

# Adaptive Cross-Network Cross-Layer Design in Heterogeneous Wireless Networks

Honghao Ju, Ben Liang, *Senior Member, IEEE*, Jiandong Li, *Senior Member, IEEE*, Yan Long, and Xiaoniu Yang

**Abstract**—A cross-network cross-layer design method is proposed to exploit the trunking, diversity, and best service assignment gains available in a heterogeneous wireless network (HWN), consisting of orthogonal radio access networks (RANs) and interference-limited RANs. Accounting for traffic-level dynamics and channel fading, we jointly design the distribution strategy for elastic and inelastic traffics, and the radio resource management strategy for RANs, in a network-separable control architecture. Optimal and quantified near-optimal radio allocation schemes are proposed for each type of RANs, which are combined into an on-line design framework that over time provides asymptotically optimal performance, maximizing the sum throughput utility for elastic traffic while guaranteeing the throughput requirements of inelastic traffic. Extensive simulation results demonstrate substantial performance improvement against suboptimal alternatives.

## I. INTRODUCTION

The current popularity of mobile Internet not only leads to a drastic increase in the traffic amount but also creates a heterogeneous traffic environment, where both elastic traffic, such as web browsing, and inelastic traffics, such as video streaming, coexist. Furthermore, the concurrent deployment of various radio access technologies, such as 3G, LTE, WiMAX, and WiFi, has led to heterogeneous wireless networks (HWNs), where multiple radio access networks (RANs) overlap with each other. Meanwhile, popular user equipment (UE) such as smart phones and tablets are becoming more powerful with their multi-homing capability. This creates a new dimension of flexibility, i.e., multiple RAN selection, to serve the user traffic, in addition to the allocation of radio resource within each RAN. A central question is on how to best match the heterogeneous traffic to the heterogeneous RANs.

An optimal match between traffic and RANs should fully exploit the benefits of HWN over traditional homogeneous networks. First, the trunking gain of multiplexing multiple traffic flows can be amplified with joint optimization of the UE traffic and the radio resource of RANs. Second, spatial transmission diversity and multi-user diversity can be enhanced by careful scheduling within each RAN to account for large- and small-scale channel fading, respectively. Third, the ramification of multi-radio transmission gain, by allowing the

UE to simultaneously transmit over multiple RANs, should be properly planned. A best service assignment principle is required to distribute to each RAN its most suitable traffic. Since the above considerations span multiple RANs and multiple network layers, a joint resource management framework is required for optimal HWN operation.

There are many technical challenges in the resultant heterogeneous joint resource management problem across RANs and network layers. First, different types of RANs each brings its own unique design issues. In particular, in an orthogonal RAN, the radio resource granularity leads to a mixed-integer non-linear optimization problem, while in an interference-limited RAN, because of the non-convexity of the SINR term, the optimal power allocation problem is NP-hard [1]. In addition, due to different hardware and legal constraints, both the average power and the maximum power of RANs may be constrained at different levels. Second, to explore the best-service assignment gain, a deliberate heterogeneous traffic management scheme has to be designed to explore both the diverse traffic requirements and the RAN characteristics. Third, traffic-level dynamics and channel fading should be considered in a realistic model of the network. These combined factors present difficult challenges to solve the joint resource management problem. As explained in Section II, no existing methods are directly applicable to provide a tractable solution.

In this paper, we specifically focus on how to jointly design the distribution of heterogeneous traffic among heterogeneous RANs and the radio resource allocation strategy in an HWN that consists of orthogonal RANs and interference-limited RANs. We propose an on-line cross-network cross-layer (CNCL) network control and resource allocation framework, which is adaptive to both traffic-level dynamics and channel fading in the HWN and is shown to offer asymptotically optimal throughput utility for elastic traffic while guaranteeing the required throughput for inelastic traffic. The main contributions of this paper are summarized as follows:

- The CNCL design introduces several Lyapunov-typed control techniques [2], [3], [4] into the HWN environment. We show that the joint resource management problem can be divided, without loss in optimality, into two components: *traffic admission control*, responsible for traffic admission and distribution among and within RANs according to the RAN load and the traffic requirement; and *RAN-level radio resource management*, responsible for radio resource allocation in orthogonal and interference-limited RANs.
- In an orthogonal RAN, the allocation of resource blocks (RBs) as the basic resource granularity leads to a mixed-integer non-linear program, while in an interference-limited RAN, the power allocation problem is NP-hard.

This work was supported by the National Natural Science Foundation of China (61231008), the National Nature Science Foundation of China (61101143), and the Program for Changjiang Scholars and Innovative Research Team in University (PCSIRT) (IRT0852), and and 111 Project under grant B08038.

Honghao Ju, Jiandong Li, Yan Long, and Xiaoniu Yang are with the Information Science Institute, State Key Laboratory of Integrated Service Networks, Xidian University, Xi'an, Shaanxi, China. Ben Liang is with the Department of Electrical and Computer Engineering, University of Toronto, Ontario, Canada. Honghao Ju was a visiting student at the University of Toronto supported by the China Scholarship Council.

Solutions are given for both problems, providing optimal or quantified near-optimal performance at each time slot. Both average and maximum power constraints are accommodated in all RANs.

- Over multiple time slots, CNCL adaptively handles traffic-level dynamics and channel fading. The proposed on-line design is shown to achieve asymptotically optimal performance, maximizing the sum elastic throughput utility while guaranteeing the inelastic throughput requirements.

The organization of the rest of this paper is as follows: In Section II, we summarize the related work. We give the HWN model in Section III. The CNCL design method is proposed in Section IV, along with analytical evaluation leading to a quantified performance bound. In Section V, we verify the performance of CNCL through simulation. We summarize this paper in Section VI.

## II. RELATED WORK

The trunking gain, spatial transmission diversity, multi-user diversity, multi-radio transmission gain, and best service assignment gain available in HWNs have been defined in [5], [6], [7], [8], [9], [10]. Even though these works promote the above benefits of HWNs to improve UE and network performance. They do not discuss how to derive these benefits by joint resource optimization in an HWN, which is studied in this paper.

Some prior works focus on deriving the best service assignment gain with a given transmission rate for each link. Given the capacity region of each RAN, the optimal service allocation strategy in an HWN is studied in [11]. The Erlang capacity of an HWN based on the M/M/m queue is given in [12]. The fairness issue regarding traffic admission control in an HWN is discussed in [13]. Game theory has been used to study the admission control problem in an HWN in [14], [15], [16]. In [17], an on-line traffic admission control algorithm is proposed for heterogeneous flows. In [18], a Markov decision approach is used to study on-line admission control in an HWN to minimize the time average blocking cost. However, none of the studies above discusses how to design the radio resource allocation strategy with respect to the RAN transmission rate, which is a main focus of this paper.

Call admission control and load balancing in the cellular-WLAN integrated network are studied in [19], [20], [21], [22], [23]. However, different from the cellular-WLAN architecture, we study joint resource management in an HWN consisting of orthogonal RANs and interference-limited RANs. Particularly, we study practical network constraints, such as resource granularity in orthogonal RANs and power optimization in both orthogonal and interference RANs.

There are prior studies that discuss the relation between the RAN transmission rate and the radio allocation strategy. Joint power and bandwidth allocation is considered in [24]. However, it only focuses on the orthogonal RAN and does not consider the resource granularity (e.g., resource block) constraint. In [25], the UE outage probability is minimized in an HWN consisting of a CDMA-based RAN and a TDMA-based RAN, and the RAN selection scheme is discussed in a

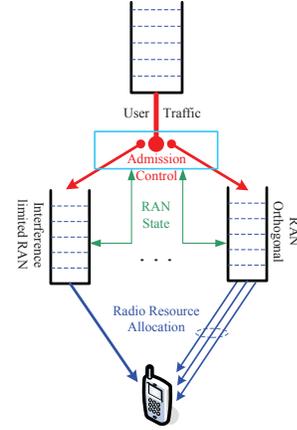


Fig. 1. UE served by a heterogeneous wireless network

CDMA/WLAN heterogeneous wireless network in [26]. Considering both elastic and inelastic traffic, joint radio resource management is studied in [27] in an HWN with orthogonal and interference-limited RANs. Decentralized resource allocation schemes in the HWN are proposed in [28], [29]. Resource management for allowing UE to transmit over both single RAN and multiple RANs is investigated in [30]. However, all of these works study the HWN performance in a deterministic fluid-flow traffic model. In contrast, our work considers adaptive HWN optimization with stochastic modeling and dynamic control.

## III. HETEROGENEOUS WIRELESS NETWORK MODEL

We consider the downlink of an HWN constituted of  $N$  RANs, denoted by the set  $\mathcal{R}$ . They may contain a mixture of interference-limited RANs (e.g., CDMA based) and orthogonal RANs (e.g., TDMA or OFDMA based). We assume that the RANs use orthogonal radio resources, and hence do not interfere with each other. For an interference-limited RAN, the intra-cell interference is considered, while for the orthogonal RAN, due to channel orthogonality, there is no interference between transmissions. We denote the set of interference-limited RANs and the set of the orthogonal RANs as  $\mathcal{R}_I$  and  $\mathcal{R}_O$  where  $\mathcal{R}_I \cup \mathcal{R}_O = \mathcal{R}$ , and let  $N_I$  and  $N_O$  be the number of RANs within each set, respectively.

The HWN serves  $M$  UEs, denoted by the set  $\mathcal{M}$ . To fully exploit the benefit introduced by network heterogeneity, the UEs have multi-homing capability and can access one or multiple RANs based on the UE traffic requirement and the RAN load. As shown in Fig. 1, a traffic admission controller distributes the UE traffic among RANs based on the feedback of the RAN state, while each RAN matches its radio resource to service the admitted UE traffic. For simplicity of illustration, we assume here that all UEs may freely connect with all RANs. The proposed model and analysis can be easily extended to the case where each UE is able to connect with only a limited subset of RANs, by setting the channel gain between some RANs and the UE to zero. The important notations used throughout this paper are summarized in Table I.

TABLE I  
NOTATION TABLE

$\mathcal{M}$	UE set
$\mathcal{R}, \mathcal{R}_I, \mathcal{R}_O$	Sets for RAN, interference-limited RAN, and orthogonal RAN
$n_I, n_O$	Index of the interference-limited RAN and the orthogonal RAN
$W_{n_I}, W_{n_O}$	Symbol rate of the interference-limited RAN $n_I$ and the orthogonal RAN $n_O$
$\sigma^2$	Noise power
$\Gamma$	Channel capacity gap from Shannon bound
$\Upsilon_{m,n_I}$	Signal processing gain for UE $m$ in interference-limited RAN $n_I$
$g_{m,n_I}(t), g_{m,n_O,i}(t)$	Channel gain between UE $m$ and interference-limited RAN $n_I$ , and between UE $m$ and orthogonal RAN $n_O$ on RB $i$ in time slot $t$
$\mu_{m,n_I}(t), \mu_{m,n_O,i}(t)$	Transmission rate for UE $m$ in interference-limited RAN $n_I$ , and for UE $m$ over RB $i$ in orthogonal RAN $n_O$ in time slot $t$
$p_{m,n_I}(t), p_{m,n_O,i}(t)$	Transmission power of interference-limited RAN $n_I$ for UE $m$ , and of orthogonal RAN on RB $i$ for UE $m$ in time slot $t$
$P_{n_I}, P_{n_O}$	Average power constraints for interference-limited RAN $n_I$ and orthogonal RAN $n_O$
$\hat{P}_{n_I}, \hat{P}_{n_O}$	Instantaneous power constraints for interference-limited RAN $n_I$ and orthogonal RAN $n_O$
$a_{m,n_O,i}(t)$	RB allocation decision for RB $i$ in orthogonal RAN $n_O$ regarding UE $m$ in time slot $t$
$\mathcal{J}_{n_O}$	RB set in orthogonal RAN $n_O$
$\hat{\pi}_{n_I}$	Joint channel gain distribution of UEs over interference-limited RAN $n_I$
$\hat{W}_{n_I}, \hat{W}_{n_O}$	Maximal transmission rates for UEs in interference-limited RAN $n_I$ and orthogonal RAN $n_O$
$\hat{W}$	$\hat{W} = \max\{\max_{n_I}\{\hat{W}_{n_I}\}, \max_{n_O}\{\hat{W}_{n_O}\}\}$
$\hat{P}$	$\hat{P} = \max\{\max_{n_I}\{\hat{P}_{n_I}\}, \max_{n_O}\{\hat{P}_{n_O}\}\}$
$\mathcal{T}_E, \mathcal{T}_I$	Sets for elastic traffic and inelastic traffic
$T_E^k, T_I^k$	Index of elastic traffic $k$ and inelastic traffic $k$
$\beta_k$	Mean throughput constraint for inelastic traffic $T_I^k$
$\mathcal{M}_E, \mathcal{M}_I$	Sets for elastic UEs and inelastic UEs
$m_I^k$	Index of an inelastic UE requesting inelastic traffic $T_I^k$
$m_E$	Index of an elastic UE
$r_{m_I^k}^n(t)$	Admitted data for inelastic UE $m_I^k$ to RAN $n$ in time slot $t$
$r_{m_I^k}(t)$	Sum admitted data for inelastic UE $m_I^k$ in time slot $t$
$\bar{r}_{m_I^k}$	Time average admitted data for inelastic UE $m_I^k$
$\hat{r}_k$	Maximal admitted data constraint for inelastic traffic $T_I^k$
$\chi_{m_E,k}(t)$	Arrival packet number of elastic traffic $T_E^k$ for elastic UE $m_E$ in time slot $t$
$X_k$	Maximal arrival packet number of elastic traffic $T_E^k$
$f_k$	Packet size of elastic traffic $T_E^k$
$\gamma_{m_E,k}^n(t)$	Admitted packet number of elastic traffic $T_E^k$ for elastic UE $m_E$ in RAN $n$ in time slot $t$
$\bar{\gamma}_{m_E,k}^n$	Time-averaged admitted packet number of elastic traffic $T_E^k$ for elastic UE $m_E$ in RAN $n$
$R$	$\max\{\sum_{k \in \mathcal{T}_E} X_k f_k, \hat{r}_k, \hat{W}_{n_I}\}$
$Y_{m_I^k}(t)$	Virtual queue length for the time-averaged rate constraint of inelastic UE $m_I^k$ in time slot $t$
$\Theta_{n_I}(t), \Theta_{n_O}(t)$	Virtual power queue length for the average power constraint of interference-limited RAN $n_I$ and orthogonal RAN $n_O$ in time slot $t$
$U_{m,n}(t)$	Queue length of UE $m$ in RAN $n$ in time slot $t$
$\rho_{n_I}$	Minimal probability of $\mathbf{p}_{n_I}^{\text{NO}}(t)$ taking $\mathbf{p}_{n_I}^{\text{opt}}(t)$

### A. Interference-Limited RAN Model

In an interference-limited RAN, all UEs shared the same channel, and the basic radio resource is power. In particular, for RAN  $n_I$ , let  $g_{m,n_I}(t)$  be the channel gain for UE  $m$  in time slot  $t$ . Let  $p_{m,n_I}(t)$  be the RAN transmission power allocated

to UE  $m$  at time  $t$ . Then the service rate to UE  $m \in \mathcal{M}$  due to RAN  $n_I$  is

$$\mu_{m,n_I}(t) = W_{n_I} \log \left( 1 + \frac{\Gamma \Upsilon_{m,n_I} g_{m,n_I}(t) p_{m,n_I}(t)}{\sum_{k \in \mathcal{M} \setminus m} g_{m,n_I}(t) p_{k,n_I}(t) + \sigma^2} \right),$$

where  $W_{n_I}$  is the symbol rate of RAN  $n_I$ ,  $\Gamma$  is the capacity gap from the Shannon channel capacity,  $\Upsilon_{m,n_I}$  is the processing gain for the UE  $m$  in RAN  $n_I$ , and  $\sigma^2$  is the noise power.

Due to hardware limitation and regulatory requirements, both the average power and instantaneous power may be constrained. Let  $P_{n_I}$  and  $\hat{P}_{n_I}$  be the average and instantaneous power constraints, respectively. We must guarantee that  $\bar{p}_{n_I} = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{u=1}^t \sum_{m \in \mathcal{M}} p_{m,n_I}(u) \leq P_{n_I}$  and  $\sum_{m \in \mathcal{M}} p_{m,n_I}(t) \leq \hat{P}_{n_I}, \forall t$ .

### B. Orthogonal RAN Model

For an orthogonal RAN, borrowing from the OFDM terminology, the radio resources include resource blocks (RBs) and power. Let  $\mathcal{J}_{n_O}$  be the RB set of the orthogonal RAN  $n_O$ , and its number is written as  $J_{n_O}$ . Let binary variable  $a_{m,n_O,i}(t)$  be the indicator function for the allocation of RB  $i$  in RAN  $n_O$  to UE  $m$ . For RB  $i$  in RAN  $n_O$  at time  $t$ , let  $g_{m,n_O,i}(t)$  be the channel gain for UE  $m$  and  $p_{m,n_O,i}(t)$  be the RAN transmission power allocated to UE  $m$ . Then the transmission rate over RB  $i$  in RAN  $n_O$  is

$$\mu_{m,n_O,i}(t) = W_{n_O} \log \left( 1 + \Gamma \frac{g_{m,n_O,i}(t) p_{m,n_O,i}(t)}{\sigma^2} \right),$$

where  $W_{n_O}$  is the symbol rate of RAN  $n_O$ ,  $\Gamma$  is the capacity gap from the Shannon channel capacity, and  $\sigma^2$  is the noise power. The overall service rate to UE  $m$  due to RAN  $n_O$  is  $\mu_{m,n_O}(t) = \sum_{i \in \mathcal{J}_{n_O}} a_{m,n_O,i}(t) \mu_{m,n_O,i}(t)$ . Since each RB is allocated orthogonally in orthogonal RAN  $n_O$ , we have

$$\sum_{m \in \mathcal{M}} a_{m,n_O,i}(t) \leq 1, \forall i \in \mathcal{J}_{n_O}.$$

Similarly to the interference-limited RAN case, let  $P_{n_O}$  and  $\hat{P}_{n_O}$  be the average and instantaneous power constraints, respectively. We must guarantee that  $\bar{p}_{n_O} = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{u=1}^t \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_O}} a_{m,n_O,i}(u) p_{m,n_O,i}(u) \leq P_{n_O}$  and  $\sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_O}} a_{m,n_O,i}(t) p_{m,n_O,i}(t) \leq \hat{P}_{n_O}, \forall t$ .

### C. Channel State Processes

We assume that the channel gains,  $g_{m,n_I}(t)$  and  $g_{m,n_O,i}(t)$ , are discretized and can take value from a finite set  $\mathcal{G}_m \triangleq \{G_m^1, G_m^2, \dots, G_m^H\}$  of size  $H$ . We also suppose that the channel gains are independent among UEs. For simplify, we further denote  $\mathbf{g}_{n_O}(t) = [g_{m,n_O,j}(t)]_{M \times J_{n_O}}$ , and  $\mathbf{g}_{n_I}(t) = [g_{m,n_I}(t)]_{1 \times M}$ .

Given the above channel gain assumption, for the interference-limited RAN, we can find a positive constant  $\hat{W}_{n_I} = \max_m \{W_{n_I} \log(1 + \Gamma \Upsilon_{m,n_I} \frac{\max_j \{G_m^j\} P_{n_I}}{\sigma^2})\}$ , which satisfies  $\mu_{m,n_I}(t) \leq \hat{W}_{n_I}, \forall m \in \mathcal{M}, n_I \in \mathcal{R}_I$ . Similarly, for the orthogonal RAN, we also can have a positive constant  $\hat{W}_{n_O} = \max_m \{J_{n_O} W_{n_O} \log(1 + \Gamma \frac{\max_j \{G_m^j\} \hat{P}_{n_O}}{\sigma^2})\}$ , such that  $\mu_{m,n_O}(t) \leq \hat{W}_{n_O}, \forall m \in \mathcal{M}, n_O \in \mathcal{R}_O$ . These upper bounds will be used later in our analysis.

### D. UE Heterogeneous Traffic Model

An UE may request either inelastic or elastic traffic. Let  $\mathcal{T}_I = \{T_I^1, T_I^2, \dots, T_I^{K_I}\}$  be the set of inelastic traffic types of size  $K_I$ , and  $\mathcal{T}_E = \{T_E^1, T_E^2, \dots, T_E^{K_E}\}$  be the set of elastic traffic types of size  $K_E$ . For each inelastic traffic session  $T_I^k$ , it has a mean arrival bit rate  $\beta_k$  that the HWN is required to support. An elastic traffic session has no such requirement. Instead, each elastic traffic type  $T_E^k$  is distinguished by its fixed packet size  $f_k$ .

For notational simplicity, we divide the UEs in the HWN into the set of UEs with inelastic traffic, denoted by  $\mathcal{M}_I$  of size  $M_I$ , and the set of UEs with elastic traffic, denoted by  $\mathcal{M}_E$  of size  $M_E$ , so that  $\mathcal{M}_E \cup \mathcal{M}_I = \mathcal{M}$ . We assume that each inelastic UE requests only one inelastic traffic type from  $\mathcal{T}_I$ , so we may denote  $m_I^k$  as a UE that requests inelastic traffic type  $T_I^k$ . In contrast, each elastic UE can request one or multiple elastic traffic types from  $\mathcal{T}_E$ , and we denote an elastic UE as  $m_E$ . Such division of UEs is without loss of generality, since a UE that requests both inelastic and elastic traffic may be modeled as multiple UEs that satisfy the above conditions.

For inelastic UE  $m_I^k$ , let  $r_{m_I^k}^n(t)$  be the admitted data rate of its required inelastic traffic type  $T_I^k$  to RAN  $n \in \mathcal{R}$  (orthogonal or interference-limited) at time  $t$ . Then the sum admitted rate for  $m_I^k$  is  $r_{m_I^k}(t) = \sum_{n \in \mathcal{R}} r_{m_I^k}^n(t)$ . Since the average admitted rate for inelastic traffic should be greater than its arrival rate, we require  $\bar{r}_{m_I^k} = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{u=1}^t r_{m_I^k}(u) \geq \beta_k$ . Furthermore, to avoid dealing with the trivial case of infinite admission creating queue instability, we assume an upper bound  $\hat{r}_k$  to the admitted data, sufficiently large such that  $r_{m_I^k}(t) \leq \hat{r}_k, \forall t, m_I^k \in \mathcal{M}_I, k \in \mathcal{T}_I$ .

For elastic UE  $m_E$ , let  $\chi_{m_E,k}(t)$  be the number of arriving packets of elastic traffic type  $T_E^k$  at time  $t$ . We assume that  $\chi_{m_E,k}(t)$  is random but upper bounded by  $X_k$ , and that the packets are of constant size  $f_k$ . Let  $\gamma_{m_E,k}^n(t)$  be the number of packets of type  $T_E^k$  admitted to RAN  $n \in \mathcal{R}$  (orthogonal or interference-limited) for UE  $m_E$  at time  $t$ . Then, the total amount of admitted data for UE  $m_E$  to RAN  $n \in \mathcal{R}$  is  $\sum_{k \in \mathcal{T}_E} \gamma_{m_E,k}^n(t) f_k$ . Note that the admitted amount of traffic should be less than the arriving of traffic in each time slot, so we have  $\sum_{n \in \mathcal{R}} \gamma_{m_E,k}^n(t) \leq \chi_{m_E,k}(t), \forall t, m_E \in \mathcal{M}_E, k \in \mathcal{T}_E$ .

### E. Queue Updating at RANs

A separate queue is maintained for each UE  $m \in \mathcal{M}$  in RAN  $n \in \mathcal{R}$ . Let  $U_{m,n}(t)$  be its queue length at time slot  $t$ . Then, the queue updating function for inelastic UE  $m_I^k$  in RAN  $n \in \mathcal{R}$  is expressed as

$$U_{m_I^k,n}(t+1) = \max\{U_{m_I^k,n}(t) - \mu_{m_I^k,n}(t), 0\} + r_{m_I^k}^n(t), \quad (1)$$

and the queue updating function for elastic UE  $m_E$  is

$$U_{m_E,n}(t+1) = \max\{U_{m_E,n}(t) - \mu_{m_E,n}(t), 0\} + \sum_{k \in \mathcal{T}_E} \gamma_{m_E,k}^n(t) f_k. \quad (2)$$

### F. Problem Statement

We aim to properly leverage the benefits provided by an HWN through jointly optimal traffic control and radio resource

allocation within each RAN and across multiple RANs. Our objective is to maximize the average utility for elastic traffic, while guaranteeing the time-averaged throughput for inelastic traffic. Summarizing the system model presented in the previous subsections, this problem is formulated as follows:

$$\max_{\mathbf{p}_{n_0}(t), \mathbf{p}_{n_1}(t), \mathbf{r}(t), \mathbf{a}_{n_0}(t), \gamma(t)} \left\{ \sum_{m_E \in \mathcal{M}_E} \sum_{k \in \mathcal{T}_E} \sum_{n \in \mathcal{R}} \bar{\gamma}_{m_E,k}^n \Psi(f_k) \right\} \quad (3)$$

s. t.

$$\bar{r}_{m_I^k} \geq \beta_k, \quad \forall m_I^k \in \mathcal{M}_I, k \in \mathcal{T}_I, \quad (4)$$

$$\bar{p}_{n_0} \leq P_{n_0}, \quad \forall n_0 \in \mathcal{R}_O, \quad (5)$$

$$\bar{p}_{n_1} \leq P_{n_1}, \quad \forall n_1 \in \mathcal{R}_I, \quad (6)$$

$$\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{u=1}^t \left[ \sum_{m \in \mathcal{M}} \sum_{n \in \mathcal{R}} \mathbb{E}\{U_{m,n}(u)\} \right] < \infty, \quad (7)$$

$$\sum_{n \in \mathcal{R}} r_{m_I^k}^n(t) \leq \hat{r}_k, \quad \forall m_I^k \in \mathcal{M}_I, k \in \mathcal{T}_I, \quad (8)$$

$$r_{m_I^k}^n(t) \geq 0, \quad \forall n \in \mathcal{R}, m_I^k \in \mathcal{M}_I, k \in \mathcal{T}_I, \quad (9)$$

$$\sum_{n \in \mathcal{R}} \gamma_{m_E,k}^n(t) \leq \chi_k(t), \quad \forall m_E \in \mathcal{M}_E, k \in \mathcal{T}_E, \quad (10)$$

$$\gamma_{m_E,k}^n(t) \geq 0, \quad \forall m_E \in \mathcal{M}_E, n \in \mathcal{R}, k \in \mathcal{T}_E, \quad (11)$$

$$\sum_{m \in \mathcal{M}} p_{m,n_1}(t) \leq \hat{P}_{n_1}, \quad \forall n_1 \in \mathcal{R}_I, \quad (12)$$

$$p_{m,n_1}(t) \geq 0, \quad \forall m \in \mathcal{M}, n_1 \in \mathcal{R}_I, \quad (13)$$

$$\sum_{m \in \mathcal{M}} a_{m,n_0,i}(t) \leq 1, \quad \forall i \in \mathcal{J}_{n_0}, n_0 \in \mathcal{R}_O, \quad (14)$$

$$a_{m,n_0,i} \in \{0, 1\}, \quad \forall m \in \mathcal{M}, i \in \mathcal{J}_{n_0}, n_0 \in \mathcal{R}_O, \quad (15)$$

$$\sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} a_{m,n_0,i}(t) p_{m,n_0,i}(t) \leq \hat{P}_{n_0}, \quad \forall n_0 \in \mathcal{R}_O, \quad (16)$$

$$p_{m,n_0,i}(t) \geq 0, \quad \forall m \in \mathcal{M}, i \in \mathcal{J}_{n_0}, n_0 \in \mathcal{R}_O, \quad (17)$$

where  $\bar{\gamma}_{m_E,k}^n = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{u=1}^t \gamma_{m_E,k}^n(u)$ ,  $\mathbf{r}(t) = [r_{m_I^k,n}^k(t)]_{M_I \times N}$ ,  $\mathbf{p}_{n_0}(t) = [p_{m,n_0,i}(t)]_{M \times J_{n_0}}$ ,  $\mathbf{p}_{n_1}(t) = [p_{m,n_1}(t)]_{M \times 1}$ ,  $\mathbf{a}_{n_0}(t) = [a_{m,n_0,i}(t)]_{M \times J_{n_0}}$ , and  $\gamma(t) = [\gamma_{m_E,k}^n(t)]_{M_E \times N \times K_E}$ .  $\Psi(f_k)$  is the utility for servicing one packet of elastic traffic type  $T_E^k$ . Given  $f_k$ ,  $\Psi(f_k)$  is a constant, which reflects the preference of the network operator towards elastic traffic type  $T_E^k$ .

Note that (4), (5), (6), and (7), respectively, are constraints on the inelastic traffic rate, orthogonal RAN power, interference-limited RAN power, and queue stability, all in the time average sense. Inequalities (8) and (9) are instantaneous rate constraints on inelastic traffic, and (10) and (11) are instantaneous packet constraints on elastic traffic. Inequalities (12) and (13) are instantaneous constraints on the radio resource of interference-limited RANs, while (14), (15), (16), and (17) are instantaneous constraints on the radio resource of orthogonal RANs.

The above optimization is a challenging cross-network and cross-layer problem, because of the high levels of correlation among the different types of RANs and among the different types of traffic flows in an HWN. Furthermore, due to the binary RB allocation variable  $a_{m,n_0,i}(t)$  in orthogonal RANs and the non-convexity of the SINR term in interference-limited

RANs, this problem is a non-convex mixed-integer program, which generally has prohibitive computational complexity. However, we show next that a network-separable joint optimization approach can efficiently solve the problem.

#### IV. ADAPTIVE NETWORK-SEPARABLE METHOD FOR CROSS-NETWORK CROSS-LAYER OPTIMIZATION

To solve the above optimization problem, we propose an on-line adaptive network-separable method to jointly design the traffic admission control strategy and the radio resource allocation scheme within and cross RANs. It allows each RAN to design its own resource allocation independently of the other RANs, and only requires the RANs to report queue lengths to the traffic admission controller. Furthermore, we observe that the computationally prohibitive non-convex optimization sub-problem in interference-limited RANs, do not need to be solved optimally in every time slot. Instead, we propose novel methods to find its optimal solution with a positive probability, which is then used to derive a globally optimal solution.

##### A. Network-Separable Reformulation

We first demonstrate how the original problem can be reformulated into a network-separable form, so that each RAN may design its own resource allocation independently of other parts of the system.

We adapt a Lyapunov optimization framework [2] to our problem, by constructing three virtual queues to accommodate the constraints (4), (5) and (6). For the average-rate constraint (4), the corresponding virtual queue has arrival and departure rates  $\beta_k$  and  $r_{m_1^k}(t)$  respectively, and its queue length in time slot  $t$  is denoted as  $Y_{m_1^k}(t)$ . The queue updating function is then expressed as

$$Y_{m_1^k}(t+1) = \max\{Y_{m_1^k}(t) - r_{m_1^k}(t), 0\} + \beta_k.$$

For the average-power constraints (5) and (6), the corresponding virtual queue has arrival rates  $p_{n_0}(t)$  and  $p_{n_1}(t)$ , and departure rates  $P_{n_0}$  and  $P_{n_1}$  respectively, and their queue lengths in time slot  $t$  are denoted as  $\Theta_{n_0}(t)$  and  $\Theta_{n_1}(t)$  respectively. The queue updating function is then expressed as

$$\Theta_{n_0}(t+1) = \max\{\Theta_{n_0}(t) - P_{n_0}, 0\} + p_{n_0}(t) \quad (18)$$

and

$$\Theta_{n_1}(t+1) = \max\{\Theta_{n_1}(t) - P_{n_1}, 0\} + p_{n_1}(t). \quad (19)$$

The constraints (4), (5), and (6) are satisfied if these three queues are mean rate stable, which indicates that the input rate is below the service rate.

We then construct a Lyapunov drift-plus-penalty function [2] with respect to the network utility objective (3) and the

constraints (4), (5), (6), and (7) as follows:

$$\begin{aligned} \Delta(\mathbf{Q}(t)) = & \mathbb{E}\left\{\frac{1}{2} \sum_{m \in \mathcal{M}} \sum_{n \in \mathcal{R}} [U_{m,n}(t+1)^2 - U_{m,n}(t)^2] + \right. \\ & \frac{1}{2} \sum_{m_1^k \in \mathcal{M}_1} [Y_{m_1^k}(t+1)^2 - Y_{m_1^k}(t)^2] + \frac{1}{2} \sum_{n_0 \in \mathcal{R}_0} [\Theta_{n_0}(t+1)^2 - \\ & \Theta_{n_0}(t)^2] + \frac{1}{2} \sum_{n_1 \in \mathcal{R}_1} [\Theta_{n_1}(t+1)^2 - \Theta_{n_1}(t)^2] | \mathbf{Q}(t) \Big\} - \\ & V \mathbb{E}\left\{ \sum_{m \in \mathcal{M}_E} \sum_{k \in \mathcal{T}_E} \sum_{n \in \mathcal{R}} \gamma_{m_E, k}^n \Psi(f_k) | \mathbf{Q}(t) \right\}, \end{aligned}$$

where  $\mathbf{U}(t) = [U_{m,n}(t)]_{M \times N}$ ,  $\mathbf{Y}(t) = [Y_{m_1^k}(t)]_{M_1 \times 1}$ ,  $\mathbf{Q}(t) = [\mathbf{U}(t), \mathbf{Y}(t), \Theta_{n_1}(t), \Theta_{n_0}(t)]$ , and  $V$  is an arbitrary positive constant.

Similar to the general derivations given in [2], it is easy to show that the above is upper bounded by

$$\begin{aligned} & \mathbb{E}\left\{ \sum_{m_1^k \in \mathcal{M}_1} \sum_{n \in \mathcal{R}} (U_{m_1^k, n}(t) - Y_{m_1^k}(t)) r_{m_1^k, n}(t) + \right. \\ & \sum_{m \in \mathcal{M}_E} \sum_{n \in \mathcal{R}} \sum_{k \in \mathcal{T}_E} (U_{m_E, n}(t) f_k - V \Psi(f_k)) \gamma_{m_E, k}^n(t) + \\ & \sum_{n_0 \in \mathcal{R}_0} \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} a_{m, n_0, i}(t) \{\Theta_{n_0}(t) p_{m, n_0, i}(t) - U_{m, n_0}(t) \mu_{m, n_0, i}(t)\} \\ & + \sum_{n_1 \in \mathcal{R}_1} \sum_{m \in \mathcal{M}} \{\Theta_{n_1}(t) p_{m, n_1}(t) - U_{m, n_1}(t) \mu_{m, n_1}(t)\} | \mathbf{Q}(t) \Big\} - \\ & \sum_{n \in \mathcal{R}_0} \Theta_{n_0}(t) P_{n_0} - \sum_{n \in \mathcal{R}_1} \Theta_{n_1}(t) P_{n_1} + \sum_{m_1^k \in \mathcal{M}_1} Y_{m_1^k}(t) \beta_k + B, \end{aligned} \quad (20)$$

where  $B$  is a constant given by  $B = MN\hat{W}^2 + N\hat{P}^2$  with  $\hat{P} = \max\{\max_{n_1} \{\hat{P}_{n_1}\}, \max_{n_0} \{\hat{P}_{n_0}\}\}$  and  $\hat{W} = \max\{\max_{n_1} \{\hat{W}_{n_1}\}, \max_{n_0} \{\hat{W}_{n_0}\}\}$ . The derivations are omitted for brevity.

Furthermore, it has been proven [2] that any strategy that minimizes this upper bound, given  $\mathbf{Q}(t)$ , also solves the original optimization problem, with  $V$  as a tuning parameter that determines the tradeoff between optimization performance and queueing delay. Therefore, the goal of our joint admission control strategy and radio resource allocation scheme is transformed to adaptively minimizing (20) given  $\mathbf{Q}(t)$  in each time slot, subject to (8)-(17). By doing so, our network control method does not require the statistical information of elastic traffic arrivals and channel fading, and can adapt to changes in network statistics.

Here we observe some important features of (20). First, the last four terms are constants given  $\mathbf{Q}(t)$ . Second, to minimize the conditional expectation, it suffices to minimize what is inside the expectation operator for each random realization of the system and given  $\mathbf{Q}(t)$ . Third, and most importantly, each of the four terms inside the expectation operator can be minimized independently of the other terms. Furthermore, we note that these four terms have the following physical interpretation:

- The first term corresponds to the admission of inelastic traffic.

- The second term corresponds to the admission of elastic traffic.
- The third term can be minimized separately by each orthogonal RAN.
- The fourth term can be minimized separately by each interference-limited RAN.

Therefore, the proposed reformulation achieves complete network separation, permitting a distributed approach to solve the original optimization problem. The handling of each of these terms are described in the following subsections.

### B. Optimal Cross-Network Heterogeneous Traffic Admission Control Strategy

The minimization of the first two terms of (20) is performed by the traffic admission controller. We need to deliberately satisfy the UE traffic requirement as well as to balance the RAN load. Intuitively, the optimal traffic admission control strategy is to exploit the cross-network design to obtain the best-service assignment gain.

1) *Inelastic traffic*: The optimal admission control decision for inelastic traffic is based on minimizing the first term in (20), which reflects the inelastic traffic requirement and the RAN load through  $Y_{m_1^k}(t)$  and  $U_{m_1^k,n}(t)$  separately. In each time slot  $t$ , it suffices to minimize  $(U_{m_1^k,n}(t) - Y_{m_1^k}(t))r_{m_1^k,n}^*(t)$ . Hence, the following amount of traffic for  $m_1^k$  is admitted for RAN  $n$ :

$$r_{m_1^k,n}^*(t) = \begin{cases} \hat{r}_k, & \text{if } n = \operatorname{argmin}_{n'} \{U_{m_1^k,n'}(t) - Y_{m_1^k}(t)\} \\ & \text{and } U_{m_1^k,n}(t) - Y_{m_1^k}(t) < 0 \\ 0, & \text{otherwise} \end{cases}. \quad (21)$$

We can take  $U_{m_1^k,n}(t)$  and  $Y_{m_1^k}(t)$  as the congestion cost and income of RAN  $n$  for serving one bit of inelastic traffic  $T_1^k$ . Then, (21) ensures that the inelastic traffic is injected to the RAN with the highest net income. Thus, the above inelastic traffic distribution scheme aims to achieve the minimum congestion cost for serving inelastic UEs.

2) *Elastic traffic*: The optimal admission control strategy for elastic UE  $m_E$  is based on minimizing the second term in (20) under constraints (10) and (11). Since the corresponding optimization problem is linear, its optimal solution can be written as

$$\gamma_{m_E,k}^{n,*} = \begin{cases} \chi_{m_E,k}(t), & \text{if } n = \operatorname{argmin}_{n'} \{U_{m_E,n'}(t)f_k - V\Psi(f_k)\} \\ & \text{and } U_{m_E,n}(t)f_k - V\Psi(f_k) < 0 \\ 0, & \text{otherwise} \end{cases}. \quad (22)$$

For the above elastic traffic admission control strategy, we can take  $U_{m_E,n}(t)f_k$  and  $V\Psi(f_k)$  as the congestion cost and income of RAN  $n$  for serving one packet of elastic traffic  $T_E^k$ . Then, (22) ensures that elastic traffic is injected to the RAN with the highest net income. Thus, the above elastic traffic distribution scheme aims to achieve the minimum congestion cost for serving elastic UEs.

### C. Optimal Radio Resource Allocation in Orthogonal RANs at Each Time Slot

The minimization of the third term of (20) is performed by orthogonal RANs through power and RB allocation. For each orthogonal RAN  $n_O$ , we rewrite this term as

$$F_{n_O}(\mathbf{a}_{n_O}(t), \mathbf{p}_{n_O}(t) | \mathbf{U}_{n_O}(t), \Theta_{n_O}(t)) \triangleq \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_O}} a_{m,n_O,i}(t) \left[ \Theta_{n_O}(t) p_{m,n_O,i}(t) - U_{m,n_O}(t) W_{n_O} \log\left(1 + \Gamma \frac{g_{m,n_O,i}(t) p_{m,n_O,i}(t)}{\sigma^2}\right) \right].$$

Then, to determine the optimal power and RB allocation decision, we need to solve the following optimization problem:

$$\begin{aligned} & \min_{\substack{\mathbf{a}_{n_O}(t) \\ \mathbf{p}_{n_O}(t)}}} F_{n_O}(\mathbf{a}_{n_O}(t), \mathbf{p}_{n_O}(t) | \mathbf{U}_{n_O}(t), \Theta_{n_O}(t)) \\ \text{s. t. } & \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_O}} a_{m,n_O,i}(t) p_{m,n_O,i}(t) \leq \hat{P}_{n_O}, \\ & a_{m,n_O,i}(t) \in \{0, 1\}, \quad \forall m \in \mathcal{M}, i \in \mathcal{J}_{n_O}, \\ & \sum_{m \in \mathcal{M}} a_{m,n_O,i}(t) \leq 1, \quad \forall i \in \mathcal{J}_{n_O}, \\ & p_{m,n_O,i}(t) \geq 0, \quad \forall m \in \mathcal{M}, i \in \mathcal{J}_{n_O}. \end{aligned} \quad (23)$$

Because of the binary RB allocation variables  $a_{m,n_O,i}(t)$ , the radio resource allocation problem in the orthogonal RAN is a mixed-integer non-linear program, which typically has prohibitive computational complexity. Similar joint power and RB allocation problem can be found in [31] without consideration of the resource granularity, and in [32] with RB and power being optimized separately. Our approach is based on the dual decomposition technique in [33], [34], [35], [36]. However, due to the average power constrain in the orthogonal RAN, which is expressed in terms of  $\Theta_{n_O}$ , the method in [33], [34], [35], [36] cannot be directly applied. Fortunately, the special structure of our optimization problem allows us to derive an optimal solution.

1) *Continuity relaxation and convexification*: We first relax  $a_{m,n_O,i}(t)$  to the continuous interval  $[0, 1]$  and further introduce a new variable  $s_{m,n_O,i}(t) = a_{m,n_O,i}(t)p_{m,n_O,i}(t)$  for each UE  $m$  and RB  $i$  in RAN  $n_O$ . Then we can rewrite (23) as <sup>1</sup>

$$\min_{\substack{\mathbf{a}_{n_O} \\ \mathbf{s}_{n_O}}} \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_O}} \left\{ s_{m,n_O,i} \Theta_{n_O} - U_{m,n_O} W_{n_O} a_{m,n_O,i} \times \log\left(1 + \Gamma \frac{g_{m,n_O,i} s_{m,n_O,i}}{\sigma^2 a_{m,n_O,i}}\right) \right\} \quad (24)$$

$$\text{s. t. } \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_O}} s_{m,n_O,i} \leq \hat{P}_{n_O}, \quad (25)$$

$$s_{m,n_O,i} \geq 0, \quad \forall m \in \mathcal{M}, i \in \mathcal{J}_{n_O}, \quad (26)$$

$$0 \leq a_{m,n_O,i} \leq 1, \quad \forall m \in \mathcal{M}, i \in \mathcal{J}_{n_O}, \quad (27)$$

$$\sum_{m \in \mathcal{M}} a_{m,n_O,i} \leq 1, \quad \forall i \in \mathcal{J}_{n_O}, \quad (28)$$

where  $\mathbf{s}_{n_O} = [s_{m,n_O,i}]_{M \times J_{n_O}}$ .

<sup>1</sup>For notation simplicity, the time index  $t$  is omitted when it is clear from the context.

It is easy to verify that the above optimization problem is convex, since it is the sum of a linear function and the perspective function of a convex log function. Furthermore, since all the constraints are affine functions, the Slater's condition is always satisfied, leading to a zero Lagrange duality gap [37]. In the following, for notation simplicity, we take  $\eta_{m,n_0} = U_{m,n_0} W_{n_0}$  and  $\xi_{m,n_0,i} = \frac{\Gamma_{m,n_0,i}}{\sigma^2}$ .

2) *Lagrange dual solution*: We relax the power constraint (25) by introducing the dual variable  $\lambda_{n_0}$ , obtaining the following Lagrangian:

$$L(\mathbf{s}_{n_0}, \mathbf{a}_{n_0}, \lambda_{n_0}) = \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} \left\{ s_{m,n_0,i} \Theta_{n_0} - \eta_{m,n_0} a_{m,n_0,i} \times \log\left(1 + \xi_{m,n_0,i} \frac{s_{m,n_0,i}}{a_{m,n_0,i}}\right) \right\} + \lambda_{n_0} \left( \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} s_{m,n_0,i} - \hat{P}_{n_0} \right), \quad (29)$$

and the corresponding Lagrange dual:

$$d(\lambda_{n_0}) = \min_{\mathbf{s}_{n_0}, \mathbf{a}_{n_0}} L(\mathbf{s}_{n_0}, \mathbf{a}_{n_0}, \lambda_{n_0}) \quad (30)$$

s. t. (26) – (28).

A standard solution approach is to solve the following the Lagrange dual problem

$$\max_{\lambda_{n_0}} d(\lambda_{n_0}) \quad (31)$$

s. t.  $\lambda_{n_0} \geq 0$ .

It is well known that the maximum of (31) equals the minimum of (24)-(28) [37]. Therefore, it remains to solve (30) and (31) to find the optimum  $\lambda_{n_0}$ ,  $s_{n_0}$ , and  $a_{n_0}$ .

To derive the solution to (30), we first give the optimal relation between  $\mathbf{a}_{n_0}$  and  $\mathbf{s}_{n_0}$  in Lemma 1.

**Lemma 1**: The optimal relation between  $s_{m,n_0,i}$  and  $a_{m,n_0,i}$  can be written as

$$s_{m,n_0,i} = \left[ \frac{\eta_{m,n_0}}{(\Theta_{n_0} + \lambda_{n_0}) \ln 2} - \frac{1}{\xi_{m,n_0,i}} \right]^+ a_{m,n_0,i}, \quad (32)$$

where  $[x]^+ \triangleq \max\{x, 0\}$ .

*Proof*: See Appendix A. ■

Then, substituting (32) into (30) and letting

$$\Lambda_{m,n_0,i}(\lambda_{n_0}) = \left[ \frac{\eta_{m,n_0}}{(\Theta_{n_0} + \lambda_{n_0}) \ln 2} - \frac{1}{\xi_{m,n_0,i}} \right]^+ \Theta_{n_0} - \eta_{m,n_0} \log \left( 1 + \xi_{m,n_0,i} \left[ \frac{\eta_{m,n_0}}{(\Theta_{n_0} + \lambda_{n_0}) \ln 2} - \frac{1}{\xi_{m,n_0,i}} \right]^+ \right) + \lambda_{n_0} \left[ \frac{\eta_{m,n_0}}{(\Theta_{n_0} + \lambda_{n_0}) \ln 2} - \frac{1}{\xi_{m,n_0,i}} \right]^+,$$

we can rewrite (30) as follows:

$$d(\lambda_{n_0}) = \min_{\mathbf{a}_{n_0}} \left\{ \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} a_{m,n_0,i} \Lambda_{m,n_0,i}(\lambda_{n_0}) - \lambda_{n_0} \hat{P}_{n_0} \right\} \quad (33)$$

s. t.

$$\sum_{m \in \mathcal{M}} a_{m,n_0,i} \leq 1, \forall i \in \mathcal{J}_{n_0},$$

$$0 \leq a_{m,n_0,i} \leq 1, \forall m \in \mathcal{M}, i \in \mathcal{J}_{n_0}.$$

This is a classical linear assignment problem. It is easy to see that an optimal solution is

$$a_{m,n_0,i} = \begin{cases} 1, & m = \operatorname{argmin}_l \{\Lambda_{l,n_0,i}(\lambda_{n_0})\} \text{ and} \\ & \Lambda_{m,n_0,i}(\lambda_{n_0}) < 0 \\ 0, & \text{otherwise} \end{cases}. \quad (34)$$

To solve the dual optimization problem (31), we note that the dual function  $d(\lambda_{n_0})$  is concave (since it is a minimum of linear functions) and there is only one dual variable. Thus, we could exploit an efficient one-dimensional search method to obtain its optimal solution. Since we can efficiently calculate the dual function for any given  $\lambda_{n_0}$ , we may take the golden-section search method [38]. We first give an upper bound for the optimal  $\lambda_{n_0}^*$ .

**Lemma 2**: The optimal  $\lambda_{n_0}^*$  is upper bounded by  $\max\{\max_{m \in \mathcal{M}, i \in \mathcal{J}_{n_0}} \left\{ \frac{\xi_{m,n_0,i} \eta_{m,n_0}}{\ln 2} - \Theta_{n_0} \right\}, 0\}$ .

*Proof*: See Appendix B. ■

Given the above upper bound for  $\lambda_{n_0}^*$ , for a fixed constant  $\epsilon$ , if  $\max_{m \in \mathcal{M}, i \in \mathcal{J}_{n_0}} \left\{ \frac{\xi_{m,n_0,i} \eta_{m,n_0}}{\ln 2} - \Theta_{n_0} \right\} > 0$ , the number of iteration steps required to ensure  $|\lambda_{n_0} - \lambda_{n_0}^*| < \epsilon$  by the golden-section search method is  $\log_{0.618} \frac{\epsilon}{\max_{m \in \mathcal{M}, i \in \mathcal{J}_{n_0}} \left\{ \frac{\xi_{m,n_0,i} \eta_{m,n_0}}{\ln 2} - \Theta_{n_0} \right\}}$ . Otherwise, if

$\max_{m \in \mathcal{M}, i \in \mathcal{J}_{n_0}} \left\{ \frac{\xi_{m,n_0,i} \eta_{m,n_0}}{\ln 2} - \Theta_{n_0} \right\} < 0$ , we directly have  $\lambda_{n_0}^* = 0$ .

3) *Tie-breaking and primal recovery*: Given the optimal dual variable  $\lambda_{n_0}^*$ , it remains non-trivial to find optimal solutions to (24)-(28),  $\mathbf{a}_{n_0}^*$  and  $\mathbf{s}_{n_0}^*$ . We define  $\Lambda_{n_0,i}^*(\lambda_{n_0}^*) = \min_{l \in \mathcal{M}} \{\Lambda_{l,n_0,i}(\lambda_{n_0}^*)\}$ . Let us denote  $\mathcal{C}_{n_0,i} = \{m : \Lambda_{m,n_0,i}(\lambda_{n_0}^*) = \Lambda_{n_0,i}^*(\lambda_{n_0}^*)\}$ . If  $|\mathcal{C}_{n_0,i}| = 1$ , we have exactly one minimizer for RB  $j$  of RAN  $n_0$  in (33). Then, to minimize (33), we only need to allocate RB  $i$  to UE  $m \in \mathcal{C}_{n_0,i}$ , i.e., setting  $a_{m,n_0,i}^* = 1$ . Since the optimal solution  $\mathbf{a}_{n_0}^*$  and  $\mathbf{s}_{n_0}^*$  to (30) is unique, they are also optimal for (24)-(28) [37].

However, we may have  $|\mathcal{C}_{n_0,i}| > 1$ . In other words, we need to break ties for the allocation of RB  $i$ . Specifically, for the solution to (33) to also optimize (24)-(28), the tie-breaking rule must comply with the following two conditions:

1) **Primal Feasibility**:  $\sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} s_{m,n_0,i}(t) \leq \hat{P}_{n_0}$ ;

2) **Complementary Slackness Condition**:

$$\lambda_{n_0}^* (\sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} s_{m,n_0,i} - \hat{P}_{n_0}) = 0.$$

Finding such a tie-breaking rule can be computational intractable, since fractional RB usage is allowed in (33). Here we exploit the maximal and minimal subgradient of  $d(\lambda_{n_0})$  at  $\lambda_{n_0}^*$  as in [34], [39].

We first denote  $\mathcal{X}_{n_0} = \{i \in \mathcal{J}_{n_0} : |\mathcal{C}_{n_0,i}| = 1\}$ , and let  $\mathcal{X}_{n_0}^c = \mathcal{J}_{n_0} \setminus \mathcal{X}_{n_0}$ . For RB  $i \in \mathcal{X}_{n_0}^c$ , we could break the tie by exhaustively and exclusively allocating RB  $i$  to exactly one UE  $m \in \mathcal{C}_{n_0,i}$ , and this will lead to  $\prod_{i \in \mathcal{X}_{n_0}^c} |\mathcal{C}_{n_0,i}|$  RB allocation results, whose convex combination can be further used to express all the other RB and power allocation results. Note that the above  $\prod_{i \in \mathcal{X}_{n_0}^c} |\mathcal{C}_{n_0,i}|$  allocation results may include the optimal power and RB allocation decision in terms of satisfying the primal feasibility and complementary slackness condition. If this is true, we just simply use the corresponding optimal allocation decision, and we are done.

However, the optimal allocation decision may not be among them, since the optimal RB allocation may be fractional. In this case, we can find two allocation decisions among these allocation results,  $\{\mathbf{a}_{n_0}^{\max}(t), \mathbf{p}_{n_0}^{\max}(t)\}$  and  $\{\mathbf{a}_{n_0}^{\min}(t), \mathbf{p}_{n_0}^{\min}(t)\}$ , which maximizes and minimizes the subgradient of  $d(\lambda_{n_0})$  at  $\lambda_{n_0}^*$ ,  $\Omega(\lambda_{n_0}^*) = \sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} s_{m,n_0,i} - \hat{P}_{n_0}$ , respectively. We denote the maximal and the minimal subgradients as  $\Omega_{\max}(\lambda_{n_0}^*)$  and  $\Omega_{\min}(\lambda_{n_0}^*)$  respectively.

Note that at the optimal  $\lambda_{n_0}^*$ ,  $\Omega_{\max}(\lambda_{n_0}^*)$  must be positive and  $\Omega_{\min}(\lambda_{n_0}^*)$  must be negative. This is because, for optimal  $\lambda_{n_0}^*$ , there always exists a fractional RB allocation such that all available power is used, i.e.,  $\sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} s_{m,n_0,i} - \hat{P}_{n_0} = 0$ , and such a fractional RB allocation is necessarily a linear combination of  $\Omega(\lambda_{n_0}^*)$  with binary  $\mathbf{a}_{n_0}$ . Then, a time sharing scheme can be adopted to find a zero subgradient based on  $\Omega_{\max}(\lambda_{n_0}^*)$  and  $\Omega_{\min}(\lambda_{n_0}^*)$ . That is we can find  $0 < z_{n_0} < 1$  such that

$$z_{n_0} \Omega_{\max}(\lambda_{n_0}^*) + (1 - z_{n_0}) \Omega_{\min}(\lambda_{n_0}^*) = 0.$$

Then we use allocation decisions,  $\{\mathbf{a}_{n_0}^{\max}(t), \mathbf{p}_{n_0}^{\max}(t)\}$  and  $\{\mathbf{a}_{n_0}^{\min}(t), \mathbf{p}_{n_0}^{\min}(t)\}$ , with the fraction of time  $z_{n_0}$  and  $1 - z_{n_0}$  separately, at time slot  $t$ .

#### 4) Optimality with Respect to (23):

**Lemma 3:** An optimal RB and power allocation to (23) is  $\mathbf{a}_{n_0}^*$  and  $\mathbf{p}_{n_0}^* = \mathbf{s}_{n_0}^*$ .

*Proof:* See Appendix C. ■

Intuitively, the above radio resource allocation strategy in the orthogonal RAN is to optimally match the RBs to the UEs under the sum power constraint, which is to exploit the multi-user diversity gain and the spatial transmission diversity gain.

### D. Probabilistically Optimal Radio Resource Allocation in Interference-Limited RANs at Each Time Slot

The minimization of the fourth term of (20) is performed by interference-limited RANs through power allocation. Due to the non-convex nature of this subproblem, we first provide a suboptimal solution and later in Section IV-E show how it can be used to achieve optimality in the overall problem.

For each interference-limited RAN  $n_1$ , we rewrite the third term of (20) as

$$F_{n_1}(\mathbf{p}_{n_1}(t) | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t)) \triangleq \sum_{m \in \mathcal{M}} \left\{ \Theta_{n_1}(t) p_{m,n_1}(t) - U_{m,n_1}(t) W_{n_1} \log(1 + \Gamma \Upsilon_{m,n_1} \text{SINR}_{m,n_1}(t)) \right\},$$

where

$$\text{SINR}_{m,n_1}(t) = \frac{g_{m,n_1}(t) p_{m,n_1}(t)}{\varphi_m(\mathbf{p}_{n_1}(t))}$$

with  $\varphi_m(\mathbf{p}_{n_1}(t)) = \sum_{k \in \mathcal{M} \setminus m} g_{m,n_1}(t) p_{k,n_1}(t) + \sigma^2$ . Then, to determine the optimal power allocation in each time slot  $t$ , we need to solve the following optimization problem:

$$\begin{aligned} & \min_{\mathbf{p}_{n_1}(t)} F_{n_1}(\mathbf{p}_{n_1}(t) | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t)) \\ \text{s. t. } & \sum_{m \in \mathcal{M}} p_{m,n_1}(t) \leq \hat{P}_{n_1}, \\ & p_{m,n_1}(t) \geq 0, \forall m \in \mathcal{M}. \end{aligned} \quad (35)$$

We denote its global optimizer as  $\mathbf{p}_{n_1}^{\text{opt}}(t)$ .

From the above formulation, it is clear that the optimal radio resource allocation in each RAN is determined by the queue lengths  $\mathbf{U}_{n_1}(t)$  and  $\Theta_{n_1}(t)$ . We can see that  $\mathbf{U}_{n_1}(t)$  coordinates the admission control strategy and the radio resource allocation strategy in the interference-limited RANs.

Due to the non-convexity of the SINR term, the above optimization problem is NP-hard [1]. A typical technique is to approximate  $\log(1 + \text{SINR}_{m,n_1})$  as  $\log(\text{SINR}_{m,n_1})$  [40], but it is appropriate only in the high SINR regime. Other known heuristics include iterative water filling [41] and asymptotic Lagrange duality either with small tone spacing [42] or with an infinite number of sub-carriers [43]. However, they are not applicable to our system as only one carrier is assumed in each interference-limited RAN. In [44], a globally optimal solution is obtained with the prismatic branch and bound method, but it has exponential computational complexity.

Rather than directly computing the global optimum of (35), we first adopt a successive convex optimization method to derive a local optimum, which is later used to find a globally optimal solution to the overall problem (3)-(17). Existing successive convex optimization techniques include single condensation [40] and logarithmic approximation [45]. Here we use a more computationally efficient approach.

We first write the term  $-\log(1 + \Gamma \Upsilon_{m,n_1} \text{SINR}_{m,n_1}(t))$  as the difference of two convex functions:

$$\log(\varphi_m(\mathbf{p}_{n_1}(t))) - \log(\Gamma \Upsilon_{m,n_1} g_{m,n_1}(t) p_{m,n_1}(t) + \varphi_m(\mathbf{p}_{n_1}(t))).$$

Then, with some initial  $\mathbf{p}_{n_1}^{(0)}(t)$ , we derive a local optimum as follows:

- 1) Approximate  $\log(\varphi_m(\mathbf{p}_{n_1}(t)))$  with its tangent line at  $\mathbf{p}_{n_1}^{(\kappa)}(t)$ .
- 2) Solve the resultant optimization problem and assign its minimizer to  $\mathbf{p}_{n_1}^{(\kappa+1)}(t)$ .
- 3) Let  $\kappa = \kappa + 1$  and repeat from 1) until convergence.

By approximating  $\log(\varphi_m(\mathbf{p}_{n_1}(t)))$  with its tangent line at  $\mathbf{p}_{n_1}^{(\kappa)}(t)$  in step 1), the approximated optimization problem is convex, and can be efficiently calculated. The sequence  $\mathbf{p}_{n_1}^{(\kappa)}(t)$  converges to a local minimizer of problem (35) [46], [47]. We denote it as  $\mathbf{p}_{n_1}^{\text{NO}}(t)$ .

Most importantly, it can be shown that  $\mathbf{p}_{n_1}^{\text{NO}}(t)$  is globally optimal with a positive probability. This is formalized in the following lemma, which will be used in Section IV-E to derive a globally optimal solution.

**Lemma 4:** With an uniformly randomly picked  $\mathbf{p}_{n_1}^{(0)}(t)$ , there exists a positive constant  $0 < \rho_{n_1} < 1$  such that  $\Pr\{\mathbf{p}_{n_1}^{\text{NO}}(t) = \mathbf{p}_{n_1}^{\text{opt}}(t)\} \geq \rho_{n_1}, \forall n_1 \in \mathcal{R}_1$ .

*Proof:* See Appendix D. ■

### E. From Near-Optimal Per-timeslot Allocation to Optimal Allocation

With the radio resource allocation strategy introduced in Subsection IV-D, we can only obtain with a positive probability an optimal resource allocation in the interference-limited RAN for each time slot. Despite this sub-optimality in each

---

**Algorithm 1: CNCL Resource Allocation Method**


---

**Output:** HWN control decisions  $\mathbf{r}(t)$ ,  $\gamma(t)$ ,  $\mathbf{a}_{n_o}(t)$ ,  $\mathbf{p}_{n_o}(t)$ ,  $\mathbf{p}_{n_i}(t)$  in each time slot  $t$

- 1 Observe the network state information  $\mathbf{U}(t)$ ,  $\mathbf{Y}(t)$ ,  $\Theta_{n_i}(t)$ ,  $\Theta_{n_o}(t)$ , and  $\mathbf{g}_{n_i}(t)$  in time slot  $t$ ;
  - 2 Decide the admitted inelastic traffic according to (21);
  - 3 Obtain the admitted elastic traffic according to (22);
  - 4 Derive the RB and power allocation decisions,  $\mathbf{a}_{n_o}(t)$  and  $\mathbf{p}_{n_o}(t)$ , for the orthogonal RAN as in Section IV-C;
  - 5 Derive the power allocation decision  $\mathbf{p}_{n_i}(t)$  for the interference-limited RAN as in Section IV-D;
  - 6 Compare to force the near-optimum to the optimum for the interference-limited RAN as in Section IV-E;
- 

time slot, in this section, we present a pick-and-compare method that allows us to obtain a time-averaged utility that is arbitrarily close to the optimal objective in (3).

The pick-and-compare method was first proposed in [48] for a static network, and it was later extended to wireless networks in [49], [50], [51]. However, such traditional pick-and-compare method cannot directly apply to our resource allocation problem due to the uncountable network control space. Instead, assisted by the near optimal solution in Subsection IV-D, our radio resource allocation strategy in the interference-limited RAN is taken as follows:

- 1) Derive the near optimal radio resource allocation strategy,  $\mathbf{p}_{n_i}^{\text{NO}}(t)$ , with the technique described in Subsection IV-D.
- 2) Define  $\tau_{n_i,t} = \max_{\tau_{n_i}} \{\tau_{n_i} : \mathbf{g}_{n_i}(\tau_{n_i}) = \mathbf{g}_{n_i}(t)\}$ . Then  $\mathbf{p}_{n_i}^{\text{NO}}(t)$  is compared with  $\mathbf{p}_{n_i}^*(\tau_{n_i,t})$ . The one that gives a smaller value to  $F_{n_i}(\mathbf{p}_{n_i}(t) | \mathbf{U}_{n_i}(t), \Theta_{n_i}(t))$  for the current queue length,  $\mathbf{U}_{n_i}(t)$  and  $\Theta_{n_i}(t)$ , is chosen as  $\mathbf{p}_{n_i}^*(t)$ .

We demonstrate the performance optimality of this strategy in Subsection IV-F.

### F. CNCL Summary and Performance Bound

1) *Overall CNCL Description:* We first summarize the proposed CNCL resource allocation method in Algorithm 1.

2) *Complexity Analysis of CNCL Method:* We analyze the complexity of our CNCL method. For the inelastic traffic control strategy (21) and elastic traffic control strategy (22), since they are closed-form expressions, they both have constant complexity.

For the joint RB and power allocation problem (23) in orthogonal RANs, given the dual variable  $\lambda_{n_o}$ , the complexity for primal recovery is polynomial. Using the golden-section search method, the iteration step required to ensure  $\epsilon$ -optimality is  $\log_{0.618} \frac{\epsilon}{\max_{m \in \mathcal{M}, i \in \mathcal{I}_{n_o}} \left\{ \frac{\xi_{m,n_o,i} \eta_{m,n_o}}{\ln 2} - \Theta_{n_o} \right\}}$ .

For the power allocation problem (35) in interference-limited RANs, using the successive convex optimization method, a convex optimization problem has to be solved in each iteration, which has a polynomial complexity. After a polynomial number of iterations, a local optimum to (35) can be achieved. To force the local optimum of (35) to global

optimum, a modified pick-and-compare algorithm is used as explained in Section IV-E. The pick-and-compare procedure is run only once in each time slot and has constant time complexity. However, since the previous power allocation result has to be recorded for each channel state, this algorithm requires a storage space of  $O(H^M)$ , where  $H$  is the number of discretized channel gains, and  $M$  is the UE number in the HWN. Noting that the power allocation problem (35) in interference-limited RANs is NP-hard, we see that such storage space requirement is necessary to achieve near optimality. To reduce the required storage space, an improved pick-and-compare algorithm [52] may be used at the cost of some additional performance loss.

3) *Performance Bound of CNCL Method:* We further quantify the performance of our proposed CNCL method in Theorem 1.

**Theorem 1:** Under the proposed CNCL method, if  $\mathbf{g}_{n_i}(t)$  and  $\mathbf{g}_{n_o}(t)$  are i.i.d. in time, we have

$$\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{u=1}^t \left[ \sum_{m \in \mathcal{M}} \sum_{n \in \mathcal{R}} \mathbb{E}\{U_{m,n}^*(u)\} \right] \leq \frac{B+C}{\delta} + V \frac{\varepsilon}{\delta}, \quad (36)$$

$$\lim_{t \rightarrow \infty} \frac{\sum_{m_i^k \in \mathcal{M}_i} \mathbb{E}\{Y_{m_i^k}^*(t)\}}{t} = 0, \quad (37)$$

$$\lim_{t \rightarrow \infty} \frac{\mathbb{E}\{\Theta_{n_o}^*(t)\}}{t} = 0, \forall n_o \in \mathcal{R}_O \quad (38)$$

$$\lim_{t \rightarrow \infty} \frac{\mathbb{E}\{\Theta_{n_i}^*(t)\}}{t} = 0, \forall n_i \in \mathcal{R}_I \quad (39)$$

$$\lim_{t \rightarrow \infty} \sum_{u=1}^t \mathbb{E}\left\{ \sum_{m_E \in \mathcal{M}_E} \sum_{k \in \mathcal{T}_E} \sum_{n \in \mathcal{R}} \gamma_{m_E,k}^{n,*}(u) \Psi(f_k) \right\} \geq \sum_{m_E \in \mathcal{M}_E} \sum_{k \in \mathcal{T}_E} \sum_{n \in \mathcal{R}} \tilde{\gamma}_{m_E,k}^{n,\text{opt}} \Psi(f_k) - \varepsilon - \frac{B+C}{V}, \quad (40)$$

where  $R = \max\{\sum_{k \in \mathcal{T}_E} X_k f_k, \hat{r}_k, \hat{W}_{n_i}\}$ ,  $C = \sum_{n_i \in \mathcal{R}_I} \frac{2M\hat{P}_{n_i}^2 + 2MR\hat{W}_{n_i}}{\min\{\hat{\pi}_{n_i}\}\rho_{n_i}}$ ,  $\hat{\pi}_{n_i}$  is the joint steady state probability of  $\mathbf{g}_{n_i}(t)$ ,  $\tilde{\gamma}_{m_E,k}^{n,\text{opt}}$  is the global optimum of (3)-(17), and  $\delta$  and  $\varepsilon$  are small positive constants.

*Proof:* See Appendix E. ■

Note that the term  $C$  is the induced cost in exchange for reducing the computational complexity, where the constant  $\rho_{n_i}$  plays an important role. Specifically,  $\frac{2M\hat{P}_{n_i}^2 + 2MR\hat{W}_{n_i}}{\min\{\hat{\pi}_{n_i}\}\rho_{n_i}}$  is the cost introduced by interference-limited RAN  $n_i$ . Intuitively, higher computational ability will lead to larger  $\rho_{n_i}$ , and this reflects the tradeoff between the HWN computational ability and the HWN performance. Moreover, (36) gives an upper bound for the time-averaged data queue length in the HWN, which is proportional to the queueing delay based on Little's Theorem. Thus, our proposed CNCL method can achieve a tradeoff between  $1 - O(\frac{1}{V})$  network utility and  $O(V)$  queueing delay. Particularly, by forcing  $V \rightarrow \infty$ , our CNCL method can obtain the optimal objective in (3).

## V. SIMULATION AND COMPARISON

In addition to the analytical bound provided above, in this section, we further present a numerical study on the CNCL method. Simulation is conducted in MATLAB, while the local optimum in Section IV-D is calculated using LINGO and then passed to MATLAB. Comparisons are made with suboptimal alternatives termed SubTraffic, DistanceS, and SubResource. Specifically, in the SubTraffic method, the UE traffic is uniformly randomly distributed among RANs, and in the DistanceS method, the UE only accesses its closest RAN. Moreover, in both SubTraffic and DistanceS methods, the resource allocation within each RAN is the same as that in our proposed CNCL method, i.e., we optimally solve the joint RB and power allocation problem in orthogonal RANs and asymptotically optimally solve the power allocation problem in interference-limited RANs. In the SubResource method, the traffic distribution scheme is the same as that in our proposed CNCL method. However, in the SubResource method, power is equally allocated among RBs in orthogonal RANs [29] and randomly picked for UEs in interference-limited RANs [49].

### A. Simulation Setup

We consider an example HWN constituted of an interference-limited RAN located at (700m, 1000m), and an orthogonal RAN located at (500m, 300m). The symbol rate is 0.2 MBaud/s for the interference-limited RAN. For the orthogonal RAN, it has 40 RBs, and each has symbol rate of 0.18 MBaud/s. The noise power spectral density is  $-174$  dbm/Hz. The mean power constraints for the interference-limited RAN and the orthogonal RAN are  $P_{n_i} = 20$  Watt and  $P_{n_o} = 30$  Watt, while the maximum power constraints for them are  $\hat{P}_{n_i} = 40$  Watt and  $\hat{P}_{n_o} = 60$  Watt [53].

We assume that UEs in the HWN are randomly allocated within a square of side-length 2000m centered at the origin. For inelastic traffic, we assume that its mean throughput requirement is 1 Mbits/s. While elastic traffic consists of three types of packets, and they all follow a Poisson arrival process with mean packet arrival intensity  $\mathbf{Z} = \zeta \times \{2, 3, 1\}$  packets/s, where  $\zeta$  is used to scale the mean packet arrival intensity. The packet sizes are set to  $\{800, 1200, 1500\}$  bytes respectively for these three arrival intensities. The processing gains for UEs in interference-limited RANs are uniformly randomly chosen from [5, 15], and are set as 1 in orthogonal RANs.

We model the channel gain by both large-scale fading and small-scale fading. For large-scale fading, we assume that it is only determined by the UE-RAN distance with the path loss exponent being 4. While for small-scale fading, we model its amplitude as Rayleigh random variables with unit average power, which are further quantized into six equal probability states  $\{0.280, 0.535, 0.734, 0.937, 1.183, 1.649\}$  [54]. We assume the Shannon capacity gap is 0.7. We further set the instantaneous admitted rate constraints for both the interference-limited RAN and the orthogonal RAN as  $\hat{r}_k = 2$  Mbits/s. We treat the utility of serving an elastic packet as its packet size to indicate the elastic throughput. We take  $t = 1000$ , and the numerical results are averaged over 1000 time slots.

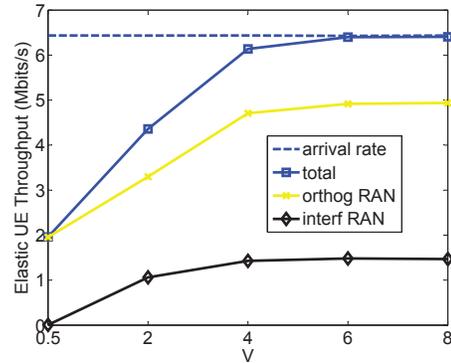


Fig. 2. Average elastic UE throughput versus  $V$ . The legend labels respectively indicate the arrival rate of elastic traffic, the total throughput, throughput from the orthogonal RAN, and throughput from the interference-limited RAN.

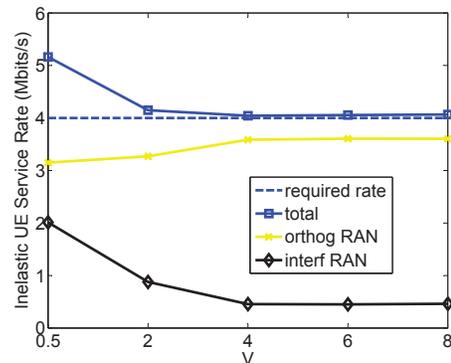
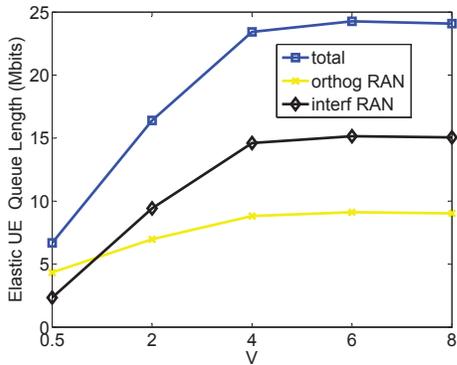
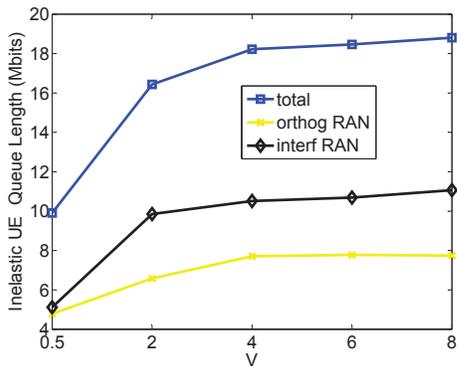


Fig. 3. Average inelastic UE service rate versus  $V$ . The legend labels respectively indicate the inelastic traffic data rate, the total service rate, service rate by the orthogonal RAN, and service rate by the interference-limited RAN.

### B. Numerical Results

We fix the elastic traffic intensity  $\zeta = 30$ , and set the elastic and inelastic UE numbers both as 4 in Fig. 2 to Fig. 5. We first study the average throughput of elastic and inelastic UEs versus  $V$  in Fig. 2 and Fig. 3. In Fig. 2, we can see that increasing  $V$  leads to an increase of elastic UE average throughput when  $V < 6$ . While  $V \geq 6$ , the average throughput of elastic UEs is equal to that of the arrival elastic traffic, indicating that the optimal objective value has been achieved. This suggests that optimal performance, which is guaranteed asymptotically by Theorem 1, can be achieved even for a moderate  $V$  value in reality. Fig. 3 illustrates that the average service rate of inelastic UEs will reduce to slightly above the inelastic UE throughput requirement while  $V$  increases. This is because more HWN resource is allocated to elastic UEs as  $V$  increases.

We also show the average queue length of elastic and inelastic UEs in Fig. 4 and Fig. 5 respectively. We see from Fig. 4 that increasing  $V$  leads to an increase in the average queue length of elastic UEs. This is because increasing  $V$  is equivalent to admitting more elastic data as shown in (22). Increasing  $V$  also increases the average queue length of inelastic UEs. This matches our observation from Fig. 3

Fig. 4. Average elastic queue length versus  $V$ Fig. 5. Average inelastic queue length versus  $V$ 

that the inelastic service rate approaches the required rate as  $V$  increases. We see that the queue lengths level off quickly, as a larger and large portion of the elastic traffic is served.

Varying the elastic traffic arrival intensity  $\zeta$  and fixing the elastic and inelastic UE numbers both as 4, we study the average elastic UE sum throughput in Fig. 6 with  $V = 2$ . We can see that the average elastic UE sum throughput gap between the CNCL method and the other schemes becomes larger as  $\zeta$  increases. This is attributed to the optimal design in traffic admission control and resource allocation within each RAN. Comparing with SubTraffic and SubResource methods, we observe that optimal design in resource allocation within each RAN brings larger throughput improvement than optimal design of admission control. Moreover, we observe that the DistanceS method outperforms the SubResource method, which further indicates the importance of optimal design in resource allocation within each RAN. Furthermore, the throughput improvement of SubTraffic over DistanceS shows that exploring multiple RANs, even sub-optimally, is of critical importance for HWN performance.

Fixing  $\zeta = 30$  and  $V = 2$ , we study the elastic UE average sum throughput versus the number of inelastic UEs in Fig. 7. We fix the total number of UEs to 8. We observe that the increase of inelastic UE number decreases the average throughput of elastic UEs, since more resource is allocated to inelastic UEs. Similar to Fig. 6, optimal design in resource allocation within each RAN brings larger throughput im-

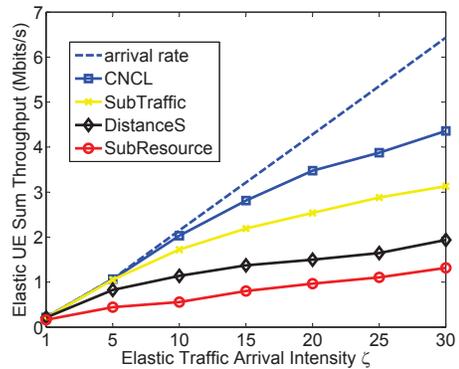
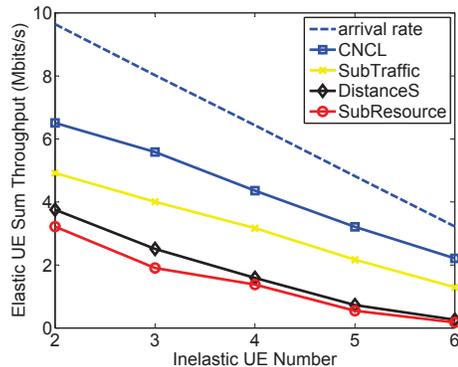
Fig. 6. Comparison of average elastic sum throughput versus  $\zeta$ 

Fig. 7. Comparison of average elastic sum throughput versus inelastic UE number

provement compared with optimal design in traffic admission control. Both Fig. 6 and Fig. 7 show the importance of joint design in traffic admission control and resource allocation in the HWN.

## VI. CONCLUSIONS

Considering both heterogeneous wireless networks and heterogeneous UE traffic, we propose a cross-network cross-layer design method to jointly optimize UE traffic distribution among RANs and RAN radio resource allocation. Our design goal is to asymptotically maximize the throughput of elastic traffic while satisfying inelastic traffic rate constraints. Efficient algorithms are proposed to handle the resultant non-linear NP-hard optimization problem in interference-limited RANs and the mixed-integer non-linear optimization program in orthogonal RANs. In addition, the proposed HWN control framework can adapt to the dynamics at both the packet level and the channel level, and it permits separable RAN control. We derive a performance bound for the proposed cross-network cross-layer design method, and verify its performance through simulation.

Interesting future research directions may include designing more efficient power allocation methods in interference-limited RANs, extending the RAN architecture from a single cell to multiple cells, and imposing more realistic constraints

such as per-RB power limitation in orthogonal RANs and limiting the number of RANs to which a UE can simultaneously connect.

APPENDIX A  
PROOF OF LEMMA 1

We write the KKT conditions of (30) as follows:

$$-\eta_{m,n_0} \log \left( 1 + \xi_{m,n_0,i} \frac{s_{m,n_0,i}}{a_{m,n_0,i}} \right) + \frac{\eta_{m,n_0} \xi_{m,n_0,i} s_{m,n_0,i}}{a_{m,n_0,i} + \xi_{m,n_0,i} s_{m,n_0,i}} + u_{m,i} - v_{m,i} + \varsigma_i = 0, \forall m, i, \quad (41)$$

$$\Theta_{n_0} - \frac{\eta_{m,n_0} a_{m,n_0,i} \xi_{m,n_0,i}}{(a_{m,n_0,i} + \xi_{m,n_0,i} s_{m,n_0,i}) \ln 2} + \lambda_{n_0} - \theta_{m,i} = 0, \forall m, i, \quad (42)$$

$$u_{m,i}(a_{m,n_0,i} - 1) = 0, \forall m, i, \quad (43)$$

$$v_{m,i} a_{m,n_0,i} = 0, \forall m, i, \quad (44)$$

$$\theta_{m,i} s_{m,n_0,i} = 0, \forall m, i, \quad (45)$$

$$\varsigma_i \left( \sum_{m \in \mathcal{M}} a_{m,n_0,i} - 1 \right) = 0, \forall i, \quad (46)$$

$$u_{m,i} \geq 0, v_{m,i} \geq 0, \theta_{m,i} \geq 0, \varsigma_i \geq 0, \forall m, i. \quad (47)$$

We can derive  $s_{m,n_0,i} = \left[ \frac{\eta_{m,n_0}}{(\Theta_{n_0} + \lambda_{n_0} - \theta_{m,i}) \ln 2} - \frac{1}{\xi_{m,n_0,i}} \right] a_{m,n_0,i}$  from (42). It follows from (45) that if  $s_{m,n_0,i} > 0$ ,  $\theta_{m,i} = 0$ . Then we have  $s_{m,n_0,i} = \left[ \frac{\eta_{m,n_0}}{(\Theta_{n_0} + \lambda_{n_0}) \ln 2} - \frac{1}{\xi_{m,n_0,i}} \right] a_{m,n_0,i}$ . Furthermore, based on (47), if  $s_{m,n_0,i} = 0$ , we have  $\theta_{m,n_0,i} \geq 0$ , thus  $\frac{\eta_{m,n_0}}{(\Theta_{n_0} + \lambda_{n_0}) \ln 2} < \frac{1}{\xi_{m,n_0,i}}$ . Hence, we can express the optimal relation between  $s_{m,n_0,i}$  and  $a_{m,n_0,i}$  with respect to (30) as

$$s_{m,n_0,i} = \left[ \frac{\eta_{m,n_0}}{(\Theta_{n_0} + \lambda_{n_0}) \ln 2} - \frac{1}{\xi_{m,n_0,i}} \right]^+ a_{m,n_0,i}, \quad (48)$$

where  $[x]^+ \triangleq \max\{x, 0\}$ .

APPENDIX B  
PROOF OF LEMMA 2

The optimal  $\mathbf{s}_{n_0}^*$  can be divided into two cases, which are  $\mathbf{s}_{n_0}^* = \mathbf{0}_{M \times J_{n_0}}$  and  $\mathbf{s}_{n_0}^* \neq \mathbf{0}_{M \times J_{n_0}}$ .

In the former case, we have  $\sum_{m \in \mathcal{M}} \sum_{i \in \mathcal{J}_{n_0}} s_{m,n_0,i}^* - \hat{P}_{n_0} = \hat{P}_{n_0}$ . Based on the complementary slackness condition, we have  $\lambda_{n_0}^* = 0$ .

While in the latter case, from (32), we can see that if  $\lambda_{n_0}^* \geq \max_{m \in \mathcal{M}, i \in \mathcal{J}_{n_0}} \left\{ \frac{\xi_{m,n_0,i} \eta_{m,n_0}}{\ln 2} - \Theta_{n_0} \right\}$ , all  $s_{m,n_0,i}^*$ 's are equal to zero, which contradicts  $\mathbf{s}_{n_0}^* \neq \mathbf{0}_{M \times J_{n_0}}$ .

APPENDIX C  
PROOF OF LEMMA 3

If there is no tie in the allocation of RB  $i$ , based on the above algorithm, we just allocate RB  $i$  to the UE  $m$  with  $\Lambda_{m,n_0,i}(\lambda_{n_0}) = \Lambda_{n_0,i}^*(\lambda_{n_0})$ . Otherwise, if ties occur in the allocation of RB  $i$ , our tie-breaking rule guarantees that the primal feasibility and complementary slackness condition hold. This gives an optimal solution to (24)-(28). Since (24)-(28) is

a relaxed version of (23), such optimum also provides a lower bound to (23). In addition, note that the minimizer to (24)-(28) also satisfies constraints of (23), and thus, our algorithm provides an optimal solution to (23).

APPENDIX D  
PROOF OF LEMMA 4

Lemma 4 is similar to a special case of Theorem 1 in [4], except that the objective function  $F_{n_1}(\mathbf{p}_{n_1}(t) | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t))$  now contains an additional linear term  $\sum_{m \in \mathcal{M}} \Theta_{n_1}(t) p_{m,n_1}(t)$ . Since such a linear term does not change the general curvature of the objective function, the same conclusion holds. We omit further proof details and refer interested readers to [4].

APPENDIX E  
PROOF OF THEOREM 1

Recalling the queue updating functions for the inelastic UE  $m_i^k$  in (1), for the elastic UE  $m_E$  in (2), and for the average power constraint of interference-limited RAN  $n_1$  in (19), we have  $U_{m_i^k, n_1}(t) \leq U_{m_i^k, n_1}(t-1) + \hat{r}_k$ , and  $U_{m_E, n_1}(t) \leq U_{m_E, n_1}(t-1) + \sum_{k \in \mathcal{T}_E} X_k f_k$ , and  $\Theta_{n_1}(t) \geq \Theta_{n_1}(t-1) - P_{n_1}$ .

Let us denote  $\mathbf{p}_{n_1}^{\text{opt}}(t)$  as the global minimizer to (35), and  $\tau_{\kappa, t} < \tau_{\kappa-1, t} < \dots < \tau_{1, t}$ , where  $\mathbf{g}_{n_1}(\tau_{i, t}) = \mathbf{g}_{n_1}(t), \forall i \leq \kappa, i \in \mathbb{N}^+$ . From Lemma 4, there must exist a  $\kappa \in \mathbb{N}^+$  such that  $\mathbf{p}_{n_1}^*(\tau_{\kappa, t}) = \mathbf{p}_{n_1}^{\text{opt}}(\tau_{\kappa, t})$ . Then we have

$$\begin{aligned} & \sum_{m \in \mathcal{M}} \left\{ \Theta_{n_1}(t) p_{m,n_1}^{\text{opt}}(t) - U_{m,n_1}(t) \mu_{m,n_1}(\mathbf{p}_{n_1}^{\text{opt}}(t)) \right\} \\ & \geq \sum_{m \in \mathcal{M}} \left\{ (\Theta_{n_1}(\tau_{\kappa, t}) - (t - \tau_{\kappa, t}) P_{n_1}) p_{m,n_1}^{\text{opt}}(t) \right\} - \\ & \sum_{m_E \in \mathcal{M}_E} \left\{ (U_{m_E, n_1}(\tau_{\kappa, t}) + (t - \tau_{\kappa, t}) \sum_{k \in \mathcal{T}_E} X_k f_k) \mu_{m_E, n_1}(\mathbf{p}_{n_1}^{\text{opt}}(t)) \right\} \\ & - \sum_{m_i^k \in \mathcal{M}_I} \left\{ (U_{m_i^k, n_1}(\tau_{\kappa, t}) + (t - \tau_{\kappa, t}) \hat{r}_k) \mu_{m_i^k, n_1}(\mathbf{p}_{n_1}^{\text{opt}}(t)) \right\} \\ & \geq \sum_{m \in \mathcal{M}} \left\{ \Theta_{n_1}(\tau_{\kappa, t}) p_{m,n_1}^{\text{opt}}(\tau_{\kappa, t}) - U_{m,n_1}(\tau_{\kappa, t}) \mu_{m,n_1}(\mathbf{p}_{n_1}^{\text{opt}}(\tau_{\kappa, t})) \right\} \\ & - (t - \tau_{\kappa, t}) \left[ \sum_{m_E \in \mathcal{M}_E} \sum_{k \in \mathcal{T}_E} X_k f_k \hat{W}_{n_1} + \sum_{m_i^k \in \mathcal{M}_I} \hat{r}_k \hat{W}_{n_1} + \sum_{m \in \mathcal{M}} P_{n_1} \hat{P}_{n_1} \right]. \quad (49) \end{aligned}$$

Moreover, from the queue updating functions (1), (2), and (19), we have  $\Theta_{n_1}(t) \leq \Theta_{n_1}(t-1) + \hat{P}_{n_1}$  and  $U_{m,n_1}(t) \geq U_{m,n_1}(t-1) - \hat{W}_{n_1}$ . It then follows that

$$\begin{aligned} & \sum_{m \in \mathcal{M}} \left\{ \Theta_{n_1}(t) p_{m,n_1}^*(t) - U_{m,n_1}(t) \mu_{m,n_1}(\mathbf{p}_{n_1}^*(t)) \right\} \\ & \stackrel{1)}{\leq} \sum_{m \in \mathcal{M}} \left\{ \Theta_{n_1}(t) p_{m,n_1}^*(\tau_{1, t}) - U_{m,n_1}(t) \mu_{m,n_1}(\mathbf{p}_{n_1}^*(\tau_{1, t})) \right\} \\ & \leq \sum_{m \in \mathcal{M}} \left\{ (\Theta_{n_1}(\tau_{\kappa, t}) + (t - \tau_{\kappa, t}) \hat{P}_{n_1}) p_{m,n_1}^*(\tau_{\kappa, t}) - \right. \\ & \left. (U_{m,n_1}(\tau_{\kappa, t}) - (t - \tau_{\kappa, t}) \hat{W}_{n_1}) \mu_{m,n_1}(\mathbf{p}_{n_1}^*(\tau_{\kappa, t})) \right\} \end{aligned}$$

$$\leq \sum_{m \in \mathcal{M}} \left\{ \Theta_{n_1}(\tau_{\kappa,t}) p_{m,n_1}^{\text{opt}}(\tau_{\kappa,t}) + (t - \tau_{\kappa,t}) \hat{P}_{n_1}^2 - U_{m,n_1}(\tau_{\kappa,t}) \mu_{m,n_1}(\mathbf{p}_{n_1}^{\text{opt}}(\tau_{\kappa,t})) + (t - \tau_{\kappa,t}) \hat{W}_{n_1}^2 \right\}, \quad (50)$$

where the inequality 1) is from the compare algorithm in Section IV-E.

Based on (49) and (50), we have

$$\begin{aligned} & \mathbb{E} \left\{ \sum_{m \in \mathcal{M}} \left\{ \Theta_{n_1}(t) p_{m,n_1}^{\text{opt}}(t) - U_{m,n_1}(t) \mu_{m,n_1}(\mathbf{p}_{n_1}^{\text{opt}}(t)) \right\} | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t) \right\} \\ & \geq \mathbb{E} \left\{ \sum_{m \in \mathcal{M}} \left\{ \Theta_{n_1}(t) p_{m,n_1}^*(t) - U_{m,n_1}(t) \mu_{m,n_1}(\mathbf{p}_{n_1}^*(t)) \right\} | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t) \right\} \\ & - \mathbb{E} \left\{ (t - \tau_{\kappa,t}) | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t) \right\} \left( \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{T}_E} X_k f_k \hat{W}_{n_1} + \sum_{m_1^k \in \mathcal{M}_1} \hat{r}_k \hat{W}_{n_1} + \sum_{m \in \mathcal{M}} P_m \hat{P}_{n_1} + \sum_{m \in \mathcal{M}} \hat{P}_{n_1}^2 + \sum_{m \in \mathcal{M}} \hat{W}_{n_1}^2 \right) \end{aligned} \quad (51)$$

We denote the probability of channel gain  $g_{m,n_1}(t)$  being  $G_m^j$  as  $\pi_{m,n_1}^j$ , and write the joint steady state probability of  $\mathbf{g}_{n_1}(t)$  as  $\hat{\pi}_{n_1}$ . Since  $\mathbb{E} \left\{ (t - \tau_{\kappa,t}) | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t) \right\} \leq \frac{1}{\min\{\hat{\pi}_{n_1}\} \rho_{n_1}}$  and  $R = \max\{\sum_{k \in \mathcal{T}_E} X_k f_k, \hat{r}_k, \hat{W}_{n_1}\}$ , we then have

$$\begin{aligned} & \mathbb{E}[F_{n_1}(\mathbf{p}_{n_1}^*(t) | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t))] - \mathbb{E}[F_{n_1}(\mathbf{p}_{n_1}^{\text{opt}}(t) | \mathbf{U}_{n_1}(t), \Theta_{n_1}(t))] \\ & \leq \frac{1}{\min\{\hat{\pi}_{n_1}\} \rho_{n_1}} \left( 2M \hat{P}_{n_1}^2 + 2MR \hat{W}_{n_1} \right). \end{aligned} \quad (52)$$

For the orthogonal RAN  $n_0$ , we can optimally minimize  $\mathbb{E}[F_{n_0}(\mathbf{a}_{n_0}(t), \mathbf{p}_{n_0}(t) | \mathbf{U}_{n_0}(t), \Theta_{n_0}(t))]$ . Then based on the Theorem 4.8 in [2], we can derive the above theorem. Details are omitted to avoid redundancy.

## REFERENCES

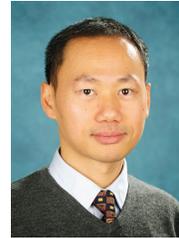
- [1] Z.-Q. Luo and S. Zhang, "Dynamic spectrum management: Complexity and duality," *IEEE Journal of Selected Topics in Signal Processing*, vol. 2, no. 1, pp. 57–73, Feb. 2008.
- [2] M. Neely, *Stochastic network optimization with application to communication and queueing systems*. Morgan & Claypool Publishers, 2010.
- [3] A. Eryilmaz, R. Srikant, and J. Perkins, "Stable scheduling policies for fading wireless channels," *IEEE/ACM Transactions on Networking*, vol. 13, no. 2, pp. 411–424, 2005.
- [4] H. Ju, B. Liang, J. Li, and X. Yang, "Dynamic power allocation for throughput utility maximization in interference-limited networks," *Wireless Communications Letters, IEEE*, no. 99, pp. 1–4, 2012.
- [5] J. Sachs and P. Magnusson, "Assessment of the access selection gain in multi-radio access networks," *European Transactions on Emerging Telecommunications Technologies*, vol. 20, no. 3, pp. 265–279, 2009.
- [6] G. Koudouridis, H. Karimi, and K. Dimou, "Switched multi-radio transmission diversity in future access networks," in *Proc. IEEE VTC*, vol. 1, Sept. 2005, pp. 235–239.
- [7] F. Berggren and R. Litjens, "Performance analysis of access selection and transmit diversity in multi-access networks," in *Proc. ACM MobiCom*, 2006, pp. 251–261.
- [8] K. Piamrat, A. Ksentini, J. Bonnin, and C. Viho, "Radio resource management in emerging heterogeneous wireless networks," *Elsevier Computer Communications*, vol. 34, no. 9, pp. 1066–1076, 2011.
- [9] D. Cavalcanti, D. Agrawal, C. Cordeiro, B. Xie, and A. Kumar, "Issues in integrating cellular networks WLANs, and MANETs: a futuristic heterogeneous wireless network," *IEEE Wireless Communications*, vol. 12, no. 3, pp. 30–41, 2005.
- [10] R. Veronesi, "Multiuser scheduling with multi radio access selection," in *2nd International Symposium on Wireless Communication Systems*. IEEE, 2005, pp. 455–459.
- [11] A. Furskar and J. Zander, "Multiservice allocation for multiaccess wireless systems," *IEEE Transactions on Wireless Communications*, vol. 4, no. 1, pp. 174–184, 2005.
- [12] I. Koo, A. Furskar, J. Zander, and K. Kim, "Erlang capacity of multiaccess systems with service-based access selection," *IEEE Communications Letters*, vol. 8, no. 11, pp. 662–664, 2004.
- [13] M. Neely, E. Modiano, and C. Rohrs, "Dynamic power allocation and routing for time-varying wireless networks," *IEEE Journal on Selected Areas in Communications*, vol. 23, no. 1, pp. 89–103, 2005.
- [14] D. Niyato and E. Hossain, "A cooperative game framework for bandwidth allocation in 4G heterogeneous wireless networks," in *Proc. ICC*, vol. 9, 2006, pp. 4357–4362.
- [15] K. Zhu, D. Niyato, and P. Wang, "Optimal bandwidth allocation with dynamic service selection in heterogeneous wireless networks," in *Proc. GLOBECOM*, 2010, pp. 1–5.
- [16] D. Niyato and E. Hossain, "A noncooperative game-theoretic framework for radio resource management in 4G heterogeneous wireless access networks," *IEEE Transactions on Mobile Computing*, vol. 7, no. 3, pp. 332–345, 2008.
- [17] R. Li, A. Eryilmaz, L. Ying, and N. Shroff, "A unified approach to optimizing performance in networks serving heterogeneous flows," *IEEE/ACM Transactions on Networking*, vol. 19, no. 1, pp. 223–236, 2011.
- [18] J. Buhler and G. Wunder, "Traffic-aware optimization of heterogeneous access management," *IEEE Transactions on Communications*, vol. 58, no. 6, pp. 1737–1747, 2010.
- [19] W. Song, H. Jiang, and W. Zhuang, "Performance analysis of the wlan-first scheme in cellular/wlan interworking," *IEEE Transactions on Wireless Communications*, vol. 6, no. 5, pp. 1932–1952, 2007.
- [20] W. Song, Y. Cheng, and W. Zhuang, "Improving voice and data services in cellular/wlan integrated networks by admission control," *IEEE Transactions on Wireless Communications*, vol. 6, no. 11, pp. 4025–4037, 2007.
- [21] W. Song and W. Zhuang, "Multi-service load sharing for resource management in the cellular/wlan integrated network," *IEEE Transactions on Wireless Communications*, vol. 8, no. 2, pp. 725–735, 2009.
- [22] P. Xue, P. Gong, J. Park, D. Park, and D. Kim, "Radio resource management with proportional rate constraint in the heterogeneous networks," *IEEE Transactions on Wireless Communications*, vol. 11, no. 3, pp. 1066–1075, 2012.
- [23] D. K. Kim, D. Griffith, and N. Golmie, "A new call admission control scheme for heterogeneous wireless networks," *IEEE Transactions on Wireless Communications*, vol. 9, no. 10, pp. 3000–3005, 2010.
- [24] Y. Choi, H. Kim, S. Han, and Y. Han, "Joint resource allocation for parallel multi-radio access in heterogeneous wireless networks," *IEEE Transactions on Wireless Communications*, vol. 9, no. 11, pp. 3324–3329, 2010.
- [25] J. Pérez-Romero, O. Salient, and R. Agusti, "On the optimum traffic allocation in heterogeneous CDMA/TDMA networks," *IEEE Transactions on Wireless Communications*, vol. 6, no. 9, pp. 3170–3174, 2007.
- [26] X. Pei, T. Jiang, D. Qu, G. Zhu, and J. Liu, "Radio-resource management and access-control mechanism based on a novel economic model in heterogeneous wireless networks," *IEEE Transactions on Vehicular Technology*, vol. 59, no. 6, pp. 3047–3056, 2010.
- [27] I. Blau, G. Wunder, I. Karla, and R. Sigle, "Decentralized utility maximization in heterogeneous multicell scenarios with interference limited and orthogonal air interfaces," *EURASIP Journal on Wireless Communications and Networking*, vol. 2009, pp. 1–12, 2009.
- [28] M. Ismail, A. Abdrabou, and W. Zhuang, "Cooperative decentralized resource allocation in heterogeneous wireless access medium," *IEEE Transactions on Wireless Communications*, vol. 12, no. 2, pp. 714–724, 2013.
- [29] M. Ismail and W. Zhuang, "A distributed multi-service resource allocation algorithm in heterogeneous wireless access medium," *IEEE Journal on Selected Areas in Communications*, vol. 30, no. 2, pp. 425–432, 2012.
- [30] M. Ismail, W. Zhuang, and M. Yu, "Radio resource allocation for single-network and multi-homing services in heterogeneous wireless access medium," in *Proc. VTC Fall*, 2012, pp. 1–5.
- [31] K. Kumaran and H. Viswanathan, "Joint power and bandwidth allocation in downlink transmission," *IEEE Transactions on Wireless Communications*, vol. 4, no. 3, pp. 1008–1016, 2005.
- [32] C. Mohanram and S. Bhashyam, "Joint subcarrier and power allocation in channel-aware queue-aware scheduling for multiuser OFDM," *IEEE*

*Transactions on Wireless Communications*, vol. 6, no. 9, pp. 3208–3213, 2007.

- [33] K. Kim, Y. Han, and S. Kim, “Joint subcarrier and power allocation in uplink OFDMA systems,” *IEEE Communications Letters*, vol. 9, no. 6, pp. 526–528, 2005.
- [34] J. Huang, V. Subramanian, R. Agrawal, and R. Berry, “Downlink scheduling and resource allocation for OFDM systems,” *IEEE Transactions on Wireless Communications*, vol. 8, no. 1, pp. 288–296, 2009.
- [35] —, “Joint scheduling and resource allocation in uplink OFDM systems for broadband wireless access networks,” *IEEE Journal on Selected Areas in Communications*, vol. 27, no. 2, pp. 226–234, 2009.
- [36] M. Hajiaghayi, M. Dong, and B. Liang, “Jointly optimal channel and power assignment for dual-hop multi-channel multi-user relaying,” *IEEE Journal on Selected Areas in Communications*, vol. 30, no. 9, pp. 1806–1814, 2012.
- [37] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge University Press, March 2004.
- [38] S. S. Rao and S. Rao, *Engineering optimization: theory and practice*. John Wiley & Sons, 2009.
- [39] V. G. Subramanian, R. A. Berry, and R. Agrawal, “Joint scheduling and resource allocation in CDMA systems,” *IEEE Transactions on Information Theory*, vol. 56, no. 5, pp. 2416–2432, 2010.
- [40] M. Chiang, C. Tan, D. Palomar, D. O’Neill, and D. Julian, “Power control by geometric programming,” *IEEE Transactions on Wireless Communications*, vol. 6, no. 7, pp. 2640–2651, 2007.
- [41] W. Yu, G. Ginis, and J. Cioffi, “Distributed multiuser power control for digital subscriber lines,” *IEEE Journal on Selected Areas in Communications*, vol. 20, no. 5, pp. 1105–1115, 2002.
- [42] R. Cendrillon, W. Yu, M. Moonen, J. Verlinden, and T. Bostoen, “Optimal multiuser spectrum balancing for digital subscriber lines,” *IEEE Transactions on Communications*, vol. 54, no. 5, pp. 922–933, 2006.
- [43] W. Yu and R. Lui, “Dual methods for nonconvex spectrum optimization of multicarrier systems,” *IEEE Transactions on Communications*, vol. 54, no. 7, pp. 1310–1322, 2006.
- [44] Y. Xu, T. Le-Ngoc, and S. Panigrahi, “Global concave minimization for optimal spectrum balancing in multi-user DSL networks,” *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 2875–2885, 2008.
- [45] J. Papandriopoulos and J. Evans, “SCALE: A low-complexity distributed protocol for spectrum balancing in multiuser DSL networks,” *IEEE Transactions on Information Theory*, vol. 55, no. 8, pp. 3711–3724, 2009.
- [46] S. Boyd, “Convex optimization II,” Lecture Notes, 2008. [Online]. Available: <http://www.stanford.edu/class/ee364b/lectures.html>
- [47] B. Marks and G. Wright, “A general inner approximation algorithm for nonconvex mathematical programs,” *Operations Research*, pp. 681–683, 1978.
- [48] L. Tassiulas, “Linear complexity algorithms for maximum throughput in radio networks and input queued switches,” in *Proc. IEEE INFOCOM*, vol. 2, 1998, pp. 533–539.
- [49] A. Eryilmaz, R. Srikant, and J. Perkins, “Stable scheduling policies for fading wireless channels,” *IEEE/ACM Transactions on Networking*, vol. 13, no. 2, pp. 411–424, 2005.
- [50] P. Chaporkar and S. Sarkar, “Stable scheduling policies for maximizing throughput in generalized constrained queueing systems,” *IEEE Transactions on Automatic Control*, vol. 53, no. 8, pp. 1913–1931, 2008.
- [51] H.-W. Lee, E. Modiano, and L. B. Le, “Distributed throughput maximization in wireless networks via random power allocation,” *IEEE Transactions on Mobile Computing*, vol. 11, no. 4, pp. 577–590, 2012.
- [52] M. Lotfinezhad, B. Liang, and E. Sousa, “On stability region and delay performance of linear-memory randomized scheduling for time-varying networks,” *IEEE/ACM Transactions on Networking*, vol. 17, no. 6, pp. 1860–1873, 2009.
- [53] “3GPP TS 36.300 V10.7.0. evolved universal terrestrial radio access (E-UTRA) and evolved universal terrestrial radio access network (E-UTRAN); overall description; stage 2,” 2012. [Online]. Available: <http://www.3gpp.org/ftp/Specs/html-info/36300.htm>
- [54] A. Goldsmith, *Wireless Communications*. Cambridge University Press, 2005.



**Honghao Ju** is currently a Ph.D. student at Xidian University, Xi’an, China, where he received his BE degree in telecommunication engineering in 2009. During 2012–2013, he was a visiting student at the University of Toronto. His research interests are stochastic network optimization, network architecture analysis, and optimization techniques in the network design.



**Ben Liang** received honors-simultaneous B.Sc. (valedictorian) and M.Sc. degrees in Electrical Engineering from Polytechnic University in Brooklyn, New York, in 1997 and the Ph.D. degree in Electrical Engineering with Computer Science minor from Cornell University in Ithaca, New York, in 2001. In the 2001–2002 academic year, he was a visiting lecturer and post-doctoral research associate at Cornell University. He joined the Department of Electrical and Computer Engineering at the University of Toronto in 2002, where he is now a Professor. His current research interests are in mobile communications and networked systems. He has served as an editor for the *IEEE Transactions on Communications* and the *IEEE Transactions on Wireless Communications*, and an associate editor for the *Wiley Security and Communication Networks* journal, in addition to regularly serving on the organizational and technical committees of a number of conferences. He is a senior member of IEEE and a member of ACM and Tau Beta Pi.



**Jiandong Li** received the BE, MS, and Ph.D. degrees from Xidian University, Xi’an, China, in 1982, 1985, and 1991, respectively, all in electrical engineering. He has been a faculty member of Telecommunications Engineering at Xidian University since 1985, where he is currently a professor and director of State Key Laboratory of Integrated Service Networks. Prof. Li is a senior member of IEEE. He was a visiting professor to the Department of Electrical and Computer Engineering at Cornell University from 2002–2003. He was a member of Personal Communication Networks (PCN) specialist group for China 863 Communication High Technology Program during 1993–1994 and again 1999–2000. He was awarded as Distinguished Young Researcher and Changjiang Scholar from Ministry of Science and Technology, China. His major research interests are wireless communication theory, cognitive radio, and signal processing.



**Yan Long** received the BSc degree in Electrical and Information Engineering from Xidian University, China, in 2009. She has been working towards the Ph.D. degree in the Department of Telecommunications Engineering at Xidian University, since 2010. From September 2011 to March 2013, she was a visiting student in the Department of Electrical and Computer Engineering, University of Florida. Her research interests include cognitive radio networks, heterogeneous networks, multi-radio multi-channel networks, wireless resource allocation, and cross-

layer optimization.



**Xiaoni Yang** is a research scientist at NO.36 research institute of CETC, Jiaxing, China. He got his BE and MS degrees in electrical engineering from Xidian University, Xi'an, China, in 1982 and 1988, where he currently holds an adjunct professor position. His research interests are software-defined radio and signal processing.