

MODELLING OF MULTI-RATE MULTI-SERVICE TRAFFIC IN THE PRESENCE OF MOBILITY

Jesús M. Juárez, Rui R. Paulo, Eva Reguera San José, Fernando J. Velez

IT-DEM, University of Beira Interior, Calçada Fonte do Lameiro, 6201-001 Covilhã, Portugal
jejuva@iespana.es, rfrpaulo@e-projects.ubi.pt, reguerasanjo@yahoo.com, fjv@ubi.pt

Abstract - A multi-service traffic model is presented, and its validation is achieved in a vehicular scenario by using event-based simulation results. In the single-service case, theoretical and experimental results for ON-OFF blocking probability are close to each other, and there is an almost perfect concordance between theoretical and simulation values when the average sojourn time in cells is equal to the average holding time. The challenges presented by the multi-service case are identified, and a question related with a possible correspondence between blocking/handover failure probabilities and ON-OFF blocking probability are presented.

Keywords – multi-service traffic, simulation, handover, bursty behaviour, validation, mobility.

I. INTRODUCTION

Multi-rate multi-service traffic modelling is one of the most challenging issues involved in the dimensioning process of mobile multimedia communication systems, e.g., UMTS (Universal Mobile Telecommunications System) and its enhancements. On the one hand, models that consider queuing of data calls can be considered, which allows for obtaining results for delay. On the other hand, for real-time (and, in most of the cases, time-based applications), a simpler model can be used, which only models blocking and handover failure probabilities. This is the case of the Bernoulli-Poisson-Pascal/Markov-modulated Poisson Processes model from [1]. In Section II, after presenting the basis of the model, the user model and equivalent user are described, and details on how to compute the arrival rate are given. Section III presents the physical and mobility scenarios considered in the simulations, the concepts associated with the generation and termination of new and handover sessions, and the definitions of the quality of service parameters. Section IV presents results for blocking and handover failure probabilities, and also for the ON-OFF blocking probability, a measure for the blocking of bursts of traffic during multimedia traffic sessions. A comparison between theoretical and simulation results is performed, and challenges are identified. Finally, some questions are presented in Section V.

II. TRAFFIC MODEL

A. Basis of the Model

In the general model of a loss system with $R_e=1$ type of resources shared by J classes (i.e., service components), a customer arrival at the resources follows a specific random process. Each customer, i.e., service components users, requests a fixed number of resource units, i.e., channels, which are granted if available. If not, the request is cleared and the customer is blocked. The classification of customers is done on the basis of their arrival process, capacity requirement and mean holding time [1]. In this work, the performance measure that one is interested in the probability that an arriving customer is blocked, i.e., the customer or connection blocking probability, P_b . Besides, as one is dealing with bursty traffic, with active and inactive, i.e., ON and OFF, periods, one needs to address ON-OFF blocking probability. One is also interested in the handover failure probability, whose limitation directly results from the existence of a threshold for the call-dropping probability.

The capacity of the resource facilities is partitioned into capacity units. A customer is assumed to need a given number of units of each facility, and the demand is granted on a first come first served basis. If a customer demand cannot be granted, it is cleared and the new customer is blocked.

One considers J customer classes, each with different spatial and temporal requirements, and c is the total of available channels. The resource capacity vector is defined as $c_v = [c_1, \dots, c_{R_e}]$ but, since $R_e = 1$, one has $c_1 = c$. The class j (the term ‘class’ referring, in the context of this work, to ‘service component’) capacity demand of channels per customer, a_j , $j \in \mathcal{G}$, where $\mathcal{G} = \{1, \dots, J\}$, and $a_j \in \mathbb{IN}$. Besides, the time that these channels, once granted, will be held by the service component (or class) j customer is i.i.d., it being specified by its mean value, whose specific distribution has no influence on the calculations that one is pursuing [1], [2]. Thus, given these considerations, the capacity vector, \mathbf{A} , is a vector of the following type

$$\mathbf{A} = [a_j]_{j=1, \dots, J} . \quad (1)$$

Let the number of class j active customers, i.e., that hold their a_j resources at time t , be represented by the random

variable $N_j(t)$. One can then express the state of the system by

$$\mathbf{N}(t) = (N_1(t), \dots, N_J(t)) \quad (2)$$

and $Y(t)$, the current resource occupancy vector as a function of the system state variable

$$Y(t) = \mathbf{N}(t) \cdot \mathbf{A}. \quad (3)$$

The set of possible states \mathcal{N} is bounded as a result of the finite resource capacity

$$\mathcal{N} = \left\{ \mathbf{n} \in \mathbb{N}^J : [n_1, \dots, n_J] \bullet \begin{bmatrix} a_1 \\ \dots \\ a_J \end{bmatrix} \leq c \right\} \quad (4)$$

where $\mathbf{n} = (n_1, n_2, \dots, n_J)$ is the state of the system (defining the number of each component active requests). In the limit, if there are more users from an application than from another, the other can use fewer channels.

If the state vector $\mathbf{N}(t) \in \mathcal{N}$, then the service occupancy vector $Y(t) \in \mathcal{Y}$, where \mathcal{Y} is simply defined by

$$\mathcal{Y} = \{y \in \mathbb{N} : y \leq c\}. \quad (5)$$

The equilibrium *pmf* of the state $N(t)$ and the occupancy $Y(t)$ are defined as follows

$$p(\mathbf{n}) = \lim_{t \rightarrow \infty} \text{Prob}\{\mathbf{N}(t) = \mathbf{n}\}, \quad (6)$$

$$q(y) = \lim_{t \rightarrow \infty} \text{Prob}\{Y(t) = y\}. \quad (7)$$

When the system is in state $N(t) = \mathbf{n}$, the time until the next arrival of a class j customer's demand is exponentially distributed with parameter $\lambda_j(n_j)$. This parameter is normalized with respect to the average class j holding time, thus, a different time unit is introduced for each customer class. Blocking takes place if a request cannot be granted entirely, i.e., a class j request arrives when the system is in the set

$$\mathcal{B}_j = \{\mathbf{n} \in \mathcal{N} : \mathbf{n} \cdot \mathbf{A} + a_j > c\}. \quad (8)$$

For exponential holding times, the BPP process can be modelled by a Markov chain, although this model allows for considering more general distributions for the holding times.

While the equilibrium *pmf* of the state $N(t)$, $p(\mathbf{n})$, has a product form, in [1] an algorithm is proposed to compute the occupancy *pmf*, $q(y)$, that there is an algorithm that is economic in terms of computation time and storage space – as long as the number of resources is not too high. This algorithm assumes a BPP arrival process.

BPP processes are those whose arrival intensity (corresponding to an exponential distribution of the inter-arrival times), conditioned to n_j customers being in the system, is of the form

$$\lambda_j(n_j) = \alpha_j + n_j \beta_j, \text{ with } \alpha_j > 0 \quad (9)$$

where $(-\beta_j)$ is the activation rate and α_j is the arrival rate. In the Poisson case, as $\beta_j=0$, the *pmf* of the number of active customers in an infinite resource is [1] $\lambda_j(n_j) = \alpha_j$.

For exponential holding times, a BPP process can be modelled by a Markov chain. Due to the normalization, mean holding times are unitary, thus, death rates are integer values.

The description and the pseudo-code for the algorithm for the computation of time and call blocking probabilities, P_{bt} and P_b , respectively, are presented in [1], [3].

B. User Model and Equivalent User

There are a total number of c available resources (or channels) in each cell, being used by a total number of equivalent users, M_T . Furthermore, one is considering two applications, voice (VOI) and Video-telephony (VTE), i.e., a total of $K_{app} = 2$ applications. The index k , $k = \text{VOI}, \text{VTE}$, refers to these applications. Given this traffic mixture, the model for applications activation by users is presented in Fig. 1. Each user can be either in an idle state or using one of the two applications, with generation rate, Λ_k , and total service rate, H_k , respectively.

Once application k is active, the service components are activated with rate $\Lambda_{j/k}$ and extinguished with total service rate $H_{j/k}$, $j = 1, \dots, J$, Fig. 2; they can be simultaneously active, or not, and some can even not be activated for a given application.

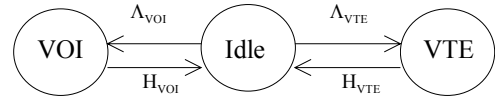


Fig. 1. Applications activation.

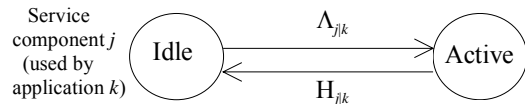


Fig. 2. Service components activation.

This is a loss system, whose performance can be measured by the blocking probability of each service component, which simplifies the analysis (because one only needs to consider the service components and not each application).

From Fig. 1 it is straightforward to derive the probability of a user having application k active

$$p_k = \frac{\Lambda_k / H_k}{1 + \sum_{i=1}^{K_{app}} \Lambda_i / H_i}. \quad (10)$$

The total traffic is

$$\rho = \sum_{k=1}^{K_{app}} (\Lambda_k / H_k). \quad (11)$$

Hence, the traffic generated by a specific application can be given by

$$\rho_k = \frac{\Lambda_k}{H_k} = prop_k \cdot \rho, \quad (12)$$

where $prop_k$ is the proportion of application k active users (numerically given by its usage).

A normalised generation rate is introduced

$$\Lambda_k^* = prop_k \cdot H_k, \quad (13)$$

such that

$$\sum_{k=1}^{K_{app}} \Lambda_k^* / H_k = 1 \quad (14)$$

(one has), reducing (10) to

$$p_k = \frac{\rho}{1+\rho} \cdot \left(\frac{\Lambda_k^*}{H_k} \right) = \frac{\rho}{1+\rho} \cdot prop_k. \quad (15)$$

The fraction of active users is

$$f = \frac{\rho}{1+\rho}. \quad (16)$$

The product $prop_k \cdot b_k$ gives the so-called maximum load per application, which represents the average load per application.

The Poisson case of the BPP model is used [1]. In the context of this model, each equivalent user (of the applications) origins a given number of actual users, one for each service component, using resources with data rates B_{sj} , during a time exponentially distributed with average $1/H_j$, $j=1, \dots, J$ (J is the total number of service components; for the service component j , one has a total service rate $H_j = \mu_j$ in the static case, and $H_j = \mu_j + \eta_j$, otherwise, where μ_j is the service rate, and η_j is the cross-over rate); so the number of users accessing each individual service component is also M_T .

The activation rate of each component j , given an application k , is defined by [3]

$$\begin{aligned} \Lambda_{j|k} &= \frac{E[\text{number of service component } j \text{ requests}]}{E[\text{duration of application } k]} = \\ &= \frac{n_{j|k}}{1/\mu_k} = n_{j|k} \cdot \mu_k \end{aligned} \quad (17)$$

where $n_{j|k}$ is the number of times the service component j is activated during application k , it being unitary for permanent service components. When the application is static (due to absence of mobility, or because it is not considered), $H_{j|k}$ is given by $\mu_j = 1/\tau$ if the service component is permanent (τ

being application k average duration), or by $1/\tau_s$ if the application is not permanent, τ_s being the average service component duration. Although the values of $\Lambda_{j|k}$ and $H_{j|k}$ change when the influence of mobility of terminals is considered (by a factor associated with the service and the cross-over rates of application k), their ratio is maintained constant, and no change will exist for traffic analysis purposes, except for the blocking/handover failure thresholds.

The number of active users of component j using their resources at time t , a_j , are represented by the random variable $N_j(t)$. As it was already pointed out, when the system is in state $N(t) = n$, the time until next user class j arrival is exponentially distributed with parameter $\lambda_j(n_j)$ (9).

C. Arrival Rate

The arrival rate is normalised with respect to the total service rate of service component j

$$\alpha_j^{norm} = \alpha_j / H_j, \quad (18)$$

meaning that different time scales are introduced for each service component.

The symmetric of the arrival rate of each service component is given by multiplying the expectation of $\Lambda_{j|k}$ by N_j , leading to [3], [4],

$$(-\alpha_j) = N_j \cdot \sum_{k=1}^{K_{app}} \Lambda_{j|k} \cdot p_k = N_j \cdot \frac{\rho}{1+\rho} \cdot \sum_{k=1}^{K_{app}} \Lambda_{j|k} \cdot prop_k \quad (19)$$

If the system is stationary, the average occupancy of component j multiplied by N_j is given by the following ratio

$$(-\alpha_j^{norm}) = \frac{-\alpha_j}{H_j} = N_j \cdot \sum_{k=1}^{K_{app}} \frac{\Lambda_{j|k}}{H_{j|k}} \cdot p_k, \quad (20)$$

here called (the symmetric of) the normalised arrival rate, meaning that the service rate of service component j is

$$H_j = \sum_{k=1}^{K_{app}} \Lambda_{j|k} \cdot \frac{\Lambda_k^*}{H_k} \left/ \left(\sum_{k=1}^{K_{app}} \frac{\Lambda_{j|k}}{H_{j|k}} \cdot \frac{\Lambda_k^*}{H_k} \right) \right. \quad (21)$$

This does depend on mobility, because of the dependence of the numerator on it. If terminal mobility is considered, $\Lambda_{j|k}$ has to be replaced by $\Lambda_{j|k}$ times a factor $(\mu_k + \eta_k) / \mu_k$, where μ_k and η_k are the service and the cross-over rates associated with application k . The holding times for every service component should be i.i.d.. An example is the particular case of having exponential distributed holding times.

The data rate associated to each application is

$$b_k = \sum_{j=1}^J \frac{n_{j|k} \cdot 1/H_{j|k}}{1/H_k} \cdot B_{sj}, \quad (22)$$

where B_{sj} is the data rate associated to service component j (note that $B_{sj} = a_j \cdot B_{s1}$, where B_{s1} is the system basic data

rate). Its value does not change with the consideration of mobility, because when it is considered, H_{jik} and H_k are both affected by the factor $(\eta_k + \mu_k)/\mu_k$, the simultaneous change being cancelled by the division.

III. PHYSICAL AND MOBILITY SCENARIO

The physical scenario has a cellular architecture composed by three cells (or ten) with the shape of a roundabout. The cellular architecture consists of a backbone network which interconnects fixed base stations, and mobile units communicating with the base stations via wireless links. Each cell has access to $c=N$ channels.

The call holding time is the average call duration if the call is not prematurely dropped, and it is assumed to be exponentially distributed with average $1/\mu$, where μ is the service rate. The transference of a mobile communication from one cell to another, while a call is in progress, is called handover (HO). If there are not enough channels available in the new cell this call will be dropped, this phenomenon is known as handover failure. The sojourn time is the time that each user stays in a cell, and it follows an exponential distribution with average $1/\eta$, where η is the cross-over rate.

The handover rate γ is given by $\gamma = \eta/\mu$, and the channel occupancy time, τ_c , is given by the minimum between the call holding time and the sojourn time. As the minimum of two variables exponentially distributed is also exponentially distributed, τ_c is exponential.

In a roundabout scenario, the traffic is homogeneous, Fig. 3. As a consequence, there is a homogeneous probability of generating new and handovers calls in the three cells with rates λ_i and η_i , respectively. Hence, $\lambda_i = \lambda \forall i$, $\eta_i = \eta \forall i$,

and $\sum_{k=1}^{N_{cells}} p_{ki} = 1 \forall i$ where p_{ki} is the probability that a call may attempt a handover from cell k to cell i , and N_{cells} is the total number cells in the geometry.

In the simulation model one uses three call generators, one for each cell, working simultaneously. Each generator models the calls of one third of the users in the entire roundabout [5].

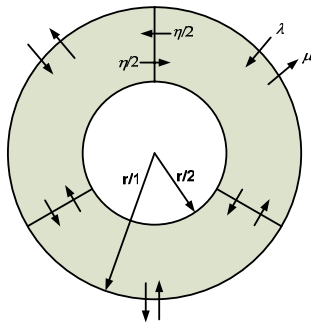


Fig. 3. Physical roundabout scenario.

The new calls are generated following a Poisson distribution with rate λ . The offered load per application is defined as $\rho = \lambda/\mu$. Packet switched traffic is commonly modelled as ON-OFF processes. Our simulator models the ON-OFF behaviour by using active/inactive time periods, according to [2]. A special model is used for real-time video-based applications like VTE due to the high level of burstiness introduced by compression techniques like MPEG-4. In simulations it was considered as having continuous occupation of channels along all the call duration.

The definition of the main concepts and parameters enables the discussion of simulation results and their comparison with other simulation results. As these parameters are general the same formulas are used for different services. The call blocking is the ratio between the number of new calls that are rejected in the process of trying to obtain channels and the total number of new calls generated. The handover failure is the ratio between the number of handovers that are rejected at the new cell in the process of trying to obtain channels, and the total number of handovers produced.

When the traffic is being modelled by ON-OFF periods, the definitions of these call level parameters will be maintained. However, new parameters are needed at the burst level. In this case, the ON-OFF blocking probability, $P_{b \text{ ONOFF}}$, is the ratio between the number of calls that are rejected at the beginning of ON periods, in the process of trying to obtain channels, and the total number of generated ON periods [6].

IV. CHALLENGES

Our simulator was used for the validation of traffic models. One performed a comparison between the theoretical values obtained by considering the Bernoulli/ Poisson/Pascal model for multi-service traffic [4], [6], and the results obtained by using the AweSim simulator [5]. Results for bursty VOI are presented in Fig. 4, where a comparison of theoretical and simulation results for $P_{b \text{ ONOFF}}$ is performed for different values of γ (VOI, $c=4$). Exponential distributions are considered for the active/inactive periods, an average session duration of 60 s is assumed, and the time intervals between arrivals are the ones presented in Table 1. For VOI, the theoretical and the experimental values of $P_{b \text{ ONOFF}}$ are close to each other, Fig. 4 ($\rho=0.2$ Erl).

Table 1
Time intervals between arrivals for single-service.

ρ	Time between calls [s]
0.05	257.14
0.10	128.57
0.15	85.71
0.20	64.29

There is an almost perfect concordance between theoretical and simulation values for $\gamma=1$, i.e, when the average sojourn time in cells is equal to the average holding time. The curves for P_b and P_{hf} follow a similar behaviour but P_{hf} takes lower values.

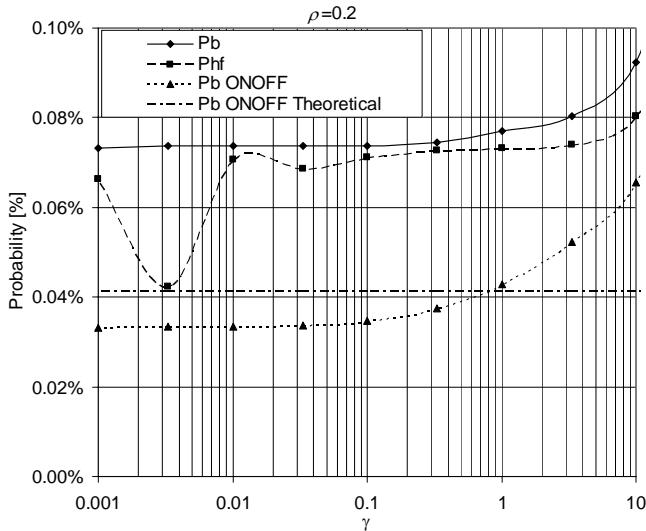


Fig. 4. P_b , P_{hf} , P_b ONOFF and theoretical P_b ONOFF as a function of γ for four channels and $\rho=0.2$ Erl and (VOI).

Besides the validation of the model for the bursty behaviour in the single-service case, one is also validating the model for the multi-service case by considering voice, VOI, video-telephony, VTE, simultaneously.

V. QUESTIONS

How will it be possible to improve the validation of the model for multi-service?

How can P_b and P_{hf} be obtained from the theoretical results for P_b ONOFF?

ACKNOWLEDGEMENTS

This work was partially funded by MULTIPLAN and CROSSNET (Portuguese Foundation for Science and Technology POSI and POSC projects with FEDER funding), and by "Projecto de Re-equipamento Científico" REEQ/1201/EEI/ 2005 (a Portuguese Foundation for Science and Technology project).

REFERENCES

[1] G. A. Awater and H. A. van de Vlag, "Exact computation of Time and Call Blocking Probabilities in Large, Multi-traffic, Multi-resource Loss Systems," *Performance Evaluation*, Vol. 25, No. 1, Mar. 1996, pp. 41-58.

- [2] J. S. Kaufman, "Blocking in a Shared Resource Environment," *IEEE Transactions on Communications*, Vol. COM-29, No. 10, Oct. 1981, pp. 1474-1481.
- [3] R. M. Carvalho and J. M. Brázio, "Multi-service Traffic Model for the Sharing of a Resource by a Homogeneous Population of Users with Stochastically Heterogeneous Demands (in portuguese)", in *Proc. of 6th Conference of the Portuguese Statistic Society*, Tomar, Portugal, June 1998.
- [4] R. M. Carvalho, *Multi-service Traffic Models for Cellular Mobile and Personal Communication Systems* (in portuguese), Graduation Report, Instituto Superior Técnico, Technical University of Lisbon, Lisbon, Portugal, Jan. 1998.
- [5] J. M. Juárez, R. R. Paulo and F. J. Velez, "Tele-Traffic Simulation for Mobile Communication Systems Beyond 3G," in *Proc. of AICT' 06 - The 2nd Advanced International Conference on Telecommunications*, Guadeloupe, French Caribbean, Feb. 2006.
- [6] J. M. Juárez, *Tele-Traffic Simulation for Mobile Communication Systems Beyond 3G*, Graduation Thesis, University of Beira Interior, Covilhã, Portugal, Sep. 2005.