant candidate user set  $\mathcal{N}_i = \mathcal{N}_i \setminus \{j\}$ ,  $\forall i \in \mathcal{N}_r^{rest}$  if  $j \neq ES$ ;
19        Update the rest computing capacity of ES  $F_{ES}^{rest} = F_{ES}^{rest} - f_{i^{max}}$  if
         $j = ES$ ;
20      Update  $\mathcal{N}_r^{un} = \mathcal{N}_r^{un} \setminus \{i^{max}\}$ ;
21    Update  $\mathcal{N}_r^{rest} = \mathcal{N}_r^{rest} \setminus \mathcal{N}_r^{in}$ ;

```

repeated until \mathcal{N}_r^{un} becomes empty, and then update $\mathcal{N}_r^{rest} = \mathcal{N}_r^{rest} \setminus \mathcal{N}_r^{in}$ and move to the next channel.

In the checking process, the reason of giving priority to the requester with more potential vacant users is because of the fact that if one requester has more potential

vacant users, then it has more chances to select the D2D link with better quality so that less interference is introduced and more requesters can be admitted in the future. In addition, if the potential serving nodes include both vacant mobile users and the ES, the priority is given to vacant users while leaving the ES to be the last for consideration. This is because i) Different from vacant mobile users, the ES can normally provide connections to most of the requesters separately due to the high-location antenna and powerful receiver processing capability of the BS; ii) The ES may need high transmit power from the requester due to the far distance, especially for cell-edge requesters, which may result in strong interference to other links using the same channel; iii) The limited cloud computing ability of the ES constrains the number of requesters that it can serve. In summary, the ES usually has the ability to serve any requester individually but it shouldn't serve too many requesters at the same time. Among multiple potential vacant users, the one with the best channel quality is selected first, since the better channel quality means the lower transmit power and less interference to other users, and thus more other users are possibly admitted into the system. The pseudo code of the greedy algorithm is summarized in Algorithm 7.

The greedy algorithm requires to explore M channels, and for each channel at most N_r requesters need to be checked. For each requester, at most $N_v + 1$ potential vacant nodes are available for serving its computation task and thus at most $N_v + 1$ checks are required for this requester. Therefore, the computational complexity of the greedy algorithm is $O(MN_rN_v)$, which is polynomial with respect to the numbers of users and channels.

Table 5.1: Main simulation parameters.

Parameter	Value
Cell radius	350 m
Path loss exponent	3
Channel bandwidth	5 MHz
Noise spectral density	-174 dBm/Hz
Receiver noise figure	9 dB
Antenna gain	MUs: 0 dBi; BS: 18 dBi
Maximum power	23 dBm
Input bits b_i	5 MBytes
Demand d_i	1 Gigacycles
Deadline T_i	4 s
Speed f_i	randomly over $\{1.0, 1.1, 1.2, 1.3, 1.4, 1.5\}$ GHz
Capacity F_j	MU: randomly over $[1,3]$ GHz
Energy per cycle θ_j	randomly over $(0, 1]$ Joule/Megacycle
Cost per energy c_j^e	ES: 1; MU: randomly over $(0,1]$ cents/Joule
Price per cycle $p_{ES}^{f_i}$	$0.5 \times 10^{-9} f_i$ cents/Megacycle;
Energy threshold	randomly over $[0.5,1]$ Joule

5.5 Simulation Results

In this section, the performance of the proposed scheme is evaluated via simulations. In simulations, users are randomly distributed in a circular cell with a radius of 350 m and the ES is located at the center. The channel gain is set as $G_{i,j}^m = D_{i,j}^{-\nu} |h_{i,j}^m|^2$, where $D_{i,j}$ denotes the distance between transmitter i and receiver j and $\nu = 3$ is the path loss exponent; $|h_{i,j}^m|^2$ captures the small-scale Rayleigh fading effect and is modeled as an exponential random variable with unit mean. Similar simulation settings for D2D communication underlaying cellular networks have been employed in [126], [128], [130]. As to computation tasks, the face recognition application as in [69], [138] is considered, where the computation input data size is $b_i = 5$ MB and the demand for the total amount of CPU cycles is $d_i = 1$ Gigacycles. The latency

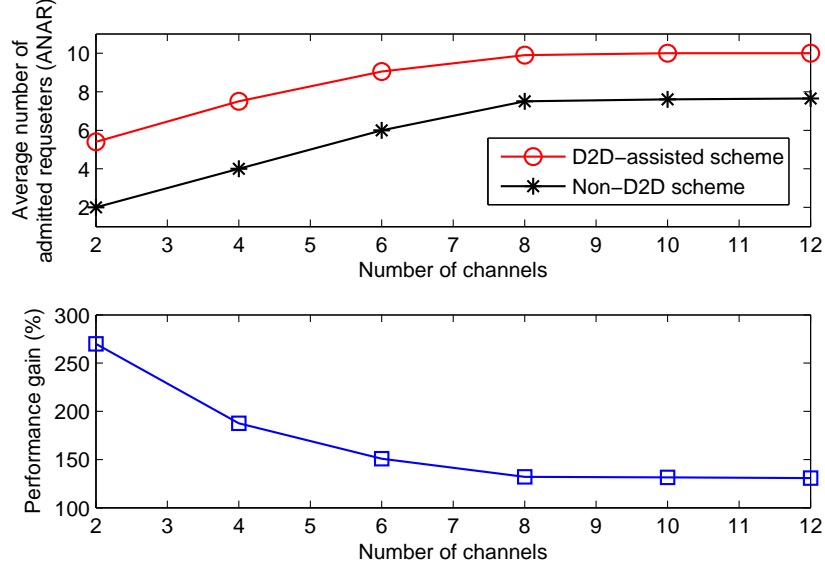


Figure 5.3: Average number of admitted requesters (ANAR) versus the number of channels when $F_{ES} = 10$ GHz, $N_v = N_r = 10$.

requirement is set as $T_i = 4$ s. Other computation model settings mainly follow references [69], [74], [120]. We fix $a_{i,j} = 1$, i.e., the objective is to maximize the number of admitted requesters. Table 5.1 summarizes the main simulation parameters.

Figs. 5.3 and 5.4 compare the performance of two schemes, i.e., mobile edge computing (MEC) with/without D2D assistance. The proposed optimal branch-and-price based algorithm is applied to solve the associated optimization problems in both schemes.

Fig. 5.3 demonstrates the effects of wireless radio resources (the number of channels) on the average number of admitted requesters (ANAR). The performance gain is defined as the ANAR ratio of the D2D-assisted scheme to the non-D2D one. From this figure, we can see that the D2D-assisted scheme always outperforms the non-D2D scheme in terms of ANAR. This is because the D2D-assisted scheme improves

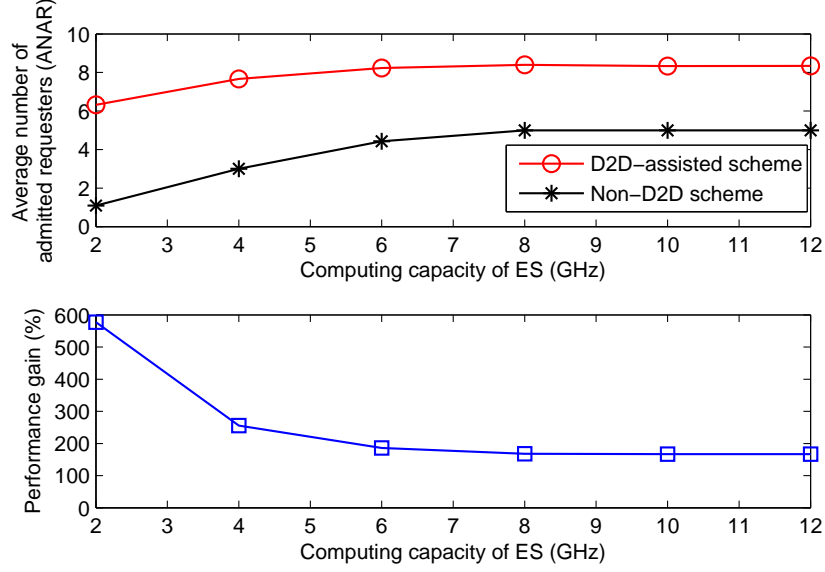


Figure 5.4: Average number of admitted requesters (ANAR) versus the computing capacity of the edge server (ES) when $|\mathcal{M}| = 5$, $N_v = N_r = 10$.

the spectral efficiency via D2D communication and at the same time expands the ES's computing capacity by using nearby vacant mobile devices. The performance gain of the D2D-assisted scheme is remarkable, especially under limited radio resources (small number of channels), and reaches around 270% when there are only 2 channels. The performance gain is decreasing as the number of channels increases, because the non-D2D scheme can be significantly enhanced with the increasing amount of radio resources. The performance gain remains almost the same at around 130% after the number of channels is larger than 8, where both D2D-assisted and non-D2D schemes are no longer limited by the transmission capacity when there are enough radio resources. Even under sufficient radio resources, however, the non-D2D scheme is still worse than the D2D-assisted scheme because the non-D2D scheme suffers from the limited computing capacity of the ES.

Fig. 5.4 illustrates the effects of cloud computing resources (computing capacity) of the ES on ANAR. From this figure, we can observe a significant performance gain with the assistance of D2D and such gain reaches around 580% when the computing capacity is 2 GHz (about 170% at 12 GHz). The performance gain is decreasing with the computing capacity of ES and remains almost invariant when the computing capacity becomes large enough (larger than 6 GHz). Even under the high computing capacity region, the performance gain still exists because the non-D2D scheme is limited by radio resources while such limitation is relieved by D2D communication in the D2D-assisted scheme. In summary, we can conclude that both the radio resource of cellular networks and the cloud computing resource of the ES have great effects on the implementation of mobile edge computing, and the D2D-assisted scheme can significantly enhance the system performance under limited radio and cloud computing resources. Note that Figs. 5.3 and 5.4 focus on the case of homogeneous channel conditions, and similar observation results can be seen for heterogeneous channel conditions, which are omitted here for conciseness.

Figs. 5.5 and 5.6 compare two optimization algorithms: i) the proposed optimal branch-and-priced based algorithm, and ii) the proposed greedy algorithm. Fig. 5.5 illustrates the average number of admitted requesters (ANAR) versus the number of channels under homogeneous channel conditions. From this figure, we can see that the performance of the proposed greedy algorithm gradually approaches the optimal one when the number of channels becomes large. Even when there are only two channels, the greedy algorithm can still achieve about 91% optimal performance. For heterogeneous channel conditions, similar observation results can be seen from Fig. 5.6.

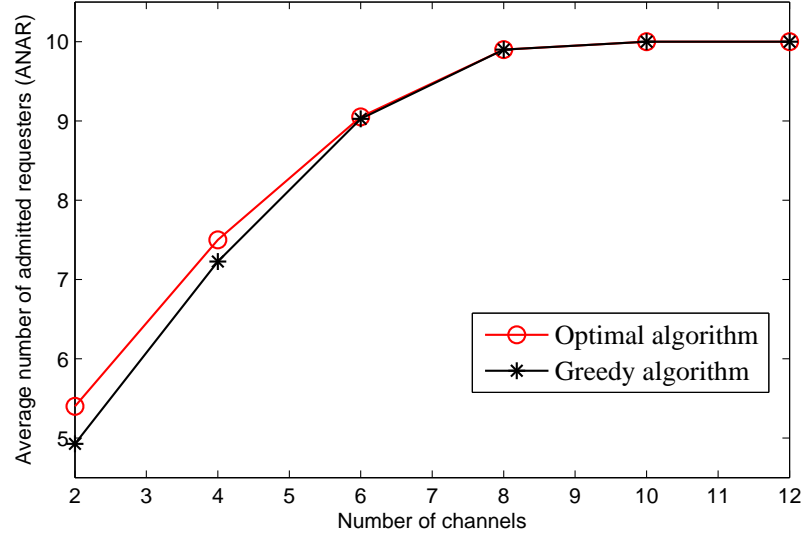


Figure 5.5: Average number of admitted requesters vs. the number of channels under homogeneous channel conditions when $N_v = N_r = 10$.

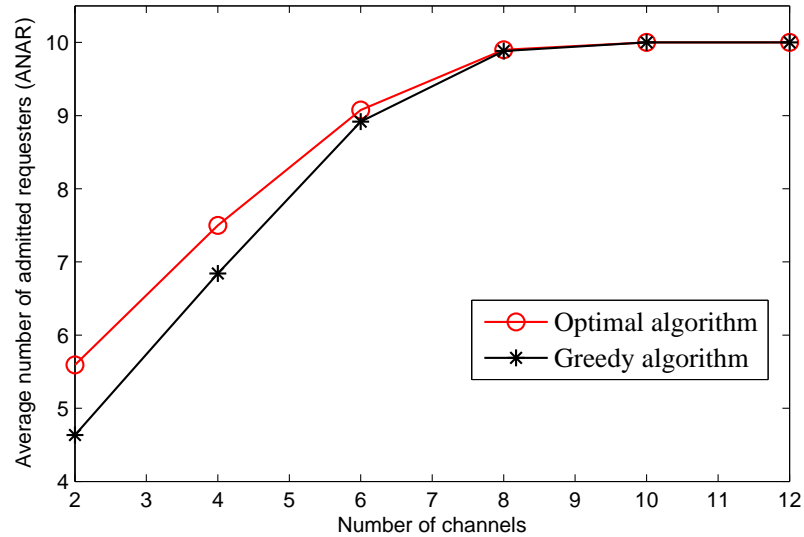


Figure 5.6: Average number of admitted requesters vs. the number of channels under heterogeneous channel conditions when $N_v = N_r = 10$.

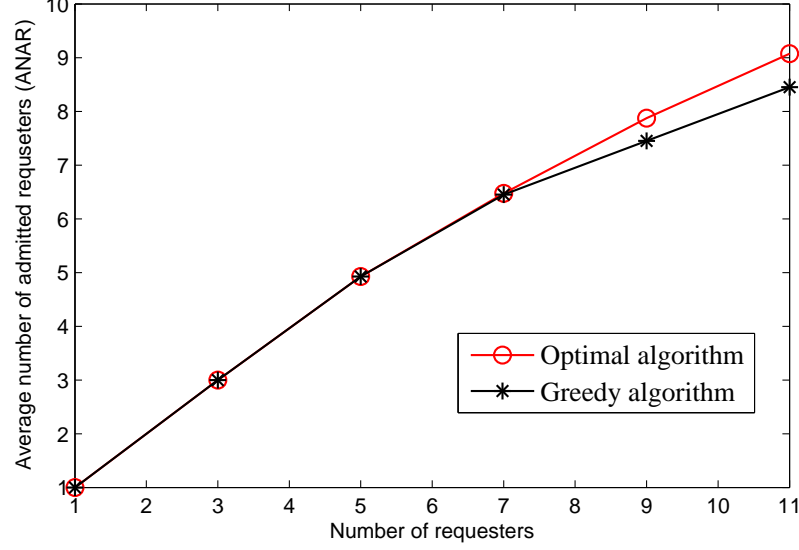


Figure 5.7: Average number of admitted requesters (ANAR) versus the number of requesters (N_r) under heterogeneous channel conditions when $N_v = N_r$ and $M = 5$.

Particularly, the proposed greedy algorithm can obtain around 85% optimality when there are two channels. Compared to the heterogeneous case, the greedy algorithm achieves closer optimality under the homogeneous case. This is because the greedy algorithm is not affected by the exploring order of channels under homogeneous channel conditions.

Fig. 5.7 illustrates the average number of admitted requesters (ANAR) versus the number of total requesters under heterogeneous channel conditions. The number of vacant users is set to be the same as the number of requesters, i.e., $N_v = N_r$, and the number of channels is set as $M = 5$. From this figure, we can see that the ANAR increases with the number of total requesters for both optimal and greedy algorithms, since the computing capacity of MEC systems is expanded by leveraging vacant users and D2D communication. The greedy algorithm can approach the optimal perfor-

Table 5.2: Average computation time of optimal and greedy algorithms.

$N_r = N_v$	1	3	5	7	9	11
Optimal (s)	0.012	0.013	1.044	1.908	18.165	62.908
Greedy (s)	0.001	0.008	0.002	0.003	0.004	0.005

mance when the number of requesters is small, and it can still achieve around 93% optimality when the system overload is 220% (2.2 requesters per channel on average, i.e., 11 requesters on 5 channels).

We further compare the average computation time of both the optimal and the greedy algorithms, as shown in Table 5.2. These results are obtained on a computer with a 3.6 GHz single-core central processing unit (CPU) and a 4 GB random-access memory (RAM). From this table, we can observe that the computation time of the greedy algorithm is always much smaller than the optimal algorithm. In addition, the computation time using the optimal algorithm increases substantially with the number of requesters/vacant users, while the computation time using the greedy algorithm increases only slightly. This clearly demonstrates the computational efficiency of the greedy algorithm. In summary, we can conclude that the proposed greedy algorithm can more effectively and efficiently address the joint admission control and resource management problem in practical D2D-assisted MEC systems. Similar results can be observed for homogeneous channel conditions, which are omitted here for conciseness.

Chapter 6

Conclusions and Future Work

6.1 Conclusions

In this research, key resource allocation and transmission scheduling issues for several wireless communication networks have been investigated. Specifically, transmit power allocation for amplify-and-forward relaying has been studied in Chapter 2, followed by an analysis of link scheduling for enhanced buffer-aided decode-and-forward relaying under time-correlated fading channels in Chapter 3. In Chapter 4, priority-aware interference-avoidance scheduling for dense multi-user coexisting networks is studied, while in Chapter 5, a joint admission control, link scheduling and resource management issue for D2D-assisted mobile edge computing (MEC) systems is discussed. Via optimizing resource allocation and link scheduling, network performance is shown to be greatly improved in terms of power consumption, queueing delay, system throughput, etc.

In Chapter 2, we formulate three power allocation strategies, denoted as Strategies

I, II and III, with different amount of overheads, so that the total transmit power consumption is minimized under a minimum rate constraint or a non-violation probability constraint. In Strategy I, both the source and relay powers are adjusted according to instantaneous CSI, and a closed-form solution is obtained, while in Strategy II, both the source and relay powers are determined based on statistical CSI. In Strategy III, the source power is adjusted based on statistical CSI while the relay power based on instantaneous CSI, which is formulated as a two-stage stochastic programming problem. An N -section method is proposed to solve optimization problems for Strategies II and III. Extensive simulation results show that the proposed algorithms obtain the same numerically optimal solution as the enumeration method. In addition, Strategy I has the lowest total power consumption but the most overheads, while Strategy II has the least overheads but results in a large performance loss (the total power consumption increases by about 8 dB (6.3 times) compared with Strategy I). Compared with Strategy I, Strategy III can cut overheads by half and only leads to a small increase in power consumption when the source-relay channel condition is better than the relay-destination channel.

In Chapter 3, a framework is proposed for analyzing the performance of buffer-aided relaying under time-correlated fading channels. An aggregate Markov chain integrating both buffer state $Q(i)$ and channel state $\beta(i)$ is established and the queueing behavior of packets in the relay's buffer is investigated. Two delay-controllable link scheduling/selection policies, denoted as Policies I and II, are considered with respect to infinite and finite buffer size cases, respectively. For Policy I, an aggregate QBD chain with an infinite length is established and a traditional matrix-geometric method

is used for solving the stationary distribution. For Policy II, the established QBD chain with a finite length is analyzed by a modified matrix-geometric method. Numerical results show that the infinite length buffer results in higher average throughput than the finite one under the same average delay constraint. For both policies, correlated fading can result in the almost same throughput as i.i.d. fading only for cases with loose delay constraints. Otherwise, for stringent delay requirements, correlated fading causes a great reduction in throughput compared with i.i.d. fading especially in low fading margins. Based on these observations, some insights on performance degradation and guidelines on performance improvement under correlated fading are provided.

In Chapter 4, the priority-aware interference-avoidance scheduling for multi-user coexisting wireless networks is addressed, where we have investigated both admission control and throughput maximization of admitted users. For the purpose of practical implementation, a sequential solution framework is proposed, at each step of which a large-scale linear subproblem with a large number of variables is formulated. To solve these large-scale linear subproblems, we introduce the column generation method, which decomposes each subproblem into a restricted master problem (RMP) and a pricing problem (PP). To improve computational efficiency, a greedy initialization algorithm (GIA) is proposed for the RMP to warmly start the column generation process. In addition, under both optimal and approximated solutions to pricing problems, we derive the upper and lower bounds on the optimal objective function of each subproblem. Applying the GIA and bounds to the solution framework comes with an accelerated algorithm, which efficiently reduces the times for optimally solving

pricing problems and thus alleviates computational costs. Simulation results show that the proposed accelerated algorithm combined with approximated solutions to pricing problems, can correctly find the maximum number of high priority levels that the system can accommodate with a very high probability. Also, it can achieve near-optimal throughput performance for admitted users even in the extremely high user density region, in a much faster way.

In Chapter 5, a mobile edge computing (MEC) system with D2D underlaying cellular networks is presented and investigated. The joint optimization problem of admission control and resource management (including link scheduling, channel assignment, and power control) is formulated, where the objective is to maximize the number of admitted computation requests. An optimal branch-and-price based algorithm and a suboptimal greedy algorithm are proposed to solve the optimization problem. Simulation results show that, with the assistance of D2D communication, the MEC system is significantly improved in terms of the average number of admitted requesters, especially under limited radio and computing resources. The benefits come from two aspects: i) via spectrum sharing among D2D and cellular connections, the spectral efficiencies and transmission capacities of cellular networks can be greatly enhanced; and ii) by leveraging the computing abilities of vacant mobile devices, the computing capacity of the edge server (ES) can be largely increased.

6.2 Future Work

Some future research directions on the resource allocation optimization of wireless networks are outlined as follows.

- **Human in the loop:** Recently, many emerging mobile applications are directly linked to human users. Some typical examples include i) traffic monitoring and navigation systems where users can report traffic congestion information timely so that other users can detour traffic jams, ii) electronic book classification systems which can collect feedbacks from readers so as to classify books more accurately, and iii) online language translation systems which may distribute translation tasks to bilingual/multilingual users for more accurate translation. These systems closely depend on not only whether human users are willing to participate, but also the accuracy and reliability of the information provided by human users. Therefore, designing effective incentive mechanisms for encouraging people to join in these systems and providing accurate information is important. Also, it is necessary to develop cheat-proof mechanisms in order to prevent people's strategic behaviors and avoid false information. On the other hand, if too many users join in the systems, we have to address the user selection issue since generally limited resources and network capacities cannot support all users. In these cases, the joint optimization problems of resource allocation, user selection, and mechanism design should be investigated.
- **Dynamic admission control:** Future wireless networks are expected to provide services with lower delay. For example, the 5th generation mobile communication system is expected to have the performance of 1 ms latency for some industrial applications that require fast responses, such as automatic driving etc. For these applications, the quasi-static optimization of admission control and resource allocation can hardly satisfy their low delay requirements, since

user requests may arrive at any time of a scheduling period but not always at the beginning, and postponing the access requests to the next scheduling period is normally infeasible. Hence, developing dynamic admission control and efficient resource management is extremely urgent for potential delay-sensitive applications in future wireless networks.

- Massive MIMO relay networks: the multiple-input-and-multiple-output (MIMO) technology has the great potential for improving network capacities by using multiple transmit and receive antennas. However, it operates on a matrix channel where a transmitter sends multiple data streams by multiple transmit antennas to multiple receive antennas, which is linked to a channel state information (CSI) matrix. In order to optimize the resource allocation of massive MIMO relay networks, a large amount of CSI may require being fed back to the control center. This will cause a lot of overheads. Hence, it necessitates developing some resource allocation schemes with reduced overheads for massive MIMO relay networks while avoiding a large performance loss.

Bibliography

- [1] T. S. Rappaport, S. Sun, *et al.*, “Millimeter wave mobile communications for 5G cellular: It will work!” *IEEE Access*, vol. 1, pp. 335–349, May 2013.
- [2] E. C. V. D. Meulen, “Transmission of information in a T-terminal discrete memoryless channel,” Ph.D. Dissertation, University of California, Berkeley, CA, 1968.
- [3] E. C. V. D. Meulen, “Three-terminal communication channels,” *Advances in Applied Probability*, vol. 3, pp. 120–154, 1971.
- [4] T. Cover and A. E. Gamal, “Capacity theorems for the relay channel,” *IEEE Transactions on Information Theory*, vol. 25, no. 5, pp. 572–584, Sept. 1979.
- [5] A. Sendonaris, E. Erkip, and B. Aazhang, “User cooperation diversity—part I: System description,” *IEEE Transactions on Communications*, vol. 51, no. 11, pp. 1927–1938, Nov. 2003.
- [6] A. Sendonaris, E. Erkip, and B. Aazhang, “User cooperation diversity—part II: Implementation aspects and performance analysis,” *IEEE Transactions on Communications*, vol. 51, no. 11, pp. 1939–1948, Nov. 2003.
- [7] J. N. Laneman and G. W. Wornell, “Distributed space-time-coded protocols for exploiting cooperative diversity in wireless networks,” *IEEE Transactions on Information Theory*, vol. 49, no. 10, pp. 2415–2425, Oct. 2003.
- [8] J. N. Laneman, D. N. C. Tse, and G. W. Wornell, “Cooperative diversity in wireless networks: Efficient protocols and outage behavior,” *IEEE Transactions on Information Theory*, vol. 50, no. 12, pp. 3062–3080, Dec. 2004.
- [9] 3GPP TR 36.814 V9.0.0 (March 2010), “Further advancements for E-UTRA: Physical layer aspects (Release 9),” [Online], Available: <http://www.3gpp.org>.
- [10] Overview of 3GPP Release 10 V0.2.1 (Jun. 2014), [Online], Available: <http://www.3gpp.org>.

- [11] IEEE Std 802.16j-2009, "IEEE Standards for local and metropolitan area networks—Part 16: Air interface for broadband wireless access systems—Amendment 1: Multihop relay specification," Jun. 2009.
- [12] IEEE Std 802.16m-2011, "IEEE Standards for local and metropolitan area networks—Part 16: Air interface for broadband wireless access systems—Amendment 3: Advanced air interface," May 2011.
- [13] S. Huang, H. Chen, J. Cai, and F. Zhao, "Energy efficiency and spectral-efficiency tradeoff in amplify-and-forward relay networks," *IEEE Transactions on Vehicular Technology*, vol. 62, no. 9, pp. 4366–4378, Nov. 2013.
- [14] A. Zafar, R. M. Radaydeh, Y. Chen, and M.-S. Alouini, "Efficient power allocation for fixed gain amplify-and-forward relaying in Rayleigh fading," in *Proc. of IEEE International Conference on Communications Workshops*, Budapest, Hungary, pp. 402–406, Jun. 2013.
- [15] A. Zafar, R. M. Radaydeh, Y. Chen, and M.-S. Alouini, "Power allocation strategies for fixed gain half-duplex amplify-and-forward relaying in Nakagami-m fading," *IEEE Transactions on Wireless Communications*, vol. 13, no. 1, pp. 159–173, Jan. 2014.
- [16] C. Luo, G. Min, F. R. Yu, M. Chen, L. T. Yang, and V. C. M. Leung, "Energy-efficient distributed relay and power control in cognitive radio cooperative communications," *IEEE Journal on Selected Areas in Communications*, vol. 31, no. 11, pp. 2442–2452, Nov. 2013.
- [17] S. Mallick, M. M. Rashid, and V. K. Bhargava, "Joint relay selection and power allocation for decode-and-forward cellular relay network with channel uncertainty," *IEEE Trans. Wireless Commun.*, vol. 11, no. 10, pp. 3496–3508, Oct. 2012.
- [18] K. Vardhe, D. Reynolds, and B. D. Woerner, "Joint power allocation and relay selection for multiuser cooperative communication," *IEEE Transactions on Wireless Communications*, vol. 9, no. 4, pp. 1255–1260, Apr. 2010.
- [19] Y. Li, B. Vuetic, Z. Zhou, and M. Dohler, "Distributed adaptive power allocation for wireless relay networks," *IEEE Transactions on Wireless Communications*, vol. 6, no. 3, pp. 948–958, Mar. 2007.
- [20] Q. Liu, W. Zhang, X. Ma, and G. T. Zhou, "Designing peak power constrained amplify-and-forward relay networks with cooperative diversity," *IEEE Transactions on Wireless Communications*, vol. 11, no. 5, pp. 1733–1743, May 2012.

- [21] I. Maric and R. Yates, "Bandwidth and power allocation for cooperative strategies in Gaussian relay networks," *IEEE Transactions on Information Theory*, vol. 56, no. 4, pp. 1880–1889, Apr. 2010.
- [22] H. A. Tous and I. Barhumi, "Joint power and bandwidth allocation for amplify-and-forward cooperative communications using Stackelberg game," *IEEE Transactions on Vehicular Technology*, vol. 62, no. 4, pp. 1678–1691, May 2013.
- [23] G. Zhao, C. Yang, G. Y. Li, D. Li, and A. C. K. Soong, "Power and channel allocation for cooperative relay in cognitive radio networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 1, pp. 151–159, Feb. 2011.
- [24] Z. Mo, W. Su, S. Batalama, and J. D. Matyjas, "Cooperative communication protocol designs based on optimum power and time allocation," *IEEE Transactions on Wireless Communications*, vol. 13, no. 8, pp. 4283–4296, Aug. 2014.
- [25] M. Duarte, C. Dick, and A. Sabharwal, "Experiment-driven characterization of full-duplex wireless systems," *IEEE Transactions on Wireless Communications*, vol. 11, no. 12, pp. 4296–4307, Dec. 2012.
- [26] D. Bharadia, E. McMillin, and S. Katti, "Full duplex radios," in *Proc. of ACM SIGCOMM 2013 Conference, SIGCOMM'13*, Hong Kong, pp. 375–386, Aug. 2013.
- [27] S. Huang, J. Cai, H. Chen, and H. Zhang, "Transmit power optimization for amplify-and-forward relay networks with reduced overheads," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 7, pp. 5033–5044, Jul. 2016.
- [28] A. Ikhlef, D. S. Michalopoulos, and R. Schober, "Max-max relay selection for relays with buffers," *IEEE Transactions on Wireless Communications*, vol. 11, no. 3, pp. 1124–1135, Mar. 2012.
- [29] A. Ikhlef, D. S. Michalopoulos, and R. Schober, "Buffers improve the performance of relay selection," in *Proc. of IEEE Globecom 2011*, Houston, Texas, USA, pp. 1–6, Dec. 2011.
- [30] I. Krikidis, T. Charalambous, and J. S. Thompson, "Buffer-aided relay selection for cooperative diversity systems without delay constraints," *IEEE Transactions on Wireless Communications*, vol. 11, no. 5, pp. 1957–1967, May 2012.
- [31] M. Zorzi, R. R. Rao, and L. B. Milstein, "On the accuracy of a first-order Markov model for data transmission on fading channels," in *Proc. of IEEE International Conference on Universal Personal Communications 1995 Fourth*, Tokyo, Japan, pp. 211–215, 1995.

- [32] M. Zorzi, R. R. Rao, and L. B. Milstein, "ARQ error control for fading mobile radio channels," *IEEE Transactions on Vehicular Technology*, vol. 46, no. 2, pp. 445–455, May 1997.
- [33] P. Sadeghi, R. A. Kennedy, P. B. Rapajic, and R. Shams, "Finite-state Markov modeling of fading channels—A survey of principles and applications," *IEEE Signal Processing Magazine*, vol. 25, no. 5, pp. 57–80, Sept. 2008.
- [34] F. Hamidi-Sepehr, H. Pfister, and J. Chamberland, "Delay-sensitive communication over fading channels: Queueing behavior and code parameter selection," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 9, pp. 3957–3970, Sept. 2015.
- [35] P. Parag, J. Chamberland, H. Pfister, and K. Narayanan, "Code-rate selection, queueing behavior, and the correlated erasure channel," *IEEE Transactions on Information Theory*, vol. 59, no. 1, pp. 397–407, Jan. 2013.
- [36] N. Zlatanov, A. Ikhlef, T. Islam, and R. Schober, "Buffer-aided cooperative communications: Opportunities and challenges," *IEEE Communications Magazine*, vol. 52, no. 4, pp. 146–153, Apr. 2014.
- [37] S. Rashwand, J. Misic, and H. Khazaei, "Performance analysis of IEEE 802.15.6 under saturation condition and error-prone channel," in *Proc. of 2011 IEEE Wireless Communications and Networking Conference (WCNC)*, Cancun, Mexico, pp. 1167–1172, Mar. 2011.
- [38] M. Deylami and E. Jovanov, "Performance analysis of coexisting IEEE 802.15.4-based health monitoring WBANs," in *Proc. of 2012 Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, San Diego, CA, pp. 2464–2467, Aug. 2012.
- [39] A. H. Moravejosharieh and E. T. Yazdi, "Study of resource utilization in IEEE 802.15.4 Wireless body sensor networks, Part I: the need for enhancement," in *Proc. of 2013 IEEE 16th International Conference on Computational Science and Engineering*, Sydney, NSW, pp. 1226–1231, Dec. 2013.
- [40] C. Xin and M. Song, "An application-oriented spectrum sharing architecture," *IEEE Transactions on Wireless Communications*, vol. 14, no. 5, pp. 2394–2401, May 2015.
- [41] E. A. Jorswieck, L. Badia, T. Fahldieck, E. Karipidis, and J. Luo, "Spectrum sharing improves the network efficiency for cellular operators," *IEEE Communications Magazine*, vol. 52, no. 3, pp. 129–136, Mar. 2014.

- [42] M. Xia and S. Aissa, "Cooperative AF relaying in spectrum-sharing systems: outage probability analysis under co-channel interferences and relay selection," *IEEE Transactions on Communications*, vol. 60, no. 11, pp. 3252–3262, Nov. 2012.
- [43] B. Fateh and M. Govindarasu, "Joint scheduling of tasks and messages for energy minimization in interference-aware real-time sensor networks," *IEEE Transactions on Mobile Computing*, vol. 14, no. 1, pp. 86–98, Jan. 2015.
- [44] Q. Wang, D. O. Wu, and P. Fan, "Delay-constrained optimal link scheduling in wireless sensor networks," *IEEE Transactions on Vehicular Technology*, vol. 59, no. 9, pp. 4564–4577, Nov. 2010.
- [45] P. Phunchongharn, E. Hossain, L. B. Le, and S. Camorlinga, "Robust scheduling and power control for vertical spectrum sharing in STDMA wireless networks," *IEEE Transactions on Wireless Communications*, vol. 11, no. 5, pp. 1850–1860, May 2012.
- [46] P. Phunchongharn and E. Hossain, "Distributed robust scheduling and power control for cognitive spatial-reuse TDMA networks," *IEEE Journal on Selected Areas in Communications*, vol. 30, no. 10, pp. 1934–1946, Nov. 2012.
- [47] J. El-Najjar, C. Assi, and B. Jaumard, "Joint routing and scheduling in WiMAX-based mesh networks," *IEEE Transactions on Wireless Communications*, vol. 9, no. 7, pp. 2371–2381, Jul. 2010.
- [48] V. Gabale, B. Raman, P. Dutta, and S. Kalyanraman, "A classification framework for scheduling algorithms in wireless mesh networks," *IEEE Communications Surveys and Tutorials*, vol. 15, no. 1, pp. 199–222, First Quarter 2013.
- [49] M. Cao, X. Wang, S.-J. Kim, and M. Madhian, "Multi-hop wireless backhaul networks: A cross-layer design paradigm," *IEEE Journal on Selected Areas in Communications*, vol. 25, no. 4, pp. 738–748, May 2007.
- [50] S. H. Cheng and C. Y. Huang, "Coloring-based inter-WBAN scheduling for mobile wireless body area networks," *IEEE Transactions on Parallel and Distributed Systems*, vol. 24, no. 2, pp. 250–259, Feb. 2013.
- [51] S. Huang and J. Cai, "Priority-aware scheduling for coexisting wireless body area networks (Invited Paper)," in *Proc. of 2015 International Conference on Wireless Communications and Signal Processing (WCSP)*, Nanjing, China, pp. 15–17, Oct. 2015.
- [52] IEEE 802.15.6 Task Group, "IEEE standard for local and metropolitan area networks—Part 15.6: wireless body area networks," Feb. 2012.

- [53] S. Manfredi, "Congestion control for differentiated healthcare service delivery in emerging heterogenous wireless body area networks," *IEEE Wireless Communications*, vol. 21, no. 2, pp. 81–90, Apr. 2014.
- [54] S. Misra, S. Moulik, and H. Chao, "A cooperative bargaining solution for priority-based data-rate tuning in a wireless body area network," *IEEE Transactions on Wireless Communications*, vol. 14, no. 5, pp. 2769–2777, May 2015.
- [55] S. Misra and S. Sarkar, "Priority-based time-slot allocation in wireless body area networks during medical emergency situations: an evolutionary game-theoretic perspective," *IEEE Journal of Biomedical and Health Informatics*, vol. 19, no. 2, pp. 541–548, Mar. 2015.
- [56] S. Moulik, S. Misra, C. Chakraborty and M. S. Obaidat, "Prioritized payload tuning mechanism for wireless body area network-based healthcare systems," in *proc. of 2014 IEEE GLOBECOM*, Austin, TX, pp. 2393–2398, 2014.
- [57] S. Moulik, S. Misra and D. Das, "AT-MAC: Adaptive MAC-frame payload tuning for reliable communication in wireless body area networks," *IEEE Transactions on Mobile Computing*, vol. 16, no. 6, pp. 1516–1529, Jun. 2017.
- [58] S. Misra, S. Moulik and H. C. Chao, "A cooperative bargaining solution for priority-based data-rate tuning in a wireless body area network," *IEEE Transactions on Wireless Communications*, vol. 14, no. 5, pp. 2769–2777, May 2015.
- [59] M. S. Chowdhury, K. Ashrafuzzaman, and K. S. Kwak, "Saturation throughput analysis of IEEE 802.15.6 slotted ALOHA in heterogeneous conditions," *IEEE Wireless Communications Letters*, vol. 3, no. 3, pp. 257–260, Mar. 2014.
- [60] J. Zhang, L. W. Hanlen, A. Y. Wang, and X. Huang, "Superframe-level time-hopping system with variable contention access period for wireless body area communications," in *Proc. of 2011 22nd IEEE International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, Montreal, QC, Canada, pp. 2178–2182, Sept. 2011.
- [61] S. Gandham, M. Dawande, and R. Prakash, "Link scheduling in wireless sensor networks: Distributed edge-coloring revisited," *Journal of Parallel and Distributed Computing*, vol. 68, no. 8, pp. 1122–1134, Aug. 2008.
- [62] C. Zhu and M. S. Corson, "A five-phase reservation protocol (FPRP) for mobile ad hoc networks," *Wireless Networks*, vol. 7, no. 4, pp. 371–384, Aug. 2001.
- [63] L. Fu, S. C. Liew, and J. Huang, "Fast algorithms for joint power control and scheduling in wireless networks," *IEEE Transactions on Wireless Communications*, vol. 9, no. 3, pp. 1186–1197, Mar. 2010.

- [64] P. Bjorklund, P. Varbrand, and D. Yuan, "A column generation method for spatial TDMA scheduling in ad hoc networks," *Ad Hoc Networks*, vol. 2, no. 4, pp. 405–418, Oct. 2004.
- [65] M. F. Uddin, H. M. K. Alazemi, and C. Assi, "Optimal flexible spectrum access in wireless networks with software defined radios," *IEEE Transactions on Wireless Communications*, vol. 10, no. 1, pp. 314–324, Jan. 2011.
- [66] S. Kompella, J. E. Wieselthier, A. Ephremides, H. D. Sherali and G. D. Nguyen, "On optimal SINR-based scheduling in multihop wireless networks," *IEEE/ACM Transactions on Networking*, vol. 18, no. 6, pp. 1713–1724, Dec. 2010.
- [67] J. Luo, C. Rosenberg and A. Girard, "Engineering wireless mesh networks: Joint scheduling, routing, power control, and rate adaptation," *IEEE/ACM Transactions on Networking*, vol. 18, no. 5, pp. 1387–1400, Oct. 2010.
- [68] E. Ahmed, A. Gani, M. Sookhak, S. H. Ab Hamid, and F. Xia, "Application optimization in mobile cloud computing: Motivation, taxonomies, and open challenges," *Journal of Network and Computer Applications*, vol. 52, pp. 52–68, 2015.
- [69] X. Chen, L. Jiao, W. Li and X. Fu, "Efficient multi-user computation offloading for mobile-edge cloud computing," *IEEE/ACM Transactions on Networking*, vol. 24, no. 5, pp. 2795–2808, Oct. 2016.
- [70] E. PORTAL, "Mobile-edge computing-introductory technical white paper," 2014.
- [71] M. Beck, M. Werner, S. Feld and S. Schimper, "Mobile edge computing: A taxonomy," in *Proc. of the Sixth International Conference on Advances in Future Internet (AFIN)*, Lisbon, Portugal, pp. 48–54, Nov. 2014.
- [72] Cisco, "Cisco visual networking index: Global mobile data traffic forecast update," Feb. 2017.
- [73] Y. Zhang, D. Niyato and P. Wang, "Offloading in mobile cloudlet systems with intermittent connectivity," *IEEE Transactions on Mobile Computing*, vol. 14, no. 12, pp. 2516–2529, Dec. 2015.
- [74] K. Zhang, Y. Mao, et al., "Energy-efficient offloading for mobile edge computing in 5G heterogeneous networks," *IEEE Access*, vol. 4, pp. 5896–5907, 2016.
- [75] A. Shapiro, D. Dentcheva, and A. Ruszczyński, *Lectures on Stochastic Programming: Modeling and Theory*. Society for Industrial and Applied Mathematics, 2009.

- [76] J. Luo, R. S. Blum, L. Cimini, L. Greenstein, and A. Haimovich, "Power allocation in a transmit diversity system with mean channel gain information," *IEEE Communications Letters*, vol. 9, no. 7, pp. 616–618, Jul. 2005.
- [77] J. Luo, R. S. Blum, L. Cimini, L. Greenstein, and A. Haimovich, "Decode-and-forward cooperative diversity with power allocation in wireless networks," *IEEE Transactions on Wireless Communications*, vol. 6, no. 3, pp. 793–799, Mar. 2007.
- [78] R. Annavajjala, P. C. Cosman, and L. B. Milstein, "Statistical channel knowledge-based optimum power allocation for relaying protocols in high SNR regime," *IEEE Journal on Selected Areas in Communications*, vol. 25, no. 2, pp. 292–305, Feb. 2007.
- [79] M.O. Hasna and M.-S. Alouini, "Optimal power allocation for relayed transmissions over Rayleigh-fading channels," *IEEE Transactions on Wireless Communications*, vol. 3, no. 6, pp. 1999–2004, Nov. 2004.
- [80] Y. Zhao, R. Adve, and T. J. Lim, "Improving amplify-and-forward relay networks: Optimal power allocation versus selection," *IEEE Transactions on Wireless Communications*, vol. 6, no. 8, pp. 3114–3123, Aug. 2007.
- [81] F. Gao, T. Cui, and A. Nallanathan, "On channel estimation and optimal training design for amplify and forward relay networks," *IEEE Transactions on Wireless Communications*, vol. 7, no. 5, pp. 1907–1916, May 2008.
- [82] C. Y. Lee and G. U. Hwang, "Fair and minimal power allocation in a two-hop relay network for QoS support," *IEEE Transactions on Wireless Communications*, vol. 10, no. 11, pp. 3864–3873, Nov. 2011.
- [83] M. Chen, S. Serbetli, and A. Yener, "Distributed power allocation strategies for parallel relay networks," *IEEE Transactions on Wireless Communications*, vol. 7, no. 2, pp. 552–561, Feb. 2008.
- [84] T. T. Pham, H. H. Nguyen, and H. D. Tuan, "Power allocation in orthogonal wireless relay networks with partial channel state information," *IEEE Transactions on Signal Processing*, vol. 58, no. 2, pp. 869–878, Feb. 2010.
- [85] A. Zafar, R. M. Radaydeh, Y. Chen, and M.-S. Alouini, "Enhancing the efficiency of constrained dual-hop variable-gain AF relaying under Nakagami- m fading," *IEEE Transactions on Signal Processing*, vol. 62, no. 14, pp. 3616–3630, Jul. 2014.
- [86] M. A. Gatzianas, L. G. Georgiadis, and G. K. Karagiannidis, "Gain adaptation policies for dual-hop nonregenerative relayed systems," *IEEE Transactions on Communications*, vol. 55, no. 8, pp. 1472–1477, Aug. 2007.

- [87] L. Jimenez Rodriguez, N. H. Tran, A. Helmy, and T. Le-Ngoc, "Optimal power adaptation for cooperative AF relaying with channel side information," *IEEE Transactions on Vehicular Technology*, vol. 62, no. 7, pp. 3164–3174, Sept. 2013.
- [88] M. Chraiti, W. Ajib, and J.-F. Frigon, "Optimal long-term power adaption for cooperative DF relaying," *IEEE Wireless Communications Letters*, vol. 3, no. 2, pp. 201–204, Apr. 2014.
- [89] N. Zlatanov, R. Schober, and P. Popovski, "Buffer-aided relaying with adaptive link selection," *IEEE Journal on Selected Areas in Communications*, vol. 31, no. 8, pp. 1530–1542, Aug. 2013.
- [90] N. Zlatanov and R. Schober, "Buffer-aided relaying with adaptive link selection—fixed and mixed rate transmission," *IEEE Transactions on Information Theory*, vol. 59, no. 5, pp. 2816–2840, May 2013.
- [91] G. Bolch, S. Greiner, H. D. Meer, and K. S. Trivedi, *Queueing Networks and Markov Chains*, 2nd edition, Hoboken, New Jersey, USA: John Wiley and Sons, 2006.
- [92] A. S. Alfa, *Queueing Theory for Telecommunications*, New York, USA: Springer Science+Business Media, 2010.
- [93] N. Akar and K. Sohraby, "Finite and infinite QBD chains: A simple and unifying algorithmic approach," in *proc. of IEEE Infocom 1997*, Kobe, Japan, pp. 1105–1113, 1997.
- [94] V. Jamali, N. Zlatanov, A. Ikhlef and R. Schober, "Achievable rate region of the bidirectional buffer-aided relay channel with block fading," *IEEE Transactions on Information Theory*, vol. 60, no. 11., pp. 7090–7111, Nov. 2014.
- [95] V. Jamali, N. Zlatanov, and R. Schober, "Bidirectional buffer-aided relay networks with fixed rate transmission—Part I: Delay-unconstrained case," *IEEE Transactions on Wireless Communications*, vol. 14, no. 3, pp. 1323–1338, Mar. 2015.
- [96] V. Jamali, N. Zlatanov, and R. Schober, "Bidirectional buffer-aided relay networks with fixed rate transmission—Part II: Delay-constrained case," *IEEE Transactions on Wireless Communications*, vol. 14, no. 3, pp. 1339–1355, Mar. 2015.
- [97] A. Zafar, M. Shaqfeh, M. S. Alouini, and H. Alnuweiri, "Resource allocation for two source-destination pairs sharing a single relay with a buffer," *IEEE Transactions on Communications*, vol. 62, no. 5, pp. 1444–1457, May 2014.

- [98] A. Zafar, M. Shaqfeh, M. S. Alouini, and H. Alnuweiri, "Exploiting multi-user diversity and multi-hop diversity in dual-hop broadcast channels," *IEEE Transactions on Wireless Communications*, vol. 12, no. 7, pp. 3314–3325, Jul. 2013.
- [99] T. Islam, A. Ikhlef, R. Schober, and V. Bhargava, "Multisource buffer-aided relay networks: Adaptive rate transmission," in *proc. of IEEE Globecom 2013*, Atlanta, GA, USA, pp. 3577–3582, Dec. 2013.
- [100] G. Chen, Z. Tian, Y. Gong, and J. Chambers, "Decode-and-forward buffer-aided relay selection in cognitive relay networks," *IEEE Transactions on Vehicular Technology*, vol. 63, no. 9, pp. 4723–4728, Nov. 2014.
- [101] A. Zafar, M. Shaqfeh, M. S. Alouini, and H. Alnuweiri, "Cooperative overlay cognitive radio with opportunistic link selection," in *proc. of the tenth International Symposium on Wireless Communication Systems 2013, ISWCS 2013*, Ilmenau, Germany, pp. 1–5, Aug. 2013.
- [102] S. Huang, J. Cai, and H. Zhang, "Relay selection for average throughput maximization in buffer-aided relay networks," in *proc. of 2015 IEEE International Conference on Communications (ICC)*, London, UK, pp. 1994–1998, Jun. 2015.
- [103] G. Latouche and V. Ramaswami, "A logarithmic reduction algorithm for quasi-birth-death processes," *Journal of Applied Probability*, vol. 30, no. 3, pp. 650–674, Sep. 1993.
- [104] D. Bini and B. Meini, "On the solution of a nonlinear matrix equation arising in queueing problems," *SIAM Journal on Matrix Analysis and Applications*, vol. 17, no. 4, pp. 906–926, 1996.
- [105] D. Bini and B. Meini, "Improved cyclic reduction for solving queueing problems," *Numerical Algorithms*, vol. 15, no. 1, pp. 57–74, 1997.
- [106] B. Hajek, "Birth-and-death processes on the integers with phases and general boundaries," *Journal of Applied Probability*, vol. 19, no. 3, pp. 488–499, Sept. 1982.
- [107] V. D. N. Persone and V. Grassi, "Solution of finite QBD processes," *Journal of Applied Probability*, vol. 33, no. 4, pp. 1003–1010, Dec. 1996.
- [108] L. Gun and A. M. Makowski, "Matrix-geometric solution for finite capacity queues with phase-type distributions," *proc. of Performance 87*, Brussels, Belgium, Dec. 1987.
- [109] E. H. Elhafsi and M. Molle, "On the solution to QBD processes with finite state space," *Stochastic Analysis and Applications*, vol. 25, no. 4, pp. 763–779, 2007.

- [110] C. Yi and J. Cai, "A priority-aware truthful mechanism for supporting multi-class delay-sensitive medical packet transmissions in E-health networks," *IEEE Transactions on Mobile Computing*, vol. 16, no. 9, pp. 2422–2435, Sept. 1 2017.
- [111] C. Yi, A. S. Alfa and J. Cai, "An incentive-compatible mechanism for transmission scheduling of delay-sensitive medical packets in E-health networks," *IEEE Transactions on Mobile Computing*, vol. 15, no. 10, pp. 2424–2436, Oct. 2016.
- [112] M. E. Lubbecke, "Column generation," *EORMS*, Jul. 2010.
- [113] A. Kakoa, T. Onoa, T. Hirata, and M. M. Halldórsson, "Approximation algorithms for the weighted independent set problem in sparse graphs," *Discrete Applied Mathematics*, no. 157, pp. 617–626, 2009.
- [114] N. Zadeh, "What is the worst case behavior of the simplex algorithm?" *Stanford University*, Stanford, California, USA, 1980.
- [115] K. Y. Yazdandoost and K. Sayrafian-Pour, "Channel model for body area network (BAN)," *IEEE 802.15.6 Technical Contribution*, document ID: 15-08-0780-09-0006, pp. 1–61, Apr. 2009.
- [116] L. Yang, J. Cao, H. Cheng and Y. Ji, "Multi-user computation partitioning for latency sensitive mobile cloud applications," *IEEE Transactions on Computers*, vol. 64, no. 8, pp. 2253–2266, Aug. 2015.
- [117] S. Kosta, A. Aucinas, P. Hui, R. Mortier and X. Zhang, "ThinkAir: Dynamic resource allocation and parallel execution in the cloud for mobile code offloading," in *proc. of IEEE INFOCOM 2012*, Orlando, FL, pp. 945–953, 2012.
- [118] H. Liang, L. X. Cai, D. Huang, X. Shen and D. Peng, "An SMDP-based service model for interdomain resource allocation in mobile cloud networks," *IEEE Transactions on Vehicular Technology*, vol. 61, no. 5, pp. 2222–2232, Jun 2012.
- [119] S. Sardellitti, G. Scutari and S. Barbarossa, "Joint optimization of radio and computational resources for multicell mobile-edge computing," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 1, no. 2, pp. 89–103, Jun. 2015.
- [120] O. Munoz, A. Pascual-Iserte and J. Vidal, "Optimization of radio and computational resources for energy efficiency in latency-constrained application offloading," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 10, pp. 4738–4755, Oct. 2015.
- [121] R. Kaewpuang, D. Niyato, P. Wang and E. Hossain, "A framework for cooperative resource management in mobile cloud computing," *IEEE Journal on Selected Areas in Communications*, vol. 31, no. 12, pp. 2685–2700, Dec. 2013.

- [122] Y. Liu, M. J. Lee and Y. Zheng, "Adaptive multi-resource allocation for cloudlet-based mobile cloud computing system," *IEEE Transactions on Mobile Computing*, vol. 15, no. 10, pp. 2398–2410, Oct. 2016.
- [123] J. Almeida, V. Almeida; D. Ardagna, I. Cunha, C. Francalanci, M. Trubian, "Joint admission control and resource allocation in virtualized servers," *Journal of Parallel and Distributed Computing*, vol. 70, no. 4, pp.344–362, 2010.
- [124] C. Xu, L. Xiong, et al., "Efficiency resource allocation for device-to-device underlay communication systems: A reverse iterative combinatorial auction based approach," *IEEE Journal on Selected Areas in Communications*, vol. 31, no. 9, pp. 348–358, Sep. 2013.
- [125] J. Gu, S. J. Bae, S. F. Hasan, and M. Y. Chung, "Heuristic algorithm for proportional fair scheduling in D2D-cellular systems," *IEEE Transactions on Wireless Communications*, vol. 15, no. 1, pp. 769–780, Jan. 2016.
- [126] F. Wang, C. Xu, L. Song and Z. Han, "Energy-efficient resource allocation for device-to-device underlay communication," *IEEE Transactions on Wireless Communications*, vol. 14, no. 4, pp. 2082–2092, Apr. 2015.
- [127] L. Zhang, M. Xiao, G. Wu, and S. Li, "Efficient scheduling and power allocation for D2D-assisted wireless caching networks," *IEEE Transactions on Communications*, vol. 64, no. 6, pp. 2438–2452, Jun. 2016.
- [128] W. Cheng, X. Zhang and H. Zhang, "Optimal power allocation with statistical QoS provisioning for D2D and cellular communications over underlaying wireless networks," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 1, pp. 151–162, Jan. 2016.
- [129] R. Yin, G. Yu, et al., "Pricing-based interference coordination for D2D communications in cellular networks," *IEEE Transactions on Wireless Communications*, vol. 14, no. 3, pp. 1519–1532, Mar. 2015.
- [130] C. Gao, J. Tang, et al, "Enabling green wireless networking with device-to-device links: A joint optimization approach," *IEEE Transactions on Wireless Communications*, vol. 15, no. 4, pp. 2770–2779, Apr. 2016.
- [131] C. Vallati, A. Virdis, E. Mingozzi and G. Stea, "Mobile-edge computing come home connecting things in future smart homes using LTE device-to-device communications," *IEEE Consumer Electronics Magazine*, vol. 5, no. 4, pp. 77-83, Oct. 2016.
- [132] M. Jo, T. Maksymyuk, B. Strykhalyuk and C. H. Cho, "Device-to-device-based heterogeneous radio access network architecture for mobile cloud computing," *IEEE Wireless Communications*, vol. 22, no. 3, pp. 50-58, Jun. 2015.

- [133] Y. Li, L. Sun and W. Wang, "Exploring device-to-device communication for mobile cloud computing," in *Proc. of 2014 IEEE International Conference on Communications (ICC)*, Sydney, NSW, 2014, pp. 2239-2244.
- [134] S. Sanghavi, D. Shah and A. S. Willsky, "Message passing for maximum weight independent set," *IEEE Transactions on Information Theory*, vol. 55, no. 11, pp. 4822-4834, Nov. 2009.
- [135] F. Vanderbeck, "Branching in branch-and-price: A generic scheme," *Mathematics Program*, vol. 130, pp. 249-294, 2011,
- [136] F. Vanderbeck, "On Dantzig-Wolfe decomposition in integer programming and ways to perform branching in a branch-and-price algorithm," *Operations Research*, vol. 48, no. 1, pp. 111-128, 2000.
- [137] C. Barnhart, E. L. Johnson, G. L. Nemhauser, M. W. P. Savelsbergh, and P. H. Vance, "Branch-and-price: Column generation for solving huge integer programs," *Operations Research*, vol. 46, no. 3, pp. 316-329, May-Jun. 1998.
- [138] T. Soyata, R. Muraleedharan, et al., "Cloud-vision: Real-time face recognition using a mobile-cloudlet-cloud acceleration architecture," in *proc. of 2012 IEEE Symposium on Computers and Communications (ISCC)*, Cappadocia, Turkey, pp. 59-66, Jul. 2012.

Appendix A

Appendixes of Chapter 3

A.1 Transition Sub-matrices under Policy I

The transition sub-matrices of the aggregate chain $Y(i)$ are derived as

$$\begin{aligned}
 \mathbf{B}_0 &= \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \lambda_{20}^0 & \lambda_{21}^0 & \lambda_{22}^0 & \lambda_{23}^0 \\ \lambda_{30}^0 & \lambda_{31}^0 & \lambda_{32}^0 & \lambda_{33}^0 \end{bmatrix}, \\
 \mathbf{B}_1 &= \begin{bmatrix} \kappa_{00}^0 & \kappa_{01}^0 & \kappa_{02}^0 & \kappa_{03}^0 \\ \kappa_{10}^0 & \kappa_{11}^0 & \kappa_{12}^0 & \kappa_{13}^0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \\
 \mathbf{A}_0 &= \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \lambda_{20} & \lambda_{21} & \lambda_{22} & \lambda_{23} \\ \lambda_{30} & \lambda_{31} & \lambda_{32} & \lambda_{33} \end{bmatrix}, \\
 \mathbf{A}_1 &= \begin{bmatrix} \kappa_{00} & \kappa_{01} & \kappa_{02} & \kappa_{03} \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \\
 \mathbf{A}_2 &= \begin{bmatrix} 0 & 0 & 0 & 0 \\ \mu_{10} & \mu_{11} & \mu_{12} & \mu_{13} \\ 0 & 0 & 0 & 0 \\ \mu_{30} & \mu_{31} & \mu_{32} & \mu_{33} \end{bmatrix}, \tag{A.1}
 \end{aligned}$$

where we introduce the following convenient notations to represent the state transition probabilities for level $q > 0$.

$$\lambda_{ab} = \Pr\{Y(i+1) = (b, q+1) | Y(i) = (a, q)\}, \tag{A.2}$$

$$\kappa_{ab} = \Pr\{Y(i+1) = (b, q) | Y(i) = (a, q)\}, \quad (\text{A.3})$$

$$\mu_{ab} = \Pr\{Y(i+1) = (b, q-1) | Y(i) = (a, q)\}. \quad (\text{A.4})$$

In the above equations, $a, b \in \{0, 1, 2, 3\}$ and $q \in \{1, 2, \dots\}$. λ_{ab} is the forward transition probability from level q to level $q+1$, κ_{ab} is the transition probability within level q , and μ_{ab} is the backward transition probability from level q to $q-1$. Similarly, for level $q=0$ (i.e., the buffer is empty), we use the following convenient notations as

$$\lambda_{ab}^0 = \Pr\{Y(i+1) = (b, 1) | Y(i) = (a, 0)\}, \quad (\text{A.5})$$

$$\kappa_{ab}^0 = \Pr\{Y(i+1) = (b, 0) | Y(i) = (a, 0)\}. \quad (\text{A.6})$$

Based on Policy I and the corresponding arrival/departure processes described in Section 3.1, the transition probabilities are derived as

$$\begin{aligned} \lambda_{0b} &= \lambda_{1b} = 0, \quad \forall b, \\ \lambda_{20} &= (1-p_0)q_1, \quad \lambda_{21} = (1-p_0)(1-q_1), \\ \lambda_{22} &= p_0q_1, \quad \lambda_{23} = p_0(1-q_1), \\ \lambda_{30} &= (1-p_0)(1-p_1)(1-P_C), \quad \lambda_{31} = (1-p_0)p_1(1-P_C), \\ \lambda_{32} &= p_0(1-p_1)(1-P_C), \quad \lambda_{33} = p_0p_1(1-P_C), \\ \kappa_{00} &= q_0q_1, \quad \kappa_{01} = q_0(1-q_1), \\ \kappa_{02} &= (1-q_0)q_1, \quad \kappa_{03} = (1-q_0)(1-q_1), \\ \kappa_{1b} &= \kappa_{2b} = \kappa_{3b} = 0, \quad \forall b, \\ \mu_{0b} &= \mu_{2b} = 0, \quad \forall b, \\ \mu_{10} &= q_0(1-p_1), \quad \mu_{11} = q_0p_1, \\ \mu_{12} &= (1-q_0)(1-p_1), \quad \mu_{13} = (1-q_0)p_1, \\ \mu_{30} &= (1-p_0)(1-p_1)P_C, \quad \mu_{31} = (1-p_0)p_1P_C, \\ \mu_{32} &= p_0(1-p_1)P_C, \quad \mu_{33} = p_0p_1P_C, \\ \lambda_{0b}^0 &= \lambda_{1b}^0 = 0, \quad \lambda_{2b}^0 = \lambda_{2b}, \quad \lambda_{3b}^0 = \lambda_{3b} + \mu_{3b}, \quad \forall b, \\ \kappa_{0b}^0 &= \kappa_{0b}, \quad \kappa_{1b}^0 = \mu_{1b}, \quad \kappa_{2b}^0 = \kappa_{3b}^0 = 0, \quad \forall b, \end{aligned}$$

A.2 Transition Sub-matrices under Policy II

Under Policy II, the transition sub-matrices \mathbf{C}_1 and \mathbf{C}_2 are derived as

$$\mathbf{C}_2 = \begin{bmatrix} 0 & 0 & 0 & 0 \\ \mu_{10}^L & \mu_{11}^L & \mu_{12}^L & \mu_{13}^L \\ 0 & 0 & 0 & 0 \\ \mu_{30}^L & \mu_{31}^L & \mu_{32}^L & \mu_{33}^L \end{bmatrix},$$

$$\mathbf{C}_1 = \begin{bmatrix} \kappa_{00}^L & \kappa_{01}^L & \kappa_{02}^L & \kappa_{03}^L \\ 0 & 0 & 0 & 0 \\ \kappa_{20}^L & \kappa_{21}^L & \kappa_{22}^L & \kappa_{23}^L \\ 0 & 0 & 0 & 0 \end{bmatrix}, \quad (\text{A.7})$$

where κ_{ab}^L and μ_{ab}^L ($a, b \in \{0, 1, 2, 3\}$) represent the state transition probabilities associated with level L (i.e., the buffer is full) as

$$\kappa_{ab}^L = \Pr\{Y(i+1) = (b, L) | Y(i) = (a, L)\}, \quad (\text{A.8})$$

$$\mu_{ab}^L = \Pr\{Y(i+1) = (b, L-1) | Y(i) = (a, L)\}. \quad (\text{A.9})$$

According to Policy II and the corresponding arrival/departure processes described in Section 3.1, the values of κ_{ab}^L and μ_{ab}^L are derived as

$$\begin{aligned} \mu_{0b}^L &= \mu_{2b}^L = 0, \quad \mu_{1b}^L = \mu_{1b}, \quad \mu_{3b}^L = \mu_{3b} + \lambda_{3b}, \quad \forall b, \\ \kappa_{0b}^L &= \kappa_{0b}, \quad \kappa_{2b}^L = \lambda_{2b}, \quad \kappa_{1b}^L = \kappa_{3b}^L = 0, \quad \forall b. \end{aligned}$$

A.3 Proofs of (3.18) and (3.19)

The equation (3.18) can be directly derived from (3.14) via some simple matrix operations (note that (3.14b) can be written as $(\boldsymbol{\pi}_0, \boldsymbol{\pi}_1)\mathbf{U} = \mathbf{1}^T$). In (3.19), the relation between $\boldsymbol{\pi}_0$ and $\boldsymbol{\pi}_1$ ($\boldsymbol{\pi}_0 = \boldsymbol{\pi}_1 \mathbf{A}_2 (\mathbf{I} - \mathbf{B}_1)^{-1}$) can be obtained from (3.14a) and the expression of $\boldsymbol{\pi}_1$ can be derived from the stationary distribution of channel fading, i.e.,

$$\begin{aligned} \boldsymbol{\beta} &= \sum_{q=0}^{\infty} \boldsymbol{\pi}_q = \boldsymbol{\pi}_0 + \sum_{q=1}^{\infty} \boldsymbol{\pi}_q \\ &= \boldsymbol{\pi}_1 \mathbf{A}_2 (\mathbf{I} - \mathbf{B}_1)^{-1} + \boldsymbol{\pi}_1 (\mathbf{I} - \mathbf{R})^{-1}, \end{aligned} \quad (\text{A.10})$$

where $\sum_{q=1}^{\infty} \boldsymbol{\pi}_q = \boldsymbol{\pi}_1 (\mathbf{I} - \mathbf{R})^{-1}$ is from (3.17).

A.4 Proof of Theorem 3.1

Theorem 1 can be proven by verifying that the solution of the stationary distribution given by (3.30)–(3.33) satisfies both the balance equation $\boldsymbol{\pi} \mathbf{T} = \boldsymbol{\pi}$ and the normalized condition $\boldsymbol{\pi} \mathbf{1} = 1$ (note that $\boldsymbol{\pi} = \{\boldsymbol{\pi}_0, \boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_L\}$). We first prove that the balance equation can be satisfied. For level 0, we have

$$\boldsymbol{\pi}_0 \mathbf{B}_1 + \boldsymbol{\pi}_1 \mathbf{A}_2 = \boldsymbol{\pi}_0 \mathbf{B}_1 + (\mathbf{v}_1 + \mathbf{v}_2 \mathbf{R}_2^{L-2}) \mathbf{A}_2 = \boldsymbol{\pi}_0, \quad (\text{A.11})$$

where the first equality is based on (3.32) with $q = 1$ (i.e., $\pi_1 = \mathbf{v}_1 + \mathbf{v}_2 \mathbf{R}_2^{L-2}$) and the second equality is from (3.33a) (i.e., $\pi_0 \mathbf{B}_1 + \mathbf{v}_1 \mathbf{A}_2 + \mathbf{v}_2 \mathbf{R}_2^{L-2} \mathbf{A}_2 = \pi_0$). For level 1, we have

$$\begin{aligned}
& \pi_0 \mathbf{B}_0 + \pi_1 \mathbf{A}_1 + \pi_2 \mathbf{A}_2 \\
&= \pi_0 \mathbf{B}_0 + (\mathbf{v}_1 + \mathbf{v}_2 \mathbf{R}_2^{L-2}) \mathbf{A}_1 + (\mathbf{v}_1 \mathbf{R}_1 + \mathbf{v}_2 \mathbf{R}_2^{L-3}) \mathbf{A}_2 \\
&= \pi_0 \mathbf{B}_0 + \mathbf{v}_1 (\mathbf{A}_1 + \mathbf{R}_1 \mathbf{A}_2) + \mathbf{v}_2 (\mathbf{R}_2^{L-2} \mathbf{A}_1 + \mathbf{R}_2^{L-3} \mathbf{A}_2) \\
&= \pi_0 \mathbf{B}_0 + \mathbf{v}_1 (\mathbf{A}_1 + \mathbf{R}_1 \mathbf{A}_2) \\
&\quad + \mathbf{v}_2 (\mathbf{R}_2^{L-2} \mathbf{A}_1 + \mathbf{R}_2^{L-3} \mathbf{A}_2 - \mathbf{R}_2^{L-2}) + \mathbf{v}_2 \mathbf{R}_2^{L-2} \\
&= \mathbf{v}_1 + \mathbf{v}_2 \mathbf{R}_2^{L-2} = \pi_1,
\end{aligned} \tag{A.12}$$

where the first equality is based on (3.32) with $q = 1$ and $q = 2$, the fourth equality is from (3.33a) (i.e., $\pi_0 \mathbf{B}_0 + \mathbf{v}_1 (\mathbf{A}_1 + \mathbf{R}_1 \mathbf{A}_2) + \mathbf{v}_2 \mathbf{R}_2^{L-3} (\mathbf{R}_2 \mathbf{A}_1 + \mathbf{A}_2 - \mathbf{R}_2) = \mathbf{v}_1$), and the last equality is derived from (3.32) with $q = 1$. Similarly, for level $L - 1$ and level L , we have

$$\begin{aligned}
& \pi_{L-2} \mathbf{A}_0 + \pi_{L-1} \mathbf{A}_1 + \pi_L \mathbf{C}_2 \\
&= (\mathbf{v}_1 \mathbf{R}_1^{L-3} + \mathbf{v}_2 \mathbf{R}_2) \mathbf{A}_0 + (\mathbf{v}_1 \mathbf{R}_1^{L-2} + \mathbf{v}_2) \mathbf{A}_1 + \pi_L \mathbf{C}_2 \\
&= \mathbf{v}_1 (\mathbf{R}_1^{L-3} \mathbf{A}_0 + \mathbf{R}_1^{L-2} \mathbf{A}_1) + \mathbf{v}_2 (\mathbf{R}_2 \mathbf{A}_0 + \mathbf{A}_1) + \pi_L \mathbf{C}_2 \\
&= \mathbf{v}_1 \mathbf{R}_1^{L-2} + \mathbf{v}_1 (\mathbf{R}_1^{L-3} \mathbf{A}_0 + \mathbf{R}_1^{L-2} \mathbf{A}_1 - \mathbf{R}_1^{L-2}) \\
&\quad + \mathbf{v}_2 (\mathbf{R}_2 \mathbf{A}_0 + \mathbf{A}_1) + \pi_L \mathbf{C}_2 \\
&= \mathbf{v}_1 \mathbf{R}_1^{L-2} + \mathbf{v}_2 = \pi_{L-1},
\end{aligned} \tag{A.13}$$

$$\pi_{L-1} \mathbf{A}_0 + \pi_L \mathbf{C}_1 = (\mathbf{v}_1 \mathbf{R}_1^{L-2} + \mathbf{v}_2) \mathbf{A}_0 + \pi_L \mathbf{C}_1 = \pi_L. \tag{A.14}$$

For level q ($2 \leq q \leq L - 2$), we have

$$\begin{aligned}
& \pi_{q-1} \mathbf{A}_0 + \pi_q \mathbf{A}_1 + \pi_{q+1} \mathbf{A}_2 \\
&= (\mathbf{v}_1 \mathbf{R}_1^{q-2} + \mathbf{v}_2 \mathbf{R}_2^{L-q}) \mathbf{A}_0 + (\mathbf{v}_1 \mathbf{R}_1^{q-1} + \mathbf{v}_2 \mathbf{R}_2^{L-q-1}) \mathbf{A}_1 \\
&\quad + (\mathbf{v}_1 \mathbf{R}_1^q + \mathbf{v}_2 \mathbf{R}_2^{L-q-2}) \mathbf{A}_2 \\
&= \mathbf{v}_1 \mathbf{R}_1^{q-2} (\mathbf{A}_0 + \mathbf{R}_1 \mathbf{A}_1 + \mathbf{R}_1^2 \mathbf{A}_2) \\
&\quad + \mathbf{v}_2 \mathbf{R}_2^{L-q-2} (\mathbf{R}_2^2 \mathbf{A}_0 + \mathbf{R}_2 \mathbf{A}_1 + \mathbf{A}_2) \\
&= \mathbf{v}_1 \mathbf{R}_1^{q-1} + \mathbf{v}_2 \mathbf{R}_2^{L-q-1} = \pi_q,
\end{aligned} \tag{A.15}$$

where the third equality is based on (3.30) and (3.31). Now, we have already verified the stationary distribution of all levels ($\pi_0, \pi_1, \dots, \pi_L$). In other words, the balance equation is satisfied when the stationary distribution is given by (3.30)–(3.33).

Secondly, we prove that the normalized condition can be satisfied as follows

$$\pi \mathbf{1} = \left(\pi_0 + \sum_{q=1}^{L-1} \pi_q + \pi_L \right) \mathbf{1}$$

$$\begin{aligned}
&= \left(\pi_0 + \sum_{q=1}^{L-1} (\mathbf{v}_1 \mathbf{R}_1^{q-1} + \mathbf{v}_2 \mathbf{R}_2^{L-q-1}) + \pi_L \right) \mathbf{1} \\
&= \left(\pi_0 + \mathbf{v}_1 \sum_{q=1}^{L-1} \mathbf{R}_1^{q-1} + \mathbf{v}_2 \sum_{q=1}^{L-1} \mathbf{R}_2^{L-q-1} + \pi_L \right) \mathbf{1} \\
&= \left(\pi_0 + \mathbf{v}_1 \sum_{m=0}^{L-2} \mathbf{R}_1^m + \mathbf{v}_2 \sum_{n=0}^{L-2} \mathbf{R}_2^n + \pi_L \right) \mathbf{1} \\
&= 1,
\end{aligned} \tag{A.16}$$

where the second equality is based on (3.32). The fourth equality is obtained by letting $m = q - 1$ and $n = L - q - 1$. The last equality is based on (3.33b).

This completes the proof.

Appendix B

Appendixes of Chapter 4

B.1 Proof of Theorem 4.1

i) Proof of the upper bound. Since the solution space of (RMP-P2), $\hat{\mathcal{F}}'_k$, is a subspace of that of (P2), $\hat{\mathcal{F}}_k$, the optimal objective function of (RMP-P2) must be no better than that of (P2), i.e., $z_{RMP-P2}^* \geq z_{P2}^*$. This completes the proof of the upper bound.

ii) Proof of the lower bound. At any iteration n , from (PP-P2), we have

$$\begin{aligned}\Delta^* &= \min\{1 - \sum_{i=1}^{|\hat{\mathcal{N}}_k|} v_i^* f_{i,m}\} \\ \Rightarrow \Delta^* &\leq 1 - \sum_{i=1}^{|\hat{\mathcal{N}}_k|} v_i^* f_{i,m} \\ \Rightarrow \Delta^* + \sum_{i=1}^{|\hat{\mathcal{N}}_k|} v_i^* f_{i,m} &\leq 1, \quad \forall m = 1, \dots, |\hat{\mathcal{F}}_k|. \end{aligned} \quad (\text{B.1})$$

Note that $\Delta^* < 0$ before the iteration process converges to the optimal solution to (P2) and $f_{i,m} \in \{0, 1\}$. Thus, if $\mathbf{f}_m = [f_{1,m}, f_{2,m}, \dots, f_{|\hat{\mathcal{N}}_k|,m}]^\top$ is not an all-zero vector, we have

$$\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \Delta^* f_{i,m} \leq \Delta^*, \quad \forall m = 1, \dots, |\hat{\mathcal{F}}_k|. \quad (\text{B.2})$$

From (B.1) and (B.2), we have

$$\begin{aligned}\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \Delta^* f_{i,m} + \sum_{i=1}^{|\hat{\mathcal{N}}_k|} v_i^* f_{i,m} &\leq 1 \\ \Rightarrow \sum_{i=1}^{|\hat{\mathcal{N}}_k|} f_{i,m} (\Delta^* + v_i^*) &\leq 1, \quad \forall m = 1, \dots, |\hat{\mathcal{F}}_k|. \end{aligned} \quad (\text{B.3})$$

If \mathbf{f}_m is an all-zero vector, then we have $\sum_{i=1}^{|\hat{\mathcal{N}}_k|} f_{i,m} (\Delta^* + v_i^*) = 0 < 1$. Therefore, regardless of \mathbf{f}_m being all-zero or not, inequality (B.3) always holds.

On the other hand, the dual problem (DP) of (P2) can be formulated as

$$(\mathbf{DP2}) : s = \max_{v_i \geq 0} \sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i v_i \quad (\text{B.4a})$$

$$\text{s.t. } \sum_{i=1}^{|\hat{\mathcal{N}}_k|} f_{i,m} v_i \leq 1, \quad \forall m = 1, \dots, |\hat{\mathcal{F}}_k|. \quad (\text{B.4b})$$

Thus, based on inequality (B.3) and dual constraint (B.4b), we can conclude that, if $\tilde{v}_i = \Delta^* + v_i^* \geq 0$, then \tilde{v}_i is a feasible solution to (DP2) and we have

$$\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i \tilde{v}_i \leq \sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i v_{i,DP2}^* = s_{DP2}^*, \quad (\text{B.5})$$

where $v_{i,DP2}^*$ and s_{DP2}^* are the optimal solution and the optimal objective function of (DP2), respectively. If $\tilde{v}_i < 0$, it is easy to check that inequality (B.5) still holds because $\tau_i \geq 0$ and $v_{i,DP2}^* \geq 0$. Thus, inequality (B.5) is always valid regardless of $\tilde{v}_i \geq 0$ or $\tilde{v}_i < 0$.

According to the *Strong Duality Theorem*, the optimal objective function of the dual problem (DP2) equals that of the primal problem (P2), i.e., $s_{DP2}^* = z_{P2}^*$. Therefore, we have $\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i \tilde{v}_i \leq s_{DP2}^* = z_{P2}^*$, which means $\sum_{i=1}^{|\hat{\mathcal{N}}_k|} \tau_i \tilde{v}_i$ is a lower bound of the optimal objective function of (P2). This completes the proof.

B.2 Proof of Theorem 4.2

At any iteration n , from (PP-P3), we have

$$\begin{aligned} \delta^* &= \max \left\{ \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} (r_i + w_i^*) f_{i,m} - \varphi^* \right\} \\ \Rightarrow \delta^* &\geq \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} (r_i + w_i^*) f_{i,m} - \varphi^* \\ \Rightarrow \varphi^* + \delta^* - \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} w_i^* f_{i,m} &\geq \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} r_i f_{i,m}, \quad \forall m = 1, \dots, |\hat{\mathcal{F}}_{k^*}|. \end{aligned} \quad (\text{B.6})$$

In addition, the dual problem of (P3) can be formulated as

$$(\mathbf{DP3}) : s_{DP3} = \min_{\varphi, w_i \geq 0} T\varphi - \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} \tau_i w_i \quad (\text{B.7a})$$

$$\text{s.t. } \varphi - \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} f_{i,m} w_i \geq \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} r_i f_{i,m}, \quad \forall m = 1, \dots, |\hat{\mathcal{F}}_{k^*}|. \quad (\text{B.7b})$$

Based on (B.6) and (B.7b), we can conclude that $(\tilde{\varphi}, w_i^* : i = 1, \dots, |\hat{\mathcal{N}}_{k^*}|)$ is a feasible solution to (DP3) where $\tilde{\varphi} = \varphi^* + \delta^* > 0$ (note that dual variable φ^* is always no less than zero and $\delta^* > 0$ before the iteration process reaches the optimality of (P3)).

Similar to the proof procedure of Theorem 4.1, we must have $T\tilde{\varphi} - \sum_{i=1}^{|\hat{\mathcal{N}}_{k^*}|} \tau_i w_i^* \geq s_{DP3}^* = z_{P3}^*$, where s_{DP3}^* and z_{P3}^* are the optimal objective functions of (DP3) and (P3), respectively. This completes the proof.

Appendix C

Appendixes of Chapter 5

C.1 Proof of Theorem 5.1

Based on PP1, we have

$$\Delta_k^{m*} = \max \left\{ \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k}^m - \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k}^m \varphi_i^{r*} - \sum_{j \in \mathcal{N}_v} \sum_{i \in \mathcal{N}_r} y_{i,j,k}^m \varphi_j^{v*} \right. \\ \left. - \sum_{i \in \mathcal{N}_r} y_{i,ES,k}^m \psi_m^* - \sum_{i \in \mathcal{N}_r} f_i y_{i,ES,k}^m \psi_{ES}^* - \mu_m^*, \forall m \in \mathcal{M}, k \in \mathcal{K}_m \right\}, \quad (\text{C.1})$$

$$\Rightarrow \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k}^m \varphi_i^{r*} + \sum_{j \in \mathcal{N}_v} \sum_{i \in \mathcal{N}_r} y_{i,j,k}^m \varphi_j^{v*} + \sum_{i \in \mathcal{N}_r} y_{i,ES,k}^m \psi_m^* + \sum_{i \in \mathcal{N}_r} f_i y_{i,ES,k}^m \psi_{ES}^* \\ + (\mu_m^* + \Delta_k^{m*}) \geq \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k}^m, \forall m \in \mathcal{M}, k \in \mathcal{K}_m. \quad (\text{C.2})$$

In addition, the dual problem of Relaxed IP1 can be formulated as

$$[\text{DP1}] \quad S_{DP1} = \min_{\varphi_i^c, \varphi_j^v, \psi_m, \psi_{ES}, \mu_m} \sum_{i \in \mathcal{N}_r} \varphi_i^c + \sum_{j \in \mathcal{N}_v} \varphi_j^v + \sum_{m \in \mathcal{M}} \psi_m + F_{ES} \psi_{ES} + \sum_{m \in \mathcal{M}} \mu_m \\ \text{s.t.} \quad \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k}^m \varphi_i^c + \sum_{j \in \mathcal{N}_v} \sum_{i \in \mathcal{N}_r} y_{i,j,k}^m \varphi_j^v + \sum_{i \in \mathcal{N}_r} y_{i,ES,k}^m \psi_m + \sum_{i \in \mathcal{N}_r} f_i \cdot y_{i,ES,k}^m \psi_{ES} + \mu_m \\ \geq \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k}^m, \forall m \in \mathcal{M}, k \in \mathcal{K}_m, \quad (\text{C.3a})$$

$$\varphi_i^c, \varphi_j^v, \psi_m, \psi_{ES} \geq 0, \quad \forall i \in \mathcal{N}_r, j \in \mathcal{N}_v, m \in \mathcal{M}. \quad (\text{C.3b})$$

According to (C.2) and (C.3a), we can observe that $(\varphi_i^{r*}, \varphi_j^{v*}, \psi_m^*, \psi_{ES}^*, \mu_m^* + \Delta_k^{m*})$ is a feasible solution to DP1. Therefore, we must have $\sum_{i \in \mathcal{N}_r} \varphi_i^{r*} + \sum_{j \in \mathcal{N}_v} \varphi_j^{v*} + \sum_{m \in \mathcal{M}} \psi_m^* + F_{ES} \psi_{ES}^* + \sum_{m \in \mathcal{M}} (\mu_m^* + \Delta_k^{m*}) \geq S_{DP1}^* = U_{IP1}^*$ where S_{DP1}^* and U_{IP1}^* are the optimal objective function values of the dual problem DP1 and the primal problem Relaxed IP1, respectively. This completes the proof.

C.2 Proof of Theorem 5.2

Based on PP2, we have

$$\Delta_k^* = \max \left\{ \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k} - \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k} \varphi_i^{r*} - \sum_{\forall j \in \mathcal{N}_v} \sum_{i \in \mathcal{N}_r} y_{i,j,k} \varphi_j^{v*} - \sum_{i \in \mathcal{N}_r} f_i \cdot y_{i,ES,k} \psi_{ES}^* - \mu^* - \sum_{q \in \underline{\mathcal{Q}}^u} I_{q,k} \cdot \underline{\mu}_q^* + \sum_{q \in \overline{\mathcal{Q}}^u} I_{q,k} \cdot \overline{\mu}_q^*, \forall k \in \mathcal{K} \right\}, \quad (\text{C.4})$$

$$\begin{aligned} \Rightarrow & \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k} \varphi_i^{r*} + \sum_{\forall j \in \mathcal{N}_v} \sum_{i \in \mathcal{N}_r} y_{i,j,k} \varphi_j^{v*} + \sum_{i \in \mathcal{N}_r} f_i \cdot y_{i,ES,k} \psi_{ES}^* + (\mu^* + \Delta_k^*) \\ & + \sum_{q \in \underline{\mathcal{Q}}^u} I_{q,k} \cdot \underline{\mu}_q^* - \sum_{q \in \overline{\mathcal{Q}}^u} I_{q,k} \cdot \overline{\mu}_q^* \geq \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k}, \quad \forall k \in \mathcal{K}. \end{aligned} \quad (\text{C.5})$$

In addition, the dual problem of Relaxed IP2 is formulated as

$$\begin{aligned} [\text{DP2}] \quad S_{DP2} = & \min_{\varphi_i^c, \varphi_j^v, \psi_{ES}, \mu, \underline{\mu}_q, \overline{\mu}_q} \sum_{i \in \mathcal{N}_r} \varphi_i^c + \sum_{j \in \mathcal{N}_v} \varphi_j^v + F_{ES} \psi_{ES} + M\mu \\ & + \sum_{q \in \underline{\mathcal{Q}}^u} \lfloor \beta_q \rfloor \underline{\mu}_q - \sum_{q \in \overline{\mathcal{Q}}^u} \lceil \beta_q \rceil \overline{\mu}_q, \\ \text{s.t.} \quad & \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} y_{i,j,k} \varphi_i^c + \sum_{\forall j \in \mathcal{N}_v} \sum_{i \in \mathcal{N}_r} y_{i,j,k} \varphi_j^v + \sum_{i \in \mathcal{N}_r} f_i \cdot y_{i,ES,k} \psi_{ES} + \mu \\ & + \sum_{q \in \underline{\mathcal{Q}}^u} I_{q,k} \cdot \underline{\mu}_q - \sum_{q \in \overline{\mathcal{Q}}^u} I_{q,k} \cdot \overline{\mu}_q \geq \sum_{i \in \mathcal{N}_r} \sum_{j \in \mathcal{N}_i \cup \{ES\}} a_{i,j} y_{i,j,k}, \quad \forall k \in \mathcal{K}. \end{aligned} \quad (\text{C.6a})$$

Thus, according to (C.5) and (C.6a), we can observe that $(\varphi_i^{r*}, \varphi_j^{v*}, \psi_{ES}^*, \mu^* + \Delta_k^*, \underline{\mu}_q^*, \overline{\mu}_q^*)$ is a feasible solution to DP2. Therefore, we have $\sum_{i \in \mathcal{N}_r} \varphi_i^{r*} + \sum_{j \in \mathcal{N}_v} \varphi_j^{v*} + F_{ES} \psi_{ES}^* + M(\mu^* + \Delta_k^*) + \sum_{q \in \underline{\mathcal{Q}}^u} \lfloor \beta_q \rfloor \underline{\mu}_q^* - \sum_{q \in \overline{\mathcal{Q}}^u} \lceil \beta_q \rceil \overline{\mu}_q^* \geq S_{DP2}^* = U_{IP2}^*$, where S_{DP2}^* and U_{IP2}^* are the optimal objective functions of the dual problem DP2 and the primal problem Relaxed IP2, respectively. This completes the proof.

List of Publications

- [1] **Shiwei Huang**, Jun Cai, and Changyan Yi, “Joint Admission Control and Resource Management for D2D-Assisted Mobile Edge Computing,” submitted to *IEEE Transactions on Mobile Computing*.
- [2] **Shiwei Huang**, Jun Cai, Hongbin Chen, and Feng Zhao, “Low-complexity Priority-aware Interference-avoidance Scheduling for Multi-user Coexisting Wireless Networks,” *IEEE Transactions on Wireless Communications*, vol. 17, no. 1, pp. 112–126, Jan. 2018.
- [3] **Shiwei Huang** and Jun Cai, “An Analysis Framework for Buffer-Aided Relay-ing Under Time-Correlated Fading Channels,” *IEEE Transactions on Vehicular Technology*, vol. 65, no. 9, pp. 6987–6999, Sep. 2016.
- [4] **Shiwei Huang**, Jun Cai, Hongbin Chen and Hong Zhang, “Transmit Power Optimization for Amplify-and-Forward Relay Networks With Reduced Over-heads,” *IEEE Transactions on Vehicular Technology*, vol. 65, no. 7, pp. 5033–5044, Jul. 2016.
- [5] **Shiwei Huang**, Hongbin Chen, Jun Cai, and Feng Zhao, “Energy Efficiency and Spectral-Efficiency Tradeoff in Amplify-and-Forward Relay Networks,” *IEEE Transactions on Vehicular Technology*, vol. 62, no. 9, pp. 4366–4378, Nov. 2013.
- [6] Zhen Zhao, **Shiwei Huang**, and Jun Cai, “An Analytical Framework for IEEE 802.15.6-Based Wireless Body Area Networks with Instantaneous Delay Constraints and Shadowing Interruptions,” accepted for publication in the *IEEE Transactions on Vehicular Technology*.
- [7] Changyan Yi, **Shiwei Huang**, and Jun Cai “An Incentive Mechanism Integrating Joint Power, Channel and Link Management for Social-Aware D2D Content Sharing and Proactive Caching,” accepted to be published in *IEEE Transactions on Mobile Computing*.
- [8] Huijin Cao, Hongqiao Tian, Jun Cai, Attahiru S. Alfa and **Shiwei Huang**, “Dynamic Load-balancing Spectrum Decision for Heterogeneous Services Pro-

- visioning in Mutil-channel Cognitive Radio Networks,” accepted to be published in *IEEE Transactions on Wireless Communications*.
- [9] Hong Zhang, Jun Cai, Xiaolong Li and **Shiwei Huang**, “Adaptive Service Rate and Vacation Length for Energy-Efficient HeNB Based on Queueing Analysis,” *IEEE Transactions on Vehicular Technology*, vol. 65, no. 10, pp. 8696–8709, Oct. 2016.
- [10] Zhen Zhao, **Shiwei Huang**, and Jun Cai, “Energy Efficient Packet Transmission Strategies for Wireless Body Area Networks with Rechargeable Sensors (**Invited Paper**),” accepted to be published in *Proc. of IEEE Vehicular Technology Conference (VTC) 2017-Fall*, Toronto, Canada, Sep. 2017.
- [11] **Shiwei Huang**, Changyan Yi, and Jun Cai, “A Sequential Posted Price Mechanism for D2D Content Sharing Communications,” in *proc. of 2016 IEEE Global Communications Conference (GLOBECOM)*, Washington DC, pp. 1–6, Dec. 2016.
- [12] **Shiwei Huang**, Jun Cai and Hong Zhang, “Relay Selection for Average Throughput Maximization in Buffer-Aided Relay Networks,” in *proc. of 2015 IEEE International Conference on Communications (ICC)*, London, UK, pp. 1994–1998, Jun. 2015.
- [13] **Shiwei Huang** and Jun Cai, “Priority-Aware Scheduling for Coexisting Wireless Body Area Networks (**Invited Paper**),” in *proc. of 2015 International Conference on Wireless Communications and Signal Processing (WCSP)*, Nanjing, China, pp. 1–5, Oct. 2015.
- [14] Hongqiao Tian, Jun Cai, Attahiru S. Alfa, **Shiwei Huang** and Huijin Cao, “Dynamic Load-Balancing Spectrum Decision for Cognitive Radio Networks with Multi-Class Services,” in *proc. of 2015 International Conference on Wireless Communications and Signal Processing (WCSP)*, Nanjing, China, pp. 1–5, Oct. 2015.