

QoS Analysis of a Scheduling Policy for Heterogeneous Users Employing AMC Jointly with ARQ

Mario Poggioni, *Member, IEEE*, Luca Rugini, *Member, IEEE*, and Paolo Banelli, *Member, IEEE*

Abstract—This paper analyzes the quality of service (QoS) of scheduling algorithms for heterogeneous users in multi-user (MU) wireless systems that take advantage from a cross-layer design with both adaptive modulation and coding (AMC) and automatic repeat request (ARQ). By developing a general theoretical framework based on a finite-state Markov chain, we analytically evaluate the average delay, the packet-loss rate (PLR) and the throughput of a scheduling algorithm based on the channel condition, the buffer occupancy, and the number of retransmissions, of users belonging to different service classes. The key assumption of our analysis, i.e., the independence of the stationary states of different users, highly reduces the computational complexity while preserving a sufficient accuracy. To present the proposed analysis, we also suggest an effective scheduling policy suitable for users belonging to different service classes, compliant with the WiMAX standard. The good match between analytical and simulated performance validates our theoretical findings, and enables the proposed approach to be used for cross-layer optimization.

Index Terms—Adaptive modulation and coding (AMC), automatic repeat request (ARQ), cross-layer design, QoS analysis, scheduling, WiMAX.

I. INTRODUCTION

IN wireless multimedia communications, the presence of multipath fading and time-varying signal-to-noise ratio (SNR) degrades the system performance and hence highly affects the quality of service (QoS) perceived by the user. Due to the increasing request for high data-rate services, a number of techniques have been proposed to counteract the negative effects of fading channels, sometimes designed in a cross-layer fashion [1]-[4]. For instance, adaptive modulation and coding (AMC) [4]-[13] allows for a throughput increase at the physical layer, by using a transmission mode (TM) with higher bit rate whenever the SNR condition is favorable. Moreover, automatic repeat request (ARQ) [7]-[17] at the data link layer further improves the performance by reducing the packet loss rate (PLR), at expense of some additional delay.

The SNR loss caused by fading channels influences the QoS not only in single-user (SU) systems, but also in multi-user (MU) scenarios, where the users may have different QoS

requirements, such as in WiMAX [18]. In wireless MU systems, the scheduling task becomes crucial, since the negative effect of fading can be overturned by scheduling a user with good channel condition. The throughput increase provided by channel-aware schedulers is usually significant, because of the high probability that at least one user has good channel condition. This effect is known as MU diversity (MUD) [19]. Anyway, the design of the scheduling algorithm, as well as the QoS performance analysis, becomes challenging.

Herein we focus on cross-layer scheduling that combines AMC with ARQ. A smart combination of AMC and ARQ can increase the throughput and reduce the PLR [8]-[13], because the error-correcting capabilities of ARQ allow for higher modulation rates. However, the additional delay associated with the packet retransmissions can be unacceptable, e.g., for real-time (RT) applications. Thus, a main goal is to find a scheduling algorithm with a good trade-off between average delay and PLR, depending on the specific application. An equally important goal is the capability of providing a QoS performance analysis for scheduling algorithms whose decisions are based on AMC-ARQ parameters, such as the channel quality, the buffer occupancy, and the number of retransmissions, when users have different QoS requirements.

Among the cross-layer designs for SU scenarios, the authors of [11] propose a delay-constrained AMC scheme and analyzes the PLR and the throughput by means of a finite-state Markov chain [14]. In [20], the delay constraint is incorporated into the rate constraint, and the power-rate adaptation strategy is obtained using the effective capacity. The policy obtained in [20] turns out to be a trade-off between time-domain waterfilling, which is suitable when the delay constraint is loose, and truncated channel inversion, which is suitable when the delay requirement is stringent. In addition, a joint AMC-ARQ design is proposed in [13], where the TM is dynamically chosen using the instantaneous SNR. From the probability distribution of the user state, closed form expressions for the PLR, the average delay, and the throughput are obtained [13], enabling the throughput maximization via exhaustive search. However, the model of [13] only considers an SU scenario.

Delay-optimal MU designs include [21]-[24]. In [21], it is shown that the delay-optimal policy assigns the transmission to the user with longest connected queue (LCQ), i.e., to the user with longest buffer occupancy among those users whose channel condition is good enough for transmission. Even if [21] assumes the same QoS for all users, it is clear that,

Paper approved by Y. Fang, the Editor for Wireless Networks of the IEEE Communications Society. Manuscript received May 20, 2009; revised February 22, 2010.

The authors are with the Department of Electronic and Information Engineering, University of Perugia, Perugia 06125, Italy (e-mail: {mario.poggioni, luca.rugini, paolo.banelli}@diei.unipg.it).

Digital Object Identifier 10.1109/TCOMM.2010.09.090276

when the users have different QoS requirements, delay-optimal schedulers cannot disregard the queue backlog of the users [22]. A delay-optimal power control and subcarrier allocation strategy is proposed in [23] for orthogonal frequency-division multiple access (OFDMA) systems. In [24], a delay-optimal power control and precoder adaptation policy is designed for multiple-input multiple-output systems. In both [23] and [24], heterogeneous users are considered. In [12], which differently from the delay-optimal strategies [21]-[24] also includes ARQ, the MU design exploits statically pre-assigned frequencies that depend on the user conditions (average SNR, buffer size, packet arrival rate, TM, QoS). A theoretical analysis for the MU scenario is provided in [8]-[9], where a statistical characterization of the QoS is addressed by deriving the cumulative distribution functions for the throughput and the delay. Although multiple service classes are considered, the theoretical analysis in [8]-[9] assume an ARQ with infinite number of retransmissions, i.e., without loss of packets induced by the channel. Moreover, the analytical model of [8]-[9] is applied to scheduling policies, such as max-rate (MR), round-robin (RR), and weighted RR (WRR), that do not consider the buffer occupancies (and the completed number of retransmissions). Differently, delay-optimal scheduling policies like [21] and [22] are strongly based on the buffer occupancies.

In this paper, we propose a theoretical QoS analysis for a broader class of scheduling algorithms that deal with heterogeneous users with different QoS requirements. Specifically, we consider a GI/M/1 service queue, where the users have finite buffer lengths and maximum number of retransmissions. One of the main features of our analysis is its reduced complexity, which is weakly dependent from the number of users. This is obtained by exploiting the key idea that the state probabilities of different users can be assumed as independent. This approximation, which is reasonable in several scenarios, is characterized by a high level of accuracy, as confirmed by extensive simulation results. To simplify the presentation, we develop the analysis for a specific scheduling algorithm that we deem significant and representative, where the different service classes are borrowed from WiMAX (RT, non-real-time (NRT), and best-effort (BE) users) [18]. However, the extension to a broader class of scheduling policies, including WRR, is theoretically equivalent, although tedious in some cases.

The main contributions of this paper, as well as the most important differences with the previous literature, can be summarized as follows.

- Differently from [11] and [13], which deal with AMC and ARQ in the SU case, our theoretical analysis focuses on the MU case, with heterogeneous users belonging to different QoS classes. Although our approach can be considered as a generalization of [13] to the MU case, we remark that a direct generalization of [13] would require a complete characterization of the states of all the users, leading to an exponential complexity in the number of users. On the contrary, our approach avoids the exponential complexity.
- Differently from [8]-[9], we consider scheduling policies that are explicitly based also on buffer occupancies and number of retransmissions. These scheduling schemes

are of great interest because they enable several trade-offs between maximum aggregate throughput (given by MR scheduling), maximum user fairness (given by RR scheduling), and minimum average delay (given by LCQ scheduling in ON/OFF channels). While our approach is valid for a large class of scheduling algorithms, the approach in [8]-[9] can be easily applied only to a subclass of algorithms that are not based on the queue lengths, such as MR, RR, and WRR.

- Differently from [25], we provide a theoretical characterization of the QoS performance. The philosophy of our scheduler is similar to that in [25], but our algorithm reduces the buffer occupancy rather than the instantaneous packet delay: this way, by Little's Theorem [14], also the average delay is reduced. Moreover, in our algorithm, RT users are favored, in order to reduce their instantaneous delay.

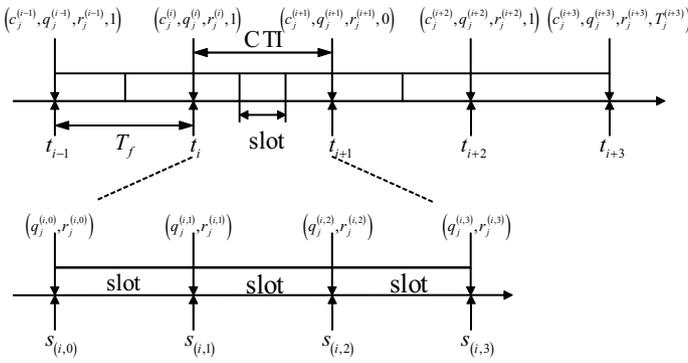
The rest of this paper is organized as follows. Section II presents the system model, and Section III the scheduling algorithm. In Section IV, we develop our theoretical QoS analysis based on the state probability distribution, whose computation is described in Section V. The simulation results of Section VI validate our analytical framework, and Section VII concludes the paper.

II. SYSTEM MODEL

We consider a wireless link shared by U users. We focus on the uplink, but it should be observed that the presented model is valid for the downlink too. The U users are divided in three WiMAX service classes [18]: 1) RT polling service, with guaranteed throughput and delay, e.g., video streaming; 2) NRT polling service, with guaranteed throughput, e.g., FTP; 3) BE, with no guarantees, e.g., e-mail. Constant bit rate users are not considered because their scheduling is trivial. We indicate with S_{RT} , S_{NRT} and S_{BE} the set of users belonging to RT, NRT and BE classes, respectively, using the class order number when opportune, e.g., $S_{RT} \equiv S_1$, $S_{NRT} \equiv S_2$ and $S_{BE} \equiv S_3$. All the users belonging to the same class have identical traffic statistics and QoS constraints. The considered uplink system can be summarized as follows: first, the receiver (base station) collects information about the user conditions (channel quality, buffer occupancy, number of retransmissions) to perform the scheduling decision; after the scheduling decision, the receiver feeds back to the scheduled user the TM to be used. In the whole process, AMC is used to increase the throughput, while ARQ is used to reduce the PLR. Anyway, our detailed assumptions are listed below.

AI: At the physical layer, time is divided in intervals of fixed length T_f seconds. In each interval, denoted by *coherence time interval* (CTI), $K \geq 1$ data-link packets of fixed size (N_P bits) can be transmitted, depending on the channel quality (AMC TM). Therefore, the CTI is divided in $K \geq 1$ slots, one for each packet, as shown in Fig. 1. For instance, in high-quality channels, the AMC selects a higher-order constellation of size M and a convolutional code with high code rate ζ , so that many data-link packets K are mapped into a single CTI, as expressed by

$$K = \frac{WT_f}{N_P} \zeta \log_2 M, \quad (1)$$


 Fig. 1. State and substate transition for user j at time epoch i .

where W is the available bandwidth (see Table I, where $W = 1.08$ MHz, $T_f = 2$ ms, and $N_P = 1080$). We assume negligible overhead, perfect time synchronization among users, and ideal power control (i.e., all the users have the same average SNR).

A2: The propagation channel is modeled by a Nakagami- m frequency-flat block-fading channel over each CTI [26]-[27], i.e., the channel is assumed constant for T_f seconds. The channel variation is captured by a finite-state Markov chain [28], where each channel state is associated to a TM. The AMC partitions the SNR range into $N + 1$ non-overlapping intervals [13] that correspond to the TMs of Table I. The boundary points of the SNR intervals are calculated using the same target packet-error rate (PER) for all the TMs (see Eqs. 1-4 in [13]), as usual in AMC systems. Notably, in OFDMA systems, (frequency-selective) multipath channels are turned into a set of parallel frequency-flat channels, one for each subcarrier. Therefore, the flat-fading channel assumption is compliant with OFDMA-based systems in frequency-selective channels, such as WiMAX [18].

A3: Depending on the instantaneous SNR at each decision epoch, the AMC selector at the receiver chooses the TM, i.e., the modulation-coding pair in Table I. After the scheduler decision, the TM is fed back to the transmitting user. The TM choice is based on perfect *channel state information* (CSI), and the feedback channel is assumed ideal, with no errors and zero latency [13].

A4: We assume Poisson packet arrival processes with packet rate λ_{S_g} packets/s, for any class $g \in \{1, 2, 3\}$. However, our model is quite general, since it can handle any GIM/1 queue, with possibly different arrival processes from class to class. The data-link layer packets, whose length is N_P bits, are mapped on physical layer time slots with different durations, according to the current TM. The size of the buffer of user j is B_j packets. Users of the same class have the same buffer size: $B_j = B_{S_1} = B_{RT}$, $\forall j \in S_{RT}$, $B_j = B_{S_2} = B_{NRT}$, $\forall j \in S_{NRT}$, and $B_j = B_{S_3} = B_{BE}$, $\forall j \in S_{BE}$.

A5: The error detection at the receiver, by means of cyclic redundancy check codes, is assumed perfect. The packets of user j are dropped either when the transmitter buffer is full, or after R_j retransmissions, where $R_j = R_{S_1} = R_{RT}$, $\forall j \in S_{RT}$, $R_j = R_{S_2} = R_{NRT}$, $\forall j \in S_{NRT}$, $R_j = R_{S_3} = R_{BE}$, $\forall j \in S_{BE}$.

 TABLE I
TRANSMISSION MODES

Transmission Mode	TM0	TM1	TM2	TM3	TM4	TM5
Channel index $c_j^{(i)}$	0	1	2	3	4	5
Modulation	-	BPSK	QPSK	QPSK	16-QAM	64-QAM
Code Rate ζ	-	1/2	1/2	3/4	3/4	3/4
Channel quality $K_{c_j^{(i)}}$ slots/CTI	1	1	2	3	6	9

III. SCHEDULING POLICY

In this section, we propose a heuristic scheduling algorithm specifically tailored to heterogeneous users. Our aim is to provide a low PLR and a low average delay to RT users, and a high throughput and a low PLR to NRT users. BE users have low priority and therefore will have high throughput and low delay only when there is enough bandwidth. In other words, we want a scheduling algorithm that gives high throughput while preserving some user fairness, summarized by the user class constraints. For instance, also the WRR scheduling [8] can handle different user classes, but it produces a waste of channel resources, because the fixed time-slot assignment does not enable any MUD [19], i.e., the diversity gathered by scheduling the user with the best channel condition. Obviously, an efficient scheduling algorithm should exploit the MUD, which boosts the system throughput with respect to the SU scenario. In this view, we propose a centralized algorithm that acts at the beginning of each CTI. For simplicity, we assume instantaneous decisions of the scheduler, based on perfect knowledge of all the user states.

The state $\psi_j^{(i)}$ of user j at time instant t_i is defined as

$$\psi_j^{(i)} = (c_j^{(i)}, q_j^{(i)}, r_j^{(i)}), \quad (2)$$

where $c_j^{(i)}$ is the channel state, $q_j^{(i)}$ is the buffer occupancy, and $r_j^{(i)}$ is the number of retransmission attempts. We associate to $c_j^{(i)}$ the channel rate $K_{c_j^{(i)}}$ as in Table I. We also define the utility function

$$\phi_j^{(i)} = K_{c_j^{(i)}} q_j^{(i)} / B_j, \quad (3)$$

which combines the channel quality $K_{c_j^{(i)}}$ and the normalized buffer occupancy $q_j^{(i)} / B_j$, thus enabling a trade-off between throughput and delay. Besides, with reference to the Venn diagram in Fig. 2, we define the following user sets:

- $S_{FB}^{(RT)} = \{j \in S_{RT} : q_j^{(i)} = B_j \wedge TM_j^{(i)} \neq TM0\}$ is the set of RT users with full buffer and a channel good enough for transmission;
- $S_{FB_r}^{(RT)} = \{j \in S_{FB}^{(RT)} : r_j^{(i)} \geq r_k^{(i)}, \forall k \in S_{FB}^{(RT)}\}$ is the subset of $S_{FB}^{(RT)}$ that groups the RT users with maximum number of retransmissions $r_j^{(i)}$;
- $S_{FB_k}^{(RT)} = \{j \in S_{FB_r}^{(RT)} : K_j^{(i)} \geq K_k^{(i)}, \forall k \in S_{FB_r}^{(RT)}\}$ represents the subset of $S_{FB_r}^{(RT)}$ with RT users with the best channel quality;
- $S_{M\Phi}^{(xRT)} = \{j \in S_{RT} \cup S_{NRT} : \phi_j^{(i)} \geq \phi_k^{(i)}, \forall k \in S_{RT} \cup S_{NRT} \wedge \phi_j^{(i)} > 0\}$ collects both RT and NRT users with maximum (nonzero) utility function $\phi_j^{(i)}$;

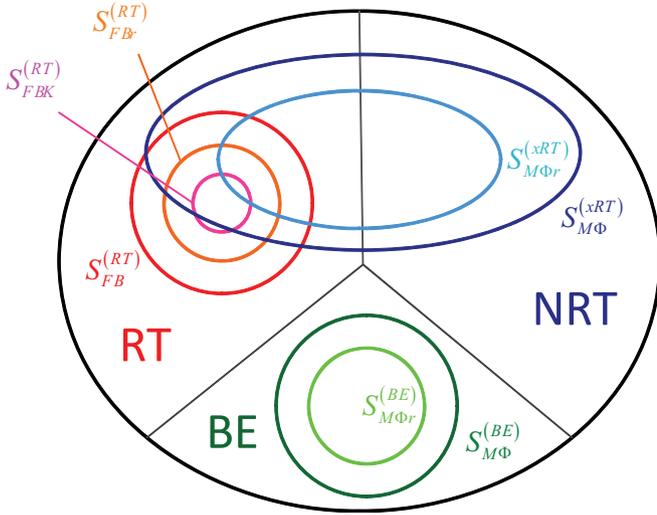


Fig. 2. Users partitioning employed by the scheduling algorithm.

- $S_{M\Phi r}^{(xRT)} = \{j \in S_{M\Phi}^{(xRT)} : r_j^{(i)}/R_j \geq r_k^{(i)}/R_k, \forall k \in S_{M\Phi}^{(xRT)}\}$ is the subset of $S_{M\Phi}^{(xRT)}$ that groups users with maximum normalized number of retransmissions $r_j^{(i)}/R_j$;
- $S_{M\Phi}^{(BE)} = \{j \in S_{BE} : \phi_j^{(i)} \geq \phi_k^{(i)}, \forall k \in S_{BE} \wedge \phi_j^{(i)} > 0\}$ represents the set of BE users with maximum and nonzero $\phi_j^{(i)}$;
- $S_{M\Phi r}^{(BE)} = \{j \in S_{M\Phi}^{(BE)} : r_j^{(i)} \geq r_k^{(i)}, \forall k \in S_{M\Phi}^{(BE)}\}$ is the subset of $S_{M\Phi}^{(BE)}$ with maximum $r_j^{(i)}$.

Using these sets, the proposed algorithm is expressed in Table II, where “%” stands for “comment.”

We remark that the proposed scheduling algorithm specifically tries to reduce the delay of RT users by first assigning the priority to RT users with full buffer, which actually are experiencing the maximum instantaneous delay. This way, also the PLR of RT users is reduced, because it is less probable that their buffers become full and packets are discarded. Since the utility function $\phi_j^{(i)}$ in (3) is proportional to the buffer occupancy, due to the Little’s Theorem [14], the average packet delays of RT and NRT users are also reduced. In addition, BE transmission is authorized only when all RT and NRT users have empty buffers or experience deep fading.

The main difference of the proposed scheduling with respect to [25] is that it considers buffer occupancies instead of instantaneous delays. Therefore, our scheduler grants service to RT users with full buffer, while in [25] users have the priority when their packet delay is higher than a temporal parameter called “time deadline,” whose choice affects the performance. A performance comparison between the two approaches is given in the simulation section.

IV. QUEUING ANALYSIS OF AMC-ARQ-BASED SCHEDULING

Our aim is to derive a probabilistic characterization of the state vector (2) of a generic user. We first summarize the SU analysis carried out in [13], which we subsequently extend to the MU case.

TABLE II
SCHEDULING ALGORITHM

```

01 if  $S_{FB}^{(RT)} \cup S_{M\Phi}^{(xRT)} \neq \emptyset$ 
02   if  $S_{FB}^{(RT)} \neq \emptyset$  % RT with full buffer
03     transmission randomly assigned to user  $\in S_{FBK}^{(RT)}$ ;
04   else
05     if  $S_{M\Phi r}^{(xRT)} \cap S_{RT} \neq \emptyset$  % RT without full buffer
06       transmission randomly assigned to user  $\in S_{M\Phi r}^{(xRT)} \cap S_{RT}$ ;
07     else % NRT transmission
08       transmission randomly assigned to user  $\in S_{M\Phi r}^{(xRT)} \cap S_{NRT}$ ;
09     end
10   end
11 else
12   if  $S_{M\Phi r}^{(BE)} \neq \emptyset$  % BE transmission
13     transmission randomly assigned to user  $\in S_{M\Phi r}^{(BE)}$ ;
14   else
15     no user transmits;
16   end
17 end

```

A. Single-User Analysis

The channel state transitions are modeled by an embedded Markov chain, described by the $(N+1) \times (N+1)$ transition matrix \mathbf{P}_C , whose elements $[\mathbf{P}_C]_{m,n} = P_{m,n}$ are derived by a level crossing rate analysis [29]. The complete SU state is described by $\psi^{(i)}$ in (2). If $c^{(i)} = n$, i.e., the AMC selector chooses TM n , then the CTI is divided in $K_{c^{(i)}} = K_n$ slots of duration T_f/K_n , as in Fig. 1. The queuing process is described by an embedded Markov chain, where the transitions of the substate $(q^{(i)}, r^{(i)})$ in each slot are described by the square transition matrix \mathbf{T}_n of size $(B+1)(R+1)$ (see Appendix A of [30]), where B is the buffer length and R is the maximum number of retransmissions. As explained in [13], the evolution of the substate $(q^{(i)}, r^{(i)})$ during the whole CTI is described by the matrix $\mathbf{T}_n^{K_n}$. The packet arrival process and the channel transitions are assumed independent, and hence the transitions of the whole state $\psi^{(i)} = (c^{(i)}, q^{(i)}, r^{(i)})$ are described by the square matrix with size $(N+1)(B+1)(R+1)$ (see Eqs. 27-28 and Prop. 1 in [13])

$$\mathbf{P} = \begin{bmatrix} P_{0,0}\mathbf{T}_0 & \cdots & P_{0,N}\mathbf{T}_0 \\ \vdots & \ddots & \vdots \\ P_{N,0}\mathbf{T}_N^{K_N} & \cdots & P_{N,N}\mathbf{T}_N^{K_N} \end{bmatrix}. \quad (4)$$

The stationary state probability vector, denoted by $\boldsymbol{\pi} = [\pi_{(0,0,0)}, \pi_{(0,0,1)}, \dots, \pi_{(N,B,R-1)}, \pi_{(N,B,R)}]$, where $\pi_{(c,q,r)}$ is the stationary probability of the state (c, q, r) , is classically calculated as [14]

$$\boldsymbol{\pi}\mathbf{P} = \boldsymbol{\pi}, \quad \sum_{c,q,r} \pi_{(c,q,r)} = 1. \quad (5)$$

From the obtained left eigenvector $\boldsymbol{\pi}$, it is possible to analytically derive the average delay τ_D , the PLR P_{PLR} , and the throughput Σ , by means of Equations 31-44 in [13], which we simply summarize by $\tau_D = f_D(\boldsymbol{\pi})$, $P_{PLR} = f_{PLR}(\boldsymbol{\pi})$, and $\Sigma = \lambda(1 - P_{PLR})$, where λ is the packet arrival rate.

B. Multi-User Analysis

From a theoretical point of view, an MU system with U users can be described by a superstate $\mathbf{s}^{(i)} = (c_1^{(i)}, q_1^{(i)}, r_1^{(i)}, \dots, c_U^{(i)}, q_U^{(i)}, r_U^{(i)})$ and by the corresponding transition matrix $\tilde{\mathbf{P}}$. This way, any scheduling policy based on memoryless utility functions $\varphi_j^{(i)} = f_j(c_j^{(i)}, q_j^{(i)}, r_j^{(i)})$, $\forall j \in \{1, \dots, U\}$, can be easily modeled. The superstate transitions are described by a finite-length Markov chain that is both irreducible (because there is a single communicating class [14]) and aperiodic (because there is at least one superstate with nonzero self-transition probability [31]). As a consequence, the stationary state probability vector $\boldsymbol{\pi}^{(s)}$ exists and is unique, and could be calculated as the left eigenvector of $\tilde{\mathbf{P}}$ associated to the unit eigenvalue, as done in (5) for SU systems. Unfortunately, this approach is impractical, because the number of superstates $L_{\boldsymbol{\pi}^{(s)}} = \prod_{g=1}^3 [(N+1)(B_{S_g}+1)(R_{S_g}+1)]^{U_g}$, where U_g is the number of users in the g th class, would be too high in real systems. Thus, the computational complexity would be exponential in the number of users $U = \sum_{g=1}^3 U_g$. To reduce complexity, we want to consider only a unique user per class, which is referred to as the *representative user* of that class, so that the number of states is greatly reduced. This SU-per-class approach is useful to describe the equilibrium condition, where the probability that user j is in state $\boldsymbol{\psi}_j = (c_j, q_j, r_j)$ is the same for all the users in the same class of user j , since users in the same class have equal traffic intensity, equal channel statistics, and equal QoS requirements.

To develop our SU-like approach for MU systems, we exploit an assumption of *independence* of the probabilities of the stationary states for all the U users, expressed by

$$\pi_{(c_1, q_1, r_1, \dots, c_U, q_U, r_U)} = \pi_{(c_1, q_1, r_1)} \dots \pi_{(c_U, q_U, r_U)}, \quad (6)$$

where $\pi_{(c_1, q_1, r_1, \dots, c_U, q_U, r_U)}$ is the stationary probability of the superstate \mathbf{s} , i.e., $\pi_{(c_1, q_1, r_1, \dots, c_U, q_U, r_U)}$ is a generic element of $\boldsymbol{\pi}^{(s)}$, while $\pi_{(c_j, q_j, r_j)}$ is the stationary probability of the state $\boldsymbol{\psi}_j = (c_j, q_j, r_j)$, i.e., $\pi_{(c_j, q_j, r_j)}$ is the corresponding element of $\boldsymbol{\pi}_j = \boldsymbol{\pi}^{(\boldsymbol{\psi}_j)}$. Obviously, $\boldsymbol{\pi}_j$ exists and is unique, and it could be theoretically obtained by marginalization of $\boldsymbol{\pi}^{(s)}$ over all the users other than j . It is worth noting that the approximation (6) is not valid for all the scheduling policies. By comparing analytical and simulation results, we have noted that (6) is more accurate for those scheduling policies that do not take into account the substates (q_j, r_j) , e.g., MR scheduling [9], and for RR-based policies. Conversely, the approximation (6) could potentially be less accurate for scheduling policies based on buffer occupancies, such as our algorithm of Table II. However, the dependence among the users' substates $\{(c_j, q_j, r_j)\}$ primarily arises when the buffers have saturated. When the scheduler is efficient and there are many users, the dependence among the users substates is generally weak, and hence (6) leads to accurate results. For the proposed scheduler, the accuracy of (6) is verified in Section VI, where we compare analytical and simulation results.

The assumption in (6) permits to focus on the evolution of the single state $\boldsymbol{\psi}_j = (c_j, q_j, r_j)$ of the user j . Basically, we want to derive, for the user j , a transition matrix $\tilde{\mathbf{P}}_j$ that plays the same role of $\tilde{\mathbf{P}}$ in (4) for SU systems. Since

the state evolution for the representative user j also depends on the channel states, buffer occupancies, and retransmission attempts of all the other users, it is clear that the desired $\tilde{\mathbf{P}}_j$ will depend on the probabilities that the other users are in specific states. At the steady-state equilibrium, these probabilities are the stationary probabilities for the states of the other users. In other words, $\tilde{\mathbf{P}}_j$ should incorporate some information about the other users. This information is represented by the dependence of $\tilde{\mathbf{P}}_j$ from $\boldsymbol{\pi}_k$, $\forall k \neq j$, which is developed in Section V. In addition, the evolution of the state $\boldsymbol{\psi}_j$ also depends on the scheduler decision. Indeed, when user j is transmitting (Tx), the state transition matrix $\tilde{\mathbf{P}}_j^{(Tx)}$ is different from the matrix $\tilde{\mathbf{P}}_j^{(no-Tx)}$ used when the user j is not transmitting (no-Tx). Both matrices, omitted for lack of space, are derived in Appendix A of [30]. However, instead of dealing with a scheduling-dependent time-varying Markov chain, we introduce (as in [8]-[9]) the extra substate T_j that represents the Tx condition: $T_j = 1$ when user j is Tx, and $T_j = 0$ when user j is no-Tx. By this approach, we can define an extended-state probability vector

$$\tilde{\boldsymbol{\psi}}_j = (\boldsymbol{\psi}_j, T_j) = (c_j, q_j, r_j, T_j), \quad (7)$$

which evolves according to a state-transition matrix $\tilde{\tilde{\mathbf{P}}}_j$ (to be determined) that depends on $\{\boldsymbol{\pi}_k, k \neq j\}$, $\tilde{\mathbf{P}}_j^{(Tx)}$, and $\tilde{\mathbf{P}}_j^{(no-Tx)}$. The state-transition matrix $\tilde{\tilde{\mathbf{P}}}_j$ is derived in the next section.

We now briefly discuss some properties of the Markov chain obtained from (7). First, it is straightforward to show that the Markov chain is irreducible, since there is a single communicating class. Indeed, it is evident that, starting from any state $\tilde{\boldsymbol{\psi}}_j^* = (c_j^*, q_j^*, r_j^*, T_j^*)$ (characterized by a channel condition c_j^* , a buffer occupancy q_j^* , a retransmission number r_j^* , and a transmission state T_j^*), the physical phenomenon induced by the memoryless scheduler does not prevent to reach any other state $\tilde{\boldsymbol{\psi}}_j^+ = (c_j^+, q_j^+, r_j^+, T_j^+)$. Note that some states, such as those of type $(c, 0, r, T)$ with $r > 0$, are not possible and hence are not considered into the Markov chain. Second, since the Markov chain is irreducible and finite, all the states are positive recurrent [14]. Moreover, it is easy to verify that the Markov chain is aperiodic. Indeed, since the transition probability from $(0, 0, 0, 0)$ to itself is positive, the state $(0, 0, 0, 0)$ is aperiodic, such as the whole Markov chain [31]. Since all the states are both aperiodic and positive recurrent, i.e., ergodic, the steady-state distribution probability vector $\tilde{\boldsymbol{\pi}}_j$ exists and is unique. Similarly to the SU case (5), this distribution can be calculated using

$$\tilde{\boldsymbol{\pi}}_j \tilde{\tilde{\mathbf{P}}}_j = \tilde{\boldsymbol{\pi}}_j, \quad \tilde{\boldsymbol{\pi}}_j \mathbf{1} = 1, \quad (8)$$

where $\mathbf{1}$ is the all-ones column vector. The stationary probability of a generic state $\tilde{\boldsymbol{\psi}}_j = (c_j, q_j, r_j, T_j)$ is expressed by $\tilde{\pi}_{(c_j, q_j, r_j, T_j)} = [\tilde{\boldsymbol{\pi}}_j]_{\tilde{k}_j}$, where $\tilde{k}_j = 2(B_j + 1)(R_j + 1)c_j + 2(R_j + 1)q_j + 2(r_j + 1) + T_j + 1$, due to the fact that the variables in (7) are used from left to right to identify the index \tilde{k}_j . Obviously, $\pi_{(c_j, q_j, r_j)} = \tilde{\pi}_{(c_j, q_j, r_j, 0)} + \tilde{\pi}_{(c_j, q_j, r_j, 1)}$. Therefore, from $\tilde{\boldsymbol{\pi}}_j$ we can simply obtain $\boldsymbol{\pi}_j$, which permits to analytically derive the average delay $\tau_{D,j}$, the PLR $P_{PLR,j}$, and the throughput Σ_j , of user j , by means of $\tau_{D,j} = f_D(\boldsymbol{\pi}_j)$, $P_{PLR,j} = f_{PLR}(\boldsymbol{\pi}_j)$, and $\Sigma_j = \lambda_j(1 - P_{PLR,j})$, where λ_j

is the packet rate of user j , as in the SU case [13]. As a result, thanks to the independence assumption (6), we have converted the MU problem into U SU-equivalent problems, one for each user. Actually, only three out of U problems are really different, since the users in the same QoS class have the same state probability distribution π_j .

V. STATE TRANSITION MATRIX AND STEADY-STATE PROBABILITY DISTRIBUTION

In this section, we explain how to calculate the state transition matrix $\tilde{\mathbf{P}}_j$ for user j used in (8) to derive the steady-state probability distribution $\tilde{\pi}_j$ for user j . The $2L\pi_j \times 2L\pi_j$ state transition matrix $\tilde{\mathbf{P}}_j$ associated to $\tilde{\psi}_j$ is defined as

$$\tilde{\mathbf{P}}_j = \begin{bmatrix} \mathbf{M}_{(0,0,0) \rightarrow (0,0,0)} & \cdots & \mathbf{M}_{(0,0,0) \rightarrow (N,B_j,R_j)} \\ \vdots & \ddots & \vdots \\ \mathbf{M}_{(N,B_j,R_j) \rightarrow (0,0,0)} & \cdots & \mathbf{M}_{(N,B_j,R_j) \rightarrow (N,B_j,R_j)} \end{bmatrix}, \quad (9)$$

$$\mathbf{M}_{\psi_j^{(i-1)} \rightarrow \psi_j^{(i)}} = \begin{bmatrix} P_{(\psi_j^{(i-1)},0) \rightarrow (\psi_j^{(i)},0)} & P_{(\psi_j^{(i-1)},0) \rightarrow (\psi_j^{(i)},1)} \\ P_{(\psi_j^{(i-1)},1) \rightarrow (\psi_j^{(i)},0)} & P_{(\psi_j^{(i-1)},1) \rightarrow (\psi_j^{(i)},1)} \end{bmatrix}, \quad (10)$$

where we have set up the compact notation $p_E = \Pr\{E\}$ to indicate the probability of an event E . The elements of (10) describe the substate transitions $\tilde{\psi}_j^{(i-1)} \rightarrow \tilde{\psi}_j^{(i)}$ for Tx and no-Tx cases. Thus, by conditional probability rules, we derive

$$P_{(\psi_j^{(i-1)},0) \rightarrow (\psi_j^{(i)},0)} = [\mathbf{P}_j^{(no-Tx)}]_{k_j^{(i-1)}, k_j^{(i)}} P_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=0 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}, \quad (11)$$

$$P_{(\psi_j^{(i-1)},0) \rightarrow (\psi_j^{(i)},1)} = [\mathbf{P}_j^{(no-Tx)}]_{k_j^{(i-1)}, k_j^{(i)}} P_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}, \quad (12)$$

$$P_{(\psi_j^{(i-1)},1) \rightarrow (\psi_j^{(i)},0)} = [\mathbf{P}_j^{(Tx)}]_{k_j^{(i-1)}, k_j^{(i)}} P_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=0 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}, \quad (13)$$

$$P_{(\psi_j^{(i-1)},1) \rightarrow (\psi_j^{(i)},1)} = [\mathbf{P}_j^{(Tx)}]_{k_j^{(i-1)}, k_j^{(i)}} P_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}, \quad (14)$$

where the index $k_j^{(i)} = (B_j + 1)(R_j + 1)c_j^{(i)} + (R_j + 1)q_j^{(i)} + r_j^{(i)} + 1$ is compliant with the substate ordering $(c_j^{(i)}, q_j^{(i)}, r_j^{(i)})$. Therefore, to compute the matrix $\tilde{\mathbf{P}}_j$ in (9), we have to compute the transition probabilities in (11)-(14) for each transition of $\psi_j^{(i)}$, based on the adopted scheduling policy. Actually, we only derive the transition probabilities with final state $T_j^{(i)} = 1$ in (12) and (14), which are simpler to obtain, and we exploit the relations $p_{(\psi_j^{(i-1)},0) \rightarrow (\psi_j^{(i)},0)} = [\mathbf{P}_j^{(no-Tx)}]_{k_j^{(i-1)}, k_j^{(i)}} - P_{(\psi_j^{(i-1)},0) \rightarrow (\psi_j^{(i)},1)}$ and $P_{(\psi_j^{(i-1)},1) \rightarrow (\psi_j^{(i)},0)} = [\mathbf{P}_j^{(Tx)}]_{k_j^{(i-1)}, k_j^{(i)}} - P_{(\psi_j^{(i-1)},1) \rightarrow (\psi_j^{(i)},1)}$ to evaluate the transition probabilities with $T_j^{(i)} = 0$ in (11) and (13). Remembering that $\mathbf{P}_j^{(Tx)}$ and $\mathbf{P}_j^{(no-Tx)}$ are available in Appendix A of [30], it turns out that the computation of $\tilde{\mathbf{P}}_j$ requests the computation of $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ and $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ in (12) and (14), respectively, $\forall \psi_j^{(i-1)}, \forall \psi_j^{(i)}$.

Since only three out of U problems (8) are different, instead of solving (8) U times, one for each user $j = 1, \dots, U$, we solve $\tilde{\pi}_{S_g} \tilde{\mathbf{P}}_{S_g} = \tilde{\pi}_{S_g}$ for $g = 1, 2, 3$, where $\tilde{\pi}_{S_g}$ and $\tilde{\mathbf{P}}_{S_g}$ represent a generic couple $\tilde{\pi}_j$ and $\tilde{\mathbf{P}}_j$ for $j \in S_g$. This fact highly reduces the complexity of the evaluation of $\tilde{\mathbf{P}}_{S_g} = \tilde{\mathbf{P}}_j$, $j \in S_g$, because $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ and $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ in (12) and (14) need to be evaluated only for the three representative users. Moreover, since the proposed scheduling policy is memoryless, $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ and $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ do not depend on the previous state $\psi_j^{(i-1)}$, but only on the current state $\psi_j^{(i)}$. This property further reduces complexity, because many transition probabilities have the same value; for instance, $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}} = p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}, \forall \psi_j^{(i-1)}$.

As it will be evident in the following, the transition probabilities (12) and (14) depend on the stationary state probabilities of all the users. Therefore, $\tilde{\mathbf{P}}_{S_g}$ depends on $\{\tilde{\pi}_{S_1}, \tilde{\pi}_{S_2}, \tilde{\pi}_{S_3}\}$. To overcome this cross-dependence, we resort to an *iterative procedure*. At the n th iteration, $\tilde{\mathbf{P}}_{S_g, n}$ is computed from the vectors $\{\tilde{\pi}_{S_1, n-1}, \tilde{\pi}_{S_2, n-1}, \tilde{\pi}_{S_3, n-1}\}$ available from the previous iteration, and then $\tilde{\pi}_{S_g, n}$ is updated using $\tilde{\mathbf{P}}_{S_g, n}$. As initialization, we select $\tilde{\pi}_{S_g, 0}$ to have a uniform distribution, excluding the impossible states. Summarizing, our iterative procedure consists of two steps.

Step 1: Computation of the state transition matrices $\{\tilde{\mathbf{P}}_{S_g, n}\}$, where $\tilde{\mathbf{P}}_{S_g, n} = f_{S_g}(\tilde{\pi}_{S_1, n-1}, \tilde{\pi}_{S_2, n-1}, \tilde{\pi}_{S_3, n-1})$ compactly summarizes (9)-(14), (20)-(21), (28)-(35), for $g \in \{1, 2, 3\}$.

Step 2: Computation of the steady-state probability distributions $\{\tilde{\pi}_{S_g, n}\}$, for $g \in \{1, 2, 3\}$, by $\tilde{\pi}_{S_g, n} = f_{LE1}(\tilde{\mathbf{P}}_{S_g, n-1})$, where $f_{LE1}(\mathbf{A})$ is the left eigenvector of \mathbf{A} associated to the unit eigenvalue.

A. Step 1: Computation of the State Transition Matrices

By (9)-(14), the state transition matrices $\{\tilde{\mathbf{P}}_{S_1, n}, \tilde{\mathbf{P}}_{S_2, n}, \tilde{\mathbf{P}}_{S_3, n}\}$ depend on the transition probabilities $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ and $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$, $\forall \psi_j^{(i-1)}$. In the following, we briefly outline how to derive these transition probabilities for a memoryless scheduling scheme. Successively, we specialize the analytical derivation to obtain the results for the proposed algorithm.

1) Outline of the Analytical Derivation in Step 1:

- To distinguish the Tx and no-Tx cases, we first introduce, for each user u , the events

$$\begin{aligned} TX0_u &:= \{T_u^{(i-1)} = 0\}, \\ TX1_u &:= \{T_u^{(i-1)} = 1\}. \end{aligned} \quad (15)$$

Then, for each set S of users (i.e., the sets defined in Section III), we derive the conditional probability that a user u belongs to the set S , conditioned on $TX0_u$, and conditioned on $TX1_u$.

- For each class S_g , we consider the representative user $j \in S_g$, and, for each state $\psi_j^{(i)}$, we calculate the conditional probability that j is scheduled. For each class, two conditional probabilities have to be found: $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ and $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$. Any single user j that can be scheduled must belong to one of the sets enabled for transmission, i.e., to one of the four sets expressed in the Lines 03, 06, 08, and 13 of Table II. Let us denote this set with \tilde{S} . All the other $U - 1$ users can be partitioned in three sets: \tilde{S}_{HP} , which contains the users with transmission priority *higher* than user j ; $\tilde{S}_{LP} = (S_1 \cup S_2 \cup S_3) \setminus (\tilde{S}_{HP} \cup \tilde{S})$, which contains the users with *equal* priority with respect to j ; and $\tilde{S}_{LP} = S_g \setminus (\tilde{S}_{HP} \cup \tilde{S})$, which contains the users with *lower* priority.
- The conditional probability $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$ is evaluated as the probability that $\tilde{S}_{HP} = \emptyset$ and user j is randomly chosen among the users in \tilde{S} . By the independence assumption (6), this probability can be computed by expressions that contain terms like

$$\sum_{u=0}^{U-1} \binom{U-1}{u} w(u) p_{EP}^u p_{LP}^{U-1-u}, \quad (16)$$

where p_{EP} is the probability that user $u \in \tilde{S}_{EP}$, p_{LP} is the probability that user $u \in \tilde{S}_{LP}$, and $w(u) = 1/(u+1)$ is a weight that takes into account the random choice among the $u+1$ users in \tilde{S} . In (16), the probability p_{HP} that user $u \in \tilde{S}_{HP}$ does not appear, because $T_j^{(i)} = 1$ implies that there are no users in \tilde{S}_{HP} , and hence the multiplicative term $p_{HP}^0 = 1$ is hidden.

- To evaluate the conditional probability $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i-1)} \rightarrow \psi_j^{(i)}}$, we denote with $u_{Tx}^{(i-1)} \neq j$ the index of the Tx user in the $(i-1)$ th CTI. Few subcases have to be considered, depending on $u_{Tx}^{(i-1)} \in \tilde{S}$ or $u_{Tx}^{(i-1)} \in \tilde{S}_{LP}$ in the i th CTI, and depending on the class of $u_{Tx}^{(i-1)}$. The conditional probability can be expressed as a weighted sum of few terms similar to (16), one for each subcase.

Now we focus on our scheduling algorithm. We evaluate the transition probabilities for RT, NRT, and BE users. For each class, we compute the probability that no other user, even of a different class, is scheduled instead of the representative user j .

2) *Transition Probabilities for RT Users:* According to our scheduling algorithm, an RT user is scheduled when Line 03 or Line 06 of Table II is executed. We distinguish between these events, identified by E1 and E2, respectively.

Event E1 (Line 03 of Table II is executed). This happens when both the conditions of Lines 01-02 are true, which leads to $\tilde{S} = S_{FBK}^{(RT)}$. In this case, the buffer $q_j^{(i)} = B_{RT}$ is full, and hence $j \in S_{FB}^{(RT)}$ when $c_j^{(i)} > 0$ ($TM_j \neq TM_0$) (The case $c_j^{(i)} = 0$ is ignored because user j would be surely no-Tx). Bearing in mind the definition of $S_{FBK}^{(RT)}$ in Section III, we define suitable conditions for $S_{FB}^{(RT)}$ and $S_{FB_r}^{(RT)}$. However, even

if $j \in S_{FB}^{(RT)}$, a user $u \in S_{FB}^{(RT)}$, with $u \neq j$, can be scheduled instead of user j : hence, to derive $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}}$ and $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}}$, we first find the probability that $u \notin S_{FB}^{(RT)}$ for any user $u \neq j$. To this aim, we define the event that user u has a lower buffer (LB) occupancy as

$$LB_u := \{u \notin S_{FB}^{(RT)}\}. \quad (17)$$

Since our scheduling policy takes ARQ into account, we have to consider also the case $u \in S_{FB}^{(RT)} \setminus S_{FB_r}^{(RT)}$, i.e., the case when user u cannot transmit because $r_u < r_j$. Thus, we define the event that user u has a lower retransmission (LR) number LR_u than the user j as

$$LR_u := \{u \in S_{FB}^{(RT)} \setminus S_{FB_r}^{(RT)}, j \in S_{FB_r}^{(RT)}\}. \quad (18)$$

Similarly, since the scheduler takes AMC into account, we consider $u \in S_{FB_r}^{(RT)} \setminus S_{FBK}^{(RT)}$, where the user u is characterized by a rate $K_u^{(i)}$ lower (or equal) than $K_j^{(i)}$, as expressed by

$$LK_u := \left\{ u \in S_{FB_r}^{(RT)} \setminus S_{FBK}^{(RT)}, j \in S_{FBK}^{(RT)} \right\},$$

$$EK_u := \left\{ u \in S_{FBK}^{(RT)}, j \in S_{FBK}^{(RT)} \right\}. \quad (19)$$

As derived in Appendix B of [30], the events in (17)-(19) allow to express $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}}$ as in (20), shown at the top of the next page, where for simplicity we omit the user index u into the conditional probabilities, which are the same for all RT users, e.g., $p_{LB_u|TX0_u} = p_{LB|TX0}$. The probabilities $p_{LB|TX0}$, $p_{LR|TX0}$, and the others, are derived in Section V.A.5.

Let us now assess the transition probability $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}}$. Since in this case user j was no-Tx during the $(i-1)$ th CTI, $u_{Tx}^{(i-1)} \neq j$ was Tx (we define $u_{Tx}^{(i-1)} = 0$ when no user was Tx in the $(i-1)$ th CTI). By distinguishing between $u_{Tx}^{(i-1)} \in S_{RT}$ and $u_{Tx}^{(i-1)} \notin S_{RT}$, the probability $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}}$, derived in Appendix B of [30], can be simplified as in (21), shown at the top of the next page, where $w = (U_{RT} - 1) p_{TX1_{RT}}$, and $p_{TX1_{RT}} = \Pr\{T_u^{(i-1)} = 1\}$ stands for p_{TX1_u} in (15) for a generic user $u \in S_{RT}$. This is consistent with (20), which indeed is equivalent to (21) when $U_{RT} = 1$. The probabilities in the right-hand side of (21), as well as p_{TX1_u} , are derived in Section V.A.5.

Event E2 (Line 06 of Table II is executed). This happens when the condition in Line 02 is false, so that the buffer $q_j^{(i)} < B_j$ is not full. In this case, $\tilde{S} = S_{M\Phi_r}^{(xRT)} \cap S_{RT}$. To evaluate the transition probabilities, we have to consider that also NRT users could transmit, depending on their utility function $\phi_u^{(i)}$.

Let us first derive $p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}}$. First of all, analogously to E1, we have to determine if a generic RT or NRT user belongs to $S_{M\Phi}^{(xRT)}$, which is a set included in the definition of $S_{M\Phi_r}^{(xRT)}$, and in the case of RT users, if they belong to $S_{FB}^{(RT)}$. We therefore define the events

$$L\Phi LB_u := \{u \in (S_{RT} \setminus S_{FB}^{(RT)}) \setminus S_{M\Phi}^{(xRT)}, j \in S_{M\Phi}^{(xRT)}\}, \quad (22)$$

$$L\Phi_u := \{u \in (S_{NRT} \setminus S_{M\Phi}^{(xRT)}) \cup (S_{BE} \setminus S_{M\Phi}^{(BE)}),$$

$$j \in S_{M\Phi}^{(xRT)} \cup S_{M\Phi}^{(BE)}\}, \quad (23)$$

$$p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} = \sum_{u=0}^{U_{RT}-1} \binom{U_{RT}-1}{u} \frac{1}{u+1} p_{EK|TX0}^u (p_{LB|TX0} + p_{LR|TX0} + p_{LK|TX0})^{U_{RT}-1-u} \quad (20)$$

$$p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} = (1-w) p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} + w \sum_{u=0}^{U_{RT}-2} \binom{U_{RT}-2}{u} p_{EK|TX0}^u \times (p_{LB|TX0} + p_{LR|TX0} + p_{LK|TX0})^{U_{RT}-2-u} \left(\frac{p_{LB|TX1} + p_{LR|TX1} + p_{LK|TX1}}{u+1} + \frac{p_{EK|TX1}}{u+2} \right) \quad (21)$$

where (22) is related to RT users and (23) to NRT and BE users (the events related to BE users will be used later on). There are also other considerations that help us in identifying useful events. If there is no user $u \neq j$ belonging to $S_{FB}^{(RT)} \cup S_{M\Phi}^{(xRT)}$, we have $S_{FB}^{(RT)} = \emptyset$ and $S_{M\Phi}^{(xRT)} = \{j\}$, and consequently user j certainly transmits. Anyway, even if some users $u \neq j$ belong to $S_{M\Phi}^{(xRT)}$, when none of them belong to $S_{M\Phi r}^{(xRT)}$, user j certainly transmits. Thus, we also define the events:

$$LRLB_u := \{u \in ((S_{RT} \cap S_{M\Phi}^{(xRT)}) \setminus S_{M\Phi r}^{(xRT)}) \setminus S_{FB}^{(RT)}, j \in S_{M\Phi r}^{(xRT)}\}, \quad (24)$$

$$ERLB_u := \{u \in S_{M\Phi r}^{(xRT)} \cap S_{RT}, j \in S_{M\Phi r}^{(xRT)}\}, \quad (25)$$

$$L\Phi R_u := \{u \in ((S_{NRT} \cap S_{M\Phi}^{(xRT)}) \setminus S_{M\Phi r}^{(xRT)}) \cup (S_{M\Phi}^{(BE)} \setminus S_{M\Phi r}^{(BE)}), j \in S_{M\Phi r}^{(xRT)} \cup S_{M\Phi r}^{(BE)}\}, \quad (26)$$

$$E\Phi R_u := \{u \in (S_{M\Phi r}^{(xRT)} \cap S_{NRT}) \cup S_{M\Phi r}^{(BE)}, j \in S_{M\Phi r}^{(xRT)} \cup S_{M\Phi r}^{(BE)}\}, \quad (27)$$

where (24)-(25) are valid for RT users, and (26)-(27) for NRT and BE users. From Appendix C of [30], we have (28), shown at the top of the next page, where $w_0 = (p_{L\Phi|TX0} + p_{E\Phi R|TX0})^{U_{NRT}}$.

We now focus on $p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}}$. In this case, we have to distinguish among the three cases: $u_{Tx}^{(i-1)} \in S_{RT}$, $u_{Tx}^{(i-1)} \in S_{NRT}$, and $u_{Tx}^{(i-1)} \notin S_{RT} \cup S_{NRT}$. In Appendix C of [30], we have derived the transmission probabilities for user j for all the cases, which lead to (29), shown at the top of the next page, where $w_1 = U_{NRT} p_{TX1_{NRT}} ((p_{L\Phi|TX1} + p_{E\Phi R|TX1}) w_0^{-1} - 1) - w_2 w_0^{-1} + 1$ and $w_2 = (U_{RT} - 1) p_{TX1_{RT}} w_0$.

3) *Transition Probabilities for NRT Users:* This case corresponds to Line 08 of Table II, i.e., an NRT user transmits. We consider again the sets $S_{FB}^{(RT)}$, $S_{FB r}^{(RT)}$, $S_{FBK}^{(RT)}$, $S_{M\Phi}^{(xRT)}$ and $S_{M\Phi r}^{(xRT)}$ defined in Section III, and the events defined in (22)-(27). When $j \in S_{M\Phi r}^{(xRT)} \cap S_{NRT}$, user j transmits only if $S_{M\Phi r}^{(xRT)} \cap S_{RT} = \emptyset$. RT users are no-Tx if $L\Phi L B_u$ or $LRLB_u$ is verified. For the other NRT users, we have to consider $L\Phi_u$ and $L\Phi R_u$ when $u \notin S_{M\Phi r}^{(xRT)} \cap S_{NRT}$, and $E\Phi R_u$ when $u \in S_{M\Phi r}^{(xRT)} \cap S_{NRT}$. From the Appendix D of [30], we obtain (30) and (31), shown in the next page, where $w_3 = (p_{L\Phi LB|TX0} + p_{LRLB|TX0})^{U_{RT}}$, $w_4 = 1 - U_{RT} p_{TX1_{RT}} - (U_{NRT} - 1) p_{TX1_{NRT}}$, $w_5 = U_{RT} p_{TX1_{RT}} (p_{L\Phi LB|TX0} + p_{LRLB|TX0})^{U_{RT}-1} (p_{L\Phi LB|TX1} + p_{LRLB|TX1})$, and $w_6 = (U_{NRT} - 1) p_{TX1_{NRT}} (p_{L\Phi LB|TX0} + p_{LRLB|TX0})^{U_{RT}}$.

4) *Transition Probabilities for BE Users:* This is the last case of our algorithm, when Line 13 of Table II is executed. We exploit the definitions of $S_{M\Phi r}^{(BE)}$ and $S_{M\Phi}^{(BE)}$ through (23), (26), and (27). We remark that BE users are scheduled for transmission only if all RT and NRT users cannot transmit, due to bad channels or empty buffers, expressed by the event $C0Q0 = \{(c_u^{(i)} = 0 \vee q_u^{(i)} = 0), \forall u \in S_{RT} \cup S_{NRT}\}$. Since $C0Q0$ depends on the class of $u_{Tx}^{(i-1)}$, we introduce the probabilities $p_{BE_1} = p_{C0Q0 | u_{Tx}^{(i-1)} \in S_{RT}}$, $p_{BE_2} = p_{C0Q0 | u_{Tx}^{(i-1)} \in S_{NRT}}$ and $p_{BE_3} = p_{C0Q0 | u_{Tx}^{(i-1)} \in S_{BE}}$. From Appendix E of [30], we obtain (32) and (33), shown in the next page, where $w_7 = U_{RT} p_{TX1_{RT}} p_{BE_1} + U_{NRT} p_{TX1_{NRT}} p_{BE_2} + (1 - U_{RT} p_{TX1_{RT}} - U_{NRT} p_{TX1_{NRT}} - (U_{BE} - 1) p_{TX1_{BE}}) p_{BE_3}$ and $w_8 = (U_{BE} - 1) p_{TX1_{BE}} p_{BE_3}$. In (33), we have also included the case $u_{Tx}^{(i-1)} = 0$.

5) *Computation of the Probabilities of Conditional Events:* The analytical expressions of the transition probabilities in (20)-(21) and (28)-(33) depend on conditional probabilities, such as $p_{LR|TX0}$ and $p_{L\Phi LB|TX1}$, which turn out to depend on the stationary state probabilities conditioned on the events $TX0_u$ and $TX1_u$. Consequently, instead of the stationary state probability vectors $\{\tilde{\pi}_{S_g, n-1}\}$, $g = \{1, 2, 3\}$, available at iteration $n-1$, we have to consider the two subvectors of $\tilde{\pi}_{S_g, n-1}$ that correspond to the stationary state probability conditioned on $T^{(i-1)} = 0$ and on $T^{(i-1)} = 1$, denoted with $\tilde{\pi}_{S_g, n-1}^{(no-Tx)}$ and $\tilde{\pi}_{S_g, n-1}^{(Tx)}$, respectively. By standard conditional probability rules, for $u \in S_g$, we obtain

$$\begin{aligned} \tilde{\pi}_{(c,q,r), S_g, n-1}^{(no-Tx)} &= \frac{\tilde{\pi}_{(c,q,r,0), S_g, n-1}}{p_{TX0_u}}, \\ \tilde{\pi}_{(c,q,r), S_g, n-1}^{(Tx)} &= \frac{\tilde{\pi}_{(c,q,r,1), S_g, n-1}}{p_{TX1_u}}, \end{aligned} \quad (34)$$

where $\tilde{\pi}_{(c,q,r, T^{(i-1)}), S_g, n-1}$, $\tilde{\pi}_{(c,q,r), S_g, n-1}^{(no-Tx)}$, and $\tilde{\pi}_{(c,q,r), S_g, n-1}^{(Tx)}$ are the elements of $\tilde{\pi}_{S_g, n-1}$, $\tilde{\pi}_{S_g, n-1}^{(no-Tx)}$, and $\tilde{\pi}_{S_g, n-1}^{(Tx)}$, respectively, corresponding to $c_u = c$, $q_u = q$ and $r_u = r$. Note that $p_{TX1_u} = \sum_{c=0}^N \sum_{q=0}^{B_u} \sum_{r=0}^{R_u} \tilde{\pi}_{(c,q,r,1), S_g, n-1}$, while $p_{TX0_u} = 1 - p_{TX1_u}$. In addition, the transition probabilities are conditioned on the transmission at the $(i-1)$ th CTI, but are calculated at the i th CTI. Consequently, we include the time evolution from $i-1$ to i by multiplying the conditional state probability vectors by the pre-computed matrices $\mathbf{P}_{S_g, n-1}^{(no-Tx)}$ and $\mathbf{P}_{S_g, n-1}^{(Tx)}$, respectively, as expressed by

$$\hat{\pi}_{S_g, n-1}^{(no-Tx)} = \tilde{\pi}_{S_g, n-1}^{(no-Tx)} \mathbf{P}_{S_g}^{(no-Tx)},$$

$$p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} = w_0 \sum_{u=0}^{U_{RT}-1} \binom{U_{RT}-1}{u} \frac{1}{u+1} p_{ERLB|TX0}^u (p_{L\Phi LB|TX0} + p_{LR LB|TX0})^{U_{RT}-1-u} \quad (28)$$

$$p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} = w_1 p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} + w_2 \sum_{u=0}^{U_{RT}-2} \binom{U_{RT}-2}{u} p_{ERLB|TX0}^u \times (p_{L\Phi LB|TX0} + p_{LR LB|TX0})^{U_{RT}-2-u} \left(\frac{p_{L\Phi LB|TX1} + p_{LR LB|TX1}}{u+1} + \frac{p_{ERLB|TX1}}{u+2} \right) \quad (29)$$

$$p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} = w_3 \sum_{u=0}^{U_{NRT}-1} \binom{U_{NRT}-1}{u} \frac{1}{u+1} p_{E\Phi R|TX0}^u (p_{L\Phi|TX0} + p_{LR\Phi|TX0})^{U_{NRT}-1-u} \quad (30)$$

$$p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} = w_4 p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} + w_5 \sum_{u=0}^{U_{NRT}-1} \binom{U_{NRT}-1}{u} \frac{p_{E\Phi R|TX0}^u (p_{L\Phi|TX0} + p_{LR\Phi|TX0})^{U_{NRT}-1-u}}{u+1} + w_6 \sum_{u=0}^{U_{NRT}-2} \binom{U_{NRT}-2}{u} p_{E\Phi R|TX0}^u (p_{L\Phi|TX0} + p_{LR\Phi|TX0})^{U_{NRT}-2-u} \left(\frac{p_{L\Phi|TX1} + p_{LR\Phi|TX1}}{u+1} + \frac{p_{E\Phi R|TX1}}{u+2} \right) \quad (31)$$

$$p_{T_j^{(i-1)}=1 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} = p_{BE3} \sum_{u=0}^{U_{BE}-1} \binom{U_{BE}-1}{u} \frac{1}{u+1} p_{E\Phi R|TX0}^u (p_{L\Phi|TX0} + p_{LR\Phi|TX0})^{U_{BE}-1-u} \quad (32)$$

$$p_{T_j^{(i-1)}=0 \rightarrow T_j^{(i)}=1 | \psi_j^{(i)}} = w_7 \sum_{u=0}^{U_{BE}-1} \binom{U_{BE}-1}{u} \frac{1}{u+1} p_{E\Phi R|TX0}^u (p_{L\Phi|TX0} + p_{LR\Phi|TX0})^{U_{BE}-1-u} + w_8 \sum_{u=0}^{U_{BE}-2} \binom{U_{BE}-2}{u} p_{ER|TX0}^u (p_{L\Phi|TX0} + p_{LR|TX0})^{U_{BE}-2-u} \left(\frac{p_{L\Phi|TX1} + p_{LR|TX1}}{u+1} + \frac{p_{ER|TX1}}{u+2} \right) \quad (33)$$

$$\hat{\boldsymbol{\pi}}_{S_g, n-1}^{(Tx)} = \hat{\boldsymbol{\pi}}_{S_g, n-1}^{(Tx)} \mathbf{P}_{S_g}^{(Tx)}. \quad (35)$$

From (35), we derive the probabilities used in (20)-(33), for Tx and no-Tx users. These probabilities correspond to the sum of the elements of subvectors of $\hat{\boldsymbol{\pi}}_{S_g, n-1}^{(Tx)}$ or $\hat{\boldsymbol{\pi}}_{S_g, n-1}^{(no-Tx)}$, depending on the conditioning on $TX1_u$ or $TX0_u$, respectively, for $u \in S_g$. By exploiting Appendix F of [30], we derive

$$p_{LRu|TX0_u} = \sum_{c=1}^N \sum_{r=0}^{\tilde{R}} \hat{\pi}_{(c, B_u, r), S_1, n-1}^{(no-Tx)}, \quad (36)$$

$$p_{LBu|TX0_u} = \sum_{c=0}^N \sum_{q=0}^{B_u-1} \sum_{r=0}^{R_u} \hat{\pi}_{(c, q, r), S_g, n-1}^{(no-Tx)} + \sum_{r=0}^{R_u} \hat{\pi}_{(0, B_u, r), S_g, n-1}^{(no-Tx)}, \quad (37)$$

$$p_{LK_u|TX0_u} = \sum_{c=1}^{c_j-1} \hat{\pi}_{(c, B_u, r_j R_u / R_j), S_1, n-1}^{(no-Tx)}, \quad (38)$$

$$p_{EK_u|TX0_u} = \hat{\pi}_{(c_j, B_u, r_j R_u / R_j), S_1, n-1}^{(no-Tx)}, \quad (39)$$

$$p_{L\Phi LB_u|TX0_u} = \sum_{c=0}^N \sum_{q=0}^{\tilde{m}} \sum_{r=0}^{R_u} \hat{\pi}_{(c, q, r), S_g, n-1}^{(no-Tx)}, \quad (40)$$

$$p_{L\Phi_u|TX0_u} = \sum_{c=0}^N \sum_{q=0}^{[\phi_j B_j / K_c]-1} \sum_{r=0}^{R_u} \hat{\pi}_{(c, q, r), S_g, n-1}^{(no-Tx)}, \quad (41)$$

$$p_{LR LB_u|TX0_u} = \sum_{c=0}^N \sum_{r=0}^{\tilde{R}} \hat{\pi}_{(c, \phi_j B_j / K_c, r), S_1, n-1}^{(no-Tx)}, \quad (42)$$

$$p_{L\Phi R_u|TX0_u} = \sum_{c=0}^N \sum_{r=0}^{\tilde{R}} \hat{\pi}_{(c, \phi_j B_j / K_c, r), S_g, n-1}^{(no-Tx)}, \quad (43)$$

$$p_{ER LB_u|TX0_u} = \sum_{c=0}^N \eta_{u, j, c} \hat{\pi}_{(c, \frac{\phi_j B_j}{K_c}, \frac{r_j R_u}{R_j}), S_1, n-1}^{(no-Tx)}, \quad (44)$$

$$p_{E\Phi R_u|TX0_u} = \sum_{c=0}^N \hat{\pi}_{(c, \phi_j B_j / K_c, r_j R_u / R_j), S_g, n-1}^{(no-Tx)}, \quad (45)$$

where, for a user $u \in S_g$, $\hat{\pi}_{(c,q,r),S_g,n-1}^{(no-Tx)}$ is the element of $\hat{\pi}_{S_g,n-1}^{(no-Tx)}$ that corresponds to $c_u = c$, $q_u = q$ and $r_u = r$, as expressed by $\hat{\pi}_{(c,q,r),S_g,n-1}^{(no-Tx)} = [\hat{\pi}_{S_g,n-1}^{(no-Tx)}]_{k_u^{(i)}}$, $k_u^{(i)} = (B_u + 1)(R_u + 1)c_u^{(i)} + (R_u + 1)q_u^{(i)} + r_u^{(i)} + 1$, $\tilde{R} = \lceil r_j R_u / R_j \rceil - 1$, and $\tilde{m} = \min \{ \lceil \phi_j B_j / K_c \rceil - 1, B_u - 1 \}$. In (44), $\eta_{u,j,c} = 1$ when $\phi_j B_j / K_c < B_u$, and $\eta_{u,j,c} = 0$ for $\phi_j B_j / K_c = B_u$. Please observe that, when $r_j R_u / R_j$ is non-integer, the probabilities (38), (39), (44), and (45) are zero. The other probabilities, i.e., those conditioned on $TX1_u$, are the same of (36)-(45), but with $\hat{\pi}_{(c,q,r),S_g,n-1}^{(Tx)}$ instead of $\hat{\pi}_{(c,q,r),S_g,n-1}^{(no-Tx)}$.

B. Step 2: Computation of the Probability Distribution of the Stationary State

The state distribution vectors $\{\tilde{\pi}_{S_g,n}\}$, $g = \{1, 2, 3\}$ are derived from the knowledge of $\tilde{\mathbf{P}}_{S_g,n}$ as follows

$$\tilde{\pi}_{S_g,n} \tilde{\mathbf{P}}_{S_g,n} = \tilde{\pi}_{S_g,n}, \quad \tilde{\pi}_{S_g,n} \mathbf{1} = 1. \quad (46)$$

Specifically, $\tilde{\pi}_{S_g,n}$ is the unique left eigenvector of $\tilde{\mathbf{P}}_{S_g,n}$ that corresponds to the unit eigenvalue, which is also maximum. Hence, the power method [32] can be used to obtain $\tilde{\pi}_{S_g,n}$. The iterative procedure is stopped when $\|1 - \tilde{\pi}_{S_g,n} \circ \tilde{\pi}_{S_g,n-1}\|_\infty \leq \varepsilon$, where \circ denotes element-wise division and ε is a suitable threshold, or when a fixed number N_{IT} of iterations have elapsed. The final vector $\tilde{\pi}_{S_g}$ is then used to calculate the average delay $\tau_{D,j}$, the PLR, and the throughput Σ_j , as in [13], for the user $j \in S_g$. In addition, from $\tilde{\pi}_{S_g}$ we can also extract the probability distribution of the substate (c, T) , used in [8]-[9] to obtain the delay probability distribution.

1) *Convergence*: The convergence of the whole iterative procedure has been investigated and verified by simulations. To speed up convergence, we split the iterative procedure into two subprocedures: the first considers RT and NRT users, while the second, which concerns BE users, exploits the results of the first one. Indeed, the performance of RT and NRT users is independent from the parameters of BE users.

C. Remarks

1) *Computational Complexity*: In order to assess the computational complexity of the analytical procedure, we estimate the number of multiplications required in Step 1 and Step 2. Actually, the total complexity is dominated by the class with the largest product $L_g = (N + 1)(B_{S_g} + 1)(R_{S_g} + 1)$. For each iteration, the number of multiplications in Step 1 is roughly estimated as $N_{M1} = 2\kappa_T U L_g$, where κ_T represents the number of terms (similar to (16)) that have to be computed. For our scheduling algorithm, $\kappa_T \leq 14$. On the other hand, the number of multiplications in Step 2 is approximately $N_{M2} = N_{PM} \kappa_S L_g$ per iteration, where N_{PM} is the number of iterations of the power method, and κ_S is the sparsity index of $\tilde{\mathbf{P}}_{S_g,n}$, herein defined as the average number of nonzero elements per row. In most cases, κ_S is very low with respect to L_g . Therefore, the total number of multiplications can be estimated as

$$N_M \approx (N_{M1} + N_{M2}) N_{IT} \approx (2\kappa_T U + N_{PM} \kappa_S) L_g N_{IT},$$

where N_{IT} is the number of iterations of our two-step iterative procedure. Hence, the complexity is linear with the number of users U . Anyway, when U is very large, accurate low-complexity approximations are possible. For instance, we can approximate (16) by exploiting the De Moivre-Laplace Theorem, which leads to

$$\sum_{u=0}^{U-1} \binom{U-1}{u} \frac{1}{u+1} x^u y^{U-u-1} \approx \frac{1}{\sqrt{2\pi U^3 x^2 y (1-y)}} \int_1^U \left(\frac{x}{1-y} \right)^z e^{-\frac{(z-U(1-y))^2}{2Uy(1-y)}} dz,$$

which is accurate when $x + y$ is close to one. Since this integral admits a closed-form solution, significant complexity reduction can be achieved. Alternatively, when $x \ll 1$, only the first terms are non-negligible, and therefore the summation can be safely truncated.

2) *Single-Class Case: RT Users Only*: On the basis of the previous results, it is interesting to examine the case of a single service class, e.g., the RT users, since this class is more demanding in terms of QoS performance. Clearly, only one matrix $\tilde{\mathbf{P}}_{S_1,n}$ and one vector $\tilde{\pi}_{S_1,n}$ must be updated at each iteration, and only RT users must be considered in Step 1 and Step 2. The transition probabilities are still expressed by (28) and (29), with the additional constraint $U_{NRT} = 0$. The obtained results correspond to those of [33], which actually neglects the very low probability $\Pr\{u_{Tx}^{(i-1)} = 0\} = 1 - (U_{RT} - 1)p_{TX1_{RT}}$.

3) *Applicability to Other Scheduling Policies*: We remark that, by the same approach, we can model other memoryless scheduling policies, by computing equations similar to (17)-(46). Indeed, although tedious, the computation of the state probability is straightforward, provided that suitable events like (17)-(19), (22)-(27) are defined. Therefore, our approach can also be used for the QoS analysis of other scheduling algorithms. In general, we expect that the accuracy level is higher for those algorithms that try to minimize the interdependencies of the user states, in such a way that the approximation (6) works well.

VI. SIMULATION RESULTS

We assume that the total bandwidth of each user link is $W = 1.08$ MHz, with packet length $N_P = 1080$ bits, and CTI of duration $T_f = 2$ ms. We use the channel model of [13] (Rayleigh fading with average SNR = 15 dB and Doppler frequency $f_d = 10$ Hz), assuming six states ($N = 5$) and choosing the TMs to guarantee a packet-error rate PER = 0.05 for each TM, as in [13]. We consider Poisson arrival processes, with packet arrival probability $p_{a|c_j^{(i)}} = (\lambda_{S_g} l_n)^a e^{-\lambda_{S_g} l_n} / a!$, $a > 0$, $j \in S_g$, $g = \{1, 2, 3\}$, where λ_{S_g} is the mean arrival rate, and a is the number of packets arriving in the i th CTI to a certain user $j \in S_g$ during a single slot of duration $l_n = T_f / K_{c_j^{(i)}}$. To quantify the throughput, we introduce the *ideal total rate* (ITR), defined as $K_{MAX}(U) = \sum_{n=0}^N p_{\{nMAX_U=n\}} K_n$, where $nMAX_U = \arg \max_{u \in S_{RT} \cup S_{NRT} \cup S_{BE}} \{c_u\}$, and $p_{\{nMAX_U=n\}}$ represents the probability of the user with the best instantaneous SNR to be in the TM n state. The value of $p_{\{nMAX_U=n\}}$ is

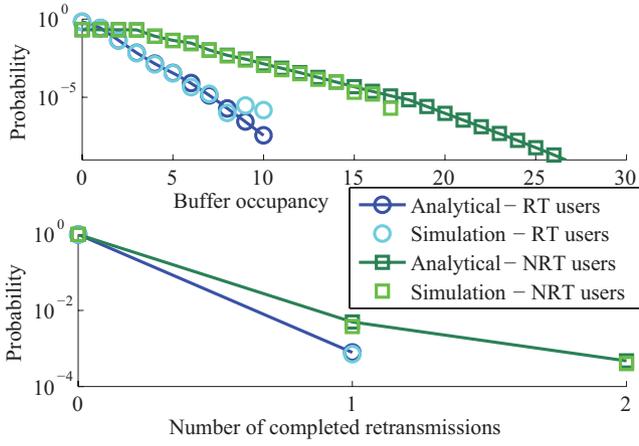


Fig. 3. Buffer occupancy and retransmission distributions for $U_{RT} = 20$, $U_{NRT} = 10$ and $U_{BE} = 20$.

obtained as the stationary state probability associated with the best-user-channel matrix $\mathbf{P}_C^{(MAX)}$ (defined similarly to \mathbf{P}_C , and derived in Appendix G of [30]). Basically, the ITR is the maximum aggregate rate, achieved by scheduling the user with the maximum rate in each CTI. In the SU case, the ITR is the average rate for infinite buffer length and no delay constraints, expressed by $K_{MAX}(1) = \sum_{n=0}^N p_{\{c_1=n\}} K_n \approx 4.9$ packets per CTI (2.65 Mb/sec), while in the MU case the ITR increases thanks to the MU diversity, up to $K_{MAX}(\infty) = K_N = 9$ packets per CTI (4.87 Mb/sec). The arrival rate $\lambda_j = \lambda_{S_g}$ is chosen as $\lambda_{S_g} = \alpha_g K_{MAX}(U) / (U_g T_f)$, where $\alpha_g \in [0, 1)$ quantifies the traffic intensity for each service class g with respect to the ITR. Clearly, by defining $A = \sum_{g=1}^3 \alpha_g$ the sum of the three intensities, $A = 1$ represents an unreachable limit for any scheduling policy in a practical system with finite buffer length. The chosen buffer lengths are $B_{RT} = 10$, $B_{NRT} = 30$, and $B_{BE} = 25$. Although the buffer lengths of NRT and BE users are quite small with respect to practical applications, the probability of higher buffer occupancies is rather small, at least for practical traffic loads such as those herein considered. Hence, short buffers do not affect the model accuracy, and reduce complexity by avoiding unmanageable matrices $\tilde{\mathbf{P}}_{S_g}$. The maximum numbers of retransmissions are $R_{RT} = 1$, $R_{NRT} = 2$, and $R_{BE} = 1$. As a consequence, the square matrices $\tilde{\mathbf{P}}_{S_1}$, $\tilde{\mathbf{P}}_{S_2}$, and $\tilde{\mathbf{P}}_{S_3}$ have sizes 264, 1116, and 624, respectively. In the iterative procedure, we fix $N_{IT} = 100$ and $\varepsilon = 10^{-3}$.

Fig. 3 illustrates the probability distribution of the buffer occupancy $q_j^{(i)}$ and of the number of retransmissions $r_j^{(i)}$, for RT and NRT users, when $U_{RT} = 20$, $U_{NRT} = 10$, and $U_{BE} = 20$. The traffic loads $\alpha_1 = 0.167$, $\alpha_2 = 0.222$, and $\alpha_3 = 0.167$, corresponding to $A = 0.55$, cause a very low transmission probability for BE users, making their QoS performance not significant. Indeed, the probability that no RT (and NRT) users have a packet to transmit is negligible, such that BE users are not able to transmit their packets and their buffers saturate. Moreover, Fig. 3 shows a very good agreement between simulated and analytical probabilities for both buffer occupancies and retransmissions. Notably, both RT and NRT users have very small probabilities of high buffer occupancies. For RT users, this reduces the average delay,

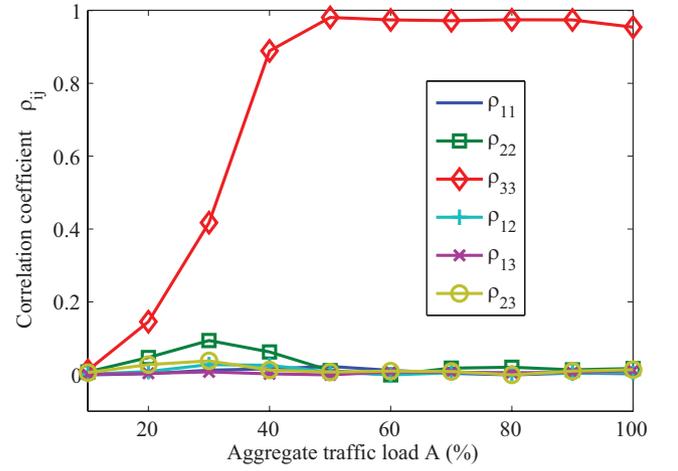


Fig. 4. Correlation coefficients among buffer occupancies of users of the classes i and j .

which is proportional to the average buffer occupancy, while for NRT users, this limits the PLR due to the low probability of exceeding R_{NRT} retransmissions. Besides, Fig. 3 justifies the choice of a short buffer also for NRT users.

In Fig. 4, we have investigated the validity of the independence assumption (6) with ad hoc simulations on the probability distribution of buffer occupancy, which is critical in the model derivation and in the QoS performance assessment. Our aim is to identify the cases when the independence approximation is not accurate. To this end, we monitor when the buffer occupancies of different users are correlated. We define the average crosscorrelation coefficient between the buffer occupancies of different users, as

$$\rho_{gg} = \sum_{i \in S_g} \sum_{j \in S_g, j \neq i} \frac{E\{q_i q_j\} - \mu_g^2}{\sigma_g^2 U_g (U_g - 1)},$$

$$\rho_{gk} = \sum_{i \in S_g} \sum_{j \in S_k, k \neq g} \frac{E\{(q_i - \mu_g)(q_j - \mu_k)\}}{\sigma_g \sigma_k U_g U_k}, \quad (47)$$

where $\mu_g = E\{q_i\}$, $\sigma_g = \sqrt{E\{q_i^2\} - \mu_g^2}$, $i \in S_g$. Fig. 4 displays the average crosscorrelation coefficient ρ_{gk} as a function of the aggregate traffic load A (α_g is assumed equal for all the three classes). It is clear that $|\rho_{33}|$ is significantly greater than 0.1, while in the other cases $|\rho_{gk}|$ is always lower than 0.1. Therefore, the buffer occupancies of BE users are correlated and hence dependent: this happens because, for high traffic loads, the buffers of BE users saturate and hence their occupancies are correlated towards the unit value. On the contrary, in the other cases, the buffer occupancies are uncorrelated. Although uncorrelatedness does not imply independence, the good accuracy obtained for RT and NRT users in the other simulations suggest that the approximation error due to the independence assumption (6) is negligible, even for high traffic loads.

Fig. 5 exhibits the average delay and the PLR for RT and NRT users as a function of U_{RT} when $U_{NRT} = 10$ and $U_{BE} = 20$. The delay of RT users, which increases linearly with U_{RT} , is much lower than for NRT users and than the maximum delay $\delta = 200$ ms defined for WiMAX [18]. The PLR of BE users is not significant, being close to one. The

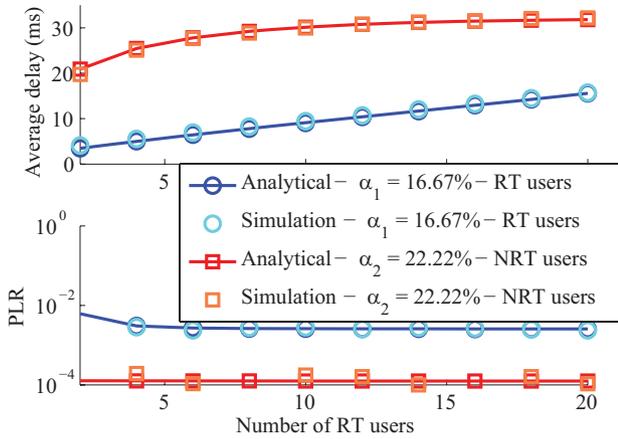


Fig. 5. Average packet delay and PLR for $U_{NRT} = 10$ and $U_{BE} = 20$.

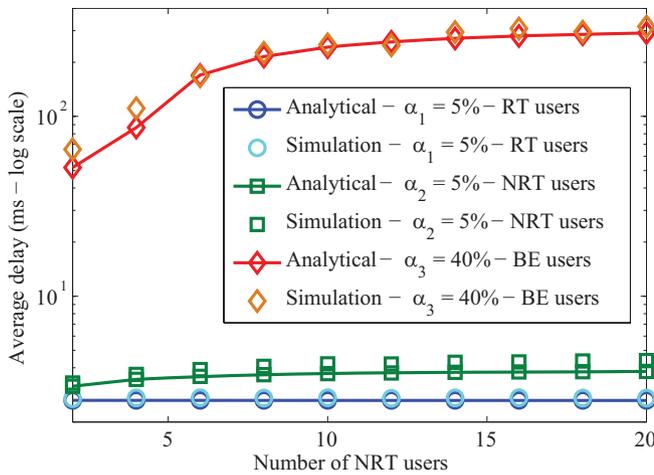


Fig. 6. Average packet delay for $U_{RT} = 20$ and $U_{BE} = 20$.

model accuracy is confirmed even when the PLR of NRT users is very low. Anyway, for both RT and NRT users, the PLR is lower than the maximum PLR defined for WiMAX (0.01 for RT and 0.001 for NRT) [18]. For these moderately high traffic loads, the accuracy is very good. However, for very high traffic loads, when many buffers are close to the saturation, the independence assumption is less accurate. Nevertheless, higher accuracies are generally obtained for higher numbers of users. Noteworthy, the complexity of the proposed analytical framework does not increase with the number of users, when De Moivre-Laplace calculations are involved. Therefore, differently from [8], [9] and [25], this work is best suited to analytically handle many users.

In Figs. 3-5, due to the high traffic loads for RT and NRT users, the probability that a BE user is scheduled is very low, so that the QoS parameters for BE users are not significant. In Fig. 6, we focus on BE users, and hence we choose $\alpha_1 = \alpha_2 = 0.05$ and $\alpha_3 = 0.4$, assuming $U_{RT} = 20$, $U_{BE} = 20$, and variable U_{NRT} . This way, BE users are scheduled quite often. Fig. 6 plots (in a logarithmic scale) the average packet delay for RT, NRT, and BE users. Clearly, in Fig. 6 the average delay for BE users is much higher than for RT and NRT users, and is quite sensitive to U_{NRT} . Fig. 6 also demonstrates a

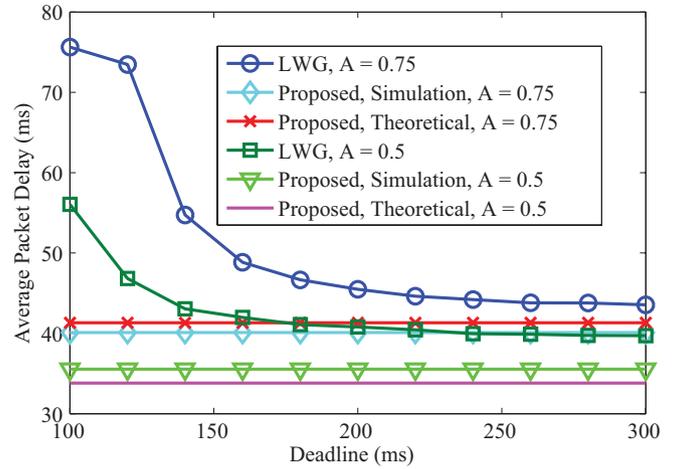


Fig. 7. Average packet delay for different values of the time deadline in LWG [25], for $U_{RT} = 20$.

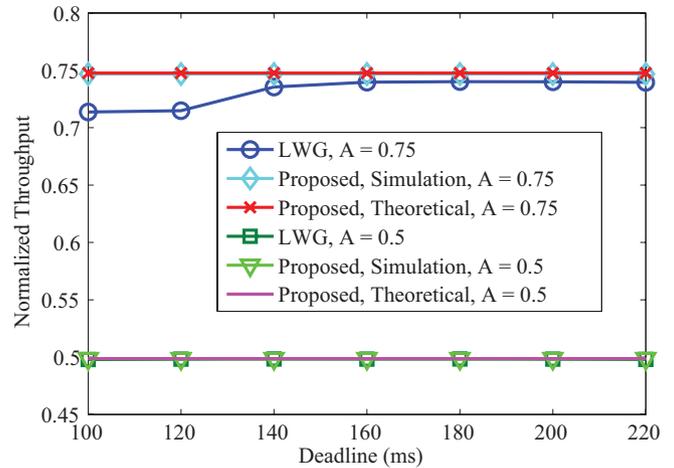


Fig. 8. Normalized throughput for different values of the time deadline in LWG [25], for $U_{RT} = 40$.

slight worse model accuracy with respect to the previous cases. Anyway, the presented model is still useful to predict the QoS performance of heterogeneous users.

Fig. 7 and Fig. 8 compare our proposed scheduler with [25], denoted with LWG. To simplify the comparison, we focus on the single-class case, where all users are RT. Fig. 7 and Fig. 8 plot the average delay and the throughput experienced by $U = 40$ RT users, for $A = 0.75$ and $A = 0.5$, as a function of the *time deadline*, i.e., the temporal parameter of [25] discussed in Section III. We observe that a longer deadline reduces the average delay of [25], which however exceeds the average delay obtained with our algorithm, even for a deadline greater than the maximum tolerable delay for voice systems ($\delta = 200$ ms). Moreover, our algorithm yields an increased throughput for $A = 0.75$, because we favor users with full buffer and therefore we reduce the packet losses.

Finally, we compare our proposed scheduler with the LCQ algorithm of [21]. Similarly to the previous case, we focus on a single class with RT users. The LCQ scheduler assigns transmission to the user with higher buffer occupancy, among the users whose channel state is not OFF, and is delay-optimal in scenarios with two TMs (ON/OFF) [21]. In order

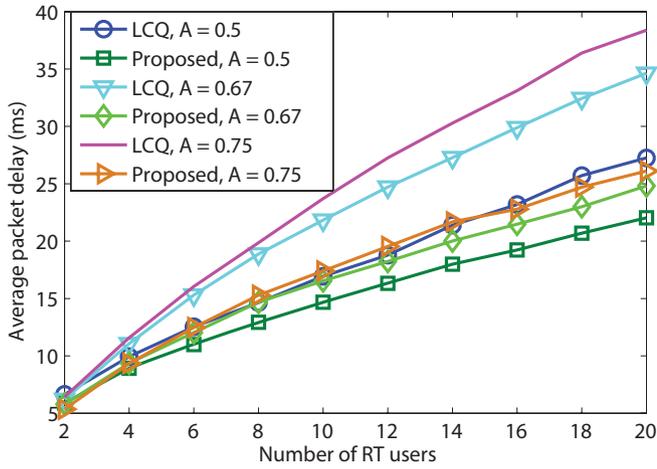


Fig. 9. Average packet delay comparison with LCQ.

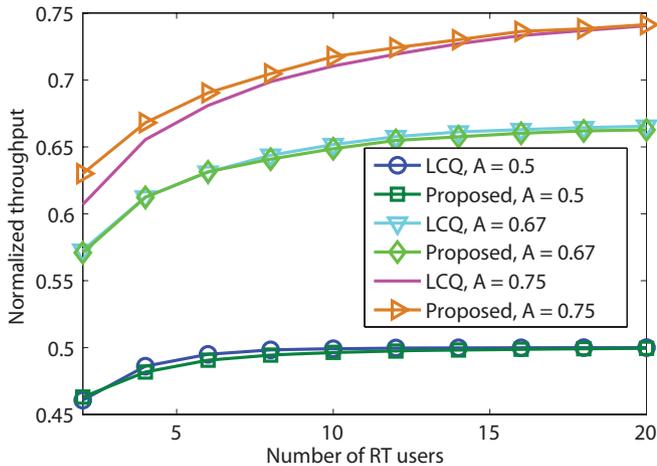


Fig. 10. Normalized throughput comparison with LCQ.

to use the LCQ algorithm in our multiple-TM scenario, we associate the OFF mode with the TM0, and the ON mode to the TM defined by the actual channel state. Therefore, we allow the LCQ algorithm to transmit multiple packets for each CTI. Fig. 9 and Fig. 10 plot the average delay and the throughput, respectively, for $A = 0.75$, $A = 0.67$, and $A = 0.5$, as a function of the number of users. Fig. 9 shows that the proposed scheduling policy produces a lower average packet delay with respect to LCQ, for all values of the traffic load A . Indeed, differently from LCQ, the proposed scheduler takes its decision depending on the value of the channel rate K_c . In practice, using our scheduler policy, usually more packets are transmitted, and this helps in reducing the buffer queues. As a result, the average delay is reduced too. Fig. 10 shows that the two schedulers produce similar throughput performances: LCQ achieves a slightly higher throughput for low traffic loads, while the proposed scheduling policy has better performance for high traffic loads. Actually, by changing the association of the ON/OFF states with our TMs, we may reduce also the average packet delay of LCQ, with a penalty in terms of throughput.

VII. CONCLUSIONS

We have presented a useful framework to characterize the theoretical QoS performance of memoryless scheduling policies for multi-user, multi-class wireless systems that combine AMC and ARQ in a cross-layer fashion. The proposed framework is able to accurately assess the PLR, the average delay, and the throughput, for a heuristic scheduling algorithm that produces a good throughput performance and a reduced average delay. The validity of the introduced approximations has been assessed by simulations, which have confirmed the good accuracy of the QoS parameters estimation, with a significantly reduced computational time with respect to extensive simulations. The proposed analytical framework can be adapted to other memoryless scheduling algorithms based on channel quality states, buffer occupancy states, and retransmission number states. As a possible future work, our framework, which is based on time-slot contention, could be extended to include frequency-slot contention policies.

REFERENCES

- [1] S. Shakkottai, T. S. Rappaport, and P. C. Karlsson, "Cross-layer design for wireless networks," *IEEE Commun. Mag.*, vol. 41, pp. 74-80, Oct. 2003.
- [2] V. Srivastava and M. Motani, "Cross-layer design: a survey and the road ahead," *IEEE Commun. Mag.*, vol. 43, pp. 112-119, Dec. 2005.
- [3] X. Lin, N. B. Shroff, and R. Srikant, "A tutorial on cross-layer optimization in wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 24, pp. 1452-1463, Aug. 2006.
- [4] V. K. N. Lau and Y.-K. R. Kwok, *Channel-Adaptive Technologies and Cross-Layer Designs for Wireless Systems with Multiple Antennas: Theory and Applications*. Hoboken, NJ: Wiley-Interscience, 2006.
- [5] S. Catreux, V. Erceg, D. Gesbert, and R. W. Heath Jr, "Adaptive modulation and MIMO coding for broadband wireless data networks," *IEEE Commun. Mag.*, vol. 40, pp. 108-115, June 2002.
- [6] L. Hanzo, J. S. Blogh, and S. Ni, *3G, HSPA and FDD versus TDD Networking: Smart Antennas and Adaptive Modulation*, 2nd edition. Chichester, UK: Wiley-IEEE Press, 2008.
- [7] S. B. Wicker, *Error Control Systems for Digital Communication and Storage*. Englewood Cliffs, NJ: Prentice-Hall, 1995.
- [8] L. B. Le, E. Hossain, and A. S. Alfa, "Service differentiation in multirate wireless networks with weighted round-robin scheduling and ARQ-based error control," *IEEE Trans. Commun.*, vol. 54, pp. 208-215, Feb. 2006.
- [9] L. B. Le, E. Hossain, and A. S. Alfa, "Delay statistics and throughput performance for multi-rate wireless networks under multiuser diversity," *IEEE Trans. Wireless Commun.*, vol. 5, pp. 3234-3243, Nov. 2006.
- [10] Q. Liu, S. Zhou, and G. B. Giannakis, "Cross-layer combining of adaptive modulation and coding with truncated ARQ over wireless links," *IEEE Trans. Wireless Commun.*, vol. 3, pp. 1746-1755, Sep. 2004.
- [11] Q. Liu, S. Zhou, and G. B. Giannakis, "Queuing with adaptive modulation and coding over wireless links: cross-layer analysis and design," *IEEE Trans. Wireless Commun.*, vol. 4, pp. 1142-1153, May 2005.
- [12] X. Wang and J. K. Tugnait, "Joint design of bandwidth distribution, truncated ARQ protocol, and adaptive modulation and coding scheme for multiple user delay sensitive traffic," in *Proc. 38th Conf. Inf. Sciences Syst.*, Mar. 2004.
- [13] X. Wang, Q. Liu, and G. B. Giannakis, "Analyzing and optimizing adaptive modulation coding jointly with ARQ for QoS-guaranteed traffic," *IEEE Trans. Veh. Technol.*, vol. 56, pp. 710-720, Mar. 2007.
- [14] L. Kleinrock, *Queueing Systems, Volume I: Theory*. New York: Wiley Interscience, 1975.
- [15] M. Zorzi and R. R. Rao, "Throughput of selective-repeat ARQ with time diversity in Markov channels with unreliable feedback," *Wireless Netw.*, vol. 2, pp. 63-75, Mar. 1996.
- [16] M. Zorzi, R. R. Rao, and L. B. Milstein, "ARQ error control for fading mobile radio channels," *IEEE Trans. Veh. Technol.*, vol. 46, pp. 445-455, May 1997.
- [17] M. Zorzi, "Some results on error control for burst-error channels under delay constraints," *IEEE Trans. Veh. Technol.*, vol. 50, pp. 12-24, Jan. 2001.

- [18] IEEE, IEEE Standard for Local and Metropolitan Area Networks – Part 16: Air Interface for Fixed and Mobile Broadband Wireless Access Systems, Feb. 2006.
- [19] P. Viswanath, D. N. C. Tse, and R. Laroia, “Opportunistic beamforming using dumb antennas,” *IEEE Trans. Inf. Theory*, vol. 48, pp. 1277-1294, June 2002.
- [20] J. Tang and X. Zhang, “Quality-of-service driven power and rate adaptation over wireless links,” *IEEE Trans. Wireless Commun.*, vol. 6, pp. 3058-3068, Aug. 2007.
- [21] L. Tassiulas and A. Ephremides, “Dynamic server allocation to parallel queues with randomly varying connectivity,” *IEEE Trans. Inf. Theory*, vol. 39, pp. 466-478, Mar. 1993.
- [22] M. J. Neely, “Order optimal delay for opportunistic scheduling in multi-user wireless uplinks and downlinks,” *IEEE/ACM Trans. Networking*, vol. 16, pp. 1188-1199, Oct. 2008.
- [23] D. S. W. Hui, V. K. N. Lau, and W. H. Lam, “Cross-layer design for OFDMA wireless systems with heterogeneous delay requirements,” *IEEE Trans. Wireless Commun.*, vol. 6, pp. 2872-2880, Aug. 2007.
- [24] V. K. N. Lau and Y. Chen, “Delay-optimal power and precoder adaptation for multi-stream MIMO systems,” *IEEE Trans. Wireless Commun.*, vol. 8, pp. 3104-3111, June 2009.
- [25] Q. Liu, X. Wang, and G. B. Giannakis, “A cross-layer scheduling algorithm with QoS support in wireless networks,” *IEEE Trans. Veh. Technol.*, vol. 55, pp. 839-847, May 2006.
- [26] M. Nakagami, “The m -distribution—a general formula of intensity distribution of rapid fading,” in *Statistical Methods in Radio Wave Propagation*. Oxford, UK: Pergamon, 1960, pp. 3-36.
- [27] G. L. Stüber, *Principles of Mobile Communications*, 2nd edition. Norwell, MA: Kluwer, 2001.
- [28] C. Comaniciu and H. V. Poor, “Jointly optimal power and admission control for delay sensitive traffic in CDMA networks with LMMSE receivers,” *IEEE Trans. Signal Process.*, vol. 51, pp. 2031-2042, Aug. 2003.
- [29] M. D. Yacoub, J. E. Vargas Bautista, and L. Guerra de Rezende Guedes, “On higher order statistics of the Nakagami- m distribution,” *IEEE Trans. Veh. Technol.*, vol. 48, pp. 790-794, May 1999.
- [30] M. Poggioni, L. Rugini, and P. Banelli, “QoS performance analysis of a heuristic scheduling algorithm for heterogeneous users employing AMC and ARQ,” Technical Report RT-002-08, Dept. Electron. Inf. Eng., University of Perugia, June 2008 [Online]. Available: <http://www.diei.unipg.it/rt/RT-002-08-Poggioni-Rugini-Banelli.pdf>
- [31] J. R. Norris, *Markov Chains*. Cambridge, UK: Cambridge Univ. Press, 1997.
- [32] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd edition. Baltimore, MD: Johns Hopkins Univ. Press, 1996.
- [33] M. Poggioni, L. Rugini, and P. Banelli, “Analyzing performance of multi-user scheduling jointly with AMC and ARQ,” in *Proc. IEEE Global Telecommun. Conf. (GLOBECOM)*, Nov. 2007, pp. 3483-3488.



Mario Poggioni (S'07-M'10) was born in Perugia, Italy, in 1979. He received the Laurea degree (magna cum laude) in electronics engineering in 2005 and the Ph.D. degree in telecommunications in 2009, from the University of Perugia. He is currently a Research Engineer with ART Srl. His research interests lie in the areas of signal processing for multicarrier communications, fast fading channels, broadcasting and cross-layer designs.



Luca Rugini (S'01-M'05) was born in Perugia, Italy, in 1975. He received the Laurea degree in electronic engineering and the Ph.D. degree in telecommunications from the University of Perugia, in 2000 and 2003, respectively. From February to July 2007, he visited Delft University of Technology, The Netherlands. He is currently an Assistant Professor with the Department of Electronic and Information Engineering at the University of Perugia. His research interests lie in the area of signal processing for multicarrier and spread-spectrum

communications.



Paolo Banelli (S'90-M'99) received the Laurea degree in electronics engineering and the Ph.D. degree in telecommunications from the University of Perugia, Perugia, Italy, in 1993 and 1998, respectively. In 2005, he was appointed Associate Professor at the Department of Electronic and Information Engineering (DIEI), University of Perugia, where he has been an Assistant Professor since 1998. In 2001, he joined the SpinComm group at the Electrical and Computer Engineering Department, University of Minnesota, Minneapolis, as a Visiting Researcher. His research

interests mainly focus on signal processing for wireless communications, with emphasis on multicarrier transmissions, and more recently on signal processing for biomedical applications, with emphasis on electrocardiography and medical ultrasounds. He has been serving as a reviewer for several technical journals, and as technical program committee member of leading international conferences on signal processing and telecommunications. In 2009, he was a General Co-Chair of the IEEE International Symposium on Signal Processing Advances for Wireless Communications.