

# Wireless Scheduling for Heterogeneous Services With Mixed Numerology in 5G Wireless Networks

Ti Ti Nguyen<sup>1</sup>, *Student Member, IEEE*, Vu Nguyen Ha<sup>2</sup>, *Member, IEEE*,  
and Long Bao Le<sup>1</sup>, *Senior Member, IEEE*

**Abstract**—This letter studies the scheduling problem which determines how time-frequency resources of different numerologies can be allocated to support heterogeneous services in 5G wireless systems. Particularly, this problem aims at scheduling as many users as possible while meeting their required service delay and transmission data. To solve the underlying integer programming (IP) scheduling problem, we first transform it into an equivalent integer linear program (ILP) and then develop two algorithms, namely Resource Partitioning-based Algorithm (RPA) and Iterative Greedy Algorithm (IGA) to acquire efficient resource scheduling solutions. Numerical results show the desirable performance of the proposed algorithms with respect to the optimal solution and their complexity-performance tradeoffs.

**Index Terms**—5G NR, wireless scheduling, mixed numerology.

## I. INTRODUCTION

THE 5G wireless network is designed to support diverse applications and use cases with different requirements including the enhanced mobile broadband (eMBB), massive machine-type communication (mMTC), and ultra-reliable and low-latency communication (uRLLC) [1]. Toward this end, the 5G wireless standard adopts the so-called mixed numerology to enable flexible configurations and assignment for different types of physical resource blocks (PRBs) and support different wireless services [2]–[4]. Particularly, while the 4G system supports only one numerology where its PRB has the bandwidth of 180 kHz and time duration of 0.5 ms, the 5G’s PRBs can have the bandwidth equal or 2, 4, 8, 16 times of 180 kHz and the time duration equal or 2, 4, 8, 16 times smaller than 0.5 ms. The scheduling problem for resource allocation (RA) of two-dimensional (2D) time-frequency resources has been studied in several recent works [5], [6]. However, the 5G-NR mixed numerology has not been fully studied in [5] while [6] assumes that only two different numerologies exist in the system. Finally, the authors of [7] study the RA problem with mixed numerology for capacity enhancement.

In this paper, we study the scheduling problem for heterogeneous services with mixed numerology which aims to maximize the number of admitted users while meeting service latency and data transmission requirements. Two algorithms, namely Resource Partitioning-based Algorithm (RPA) and

Manuscript received October 4, 2019; accepted October 29, 2019. Date of publication November 8, 2019; date of current version February 11, 2020. Research is supported by Vingroup Innovation Foundation (VINIF) in project code VINIF.2019.DA01. The associate editor coordinating the review of this letter and approving it for publication was A. G. Saaavedra. (*Corresponding author: Long Le.*)

T. T. Nguyen and L. B. Le are with the INRS-EMT, Université du Québec, Montréal, QC H5A1K6, Canada (e-mail: titi.nguyen@emt.inrs.ca; long.le@emt.inrs.ca).

V. N. Ha is with the École Polytechnique de Montréal, Montréal, QC H3T1J4, Canada (e-mail: vu.ha-nguyen@polymtl.ca).

Digital Object Identifier 10.1109/LCOMM.2019.2951375

1558-2558 © 2019 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.

See <https://www.ieee.org/publications/rights/index.html> for more information.

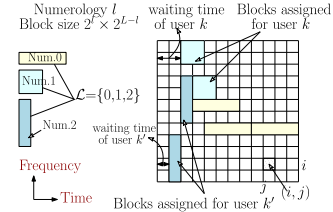


Fig. 1. PRB allocation in 5G wireless networks with mixed numerology.

Iterative Greedy Algorithm (IGA), are proposed to tackle the underlying problem. The RPA algorithm is developed by partitioning the resources and users into smaller groups, optimizing the RA for each group, then performing resource defragmentation and additional resource assignment for unallocated resources to obtain a final solution. In the IGA algorithm, we iteratively allocate PRBs to users based on an appropriate resource assignment weight to obtain an efficient scheduling solution with low computation complexity. Numerical studies are performed to demonstrate the efficacy of the proposed algorithms.

## II. SYSTEM MODEL

We consider the 5G system where the available time-frequency resource is divided into resource elements (RE). Each RE occupies the bandwidth of  $\Delta_{\min}^f$  (Hz) and the slot duration of  $\Delta_{\min}^t$  (seconds). The link/channel conditions for each subcarrier are assumed unchanged during the scheduling time. Moreover, we assume that the 2D RA is performed over each scheduling interval of  $T = N^t \Delta_{\min}^t$  (seconds) and the bandwidth of  $B^f = M^f \Delta_{\min}^f$  (Hz). Considering the scheduling problem for users where the serving base station supports multiple numerologies as shown in Fig 1. Particularly, the bandwidth of a PRB in numerology  $l$  is equal to one half of that in numerology  $l + 1$  while the time slot duration of a PRB in numerology  $l$  is twice of that in numerology  $l + 1$ . The bandwidth of a PRB in numerology  $l$  is defined as  $\Delta_l^f$  and the slot duration of a PRB in numerology  $l$  is defined as  $\Delta_l^t$ . Then, we have  $\Delta_l^t = \Delta_{l-1}^t/2$ ,  $\Delta_l^f = 2\Delta_{l-1}^f$ ,  $\Delta_{\min}^t = \min\{\Delta_l^t, \forall l\}$ , and  $\Delta_{\min}^f = \min\{\Delta_l^f, \forall l\}$ . For convenience, the numerology used by user  $k$  is denoted as  $l_k$ , the set of all users is denoted as  $\mathcal{K}$ , the set of numerologies is denoted as  $\mathcal{L}$ ,  $l_{\max} = \max\{l \in \mathcal{L}\}$ ,  $l_{\min} = \min\{l \in \mathcal{L}\}$ , and  $L = l_{\max} - l_{\min}$ , and the cardinals of set  $\mathcal{L}$  is denoted as  $|\mathcal{L}|$ .

Each user  $k$  requires a data chunk of  $d_k^{rq}$  bits be completely transmitted and the total waiting time for its data transmission must not be larger than  $\tau_k^{\max}$ . Note that  $d_k^{rq}$  can be the whole data carried by a data flow (e.g., sensing data) or a part of the data of the underlying data flow (e.g., streaming data) in  $T$  seconds.<sup>1</sup> We use  $(i, j)$  to refer to a particular RE where

<sup>1</sup>We assume that the  $d_k^{rq}$  is known or can be estimated from the QoS requirement and data properties.

its location is given as  $f \in [(i-1)\Delta_{\min}^f : i\Delta_{\min}^f]$  and  $t \in [(j-1)\Delta_{\min}^t : j\Delta_{\min}^t]$ , for  $1 \leq i \leq M^f$  and  $1 \leq j \leq N^t$ .

### A. Problem Formulation

PRBs are allocated to users where the numerology is selected in advance by each user. Moreover, each RE is allocated to only one user and associated with one numerology. We represent the mapping for one particular PRB of numerology  $l$  to REs in the 2-D frequency-time resource space as follows:

$$\mathbf{q}_{l,m,n} = \{(i,j) | m \leq i \leq m + M_l, n \leq j \leq n + N_l\}, \quad (1)$$

where  $M_l = 2^{l-l_{\min}} - 1$ ,  $N_l = 2^{l_{\max}-l} - 1$ . Assuming that the number of PRBs assigned to user  $k$  is not larger than  $C_k$  to maintain certain fairness among users. Then, we introduce binary variables  $x_{i,j}^{k,c}$ 's,  $y_{m,n}^{k,c}$ 's where  $y_{m,n}^{k,c} = 1$  if  $\mathbf{q}_{l,m,n}$  corresponds to the  $c^{\text{th}}$  assigned PRB of user  $k$ , and  $y_{m,n}^{k,c} = 0$  otherwise;  $x_{i,j}^{k,c} = 1$  if RE  $(i,j)$  is assigned for user  $k$  in its  $c^{\text{th}}$  PRB, and  $x_{i,j}^{k,c} = 0$  otherwise. For  $k \in \mathcal{K}$ , the ranges of  $c, i, j, m$ , and  $n$  are  $c = 1 : C_k$ ,  $i = 1 : M^f$ ,  $j = 1 : N^t$ ,  $m = 1 : M^f - M_{l_k}$ , and  $n = 1 : N^t - N_{l_k}$ , respectively. We impose the following constraints to ensure non-overlapping RA:

$$\sum_{i'=m:M_{l_k}+1} \sum_{j'=n:N_{l_k}+1} x_{i',j'}^{k,c} \geq 2^L y_{m,n}^{k,c}, \quad \forall k, c, m, n, \quad (2a)$$

$$\sum_{k \in \mathcal{K}} \sum_c x_{i,j}^{k,c} \leq 1, \quad \forall i, j, \quad \text{and} \quad \sum_m \sum_n y_{m,n}^{k,c} \leq 1, \quad \forall k, c. \quad (2b)$$

Specifically, (2a) implies that the number of REs per PRB remains constant and equal to  $2^L$ . Moreover, (2b) indicates that each RE should belong to only one PRB, and each PRB can be assigned to at most one user. Let  $r_{m,n}^k$  denote the transmission rate of user  $k$  on PRB  $\mathbf{q}_{l,m,n}$ . Then, the total amount of data transmitted by user  $k$  during the scheduling interval is  $d_k = \Delta_{l_k}^t \sum_{m=1}^{M^f - M_{l_k}} \sum_{n=1}^{N^t - N_{l_k}} \sum_{c=1}^{C_k} r_{m,n}^k y_{m,n}^{k,c}$ .

Each user  $k$  wants its data chunk to be completely transmitted and the total waiting time be not larger than  $\tau_k^{\max}$ . Let  $\tau_{k,0}$  denote the initial waiting time (of the data chunk) of user  $k \in \mathcal{K}$  at the beginning of the considered scheduling interval.<sup>2</sup> Then, the total waiting time until the transmission instant of user  $k$  can be written as  $\tau_k = \tau_{k,0} + \tau_{k,1}$ , where  $\tau_{k,1}$  is the additional waiting time before user  $k$  is served in the scheduling interval, which can be expressed as  $\tau_{k,1} = \Delta_{\min}^t \min\{j-1 \mid x_{i,j}^{k,c} = 1, \forall i, j, c\}$ .

Our design aims to schedule as many users as possible while meeting their data demand and latency requirements. Recall that user  $k$  wishes to transmit a data chunk of  $d_k^{\text{rq}}$  bits with the total waiting time not larger than  $\tau_k^{\max}$ . To maintain these constraints, we define a function capturing if both constraints are satisfied as  $u_k = \mathbb{1}_{d_k - d_k^{\text{rq}} \leq \tau_k^{\max} - \tau_k}$ , where  $\mathbb{1}_x$  stands for the step function, i.e.,  $\mathbb{1}_x = 1$  if  $x \geq 0$ , and  $\mathbb{1}_x = 0$ , otherwise. In fact, if a scheduling solution ensures that the amount of transmitted data and the total waiting time satisfy  $d_k \geq d_k^{\text{rq}}$  and  $\tau_k \leq \tau_k^{\max}$ , respectively, we have  $u_k = 1$ ; otherwise,

<sup>2</sup>This initial waiting time is applied to the first chunk of a new data flow when the data flow arrives in the middle of the previous scheduling interval.

$u_k = 0$ . Then, the scheduling problem can be formulated as

$$(\mathcal{P}_1) \max_{\mathbf{x}, \mathbf{y}} \sum_{k \in \mathcal{K}} u_k \quad \text{s.t.} \quad (2a), (2b), \quad \text{and} \quad \mathbf{x}, \mathbf{y} \in \{0, 1\},$$

where  $\mathbf{x} = \{x_{i,j}^{k,c} \mid \forall i, j, k, c\}$  and  $\mathbf{y} = \{y_{m,n}^{k,c} \mid \forall k, c, m, n\}$ .

### B. Problem Transformation

To deal with non-continuous and non-linear functions in problem  $(\mathcal{P}_1)$ , we first transform this problem into a standard ILP. Specifically, we can equivalently express  $u_k$  as

$$u_k \in \{0, 1\}; \quad u_k(\tau_k^{\max} - \tau_k) \geq 0; \quad u_k(d_k - d_k^{\text{rq}}) \geq 0. \quad (3)$$

Then, we introduce auxiliary variables  $z_{m,n}^{k,c}$ 's as  $z_{m,n}^{k,c} = u_k y_{m,n}^{k,c}$ . Using  $z_{m,n}^{k,c}$ 's,  $(\mathcal{P}_1)$  can be transformed into the ILP form as stated in the following proposition.

*Proposition 1:*  $(\mathcal{P}_1)$  is equivalent to the following problem

$$(\mathcal{P}_1^{\text{ILP}}) \max_{\mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{u} \in \{0, 1\}} \sum_{k \in \mathcal{K}} u_k \quad \text{s.t.} \quad (2a), (2b), \\ \sum_m \sum_{n'=1:N_k^{\text{rq}}} \sum_c z_{m,n'}^{k,c} - u_k \geq 0, \quad \forall k \in \mathcal{K}, \quad (4a)$$

$$\Delta_{l_k}^t \sum_m \sum_n \sum_c r_{m,n}^k z_{m,n}^{k,c} - u_k d_k^{\text{rq}} \geq 0, \quad \forall k \in \mathcal{K}, \quad (4b)$$

$$z_{m,n}^{k,c} \geq u_k + y_{m,n}^{k,c} - 1, \quad z_{m,n}^{k,c} \leq \min\{u_k, y_{m,n}^{k,c}\}, \\ \forall k, c, m, n, \quad (4c)$$

where  $\mathbf{z} = \{z_{m,n}^{k,c}, z_{m,n}^{k,c} \mid \forall k, c, m, n\}$  and  $\mathbf{u} = \{u_k^r, u_k^d \mid \forall k\}$ . Here,  $N_k^{\text{rq}}$  represents the maximum number of time slots (with size of  $\Delta_{\min}^t$  seconds) that user  $k$  can wait, counting from the beginning of the scheduling interval, to meet its delay constraint which is determined as  $N_k^{\text{rq}} = \lfloor (\tau_k^{\max} - \tau_{k_0}) / \Delta_{\min}^t \rfloor$ , where  $\lfloor \cdot \rfloor$  is the floor operation.

*Proof:* (4c) with  $\mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{u} \in \{0, 1\}$  are equivalent to  $z_{m,n}^{k,c} = u_k y_{m,n}^{k,c}$ ,  $\forall k, c, m, n$ . As a result, (4a) and (4b) in  $(\mathcal{P}_1)$  are equivalent to (3). Hence,  $(\mathcal{P}_1)$  is equivalent to  $(\mathcal{P}_1^{\text{ILP}})$ . ■

*Proposition 2:*  $(\mathcal{P}_1^{\text{ILP}})$  is  $\mathcal{NP}$ -hard.

*Proof:* If the number of PRBs assigned to each user is known, the remaining problem can be formulated as a well-known ‘‘Cutting Stock Problem’’ (CSP), which is strongly NP-hard [8]. Thus,  $(\mathcal{P}_1^{\text{ILP}})$  is more complex than the standard CSP since the numbers of PRBs assigned for different users should also be optimized. Therefore,  $(\mathcal{P}_1^{\text{ILP}})$  must be NP-hard. ■

## III. PROPOSED ALGORITHMS

### A. Resource Partitioning Based Algorithm (RPA)

We propose a low-complexity algorithm which solves  $(\mathcal{P}_1^{\text{ILP}})$  by decomposing it into parallel small-scale sub-problems. In fact, if we can maintain the relationship between available resources and users' demands in each sub-problem similar to that in the original problem, then solving small-scale sub-problems can return a solution as good as the one obtained by directly solving the original problem. To realize this idea, we first divide the available bandwidth into  $M_B$  sub-bands where sub-band  $m_B$  occupies the spectrum from  $(m_B - 1)[M^f / M_B] \Delta_{\min}^f$  to  $\min\{m_B[M^f / M_B] \Delta_{\min}^f, M^f \Delta_{\min}^f\}$ . Then, the following three steps are taken in RPA: 1) perform RA on each sub-band,

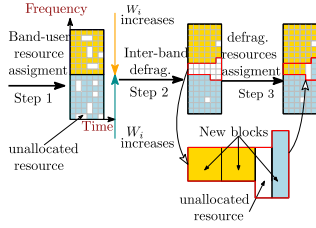


Fig. 2. Illustration of the three steps of RPA algorithm.

2) re-arrange unallocated resources for consecutive sub-bands, and 3) assign these re-arranged resources to support more (unadmitted) users.

### Algorithm 1 Resource Partitioning Based Algorithm (RPA)

- 1: Initialize: Set initial value for  $M_B$ .
- 2: Partition resources into  $M_B$  sub-bands and distribute users into these sub-bands as in **Section III-A**.
- 3: **Step 1**: Solve  $(\mathcal{P}_{1,m_B}^{\text{ILP}})$  to obtain  $u_k^{S1*}$  for all  $m_B, k \in \mathcal{K}_{m_B}$ .
- 4: **Step 2**: Solve  $(\mathcal{P}_{m_B})$  to create a contiguous region of unallocated resources between two consecutive sub-bands while still satisfying the requirements of admitted users, i.e., users with  $u_k^{S1*} = 1, \forall k$ .
- 5: **Step 3**: Solve  $(\mathcal{P}_{\text{RPA}})$  to assign unallocated resources to unadmitted users, i.e., users with  $u_k^{S1*} = 0, \forall k$ .

The key steps of RPA are illustrated in Fig. 2 and it is summarized in **Algorithm 1**. In **Step 1**, we randomly distribute users into sub-bands to make resource demands on different sub-bands similar. Denote the set of users associated with sub-band  $m_B$  as  $\mathcal{K}_{m_B}$ , and the set of REs in the frequency dimension as  $\mathcal{I}_{m_B} = \{m | m = (m_B - 1) \lfloor M^f / M_B \rfloor : \min\{m_B \lfloor M^f / M_B \rfloor, M^f\}\}$ . We then solve  $(\mathcal{P}_1^{\text{ILP}})$  corresponding to each sub-band  $m_B$  and the set of users  $\mathcal{K}_{m_B}$  to obtain admission decisions, denoted as  $u_k^{S1*}$ 's. The sub-problem for sub-band  $m_B$  is named  $(\mathcal{P}_{1,m_B}^{\text{ILP}})$ . In **Step 2**, after finding  $u_k^{S1*}$ 's, we re-arrange the allocated resources so that the unallocated REs from two consecutive sub-bands can be arranged close to one another and they can be combined and mapped into PRBs of certain numerology as defined in (1). The re-arrangement of unallocated REs on sub-band  $m_B$  can be achieved by solving the following problem:

$$\begin{aligned}
 (\mathcal{P}_{m_B}) \quad & \min_{\mathbf{x}, \mathbf{y}, \mathbf{z} \in \{0,1\}} \sum_{k \in \mathcal{K}_{m_B}} \sum_{i \in \mathcal{I}_{m_B}} \sum_{j=1:N^t - N_{l_k}} \sum_{c=1:C_k} W_i x_{i,j}^{k,c} \\
 \text{s.t.} \quad & u_k = u_k^{S1*}, (2a), (2b), (4a), (4b), \\
 & \forall k \in \mathcal{K}_{m_B}, i \in \mathcal{I}_{m_B},
 \end{aligned}$$

where  $\{W_i\}$  is an increasing series, e.g.,  $W_i = 2^i$  if the sub-band index is odd and  $\{W_i\}$  is a decreasing series, e.g.,  $W_i = 2^{-i}$  if the sub-band index is even. It can be verified that after solving  $(\mathcal{P}_{m_B})$ , all unallocated REs in two consecutive sub-bands will be pushed close to one another to create a contiguous resource region as large as possible. Let  $\{\mathbf{x}_{\mathcal{P}_{m_B}}^*, \mathbf{y}_{\mathcal{P}_{m_B}}^*, \mathbf{z}_{\mathcal{P}_{m_B}}^*\}$  denote the optimal solution of  $(\mathcal{P}_{m_B})$ . In **Step 3**, we assign the unallocated resources,  $\Omega = \{(i, j) | \mathbf{x}_{\mathcal{P}_{m_B}}^* = 0, \forall m_B\}$ , to the set of unadmitted users  $\bar{\mathcal{K}} = \{k | u_k^{S1*} = 0\}$  by solving the following problem  $(\mathcal{P}_{\text{RPA}})$ :

$$\begin{aligned}
 \max_{\mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{u} \in \{0,1\}} \quad & \sum_{k \in \bar{\mathcal{K}}} u_k \\
 \text{s.t.} \quad & (2a), (2b), (4a), (4b), \forall k \in \bar{\mathcal{K}}, (i, j) \in \Omega.
 \end{aligned}$$

### Algorithm 2 Iterative Greedy Algorithm (IGA)

- 1: Initialize:  $d_{k,0}^{\text{rq}} = d_k^{\text{rq}}, c_k = 0, \mathcal{W}_{m,n} = 1, W = 10$ .
- 2: **repeat**
- 3: Compute  $\mathcal{U}_{m,n,k}$  find the largest value of  $\mathcal{U}_{m_0,n_0,k_0}$ , and perform the corresponding PRB allocation.
- 4: Update different parameters after the PRB allocation as  $c_{k_0} = c_{k_0} + 1$ , assign  $y_{m_0,n_0}^{k_0,c_{k_0}} = 1$ , update the remaining required data  $d_{k,0}^{\text{rq}} = d_{k,0}^{\text{rq}} - r_{m_0,n_0}^{k_0} \Delta_{l_k}^t$ , and  $\mathcal{W}_{m_0,n} = W, \forall n = 1 : N^t$ .
- 5: Drop all overlapped PRBs  $\mathbf{q}_{l_k,m,n}$  to PRB  $\mathbf{q}_{l_{k_0},m_0,n_0}$
- 6: **until**  $\mathcal{U}_{m,n,k} = 0, \forall m, n, k$

### B. Iterative Greedy Algorithm (IGA)

We propose another fast iterative algorithm where we greedily assign resources to users based on an assignment weight which depends on the requirements of the underlying user, the amount of data transmitted and the achieved latency if the underlying PRB is assigned to the user. In each iteration, we calculate the assignment weight for each pair of an available PRB and a user based on which the resource assignment is performed for the PRB-user pair achieving the largest weight. After that, the available PRBs and the weights of all possible PRB-user pairs are updated to prepare for further resource assignment in the next iteration. This process is repeated until there is no more available PRB or unsatisfied user.

We now define the resource assignment weight for a particular PRB-user-resource pair as follows:  $\mathcal{U}_{m,n}^k = \frac{r_{m,n}^k \Delta_{l_k}^t}{d_{k,0}^{\text{rq}}} \frac{n}{N_k^{\text{rq}}} \mathbb{1}_{n \in \mathcal{N}_k} \mathcal{W}_{m,n}$  if  $c_k \leq C_k, d_{k,0}^{\text{rq}} > 0$ , and  $\mathcal{U}_{m,n}^k = 0$ , otherwise, where  $d_{k,0}^{\text{rq}}$  is the remaining required data amount in each iteration, which is equal to  $d_k^{\text{rq}}$  in the first iteration,  $c_k$  is the current total PRBs assigned to user  $k$ , and  $\mathcal{N}_k$  is the set of REs in the time domain, which is defined as  $\mathcal{N}_k = \{n | n \leq N_k^{\text{rq}}\}$  if  $\sum_{n=1}^{N_k^{\text{rq}}} \sum_{m=1}^{M^t - M_{l_k}} \sum_{c=1}^{C_k} y_{m,n}^{k,c} = 0$ , and  $\mathcal{N}_k = \{n | n \leq N^t\}$ , otherwise, and  $\mathcal{W}_{m,n}$  is a matrix used to mitigate the resource fragmentation in the allocation process, which is updated in each iteration. The unit matrix is initially assigned to  $\mathcal{W}_{m,n}$ . In particular,  $\mathcal{U}_{m,n}^k$  is chosen based on the following criteria: 1) Users with smaller required data chunk receive higher scheduling priorities; 2) If a user is not admitted yet, a PRB at the time location closer to the slot corresponding to the allowed maximum waiting time is more prioritized; 3) Admitted users are allocated resources until their requirements are completely satisfied; and 4) The resource fragmentation is prevented to ease future PRB allocations. The IGA is summarized in **Algorithm 2**. In each iteration of IGA, we need to compute the resource assignment weights  $\mathcal{U}_{m,n,k}$  for all available  $m, n, k$ , determine the largest  $\mathcal{U}_{m_0,n_0,k_0}$  to perform RA, update different parameters, and drop all overlapped PRBs to the assigned block. The worst-case complexity of each iteration is  $\mathcal{O}(M^f N^t K)$ . Let  $N^{\text{iter}}$  be the number of iterations, which is upper bounded by  $M^f N^t |\mathcal{L}|$ , the overall worst-case complexity of IGA is  $\mathcal{O}(N^{\text{iter}} M^f N^t |\mathcal{L}| K)$ .

### IV. NUMERICAL RESULTS

We consider a wireless system with pedestrian and high moving users in a cell with radius of 500 meters. The channel path-loss  $\beta_k$  (dB) =  $128.1 + 37.6 \log_{10}(\gamma_k)$  where  $\gamma_k$  is the



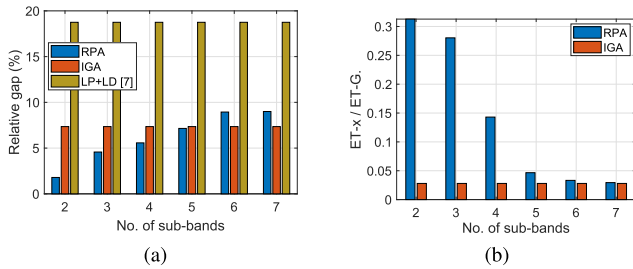


Fig. 3. Comparison of RPA and IGA with the optimum on the relative gap (a) and the execution time (b).

distance between user  $k$  and the BS (in km). For small-scale channel fading, ITU pedestrian-B channel parameters with Doppler shift of 50 Hz and ITU Vehicular-A channel parameters with Doppler shift of 500 Hz are used for pedestrian users and high moving users, respectively. We set three user groups A, B, and C and the numbers of users in these group are  $\lfloor K/3 \rfloor$ ,  $\lfloor K/3 \rfloor$  and  $K - 2\lfloor K/3 \rfloor$ , respectively. Specifically, group A adopts numerology 0 corresponding to high moving users with large data demand, group C uses numerology 2 corresponding to high moving users requiring low waiting time, and group B employs numerology 1 corresponding to pedestrian users with average requirements on the transmission data and waiting time. The transmission rate is calculated according to Shannon's capacity where the ratio of transmit power per Hz to noise power density is set equal to  $2.8 \times 10^5$ . The required data chunks  $d_k^{rq}$  over the interval  $T$  of 1 ms for users in groups A, B, and C are set randomly in  $[500 - 2000]$  (bits),  $[500 - 1000]$  (bits), and  $[180 - 500]$  (bits), respectively, and  $C_k$  is set equal to 10. Besides,  $N_k^{rq}$  defined in Proposition 1 is set equal to 8 for users in group A and randomly in  $[3-6]$  and  $[1-4]$  for users in groups B and C, respectively. All numerical results are obtained by averaging the results over 50 random realizations.

We show the relative gap in Fig. 3-(a) which is calculated as  $(\sum_k u_k^G - \sum_k u_k^{RPA/IGA}) \times 100\% / \sum_k u_k^G$  for  $M^f = 32$  and  $K = 20$  where  $u_k^{RPA/IGA/LP+LD[7]}$  and  $u_k^G$  represent the objective values for user  $k$  obtained by using RPA/IGA/LP+LD [7] and the CVX-Gurobi solver, respectively. Fig. 3-(b) shows the ratio between the average execution time (ET) of the RPA/IGA and that required to solve  $(\mathcal{P}_1^{ILP})$  by the CVX-Gurobi solver (ET-G). The "LP+LD" algorithm proposed in [7] includes two loops: an outer loop to assign PRBs to users based on the utility matrix for all PRB-user pairs, and an inner loop to determine the utility matrix. The figure shows that our proposed algorithms outperform the "LP+LD" algorithm. Besides, as shown in this figure, the relative gap due to RPA increases when the number of sub-bands  $M_B$  increases. This is because larger  $M_B$  reduces RA flexibility in each sub-band and thus resource utilization efficiency. However, the execution time of the RPA can be reduced significantly when  $M_B$  becomes larger. In contrast, IGA always explores good resource-user pairs for efficient resource utilization and IGA is not affected by  $M_B$ .<sup>3</sup>

<sup>3</sup>By implementing the code on Matlab to solve  $(\mathcal{P}_1^{ILP})$  with Gurobi using a computer equipped with CPU chipset Intel core i7-4790 and 12 GB RAM, average execution time is about 6.39 seconds when  $M^f = 32$  and  $K = 20$ .

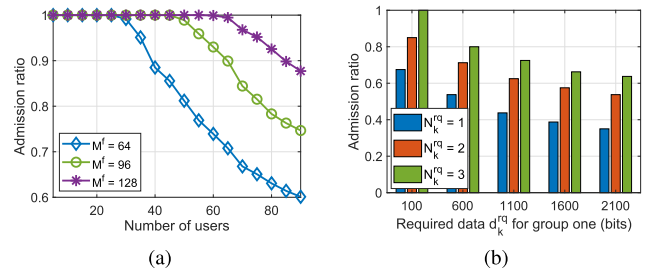


Fig. 4. Admission ratio due to IGA for different number of users (a) and for different required data (b).

Fig. 4-(a) shows the admission ratio achieved by IGA with different values of  $M^f$  where this ratio is equal to the number of admitted users divided by the total number of users. As can be seen, the admission ratio decreases with the increasing number of requesting users as expected. Moreover, the proposed IGA can admit all users when there are sufficient network resources. We study the interactions between two user groups with different requirements and  $M^f = 64$  in Fig. 4-(b). Specifically, group one has varying data transmission demand with maximum waiting time of 1ms while group two requires smaller waiting time compared to group one and each user has a data chunk of  $d_k^{rq} = 250$  bits. In addition, there are 40 users in group one selecting numerology 0 and 40 users in group two using numerology 2 with the same maximum waiting time. For group two, we consider three different values of  $N_k^{rq}$ , which are 1, 2, 3. It can be observed that the higher the data transmission demand per user of group one, the lower the admission ratio that can be achieved by IGA. Also, the improvement in the admission ratio when  $N_k^{rq}$  increases from 1 to 2 is considerably greater than that when  $N_k^{rq}$  increases from 2 to 3. This implies that strict delay requirement may have very negative impact on the system performance.

## V. CONCLUSION

We have proposed two low-complexity algorithms to tackle the scheduling problem for the 5G wireless system supporting heterogeneous services. Numerical results have revealed the desirable performance and complexity for the proposed algorithms.

## REFERENCES

- [1] P. Popovski *et al.*, "5G wireless network slicing for eMBB, URLLC, and mMTC: A communication-theoretic view," *IEEE Access*, vol. 6, pp. 55765–55779, 2018.
- [2] *Physical Channel and Modulation (Release 15)*, document 3GPP-TS-38.211, 2019.
- [3] K. I. Pedersen *et al.*, "A flexible 5G frame structure design for frequency-division duplex cases," *IEEE Commun. Mag.*, vol. 54, no. 3, pp. 53–59, Mar. 2016.
- [4] E. Dahlman *et al.*, *5G NR: The Next Generation Wireless Access Technology*. New York, NY, USA: Academic, 2018.
- [5] R. Kassab *et al.*, "Non-orthogonal multiplexing of ultra-reliable and broadband services in fog-radio architectures," *IEEE Access*, vol. 7, pp. 13035–13049, 2019.
- [6] J. Tang *et al.*, "Service multiplexing and revenue maximization in sliced C-RAN incorporated with URLLC and multicast eMBB," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 4, pp. 881–895, Apr. 2019.
- [7] L. You *et al.*, "Resource optimization with flexible numerology and frame structure for heterogeneous services," *IEEE Commun. Lett.*, vol. 22, no. 12, pp. 2579–2582, Dec. 2018.
- [8] C. McDiarmid, "Pattern minimisation in cutting stock problems," *Discrete Appl. Math.*, vol. 98, nos. 1–2, pp. 121–130, Oct. 1999.