

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

Queuing Network Models of Multiservice RANs

This is a pre print version of the following article:

Original Citation:

Availability:

This version is available <http://hdl.handle.net/2318/1965548> since 2024-03-29T23:37:00Z

Published version:

DOI:10.1145/3649307

Terms of use:

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

Queuing Network Models of Multiservice RANs

ANDREA MARIN, Università Cà Foscari di Venezia, Italy

MICHELA MEO, Politecnico di Torino, Italy

MATTEO SERENO, Università di Torino, Italy

MARCO AJMONE MARSAN, Politecnico di Torino, Italy, and IMDEA Networks Institute, Spain

In this paper, we present a new queuing network model for the analysis of a portion of a radio access network (RAN) comprising macro cell base stations (BSs) and small cell BSs offering "streaming" and "elastic" services. Streaming services require a certain data rate for a random time. The required data rates depend on the type of service, e.g., audio and video. Elastic services require the transfer of random data volumes, and their data rate adjusts dynamically based on the capacity not utilized by the streaming services. To derive performance measures for the proposed model we develop a computationally efficient framework that exploits a new product form result for streaming services, relying on a well-known blocking policy, and an approximate product form for elastic services. Insensitivity to the distribution of service requirements holds in the case of negligible end user mobility.

We show the high accuracy of our model in predicting the performance of practical system configurations by conducting a thorough comparison between the model's results and those obtained from a detailed discrete-event simulator. Through this analysis, we uncover significant counter-intuitive behaviors that arise from the competition between streaming services with diverse demands, and our model effectively captures and predicts these behaviors.

Our computationally efficient queuing model is a useful new tool to support design and planning of multiservice RANs whose complex structures result from the coexistence of BSs of different generations in dense areas.

ACM Reference Format:

Andrea Marin, Michela Meo, Matteo Sereno, and Marco Ajmone Marsan. 2024. Queuing Network Models of Multiservice RANs. 1, 1 (February 2024), 22 pages. <https://doi.org/10.1145/nmnnnnn.nnnnnnn>

1 INTRODUCTION

The radio access networks (RANs) that today cover large metropolitan areas are quite complex. For example, the metropolitan area of Milan in Italy is covered by about 5,000 base stations (BSs) belonging to 4 infrastructured mobile network operators (MNOs). Not all BSs are equal: some transmit at high power and have a reach of 1 km or more; some use much lower power and have much shorter reach, down to tens of meters. Each BS provides services in the cell it defines, and the mix of services requested by end users is extremely diverse, including video, audio, images, text, large and small data chunks, etc. According to the 5G jargon, services are grouped in service categories, such as eMBB (enhanced mobile broadband) mMTC (massive machine-type communications) and URLLC (ultra-reliable low-latency communications), implemented over network slices, i.e., virtual networks implemented using (possibly non-disjoint) subsets of resources.

Authors' addresses: Andrea Marin, andrea.marin@unive.it, Università Cà Foscari di Venezia, Italy; Michela Meo, michela.meo@polito.it, Politecnico di Torino, Italy; Matteo Sereno, Università di Torino, Italy, matteo.sereno@unito.it; Marco Ajmone Marsan, Politecnico di Torino, Italy, and IMDEA Networks Institute, Spain.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

53 Using low power values, hence reducing the reach of BSs, has the double advantage of a lower amount of electro-
54 magnetic radiation and of a higher spatial reuse of the spectrum portions allocated to mobile communication services,
55 which leads to an increase in capacity per unit area. The latter aspect, known under the name *spectrum reuse* or *network*
56 *densification* has been the key factor that allowed RANs to dramatically increase their capacity over the years, thus
57 serving more and more users at higher and higher data rates. It is expected that further densification will take place
58 with the diffusion of 5G networks, and even more when 6G will arrive.

59 The effective deployment of over one thousand BSs by a MNO in a metropolitan area calls for sophisticated design
60 and planning tools to forecast both coverage and capacity, so as to provide customers with their chosen services in the
61 most efficient and cost-effective way.

62 The planning problem is further complicated by the fact that several RAN generations typically coexist, at least for
63 some periods of time, so that users can access services through BSs of different generations, with different capacity,
64 different reach, and different technology. The resulting RAN architectures include overlapping cells of different sizes.
65 Today, the most visible aspect of the presence of cells of different types is due to the coexistence of 4G and 5G, with the
66 residual presence of older generations for some types of narrow band Internet of things (IoT) services.

67 The complexity of the design and planning problem is such that only rather small portions of a RAN can be considered
68 at a time (accounting for hundreds or thousands of BSs in one study is not feasible), but even considering few cells at a
69 time raises significant issues.

70 Simulation is typically very cumbersome, and can be effective to verify the indications of analytical models. The
71 solution of analytical models is also far from trivial, because the mix of different service types and the limited capacity
72 of BSs make the analysis difficult. On the one hand, using a direct analysis approach based on continuous-time Markov
73 chains (assuming exponential assumptions are acceptable) is not feasible due to the humongous size of the resulting
74 state space; on the other hand, the standard queuing network modeling approach is also, in general, intractable. Indeed,
75 product forms allowing an efficient numerical analysis of the systems are not known, unless additional assumptions are
76 incorporated into the model.

77 In a recent paper [17], we derived the conditions under which the queuing model of one BS offering multiple services
78 admits a product form solution for the limiting joint probability distribution of the numbers of active services of
79 different types, even allowing, under some conditions, non-exponential distributions for service durations or volumes.
80 The considered service mix consists of streaming and elastic services. Streaming services require a constant data rate
81 for their entire duration, while elastic services can adapt their data rate to the capacity available in the RAN. Streaming
82 and elastic services can be viewed as the basic building blocks of the 5G service categories that we mentioned above.

83 The results of [17] showed that product form is achieved through a state-dependent admission control algorithm for
84 elastic services, and in addition proved that the conditions for the existence of a product form solution, in the case of
85 negligible mobility of end users, also lead to the insensitivity of the result to the distribution of service requirements,
86 thus allowing the model to depart from exponential assumptions.

87 In this paper, we use an approach inspired by [17] to study the queuing network model of a portion of a multiservice
88 RAN comprising several BSs.

89 We prove that the queuing network model admits a product form solution for the limiting joint probability distribution
90 of the numbers of active streaming services of different types at different BSs under a state-dependent routing of service
91 requests for both types of services. In addition, we develop an approximate but accurate product form solution for
92 elastic customers, and we prove that in case of negligible user mobility the results of the product form are insensitive to
93 service time distributions.

105 While a few papers in the literature have discussed analytical models for the case of one BS loaded with traffic
 106 resulting from a mix of streaming and elastic services (see for example [3, 5, 6, 8–10, 14, 16, 21–23]) – but the possibility
 107 of a product form solution and of an insensitivity to service duration was never proved before [17] – to the best of our
 108 knowledge, no paper has yet tackled the modeling of groups of several multiservice BSs.

109 The solution technique that we propose in this paper can scale to portions of RANs comprising some tens of cells
 110 and few thousands users.
 111
 112
 113

114 2 THE RADIO ACCESS NETWORK

115 We consider a portion of a radio access network (RAN) comprising a number $N^{(BS)}$ of base stations (BSs) and a number
 116 $N^{(UE)}$ of end user terminals (called user equipments – UEs). Each BS k , with $k = 1, 2, \dots, N^{(BS)}$, has a user plane data
 117 rate C_k bit/s.
 118
 119

120 UEs roam over the area served by the RAN and in different time periods they may be associated with different BSs
 121 (while at any time being associated with only one BS, except for short handover transients). The association of UEs
 122 with BS k , i.e., the dwell times of UEs at the BS, have durations described by the i.i.d. random variables δ_k .

123 UEs may request service from the BS with which they are currently associated. Services requested by UEs are of two
 124 different types: streaming and elastic.
 125

126 Streaming services are characterized by a fixed data rate and a random duration. Streaming services can belong to
 127 different classes, depending on their data rate requirement and/or duration. Examples of streaming services are the
 128 access to real time videos (e.g., to watch a live sport event or a newscast) or conference calls including audio and video
 129 streams from remote partners. We denote the number of streaming service classes with S . Class i ($i = 1, 2, \dots, S$) is
 130 characterized by a fixed data rate $R_i^{(s)}$ and by service durations described by i.i.d. random variables. The duration of a
 131 generic class i service is denoted by $\tau_i^{(s)}$.
 132
 133

134 Elastic services are characterized by a random volume of data to be transferred and by the possibility to use a variable
 135 data rate (that may drop temporarily to zero), according to the available data rate at the BS. Elastic services can belong
 136 to different classes, depending on the volume of data to be transferred. Examples of elastic services are the download
 137 of web pages or recorded videos, and the transmission of messages possibly including audio, video and images. We
 138 denote with E the number of elastic service classes. Class j ($j = 1, 2, \dots, E$) is characterized by volumes of data to be
 139 transferred that are described by i.i.d. random variables. The volume to be transferred by a generic class j service is
 140 denoted by $\varphi_j^{(e)}$.
 141
 142

143 A BS (say BS k) allocates resources for a service either because one of the UEs associated with BS k issues a new
 144 service request, or because a UE with an active service that was previously associated with a different BS becomes
 145 associated with BS k due to roaming. A streaming service request of class i is accepted by BS k only if the BS can
 146 allocate the required data rate $R_i^{(s)}$ (possibly reducing the data rate used by elastic services); otherwise, the request is
 147 blocked. Elastic service requests use the BS data rate not used by streaming customers (we will refer to such data rate
 148 with the term *residual data rate*). Elastic service requests of class j are always accepted by BS k , due to the fact that they
 149 have no minimum data rate requirement, until the maximum number of admissible simultaneous class j services, $N_{j,k}^{(e)}$,
 150 is reached. All active elastic services evenly share the residual data rate, regardless of their class. BS k can reserve a
 151 portion of its user plane data rate C_k to be used by elastic services, so as to avoid that the data rate available to elastic
 152 services drops to zero. We denote the data rate reserved to elastic services in BS k as $C_k^{(e)}$.
 153
 154
 155
 156

Summing up, streaming services are blocked if at the epoch of their arrival there is not enough free bandwidth to begin their service, while elastic services are dropped if at their arrival epoch, they find the buffer for elastic services saturated. Furthermore, when a request is blocked, it may still be served by an adjacent BS if this has free resources and the connection is physically possible, otherwise it is dropped.

Active services leave a BS (say BS k) for one of three possible reasons: (i) the service reaches completion – in case of streaming services this means that the total amount of time of service (accounting for service times in all cells visited during service) reaches the value sampled for $\tau_i^{(s)}$; in case of elastic services completion means that the total amount of transferred data (accounting for all cells visited during service and the available data rate) reaches the value sampled for $\varphi_j^{(e)}$; (ii) the UE modifies its association to a BS different from k (e.g., because of roaming); (iii) the active service request that reaches BS k is blocked upon arrival.

The changes of association of a UE (say UE u) from one to another of the BSs present in the service area because of roaming are described by the routing probabilities of its services. In this case, when UE u requests to change its association from BS k to a different BS, the service request reaches BS ℓ , with probability $p_{k,\ell}^{(M)}$. Note that this probability is conditional on the UE leaving the BS due to mobility, not to service completion. Obviously, we have

$$\sum_{\ell=1}^{N^{(BS)}} p_{k,\ell}^{(M)} = 1 \quad \forall k \in 1, 2, \dots, N^{(BS)}, k \neq \ell. \quad (1)$$

3 THE QUEUING NETWORK MODEL

The RAN portion under investigation is modeled with a network of queues in which each BS corresponds to a queue. Customers moving over the queuing network correspond to service instances. Hence, customers are said to be either streaming or elastic of the appropriate class, depending on the type of service they correspond to.

In the case of streaming customers of class i , we assume that service times at queue k are modeled by i.i.d. exponentially distributed random variables with rate $\mu_{k,i}^{(s)}$. In the case of elastic customers of class j , we assume that the amounts of data to be transferred are modeled by i.i.d. exponentially distributed random variables with rate $v_{k,j}^{(e)}$. Insensitivity to service data/volume distributions will be discussed later on. For all customers, we assume that dwell times are modeled by i.i.d. exponentially distributed random variables with rate $\mu_{k,H}$.

The service of (non-blocked) streaming customers of class i at queue k follows a multiserver paradigm, since services proceed in parallel, each one for a time duration equal to the minimum between the time to completion and the dwell time, hence for exponentially distributed times with rate $\mu_{k,i}^{(s)} + \mu_{k,H}$.

The service of (non-blocked) elastic customers follows a state dependent processor sharing paradigm, since elastic services proceed in parallel, fairly sharing the BS residual data rate. The service time duration for elastic customers of class j results from the combination of the volume of data to be transmitted, the dwell time and the residual data rate.

One additional infinite-server queue (called *delay station*) models the time between the end of a (streaming or elastic) service instance and the generation of a new service request (of the same type) from the same UE. At this queue, that will be numbered 0 (all other queues inherit the numbering of the corresponding BS), service times are modeled by arbitrary i.i.d. random variables with means $(\mu_{0,i}^{(s)})^{-1}$ and $(\mu_{0,j}^{(e)})^{-1}$ in the case of streaming (or elastic) customers of class i (j).

In case a customer reaches a queue where it cannot be accepted (because the available data rate is not sufficient in the case of a streaming customer, or, in the case of an elastic customer, because the maximum admissible number

has already been reached), we assume it is transferred to queue ℓ with probability $p_{k,\ell}^{(B)} = p_{k,\ell}^{(M)}$ for all k and ℓ , which implies that a blocked service request behaves exactly as a non-blocked request in terms of routing. This assumption is known with the name of *skipping* in the literature, and is instrumental to our modeling approach. The assumption is reasonable in the case of further handovers, since a request that cannot be accepted at a BS tries to move to another one of the neighboring BSs, but these are the same toward which handovers are possible. Similarly, customers that cannot be accepted at a queue are moved to the delay station with the same probability with which service would complete at the queue. This is interpreted as a forced completion, hence a true blocking of the request.

4 STATIONARY ANALYSIS OF THE QUEUING NETWORK

We are interested in computing the limiting joint distribution of the numbers of customers of the different streaming and elastic service classes at all queues. Although the CTMC underlying the model is finite, the very high cardinality of the state space for realistic scenarios makes the solution of the system of global balance equations intractable. For this reason, we investigate the existence of a product form solution allowing the definition of a numerically stable and efficient algorithm for the computation of the performance indices.

More formally, we are interested in computing the joint probabilities

$$P\{\mathcal{N}_0 = \mathbf{N}_0, \mathcal{N}_1 = \mathbf{N}_1, \dots, \mathcal{N}_{N(BS)} = \mathbf{N}_{N(BS)}\} \quad (2)$$

with

$$\mathcal{N}_k = [\mathcal{N}_k^{(s)}, \mathcal{N}_k^{(e)}], \quad \mathcal{N}_k^{(s)} = [\mathcal{N}_{k,1}^{(s)}, \dots, \mathcal{N}_{k,S}^{(s)}], \quad \mathcal{N}_k^{(e)} = [\mathcal{N}_{k,1}^{(e)}, \dots, \mathcal{N}_{k,E}^{(e)}],$$

where the random variables $\mathcal{N}_{k,i}^{(s)}$ and $\mathcal{N}_{k,j}^{(e)}$ indicate the numbers of streaming services of class i and elastic services of class j in progress at BS k , respectively.

In addition,

$$\mathbf{N}_k = [\mathbf{N}_k^{(s)}, \mathbf{N}_k^{(e)}], \quad \mathbf{N}_k^{(s)} = [n_{k,1}^{(s)}, \dots, n_{k,S}^{(s)}], \quad \mathbf{N}_k^{(e)} = [n_{k,1}^{(e)}, \dots, n_{k,E}^{(e)}],$$

where $n_{k,i}^{(s)}$ and $n_{k,j}^{(e)}$ are the values assumed by the random variables $\mathcal{N}_{k,i}^{(s)}$ and $\mathcal{N}_{k,j}^{(e)}$, respectively.

The existence and uniqueness of the limiting distribution is ensured by the ergodicity of the CTMC underlying the queuing network that, in turns, is implied by the observation that the state space is finite and irreducible.

Observe that the marginal distribution for the network of streaming services can be obtained without considering the network of elastic services. More precisely, the stochastic process underlying the network of streaming services is a modulating process (or environment) for the Markov modulated process underlying the network of elastic services.

As a consequence, we propose a solution method that consists of two steps. First, we describe only the network of streaming services, we derive the conditions for its product form solution, and we give a convolution-based algorithm for the efficient computation of the stationary performance indices. In this step, we also study the conditions for the insensitivity of the model to the distribution of the service times. In the second step, we resort to a standard approximation for the analysis of Markov modulated processes. In fact, it is well-known that, in general, exact solutions for such models are intractable [18], thus we develop a computationally efficient method to obtain accurate estimates of the stationary performance indices including dropping probabilities. The efficiency of this method relies on the product form solution of the network of elastic services conditioned to the state of the streaming network.

4.1 Product form solution for the network of streaming services

Let us consider the network of streaming services with the service and blocking policy defined in Section 3. The following theorem shows that there exists a product form solution both among the stations that form the streaming network and within each station among the classes of service.

THEOREM 1 (PRODUCT FORM FOR THE STREAMING NETWORK). *The streaming network with the blocking policy defined in Section 3 has product form stationary distribution, i.e.:*

$$P\{\mathcal{N}_0^{(s)} = \mathbf{N}_0^{(s)}, \dots, \mathcal{N}_{N^{(BS)}}^{(s)} = \mathbf{N}_{N^{(BS)}}^{(s)}\} = \frac{1}{G^{(s)}} \prod_{k=0}^{N^{(BS)}} g_k(\mathbf{N}_k^{(s)}), \quad (3)$$

moreover, we have that:

$$g_k(\mathbf{N}_k^{(s)}) = \prod_{t=1}^S g_{k,t}(n_{k,t}^{(s)}) \quad (4)$$

at every queue k , where $k = 0, \dots, N^{(BS)}$ and $G^{(s)}$ is the normalizing constant. Functions $g_{k,t}(n_{k,t}^{(s)})$ are defined as:

$$\forall k, t \quad g_{k,t}(n_{k,t}^{(s)}) = \frac{1}{n_{k,t}^{(s)}!} \left(\frac{\sigma_{k,t}}{\mu_{k,H} + \mu_{k,t}} \right)^{n_{k,t}^{(s)}}$$

and $\sigma_{k,t}$ is any non-trivial solution of the BCMP-like system of traffic equations [4] on the routing probabilities defined for streaming service t as:

$$p_{k,\ell}^{t*} = \frac{\mu_{k,H}}{\mu_{k,H} + \mu_{k,t}^{(s)}} p_{k,\ell}^{(M)} + \frac{\mu_{k,t}^{(s)}}{\mu_{k,H} + \mu_{k,t}^{(s)}} 1_{\ell=0} \quad k > 0,$$

where 1 is the indicator function, and $p_{0,\ell}^{t*}$ is the probability that a service is initiated at cell ℓ .

The routing matrices $[p_{k,\ell}^{t*}]$ used to solve the model are obtained as a mixture of routing probabilities due to mobility and to service completion. In other words, routing probabilities in the queuing network do not coincide with roaming probabilities in the system. Indeed, since the residence times of all streaming jobs are independent of the states of the stations, and given the exponential assumption, the probability of terminating the service because of mobility is $\mu_{k,H} / (\mu_{k,H} + \mu_{k,t}^{(s)})$. Moreover, the residence time is still exponentially distributed because it turns out to be the minimum of two independent exponentially distributed random variables (mobility and service times).

Proof sketch. In contrast with the monolithic proofs based on the verification of the system of global balance equations of the entire streaming network as in [19], our proof follows the ideas of [2], or alternatively we may resort to the extended quasi-reversibility result of [13]. We consider a single isolated queue k serving streaming services. If the queue can accommodate $n_{k,1}^{(s)}, \dots, n_{k,S}^{(s)}$ customers, then class t customers are served with rate $n_{k,t}^{(s)} \mu_{k,t}^{(s)}$ and leave the station for mobility with rate $n_{k,t}^{(s)} \mu_{k,H}$. This corresponds to a situation in which customers of class t are served with rate $\mu_{k,H} + \mu_{k,t}^{(s)}$ and at the service completion, we use the mobility routing with probability $\mu_{k,H} / (\mu_{k,H} + \mu_{k,t}^{(s)})$ and the service completion routing with its complement. This implies that the process underlying a single isolated queue is reversible (but not the process underlying the joint process!), thus the reversed rate of the service completion events (due either to mobility or actual service completion) are equal to the intensity of the arrival process at station k for each class. At the arrival of a streaming job at a saturated station, we use the skipping policy. According to the the Reversed

Compound Agent Theorem (see, e.g., [2]), this ensures the product form of the model. For what concerns the multi class delay station, its product form properties are well-known from the literature (see, e.g., [4]). \square

Notice that, although we inherit the skipping policy from [2, 19], the proof of the product form for networks of Erlang-B stations is, to the best of our knowledge, new.

COROLLARY 1 (INSENSITIVITY). *If the effect of mobility is negligible (i.e., mobility is very slow with respect to service completion), the model is insensitive to moments of the service time distribution higher than the first one.*

Proof sketch. Notice that the topology of the network imposes that, after completion, service requests return to the delay station. If we neglect mobility, we can represent the service time distribution with a Coxian distribution and observe that the proof of Theorem 1 still holds true. Coxian distributions are dense in the domain of non-negative continuous distributions, and this suffices to conclude the proof. \square

For what concerns the algorithm for the computation of the normalizing constant $G^{(s)}$, the standard convolution algorithm for multiclass queuing networks (see [11, 12] for details) cannot be applied in our model because of the finite capacity of the queues. The relations between the convolution constants and the average stationary performance indices are different from those well-known from the literature because of the finite capacity of some stations, but their derivation is totally analogous [11, 12]. Notice that, the choice of resorting to a convolution algorithm rather than a mean value analysis will be clear in Section 4.3 since we will use the convolution table to sample a subset of states of the network of streaming services.

THEOREM 2 (CONVOLUTION FOR STREAMING SERVICES). *For the product form network of streaming services, the following relation between the normalizing constants holds:*

$$G_{\Omega, \mathbf{m}^{(s)}}^{(s)} = \sum_{\mathbf{n} \in \mathcal{S}_k(\mathbf{m}^{(s)})} g_k(\mathbf{n}) G_{\Omega \setminus \{k\}, \mathbf{m}^{(s)} - \mathbf{n}}^{(s)},$$

where Ω is a non-empty (sub)set of stations and k an arbitrary element, $\mathbf{m}^{(s)} = (m_1^{(s)}, \dots, m_S^{(s)})$ is the population of streaming services, and

$$\mathcal{S}_k(\mathbf{m}^{(s)}) = \{(n_1, \dots, n_S) : n_i \leq m_i^{(s)} \wedge \sum_{i=1}^S n_i R_i^{(s)} \leq C_k\}.$$

Proof sketch. The proof of this convolution for finite capacity queues follows from algebraic manipulations for the product form expression (3). \square

Theorem 2 allows us to define a recursive algorithm to compute the normalizing constant $G^{(s)} = G_{\{0, \dots, N^{(BS)}\}, \mathbf{m}^{(s)}}^{(s)}$. The complexity of this algorithm is bounded by $O((N^{(BS)} + 1)H^2)$, $H = \prod_{i=1}^S (m_i^{(s)} + 1)$.

4.2 Product Form for the network of elastic services conditioned on the state of streaming networks

Assume that the network of streaming services is in a certain state $\mathbf{N}^{(s)}$ and persists in this state for a sufficient long time to let the stochastic process underlying the network of elastic services reach its limiting distribution. Clearly, this distribution depends on $\mathbf{N}^{(s)}$ since this state determines the residual data rate for elastic services. Henceforth, in this subsection, we will reason under this assumption.

Similarly to the case of streaming services, the network of elastic services presents a product form solution. However, in this case the stations have constant service capacity (with the exception of the delay station) and processor sharing discipline. Given the small amount of data corresponding to elastic services, we ignore the effects of mobility. Therefore, Theorem 1 reduces to the result of [19], i.e., functions $g_{k,t}$ must be replaced by:

$$h_{k,t}(n_{k,t}^{(e)}) = \begin{cases} \left(\frac{\rho_{0,k}}{\mu_{k,t}^{(e)}} \right)^{n_{k,t}^{(e)}} & \text{if } 0 < k \leq N^{(BS)} \\ \frac{1}{n_{0,t}^{(e)}!} \left(\frac{1}{\mu_{0,t}^{(e)}} \right)^{n_{0,t}^{(e)}} & \text{if } k = 0 \end{cases}$$

where, for $k > 0$, we define the service rate for class t elastic service as the proportion of the residual data rate used for serving the requests:

$$\mu_{k,t}^{(e)} = \frac{C_k - \sum_{i=1}^S R_i^{(s)} N_{k,i}^{(s)} + C_k^{(e)}}{\varphi_t^{(e)}} \frac{n_{k,t}^{(e)}}{\sum_{i=1}^E n_{k,i}^{(e)}}.$$

Also the network of elastic services is insensitive to moments higher than the first of the data to be transferred.

Finally, although convolution Theorem 2 holds also for the elastic network, we can resort to the mean value analysis proposed in [24] since the network satisfies the assumptions of the model with finite capacity stations analyzed in this work. In this case, the computational complexity is bounded by $O((N^{(BS)} + 1)^2 H)$, $H = \prod_{i=1}^E (m_i^{(e)} + 1)$, where $m_i^{(e)}$ denotes the total number of elastic services of class i since all stations with the exception of the delay station have fixed service rate.

4.3 Approximate solution of the joint model between streaming and elastic services

The standard approach to the approximate solution of Markov modulated processes is that of computing the stationary probabilities of the modulated process (elastic service network) conditioned to each state of the environment (streaming service network) assuming that this does not change (see, e.g., [15]). The underlying assumption is that the modulated process reaches the stationary behavior conditioned to the environment state after a short transient and hence the marginal stationary probabilities can be obtained by averaging these values. More formally, let $\pi(e, m)$ be the joint stationary probability and $\pi_e(e)$, $\pi_m(m)$ the marginal stationary distributions of the environment and the modulated process, respectively. Then, it holds that:

$$\pi_m(m) = \sum_e \pi(m|e) \pi_e(e) \simeq \sum_e \pi^*(m|e) \pi_e(e), \quad (5)$$

where $\pi^*(m|e)$ is the stationary distribution of the modulated process conditioned to a non-changing environment in state e .

In our case, $\pi_e(e)$ can be computed efficiently thanks to Theorems 1 and 2 and the analogues for the elastic service network allow us to compute efficiently $\pi^*(m|e)$. However, we have to tackle the problem of the cardinality of the state space of the streaming services that makes a direct use of (5) unfeasible.

We propose a novel method to solve Markov modulated processes when the cardinality of the state space of the environment is very high. The idea is to sample L states from the space of the streaming network according to the technique of *perfect sampling*, i.e., the samples are chosen with a probability that exactly corresponds to their stationary probability (see, e.g. [7]).

Perfect sampling is usually achieved thanks to an algorithm called *coupling from the past* [20], but in our case we must consider that the streaming network has already been solved thanks to the convolution algorithm. This generates a table of normalizing constants from which we can draw the samples according to their stationary distribution. Suppose we have applied the algorithm to the stations in ascending order $0, \dots, N^{(BS)}$. In accordance with convolution theory on queuing networks [11, 12], the marginal distribution of station $N^{(BS)}$ is given by:

$$P\{\mathcal{N}_{N^{(BS)}} = \mathbf{N}_{N^{(BS)}}\} = g_{N^{(BS)}}(\mathbf{N}_{N^{(BS)}}) \frac{G_{\Omega \setminus \{N^{(BS)}\}, \mathbf{m}^{(s)} - \mathbf{N}_{N^{(BS)}}}^{(s)}}{G_{\Omega, \mathbf{m}^{(s)}}^{(s)}}.$$

At this point we can sample component $N^{(BS)}$ of state vector. The step can be recursively repeated for station $N^{(BS)} - 1$ conditioned on the choice made for station $N^{(BS)}$ and so on until we reach station 0. Therefore, with $N^{(BS)}$ random variates per sample and the convolution table we can construct our perfect sample set \mathcal{P} very efficiently. Thus, we can approximate Equation (5) as follows:

$$\pi_m(m) \simeq \frac{1}{L} \sum_{\mathbf{m}^{(e)} \in \mathcal{P}} \pi^*(m|e)$$

whose evaluation is fast, thanks to the low complexity of the convolution algorithm for elastic services.

The number of samples contained in \mathcal{P} is clearly related to the accuracy that one desires to achieve. However, (5) is already an approximation and represents the main source of error in the results. According to our experiments, a few hundred samples are sufficient to have very accurate results as discussed in Section 5.

4.4 Computations of the blocking probabilities

A streaming service is blocked in two cases: (i) it exits from the delay station and returns to the same station without being served by any base station or (ii) its service at a base station is interrupted because of mobility, but the job returns to the delay station. These events occur with a state-dependent probability following the skipping routing described before. To compute the blocking probability, we use a relation among the network station throughputs that we obtain from the convolution table (see, e.g., [11, 12]). More specifically, let $X_{0,t}^{(s)}$ be the throughput of the delay station for class t , and let $X_{k,t}^{(s)}$ be the throughput of station $k > 0$ for class t . Among the jobs that leave station k , only $\mu_{k,t}^{(s)} / (\mu_{k,H} + \mu_{k,t}^{(s)})$ complete their service, while the others try to access another base station. Thus, the total throughput of the network for class t can be expressed by:

$$X_t^{(s)} = \sum_{k=1}^{N^{(BS)}} \frac{\mu_{k,t}^{(s)}}{\mu_{k,H} + \mu_{k,t}^{(s)}} X_{k,t}^{(s)}.$$

Since the total service demand is $X_{0,t}^{(s)}$, the blocking probability for class t is $1 - X_t^{(s)} / X_{0,t}^{(s)}$. A similar reasoning (actually simpler, thanks to the assumption of no mobility) holds also for elastic services.

4.5 Methodological contributions

Product forms have been widely applied for the performance evaluation of telecommunication systems. This paper relies on this theory to develop a novel efficient algorithm for the solution of a queuing network with multiple classes of customers that interact in a peculiar way. To the best of our knowledge, this is the first time that the *skipping* policy, firstly introduced in [19] and still widely studied as, e.g., in [24], has been applied to networks of multiclass queues of Erlang type with alternative routing for impatient customers. This new product-form model has required a specific convolution algorithm for the computation of the normalizing constant that is a non-trivial generalization of the one

presented in [11, 12]. The performance indices obtained from the convolution output are original since they must account for the finite capacity and the competition of the customer classes. Indeed, we do not have a traditional ‘queue capacity’ that can accommodate a certain number of customers for each class, but the availability of the data rate gives the condition of blocking. Finally, we are not aware of any other work using a perfect sampling algorithm based on the convolution table to study Markov modulated processes. This approach is flexible in the sense that it allows one to balance the desired accuracy with the computational efforts by choosing the amount of samples to consider.

5 EXPERIMENTS

We present in this section two different sets of numerical results. The first one refers to a hypothetical RAN layout, while the second considers a real placement of macro and small cell base stations.

5.1 A Road Segment

We consider a portion of a RAN offering services over a road segment with one macro cell BS (named BS_M , and indexed 1) and 3 small cell BSs (named BS_{S1} , BS_{S2} , BS_{S3} , and indexed 2, 3, 4). The reach of BS_M includes the whole road segment, while the coverage of each of the three small cells includes only a portion of the road segment, and small cells are adjacent to one another. We assume that each small cell has an area equal to $1/9$ of the area reached by the macro cell, so that $1/3$ of the considered area is served by the three small cells and the remaining $2/3$ are served by the macro cell.

Since the density of UEs is assumed to be constant over time, whenever a UE requests a new service after an idle period, the request is directed to a BS with a probability proportional to the cell area, hence $2/3$ for the macro cell BS and $1/9$ for each small cell.

The requests for handover from cell i to cell j have the following values:

$$\mathbf{p}^M = [P_{i,j}^M] = \begin{bmatrix} 0 & 3/8 & 2/8 & 3/8 \\ 3/4 & 0 & 1/4 & 0 \\ 1/2 & 1/4 & 0 & 1/4 \\ 3/4 & 0 & 1/4 & 0 \end{bmatrix} \quad (6)$$

UEs can request access to the streaming and elastic services offered by the RAN. A streaming service instance can belong to one of two classes: either audio or video. A video service requires a data rate equal to 10 Mb/s, while an audio service requires 100 kb/s. Video streaming service durations are i.i.d. exponentially distributed random variables with average 1800 s. Audio durations are i.i.d. exponentially distributed random variables with average 600 s. Elastic services require the transfer of i.i.d. exponentially distributed random amounts of data with average 1 Mb. The capacity of each BS (for macro as well as small cells) is 300 Mb/s, of which 5 are reserved for elastic services. Streaming services are accepted at a BS as long as data rate is available (without consuming the 5 Mb/s reserved for elastic services). Elastic services are accepted up to a maximum number set equal to 50. The dwell times of UEs in BS_{S1} , BS_{S2} , BS_{S3} are i.i.d. exponentially distributed random variables with average 400 s, while the dwell times in BS_M are i.i.d. exponentially distributed random variables with average 600 s. Note that these numbers are derived from the specified user density in the cell and the cell area. The number of small cells we consider within the macro cell area is typical of today’s RANs (as we see in the real deployment example that we discuss next) and the fact that all BSs use the same data rate stems from the fact that operators often deploy the same equipment to implement macro and small cells, changing just the emitted power. Of course, our model copes equally well with different BS data rates.

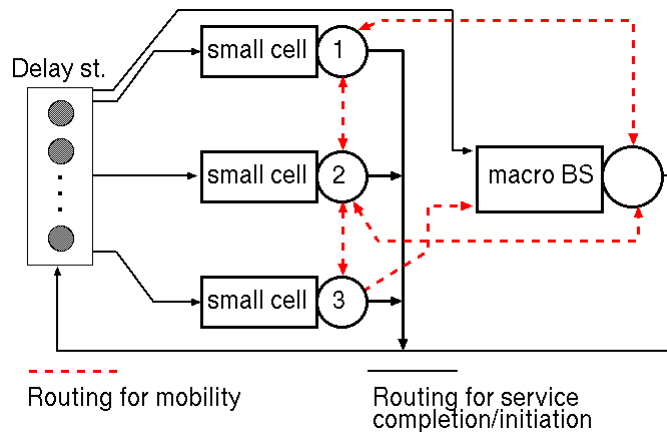


Fig. 1. Sketch of the queuing network modeling the RAN portion comprising one macro cell BS and three small cell BSs.

After a streaming service completion, UEs remain idle for i.i.d. exponentially distributed random times with average 200 s before issuing a new service request of the same type. Instead, the dynamic of elastic services is much faster. Following completion, a new request is issued after i.i.d. exponentially distributed random times with average 5 s. We consider the case of a fixed number $N^{(UE)}$ of UEs in the road segment, and we study the network performance for variable values of $N^{(UE)}$. For every video user in the area, we have 2 audio and 2 elastic users. We study the system for up to 200 video users, hence one thousand users in total.

Figure 1 shows the structure of the queuing network model used to study the described portion of a RAN.

The first two plots in Figure 2 report the blocking probabilities for video and audio streaming services and for elastic services, as a function of the number of video streaming customers in the network of queues.

In Figure 2a we report analytical predictions and simulation results for the blocking probability of audio and video services. Observe the extremely good match between analytical and simulation results (the analytical model is indeed exact for these services, so that the good match is a validation of the simulator). The other evident and surprising element of the plots is the oscillation in the blocking probability of audio services. This phenomenon is in line with what was already observed in [17] in the case of a single cell loaded by streaming and elastic services, and results from the interplay of video services with high data rate and audio services whose data rate is two orders of magnitude smaller. Blocking one video service makes room for 100 voice services, and when the probability of blocking one more video request becomes significant, the probability of blocking for audio services decreases. It should be observed that a service level agreement (SLA) with a blocking probability limit for audio services equal to 10^{-3} is respected up to about 60 video customers, violated for a number of video customers between about 60 and 100, respected again for a number of video customers between about 100 and 130, and finally violated when the number of video customers is over about 130. This behavior is quite unexpected, and should be known to an MNO, who can activate appropriate measures to avoid user dissatisfaction (for example setting a limit to the maximum number of simultaneously active video services which is compatible with the first crossing of a desired threshold for the blocking probability of audio services, as predicted by our model – if such threshold is set to 10^{-3} , our model says that the number of active video services should not exceed 60; this is easily implemented through an admission control algorithm).

573

574

575

576

577

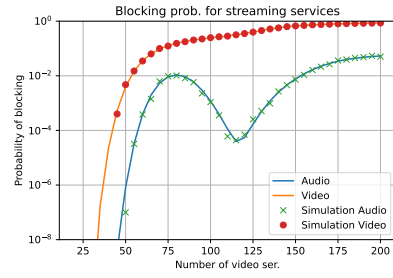
578

579

580

581

582



(a) Blocking probabilities for streaming services.

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

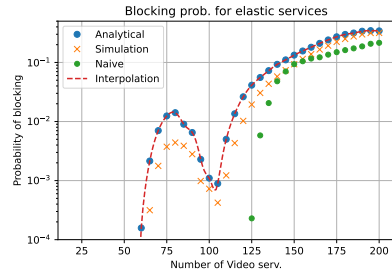
620

621

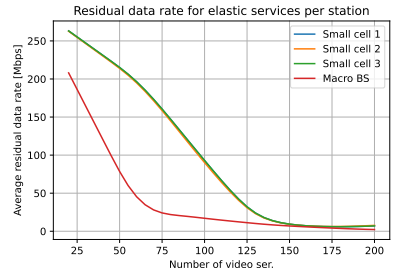
622

623

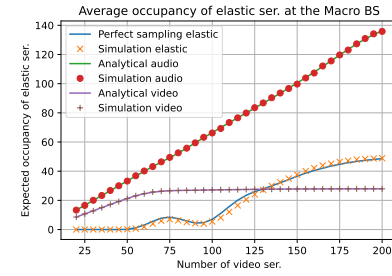
624



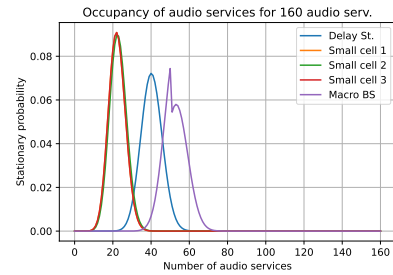
(b) Blocking probabilities for elastic services.



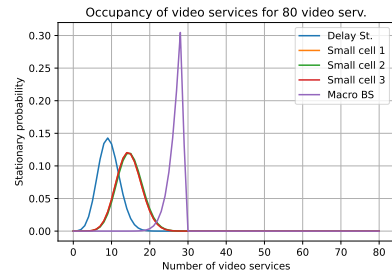
(c) Available data rate for elastic services



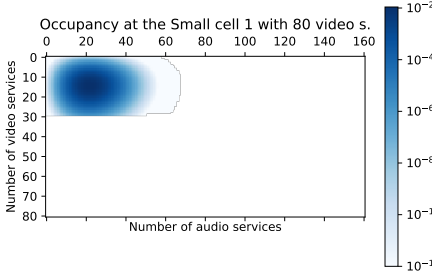
(d) Average number of services at the macro BS



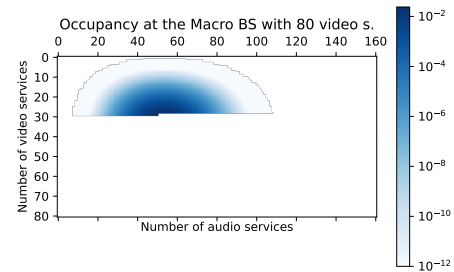
(e) Distribution of audio services at the stations



(f) Distribution of video services at the stations



(g) Distribution of streaming services at the small cells 1-3



(h) Distribution of streaming services at the macro BS

Fig. 2. Numerical results for the case of one macro cell BS with user density 1 and three small cell BSs with user density 2.

In Figure 2b we report analytical predictions and simulation results for the blocking probability of elastic services. We can see that the oscillations in blocking probability are present also in this case, and that the analytical model provides predictions which are pessimistic with respect to simulation. It is quite interesting to observe that the careful model we presented in this paper allows much better results to be obtained with respect to a "naive" approach that only considers the average data rate that remains for elastic services after the allocation to streaming services. In this case, analytical predictions (green dots in Figure 2b) come close to simulation only when blocking probabilities approach 10%, and are unreasonably optimistic otherwise.

Figure 2c reports the residual data rate as a function of the number of video streaming customers in the network of queues at the macro BS (red curve) and at the three small cell BSs (the curves overlap because the cells are assumed to have the same size and the same user density). We see that in this case the macro BS carries more load deriving from streaming services, which is due to the fact that its area is 6 times larger.

Figure 2d shows the average numbers of active audio, video and elastic services at the macro BS as well as in each small cell. It is important to note the oscillation in the number of active elastic services, as well as the very good match between simulation results and results obtained with perfect sampling.

Figures 2e and 2f report the probability density functions of the numbers of audio and video services in progress at the four BSs, as well as the number of customers that are experiencing a pause between a service and the next one, in the case in which 80 video customers, 160 audio and 160 elastic customers are present in the queuing network. Again, because of symmetry, the curves of the three small cell BSs overlap. We can see that the number of active video services at the macro cell BS is with high probability close to 29, which is the maximum possible number (remember that each BS can allocate up to 295 Mb/s to streaming services, and each video service requires 10 Mb/s). The spike in the distribution of the number of audio services in progress at the macro BS for 50 audio services is quite interesting. This value corresponds to the allocation of the remaining 5 Mb/s (of the 300 available, after 290 have been allocated to video and 5 are reserved to elastic) to audio services. The same phenomenon is also visible in Figures 2g and 2h, which represent through heat maps the joint probability distributions of the numbers of active video and audio services at small cells (equal for the three small cells because of symmetry) and at the macro. We can see that the probability distribution at the macro peaks around the discontinuity at 29 videos and 50 audio.

To see what happens when the three small cells are not symmetric, we also present results for the case in which the macro cell, as well as small cells 1 and 3 have user density equal to 1, and only small cell 2 has user density equal to 2.

This leads to handover probabilities from cell i to cell j with the following values:

$$\mathbf{P}^M = [P_{i,j}^M] = \begin{bmatrix} 0 & 1/5 & 3/5 & 1/5 \\ 3/4 & 0 & 1/4 & 0 \\ 1/2 & 1/4 & 0 & 1/4 \\ 3/4 & 0 & 1/4 & 0 \end{bmatrix} \quad (7)$$

and to dwell times with average equal to 600 s for all cells.

Results for this case are shown in Figure 3. In Figures 3a and 3b we report blocking probabilities for streaming and elastic services, and we again see the oscillations already observed in Figure 2a, as well as the superior accuracy of our modeling approach with respect to the naive analysis. In Figure 3c we see that in this case the residual data rate is lower for small cell 2 with respect to small cells 1 and 3, due to the higher user density. Finally, in Figure 3d we report the average numbers of active audio, video and elastic services at small cell 2. Also in this case, the match between simulation results and results obtained with perfect sampling is very good.

677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728

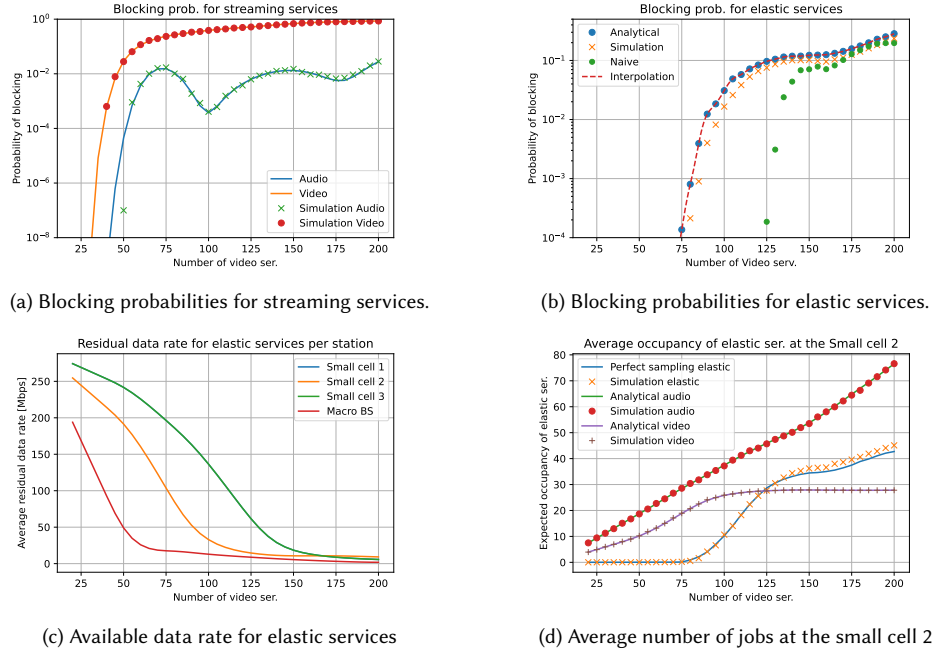


Fig. 3. Numerical results for the case of one macro cell BS with user density 1 and three small cell BSs with user density 1, 2 and 1.

5.2 A University Campus

Our second numerical example considers a real deployment of BSs over a University Campus. Two macro cell BSs and four small cell BSs provide services over the campus area, and are positioned as shown in the map of Fig. 4 (taken from the web site lteitaly.it that reports BS positions in Italy), in which grey and purple markers represent macro cells and small cells, respectively. The macro cell BS in position 1 (called macro BS 1) has within its coverage area the small cell BSs in positions 3, 4 and 5 (called small BS 3, 4 and 5). The small BS 5 is also within the coverage area of macro BS 2, and so are small BSs 3 and 6.

The parameters that characterize the system are the same as for the previous experiments, except for the fact that macro BS 1 has capacity 300 Mb/s and all other BSs have capacity 100 Mb/s (with 5 Mb/s reserved to elastic services at all BSs), the average size of the files to be transmitted by elastic services is 100 kb, the average think time for elastic services is 500 ms, and dwell times in small BSs 3, 4, 5 and 6 are, respectively, 702, 583, 400, 654 s. The handover probabilities (for streaming services) from cell i to cell j have the following values:

$$P^M = [P_{i,j}^M] = \begin{bmatrix} 0.0 & 0.1 & 0.187 & 0.613 & 0.1 & 0.0 \\ 0.1 & 0.0 & 0.319 & 0.0 & 0.1 & 0.481 \\ 0.465 & 0.516 & 0.0 & 0.0 & 0.019 & 0.0 \\ 0.284 & 0.0 & 0.0 & 0.0 & 0.716 & 0.0 \\ 0.1 & 0.1 & 0.29 & 0.324 & 0.0 & 0.186 \\ 0.0 & 0.349 & 0.0 & 0.0 & 0.651 & 0.0 \end{bmatrix} \quad (8)$$

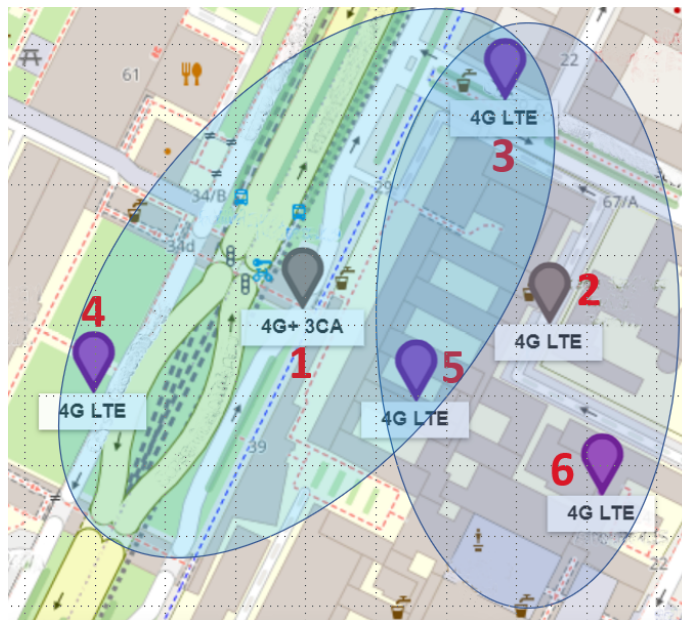


Fig. 4. Layout of the BSs covering a University Campus. Macro cell BSs are in positions 1 and 2; small cell BSs are in positions 3, 4, 5, and 6.

The average residence times and the routing probabilities have been estimated taking into account the sizes of the areas covered by the cells and by balancing the mobility flows of users among these areas. To obtain a closed system, we have assumed that users move among the considered cells neglecting incoming (departure) flows from (to) the outside [1].

In this setting, the blocking probabilities for streaming and elastic services are reported in Fig. 5 as a function of the number of video users in the system (for each video user there are also 2 audio users and 2 elastic users in the system; hence, with 200 video users, the total number of users in the system is 1000). In Fig. 5b, which shows the blocking probabilities for elastic services, blue and green dots represent the analytical results of our model and the naive model, respectively, while orange crosses correspond to simulation results. The curves once more prove the accuracy of our approach and its superiority to the naive approach based on the use of the average residual data rate for the computation of the blocking probability of elastic services.

The average residual data rate for elastic services for all BSs is presented in Fig. 6. Macro BS 1 exhibits higher values of the average residual data rate with respect to all other BSs because it has three times more capacity. All other BSs, have equal capacity and very similar values of average residual data rate. The plot on the right zooms in the behavior with large numbers of users (from 600 to 1000 in total), allowing us to see that the data rate available to elastic services (in addition to the reserved 5 Mb/s) is significantly larger at Macro BS 1 than at all other BSs. For example, with 200 video users (1000 total users) the data rate available to elastic services is about 7 Mb/s (5 reserved plus about 2, as we see on the plot) at all BSs, except Macro BS 1, where it is about 9.5 Mb/s.

The average number of services in progress at each one of the small cell BSs is reported in Fig. 7. The behavior of the four small cell BSs is very similar, as expected, since differences are only due to (small) asymmetries in routing. The maximum average number of simultaneously active video services is 9, thus consuming 90 of the available 100

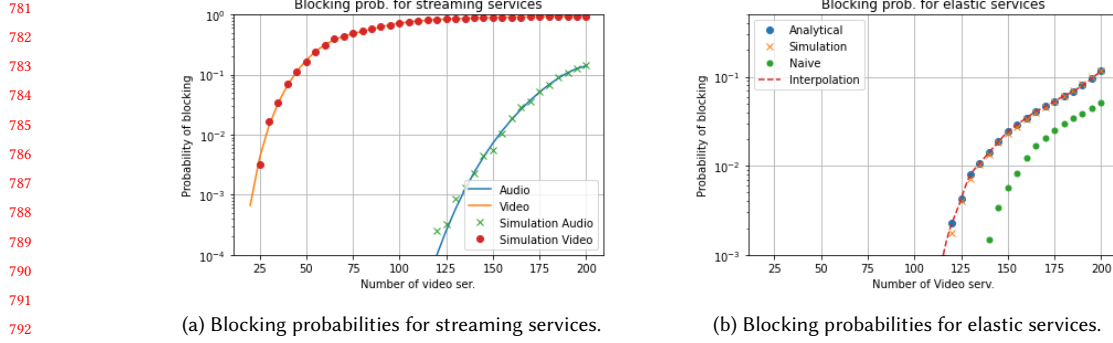


Fig. 5. Blocking probability for streaming and elastic services versus the number of video users in the real base station layout. In the system, for each video user we also have 2 audio users and 2 elastic users

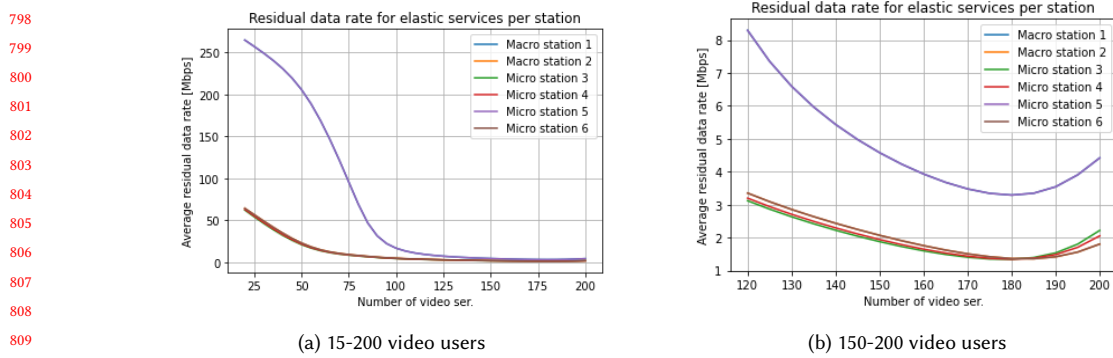


Fig. 6. Average residual data rate at all BSs in the University Campus scenario: plot (b) shows a zoom on the tail of plot (a)

Mb/s, so that 5 Mb/s remain for audios, in addition to the 5 Mb/s reserved for elastic services. The average number of simultaneously active audio services grows up to 50, thus bringing the data rate allocated to streaming services to 95 Mb/s, hence respecting the reservation of 5 Mb/s for elastic services. The average number of elastic services in progress is very low when the number of users is small, because the residual data rate is high, so that each elastic service is allocated a high data rate and thus completes very quickly (transferring 1 Mb at 5 Mb/s requires 200 ms only). Instead, when streaming services consume a large portion of the data rate available at the BS (we can see in Fig. 6 that this begins to happen for a number of video users between 75 and 100), the average residual data rate becomes small, the duration of each elastic service grows, and so does the average number of elastic services in progress. Note that the average number of simultaneous elastic services remains far from the maximum value permitted, namely 50.

The average number of services in progress at each one of the macro cell base stations is reported in Fig. 8. The results for the two macro BSs are obviously different, since macro BS 1 has a data rate three times higher than macro BS 2. For this reason we see that in macro BS 1 the average number of simultaneously active video services grows up to 29, hence consuming 290 Mb/s out of the 295 available for streaming services. On the contrary, at macro BS 2 the number of active video services grows only up to 9, hence consuming 90 Mb/s out of the 95 available. In both BSs the

Manuscript submitted to ACM

833 average number of simultaneously active audio services grows up to 50, consuming the 5 Mb/s that are not used by
834 video and not reserved to elastic services.

835 The average number of simultaneously active elastic services, like in the previous plots, is very small before the
836 available data rate saturates, and also in this case does not reach the maximum value of 50.

837 Quite interesting is the fact that if we increase the average file size to be transferred by elastic services to 1 Mb, and
838 we simultaneously increase the average time between the end of an elastic service and the issue of a new request to 5 s
839 (the same parameters we used for the case of one macro cell BS and three small cell BSs in Section 5.1) we obtain very
840 accurate results for all cells, except Macro cell BS 1, whose results are reported in Fig.9. In this figure, we can observe a
841 significant discrepancy between the average number of simultaneously active elastic services predicted by our model
842 and the number computed by detailed simulations. The reason of this discrepancy is in the transient behavior in the
843 dynamics of elastic services between state changes of streaming services. Indeed, our model computes the steady state
844 distribution of the number of elastic services for every configuration of streaming services and averages the results.
845 This implies that our model is oblivious to the transients in the number of elastic services generated at every change of
846 state of streaming services. Consider for example a change of state induced by the termination of one video service.
847 When this happens, all of a sudden 10 Mb/s become available, and are immediately exploited by elastic services, whose
848 number goes down very fast, toward the steady state associated with the new value of available data rate. Conversely,
849 when a new video service starts, the available data rate is reduced by 10 Mb/s, but the number of active elastic services
850 only grows toward the new steady state value with the dynamics of their arrival rate combined with the new (slower)
851 service rate. This means that the reaction of the number of elastic services to a video completion is faster than the
852 reaction to a video start, so that the real average number of simultaneously active elastic services results lower than
853 what our model predicts, as we see in Fig. 9. A confirmation of this effect is given by the fact that considering the faster
854 dynamics of elastic services that we used to generate the results in Fig. 8, the inaccuracy becomes negligible. The reason
855 why the inaccuracy is visible only for Macro cell BS 1 is in its higher data rate with respect to all other BSs (300 Mb/s
856 instead of 100). This leads to the possibility of 29 simultaneous active video services (instead of 9), which makes the
857 dynamics of the number of active video services over three times faster for this BS, hence generating a larger number
858 of transients for elastic services.

859 5.3 Comparison with a routing policy without skipping

860 Up to now, we discussed the accuracy of the analytical model by comparing against simulations that implement the same
861 routing that is necessary to obtain the product form solution for the network of queues that describes the dynamics of
862 streaming services, i.e., the skipping routing policy. Indeed, we observed that the extremely good match between the
863 analytical and simulation results for streaming services is a validation of the simulator rather than the model, while the
864 very good accuracy of the results for elastic services is an indication of the accuracy of the approximations introduced in
865 their analytical model solution. However, the actual behavior of service requests in RANs does not follow the skipping
866 policy, so that it is important to assess the impact of such approximation with respect to results obtained for realistic
867 service request routing. In other words, up to now we discussed the validity of the approximations introduced in the
868 model *solution*, while now we will discuss the impact of the approximations introduced in the model *construction*.

869 In Fig. 10, we show the average number of video, audio and elastic services simultaneously in progress at the Macro
870 cell BS 1 and 2 as functions of the number of video users in the system. In Fig. 11 we report the same quantities for
871 the four small cell BSs. Results show that the model provides acceptable estimates, especially in the regions before
872 saturation of resources, i.e., in regions where the system is expected to operate most of the time.

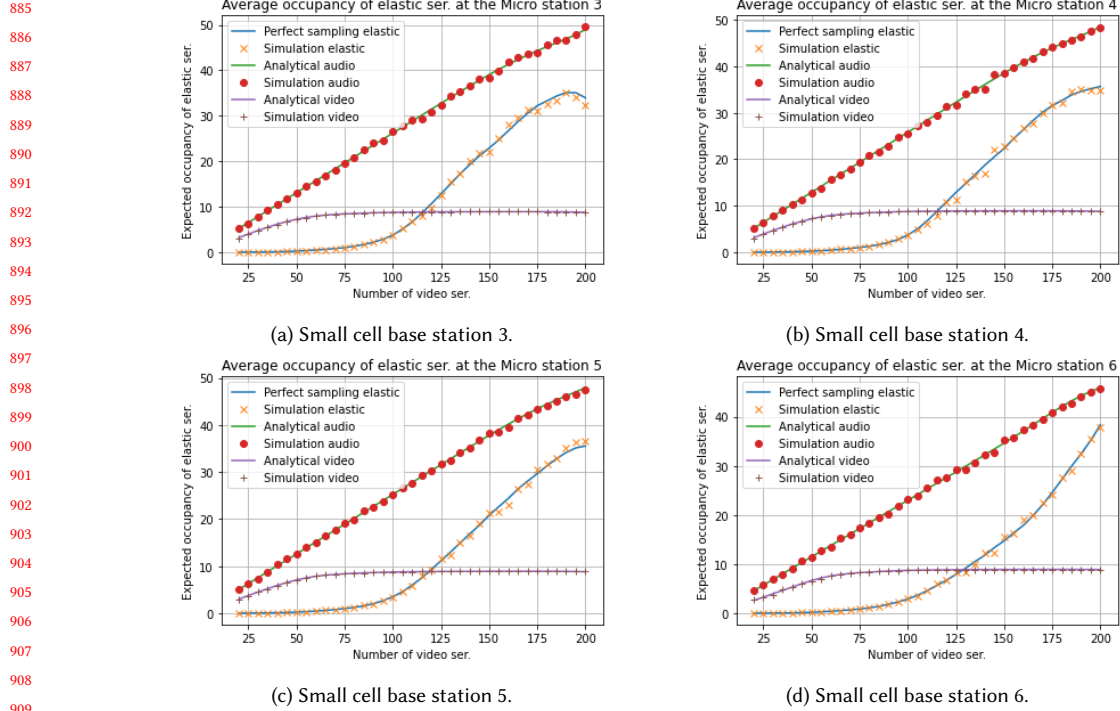


Fig. 7. Average number of streaming and elastic services active at each of the small cell base stations in the real base station layout

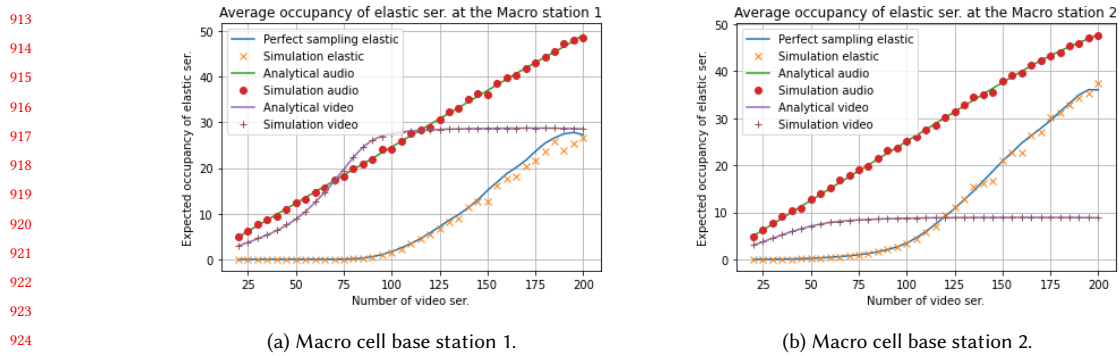


Fig. 8. Average number of streaming and elastic services active at each of the macro cell base stations in the real base station layout

Blocking probabilities for streaming and elastic service requests are presented in Fig. 12. We can see that the skipping routing policy leads to pessimistic estimated for the service request blocking probability, as can be expected by considering that skipping implies a possibility of blocking at each visited BS, while the implemented routing produces blocking only when all considered BSs have no available resources. The fact that the efficient solution technique leads to pessimistic estimates is quite positive for applications in network planning and design, since it provides conservative estimates, and leads to some overdimensioning, as customary in the networking domain.

937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988

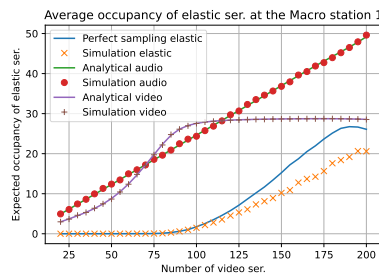


Fig. 9. Average number of streaming and elastic services active at the macro cell base station 1 in the real base station layout with 1 Mb file size for elastic services

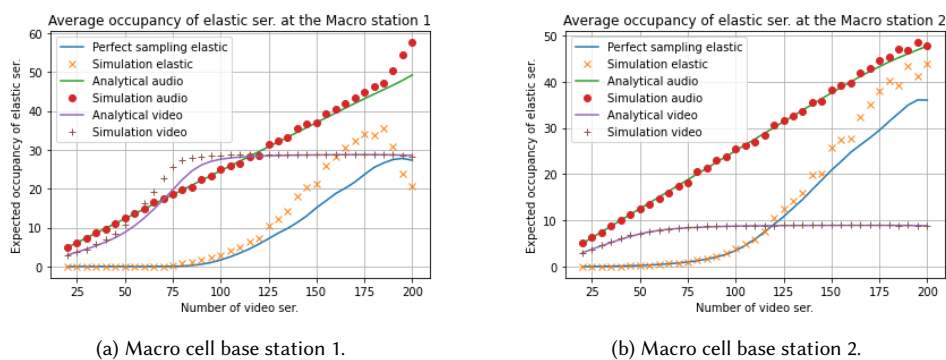


Fig. 10. Average number of streaming and elastic services active at each of the macro cell base stations in the real base station layout; model results with skipping, simulation results without skipping

6 STRENGTHS AND LIMITATIONS OF THE APPROACH

The study by simulation and the exact analysis of multiservice RANs is extremely complicated. The cardinality of the system state space makes a direct solution of the underlying stochastic process unfeasible even for quite small networks. As discussed in Section 4, approximate methods are available to study Markov modulated processes such as those used in [15, 18], but they are inapplicable to our models because of the extremely large state space of the modulating process, i.e., the network of streaming services.

One simple way to overcome these problems could be to study the network of elastic services conditioned to the *average* residual data rate not used by streaming services. This method has been compared to ours in Section 5 under the name of *naive*. Unfortunately, despite its low complexity, this approach provides extremely inaccurate and optimistic estimates of the blocking probabilities. This is due to the non-linearity of the blocking probabilities as function of the residual data rate. Indeed, blocking probabilities tend to grow only when the BS is close to saturation, thus averaging the residual data rate leads to very optimistic estimations of the performance.

Simulation is another tool that could be used to study multiservice RANs and, indeed, we have resorted to this approach to validate our solution method. However, discrete event simulations have serious problems for realistic scenarios. The most important is the fact that the dynamics of elastic services are much faster than those of streaming services and this implies that most of the events of the simulation are devoted to handle elastic services, making the

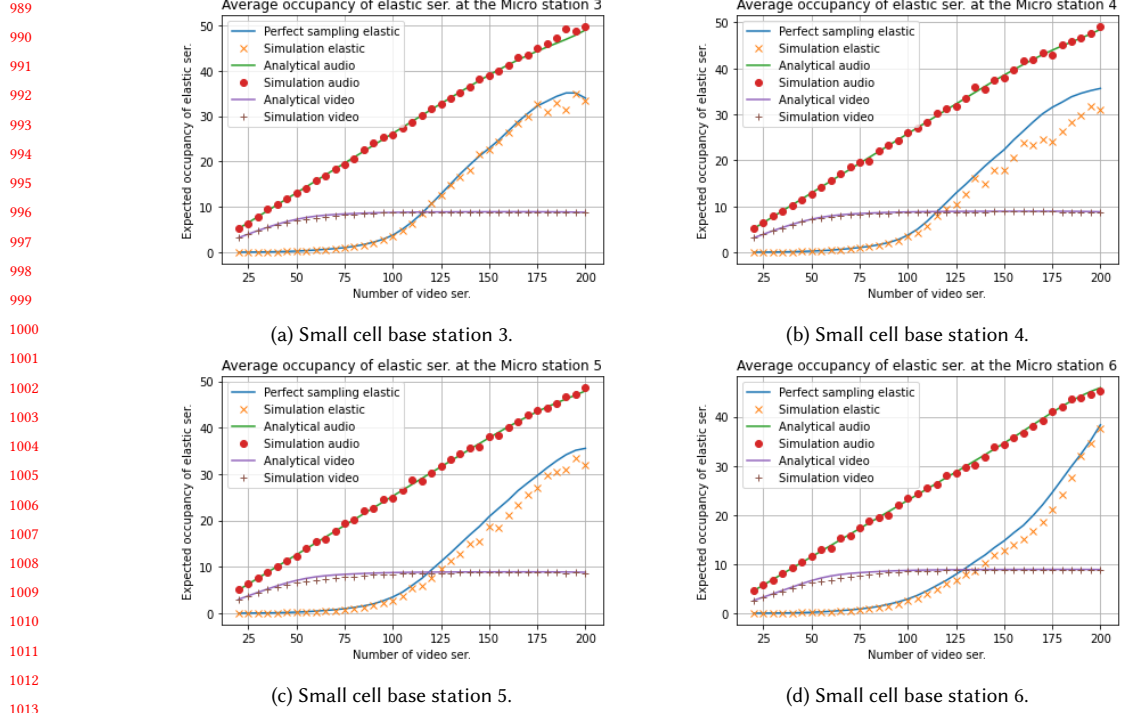


Fig. 11. Average number of streaming and elastic services active at each of the small cell base stations in the real base station layout; model results with skipping, simulation results without skipping

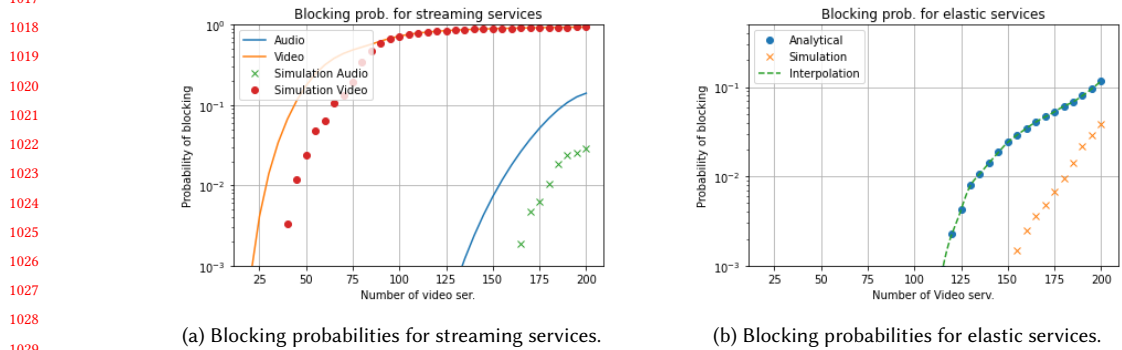


Fig. 12. Blocking probability for streaming and elastic services versus the number of video users in the real base station layout; model results with skipping, simulation results without skipping

time required to skip the transient phase quite long. Moreover, an accurate estimation of the blocking probability in moderate load for audio services (e.g., $\sim 10^{-4}$) requires to process millions of events even with a small number of cells.

It is interesting to observe that, while the problem of the different speeds of the dynamics of elastic and streaming services undermines the applicability of discrete event simulations, it is the characteristic that allows our approach to work properly. This makes the two methods somehow complementary.

1041 Finally, observe that the procedure proposed to perfectly sample a streaming network state can also be used to
1042 randomly choose the initial state of the simulation model, thus reducing the duration of the simulation transient to that
1043 required by elastic services to reach their stationary behavior.
1044

1045 The proposed approach has been shown to work for networks consisting of 7 stations and 1,000 users (400 audio,
1046 200 video and 400 elastic), which would be unfeasible with other approaches since the resulting models have state
1047 space cardinality larger than $O(10^{20})$. Notice from the asymptotic complexity of the convolution algorithm that the
1048 number of stations can grow to a dozen without creating numerical issues, while higher numbers of customers and
1049 classes negatively affect the applicability of the method.
1050

1051 7 CONCLUSIONS

1052 In this paper, we presented the first tractable analytical model for the performance evaluation of a group of base stations
1053 of a radio access network offering a mix of elastic and streaming services. The model is based on a closed network
1054 of queues, in which customers represent service instances that roam over the area served by the base stations, until
1055 either service completion or blocking. The queuing network model is exact for what concerns streaming customers,
1056 and provides an approximation for the behaviour of elastic customers. The model can be solved in product form under
1057 some specific conditions concerning customer routing, and in addition exhibits insensitivity to service distributions,
1058 provided that mobility can be neglected.
1059

1060 Numerical results show quite good accuracy of the elastic services model, and highlight unexpected oscillating
1061 behaviours that were already observed in the literature in the case of just one base station, but are described in this
1062 paper for the first time in the case of a group of base stations.
1063

1064 Capturing such behaviours in an analytical model is very important, since it allows a mobile network operator to
1065 understand the dynamics of its network, and to correct undesired behaviours with simple admission control algorithms
1066 and/or traffic management approaches.
1067

REFERENCES

- 1093
1094
1095 [1] M. Ajmone Marsan, M. Meo, and M. Sereno. Modeling simple hetnet configurations with mixed traffic loads. In *22nd IEEE International Symposium on a World of Wireless, Mobile and Multimedia Networks, WoWMoM*, pages 119–128, 2021.
- 1096 [2] S. Balsamo, P. G. Harrison, and A. Marin. A unifying approach to product-forms in networks with finite capacity constraints. In *Proc. of ACM SIGMETRICS*, pages 25–36, 2010.
- 1097 [3] G. P. Basharin and T. V. Aterekova. Analytical model of streaming and elastic traffic with dynamic channel allocation scheme. In *Proc. of International Congress on Ultra Modern Telecommunications and Control Systems*, pages 1086–1090, 2010.
- 1098 [4] F. Baskett, K.M. Chandy, R.R. Muntz, and F.G. Palacios. Open, closed and mixed networks of queues with different classes of customers. *J. of ACM*, 22:248–260, 1975.
- 1099 [5] N. Benameur, S. B. Fredj, F. Delcoigne, S. Oueslati-Boulahia, and J. Roberts. Integrated admission control for streaming and elastic traffic. In *Proc. of QoSIS, LNCS 2156*, pages 69–81, 2001.
- 1100 [6] E. Bernal-Mor, V. Pla, and J. Martinez-Bauset. Robust admission control for streaming and elastic services in cellular networks. In *Proc. of The IEEE symposium on Computers and Communications*, pages 372–374, 2010.
- 1101 [7] J. Blanchet and X. Chen. Perfect sampling of generalized Jackson networks. *Mathematics of Operations Research*, 44(2), 2019.
- 1102 [8] S. Borst and N. Hegde. Integration of streaming and elastic traffic in wireless networks. In *Proc. of IEEE INFOCOM*, pages 1884–1892, 2007.
- 1103 [9] O. J. Boxma, A. Bumb, R. Nunez-Queija, and H.-P. Tan. Integration of streaming and elastic traffic in a single UMTS cell: Modeling and performance analysis. Technical report, EURANDOM Report, 2005.
- 1104 [10] O. J. Boxma, A. F. Gabor, R. Nunez-Queija, and H.-P. Tan. Performance analysis of admission control for integrated services with minimum rate guarantees. In *Proc. of 2nd Conference on Next Generation Internet Design and Engineering (NGI)*, pages 7–47, 2006.
- 1105 [11] S.C. Bruell and G. Balbo. *Computational algorithms for closed queueing networks*. North Holland, 1980.
- 1106 [12] J. P. Buzen. Computational algorithms for closed queueing networks with exponential servers. *Communications of the ACM*, 16(9):527–531, 1973.
- 1107 [13] X. Chao, M. Miyazawa, and M. Pinedo. *Queueing Networks - Customers, Signals, and Product Form Solutions*. John Wiley and Sons, 1999.
- 1108 [14] D. Garcia, J. Martinez, and V. Pla. Admission control policies in multiservice cellular networks: Optimum configuration and sensitivity. In *Proc. of Wireless Systems and Mobility in Next Generation Internet, LNCS 3427*, pages 121–135, 2005.
- 1109 [15] E. Gelenbe and C. Rosenberg. Queues with slowly varying arrival and service processes. *Management Science*, 36(8):928–937, 1990.
- 1110 [16] S. Hanczewski, Stasiak M., and J. Weissenberg. A model of a system with stream and elastic traffic. *IEEE Access*, 9:7789–7796, 2021.
- 1111 [17] A. Marin, M. Meo, M. Sereno, and M. Ajmone Marsan. Modeling service mixes in access links: Product form and oscillations. In *Proc. of WOWMOM 2022, 2022*.
- 1112 [18] I. Mitrani. Spectral expansion solutions for markov-modulated queues. In *Proc. of IFIP Performance, LNCS 2459*, pages 17–35, 2002.
- 1113 [19] B. G. Pittel. Closed exponential networks of queues with saturation: The Jackson-type stationary distribution and its asymptotic analysis. *Math. Oper. Res.*, 4(4):357–378, 1979.
- 1114 [20] J. G. Propp and D. B. Wilson. Exact sampling with coupled Markov chains and applications to statistical mechanics. In *Proc. of Int. Conference on Random Structures and Algorithms*, pages 223–252, 1996.
- 1115 [21] R. Ramjee, R. Nagarajan, and D. Towsley. On optimal call admission control in cellular networks. *Wireless Networks J.*, 3(1):29–41, 1997.
- 1116 [22] W. Song and W. Zhuang. Multi-class resource management in a cellular/WLAN integrated network. In *Proc. of IEEE Wireless Communications and Networking Conference*, pages 3070–3075, 2007.
- 1117 [23] W. Song and W. Zhuang. Resource allocation for conversational, streaming, and interactive services in cellular/WLAN interworking. In *Proc. of IEEE GLOBECOM*, pages 4785–4789, 2007.
- 1118 [24] J. van der Gaast, R.B.M. de Koster, I. J. B. F. Adan, and J. A. C. Resing. Capacity analysis of sequential zone picking systems. *Operations Research*, 68(1):161–179, 2020.