# Performance Optimization of Cloud Radio Access Networks

Mohamed Labana

**A Thesis**

**in**

**The Department**

**of**

**Electrical and Computer Engineering**

**Presented in Partial Fulfillment of the Requirements for the Degree of**

**Master of Applied Science (Electrical and Computer Engineering) at**

**Concordia University**

**Montréal, Québec, Canada**

**March 2021**

## CONCORDIA UNIVERSITY

### School of Graduate Studies

This is to certify that the thesis prepared

By:              Mohamed Labana

Entitled:       **Performance Optimization of Cloud Radio Access Networks**

and submitted in partial fulfillment of the requirements for the degree of

**Master of Applied Science (Electrical and Computer Engineering)**

complies with the regulations of this University and meets the accepted standards with respect to originality and quality.

Signed by the Final Examining Committee:

_____ Chair
*Dr. D. Qiu*

_____ External Examiner
*Dr. A. Youssef (CIISE)*

_____ Examiner
*Dr. D. Qiu*

_____ Supervisor
*Dr. Walaa Hamouda*

Approved by      _____
Dr. Y.R. Shayan, Chair
Department of Electrical and Computer Engineering

March, 2021          _____
Dr. Mourad Debbabi, Interim Dean
Gina Cody School of Engineering and Computer Science

# Abstract

Performance Optimization of Cloud Radio Access Networks

Mohamed Labana

The exponential growth of cellular data traffic over the years imposes a hard challenge on the next cellular generations. The cloud radio access network (CRAN) is an emerging cellular architecture that is expected to face that challenge effectively. The main difference between the CRAN architecture and the conventional cellular architecture is that the baseband units (BBUs) are aggregated at a centralized baseband unit pool, hence, enabling statistical multiplexing gains. However, to acquire the several advantages offered by the CRAN architecture, efficient optimization algorithms and transmission techniques should be implemented to enhance the network performance. Hence, in this thesis, we consider jointly optimizing user association, resource allocation and power allocation in a two tier heterogeneous cloud radio access network (H-CRAN). Our objective is to utilize all the network resources in the most efficient way to maximize the network average throughput, while keeping some constraints such as the quality of service (QoS), interference protection to the devices associated with the Macro remote radio head (MRRH), and fronthaul capacity. In our system, we propose using coordinated multi-point (CoMP) transmissions to utilize any excess resources to maximize the network performance, in contrast to the literature, in which CoMP is usually used only to support edge users. We divide our joint problem into three sub-problems: user association, radio resource allocation, and power allocation. We propose matching game based low complexity algorithms to tackle the first two sub-problems. For the power allocation sub-problem, we propose a novel technique to convexify the non-convex original problem to obtain the optimal solution. Given the conducted simulations, our proposed algorithms proved to enhance the network average weighted sum rate significantly, compared to the state of the art algorithms in the literature.

The high computational complexity of the optimization techniques currently proposed in the

literature prevents from totally reaping the benefits of the CRAN architecture. Learning based techniques are expected to replace the conventional optimization techniques due to their high performance and very low online computational complexity. In this thesis, we propose tackling the power allocation in CRAN via an unsupervised deep learning based approach. Different from the previous works, user association is considered in our optimization problem to reflect a real cellular scenario. Additionally, we propose a novel scheme that can enhance the deep learning based power allocation approaches, significantly. We provide intensive analysis to discuss the trade-offs faced when employing our deep learning based approach for power allocation. Simulation results prove that the proposed technique can obtain a very close to optimal performance with negligible computational complexity.

# Acknowledgments

First and foremost, I thank Allah The Almighty for giving me this opportunity and strength to accomplish my degree.

I would like to express my deepest gratitude to my supervisor, Prof. Walaa Hamouda, for his patience and continuous support. Without his guidance and care, I would not have been able to complete this thesis.

I would also like to thank my colleagues in the Lab for the wonderful times and interesting discussions we had over the year. Their companionship and comradeship have significantly enriched my experience at Concordia University.

Finally, I have and would love to thank my beloved parents. I owe them everything in life. Their support and motivation enabled me to complete my research.

# Contents

# List of Figures

# List of Tables

# List of Abbreviations

| | |
|---|---|
| BBU | Baseband unit |
| BS | Base Station |
| CoMP | Coordinated multi-point |
| CRAN | Cloud Radio Access Network |
| CSI | Channel State Information |
| DAS | Distributed Antenna Systems |
| DBN | Deep Belief Network |
| DNN | Deep Neural Network |
| GA | Genetic Algorithm |
| H-CRAN | Heterogeneous Cloud Radio Access Network |
| HD | High Definition |
| IoT | Internet of Things |
| MC-NOMA | Multi-Carrier Non-Orthogonal Multiple Access Multiple Access |
| MRRH | Macro Remote Radio Head |
| PRRH | Pico Remote Radio Head |
| PSD | Power Spectral Density |
| QoE | Quality of Experience |
| QoS | Quality of Service |
| RAU | Remote Access Unit |
| RF | Radio Frequency |
| RRH | Remote Radio Head |
| SCA | Successive Convex Approximation |

| | |
|---|---|
| SDN | Software Defined Network |
| SINR | Signal-to-Interference-plus-Noise Ratio |
| SNR | Signal-to-Noise Ratio |
| SVMs | Support Vector Machines |
| TI | Tactile Internet |
| URLLC | Ultra Reliable Low Latency Communications |
| VoIP | Voice over IP |
| WMMSE | Weighted Minimum Mean Square Error |
| WSR | Weighted Sum Rate |

# Notations

| | |
|---|---|
| MU | Macro RRH User |
| PA | Power Allocation |
| PU | Pico RRH User |
| RA | Resource Allocation |
| RB | Resource Block |
| SU | Small RRH User |
| UA | User Association |
| UE | User Equipment |
| s.t. | Subject to |
| $\sigma^2$ | Noise Variance |
| $\odot$ | Hadamard Product (Element-wise Product) |
| $\mathbb{E}\left[\cdot\right]$ | Expected Value |
| $\left[\cdot\right]^T$ | Transpose |

# Chapter 1

# Introduction

## 1.1 Cloud Radio Access Networks

There is a massive evolution in communications industry. New applications that require reliable network connectivity, such as high definition (HD) videos, live video streaming, voice over IP (VoIP), and virtual reality are introduced everyday. Moreover, billions of internet of things (IoT) devices are expected to be fully operating by 2022 [1]. Thus, these IoT devices will also have their share from the cellular networks resources. The diverse types of applications also need different network requirements, and reliability metrics. For instance, IoT devices may be related to fire alarm systems or medical applications, so they will require ultra low latency requirements. Consequently, with this diversity in applications and different requirements, current research have been focused on developing new cellular architectures and network technologies that can face these challenges. Owing to these research efforts, the cloud radio access network (CRAN), an evolutionary cellular architecture was proposed, and is expected to play a key role in the upcoming cellular generations [2].

The CRAN architecture divides the base stations in the traditional cellular networks into remote radio heads (RRHs), which usually perform only the basic radio frequency functions; and baseband units (BBUs) that are usually aggregated in a centralized position [3]. This architecture enables supporting higher capacities with much lower costs. Additionally, the centralized processing at the BBU pool enables achieving much better network performance and lower power consumption, thanks to the statistical multiplexing gain.

## 1.2 Coordinated Multi-point Transmissions

Coordinated multipoint (CoMP) is a technique that enables serving a user by more than one base station in a cellular network. It helps to improve the network performance, such as enhancing the network throughput, increasing the coverage probability, and improving the quality of service (QoS) for the cell edge users [4].

In one of our contributions, we consider combining both emerging technologies, CRAN and CoMP. In particular, we consider a downlink heterogeneous H-CRAN system while implementing CoMP. We focus on jointly optimizing the user association, resource allocation, and power allocation to maximize the network throughput and utilize any available resources to enhance the network performance in the most efficient way.

## 1.3 Optimization in CRAN

Optimization techniques are essential for cellular networks. The main aim of optimization is to ensure that the network resources are utilized in the most efficient way to maximize the network performance. Several metrics can be considered for optimization according to the type of users of the network. The most important metrics that are usually considered by researchers as the optimization objectives are maximizing the network weighted sum rate, minimizing power consumption, or minimizing network latency. The ultimate goal of optimization is to achieve the optimal performance, according to the metric being optimized, subject to some constraints that must be satisfied. For instance, in green networks, the objective is usually to minimize the power consumption, subject to a condition that each user must achieve a minimum data rate that is determined according to the user's target quality of service (QoS).

## 1.4 Motivations

The CoMP transmission technique was initially limited to serve the network edge users, which usually suffer high interference. Following that, researchers started to remove this limitation and use CoMP transmissions to serve generic users in order to improve the network performance. However,

the number of remote radio heads (RRHs) serving a specific user was usually limited (or fixed) in the proposed works. Differently, we propose using generalized CoMP transmissions to utilize any excess resources in the network to obtain higher throughput. In our work, the number of RRHs serving some user is not limited, such that it is decided according to the network status to obtain the best possible performance.

One of the main challenges that prevents from totally reaping the benefits of the centralized optimization is the high computational complexity of the conventional optimization algorithms introduced in the literature. This high complexity results in increased costs in the required hardware; high latency, since these algorithms require large time to be processed; and less power efficiency. Moreover, in ultra dense network settings, where the network might have several thousands of small cells and huge number of users, the implementation of these algorithms might be infeasible.

Due to the aforementioned drawbacks of the conventional optimization techniques, recent research has been focused on developing novel low complexity optimization techniques that can offer performance close to, or outperform the conventional ones. Game theoretic based algorithms proved to be efficient in solving optimization problems with integer variables that are usually intractable and can be solved optimally, solely, through the exhaustive search technique [5,6]. One of the most promising solutions is to opt for techniques based on machine learning and artificial intelligence. The main advantage of machine learning based optimization is that its computational complexity is only experienced during the model training phase [7]. Once a machine learning model is trained successfully, it can be used for prediction with negligible computational complexity. Thus, utilizing machine learning based optimization techniques in cellular networks can alleviate the huge burden on the cellular baseband units [8], and reduce the latency significantly [9]. Particularly, reducing the latency is a superior advantage for networks serving IoT users [10]. Moreover, the hardware requirements, and accordingly the costs of construction of baseband units will highly decrease, as the training of the machine learning models does not have to be done in each BBU. Rather, these models can be trained in a different location then transmitted to be utilized in several BBUs.

Tackling the power allocation optimization problem in wireless networks via deep neural networks (DNNs) was proposed in the literature. However, a cellular network scenario where users have the possibility to associate with several RRHs was not considered. Here and different from

previous works, we consider optimizing the power allocation in a CRAN with DNNs, while considering user association to reflect a real cellular scenario, where users have the possibility to associate with several RRHs, and the locations of the users are highly random.

## 1.5 Thesis Contributions

Based on the previous discussions, and motivated by the void in the literature, the contributions of this thesis can be divided into two main parts. Firstly, we consider solving a joint user association, resource allocation, and power allocation optimization problem in H-CRAN, while implementing CoMP. To the best of our knowledge, no research work considered a similar problem while utilizing CoMP transmissions. Our contributions in this part can be summarized as follows:

- We propose using generalized CoMP transmissions to utilize any excess resources in the network to obtain higher throughput.

- We propose matching game based low complexity algorithms to tackle the user association and resource allocation sub-problems. Matching game was also proposed in [11] to solve a similar optimization problem. However, CoMP was not supported in their proposed algorithms.

- Our user association algorithm proved to realize a good tradeoff between the cooperation gain and fronthaul consumption, to achieve a cooperation gain even in case of tight fronthaul constraints.

- Furthermore, we propose a novel approach to convexify the power allocation sub-problem to obtain an optimal solution.

Regarding the second part, we aim to combine the benefits of both, the CRAN architecture and the deep learning based optimization algorithms. Specifically, we consider jointly optimizing the user association and power allocation in CRAN downlink transmission, where the power allocation sub-problem will be tackled via a deep neural network (DNN) based algorithm. The contributions of this part can be summarized as follows:

- We optimize the power allocation in a CRAN with DNNs, while considering user association to reflect a real cellular scenario.

- We provide intensive simulations to discuss the trade-off between the ability of the DNNs to obtain high performance, and their ability to predict solutions that satisfy the optimization problem constraints (in our case, QoS constraints).

- Furthermore, we propose a novel approach to enhance the ability of the DNNs to obtain higher performance and better constraint preservation capability.

## 1.6  Thesis Organization

The rest of this thesis will be organized as follows:

In chapter 2, we give a background about the CRAN architecture and coordinated multi-point transmissions. Moreover, we explain the optimization objectives and optimization parameters targeted by researchers in CRAN. Finally, we review some related works.

We discuss our proposed generalized CoMP scheme in chapter 3. A detailed explanation of the system model, problem formulation and proposed algorithms will be provided. Furthermore, we prove by simulations the superiority of the proposed algorithms, compared to the literature.

In chapter 4, our proposed DNNs based power allocation scheme will be introduced. We provide a detailed explanation of the proposed DNNs architecture, training process, and prediction process. Additionally, we provide extensive simulations to explain the trade-offs associated with utilizing our proposed algorithm. Moreover, we compare the proposed algorithm to the state-of-the-art conventional optimization algorithms available in the literature in terms of the achieved performance and the computational complexity. We prove that the proposed algorithm clearly outperforms the conventional ones.

Finally, in chapter 5, we draw our conclusions; discuss the challenges facing the CRAN performance optimization, and the possible research directions.

# Chapter 2

# Background and Literature Review

## 2.1 Introduction

The evolution in the telecommunications industry has led to developing novel cellular technologies and architectures. One of the proposed promising technologies is the software defined network (SDN), in which the network control plane is disassociated from the network devices and transferred to a logical control center to be centrally managed using software applications. Integrating the concept of SDNs with cellular networks allows the network resources to be centrally managed and allocated to the users based on a comprehensive view of the network, which helps significantly in improving the users' QoE (quality of experience). The authors in [12] gave a comprehensive overview on SDNs, and proposed a method to enhance the QoE of software defined multi-tier Long Term Evolution-Advanced (LTE-A) networks' users. In addition to SDNs, the cloud radio access network (CRAN), an evolutionary cellular architecture was proposed, and is expected to have a key role in the upcoming cellular generations [2].

The main contribution in the CRAN architecture was that the regular basestations in the ordinary cellular networks were decoupled into two units. The first unit is the remote radio head, which performs mainly the RF functions, while the second is the baseband unit (BBU) that performs the baseband signal processing functions. Moreover, the baseband units (BBUs) of numerous RRHs are aggregated at a centralized position, namely, baseband unit pool (BBU pool) [11].

There are several perks for the CRAN architecture, due to the statistical multiplexing gain.

6

Firstly, there is a large reduction in construction costs when compared to the conventional cellular networks, as the BBUs are aggregated at one place; hence, they need only a single cooling unit. This also results in a significant reduction in the power consumption. Furthermore, since the BBU pool controls numerous RRHs, more complicated tasks can be synchronized and performed co-operatively and smoothly between the RRHs, such as coordinated multi-point transmissions, and multi-connectivity, in which a user can be served by multiple RRHs.

To reap the benefits of the CRAN architecture, high performance optimization techniques must be used to obtain the important system parameters. Several research works were proposed in the literature focusing on CRAN performance optimization, some of which can obtain significant performance improvements.

In this chapter, we aim to give an overview of the CRAN architecture. Then, we explain the CoMP technique and discuss how it significantly improves the network performance. Moreover, we discuss the important optimization objectives targeted by researchers in CRAN, which reflect the key performance metrics of the network. Additionally, we discuss the important parameters that should be optimized, and how they affect the CRAN performance. Furthermore, we review some of the important research works on CRAN performance optimization with single association, and coordinated multi-point techniques. Finally, we review some research works utilizing machine learning based techniques for CRAN performance optimization, where it is known that learning based techniques have very low computational complexity.

## 2.2 H-CRAN Architecture

A CRAN usually consists of multiple network tiers. The first tier is formed of a macro RRH (MRRH), while the inner tiers can include micro, femto or pico RRHs. Fig. 4.2 shows a two-tier H-CRAN. Note that the macro RRH users are denoted as MUs, while the small RRHs users are denoted as SUs in the figure. It can be seen in Fig. 4.2 that BBUs are decoupled from the RRHs and aggregated in the BBU pool. Fronthaul links are used to connect the RRHs to the BBU pool. The most common method to deploy these links is through optical fibers [13]. However, the fronthaul links can be also in the form of microwave wireless links, or deployed using the Ethernet standard

to be shared among different RRHs [13]. The presence of the inner tiers allows the network to serve much more users, compared to the number served by the MRRH solely; hence, significantly improving the network coverage and spectral efficiency. Moreover, the inner-tier RHHs can be turned off when not serving any user, resulting in higher power efficiency. However, the presence of multiple network tiers causes inter-tier interference, if the radio resources are not orthogonal. Accordingly, optimization techniques are needed to reap the benefits of the H-CRAN architecture, while mitigating the inter-tier interference.

Fortunately, the CRAN architecture, which aggregates the baseband units (BBUs) together in the BBU pool, gives the opportunity to implement more centralized optimization techniques. This allows achieving much higher performance gains. Additionally, centralized optimization facilitates the implementation of interference mitigation techniques such as coordinated multi-point transmission.



Figure 2.1: Two-tier H-CRAN Architecture

8

## 2.3 Coordinated Multi-point Technique

Coordinated multi-point (CoMP) technique is expected to be a key factor in 5G systems and beyond. It allows a user to be served by numerous basestations (RRHs in our case) simultaneously using the same frequency block. This helps to significantly improve the QoS, especially for the network edge users. Consequently, utilizing CoMP can enhance the network performance greatly [4]. Fig. 2.2 shows a typical implementation of CoMP, in which the edge users, who usually suffer high interference are served with multiple RRHs. Noting that a user equipment is denoted as UE in the figure.



Figure 2.2: Coordinated multi-point transmission

To illustrate more how CoMP improves the system performance, consider a network consisting of numerous RRHs which belong to the set $J = \{1, 2, 3...M\}$ and users belonging to the set $U = \{1, 2, 3...K\}$. If CoMP is not used, the signal to interference plus noise ratio ($SINR$) received by user $i$ served by RRH $j$ can be expressed as:

$$SINR_i = \frac{P_j g_{ji}}{\sum_{m \in J/\{j\}} P_m g_{mi} + \sigma^2},$$ 

(2.1)

where $P_j$ is the transmission power of RRH $j$, $g_{ji}$ is the channel gain between user $i$ and RRH $j$, and $\sigma^2$ is the noise variance. Now, assume that CoMP is adopted, and RRHs $j_1, j_2, j_3$ are serving user $i$. Hence, $SINR_i$ can be rewritten as:

$$SINR_i = \frac{P_{j_1} g_{j_1 i} + P_{j_2} g_{j_2 i} + P_{j_3} g_{j_3 i}}{\sum_{m \in J/\{j_1, j_2, j_3\}} P_m g_{mi} + \sigma^2}, \qquad (2.2)$$

Obviously, it can be seen that the received useful signal power at $i$ is much higher at the second case. Moreover, the interference power has been significantly decreased. Thus, the CoMP technique simply maps the interference signals into useful signals. It can be easily deduced that receiving higher $SINR$ means larger data rates received by user $i$.

$$R_i = F \log_2(1 + SINR_i), \qquad (2.3)$$

where $F$ is the bandwidth of the utilized radio resource block.

## 2.4 Optimization Objectives

The optimization metrics are usually chosen according to the users. For instance, if the users are mainly IoT devices related to critical applications, such as fire detection, then minimizing the network latency is the most important optimization objective to consider. Another important metrics that are usually considered by researchers as the optimization objectives are maximizing the network weighted sum rate and minimizing power consumption. This section reviews some of the popular optimization objectives introduced in the literature.

- Maximizing the network weighted sum rate (throughput): with the evolution of the new cellular applications that require very high date rates such as high definition (HD) and 4K videos, it is logical that one of the most important metrics to optimize is the total rate achieved by the network users. Thus, if the CRAN or the cellular network in general is intended to mainly serve these types of applications, then maximizing the weighted sum rate is the best optimization objective to consider, in order to guarantee a good quality of experience (QoE).

- Minimizing network latency: some applications are delay critical, such as the medical communications, or vehicular security critical messages...etc. In general, communications including internet of things (IoT) devices are usually sensitive to the network delay, and do not require high data rates. Consequently, networks serving these type of devices should consider minimizing the network latency as the optimization metric.

- Minimizing total transmission power: using less power means less recurring costs for the network operator; and hence more profit. Additionally, reducing power consumption and energy saving is becoming a trend in many manufacturing sectors. Thus, green communications has received much importance recently, in which the main concern is reducing the energy consumption while maintaining the QoS.

- Maximizing the ratio between weighted sum rate and total transmission power: this objective function achieves a very good trade-off between the network throughput and the power consumption. Thus, the task here is to obtain the best performance with the lowest cost in terms of transmission power.

- Maximizing fairness among users: this can be achieved through several methods. One of the most popular optimization objective functions that ensure fairness is the max-min fairness function, in which the minimum achievable signal to interference plus noise ratio (SINR) for each user is maximized.

These are the most important optimization objectives usually considered in the literature, as they play a pivotal role in determining the network performance. In the next section, we discuss the most important parameters that should be optimized in order to achieve the optimization objectives.

## 2.5   Optimization Parameters

To achieve a superior performance, several parameters need to be optimized in a cellular network, according to the required performance metrics explained in the last section. In this section, we discuss the most important optimization parameters that should be taken into consideration in CRAN.

### 2.5.1  RRH-User Association

This is one of the most important parameters [14]. Users should be associated to the RRHs that can provide them with the best performance. One of the most popular user association techniques is to simply associate each user with the RRH providing the highest signal power. This technique can achieve acceptable performance with low computational complexity. Also, many algorithms based on game theory were proposed for the same purpose. Moreover, norm approximation methods proved to be efficient in solving optimization problems with integer values, in which RRH-user association is one of them. The idea of norm approximation methods is to relax the integer values in the combinatorial optimization problem [15], then iteratively solve the relaxed problem. However, algorithms based on game theory are the best candidates for the user association optimization problems, due to their low complexity and high performance. Game theoretic based algorithms converge after reaching the Nash Equilibrium [16], in which the players reach the maximum possible utility according to the game scenario. The concerned players can be the base stations (RRHs in our case) [17], users [18], or both [19]. One of the important gaming models introduced in the literature is the bargaining game [20,21], where players negotiate with each others to obtain the highest possible mutual gains. Another gaming model is the matching game considered in [11], where each user should be matched to the RRH offering the highest data rate, while taking the fronthaul capacity into consideration.

### 2.5.2  Radio Resource Allocation

Optimization problems including radio resource allocation usually include integer values. Thus, they can be tackled via the classic non-convex optimization, where the binary variables can be relaxed with numerous techniques including the norm approximation methods [15, 22]. In general, a user should be assigned the radio resource with highest channel gain to obtain the best performance. Another way to tackle the radio resource allocation optimization problems is through the game theoretic based techniques. The gaming models utilized in the literature are similar to the used models for user association, such as the bargaining game [20, 21], matching game [11], and coalition game [23]. While the matching game and bargaining game were discussed in the previous

subsection, the idea of the coalition game is to divide the players into different sets (coalitions). The players in each coalition form a cooperative team that aims to maximize their gains during the game, while competing against the other coalitions.

### 2.5.3 Power Allocation

Optimizing the values of the signal powers allocated from each RRH to each user can help in reducing the intra-tier and inter-tier interference significantly. Consequently, power allocation is one of the critical parameters to be optimized in CRAN. Nevertheless, power allocation is usually a non-convex optimization problem. Thus, it is usually tackled with some high complexity algorithms. Some of the popular algorithms used for power allocation are: genetic algorithm [24], successive convex approximation [25], and weighted minimum mean square error [26]. These algorithms are highly iterative with significant computational complexity, such that their application in networks with large number of users and RRHs (ultra-dense networks) might be infeasible.

### 2.5.4 RRH Activation

In optimization problems with the objective of minimizing power consumption, or maximizing the power efficiency, switching off some RRHs that are not going to highly affect the performance can cause significant power efficiency gain. Thus, the main task of RRH activation optimization is to save power while causing a minimal effect on the network performance, and maintaining the quality of service (QoS) constraints. Numerous approaches were proposed to optimize the RRH activation in CRAN. In [27], authors proposed using norm approximation methods to jointly optimize RRH activation and beamforming. Researchers in [28] considered jointly optimizing RRH activation and user association, where they adopted the max-sum algorithm to efficiently tackle their problem. Another work in [29] considered utilizing Benders decomposition to optimally tackle a joint optimization problem considering RRH activation, user association, and power allocation.

### 2.5.5 RRH-BBU Mapping

BBUs are aggregated in the BBU pool. Each BBU can be modeled as a virtual machine in the optimization process and mapped to serve multiple RRHs. Optimizing RRH-BBU mapping can

help reduce the latency significantly. Moreover, some BBUs can be turned off to increase system power efficiency, if this will not affect the performance subject to the system constraints. Some works were proposed to deal with RRH-BBU mapping optimization. There are several optimization objectives to consider for problems involving RRH-BBU mapping such as load balancing [30], or minimizing power consumption [31]. In [30], authors proposed a particle swarm technique for the RRH-BBU mapping in order to minimize the number of blocked calls and balance the load between BBUs. Load balancing is a very important metric as it can prevent the network congestion [32]. Additionally, researchers in [32] proposed a borrow-lend approach for RRH-BBU mapping that can improve the load balancing and prevent network congestion. Regarding power efficiency, authors in [31] formulated RRH-BBU mapping as a bin-packing problem, then proposed a heuristic algorithm to minimize the power consumption while maintaining users' QoS.

### 2.5.6  Cache Allocation

The types of content that are frequently requested by mobile network users should be cached in the CRAN cloud to reduce the burden on the core network and backhaul links. Moreover, since the cached content becomes closer to the users, the network performance can be improved significantly in terms of latency and QoS. Therefore, cache allocation optimization has significant attention lately from researchers. There are important factors to consider when optimizing cache allocation such as the associated user preferences and mobility, as well as the available hardware resources for caching [33].

As it can be noticed, most of the previously mentioned parameters are related to each other, and the way one of them is tuned will affect other parameters, and the whole network performance. Accordingly, numerous researches focus on solving joint optimization problems that consider several optimization parameters. Jointly considering these parameters improves the network performance significantly. Several examples will be discussed in the upcoming sections.

## 2.6 Optimization in CRAN With Single Association

Communication networks in which each user can only be associated with only one basestation are the simplest type of networks. The main advantage of this type of network operation is that the complexity of the optimization algorithms is relatively lower. In this section, we review some of the recent research works that focused on optimizing the CRAN performance under the single association constraint.

Authors in [11] aimed to maximize the weighted sum rate in a heterogeneous CRAN. They proposed jointly optimizing the RRH-user association, radio resource allocation, and power allocation. They proposed a novel scheme to tackle the joint problem, in which it was divided into three sub-problems. The first two sub-problems, namely, RRH-user association, and radio resource allocation were tackled through matching game based algorithms. The third sub-problem, power allocation, was relaxed via the high SINR regime to a convex problem and solved by updating the Lagrangian multipliers iteratively.

With a different optimization objective to minimize the network latency, researchers in [2] proposed solving a joint optimization problem including communication and computing resource allocation. The communication resources taken into consideration in their model included RRH-user association, radio resource block allocation, and power allocation. Regarding the computing resources, each BBU in their system was divided into virtual machines (VMs), each of them can serve one RRH. Thus, the computing resources considered in their optimization were namely, the RRH-BBU mapping, and VM allocation for the RRHs. To overcome the intractability of their optimization problems, they proposed new algorithms inspired from the auction theory. Tactile Internet (TI) is one of the services that are expected to be provided in the upcoming cellular generations. However, to provide such a service, URLLC (ultra reliable low latency communication) is required. Fortunately, the CRAN architecture allows researchers to look forward to achieving this type of communications. Thus, some research aimed to reap the benefits of CRAN architecture to satisfy the strict requirements of the TI applications. Authors of [34] considered a system model in which a CRAN is serving several pairs of tactile users, each of which are communicating together. They proposed a queuing model that is especially designed to fit the TI applications. The objective of

their optimization problem was formulated to minimize the overall power consumption, while satisfying the stringent delay constraints for the TI. The optimization parameters being considered in their system were mainly power allocation and radio resource allocation. In their model, they considered the queuing delay in addition to the fronthaul delay. In order to tackle the NP-hardness and non-convex nature of their formulated problem, they utilized techniques based on successive convex approximation (SCA), and the difference between two convex functions. Their algorithms proved to reduce the network power consumption while satisfying the strict latency requirements of the Tactile Internet.

These aforementioned works in CRAN optimization have improved the performance significantly, and opened the door for more developments. However, restricting the system to single association schemes prevents the network from achieving better performances that can be achieved using multiple association techniques. Especially that the architecture of the CRAN facilitates centralized optimization, which is the ideal environment to perform these techniques. In the next section, we discuss the recent advances in performance optimization of coordinated multi-point enabled CRAN.

## 2.7 Optimization in CoMP-Enabled CRAN

The high ability of CoMP to cancel the inter-tier and intra-tier interference, and transferring the interfering signals into useful ones, has driven many research works to start exploring how to use CoMP in the most efficient way. The authors in [35] proposed optimizing the resource and power allocation problems in a fronthaul constrained H-CRAN. They used a price-based outer iteration scheme to control the fronthaul capacity, and weighted minimum mean square error-based inner iteration approach to obtain the power allocation. In [20], the authors solve the clustering problem through a cooperative bargaining game, while ensuring fairness amongst users. They utilized CoMP transmission in order to mitigate the interchannel interference (ICI). Nevertheless, their solution is only applicable to CRANs with high capacity fronthaul links. The clustering problem was also tackled in [36], in which a CoMP heuristic user association was proposed to maximize the energy efficiency. Their algorithm showed enhancement in the system energy performance when compared to the baseline nearest RRH user association scheme. However, they only focused on a single

tier CRAN. Similar to [20], authors in [37] jointly optimized the user association, precoding and power allocation in a single tier CRAN, with an objective function aiming at minimizing the power consumption of the network. They proposed an iterative linear relaxed algorithm to approximate the fronthaul capacity constraint function then solving the corresponding linear problem, and hence obtaining a sub-optimal solution. In [5], authors proposed an algorithm to optimize the CoMP selection and resource allocation joint problem, where they considered only a single tier CRAN. The authors in [38] also considered the downlink power allocation problem in a single tier CRAN. They managed to reach the optimal solution but in a system model in which only one UE is served by multiple RRHs. Thus, there was no interference in their model, and the performance was only related to SNR.

In order to tackle the void in the literature, in the next chapter, we propose a generalized CoMP scheme that aims to maximize the network weighted sum rate. A two tier H-CRAN is considered. Thus, constraints to limit inter-tier interference are added, in addition to the QoS constraints. We consider all possibilities regarding the fronthaul links, which can be individual or shared links with a limited capacity, or high capacity links.

## 2.8   Deep Learning-Based Optimization in CRAN

As mentioned before, many optimization techniques are highly complex. This can cause a huge burden on the hardware performing the optimization in the cellular sites (in our case BBU pool). Additionally, it can lead to high latency, as the processing may take time to reach acceptable solutions. Consequently, more research should be carried in order to have less complex optimization algorithms that can perform accurately. The application of machine learning and deep learning algorithms in optimization can alleviate significantly its complexity. The key idea behind using machine learning in optimization is that the computational complexity due to machine learning comes only during the model training phase, which is performed offline. The machine learning model is said to be trained successfully if a target prediction accuracy can be reached. After the model training is accomplished, the online computational complexity due to prediction is very small and can be neglected. As we propose a deep learning based power allocation scheme in chapter 4,

we review here some related works.

Tackling the power allocation optimization problem in wireless networks via DNNs was proposed in the literature. Authors in [39] considered a joint problem of power allocation, scheduling, and flow allocation in a static wireless network. Their optimization objective was to maximize the weighted sum rate of the network under fairness constraints. To tackle their problem, they utilized supervised machine learning, in which they aimed to train their learning model to approximate the optimal solutions obtained from offline computations. Their learning based algorithm consisted of both support vector machines (SVMs) and deep belief networks (DBNs). The performance of their model proved to be superior to some of the other popular optimization algorithm. However, their static wireless network system model and proposed algorithms cannot be applied to a cellular environment where the users' locations are always varying, and users have many possibilities to associate with several base stations. Researchers in [40] also considered the power allocation problem in distributed antenna systems (DAS). Their system model consisted of several remote access units (RAUs) connected to their corresponding users on orthogonal radio resources. Thus, there was no interference in their model. They used DNNs trained through supervised learning to approximate the optimal solutions obtained offline via the sub-gradient technique. Their algorithm showed a high ability in obtaining performance close to the optimal. However, the fact that their model did not include interference effects prevents it from being considered in real cellular scenarios. A joint resource allocation and power allocation problem was considered in [41], with an objective function to minimize the total transmission power, while maintaining the QoS of users. In [41], the system model consisted of a single base station connected to numerous users through the multi-carrier non-orthogonal multiple access (MC-NOMA) technique. However, the considered system model does not reflect a practical scenario for a cellular network, since it is based only on one base station. Authors in [42] considered a power allocation problem in a generic wireless network consisting of multiple pairs of transmitters/receivers. They considered the two cases in which the transmitters/receiver pairs can have either fixed or random locations. They utilized unsupervised DNNs to tackle their power allocation problem. Their simulation results showed that unsupervised learning technique can outperform the traditional optimization algorithms. However, their generic system model cannot be directly applied to a cellular network, since they did not consider user association.

Different from the previous works, in this thesis, we consider a deep learning based power allocation algorithm while taking into account the fact that a user has the possibility to associate with different RRHs. Moreover, to reflect a real cellular scenario, our proposed deep learning models are trained with independent channel realizations, in which the users and RRHs have random locations in each realization.

## 2.9   Conclusions

Cloud radio access networks are expected to play a great role in the next generation cellular systems. They have a unique architecture in which BBUs are aggregated in a centralized BBU pool, enabling much lower power consumption, less costs, and more efficient optimization. In this chapter, we explained the CRAN architecture and how CoMP transmissions can highly enhance the network performance. Furthermore, we discussed the main optimization objectives and optimization parameters in CRAN. Several research works proposed techniques for CRAN performance optimization. Some of the proposed techniques can reach very high efficiencies and significant performance gains. However, there are still challenges that face improving the conventional performance optimization techniques and reducing their complexity. As noted, machine learning based techniques are highly qualified candidates that have negligible computational complexity.

# Chapter 3

# Joint User Association and Resource Allocation in CoMP-Enabled Heterogeneous CRAN

## 3.1 Introduction

As previously discussed, the exponential growth of the cellular traffic over the years has led the researchers to think about novel architectures that can handle this hard challenge. Hence, the CRAN was proposed. In typical cellular architectures, performing CoMP is considered a challenge, as high synchronization is required between the serving basestations (BSs). Nevertheless, in the CRAN architecture, synchronization is no more a severe issue, as the the centralized BBU pool controls numerous RRHs. Thus, synchronization can be performed smoothly [43].

Since CoMP transmission technique and CRAN architecture are well suited together, we combine both technologies in this chapter. Related works were reviewed in section 2.7. Different from the previous works, we propose using generalized CoMP transmissions to utilize any excess resources in the network to obtain higher weighted sum rate. A joint optimization problem is considered, where the optimization parameters are namely, user association, resource allocation, and power allocation. We tackle the user association and resource allocation sub-problems via matching

game based low complexity algorithms. Our user association algorithm proved to realize a good balance between cooperation gain and fronthaul consumption, such that a cooperation gain is achieved even in case of tight fronthaul constraints. Moreover, we propose a novel approach to convexify the power allocation sub-problem to obtain an optimal solution.

The rest of this chapter is organized as follows: in section 3.2, we discuss the system model and formulate our problem; our proposed algorithms are presented in section 3.3; the performance analysis and numerical results are investigated in section 3.4; and finally, our conclusions are drawn in section 3.5.

## 3.2   System Model and Problem Formulation

We consider the downlink transmission in a H-CRAN with the architecture shown in figure 3.1. Our system model includes a Macro remote radio head (MRRH) associated with some users denoted as MUs, and belong to the set $U^{MUs}$. Within the area of the MRRH coverage, several Pico remote radio heads (PRRHs) are deployed to serve a set of devices denoted as PUs, and belong to the set $U^{PUs}$. The MRRH and the PRRHs are assigned the same orthogonal radio resources from the set $\mathcal{N} = 1, 2, 3..N$. All the RRHs are connected to a baseband unit (BBU) pool via fronthaul links.

It can be noticed that the devices served by different PRRHs will suffer high interference from the MRRH and the non-serving PRRHs that use the same radio resources. Consequently, we implement the CoMP transmission technique to ensure that the QoS is satisfied for the PRRHs users (PUs), and to utilize the additional resources to optimize the network performance.

One of the great advantages of the CRAN architecture is that RRHs are connected to a central BBU pool, where computations can be done in a centralized manner; and hence, obtaining significant performance gain. Thus, in our model, the processing of all the proposed algorithms is performed at the BBU pool assuming perfect channel state information (CSI). The rate achieved by PU $i$ on radio resource $n$:

$$R_i^n = B^n \log_2(1 + SINR_i^n), \tag{3.1}$$

where $B^n$ is the bandwidth of radio resource $n$, $SINR_i^n$ is the signal to interference plus noise ratio
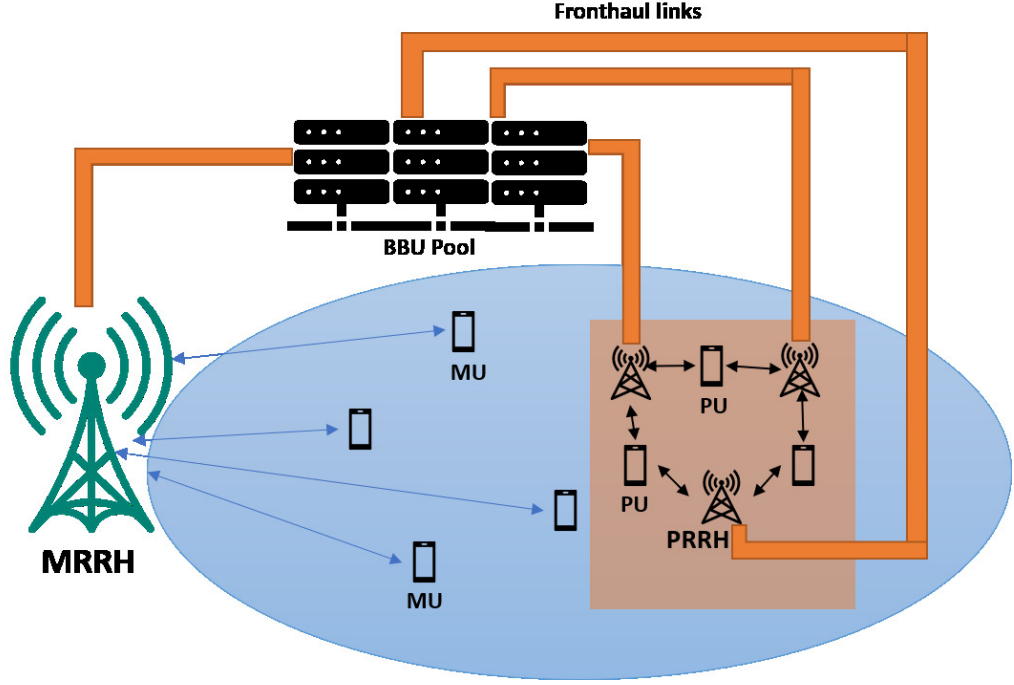
Figure 3.1: Our Network Architecture

received by PU $i$ at $n$,

$$SINR_i^n = \frac{\sum_{j \in \Upsilon} y_{ij} P_j^n g_{ji}^n}{P_M^n g_{Mi}^n + \sum_{k \in \rho_n^i} P_k^n g_{ki}^n + \sigma^2}, \tag{3.2}$$

where $i$, $j$ and $n$ are the indices for PUs, PRRHs and radio resources, respectively; $y_{ij}$ is the association coefficient, where $y_{ij} = 1$ if PU $i$ is served by PRRH $j$, and $y_{ij} = 0$ otherwise; $\Upsilon$ is the set of PRRHs; the index $M$ is used to denote the MRRH; $P_j^n$ and $P_M^n$ are the transmission powers of PRRH $j$ and the MRRH on radio resource $n$, respectively; $\rho_n^i$ is the set of PRRHs which use radio resource $n$ and do not serve PU $i$. Finally, $g_{ji}^n$, and $g_{Mi}^n$ are the channel gains on radio resource $n$ between PU $i$ and PRRH $j$ or the MRRH, respectively.

To guarantee the required QoS, a constraint must be defined. Accordingly, the data rate assigned for each PU must be greater than or equal to a predefined threshold $R_i^{min}$. Thus,

$$\sum_{n \in \aleph} \propto_i^n R_i^n \geq R_i^{min}, \forall i \in U^{PUs}, \tag{3.3}$$

where $\propto_i^n$ is the binary resource allocation coefficient. $\propto_i^n = 1$ only if the radio resource $n$ is assigned to PU $i$. Another constraint must be set in order to protect the MRRH users (MUs) from

the interference caused by the PRRHs. Consequently, the interference on the MU $m$ that is allocated radio resource $n$ must be less than a predefined threshold $I_{\mathrm{m}}^n$.

$$\sum_{j \in \Upsilon} P_j^n g_{jm}^n \leq I_{\mathrm{m}}^n, \ \forall n \in \aleph, \ m \in U^{MUs}. \tag{3.4}$$

Additionally, we set a constraint to ensure that the sum of the achievable rates of the PUs associated with each PRRH $j$ is less than the fronthaul capacity $Ca_j^{PRRH}$.

$$\sum_{n \in \aleph} \sum_{i \in U^{PUs}} \propto_i^n y_{ij} R_i^n \leq Ca_j^{PRRH}, \forall j \in \Upsilon. \tag{3.5}$$

Consequently, with an objective function to maximize the weighted sum rate of the network, our joint user association, resource allocation, and power allocation optimization problem can be formulated as follows:

$$\mathbf{P1}: \max_{(\mathbf{y}, \propto, \mathbf{P})} \cdot \sum_{j \in \Upsilon} \sum_{i \in U^{PUs}} \sum_{n \in \aleph} \propto_i^n R_i^n \tag{3.6}$$

$$\text{s.t.} \quad (3.3), (3.4), (3.5),$$

$$\sum_{n \in \aleph} \propto_i^n \leq 1, \forall i \in U^{PUs} \tag{3.7}$$

$$P_j^{\min} \leq \sum_{n \in \aleph} P_j^n \leq P_j^{\max}, \ \forall j \in \Upsilon \tag{3.8}$$

$$y_{ij} = \{0, 1\} \tag{3.9}$$

$$\propto_i^n = \{0, 1\} \tag{3.10}$$

Note that the constraint (3.7) indicates that a user can only be assigned one resource block (RB). Nevertheless, such a RB can be served by numerous PRRHs, according to the network status. On the other hand, constraint (3.8) limits the total power of PRRH $j$ to some maximum value. Finally, constraints (3.9), and (3.10) indicate that $y_{ij}$ and $\propto_i^n$ are binary coefficients. Thus we can easily conclude that $\mathbf{P1}$ is NP-Hard problem, which is computationally intractable. Therefore, we consider a sub-optimal solution by dividing $\mathbf{P1}$ to three sub-problems considering the user association $(\mathbf{P1 - UA})$, resource allocation $(\mathbf{P1 - RA})$, and power allocation $(\mathbf{P1 - PA})$, respectively.

$$\mathbf{P1-UA} \quad \max_{\mathbf{y}} . \quad \sum_{j \in \Upsilon} \sum_{i \in U^{PUs}} R_i^n \tag{3.11}$$

$$\text{s.t.} \quad (3.5), (3.9).$$

The first sub-problem will consider the user association using a many-to-many matching game [44] between PUs and PRRHs, while considering the fronthaul capacity constraint.

$$\mathbf{P1-RA} \quad \max_{\propto} . \quad \sum_{j \in \Upsilon} \sum_{i \in U^{PUs}} \sum_{n \in \aleph} \propto_i^n R_i^n \tag{3.12}$$

$$\text{s.t.} \quad (3.4), (3.7), (3.10).$$

The second sub-problem that considers the resource allocation is solved using a many-to-many matching game between PUs and radio resources, while taking the interference protection on MUs into account

$$\mathbf{P1-PA} \quad \max_{\mathbf{P}(\omega_{\mathbf{n}})} . \quad \sum_{i \in \omega_n} R_i^n. \tag{3.13}$$

$$\text{s.t.} \quad (3.3), (3.4), (3.5), (3.8).$$

The third sub-problem is the power allocation, where $\omega_n$ is the set of PUs using radio resource $n$ and their serving PRRHs. The power allocation sub-problem is in general non-convex. However, we will introduce an additional constraint to transform the problem to be a convex one. The additional constraint and the proof of convexity of the new problem will be discussed later on.

The fronthaul communications in CRAN architectures can be based on several technologies. For instance, wireless fronthaul networks were proposed based on microwave links, or WiFi standard in indoor environments [13]. Optical fiber based networks are always very efficient candidates due to their large capacities [13]. In many of these fronthaul network architectures, the transmission medium might be shared, especially among the RRHs serving at the same geographical area.

In what follows, it will be more convenient to reformulate the fronthaul capacity constraint as a

sum fronthaul constraint. Thus, our new constraint can be represented as follows:

$$\sum_{i \in U^{PUs}} \sum_{n \in \aleph} \sum_{j \in \Upsilon} \propto_i^n y_{ij} R_i^n \leq Ca^{Total}, \tag{3.14}$$

where $Ca^{Total}$ is the total capacity of the fronthaul network.

## 3.3 Proposed Algorithms

In CRAN, all RRHS are connected to a centralized BBU pool, where our algorithms are assumed to be implemented, with PRRHs and PUs being mapped to virtual nodes. In the rest of this, section we will discuss the matching game algorithms used to solve the user association and resource allocation sub-problems. Finally, we introduce our novel approach to convexify the power allocation sub-problem.

### 3.3.1 Many-to-many Matching Game Based User Association

We use a matching game based on the deferred acceptance scheme [44]. Algorithm 3.1 simply works as follows: each PU $i$ proposes to be matched to its preferred PRRH $j$, which initially accepts the proposal if there is enough capacity to serve $i$ in its fronthaul link, according to the minimum data rate needed for $i$, $R_i^{min}$. If there is not enough capacity, $j$ starts to sequentially reject the previously initially accepted PUs which are less preferred than $i$ until there is enough fronthaul capacity to admit $i$ or there are not other PUs to reject. If the latter case occurs, $i$ and all the rejected PUs will have to remove $j$ from their preference lists, and the same should be done by $j$. These steps should be repeated multiple times until convergence is reached. Additionally, the whole previously mentioned steps should also be repeated while updating the preference lists to ensure that the already matched PUs will not be considered again for matching with PRRHs they are associated with. The outputs are the user association sets $\Theta_j$, and $\Theta_i$.

To determine the preferred PRRHs for each PU, the utility of each PU $i$ with respect to each PRRH $j$ ($u_i^j$) and the utility of each PRRH $j$ with respect to each PU $i$ ($u_j^i$) must be calculated. The

result is applied as an input to the algorithm, where,

$$u_i^j = \log_2(1 + \sum_{n \in \aleph} SINR_{ij}^n) \tag{3.15}$$

$$u_j^i = \sum_{n \in \aleph} P_j^{max} g_{ji}^n. \tag{3.16}$$

Consequently, the preference list for each PU $i$ $(PL_i)$ is simply calculated by arranging the PRRHs in a descending order according to the values of $u_i^j$. The same process is done to calculate the preference list of each PRRH $j$ $(PL_j)$. Hence, each PU ranks PRRHs according to which will serve it with the highest rate, averaged on all radio resources. Also, each PRRH ranks PUs according to which will receive the highest power form it, averaged over all radio resources. During the user association and the radio resource allocation phases of our problem, the PRRHs will be virtually assumed to be operating with the maximum power $P_j^{max}$

When the fronthaul capacity is limited, we can approximately assume that the rates achieved by the PUs connected to a specific PRRH is tied by the capacity of its fronthaul. Thus, one can roughly approximate the initial rate that will be achieved by PU $i$ from PRRH $j$ as:

$$R_i^j = max(R_i^{min}, \frac{Ca_j^{PRRH}}{|\Theta_j|}) \tag{3.17}$$

where $|\Theta_j|$ is the number of PUs associated to PRRH $j$. Now assume that $i$ is associated to some PRRHs $j_1, j_2, j_3...etc$. The rate achieved by $i$ will be tied by the lowest rate it can achieve at $j_1, j_2, j_3...etc$, according to their fronthaul capacities, and the number of PUs associated to each. This rate can be initially approximated as:

$$R_i^{init} = min(R_i^{j_1}, R_i^{j_2}, R_i^{j_3}, ...). \tag{3.18}$$

Consequently, we can avoid cooperations that will actually lead to data rate loss based on the approximate equations (3.17), and (3.18). In algorithm 3.1, the term $Co$ is true only if a cooperation gain is expected. Thus, no PU will be associated with more than one PRRH, unless cooperation gain is guaranteed, or in the worst case, no performance degradation will occur.

**Algorithm 3.1:** Many to many matching user association

---

**Input:** $u_i^j$, $u_j^i$, $PL_i$, $PL_j$, $\forall j \in \Upsilon, i \in U^{PUs}$

1 **Initialize:** $t_1 = 0, t_2 = 0, \Theta_j(0) = \emptyset \,\forall j \in \Upsilon, Ca_j^{av} = Ca_j^{PRRH}$

2 **do**

3     $t_1 \leftarrow t_1 + 1$

4     **for** $j \in \Upsilon$ **do**

5       **for** $i \in U^{PUs}$ **do**

6         **if** $i \in \Theta_j(t_1)$ **then**

7           $u_i^j(t_1) \leftarrow 0, u_j^i(t_1) \leftarrow 0$

8         **else**

9           $u_i^j(t_1) \leftarrow u_i^j(t_1 - 1)\,, u_j^i(t_1) \leftarrow u_j^i(t_1 - 1)$

10     $PL_i(t_1) \leftarrow update(PL_i(t_1 - 1))$

11     $PL_j(t_1) \leftarrow update(PL_j(t_1 - 1))$

12     **do**

13       $\Psi_j \leftarrow \emptyset \,\forall j \in \Upsilon$

14       $t_2 \leftarrow t_2 + 1$

15       **for** $j \in \Upsilon$ **do**

16         **for** $i$ $with$ $j$ $as$ $its$ $most$ $preferred$ $in$ $PL_i$ **do**

17           **while** $i \notin \Psi_j$ **do**

18             **if** $Ca_j^{av} \geq R_i^{min}$ **then**

19               $\Psi_j(t_2) \leftarrow \Psi_j(t_2) \cup i, \;\; Ca_j^{av} \leftarrow Ca_j^{av} - R_i^{min}$

20             **else**

21               $PL_j'(t_2) \leftarrow \{i' \in \Psi_j(t_2) \,|i \succ_j i'\}$

22               Remove least preferred $i' \in PL_j'(t_2)$ from $\Psi_j(t_2)$ till $(PL_j'(t) = \emptyset)$
               Or $(Ca_j^{av} \leq R_i^{min})$

23               **if** $Ca_j^{av} \geq R_i^{min}$ **then**

24                do step 19

25               **else**

26                $D_{Lp} \leftarrow i$

27                $z_j \leftarrow \{z \in PL_j(t_2) \,|D_{Lp} \succ_j z\} \cup D_{Lp}$

28                **for** $z \in z_j$ **do**

29                  $PL_i(t_2) \leftarrow PL_i(t_2) \setminus \{j\}$

30                  $PL_j(t_2) \leftarrow PL_j(t_2) \setminus \{z\}$

31       **while** $\Psi_j(t_2) \neq \Psi_j(t_2 - 1), \,\forall j \in \Upsilon$

32       **for** $j \in \Upsilon$ **do**

33         **for** $i \in U^{PUs}$ **do**

34           **if** $(i \in \Psi_j(t_1)) \cap (i \notin \Theta_j(t_1), \,\forall j \in \Upsilon)$ **then**

35             $\Theta_j(t_1) \leftarrow \Theta_j(t_1) \cup \{i\}$

36           **else if** $(i \in \Psi_j(t_1)) \cap (i \notin \Theta_j(t_1)) \cap Co$ **then**

37             $\Theta_j(t_1) \leftarrow \Theta_j(t_1) \cup \{i\}$

38 **while** $\Theta_j^i(t_1) \neq \Theta_j^i(t_1 - 1)$

    **Output:** $\Theta_j, \Theta_i, \,\forall j \in \Upsilon, i \in U^{PUs}$

### 3.3.2 Many-to-one Matching Game Based Resource Allocation

Regarding the resource allocation (RA) algorithm, it will follow procedures similar to the user association (UA) algorithm that is already explained in details in Algorithm 3.1, and it is also based on deferred acceptance [44]. The utilities that will be input to the algorithm are:

$$u_n^i = \sum_{j \in \Theta_i^j} P_j^{max} g_{jm}^n \tag{3.19}$$

$$u_i^n = \log_2(1 + \sum_{j \in \Theta_i^j} SINR_{ij}^n) \tag{3.20}$$

$$u_n^j = P_j^{max} g_{jm}^n \tag{3.21}$$

$$u_i^{(n)(j)} = \log_2(1 + SINR_{ij}^n) \tag{3.22}$$

where $u_n^i$ is the utility of each $n$ with respect to each PU $i$ which is equal to the interference caused on MU $m$ (allocated radio resource $n$) by the PRRHs associated with $i$. $u_i^n$ is, similarly, the utility of each PU $i$ with respect to each radio resource $n$ and is equal to the rate achieved on $n$ using CoMP transmissions from the associated PRRHs. $u_n^j$ is the interference caused on MU $m$ by PRRH $j$. $u_i^{(n)(j)}$ is the rate achieved by PU $i$ when allocated radio resource n and associated with PRRH $j$.

In addition to the utilities, $PL_n$, $PL_i$, $\Theta_i$ , $\Theta_j$ should also be input to the resource allocation (RA) algorithm, in which $PL_n$, $PL_i$ are the preference lists that can be obtained as explained before, and $\Theta_i$, $\Theta_j$ are the user association sets obtained from Algorithm 3.1.

Now, we illustrate the operation of the RA algorithm. Firstly, matching is done between PUs and radio resources based on the interference caused by the PRRHs associated with $i$, according to the predefined interference threshold $I_m^n$. If a PU is initially accepted on a radio resource $n$, it competes with the PUs that were previously initially accepted on the same $n$ and associated with the same PRRHs. The PU that achieves higher rate from each PRRH will be associated with it on radio resource $n$. If a PU losses all its associated PRRHs on a specific radio resource $n$, the radio resource will be removed from its preference list, and this PU will propose its second preferred $n$ in the subsequent cycle. Additionally, all the utilities and preference lists should be updated after any

change in the user association set. These steps will be repeated until convergence. The output from the RA algorithm will be the set $\kappa_n$ containing the PUs assigned each radio resource $n$, and their serving PRRHs.

### 3.3.3 Power Allocation (PA) Algorithm

Given that the set $\kappa_n$ was obtained by applying the user association and resource allocation algorithms, the power allocation problem can now be solved. We can write the objective function of our power allocation sub-problem as:

$$
\sum_{n \in \aleph} \sum_{i \in \kappa_n} R_i^n = \sum_{n \in \aleph} \sum_{i \in \kappa_n} B^n \log_2(1 + SINR_i^n) =
$$
$$
\sum_{n \in \aleph} \sum_{i \in \kappa_n} B^n \log_2(1 + \frac{\sum_{j \in \tau_i} P_j^n g_{ji}^n}{P_M^n g_{Mi}^n + \sum_{j \notin \tau_i} P_j^n g_{ji}^n + \sigma^2}) \tag{3.23}
$$

where $\tau_i$ is the set of PRRHs serving PU $i$. To transform $(\mathbf{P1 - PA})$ to a convex problem, we introduce an additional constraint. Thus, the resultant modified problem $(\mathbf{P1 - PA})$ can be stated as follows:

$$
\mathbf{P1 - PA} \quad \max_{(\mathbf{P}(\kappa_{\mathbf{n}}))}. \quad \sum_{n \in \aleph} \sum_{i \in \kappa_n} R_i^n \tag{3.24}
$$

$$
\text{s.t.} \quad (3.3), (3.4), (3.8),
$$

$$
P_M^n g_{Mi}^n + \sum_{j \notin \tau_i} P_j^n g_{ji}^n + \sigma^2 \leq \sum_{j \in \tau_i} P_j^n g_{ji}^n \quad \forall \, i \in \kappa_n, n \in \aleph. \tag{3.25}
$$

Constraint (3.25) is simply stating that $SINR \geq 1$ for all the associated PUs. The proof of the convexity of $(\mathbf{P1 - PA})$ is provided in appendix A. Since our new problem is convex, it can be solved with any of the well known convex optimization tools to obtain the optimal solution.

For the fronthaul capacity constraint in (3.5) to be satisfied, after applying the UA, RA, and PA algorithms, each PRRH checks the total rate of its associated PUs. Then PRRHs start to reject the ones with the smallest achievable rates until the constraint is satisfied.

It is important to note that, for the special cases when the fronthaul links are shared among the RRHs or the fronthaul capacity is unlimited, the same algorithms and steps are employed to

solve the problem. However, in the user association algorithm, PUs can freely be associated with numerous PRRHs, as using CoMP will not result in overloaded fronthaul links. In contrast, it will help achieving much higher network throughput.

## 3.4   Simulation Results

Our simulation model considers six PRRHs of a maximum transmission power of $20dBm$ placed inside the coverage area of a single MRRH with a constant transmission power of $46dBm$. The PRRHs and PUs are uniformly distributed at an indoor area of $300m^2$. The MUs are placed outdoors directly beside the $300m^2$ indoor area, in order to have the worst case inter-tier interference scenario. The MRRH and the PRRHs use the same orthogonal six radio resources, each of which is $180KHz$ bandwidth. We assume that each radio resource is already allocated to a MRRH user (MU), and that the interference threshold of each of the six MUs is $-100dBm$. The noise power spectral density (PSD) is $-174\ dBm/Hz$. The wireless channel follows a Rayleigh fading model, with the path-loss and shadowing models implemented as [45]. The distance between the MRRH and the indoor area is $600m$.

In Fig. 3.2, we consider that our network has fronthaul links with unlimited capacity. Thus, constraints (3.5) or (3.14) are not taken into account. To assess the performance of our proposed algorithms, we compare the performance of our network while implementing three different settings:

- Our proposed generalized CoMP with matching game UA and RA algorithms

- No CoMP, user association is done with the high SINR algorithm, and resource allocation with greedy algorithm.

- No CoMP, with matching game UA and RA algorithms [11]

In the three settings, the optimal power allocation is obtained.

We can see that significantly higher throughput could be obtained with CoMP, as the excess resources are utilized to achieve much higher rate for the connected PUs. The network performance with CoMP tends to become closer to the non-CoMP algorithms as the number of the served PUs increases, due to the decrease in the available excess resources. Thus, when the network is overloaded,
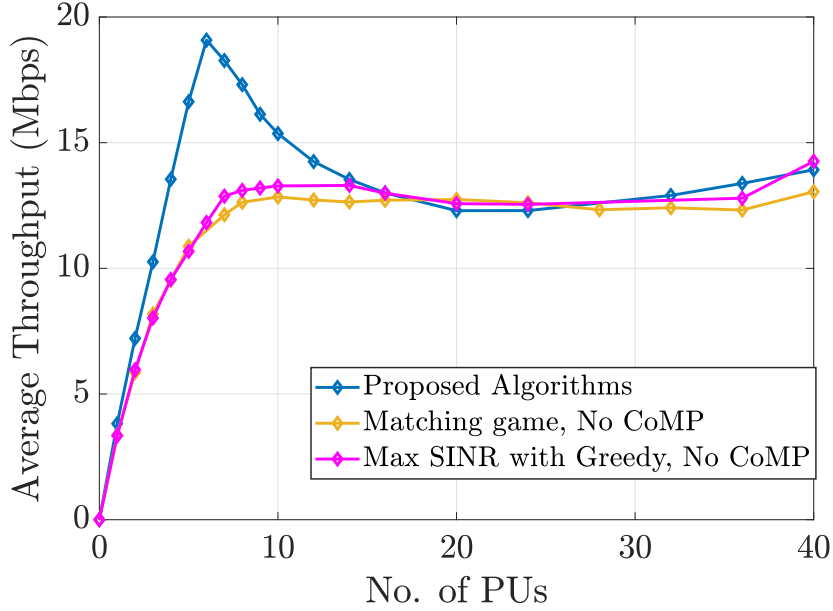
Figure 3.2: Weighted sum rate of the network versus number of PUs- large fronthaul capacity

our generalized CoMP algorithm is equivalent to non-CoMP algorithms.

To assess the performance of our network in the case of limited capacity fronthaul, we consider both cases of either individual or shared fronthaul links. The capacity of the shared fronthaul $Ca^{Total} = 9Mbps$. On the other hand, the capacity of each individual fronthaul link $Ca_j^{PRRH} = 1.5Mbps$. Fig. 3.3 represents the weighted sum rate of the network versus the number of served PUs. We can see that our proposed CoMP algorithms can achieve considerable gains even in the case of tight individual fronthaul constraints. This is due to the fact that our user association algorithm can create a good tradeoff between fronthaul consumption and cooperation gain.

To further assess the performance of our algorithms in more random environments, our parameters will be changed such that the PRRHs will be Poisson distributed with density $\lambda_{PRRH} = 6PRRHs$, and the PUs will also follow the same distribution with $\lambda_{PU} = 16PUs$. The efficiency of our generalized CoMP algorithm can be more realized from Fig. 3.4, in which the weighted sum rate of the network is plotted against the number of radio resources available for each PRRH. It can be seen that the performance of our CoMP algorithm is close to the non-CoMP algorithms when the number of available resources is small. However, as the number of available resources, and accordingly, excess resources increase, the throughput of the network becomes much higher when
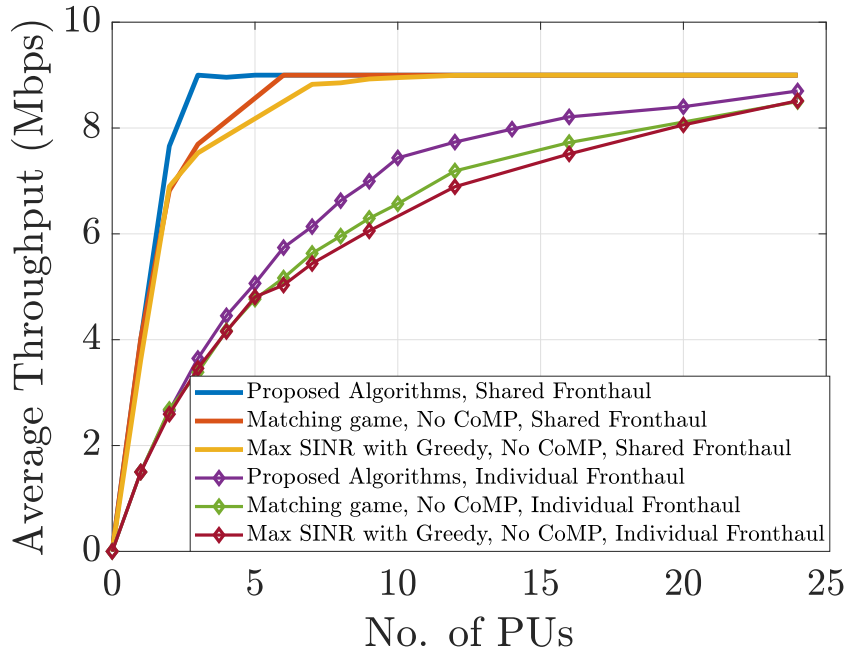
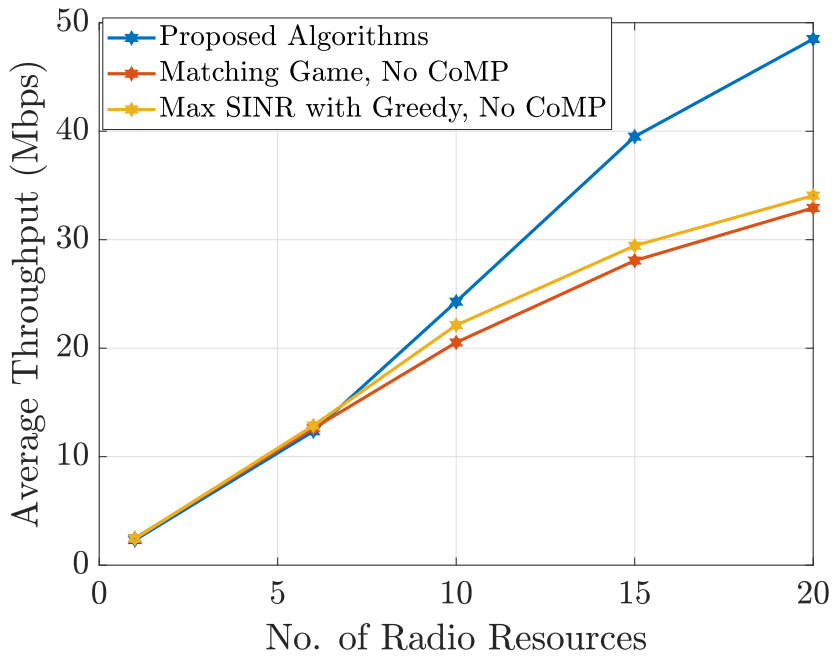Figure 3.3: Weighted sum rate of the network versus number of PUs- limited fronthaul capacity



Figure 3.4: Weighted sum rate of the network versus number of radio resources- large fronthaul capacity

employing the proposed algorithms.

## 3.5 Conclusions

We proposed generalized CoMP in order to utilize any excess resources in the network to improve the throughput. The simulation results proved the significant improvements in the network performance when utilizing our proposed algorithms. The performance of our generalized CoMP scheme becomes more superior when excess radio resources and high capacity fronthaul links are available in the network. Moreover, our proposed user association algorithm proved to achieve cooperation gains even with very tight fronthaul constraints.

# Chapter 4

# Unsupervised Deep Learning Approach for Near Optimal Power Allocation in CRAN

## 4.1 Introduction

Most of the research works targeting CRAN performance optimization utilize the conventional optimization techniques. The joint problem of user association to each RRH, BBU, and fronthaul link, in addition to power allocation is studied in [46]. To solve their optimization problem, they propose a two level iterative algorithm, in which at the first level the problem is solved two times with one of the variables fixed at each, while at the lower level the problem is divided into many sub-problems and solved via the successive convex approximation technique. Their algorithms could achieve an enhancement in the network throughput, but with high computational complexity. Regarding the energy efficiency and the green CRAN concept, authors in [26, 47] considered minimizing the power consumption in CRANs, while maintaining some constraints such as the QoS. Nevertheless, to tackle their problems, they used online optimization techniques with relatively high computational complexities, such as the weighted minimum mean square error (WMMSE) and Lagrangian relaxation. Another research in [48] considered jointly optimizing user association and

precoding in CRAN. To deal with the NP-hardness of their problem, they utilized the successive convex approximation technique to approximate the non-convex constraints in their problem. However, since their algorithm is based on solving a convex problem multiple times, it suffers from high computational complexity.

To tackle the problem of the high computational complexity of the conventional optimization algorithms, in this chapter, we propose a deep learning based power allocation algorithm. Particularly, we consider a joint optimization problem of user association and power allocation in CRAN downlink transmission, with an objective to maximize the network weighted sum rate, while maintaining the users' QoS constraints. We divide our joint problem into two sub-problems. The first is the user association, which will be tackled through a low complexity matching game based algorithm [11, 49], while the second sub-problem is the power allocation which will be solved through a deep neural network (DNN) based algorithm. The contributions of this chapter can be stated as follows:

- we optimize the power allocation in a CRAN with DNNs, while taking user association into account to reflect a real cellular scenario, where users have the possibility to associate with several RRHs, and the locations of the users are highly random.

- In this work, we consider the unsupervised learning approach for power allocation employed in [42]. Also, we add the QoS constraints to the DNN's loss function through a $ReLU$ penalty term as in [42, 50]. However, there is a trade-off between the ability of the DNN to obtain high data rates, and its ability to maintain the QoS of the users [42]. There is a hyperparameter that controls this trade-off. We call this parameter the QoS coefficient. In this work, we provide intensive simulations to show how the QoS coefficient can affect the DNN performance significantly. Moreover, we define a novel performance metric to measure the QoS preservation capability of the DNNs based scheme.

- Furthermore, we propose a novel approach to enhance the ability of the DNNs to obtain higher data rates with better QoS preservation capability, through directly inputting the QoS requirements of the users to the DNN.

The rest of this chapter is organized as follows: in section 4.2, we present our system model

35

and formulate the optimization problem; section 4.3 presents our proposed algorithms; numerical analysis is presented in section 4.4; finally, we draw our conclusions in section 4.5.

## 4.2   System Model and Problem Formulation

### 4.2.1   System Model

A downlink scenario in a CRAN is considered. Our system architecture (Fig.4.1) consists of $M$ RRHs and $R$ users, in which RRHs communicate with their associated users utilizing the same radio resource block. Hence, high interference is anticipated and efficient optimization algorithms are required to achieve acceptable performance. The RRHs are connected to the centralized BBU pool via high capacity fronthaul links. In our system, we consider the single association scheme, where each user can be served by maximally one RRH. The RRHs belong to the set $J = 1, 2, 3...M$ and denoted by $j$, while the users belong to the set $U = 1, 2, 3...R$ and denoted by $i$. To cohere with the CRAN architecture, which enables efficient centralized optimization, we assume that our optimization algorithms are implemented in the BBU pool, and that perfect channel state information (CSI) is available.
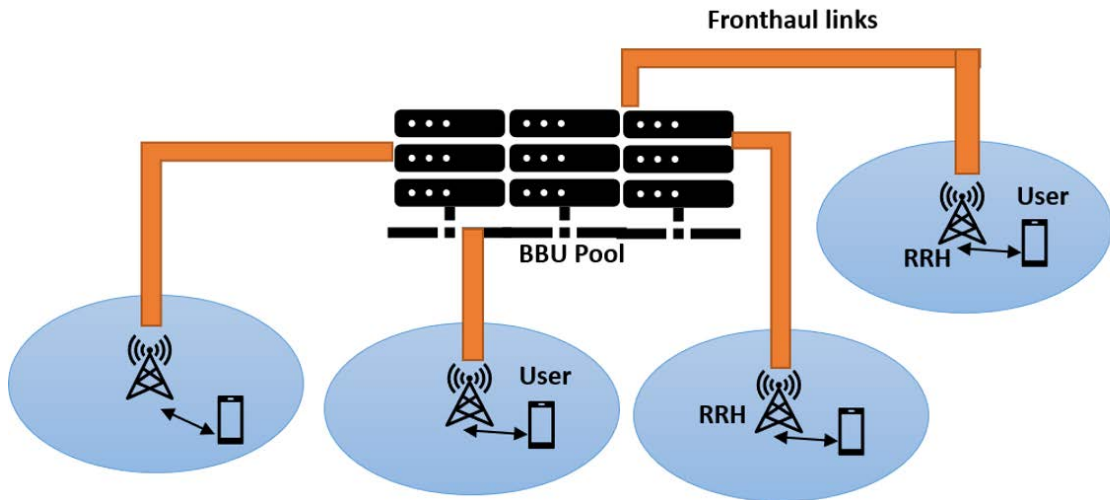


Figure 4.1: Network Architecture

### 4.2.2 Problem Formulation

Recall that the data rate achieved by user $i$ associated to RRH $j$ can be calculated through (2.3) and (2.1). In order to achieve the QoS requirements of each user, a constraint must be defined, similar to our approach in chapter 3:

$$\sum_{j \in J} y_i^j R_i^j \geq R_i^{min}, \forall i \in U, \tag{4.1}$$

where $R_i^{min}$ is the minimum required data rate for user $i$, and $y_i^j$ is the user association coefficient. $y_i^j = 1$ if the user $i$ is associated to RRH $j$, and $y_i^j = 0$ otherwise. Thus, (4.1) states that each user must obtain at least its minimum required data rate ($R_i^{min}$).

Considering user association and power allocation, our weighted sum rate maximization problem can be formulated as follows:

$$\mathbf{P2}: \quad \max_{(\mathbf{y},\mathbf{P})} . \quad \sum_{j \in J} \sum_{i \in U} y_i^j R_i^j \tag{4.2}$$

$$\text{s.t.} \quad (4.1),$$

$$\sum_{j \in J} y_i^j \leq 1, \forall i \in U \tag{4.3}$$

$$P_j^{\min} \leq P_j \leq P_j^{\max}, \forall j \in J \tag{4.4}$$

$$y_i^j = \{0, 1\} \forall i \in U, j \in J \tag{4.5}$$

The constraint (4.3) restricts the system to the single association scheme. On the other hand, constraint (4.4) limits the transmission power of RRH $j$ to some maximum value $P_j^{\max}$. Lastly, constraint (4.5) restricts the association coefficient $y_i^j$ to binary values. Accordingly, it can be observed that $\mathbf{P2}$ is a mixed integer programming NP-Hard problem that can be tackled optimally, solely, via the exhaustive search technique which has exponential computational complexity. Therefore, we consider sub-optimal solutions. To tackle our problem $\mathbf{P2}$, we firstly divide it into two sub-problems namely, $(\mathbf{P2} - \mathbf{UA})$ which considers the user association, and $(\mathbf{P2} - \mathbf{PA})$ considering the power allocation.

### 4.2.3 Sub-problem Formulation

The first sub-problem considers the user association, and will be tackled using the matching game algorithm [49], while restricting the system to the single association scheme.

$$\mathbf{P2 - UA}: \quad \max_{\mathbf{y}}. \quad \sum_{j \in J} \sum_{i \in U} y_i^j \, R_i^j \tag{4.6}$$

$$\text{s.t.} \quad (4.1), (4.3), (4.5).$$

The second sub-problem considers the power allocation. The set $\theta$ represents the active RRHs, while $\beta_j$ is the user associated with RRH $j$. Generally, $\mathbf{P2 - PA}$ is a non-convex optimization problem that is usually tackled in the literature through high complexity techniques. Nevertheless, we utilize DNN models, that are trained offline, to predict the power profile $P(\theta)$. The online computational complexity of DNNs is negligible compared to the other optimization algorithms. In the upcoming sections, we will discuss thoroughly how to train the DNNs, such that they reach performance that can be superior to the high complexity optimization algorithms. Furthermore, we provide an extensive review of the trade-offs faced during DNNs training, and how to tune the different parameters according to our desired performance.

$$\mathbf{P2 - PA}: \quad \max_{\mathbf{P}(\theta)}. \quad \sum_{j \in \theta} R_{\beta_j}^j \tag{4.7}$$

$$\text{s.t.} \quad (4.1), (4.4).$$

## 4.3 Proposed Algorithms

### 4.3.1 Matching Game Based User Association

As previously mentioned, the user association sub-problem is undertaken with the matching game algorithm. The objective is simply to match each user with the RRH expected to provide the highest data rate. Our matching game algorithm is presented in Algorithm 4.1. Note that Algorithm 4.1 is similar to Algorithm 3.1, but customized to match our system model presented in this chapter. The inputs to the algorithm are the preference lists $(L_j)$, $(L_i)$ of the RRHs and users, respectively.

The preference list of each user represents the RRHs providing the highest $SINR$ arranged from the best to the worst, while the preference list of each RRH also represents the users receiving the highest $SINR$ from that RRH arranged in a descending order. Now, we illustrate briefly the operation of the matching game based user association. In each iteration, every user $i$ proposes to be matched with its best preferred RRH $j$. The proposal of user $i$ is initially accepted if no users are associated with RRH $j$, or if the user $i$ precedes in the preference list $L_j$ the user that is associated with $j$ at the current iteration. If none of the aforementioned conditions is true, $j$ will reject user $i$, which will have to remove $j$ from its preference list $L_i$. Then, user $i$ will propose its second preferred RRH in the subsequent iteration, in which the previous steps will be repeated. The operation continues for all users and RRHs until convergence is reached. The outputs of the matching game algorithm are the sets $\theta = \{\theta_1, \theta_2, \theta_3, ...\theta_K\}$ and $\beta = \{\beta_{\theta_1}, \beta_{\theta_2}, \beta_{\theta_3}, ...\beta_{\theta_K}\}$ representing the active RRHs, and the user associated with each of them, respectively. Accordingly, we will have $K$ RRH-user pairs.

---

**Algorithm 4.1:** Matching game based user association

---

    **Input:** $L_i,\ L_j,\ \ \forall j \in J, i \in U$

1  **Initialize:** $t = 0$, $\beta_j(0) = \emptyset\ \forall j \in J$

2  **do**

3      $t \leftarrow t + 1$

4      **for** $j \in J$ **do**

5         $i' \leftarrow \beta_j(t)$

6         **for** $i\ with\ j\ as\ its\ most\ preferred\ in\ L_i$ **do**

7            **if** $(i' = \emptyset) \cup (i\ precedes\ i'\ in\ L_j)$ **then**

8               $\beta_j(t) \leftarrow i$

9            **else**

10               $L_i(t) \leftarrow L_i(t) \setminus \{j\}$

11              $L_j(t) \leftarrow L_j(t) \setminus \{i\}$

12 **while** $\beta(t) \neq \beta(t-1)$

    **Output:** $\theta,\ \beta$

---

Given the sets $\theta$ and $\beta$ have been obtained, the power profile can be calculated with a trained DNN model. The input vector to the DNN can be expressed as:

$$V_{DNN} = [g_{\theta_1 \beta_{\theta_1}}, g_{\theta_2 \beta_{\theta_1}}, ....g_{\theta_K \beta_{\theta_1}}, g_{\theta_1 \beta_{\theta_2}}, g_{\theta_2 \beta_{\theta_2}}, ....g_{\theta_K \beta_{\theta_2}}, ....g_{\theta_1 \beta_{\theta_K}}, g_{\theta_2 \beta_{\theta_K}}, ....g_{\theta_K \beta_{\theta_K}},$$

$$R_{\beta_{\theta_1}}^{min}, R_{\beta_{\theta_2}}^{min}, ....R_{\beta_{\theta_K}}^{min}]^T \quad (4.8)$$

where $g_{\theta_1 \beta_{\theta_1}}$ is the main channel gain from RRH $\theta_1$ to user $\beta_{\theta_1}$, $g_{\theta_2 \beta_{\theta_1}}$ is the interference channel gain from RRH $\theta_2$ to user $\beta_{\theta_1}$, and so forth. Additionally, we adopt a novel approach to input the QoS requirements of each user $R_i^{min}$ to the DNN, in which we will prove later on that it enhances the DNNs' ability to maintain the QoS constraints and achieve higher data rates, significantly.
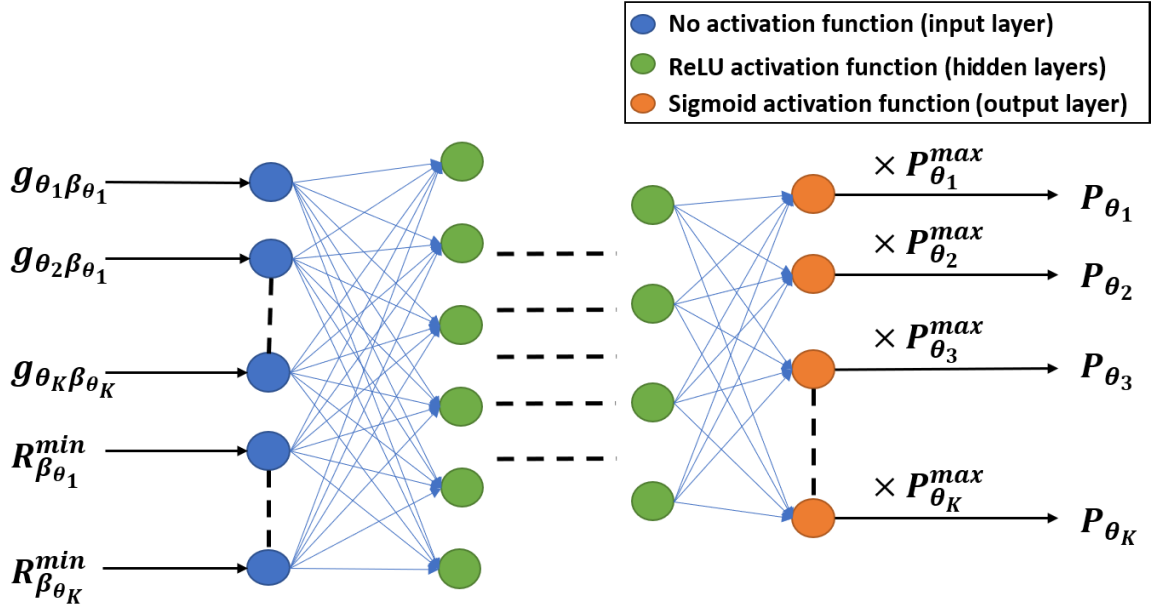
### 4.3.2 DNNs Architecture



Figure 4.2: The structure of DNN

We consider the conventional fully connected deep neural networks. Our proposed DNN architecture (Fig. 4.2) consists of an input layer, $(D-2)$ hidden layers, and an output layer. The input layer consists of $(K^2 + K)$ nodes, such that the first $K^2$ nodes correspond to the channel coefficients, while the other $K$ nodes correspond to the QoS minimum requirements for the users $(R_i^{min})$. The output layer consists of $K$ nodes corresponding to the normalized power profile. Our DNN layers are denoted by $x = \{1, 2, 3, 4...D\}$, where layer 1 is the input layer, and layer $D$ is the output layer. The number of nodes in each layer is denoted by $n_x$, such that $n_1 = K^2 + K$.

We adopt the $ReLU(.)$ as the activation function for the hidden layers of our DNNs. Activation functions are critical as they introduce non-linearity to the deep learning models [51], otherwise the resultant DNN models will be equivalent to linear regression. Hence, the output of an arbitrary

hidden layer $x$ can be calculated as:

$$B_x = ReLU(W_x B_{x-1} + b_x), \tag{4.9}$$

where $ReLU(u) = max(0, u)$; $B_x$ is the output vector of layer $x$ and has a size of $n_x \times 1$, such that each element in $B_x$ is the output of its corresponding node; $W_x$ is weights matrix with dimensions $n_x \times n_{x-1}$; and $b_x$ is the $n_x \times 1$ bias vector.

Considering that the output of the $ReLU(.)$ activation is unbounded, while the transmission power of RRH $j$ must be confined to a maximum value $P_j^{max}$, we utilize the sigmoid activation function for the output layer. Consequently, the DNN output vector can be expressed as:

$$B_D = Sig(W_D B_{D-1} + b_D), \tag{4.10}$$

in which:

$$Sig(u) = \frac{1}{1 + e^{-u}}. \tag{4.11}$$

Thus, the elements of the output vector of our DNN $(B_D)$ are always in the range $[0, 1]$. Finally, we can calculate the allocated powers vector $P = [P_{\theta_1}, P_{\theta_2}, P_{\theta_3}, ....P_{\theta_K}]^T$ as:

$$P = P^{max} \odot B_D, \tag{4.12}$$

in which $P^{max} = [P_{\theta_1}^{max}, P_{\theta_2}^{max}, ....P_{\theta_K}^{max}]^T$ represents the maximum powers of the active RRHs. Hence, (4.12) is simply the element-wise product between $P^{max}$ and $B_D$.

Before the formed DNN can be used for prediction, it must be trained with a sufficient number of data points in order to obtain high performance. This is discussed in details in the upcoming subsection.

### 4.3.3 Training DNNs

The training process of DNNs is crucial to obtain high performance. Previous research works in [52] and [53] utilized supervised training to train their DNNs. The main idea of the supervised

training is to generate labeled data, which acts as the ground truth, via one of the well known online optimization algorithms such as WMMSE, as adopted in [52]. Then, the DNN training process is done on the basis of minimizing an error function between the predicted power profile and the data labels. However, this approach has several limitations. Considering that DNNs usually require very large data sets to be trained accurately, generating such data sets of labeled training data will be an expensive process in terms of computational complexity, required hardware, and required time to accomplish the process. Moreover, the trained DNNs cannot outperform the optimization algorithm used to generate the ground truth data labels. Thus, we opt for the unsupervised learning. For the unsupervised training, we directly utilize the objective function in problem **P2-PA** as the loss function of our DNN, in addition to embedding the QoS constraints (4.1) in that loss function via ReLU penalty terms, as in [42] and [54]. To this point, we can express our loss function as:

$$Lo_{DNN} = \mathbb{E}\Big[ - R(V_{DNN}, W, b) + q \sum_{i \in \beta} ReLU(R_i^{min} - R_i(V_{DNN}, W, b)) \Big], \qquad (4.13)$$

where $q$ is the QoS coefficient; $R(V_{DNN}, W, b)$ and $R_i(V_{DNN}, W, b)$ are the weighted sum rate of the network, and the data rate achieved by user $i$, respectively. The terms $R(V_{DNN}, W, b)$ and $R_i(V_{DNN}, W, b)$ are functions of the DNN's weights matrix ($W$), biases vector ($b$), and vector $V_{DNN}$ representing a specific channel realization along with the QoS requirements of the users. Thus, each ReLU penalty term will have a positive value only if the QoS requirement of the corresponding user $i$ is violated, enforcing the training process towards satisfying the given requirement. Otherwise, the value of the ReLU term will be zero and will not affect the DNN training process. To be noted that the parameter $q$ can be considered as a critical hyperparameter that needs to be tuned according to our target performance. In general, there is no analytical method to calculate the optimal values of the hyperparameters related to the DNNs' training process. Therefore, we provide intensive experimental results to discuss how to tune $q$ according to our desired performance. Broadly speaking, when the value of $q$ is chosen to be very high, the DNN training process will be highly biased towards maintaining QoS constraints, rather than maximizing the weighted sum rate. On the other hand, if $q$ is set to have a very small value, the opposite occurs and the training process will have much higher tendency to increase the sum rate on the cost of violating some users' QoS

constraints. Thus, the manner in which $q$ is tuned provides a tradeoff between the total achieved sum rate and QoS satisfaction for users.

It is to be noted that the loss function $Lo_{DNN}$ in (4.13) represents an expected value that is related to the distribution of the elements in the DNN input vector $V_{DNN}$ in (4.8). Thus, it is quite challenging to directly compute $Lo_{DNN}$. We adopt the widely used approach, mini-batch stochastic gradient descent [55]. Hence, the generated offline training data will be divided into numerous equal sized mini-batches. Assume that the channel realizations related to a specific mini-batch belong to the set $\psi$. Thus, the number of $V_{DNN}$ vectors in each mini-batch is equal to $|\psi|$, which is called the batch size. The loss function in (4.13) can be approximated as:

$$Lo_{DNN} \approx \frac{1}{|\psi|} \sum_{V_{DNN} \in \psi} \Big[ - R(V_{DNN}, W, b) + q \sum_{i \in \beta} ReLU(R_i^{min} - R_i(V_{DNN}, W, b)) \Big]. \quad (4.14)$$

In each training step, the approximate (4.14) is used instead of (4.13) to calculate the loss function for each mini-batch, and then update the DNN's weights and biases. One training epoch is completed after a training step is done on each of the formed mini batches. The batch size is an important hyperparameter for the DNNs' training process. In deep learning, a common practice is to choose small batch sizes, as this improves the generalization capability of the DNN models [56]. However, choosing a very small batch size can result in a dramatic increase in the DNNs needed training time due to the loss of vectorization effect [57]. Thus, the choice of the batch size should balance the complexity and accuracy.

It is important to note: there is no guarantee that the power profiles predicted through the DNN will always be feasible, regardless of the value of $q$. This is due to the fact that the DNN mainly makes predictions according to the input channel realizations and QoS requirements, unlike the other optimization algorithms which solve the problem online. Moreover, for some channel realizations and QoS requirements, the problem may initially be infeasible. For instance, in case some users have relatively high QoS requirements that cannot be satisfied because of the interference level.

To investigate more the effect of the way $q$ is tuned on the DNNs performance, we introduce a

new performance metric:

$$P_{QV} = Pr[QoS\ is\ Violated].\qquad(4.15)$$

It can be easily deduced that $P_{QV}$ is the probability that the QoS requirements of some user in the network will be violated. Unfortunately, we currently do not have an analytical approach to calculate $P_{QV}$. However, in our simulations, we provide intensive analysis on how the method $q$ is tuned can affect $P_{QV}$. Moreover, we prove that our novel approach to input the QoS requirements to the DNN can achieve high data rates with significantly reduced values of $P_{QV}$, regardless of the value of $q$.

After the DNN models are trained offline, according to the target performance, they can be effectively used for online prediction. Our complete online scheme is presented in Algorithm 4.2.

Regarding the training process, several thousands of channel realizations and QoS requirements should be generated offline. Then steps from 1 to 4 in Algorithm 4.2 should be used to calculate $V_{DNN}$ corresponding to each of the generated realizations. Finally, the calculated vectors should be used for the offline training process.

---

**Algorithm 4.2:** The complete online algorithm

---

**Input:** $U$, $J$, $R_i^{min}\ \forall i \in U$, and $g_{ji}\ \forall j \in J,\ i \in U$

1 Calculate $SINR_i^j\ \forall j \in J,\ i \in U$
2 Calculate $L_i,\ L_j,\ \forall j \in J, i \in U$
3 Run Algorithm 4.1 to get the outputs $\theta$ and $\beta$
4 Calculate $V_{DNN}$
5 Load the appropriate trained DNN model, according to the value of $K$
6 Use the loaded model to calculate the normalized powers vector $B_D$
7 Calculate the power profile vector $P$ from (4.12)

**Output:** $P$

---

## 4.4   Numerical Analysis

In this section, we present our numerical results to verify the great enhancement in the network performance when employing the proposed DNNs for power allocation. Our simulations related to training DNNs and validating their performance were implemented with Python programming language, utilizing TensorFlow, the open source library. The rest of our simulations related to the

other optimization techniques, we compare against, were accomplished via MATLAB. Without loss of generality, all the trained models in our simulations had $8$ hidden layers, where each hidden layer was composed of $n_x = 100\ nodes$. To obtain DNN models with high performance for each scenario, several hundreds of training epochs were run, and the best performing models on the validation data, according to our target performance were saved. Adam optimizer with initial learning rate of $\eta = 0.001$ was used in the training due to its well known high performance [58]. For each model, we used $150,000$ data points for training, and $50,000$ for validation. The data points were divided into equal sized mini-batches. To provide a good balance between the training time and complexity, we chose $batch\ size = 100$ for all our trained models. The reason behind using a large number of data points for validation is to verify that our trained models can perfectly generalize over any random inputs they have not experienced. All the values provided in our numerical analysis are based on the performance of the DNN models on the validation data.

Independent channel realizations and QoS requirements were used to generate each data vector $V_{DNN}$, for both the training and validation data. The users and RRH were assumed to be placed randomly with a uniform distribution at an area of $600m^2$. The QoS requirements of the users were randomly selected from a uniform distribution in the range $[0, 300]$ Kbps. The Rayleigh fading channel model was assumed, with a path-loss and shadowing models as deployed in [11], [45], according to the 3GPP standards. A $600KHz$ radio resource was utilized in our simulations. The same settings were used for the schemes we compare against.

Note that the matching game based user association in Algorithm 4.1 was applied before our proposed DNNs based power allocation scheme, and before the power allocation schemes we compare against. Hence, we will have: number of RRHs=number of users=number of RRH-user pairs=$K$.

To validate the efficiency of our power allocation algorithm based on DNNs, we compare against three different power allocation schemes:

- genetic algorithm (GA),

- global search (optimal power allocation),

- the conventional deep learning based scheme [42, 50, 59, 60].

The idea of the global search method is to solve the power allocation problem numerous times

through one of the convex optimization techniques, where each time the algorithm uses a different initial point. Following that, the best achieved solution is selected. The number of examined initial points is chosen to be sufficiently large such that an optimal solution is guaranteed, with a high computational complexity. The genetic algorithm is a random based heuristic algorithm [24]. It belongs to a class named "evolutionary algorithms". The idea of the genetic algorithm is to generate numerous feasible solutions for the problem being tackled each iteration, in which this set of points is called "population". From the produced population, higher quality solutions can be reproduced utilizing some biologically inspired operations, replacing the previous population. Each population of solutions can be called a generation. These generations continue evolving till an acceptable solution is produced. The genetic algorithm can reach high quality solutions for optimization problems, but an optimal solution is not guaranteed. There are two parameters that play a key role in the performance of GA, the population size ($PS$) and the maximum number of generations. The population size is simply the number of solutions produced each generation. The algorithm stops after the maximum number of generations is reached. Increasing either of the population size or the maximum number of generations will result in higher quality solutions that can be close to optimal when their values are sufficiently high. However, increasing the values of the aforementioned parameters will also result in much higher computational complexity. In our simulations, we used $max\ number\ of\ generations = 100 \times K$, and different values of $PS$ for comparison. Regarding the conventional deep learning based scheme, it is implemented with almost the same settings as our proposed algorithm. The only difference is that for the conventional scheme, the QoS requirements of the users are only considered in the loss function but not fed to the DNN's architecture.

### 4.4.1 Achieved Data Rate Performance

Firstly, we compare the performance of our proposed DNNs based power allocation algorithm against the GA and the global search method. For the DNNs based algorithm, we consider three cases:

- high QoS preservation

- moderate QoS preservation

- flexible QoS preservation.

In the first case, the value of the QoS coefficient $q$ is chosen to be very high ($10^6$), such that the DNN training is extremely biased towards maintaining the QoS constraints, on the cost of the achieved data rates. Regarding the second case "moderate QoS preservation" , the values of $q$ are chosen in order to achieve a trade-off between the achieved data rates and QoS preservation performance. Lastly, for the case of "flexible QoS preservation", $q$ is tuned more flexibly to achieve higher data rates, while keeping reasonable QoS preservation performance.

In Fig. 4.3, the achieved average data rate per user is plotted versus the number of RRH-user pairs in the network. It can be seen that the proposed DNNs based power allocation scheme can highly outperform the genetic algorithm for the "flexible QoS preservation" scheme. Moreover, this scheme can achieve a performance that is very close to the optimal, since the unsupervised learning is not upper bounded, unlike the supervised technique which cannot outperform the optimization algorithm used to obtain the ground truth data labels. It can be also noticed that the "moderate QoS preservation" scheme achieves a performance that is close to or even outperforming the genetic algorithm for all population size ($PS$) values employed in our simulations. Additionally, from Fig. 4.3, we can see that the performance of the genetic algorithm deteriorates when using smaller population sizes, which is expected. Likewise, the DNN models trained to have "high QoS preservation" were expected to have poor performance in-terms of the achieved data rates. In Fig. 4.4, a bar-chart of the achieved average weighted sum rates is plotted for the same schemes as before. The results confirm our conclusions from Fig. 4.3, where our DNN based scheme outperforms the genetic algorithm and can obtain a performance close to the optimal.

As pointed before, the main drawback of the DNN based power allocation is that a feasible solution is not always guaranteed, since its performance mainly relies on prediction. In tables 4.1, 4.2, and 4.3, the utilized values of $q$, along with the achieved data rates, and the probability that a user's QoS constraint is not satisfied ($P_{QV}$) are presented. We can notice that $P_{QV}$ values were very small for the highly preservative case. Then, there was a noticeable increase in $P_{QV}$ values for the moderate case when compared to the highly preservative. Finally, there was a slight increase in $P_{QV}$ when comparing the flexible case against the moderate.
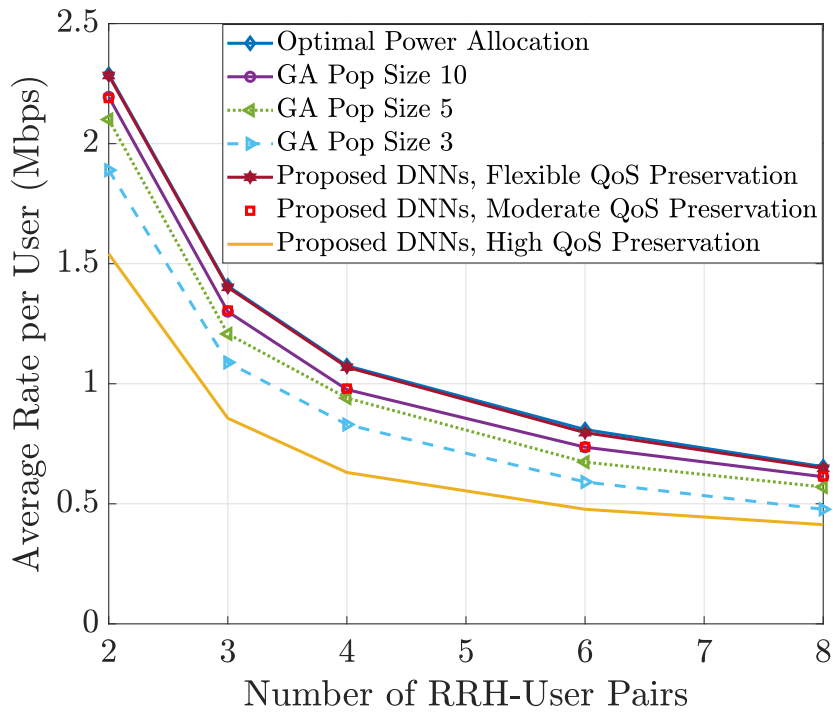
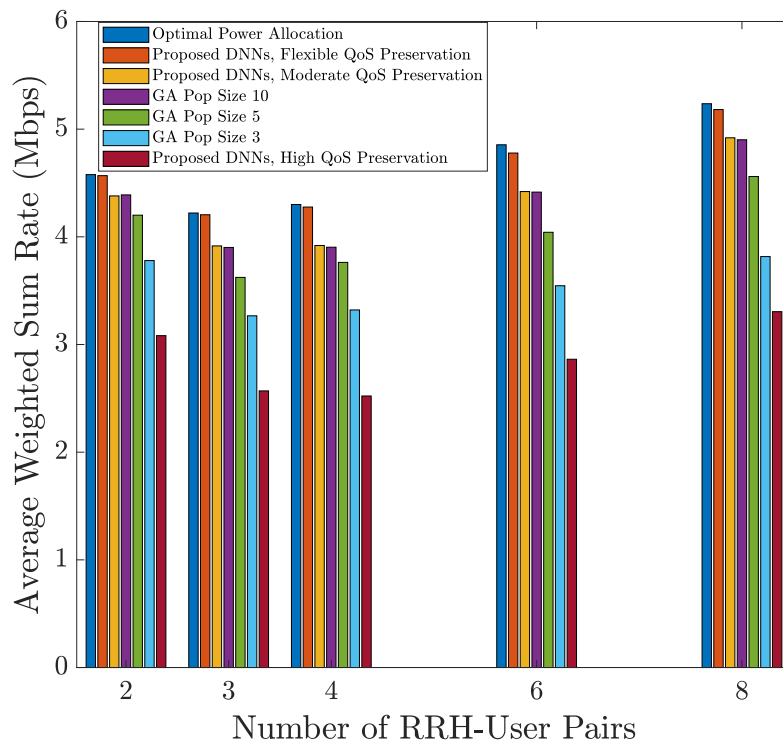Figure 4.3: Average data rate per user in Mbps Vs. Number of RRH-user pairs



Figure 4.4: Average weighted sum rate of the network Vs. Number of RRH-user pairs

48

Generally, increasing the value of $q$ results in trained DNN models that are more preservative regarding $P_{QV}$. Nevertheless, in some cases, using the same values of $q$, it is possible to achieve DNN models that have different performance regarding the achieved data rates and the resultant $P_{QV}$. For instance, in the case we had 4 RRH-user pairs, the same value of $q = 40$ was used to obtain both of the models which have the moderate QoS preservation, and flexible QoS preservation. However, the trade-off is always fixed: higher value of $P_{QV}$ is equivalent to higher achieved data rates, and the reverse. This is in fact due to the high experimental nature of deep learning, where many values of the different hyperparameters should be tested for several training epochs, and the best performing models, according to the pre-defined performance metrics should be saved for further usage.

| $K$ | $q$ | Av. Weighted Sum Rate (Mbps) | $P_{QV}$ |
|---|---|---|---|
| 2 | $10^6$ | 3.08223 | 0.00065 |
| 3 | $10^6$ | 2.56918 | 0.0189 |
| 4 | $10^6$ | 2.52191 | 0.0357 |
| 6 | $10^6$ | 2.86308 | 0.0729 |
| 8 | $10^6$ | 3.30492 | 0.1061 |

Table 4.1: High QoS preservation

| $K$ | $q$ | Av. Weighted Sum Rate (Mbps) | $P_{QV}$ |
|---|---|---|---|
| 2 | 240 | 4.3800 | 0.08279 |
| 3 | 100 | 3.9156 | 0.108 |
| 4 | 40 | 3.9196 | 0.12286 |
| 6 | 40 | 4.4204 | 0.14614 |
| 8 | 60 | 4.9196 | 0.16811 |

Table 4.2: Moderate QoS preservation

| $K$ | $q$ | Av. Weighted Sum Rate (Mbps) | $P_{QV}$ |
|---|---|---|---|
| 2 | 200 | 4.5681 | 0.09599 |
| 3 | 80 | 4.2050 | 0.1359 |
| 4 | 40 | 4.2765 | 0.1529 |
| 6 | 20 | 4.7781 | 0.16675 |
| 8 | 40 | 5.1820 | 0.17956 |

Table 4.3: Flexible QoS preservation

*Remark:* In our system model, we assumed only one radio resource for simplicity. However, our trained DNN models can be also utilized if the system has more than one radio resource, since they have the same probability distribution. Note that if the radio resources have different bandwidths, then DNN models should be trained for each bandwidth value.

### 4.4.2 Computational Complexity Analysis

It is well known that the DNNs training process is expensive in terms of computational complexity and required hardware [42]. Additionally, several values of the different hyperparameters need to be experimented such that the best performance can be obtained, which further increases the complexity of the process. However, the main motivation to use DNNs for optimization in cellular networks is that the training complexity is only experienced offline. Thus, after the DNNs models are trained, their online computational complexity becomes negligible. In this subsection, we compare the online computational complexity of the proposed DNNs for power allocation against the genetic algorithm and global search.

The DNNs online prediction process (based on forward propagation) is simply a series of vector-matrix multiplications that can be performed in a highly efficient way, thanks to the vectorization technique [57], in the modern platforms such as TensorFlow and MATLAB. With respect to other online optimization algorithms used in the literature, they are highly iterative in nature, in which each step is dependent on the previous, such that vectorization is not possible. Accordingly, it is naive and unrealistic to compare the online computational complexity of DNNs against those

algorithms only in terms of the number of floating-point operations performed. Thus, we compare practically the online computational complexity of our DNN based algorithm against the global search and genetic algorithms. For a fair comparison, we imported our trained DNN models from TensorFlow to MATLAB. Consequently, the same hardware and software platforms were used to run our online simulations to measure the computational complexity. Table 4.4 shows the average running time in seconds for the different algorithms for different values of RRH-user pairs $K$. It is clear that the global search has the highest computational complexity, which is expected. Regarding the genetic algorithm, as population size $PS$ increases, the complexity increases significantly, but we obtain higher data rates as shown in the previous subsection. We can also see that the least complex case for the genetic algorithm ($PS = 3$), which has a poor performance in terms of data rate, still has computational time that is orders of magnitude higher than the DNN models. Consequently, we can verify that our proposed DNNs can obtain superior performance that is close to the optimal with negligible computational complexity.

| Algorithm | K=2 | K=3 | K=4 | K=6 | K=8 |
|---|---|---|---|---|---|
| GL Search | 1.5962 | 4.9141 | 6.8890 | 17.5236 | 37.7127 |
| GA PS=10 | 0.5642 | 1.1226 | 1.6493 | 3.5816 | 5.8616 |
| GA PS=5 | 0.3863 | 0.6241 | 0.8465 | 1.5934 | 3.0056 |
| GA PS=3 | 0.1933 | 0.1968 | 0.2088 | 0.2702 | 0.3620 |
| DNNs | 0.00174 | 0.00175 | 0.0018 | 0.00181 | 0.00187 |

Table 4.4: Computation time in seconds

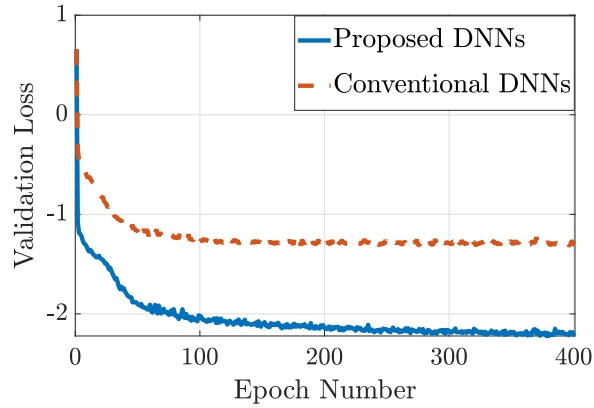### 4.4.3 Proposed Scheme Vs. Conventional Deep Learning Based Scheme

In this subsection, we discuss the added value obtained when feeding the QoS requirements of the users to the DNN structure. We prove numerically how our novel scheme outperforms the conventional deep learning based scheme. To compare against the conventional scheme, we consider the same settings employed regarding the DNNs architectures; number of data points and validation points, and their probability distribution. The only difference is that for the conventional scheme, the number of nodes in the input layer $n_1 = K^2$, which represent only the channel gains of the

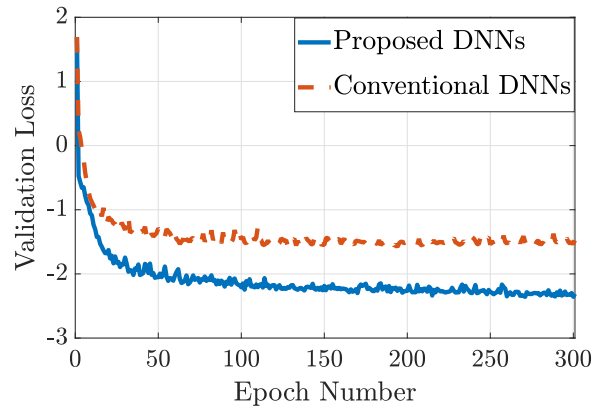network, since the QoS requirements are not inputted to the DNN structure.

In Fig. 4.5, the value of the validation loss function after each training epoch is plotted against the epoch number. It can be seen that our proposed scheme can converge to much smaller values of validation loss. Our interpretation for the improved performance is that in our proposed architecture, the DNN can have a comprehensive view of the relationship between the input QoS requirements and the loss function, since the loss function is dependent on the QoS requirements. Hence, during training, the DNN becomes able to learn better how to tune its trainable parameters (weights and biases) according to the input QoS requirements, such that a better performance can be achieved. Obviously, achieving smaller values of the loss function results in producing DNN models that can achieve higher data rates and less probability that the QoS requirement of a user is violated ($P_{QV}$). This is clearly shown in Fig. 4.6, in which the achieved data rates from the resulting DNN models are plotted against the corresponding value of $P_{QV}$. It can be noticed that the same trade-off holds as before: higher data rate results in higher $P_{QV}$. However, our proposed scheme outperforms the conventional scheme, such that significantly higher data rates can be achieved with smaller values of the resultant $P_{QV}$.

To discuss more how our proposed modification enhances the performance of the DNNs based power allocation schemes, we consider the "high QoS preservation" case, in which the value of the QoS coefficient is chosen to be very high ($q = 10^6$) such that the ultimate goal is to maintain the QoS constraints, as previously mentioned. Fig. 4.7 represents a comparison between our proposed scheme and the conventional DNNs based scheme in-terms of the minimum achieved values of $P_{QV}$ for the underlaying case. To obtain the results in Fig. 4.7, DNN models were trained for each value of $K$ (number of RRH-user pairs), such that each model's training process was run for 300 Epochs, and the model achieving the minimum value of $P_{QV}$ for the validation data was saved. It can be seen from the figure that employing our proposed scheme can offer significant reduction in the proportion of devices with violated QoS constraints, when compared to conventional scheme. The efficiency of our proposed scheme becomes more significant with the increase in the number of RRH-user pairs $K$, since it becomes more challenging for the DNNs to predict power profiles satisfying the QoS requirements of larger number of users. However, as our scheme enhances the learning capability of DNNs, this challenge can be faced more efficiently. Thus, we can confirm that
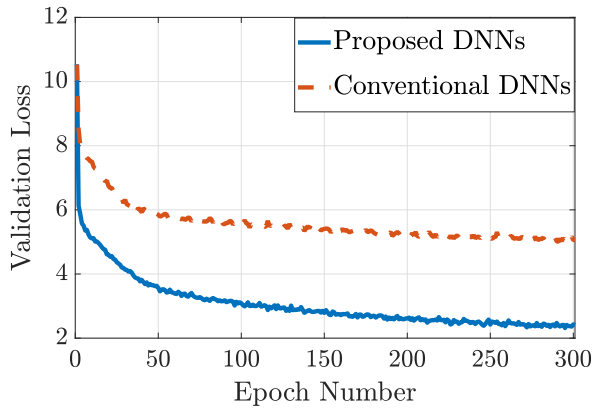
our scheme enhances the learning capability and performance of the DNNs, such that less values of $P_{QV}$ (probability QoS of a user is violated) were achieved when the training process was biased towards maintaining the QoS constraints.
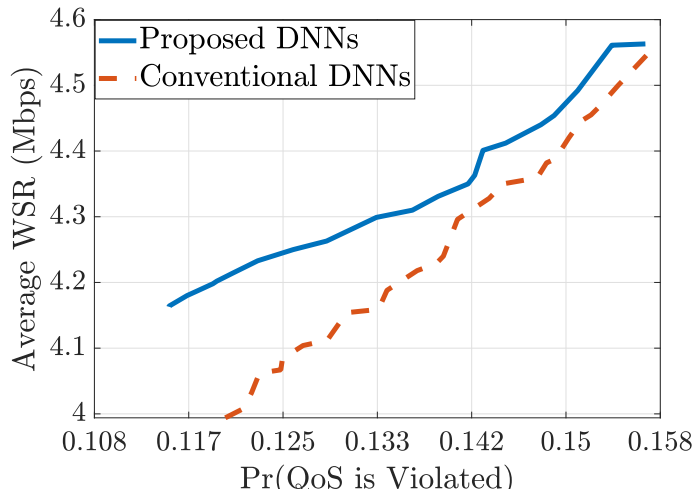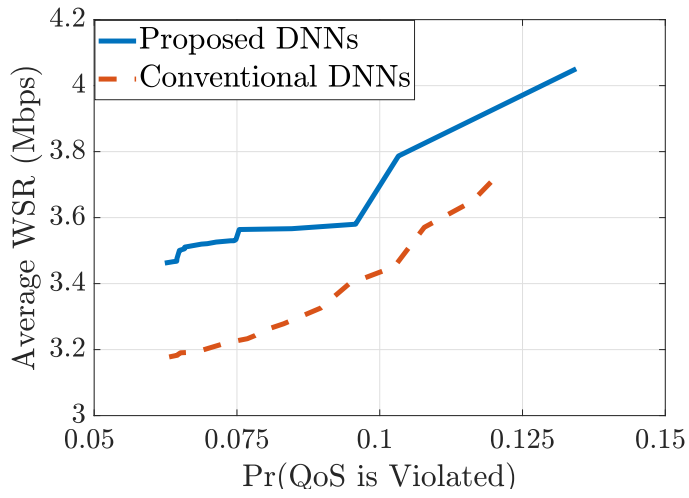


(a) $K = 6, q = 40$
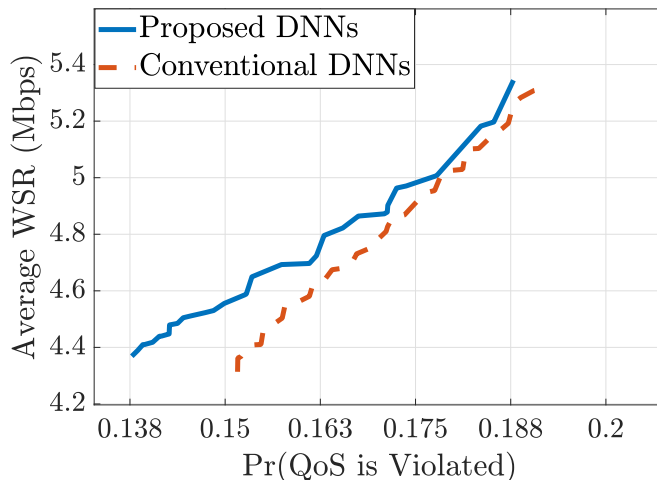
(b) $K = 4, q = 80$

(c) $K = 8, q = 80$

Figure 4.5: Validation Loss Vs Epoch Number

(a) $K = 6$, $q = 40$



(b) $K = 4$, $q = 80$



(c) $K = 8$, $q = 80$

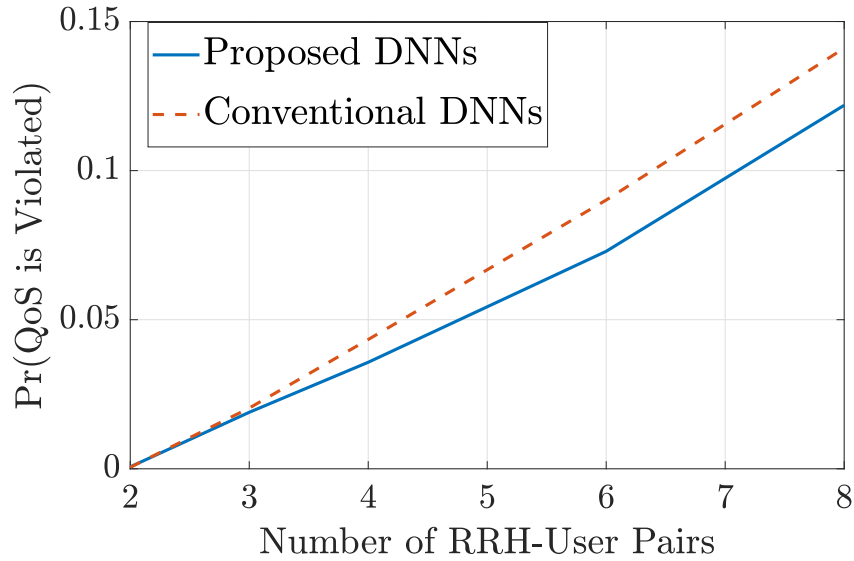Figure 4.6: Average Weighted Sum Rate Vs $P_{QV}$

Figure 4.7: $P_{QV}$ Vs. Number of RRH-user pairs, $q = 10^6$

## 4.5  Conclusions

We proposed a deep learning based power allocation scheme, while considering user associa-
tion to reflect a real cellular scenario. It was shown that the proposed algorithm could achieve a
performance close to optimal with negligible computational complexity. Moreover, we proposed
a novel approach that can enhance the learning capability of DNNs through directly feeding the
QoS requirements of the users to the DNN structure, which proved to outperform the traditional ap-
proach in the literature. However, the sole drawback associated with using the deep learning based
algorithm for power allocation is that a feasible solution is not always guaranteed. The trade-off
between achieved data rates and QoS preservation performance was discussed thoroughly in the
provided simulation results. One remedy for that issue is to opt for the traditional optimization al-
gorithms whenever the DNNs predictions are not feasible. In such case, we can reap the advantages
of deep learning, while avoiding its main drawback.

# Chapter 5

# Conclusions and Future Works

## 5.1 Conclusions

In this thesis, we proposed enhancing the performance of cloud radio access networks. We considered various optimization techniques based on game theory, and convex optimization. Moreover, to overcome the computational complexity of the conventional optimization techniques, we proposed utilizing deep neural networks for power allocation.

First, we proposed a generalized CoMP transmission technique, in which the number of RRHs serving a specific user is unlimited and determined according to the network status. We focused on optimizing the user association, radio resource allocation, and power allocation. The user association and resource allocation sub-problems were tackled via novel matching game based algorithms. A novel approach was proposed to convexify the power allocation sub-problem. We considered three scenarios for the H-CRAN fronthaul network: individual fronthaul links with limited capacity, shared fronthaul links with limited capacity, and fronthaul links with unlimited capacity. Our proposed algorithms proved to achieve significant cooperation gain with CoMP, even with very tight fronthaul constraints.

Finally, we proposed optimizing the power allocation in CRAN via unsupervised deep learning, in order to maximize the network weighted sum rate. We provided intensive simulations to discuss the trade-off between the QoS preservation performance and the achieved data rates. Also, we explained how to tune the QoS coefficient hyperparameter to obtain the desired performance.

Our proposed algorithm proved to outperform the conventional online optimization algorithms, with negligible computational complexity. Moreover, we proposed a modification to enhance the learning capability of the DNNs based power allocation algorithms. Our proposed modification proved to provide a significant performance gain regarding the achieved data rates and users' QoS preservation.

## 5.2  Future Works

CoMP can achieve magnificent network performance improvements. However, in a CRAN with limited fronthaul capacity constraints, where the network radio capacity is much higher than the fronthaul capacity, CoMP can actually cause performance deterioration. The reason is that with CoMP, the data of the served user must be sent from the BBU pool to all the serving RRHs, which causes a high burden on the fronthaul links. Consequently, the network coverage and spectral efficiency may be reduced, as the number of served users will decrease because of the limited fronthaul capacity. In chapter 3, we proposed a generalized CoMP scheme with an objective to maximize the weighted sum rate of the network, where our user association algorithm could achieve cooperation gain, even with tight fronthaul constraints. However, more research is needed to develop similar algorithms for other optimization objectives such as minimizing latency, maximizing power efficiency, and maximizing fairness, among others.

The scarcity in the available radio resources has driven the researchers to investigate the utilization of millimeter wave (mmWave) frequencies in wireless communications. However, due to their extremely small wavelengths, communications with mmWave are vulnerable, and can never be accomplished without a line of sight (LoS) between the source and destination, which cannot be guaranteed in cellular networks. Consequently, researchers proposed the concept of multi-connectivity [61], in which a user can be associated with several basestation (RRHs) while using different radio access technologies. This means that a user can be allocated mmWave and micro wave frequencies simultaneously. The usage of mmWave with multi-connectivity can enhance the CRAN performance significantly through improving the coverage, enhancing QoS, and providing much higher data rates. Currently, research has been focused on the performance optimization of

multi-connectivity enabled CRAN. For instance, authors of [62] proposed a heuristic algorithm to tackle a user association problem, with an objective to minimize the network power consumption. However, more research works are still needed to tackle the optimization problems in such CRANs. The main challenge is that with multi-connectivity, these problems become much more complex, which requires more advanced optimization algorithms. One idea is to opt-for machine learning based algorithms to deal with such problems. In chapter 4, we proposed a deep learning based power allocation for CRAN with single association. A possible expansion to our work is to consider a system model that includes multi-connectivity. Another possible expansion is to consider deep learning based power allocation algorithms for a system model utilizing CoMP transmissions.

# Appendix A

# Proof of convexity of P1-PA

To prove the convexity of our maximization optimization problem, we need to prove that the objective function is concave, and inequality constraints are convex functions. For the objective function of **P1-PA** (3.23), we can easily rewrite it as follows:

$$\sum_{n \in \aleph} B^n \sum_{i \in \kappa_n} \log_2(\sum_{j \in \tau_i} P_j^n g_{ji}^n + P_M^n g_{Mi}^n + \sum_{j \notin \tau_i} P_j^n g_{ji}^n + \sigma^2) +$$

$$\log_2(\frac{1}{P_M^n g_{Mi}^n + \sum_{j \notin \tau_i} P_j^n g_{ji}^n + \sigma^2}) \quad \text{(A.1)}$$

Apparently, the first logarithmic term in the objective function, $\log_2(\sum_{j \in \tau_i} P_j^n g_{ji}^n + P_M^n g_{Mi}^n + \sum_{j \notin \tau_i} P_j^n g_{ji}^n + \sigma^2)$ is concave, while the second logarithmic term is convex. Thus, the sum of the two terms is generally neither concave or convex. However, if the additional constraint (3.25) can be satisfied and our problem is feasible, the first logarithmic concave term will be always greater that the second logarithmic convex term, and any increase in the second term will result in a higher increase in the first. Hence, the sum of the two logarithmic terms will be monotonically increasing in the feasible region, and consequently, concave. Since the sum of concave functions is also concave, our whole objective function is concave. Regarding the constraints, (3.3) can be easily proved to be a convex function in a similar way as above; while constraints (3.4), and (3.8) are clearly affine functions. Accordingly, **P1-PA** is a convex problem.

# Appendix B

# List of Publications

- M. Labana and W. Hamouda, "Advances in CRAN Performance Optimization," in IEEE Network, doi: 10.1109/MNET.011.2000502.

- M. Labana and W. Hamouda, "Joint User Association and Resource Allocation in CoMP-Enabled Heterogeneous CRAN," GLOBECOM 2020 - 2020 IEEE Global Communications Conference, Taipei, Taiwan, 2020, pp. 1-6, doi: 10.1109/GLOBECOM42002.2020.9322501.

- M. Labana and W. Hamouda, "Unsupervised Deep Learning for Power Allocation in CRAN," ICC 2021 - 2021 IEEE International Conference on Communications (ICC), Montreal, Canada, 2021, Accepted.

- M. Labana and W. Hamouda, "Unsupervised Deep Learning Approach for Near Optimal Power Allocation in CRAN," IEEE Transactions on Vehicular Technology, Major Revision.

# Bibliography

[1] M. Elbayoumi, M. Kamel, W. Hamouda, and A. Youssef, "NOMA-assisted machine-type communications in UDN: State-of-the-art and challenges," *IEEE Communications Surveys Tutorials*, vol. 22, no. 2, pp. 1276–1304, 2020.

[2] L. Ferdouse, A. Anpalagan, and S. Erkucuk, "Joint communication and computing resource allocation in 5G cloud radio access networks," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 9, pp. 9122–9135, 2019.

[3] Y. Lee, K. Miyanabe, H. Nishiyama, N. Kato, and T. Yamada, "Threshold-based RRH switching scheme considering baseband unit aggregation for power saving in a cloud radio access network," *IEEE Systems Journal*, vol. 13, no. 3, pp. 2676–2687, 2019.

[4] D. Zeng, J. Zhang, L. Gu, S. Guo, and J. Luo, "Energy-efficient coordinated multipoint scheduling in green cloud radio access network," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 10, pp. 9922–9930, 2018.

[5] L. You and D. Yuan, "User-centric performance optimization with remote radio head cooperation in C-RAN," *IEEE Transactions on Wireless Communications*, vol. 19, no. 1, pp. 340–353, 2020.

[6] M. Labana and W. Hamouda, "Advances in CRAN performance optimization," *IEEE Network*, pp. 1–7, 2020.

[7] T. Zhang and S. Mao, "Energy-efficient power control in wireless networks with spatial deep neural networks," *IEEE Transactions on Cognitive Communications and Networking*, vol. 6, no. 1, pp. 111–124, 2020.

[8] L. Liang, H. Ye, G. Yu, and G. Y. Li, "Deep-learning-based wireless resource allocation with application to vehicular networks," *Proceedings of the IEEE*, vol. 108, no. 2, pp. 341–356, 2020.

[9] T. Koketsu Rodrigues, K. Suto, and N. Kato, "Edge cloud server deployment with transmission power control through machine learning for 6G internet of things," *IEEE Transactions on Emerging Topics in Computing*, pp. 1–1, 2019.

[10] K. Miyanabe, T. Gama Rodrigues, Y. Lee, H. Nishiyama, and N. Kato, "An internet of things traffic-based power saving scheme in cloud-radio access network," *IEEE Internet of Things Journal*, vol. 6, no. 2, pp. 3087–3096, 2019.

[11] M. K. Elhattab, M. M. Elmesalawy, T. Ismail, H. H. Esmat, M. M. Abdelhakam, and H. Selmy, "A matching game for device association and resource allocation in heterogeneous cloud radio access networks," *IEEE Communications Letters*, vol. 22, no. 8, pp. 1664–1667, 2018.

[12] J. Liu, S. Zhang, N. Kato, H. Ujikawa, and K. Suzuki, "Device-to-device communications for enhancing quality of experience in software defined multi-tier LTE-A networks," *IEEE Network*, vol. 29, no. 4, pp. 46–52, 2015.

[13] M. Waqar, "A study of fronthaul networks in CRANs - requirements and recent advancements," *KSII Transactions on Internet and Information Systems*, 10 2018.

[14] D. Liu, L. Wang, Y. Chen, M. Elkashlan, K. Wong, R. Schober, and L. Hanzo, "User association in 5G networks: A survey and an outlook," *IEEE Communications Surveys Tutorials*, vol. 18, no. 2, pp. 1018–1044, 2016.

[15] A. A. Ahmad, H. Dahrouj, A. Chaaban, A. Sezgin, T. Y. Al-Naffouri, and M. Alouini, "Distributed cloud association and beamforming in downlink multi-cloud radio access networks," in *2020 IEEE International Conference on Communications Workshops (ICC Workshops)*, 2020, pp. 1–6.

[16] R. Trestian, O. Ormond, and G. Muntean, "Game theory-based network selection: Solutions and challenges," *IEEE Communications Surveys Tutorials*, vol. 14, no. 4, pp. 1212–1231, 2012.

[17] D. Liu, Y. Chen, K. K. Chai, T. Zhang, and M. Elkashlan, "Opportunistic user association for multi-service HetNets using nash bargaining solution," *IEEE Communications Letters*, vol. 18, no. 3, pp. 463–466, 2014.

[18] V. N. Ha and L. B. Le, "Distributed base station association and power control for heterogeneous cellular networks," *IEEE Transactions on Vehicular Technology*, vol. 63, no. 1, pp. 282–296, 2014.

[19] M. Hong and Z. Luo, "Distributed linear precoder optimization and base station selection for an uplink heterogeneous network," *IEEE Transactions on Signal Processing*, vol. 61, no. 12, pp. 3214–3228, 2013.

[20] M. M. Abdelhakam, M. M. Elmesalawy, K. R. Mahmoud, and I. I. Ibrahim, "A cooperation strategy based on bargaining game for fair user-centric clustering in cloud-RAN," *IEEE Communications Letters*, vol. 22, no. 7, pp. 1454–1457, 2018.

[21] H. Zheng, S. Hou, H. Li, Z. Song, and Y. Hao, "Power allocation and user clustering for uplink MC-NOMA in D2D underlaid cellular networks," *IEEE Wireless Communications Letters*, vol. 7, no. 6, pp. 1030–1033, 2018.

[22] S. Luo, R. Zhang, and T. J. Lim, "Downlink and uplink energy minimization through user association and beamforming in C-RAN," *IEEE Transactions on Wireless Communications*, vol. 14, no. 1, pp. 494–508, 2015.

[23] Y. Sun, F. Wang, and Z. Liu, "Coalition formation game for resource allocation in D2D uplink underlaying cellular networks," *IEEE Communications Letters*, vol. 23, no. 5, pp. 888–891, 2019.

[24] G. R. Harik, F. G. Lobo, and D. E. Goldberg, "The compact genetic algorithm," *IEEE transactions on evolutionary computation*, vol. 3, no. 4, pp. 287–297, 1999.

[25] A. Mokdad, P. Azmi, N. Mokari, M. Moltafet, and M. Ghaffari-Miab, "Cross-layer energy efficient resource allocation in PD-NOMA based H-CRANs: Implementation via GPU," *IEEE Transactions on Mobile Computing*, vol. 18, no. 6, pp. 1246–1259, 2019.

[26] K. Wang, W. Zhou, and S. Mao, "On joint BBU/RRH resource allocation in heterogeneous cloud-RANs," *IEEE Internet of Things Journal*, vol. 4, no. 3, pp. 749–759, 2017.

[27] K. Zhang, W. Tan, G. Xu, C. Yin, W. Liu, and C. Li, "Joint RRH activation and robust coordinated beamforming for massive MIMO heterogeneous cloud radio access networks," *IEEE Access*, vol. 6, pp. 40 506–40 518, 2018.

[28] Z. Wu, Z. Fei, Z. Zheng, B. Li, and Z. Han, "Remote radio head activation and user association in dense C-RANs," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 10, pp. 12 216–12 228, 2020.

[29] Z. Wu, Z. Fei, Y. Yu, and Z. Han, "Toward optimal remote radio head activation, user association, and power allocation in C-RANs using Benders decomposition and ADMM," *IEEE Transactions on Communications*, vol. 67, no. 7, pp. 5008–5023, 2019.

[30] E. A. Ramos da Paixão, R. F. Vieira, W. V. Araújo, and D. L. Cardoso, "Optimized load balancing by dynamic BBU-RRH mapping in C-RAN architecture," in *2018 Third International Conference on Fog and Mobile Edge Computing (FMEC)*, 2018, pp. 100–104.

[31] K. Boulos, M. El Helou, and S. Lahoud, "RRH clustering in cloud radio access networks," in *2015 International Conference on Applied Research in Computer Science and Engineering (ICAR)*, 2015, pp. 1–6.

[32] Y. Chen, W. Chiang, and M. Shih, "A dynamic BBU–RRH mapping scheme using borrow-and-lend approach in cloud radio access networks," *IEEE Systems Journal*, vol. 12, no. 2, pp. 1632–1643, 2018.

[33] T. H. T. Le, N. H. Tran, P. L. Vo, Z. Han, M. Bennis, and C. S. Hong, "Joint cache allocation with incentive and user association in cloud radio access networks using hierarchical game," *IEEE Access*, vol. 7, pp. 20 773–20 788, 2019.

[34] N. Gholipoor, S. Parsaeefard, M. R. Javan, N. Mokari, H. Saeedi, and H. Pishro-Nik, "Cloud-based queuing model for Tactile Internet in next generation of RAN," in *2020 IEEE 91st Vehicular Technology Conference (VTC2020-Spring)*, 2020, pp. 1–6.

[35] W. Hao, O. Muta, and H. Gacanin, "Price-based resource allocation in massive MIMO H-CRANs with limited fronthaul capacity," *IEEE Transactions on Wireless Communications*, vol. 17, no. 11, pp. 7691–7703, 2018.

[36] J. Zuo, J. Zhang, C. Yuen, W. Jiang, and W. Luo, "Energy efficient user association for cloud radio access networks," *IEEE Access*, vol. 4, pp. 2429–2438, 2016.

[37] V. N. Ha, L. B. Le, and N. Đào, "Coordinated multipoint transmission design for cloud-RANs with limited fronthaul capacity constraints," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 9, pp. 7432–7447, 2016.

[38] M. Farahmand and A. Mohammadi, "Sparse power allocation in downlink transmission of cloud radio access networks," *IET Communications*, vol. 11, no. 16, pp. 2531–2538, 2017.

[39] X. Cao, R. Ma, L. Liu, H. Shi, Y. Cheng, and C. Sun, "A machine learning-based algorithm for joint scheduling and power control in wireless networks," *IEEE Internet of Things Journal*, vol. 5, no. 6, pp. 4308–4318, 2018.

[40] G. Qian, Z. Li, C. He, X. Li, and X. Ding, "Power allocation schemes based on deep learning for distributed antenna systems," *IEEE Access*, vol. 8, pp. 31 245–31 253, 2020.

[41] J. Luo, J. Tang, D. K. C. So, G. Chen, K. Cumanan, and J. A. Chambers, "A deep learning-based approach to power minimization in multi-carrier NOMA with SWIPT," *IEEE Access*, vol. 7, pp. 17 450–17 460, 2019.

[42] F. Liang, C. Shen, W. Yu, and F. Wu, "Towards optimal power control via ensembling deep neural networks," *IEEE Transactions on Communications*, vol. 68, no. 3, pp. 1760–1776, 2020.

[43] M. Peng, Y. Sun, X. Li, Z. Mao, and C. Wang, "Recent advances in cloud radio access networks: System architectures, key techniques, and open issues," *IEEE Communications Surveys Tutorials*, vol. 18, no. 3, pp. 2282–2308, thirdquarter 2016.

[44] A. E. Roth, "Deferred acceptance algorithms: History, theory, practice, and open questions," *international Journal of game Theory*, vol. 36, no. 3-4, pp. 537–569, 2008.

[45] T. Kim and J. M. Chang, "QoS-aware energy-efficient association and resource scheduling for HetNets," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 1, pp. 650–664, Jan 2018.

[46] S. Parsaeefard, V. Jumba, A. Dalili, M. Derakhshani, and T. Le-Ngoc, "User association in cloud RANs with massive MIMO," *IEEE Transactions on Cloud Computing*, pp. 1–1, 2018.

[47] J. Yao and N. Ansari, "QoS-aware joint BBU-RRH mapping and user association in cloud-RANs," *IEEE Transactions on Green Communications and Networking*, vol. 2, no. 4, pp. 881–889, 2018.

[48] M. Servetnyk and C. C. Fung, "Precoding and selection for coordinated multipoint transmission in fronthaul-constrained cloud-RAN," *IEEE Wireless Communications Letters*, vol. 9, no. 1, pp. 51–55, 2020.

[49] M. Labana and W. Hamouda, "Joint user association and resource allocation in CoMP-enabled heterogeneous CRAN," in *GLOBECOM 2020 - 2020 IEEE Global Communications Conference*, 2020, pp. 1–6.

[50] W. Lee, M. Kim, and D. Cho, "Transmit power control using deep neural network for underlay device-to-device communication," *IEEE Wireless Communications Letters*, vol. 8, no. 1, pp. 141–144, 2019.

[51] S. Sharma, "Activation functions in neural networks," *Towards Data Science*, vol. 6, 2017.

[52] H. Sun, X. Chen, Q. Shi, M. Hong, X. Fu, and N. D. Sidiropoulos, "Learning to optimize: Training deep neural networks for interference management," *IEEE Transactions on Signal Processing*, vol. 66, no. 20, pp. 5438–5453, 2018.

[53] W. Lee, M. Kim, and D. Cho, "Deep power control: Transmit power control scheme based on convolutional neural network," *IEEE Communications Letters*, vol. 22, no. 6, pp. 1276–1279, 2018.

[54] W. Lee, "Resource allocation for multi-channel underlay cognitive radio network based on deep neural network," *IEEE Communications Letters*, vol. 22, no. 9, pp. 1942–1945, 2018.

[55] Q. Qian, R. Jin, J. Yi, L. Zhang, and S. Zhu, "Efficient distance metric learning by adaptive sampling and mini-batch stochastic gradient descent (SGD)," *Machine Learning*, vol. 99, no. 3, pp. 353–372, 2015.

[56] D. Masters and C. Luschi, "Revisiting small batch training for deep neural networks," 2018.

[57] J. Ren and L. Xu, "On vectorization of deep convolutional neural networks for vision tasks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 29, no. 1, 2015.

[58] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2017.

[59] Y. Wang, S. Wang, and L. Liu, "Joint beamforming and power allocation using deep learning for D2D communication in heterogeneous networks," *IET Communications*, vol. 14, no. 18, pp. 3095–3101, 2020.

[60] C. Du, Z. Zhang, X. Wang, and J. An, "Deep learning based power allocation for workload driven full-duplex D2D-aided underlaying networks," *IEEE Transactions on Vehicular Technology*, pp. 1–1, 2020.

[61] J. Deng, O. Tirkkonen, R. Freij-Hollanti, T. Chen, and N. Nikaein, "Resource allocation and interference management for opportunistic relaying in integrated mmwave/sub-6 GHz 5G networks," *IEEE Communications Magazine*, vol. 55, no. 6, pp. 94–101, 2017.

[62] M. Saimler and S. Coleri, "Multi-connectivity based uplink/downlink decoupled energy efficient user association in 5G heterogenous CRAN," *IEEE Communications Letters*, vol. 24, no. 4, pp. 858–862, 2020.