

# On the Integration of Best-Effort and Guaranteed Performance Services

Eitan Altman, Damien Artiges, Karim-Frédéric Traore

► **To cite this version:**

Eitan Altman, Damien Artiges, Karim-Frédéric Traore. On the Integration of Best-Effort and Guaranteed Performance Services. RR-3222, INRIA. 1997. <inria-00073467>

**HAL Id: inria-00073467**

**<https://hal.inria.fr/inria-00073467>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

*On the Integration of Best-Effort and Guaranteed  
Performance Services*

Eitan ALTMAN Damien ARTIGES and Karim TRAORE

**N° 3222**

July, 1997

\_\_\_\_\_ THÈME 1 \_\_\_\_\_



*R*apport  
de recherche





## On the Integration of Best-Effort and Guaranteed Performance Services

Eitan ALTMAN\* Damien ARTIGES\*\* and Karim TRAORE\*\*

Thème 1 — Réseaux et systèmes  
Projet Mistral

Rapport de recherche n° 3222 — July, 1997 — 25 pages

**Abstract:** One of the main challenges of emerging high speed telecommunication networks is the integration of services. Both ATM as well as the INTERNET have been evolving so as to accomodate both best-effort type traffic (such as file transfer) as well as real time traffic requiring guaranteed performance. The co-existence of these two service types on the same network carries important benefits such as resource sharing between service classes, and the ability of the user to easily select an appropriate service class according to requirements and preferences. These benefits depend on efficient network management and resource sharing strategies. The aim of this paper is twofold. First, we study the performance measures of both service types, as a function of the network management strategy. This study allows us at a second phase to design efficient network management schemes. This includes questions such as whether best-effort traffic should use only bandwidth left-over by the guaranteed performance ones, or whether (and how much) bandwidth needs to be reserved to best-effort traffic. Another management issue that we study is the pricing. We allow for different classes of guaranteed-performance traffic to have different priorities with respect to the rejection probabilities.

**Key-words:** Best-effort services, Guaranteed performance services, ATM; Quality of Service; Resource sharing; Matrix Geometric analysis; Constant Bit Rate; Available Bit Rate; Call admission control, pricing.

*(Résumé : tsvp)*

\* INRIA, B.P.93, 2004 Route des Lucioles, 06902 Sophia Antipolis Cedex, France. E-mail: altman@sophia.inria.fr URL:<http://www.inria.fr:80/mistral/personnel/Eitan.Altman/me.html>

\*\* ASCOM B.P.93, C.I.C.A, 2229 route des cretes, 06560 Sophia-Antipolis Cedex, France. Email: {artiges,traore}@ascom.eurecom.fr

## Intégration de services avec et sans réservation de ressources

**Résumé :** Cet article présente un modèle de files d'attente pour étudier les performances d'un commutateur ATM traitant les services CBR et ABR. La question abordée est de savoir s'il peut être avantageux de réserver une partie de la bande passante pour le service ABR, aux dépens de la probabilité de blocage des appels CBR.

Le commutateur ATM est représenté par un ensemble de  $N$  serveurs, les appels CBR sont la superposition de  $K$  flux d'arrivées différenciés par des seuils de rejet  $N_K \leq \dots \leq N_1 \leq N$ : un appel de classe  $k$  est accepté dans le système si le nombre de serveurs occupés par des appels CBR est inférieur à  $N_k$ , et est rejeté sinon. Tous les appels ABR sont acceptés et partagent la puissance de service non utilisée par les appels CBR. Le processus d'arrivées des appels est poissonnien et la distribution de la quantité d'informations transmises par chaque appel est exponentielle. Nous nous intéressons aux mesures macroscopiques suivantes, calculées numériquement par une approche matricielle géométrique: probabilités de blocage pour CBR, et temps de connexion moyen pour ABR.

En imposant un seuil de rejet sur les appels CBR, leur probabilité de blocage est augmentée, mais le temps de connexion pour ABR est réduit. Nous observons que le bénéfice pour ABR est particulièrement sensible en situation de trafic chargé, et qu'un seuil de rejet trop élevé sur les appels CBR peut même conduire à l'instabilité du service ABR.

Nous étudions quantitativement le compromis entre CBR et ABR en attribuant une fonction d'utilité à chaque service, et nous calculons les valeurs de seuil qui maximisent le revenu global.

Enfin, en supposant qu'une partie du trafic ABR peut utiliser le service CBR si les temps de connexions ABR sont trop grands, nous montrons comment ajuster les tarifs de chaque service pour que le réseau fonctionne au point optimal.

**Mots-clé :** services ATM; qualité de service; partage de ressources; distribution matricielle géométrique; Constant Bit Rate (CBR); Available Bit Rate (ABR); tarification ATM.

## 1 Introduction

In recent years, a big effort has been made by the ITU-T and the ATM forum to standardize several service categories, so as to allow for integration of services sharing the same network. CBR (Constant Bit Rate) and rt-VBR (real-time Variable Bit Rate) were designed to support real-time services, whereas for non-real-time traffic, the nrt-VBR (non-real-time Variable Bit Rate), ABR (Available Bit Rate), UBR (Unspecified Bit Rate) and ABT (ATM Block Transfer) services categories were developed. We make the distinction between “guaranteed services”, in which a fixed amount of bandwidth is reserved for the whole duration of the session (CBR and VBR services), and “best effort services”, in which the bandwidth allocation may change dynamically according to the source’s requirement, on one hand, and on the available bandwidth in the network, on the other hand. Although both best-effort and guaranteed performance services are supported by ATM, the appearance of best-effort service is quite recent, and even today it is generally considered to have a secondary role; it is believed that this service type would be used only to fill bandwidth unused by guaranteed-performance services. For example, in ATM Forum/95-1346 1995 [5], where an improvement of the ERICA control mechanism for ABR is proposed, it says: *“If a switch provides both ABR and VBR service (and probably CBR as well), VBR takes precedence. The available bandwidth for ABR is not link rate but what ever is left over after VBR. This is handled by the following modification to the ERICA...”*

The Internet technology follows an opposite development. It was designed originally to provide only best-effort services (TCP and UDP transport protocols). However, nowadays it evolves towards the integration of different services categories. In the IPv6 version, it is possible to provide different Quality of Services (QoS) to different connections (flows) and to obtain some service guarantees by using signaling [2, 3, 7, 12].

The main goal of this paper is to study different policies for integration of guaranteed-performance services with best-effort ones. We shall propose several tools to obtain efficient operation of both services.

We analyse and compare two network management approaches:

- (i) the approach that consists of allocating best-effort traffic the bandwidth leftover by higher priority guaranteed-performance traffic, and
- (ii) the approach that consists of limiting the amount of bandwidth available for best-effort service categories, so that the best-effort traffic has some minimum pre-allocated bandwidth.

Limiting the amount of bandwidth for the guaranteed-performance (using some Call Admission Control) results in increasing the blocking probabilities of guaranteed-performance calls. On the other hand, if best-effort traffic is allocated *only* the bandwidth leftover by

higher priority guaranteed-performance traffic, there may be situations in which the *long-run* average throughput of best-effort traffic is unacceptably low, and the time for handling a best-effort session is unacceptably large. In particular, if the *rate of arrival of best-effort sessions* is larger than the *rate that the network can handle*, then the number of best-effort connections may grow without bound. In practice, this would imply that the network will not be attractable for best-effort applications.

To describe the above situation, we may define congestion at the edges of the network (as opposed to congestion in switches) as a congestion experienced by the users. Guaranteed-performance users experience congestion by suffering from call rejections. Best effort users experience congestion when they receive very low throughputs. Because of CAC and flow control mechanisms (such as the one used in the ABR service category of ATM, or the TCP/IP in the Internet), the congestion inside the network due to an excess of data sent by the users can be avoided. However, the congestion can still exist but will now be moved to the edges of the network, where data will have to wait before being able to enter the network: requests for new (guaranteed performance) connections may be blocked, and file transmissions (using best-effort services) already in progress may suffer very large delays.

Improving transmission delays for best-effort connections may not always be an important concern, as is the case for automatic file transfer during nighttime. But for other applications, it can be a valuable service to some customers to have a best-effort service where some (long-run) average bandwidth is kept at a reasonable level (this may be a much weaker requirement than having a minimum cell rate guarantee, which is possible in the ABR service category in ATM). This may be the case for transfer of files, datagrams in interactive applications, or any service where extremely long delays are bothersome.

An alternative way to limit congestion of best-effort traffic could be to use Call Admission Control also for best-effort traffic. This, however, seems undesirable, as expressed in the ATM forum traffic management specification, concerning ABR traffic (in case it has no minimum cell rate guarantee): “the CAC will not block the connection attempt because of bandwidth allocated to other connections” ([1] p. 85). The reason that ABR sessions are not refused access to the network is that ABR users implicitly agree to use the resources left unused by other services and to accept the performance of best effort services.

We propose in Section 2 a queueing model and study both approaches for handling best-effort traffic described above. We use matrix geometric tools [9] following an approach similar to [10]. We allow in our analysis for different priorities among the guaranteed bandwidth: the actions of the call admission control may depend on the type (e.g. the Virtual Path) of the arriving session, so that different types may experience different call rejection probabilities. We

compute the rejection probabilities of the guaranteed-performance sessions, and the average throughput and sojourn time of best-effort sessions, as a function of the network management policy.

Another important factor that has an impact on the integration of services, and on limiting congestion, is pricing. We assume that long-run congestion of best-effort sessions might result in transformation of some applications to guaranteed-performance ones. This phenomenon may be enhanced or discouraged by the network management by using an intelligent pricing scheme, which we propose in Section 5. In designing such a scheme, we first identify the operation conditions that are optimal for the network manager. We then compute a pricing mechanism that induces users' behavior that achieves the network optimal conditions.

The optimal network management thus combines both the pricing and the appropriate pre-allocation of bandwidth for best-effort traffic.

## 2 Definition of the Model

We consider a service station (e.g. an ATM switch) with a total of  $\Delta$  units of bandwidth. We assume, for simplicity, that each guaranteed performance (GP) session requires the same amount  $L = 1$  unit of bandwidth, and the length of (i.e. the amount of information in) GP sessions are i.i.d. and exponentially distributed with parameter  $\mu_{\text{GP}}$ . The service station can thus process at most  $\lfloor \Delta \rfloor$  GP sessions simultaneously ( $\lfloor x \rfloor$  stands for the largest integer smaller than or equal to  $x$ ). Best-effort (BE) calls arrive according to a Poisson process with rate  $\lambda_0$ . Their lengths are i.i.d. exponentially distributed with parameter  $\mu_{\text{BE}}$ . All BE calls are accepted in the system and share the bandwidth left unused by GP.

GP calls arrive as  $K$  incoming flows (streams) differentiated by acceptance thresholds  $N_K \leq \dots \leq N_1 \leq \Delta$ . A class  $k$  call is accepted into the system if the number of servers (i.e. the amount of bandwidth) already occupied by GP calls is less than  $N_k$ , and is otherwise lost.  $N_k, k = 1, \dots, K$  are thus acceptance thresholds which define a CAC. Thus the GP calls of class  $k$  receive a better service than those of class  $K + 1$  in terms of call acceptance probability: if the number of GP sessions at time  $t$  is in the interval  $[N_{k+1}, N_k)$ , an incoming call of GP class  $k$  will be accepted, but a call of class  $k + 1$  will be blocked. We assume that the GP calls which are not accepted in the system are simply lost. Once accepted in the system, each GP session occupies one unique server (i.e. one unit of the bandwidth) until service completion; the GP service times are exponential with parameter  $\mu_{\text{GP}}$  and do not depend on the GP class. A GP session in service can not be preempted by any other session, even from a higher priority class. Thus, the difference in service between classes only concern call blocking probabilities.



The largest threshold value  $N_1$  defines the maximum amount of resource occupied by GP calls. If  $N_1 = \Delta$ , then the GP service may at some times use up all the service capacity of the system. If  $N_1 < \Delta$ , then the number of servers unused by GP is at all times at least  $\Delta - N_1$ , meaning that this fraction of the service is always available for BE connections.

Class  $k$  GP sessions arrive to the system according to a Poisson process with rate  $\lambda_k$ ,  $k = 1, \dots, K$ . Denote  $\lambda_{\text{GP}} = \sum_{k=1}^K \lambda_k$ . The arrival process of BE sessions is also Poissonian with rate  $\lambda_{\text{BE}}$ . The arrival times of different GP flows as well as of BE traffic and service times are independent.

There is no call blocking for BE in our model, all incoming calls are accepted into the network. Since BE sessions share the available capacity, we model their service by a processor sharing discipline (for ATM, this is in accordance with the min-max fairness objective in the flow control for ABR traffic [1]).

The performance measures of interest concern macroscopic values in the stationary behavior. For the GP calls, we want to calculate the call rejection probability for each GP class. In our model, because the arrivals are Poisson, it follows from the PASTA property that the call blocking probability of class  $k$  is equal to the probability that the number of GP calls is larger than or equal to  $N_k$ . For the BE connections, we are mainly interested in the mean duration of a session. This represents for example the average time it takes to transfer some amount of data. Typical target values of the average connection time should be in the order of magnitude of seconds or less for file sizes of a few megabits.

Denote by  $X_1(t)$  the total number of GP sessions in the system, and by  $X_2(t)$  the total number of BE ones, at time  $t$ . GP calls of all classes have preemptive priority over BE (as long as they do not pass their respective acceptance thresholds).

The bandwidth available for BE sessions at time  $t$  is  $\Delta - X_1(t)$ . In order to obtain the average amount of workload of BE sessions, or the average number of BE sessions, we make the following observation: these *averages* remain unchanged if we consider a FIFO discipline for BE sessions (instead of a PS discipline). We shall thus study, below, the FIFO discipline instead (each BE session has length  $\mu_{\text{BE}}$ ). (Note that this transformation does alter the other moments.)

With this transformation, the dynamics of the BE traffic is that of an M/M/1 with a random environment for the service time; the instantaneous departure rate of BE sessions is  $\nu(X_1(t))$  where

$$\nu(x) := \mu_{\text{BE}}(\Delta - x). \quad (1)$$

The process  $X(t) = (X_1(t), X_2(t))$  is an irreducible Markov chain. It is ergodic if and only if the average service capacity available to the BE sessions is larger than their arrival rate  $\lambda_{\text{BE}}$ :

$$\nu_{\text{BE}} := \mu_{\text{BE}}(\Delta - \mathbb{E}X_1) > \lambda_{\text{BE}}. \quad (2)$$

We assume that this condition is satisfied. We show in Section 3.1 how to calculate the mean GP server occupancy  $\mathbb{E}X_1(t)$ . To establish the above stability condition, one may simply use the fluid limit approach. The sufficiency part of the stability condition can be obtained using the approach in [4]. The necessity of the condition follows from [6].

### 3 Numerical Approach

This section describes the numerical algorithms that we have implemented to calculate the performance statistics of the model.

#### 3.1 Distribution of the number of GP sessions

We show here how to calculate the distribution of  $X_1(t)$ , which permits to derive the call blocking probabilities of each GP class and also the mean value of  $X_1(t)$ . >From the assumption of priority of GP sessions over BE, the distribution of  $X_1(t)$  does not depend on the arrivals and departures of BE sessions, and thus we can ignore them.

When  $X_1(t) = i$  with  $N_{k+1} \leq i < N_k$  for some  $k$ , the arrival rate of GP sessions is equal to  $\lambda_1 + \dots + \lambda_k$ , as only sessions with class indice lower or equal to  $k$  are allowed into the system, and the service rate is equal to  $i\mu_{\text{GP}}$ , as exactly  $i$  servers are busy with GP sessions. Define  $\gamma_k = (\lambda_1 + \dots + \lambda_k)/\mu_{\text{GP}}$ ; then, from the balance equations of the Markov process  $X_1(t)$ , the stationary probabilities  $\pi_i = \text{P}(X_1(t) = i)$  satisfy the relation:

$$\pi_{i+1} = \pi_i \frac{\gamma_k}{i+1}. \quad (3)$$

>From the above equality, we can compute the  $(\pi_i)$  as follows. For  $N_{k+1} \leq i < N_k$ , let

$$t(i) = \frac{\gamma_K^{N_K} \gamma_{K-1}^{N_{K-1}-N_K} \dots \gamma_{k+1}^{N_{k+1}-N_{k+2}} \gamma_k^{i-N_{k+1}}}{i!},$$

and for the terms with  $k = K$ , we use the convention  $N_{K+1} = 0$ . >From (3), we have  $\pi_i = t(i)\pi_0$ . We further define

$$S_k = \sum_{N_{k+1} \leq i < N_k} t(i) \quad (4)$$

$$S = t(N_1) + \sum_{1 \leq k < K} S_k. \quad (5)$$

We have then  $\pi_i = t(i)/S$  for all  $i$  in  $[0, N_1]$ , and the blocking probability for class  $k$  sessions is the probability that  $X_1(t)$  is greater or equal to  $N_k$ , which is equal to  $1 - P(X_1(t) < N_k)$ , with

$$P(X_1(t) < N_k) = (S_K + S_{K-1} + \cdots + S_k)/S. \quad (6)$$

We also need a formula for  $\mathbb{E}X_1(t)$ , the mean number of servers occupied by GP sessions:

$$\begin{aligned} \mathbb{E}X_1(t) &= \sum_{1 \leq i \leq N_1} it(i) \\ &= \sum_{1 \leq k \leq K} \sum_{N_{k+1} < i \leq N_k} it(i) \\ &= \sum_{1 \leq k \leq K} \sum_{N_{k+1} < i \leq N_k} \gamma_k t(i-1) \\ &= \sum_{1 \leq k \leq K} \gamma_k S_k. \end{aligned} \quad (7)$$

The relations (6) and (7) above provide the means to calculate the performance values for GP traffic that we need below.

### 3.2 Distribution of the waiting time of BE sessions

In this section, we use the results on matrix-geometric distributions presented by M. F. Neuts in [9], and in particular Theorem 1.7.1 on page 32. We calculate the distribution  $\pi$  of the process  $(X_1(t), X_2(t))$  and derive from there the average waiting time of BE sessions.

The pair  $X(t) = (X_1(t), X_2(t))$  is a quasi birth and death process whose infinitesimal generator  $Q$  can be written in the following block matrix form:

$$Q = \begin{pmatrix} B & A_0 & 0 & 0 & \cdots \\ A_2 & A_1 & A_0 & 0 & \cdots \\ 0 & A_2 & A_1 & A_0 & \cdots \\ 0 & 0 & A_2 & A_1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \cdots \\ \vdots & \vdots & \vdots & \vdots & \cdots \end{pmatrix}, \quad (8)$$

In the above,  $B$ ,  $A_0$ ,  $A_1$ , and  $A_2$  are square matrices of dimension  $N_1 + 1$  such that a non-diagonal  $(n, n')$ -entry of  $Q$ , with  $n = i + (N_1 + 1)j$ ,  $n' = i' + (N_1 + 1)j'$ , and  $0 \leq n, n' \leq N_1$ , is the transition rate of the Markov process  $X(t)$  from state  $(i, j)$  to state  $(i', j')$ . Note that from the independence assumptions, the probability that two events (arrival and service completion

in the different service classes) occur simultaneously is zero. The matrix  $A_0$  corresponds to an arrival in the BE queue; this event occurs with rate  $\lambda_0$ , thus we have:

$$A_0 = \text{diag}(\lambda_0; 0 \leq i \leq N_1) \quad (9)$$

Similarly, the matrix  $A_2$  corresponds to a departure in the BE queue, which happens with rate  $\nu(X_1(t)) = (\Delta - X_1(t))\mu_{\text{BE}}$ , thus

$$A_2 = \text{diag}((\Delta - i)\mu_{\text{BE}}; 0 \leq i \leq N_1) \quad (10)$$

In the matrix  $A_1$ , the off-diagonal elements are the transition rates of the Markov process  $X_1(t)$ , and the diagonal elements are such that the elements in each row of  $Q$  sum to zero. Because no more than one GP session can leave or arrive in the system at a time, the matrix  $A_1$  is tridiagonal, containing the following elements, with  $i$  varying between the appropriate bounds.

$$A_1[i][i+1] = \sum_{1 \leq k \leq K} \lambda_k \mathbf{1}(i < N_k) \quad (11)$$

$$A_1[i][i-1] = i\mu_{\text{GP}} \quad (12)$$

$$A_1[i][i] = -\lambda_0 - (\Delta - i)\mu_{\text{GE}} - \sum_{1 \leq k \leq K} \lambda_k \mathbf{1}(i < N_k) - i\mu_{\text{GP}} \quad (13)$$

We write  $A_1 = T - D$ , where  $D$  is a diagonal matrix and  $T$  has a zero diagonal. Note that the diagonal matrix  $D$  is positive and invertible, and that  $T$  is non-negative and has the same off-diagonal elements as the infinitesimal generator of  $X_1(t)$ . The matrix  $B$  is identical to  $A_1$  except for the diagonal, which is such that the diagonal of  $Q$  is zero. We can write  $B = T - D_B$ , with

$$D_B = \text{diag} \left( \lambda_0 + \sum_{1 \leq k \leq K} \lambda_k \mathbf{1}(i < N_k); 0 \leq i \leq N_1 \right) \quad (14)$$

The balance equations for the two-dimensional chain can be written  $\pi Q = 0$ . Let  $\pi = [\pi(0), \pi(1), \dots, \pi(j), \dots]$ , with  $\pi(j) = [\pi(0, j), \pi(1, j), \dots, \pi(N_1, j)]$ . We further define

$$\pi_i = \sum_{j=0}^{\infty} \pi(i, j) \quad (15)$$

the stationary probability of having  $i$  GP sessions in the system. We have then:

$$\pi(0)B + \pi(1)A_2 = 0 \quad (16)$$

$$\pi(j)A_0 + \pi(j+1)A_1 + \pi(j+2)A_2 = 0 \quad (17)$$

We now refer to the work on matrix-geometric Markov processes presented by M. Neuts in [9]. From the stability condition (2) above, we know that when  $\nu_{\text{BE}} < 1$ , the steady-state distribution  $\pi$  exists, and from [9], it is characterized by

$$\pi(j) = \pi(0)R^j, \quad (18)$$

where the matrix  $R$  is the minimum non-negative solution of the equation

$$A_0 + RA_1 + R^2A_2 = 0. \quad (19)$$

The above equation is equivalent to  $R = (A_0 + RT + R^2A_2)D^{-1}$ . To calculate the matrix  $R$ , we use the following iterative scheme:

$$R_0 = 0 \quad (20)$$

$$R_{n+1} = (A_0 + R_nT + R_n^2A_2)D^{-1}. \quad (21)$$

We know from [9] and from the stability condition (2) that there exists a finite non-negative matrix which is a solution to equation (19). It can be shown without difficulty that the matrix sequence  $(R_n)_n$  for  $n \in \mathbb{N}$  is non-decreasing and converges to the matrix  $R$ . This is the algorithm that has been used to obtain the matrix  $R$ , with a precision equal to the smallest positive number known by our computer.

Once the matrix  $R$  is known, we get the distribution of  $X(t)$  by performing its spectral analysis, that is by calculating its  $N_1 + 1$  eigenvalues and eigenvectors. This part of the numerical analysis has been implemented with the help of the library *meschach*, a freeware package in C language for linear algebra (see reference manual [8]). Let  $(r_i)_{0 \leq i \leq N_1}$  be the eigenvalues of  $R$ . The boundary condition (16) of  $\pi$  translates into a linear system of equations that is also solved with the help of the *meschach* libraries and yields a vector  $a = [a_0, a_1, \dots, a_{N_1}]$  such that for all  $j$  in  $\mathbb{N}$

$$P(X_2(t) = j) = \sum_{0 \leq i \leq N_1} a_i (r_i)^j, \quad (22)$$

and

$$\mathbb{E}X_2(t) = \sum_{0 \leq i \leq N_1} a_i \frac{r_i}{1 - r_i^2}. \quad (23)$$

>From the mean value of  $X_2(t)$ , we finally derive the expectation of the sojourn time of a BE session through the Little's formula. This mean sojourn time, which we denote by  $T_{\text{BE}}$ , is given by  $T_{\text{BE}} = \mathbb{E}X_2(t)/\lambda_{\text{BE}}$ .

## 4 Experiment Results

This section presents and comments some numerical experiments on the performances of the model, in different situation of the traffic load: first the light traffic situation, where the total input load of GP and BE traffic is small compared to the service capacity of the station; second, the heavy traffic situation, where the total load is less than but close to the service capacity; and finally the saturating traffic situation where the total load is more than the service capacity.

In each situation (light, heavy, or saturating traffic), we present the numerical values obtained for the mean BE call duration  $T_{\text{BE}}$  and GP call blocking probability as a function of the GP acceptance threshold  $N_1$  varying from 0 to  $\Delta$ .

In each situation, we also show different curves to demonstrate the impact of GP service time on BE performances. We take constant load factors  $\rho_{\text{GP}} := \lambda_{\text{GP}}/\mu_{\text{GP}}$  and thus constant mean occupation level of servers by GP sessions, and let the mean duration  $\mu_{\text{GP}}^{-1}$  of GP sessions vary. We thus investigate the problem of of time scales, also investigated by Nunez-Queija and Boxma [10] (who restrict their model and analysis to a single GP class and to  $N_1 = \Delta$ , i.e. no bandwidth reservation for BE traffic. The results in [10] were obtained independently to ours, which have been already described in [11]). As in [10], what we observe is that for a constant GP load factor, the BE mean duration is lowest if the GP service time is lowest. If the GP call duration is very small compared to the BE connection time, then seen from an BE call, the bandwidth available can be considered constant and equal to its mean value  $\Delta - \mathbb{E}X_1(t)$ . When the mean GP service time increases, then the mean BE call duration  $\mathbb{E}T$  also increases to potentially high values, depending on the BE load and the threshold level.

Thus, in order to keep the BE service with a relatively low call duration, there is an important relation between the target value for the mean BE call duration, the mean GP service time, and the GP acceptance thresholds. For example, the sensitivity of BE mean connection to GP service time duration can be important when the traffic load is high and when the BE delays are large; in light traffic conditions, however, this sensitivity is in most cases not significant.

In the following subsections, we investigate the dependency of BE performances on GP service times and GP threshold values. In all figures, the total available bandwidth  $\Delta$  (number of service units) is 20, the BE service rate  $\mu_{\text{BE}}$  is 0.2, and the values for the GP service rate  $\mu_{\text{GP}}$  vary from 0.001 to 10.0. Each curve is a function of the GP threshold  $N_1$ . In all cases, the lowest BE call time is for  $\mu_{\text{GP}} = 10$  and the highest BE call time for  $\mu_{\text{GP}} = 0.001$ . The GP call blocking probability is independent of  $\mu_{\text{GP}}$  as the GP load factor is kept constant.

## 4.1 Light Traffic

In light traffic situations, when varying the threshold values of GP, the variations in blocking probabilities are not very important, but the variations in the BE mean call duration are even less significant. Thus, in such cases, it can be advocated not to set any threshold on GP call acceptance, at least for the class  $k = 1$ . As expected, in this situation, the performances of both services can be made very good at the same time. This is illustrated on the Figures 1 and 2, where the arrival rates are chosen such that  $\lambda_{GP}/\mu_{GP} = 4$  and  $\lambda_{BE}/\mu_{BE} = 5$ .

In summary, we obtain good performances for both services, and low sensitivity of the performance of BE sessions to the mean duration of GP services and to the thresholds of GP traffic.

## 4.2 Heavy Traffic

In heavy traffic situations, there is some degradation in the performances of the services, and a there is a tradeoff to be found for the threshold  $N_1$  if one does not want to see the BE service quality suffer too much from the high traffic load. In this case, it makes sense to set a threshold on GP acceptance, as this can greatly reduce the BE mean delays while not degrading too much the GP blocking probability. This is illustrated in Figures 3 and 4, where the arrival rates are chosen such that  $\lambda_{GP}/\mu_{GP} = 9$  and  $\lambda_{BE}/\mu_{BE} = 9$ .

Further, we can see that when the BE mean duration is high, the sensitivity to the mean duration of a GP session is also very high. In summary, there is room for some performance tradeoff between the traffic classes; high sensitivity of BE's performance to the duration of GP sessions.

When the GP and BE services are heavily loaded, it can be a good choice to define 2 distinct GP classes with different acceptance thresholds to make sure that at least one GP service has a low blocking probability. This way, the performance tradeoff is among three classes of service: one GP class with low blocking probability, one GP class with relatively high blocking probability, and one BE class with acceptable service capacity.

Performances measures with two different GP classes are shown in Figures 5 and 6. The total bandwidth  $\Delta$  is still 20, we take  $N_1 = \Delta$  and let  $N_2$  vary from 0 to  $N_1$ . The different curves on Fig. 5 are as above for different values of  $\mu_{GP}$  (0.001, 0.01, 0.1, 10.0). Other parameters are  $\text{load}(\text{GP1})=4$ ,  $\text{load}(\text{GP2})=5$ ,  $\text{load}(\text{BE})=9$ , where we call load the ratio of the arrival rate of a stream over the service rate  $\mu_{GP}$  or  $\mu_{BE}$ . Note that the BE load and the global GP load are the same as in the previous example. We can see that in this case it is possible to maintain a fairly good BE service (mean duration in the order of seconds) while keeping good performance for GP1 service.

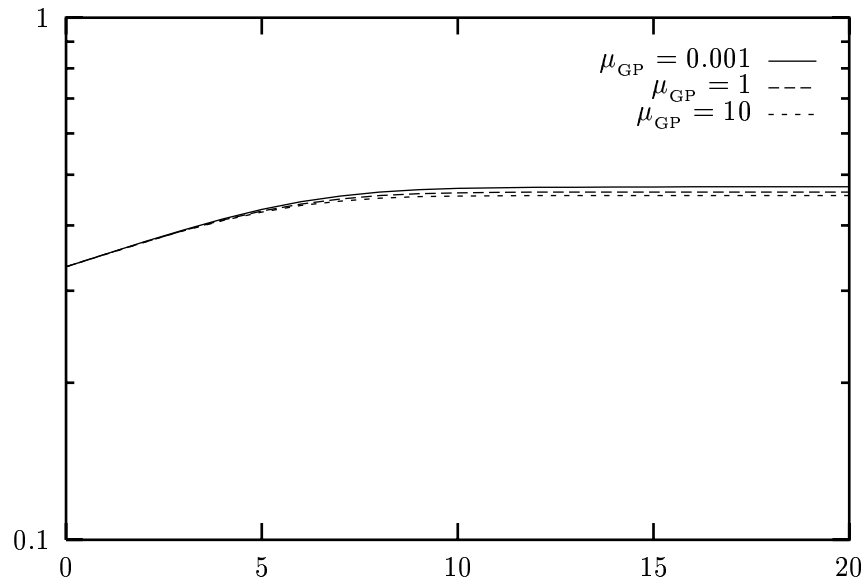


Figure 1: Mean BE call duration in light traffic.

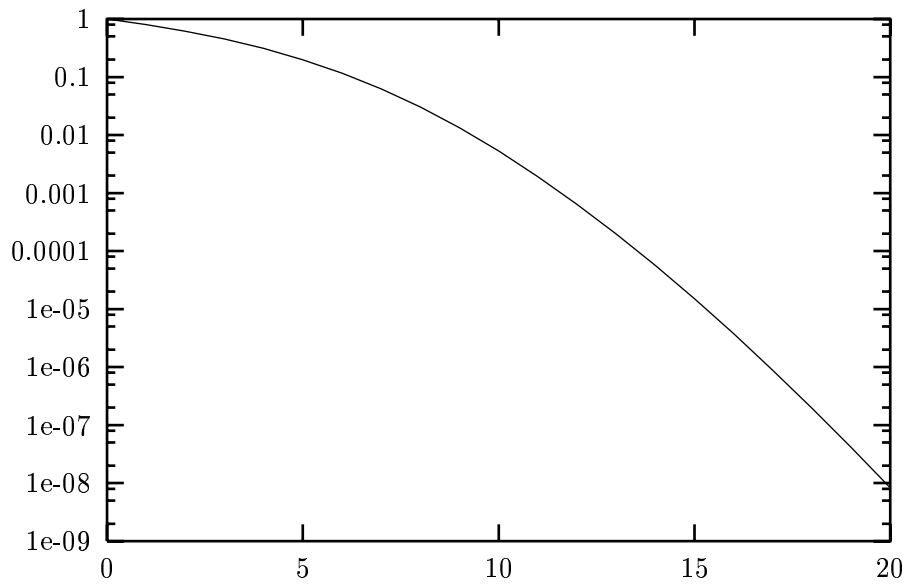


Figure 2: Blocking probability of GP sessions in light traffic.



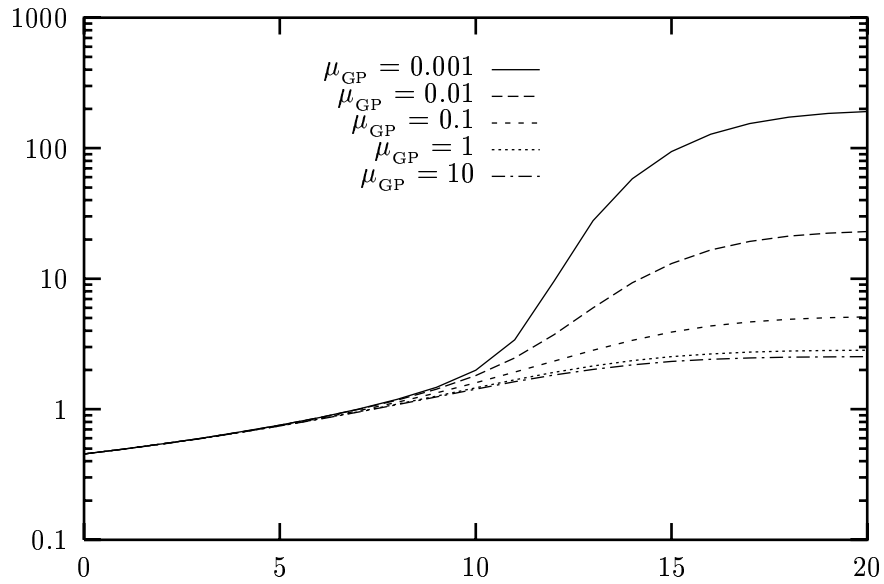


Figure 3: Mean BE call duration in heavy traffic.

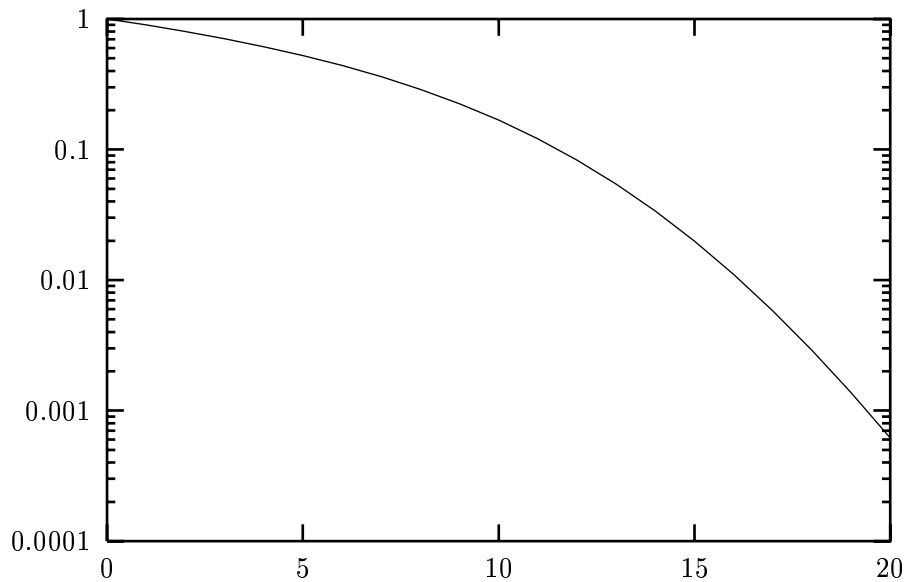


Figure 4: Blocking probability of GP sessions in heavy traffic.

### 4.3 Saturated Traffic

An important problem arises: the stability of the BE queue. When the total traffic load exceeds the service capacity of the  $\Delta$  servers, the choice of GP rejection thresholds can greatly affect the connection time of BE, up to the point of instability (infinite connection time). The latter occurs when  $\mu_{\text{BE}}(\Delta - \mathbb{E}X_1) < \lambda_{\text{BE}}$ . Thus, in some cases, the BE service is unstable if there is no acceptance threshold on GP, but can be made stable by setting a threshold. In saturated traffic, setting a low acceptance threshold on a given GP class, even though it will allow a stable BE service, will of course degrade the acceptance probability of this GP class. But in most cases it is possible to set a low threshold on some GP class, which will correspond to a service with a high blocking probability, while setting a higher acceptance threshold for some other GP class and keeping a very low blocking probability for this class.

This situation is exemplified in Fig. 7 and 8, where the arrival rates are chosen such that  $\lambda_{\text{GP}}/\mu_{\text{GP}} = 10$  and  $\lambda_{\text{BE}}/\mu_{\text{BE}} = 11$ . In summary: at least one of the traffic classes is congested, and there is high BE sensitivity to the mean duration of GP sessions.

Note that in our numerical example, the mean BE call time in steady state is infinite if the threshold  $N_1$  is more than 12.

## 5 The pricing

The issues raised by the provision of multiple classes of service over the same network infrastructure includes the objective for the service provider to offer a fair access and quality of service for all classes of services. The BE service should be allocated the network resources that are not used by GP calls.

In ATM networks, CBR and VBR services are expected to be charged at a much higher rate than ABR and UBR services. From the point of view of users who can not afford GP calls prices, a reasonable service availability and performance for BE services is important since usually the offering of telecommunication services is a multi-part tariff with at least one component which is fixed (does not depend on bandwidth, durations of calls or the volume of traffic). A BE service that provides too bad performances to be used may pose some legal issues to the service provider since users are still paying the fixed part of the tariff (e.g., the subscription fees). Since the emerging pricing are set for bundled services (e.g., wireless, internet and ATM), the service provider may also loose a large number of users by side effect.

Therefore, the first step is to establish a resource allocation policy between the different classes of services in order to guarantee the availability of a minimum amount of resources for BE services. This first step results in an optimization of the distribution of resources according to the minimum amount of resources reserved for BE services. Since the BE service does not

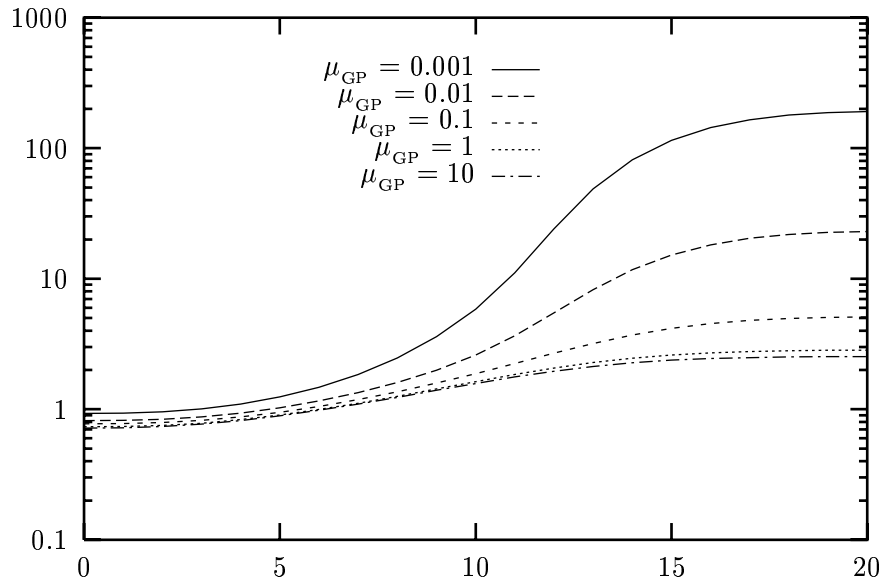


Figure 5: Mean BE call duration in heavy traffic, with two GP classes.

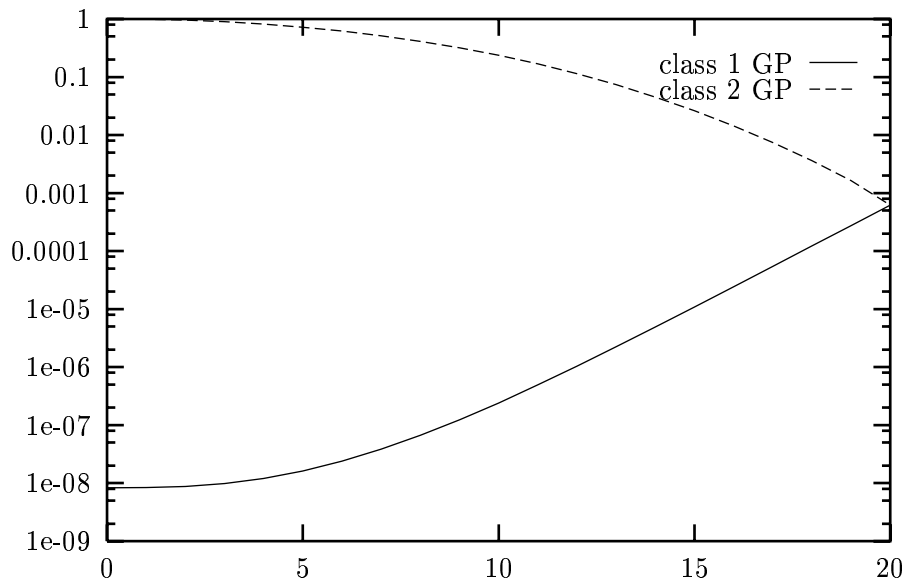


Figure 6: GP1 and GP2 call blocking probabilities in heavy traffic.

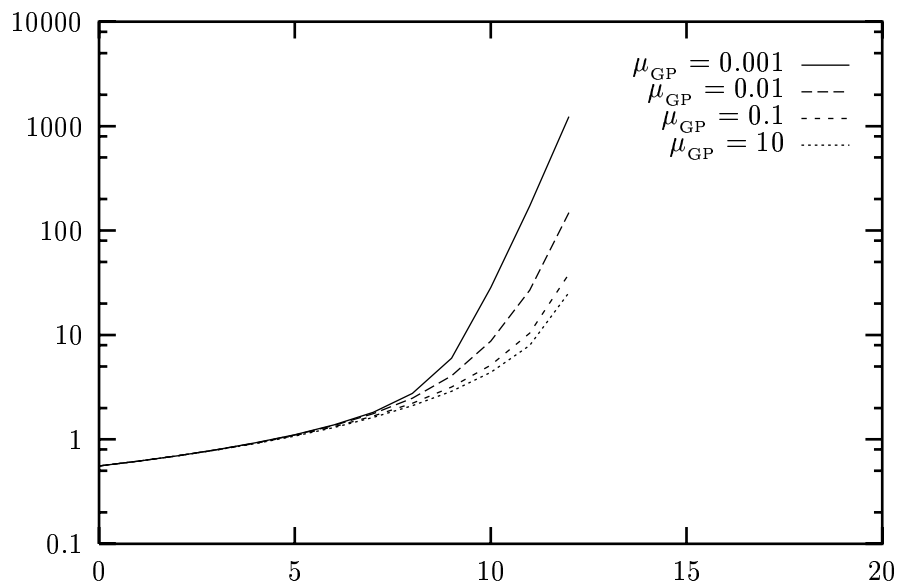


Figure 7: Mean BE call duration in saturated traffic.

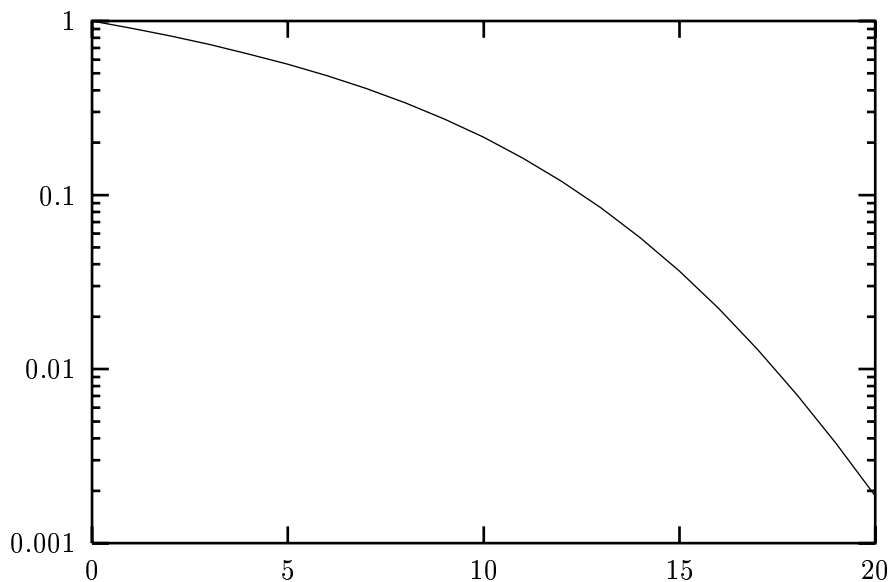


Figure 8: Blocking probability of GP sessions in saturated traffic.

implement a strict Call Acceptance Control mechanism, the resource allocation policy does not guarantee a reasonable service quality. For the BE service the service quality measured for example in terms of delay depends on the total load on this service. Economics is the last recourse left to provide reasonable service quality for BE users by influencing the demand for this service. In fact, the demand of users for a class of service is a recursive function of its service quality and cost.

The service quality tradeoffs and resource allocation objectives among GP and BE services can be efficiently reinforced in a second step by a pricing policy that encourages users not to deviate from some desirable (optimum) operation point which maximises the social welfare. We propose in this section a pricing for the model described previously for determining tariffs at the network optimal operation point.

We consider three categories of sessions: the first one is the pure GP sessions that characterise applications that require hard performance guarantees. The second category is the pure BE sessions that correspond to elastic applications and for which users can not afford (or have no reason to prefer) the high prices of the GP service. The third category encompasses sessions, called mixed sessions, for which users can purchase both GP and BE services.

Our goal is to analyse the impact of the pricing on the last category, i.e. on those sessions that may use both GP or BE services for their sessions; the pricing will add to the already existing trades off between the call blocking probability of the GP sessions and the average waiting time of BE sessions.

As is the case for any GP session, when sessions of the third category decide to use a GP service, they are subject to a Call Admission Control; we assume that calls of this category are accepted if and only if the total number of ongoing GP sessions (including the first category) is no larger than  $N_2$ . Sessions of the first category are accepted as long as the total of number of GP ongoing sessions is no larger than  $N_1$ . We take  $\Delta = N_1 \geq N_2$ ; our goal will be to design an efficient combined CAC (i.e. to determine  $N_2$ ) together with a pricing mechanism.

The total flow that the third (mixed) category submits to the network is constant and denoted  $\lambda_{\text{GP-BE}}$ . The total flows submitted by the pure GP sessions and pure BE sessions are respectively denoted  $\lambda_{\text{GP}}$  and  $\lambda_{\text{BE}}$ . Let  $\alpha \in [0, 1]$  be the fraction of mixed sessions that use GP service.

We assume that the objective of the network manager is to maximize over the possible distributions  $\alpha$  and the CAC threshold  $N_2$ , the difference  $J(\alpha, N_2)$  between the *global value* (for all ongoing sessions) and *the cost* generated by these sessions for the network (this will be made precise below).

We assume that the value of the GP service to a user is reduced with the increase of the call blocking probability. Similarly, we assume that the value of the BE service to a user decreases with the increase of the average waiting time in the network.

The total flow of GP sessions (including those originating from the mixed sessions) is:

$$\lambda_{\text{GP}} + \alpha\lambda_{\text{GP-BE}} \quad (24)$$

and the total flow submitted to the BE service is:

$$\lambda_{\text{BE}} + (1 - \alpha)\lambda_{\text{GP-BE}}. \quad (25)$$

We denote by  $u_{\text{BE}}(\alpha, N_2)$ ,  $u_{\text{GP}}^1(\alpha, N_2)$  and  $u_{\text{GP}}^{\text{mixed}}(\alpha, N_2)$  the values of the services offered to a BE session, a GP session of the first category, and GP session of the third category (i.e. originating from the mixed sessions), respectively, as function of  $\alpha$  and  $N_2$ . We assume that these functions are piecewise continuous in  $\alpha$ . The call blocking probability levels for sessions of category 1 and 3, i.e. those using GP, are monotone increasing in the total flow of these sessions, and thus in  $\alpha$ . Note that these call blocking probabilities are different for traffic of categories 1 and 3, although both use GP, since their corresponding thresholds are different.  $\alpha$  has also an impact on the mean waiting time of the BE service.

As already illustrated in Section 4.3, it may not always be desirable to allocate the network resources to the most valuable calls because in some cases the global value to users of a service can increase considerably with a small decrease of the global value for the other services.

The design of the combined CAC and pricing is performed in two steps.

### (i) Determining the social optimum value

First, the network manager finds the following social optimum over  $\alpha$  and  $N_2$  of  $J(\alpha, N_2)$ , where

$$\begin{aligned} J(\alpha, N_2) &:= (\lambda_{\text{BE}} + (1 - \alpha)\lambda_{\text{GP-BE}})u_{\text{BE}}(\alpha, N_2) - c_{\text{BE}}(\alpha) \\ &\quad + \lambda_{\text{GP}}u_{\text{GP}}^1(\alpha, N_2) + \alpha\lambda_{\text{GP-BE}}u_{\text{GP}}^{\text{mixed}}(\alpha, N_2) - c_{\text{GP}}(\alpha) \end{aligned}$$

$c_{\text{BE}}(\alpha, N_2)$  and  $c_{\text{GP}}(\alpha, N_2)$  are the costs to the network of handling respectively the GP and the BE sessions.

Let  $\alpha^*$  and  $N_2^*$  denote some optimal value of  $\alpha$  and  $N_2$ , respectively.

### (ii) Determining the prices

Our next goal is to design a pricing mechanism that would induce the mixed users to send indeed a fraction  $\alpha^*$  of their global inflow of sessions to the GP service, given that the CAC for these is set at level  $N_2^*$ .

In order to do that, we view this problem as a stackelberg game. The leading player is the network manager that determines the prices. The users having mixed sessions are the followers, and react to the prices. We describe these users as a *continuous* number of players, each sending a small fraction  $d\alpha$  of overall mixed sessions.

Let  $\Pi_{\text{BE}}(\alpha)$  and  $\Pi_{\text{GP}}(\alpha)$  be the price per session of the 3rd category (mixed sessions) that use BE and a GP, respectively, as a function of  $\alpha$ . We shall assume that they are piecewise continuous in  $\alpha$ . The average utility per mixed session when all users send a fraction  $\alpha$  of their mixed sessions to GP, is

$$V(\alpha) := \alpha[u_{\text{GP}}^{\text{mixed}}(\alpha, N_2^*) - \Pi_{\text{GP}}(\alpha)] + (1 - \alpha)[u_{\text{BE}}(\alpha, N_2^*) - \Pi_{\text{BE}}(\alpha)].$$

Assume that mixed sessions are sent according to the policy (proportion) of  $\alpha$ . Assume that a fraction  $\epsilon$  of the users deviates from  $\alpha$ , and send a fraction of  $\beta$  of their mixed sessions to GP. This results in a new overall proportion of

$$\gamma = (1 - \epsilon)\alpha + \epsilon\beta$$

of mixed sessions that are sent to GP. The average utility per session for the users that deviated to the proportion of  $\beta$  is

$$V(\alpha; [\epsilon, \beta]) := \beta[u_{\text{GP}}^{\text{mixed}}(\gamma) - \Pi_{\text{GP}}(\gamma)] + (1 - \beta)[u_{\text{BE}}(\gamma) - \Pi_{\text{BE}}(\gamma)].$$

**Definition 5.1** (i)  $\Pi_{\text{BE}}(\alpha)$  and  $\Pi_{\text{GP}}(\alpha)$  are said to be weakly stable at  $\alpha^*$  if

$$\delta(\alpha^*) = 0, \tag{26}$$

where

$$\delta(\alpha) := [u_{\text{BE}}(\alpha) - \Pi_{\text{BE}}(\alpha)] - [u_{\text{GP}}^{\text{mixed}}(\alpha) - \Pi_{\text{GP}}(\alpha)].$$

(ii)  $\Pi_{\text{BE}}(\alpha)$  and  $\Pi_{\text{GP}}(\alpha)$  are said to be stable at an  $\epsilon$ -neighborhood of  $\alpha^*$  if the following holds: Given that all users of mixed sessions use the policy of sending a portion of  $\alpha^*$  of their sessions to GP and the rest to BE, then no fraction  $0 < \epsilon_0 \leq \epsilon$  of these users can benefit by deviating from this policy. In other words, for all  $\epsilon_0 \leq \epsilon$  and all  $\beta$ ,

$$V(\alpha^*; [\epsilon_0, \beta]) \leq V(\alpha^*).$$

(iii)  $\Pi_{\text{BE}}(\alpha)$  and  $\Pi_{\text{GP}}(\alpha)$  are said to be globally stable at  $\alpha^*$  if they are stable at  $\epsilon$ -neighborhood of  $\alpha^*$  for all  $\epsilon$ .

(iv)  $\Pi_{\text{BE}}(\alpha)$  and  $\Pi_{\text{GP}}(\alpha)$  are said to be  $\nu$ -stable at an  $\epsilon$ -neighborhood of  $\alpha^*$  if the following

holds: Given that all users of mixed sessions use the policy of sending a portion of  $\alpha^*$  of their sessions to GP and the rest to BE, then no fraction  $0 < \epsilon_0 \leq \epsilon$  of these users can benefit more than  $\nu$  by deviating from this policy. In other words, for all  $\epsilon_0 \leq \epsilon$  and all  $\beta$ ,

$$V(\alpha^*; [\epsilon_0, \beta]) \leq V(\alpha^*) + \nu.$$

In order to induce the mixed sessions to send a fraction  $\alpha^*$  of their sessions as GP ones, the prices need to be stable at an  $\epsilon$ -neighborhood of  $\alpha^*$ .

**Remark 5.1** (i) *Weak stability implies that an “infinitesimal” user of mixed sessions, i.e. a user that sends an overall fraction  $d\alpha$  of mixed sessions, is indifferent between sending them as GP or as BE sessions. Hence, if all users are infinitesimal, no user has any incentive to deviate from  $\alpha^*$ . This call this a weak stability since it might be the case that if a whole positive fraction of the sessions deviate from the proportions  $\alpha^*$ , then they do benefit from that deviation.*

(ii) *One can show that if  $V(\alpha; [\epsilon, \beta]) < V(\alpha^*)$  for some  $\epsilon$  and  $\beta$ , then it also holds for some  $\epsilon' \leq \epsilon$  and for either  $\beta = 0$  or  $\beta = 1$ . Hence, one may replace “for all  $\beta$ ”, in the above definition (parts ii-iii), with “for  $\beta = 0, 1$ ”.*

**Lemma 5.1** *A necessary condition on the prices  $\Pi$  to be stable at some  $\epsilon$ -neighborhood of  $\alpha^*$  is that it is weakly stable. and  $u_{\text{BE}}(\alpha, N_2)$  is continuous in  $\alpha$  at  $\alpha^*$ .*

**Proof:** Assume that a fraction  $\alpha^*$  of the mixed sessions are sent as GP ones. If  $\delta(\alpha^*) > 0$  then some fraction of users that send an overall of  $d\alpha$  of the mixed sessions will benefit mostly by sending *all their mixed sessions* as GP ones. Similarly, if  $\delta(\alpha^*) < 0$  then they will benefit mostly by sending *all their mixed sessions* as BE ones. ■

It is not difficult to show that the above condition need not be sufficient.

Motivated by Lemma 5.1, we propose the following simple pricing mechanism:

## 5.1 Constant pricing scheme

The price per session type is constant:

$$\Pi_{\text{BE}}(\alpha) = q_{\text{BE}} \text{ and } \Pi_{\text{GP}}(\alpha) = q_{\text{GP}}.$$

Hence a user that sends  $d\alpha$  of the mixed sessions, of which a fraction  $\beta$  is GP, is charged  $(\beta q_{\text{GP}} + (1 - \beta)q_{\text{GE}})d\alpha$ .

From Lemma 5.1, it follows that this pricing should satisfy:

$$\delta^* := q_{\text{GP}} - q_{\text{GE}} = u_{\text{GP}}^{\text{mixed}}(\alpha^*) - u_{\text{BE}}^{\text{mixed}}(\alpha^*). \tag{27}$$



We obtain the following:

**Theorem 5.1** *Consider the following simple pricing:*

$$\begin{aligned} q_{\text{BE}} &= Q, \\ q_{\text{GE}} &= Q + u_{\text{GP}}(\alpha^*) - u_{\text{BE}}^{\text{mixed}}(\alpha^*), \end{aligned}$$

where  $Q$  is an arbitrary constant chosen by the network so as to obtain a desirable profit. Assume that  $u_{\text{BE}}(\alpha, N_2^*)$  is decreasing in  $\alpha$  at some interval  $[\alpha_1, \alpha^*]$  and continuous at  $\alpha^*$ . Then the above pricing is stable at some  $\epsilon$ -neighborhood of  $\alpha^*$ . If  $\alpha_1 = 0$  then the pricing scheme is globally stable.

**Proof:** We check for  $\beta = 0, 1$ , as follows from Remark 5.1. We have under the above pricing:

$$\begin{aligned} V(\alpha^*; [\epsilon, 1]) - V(\alpha^*) &= [u_{\text{GP}}^{\text{mixed}}(\gamma, N_2^*) - q_{\text{PG}}] \\ &\quad - \alpha^* [u_{\text{GP}}^{\text{mixed}}(\alpha^*, N_2^*) - q_{\text{GP}}] - (1 - \alpha^*) [u_{\text{BE}}(\alpha^*, N_2^*) - q_{\text{BE}}] \\ &= u_{\text{GP}}^{\text{mixed}}(\gamma, N_2^*) - q_{\text{PG}} - [u_{\text{BE}}(\alpha^*, N_2^*) - q_{\text{BE}}] \\ &= u_{\text{GP}}^{\text{mixed}}(\gamma, N_2^*) - u_{\text{GP}}^{\text{mixed}}(\alpha^*, N_2^*) \leq 0, \end{aligned}$$

where  $\gamma = (1 - \epsilon)\alpha^* + \epsilon$ . We made twice use above of the relation (27). The last inequality follows from the fact that  $\gamma \geq \alpha^*$ , and since the value of GP decreases in  $\alpha$  (as rejection probabilities of GP increase).

For  $\beta = 0$  we have:

$$\begin{aligned} V(\alpha^*; [\epsilon, 0]) - V(\alpha^*) &= [u_{\text{BE}}(\gamma, N_2^*) - q_{\text{BE}}] \\ &\quad - \alpha^* [u_{\text{GP}}^{\text{mixed}}(\alpha^*, N_2^*) - q_{\text{GP}}] - (1 - \alpha^*) [u_{\text{BE}}(\alpha^*, N_2^*) - q_{\text{BE}}] \\ &= u_{\text{BE}}(\gamma, N_2^*) - u_{\text{BE}}(\alpha^*, N_2^*) \leq 0, \end{aligned}$$

where  $\gamma = (1 - \epsilon)\alpha^* \leq \alpha^*$ . The last inequality follows from the hypothesis on the monotonicity of  $u_{\text{BE}}$  in  $\alpha$ . ■

The monotonicity assumption used in Theorem 5.1 often holds, in which case we may use the above constant pricing scheme. We next propose an alternative dynamic pricing scheme for the case that the monotonicity assumption does not hold.

## 5.2 Load-dependent pricing scheme

We introduce the following price per session type:

$$\Pi_{\text{BE}}(\alpha) = q_{\text{BE}} + \theta I\{\alpha \leq \alpha^* - \zeta\} \text{ and } \Pi_{\text{GP}}(\alpha) = q_{\text{GP}}, \quad (28)$$

where  $\zeta$  and  $\theta$  are some positive constants.

Hence a user that sends  $d\alpha$  of the mixed sessions, of which a fraction  $\beta$  is GP, is charged again  $(\beta q_{\text{GP}} + (1 - \beta)q_{\text{GE}})d\alpha$ , as long as  $\alpha > \alpha^* - \zeta$ . If the latter holds, it is charged an extra overload charge of  $\theta$  per unit of BE session that is sent, which amounts in an extra price of  $(1 - \beta)\theta d\alpha$ .

We may apply again Lemma 5.1 to conclude that this pricing should satisfy (27). We obtain the following:

**Theorem 5.2** *Consider the pricing scheme given in (28) with*

$$\begin{aligned} q_{\text{BE}} &= Q, \\ q_{\text{GE}} &= Q + u_{\text{GP}}(\alpha^*) - u_{\text{BE}}^{\text{mixed}}(\alpha^*), \end{aligned}$$

where  $Q$  is an arbitrary constant chosen by the network so as to obtain a desirable profit.

Assume that  $u_{\text{BE}}(\alpha, N_2)$  is continuous in  $\alpha$  at  $\alpha^*$ . Then for any  $\nu > 0$ , and  $\epsilon > 0$ , one can choose positive constants  $\theta$  and  $\zeta$  such that the above pricing is  $\nu$ -stable at the  $\epsilon$ -neighborhood of  $\alpha^*$ .

**Proof:** As in the proof of Theorem 5.1, it follows that  $V(\alpha^*; [\epsilon, \beta]) - V(\alpha^*) \geq 0$  for  $\beta = 1$ . It remains to show that  $V(\alpha^*; [\epsilon, \beta]) - V(\alpha^*) \geq -\nu$  for  $\beta = 0$ , for all  $\epsilon$  sufficiently small.

$$\begin{aligned} V(\alpha^*; [\epsilon, 0]) - V(\alpha) &\leq [u_{\text{BE}}(\gamma, N_2^*) - q_{\text{BE}} - \theta I\{\alpha \leq \alpha^* - \zeta\}] \\ &\quad - \alpha^* [u_{\text{GP}}^{\text{mixed}}(\alpha^*, N_2^*) - q_{\text{GP}}] - (1 - \alpha^*) [u_{\text{BE}}(\alpha^*, N_2^*) - q_{\text{BE}}] \\ &= u_{\text{BE}}(\gamma, N_2^*) - \theta I\{\alpha \leq \alpha^* - \zeta\} - u_{\text{BE}}(\alpha^*, N_2^*) \leq 0, \end{aligned}$$

where  $\gamma = (1 - \epsilon)\alpha^*$ .

Since  $u_{\text{BE}}(\alpha, N_2^*)$  is continuous at  $\alpha^*$ , one may choose  $\zeta$  such that  $u_{\text{BE}}(\gamma, N_2^*) - u_{\text{BE}}(\alpha^*, N_2^*) \leq -\nu$ , for all  $\epsilon'$  smaller than some  $\epsilon_1$  (recall that  $\gamma$  is a function of  $\epsilon$ ). Thus the pricing scheme is  $\nu$ -stable at an  $\epsilon_1$ -neighborhood of  $\alpha^*$ . We now show that the  $\nu$ -stability can be extended to any  $\epsilon > \epsilon_1$  as well. Indeed, this is done by choosing

$$\theta := \sup_{\epsilon' \leq \epsilon} [u_{\text{BE}}(\gamma(\epsilon'), N_2^*) - u_{\text{BE}}(\alpha^*, N_2^*)].$$

■

**Remark 5.2** *It is clear from the proof of the above Theorem that we could have taken a price of the form*

$$\Pi_{\text{BE}}(\alpha) = q_{\text{BE}} + \theta I\{\alpha \leq \alpha^*\} \text{ and } \Pi_{\text{GP}}(\alpha) = q_{\text{GP}}, \quad (29)$$

where  $q_{\text{PG}}$  and  $q_{\text{BE}}$  are given in (29). This pricing would be stable at some  $\epsilon$ -neighborhood of  $\alpha^*$  for some choice of  $\theta > 0$ . In practice, such a pricing mechanism is not desirable and should be replaced by the weaker stability notion that appears in Theorem 5.2. The reason is the following. In practice, there is no way for the network to know the exact value of  $\alpha$ , and only some estimates (e.g. the empirical average number of sessions during some finite horizon) can be used to approximate it. Thus, even if the users all send their mixed sessions at the desirable ratio of  $\alpha^*$ , we might mistakenly charge them the extra penalty of  $\theta$  per session due to measurement limitation. By using the pricing mechanism (28), we can guarantee that the probability of charging an extra penalty of  $\theta$  (due to erroneous measurements of  $\alpha$ ), while all users split their sessions according to  $\alpha^*$ , can be made arbitrarily small (by choosing the measurement interval sufficiently long).

## References

- [1] The ATM Forum Technical Committee, *Traffic Management Specification*, Version 4.0, April 1996.
- [2] A. Birman, V. Firoiu, R. Guérin, and D. Kandlur, "Provisioning of RSVP-based Services over a Large ATM Network," Research Report RC 20250, IBM Research Division, T. J. Watson Research Center, Yorktown Heights, New York, USA, 1995.
- [3] R. Braden, L. Zhang, and S. Berson, "Resource reservation protocol (RSVP) - version 1 functional specification," Internet Draft, Internet Engineering Task Force, Nov. 1995. <ftp://ds.internic.net/internet-drafts/draft-ietf-rsvp-spec-08.txt>.
- [4] J. D. Dai and S. P. Meyn, "Stability and convergence of moments for multiclass queueing networks via fluid limit models", *IEEE Trans. Automatic Control* **40**, pp. 1889-1904, 1995.
- [5] Raj Jain, Shiv Kalyanaraman, Rohit Goyal, Sonia Fahmy, and Fang Lu, "ERICA+: Extensions to the ERICA Switch Algorithm", ATM Forum Document Number: ATM Forum/95-1346 1995.
- [6] S. P. Meyn, "Transience of multiclass queueing networks via fluid limit models", *Ann. Appl. Probab.* **5**, pp. 946-957, 1995.

- 
- [7] D. J. Mitzel, D. Estrin, S. Shenker, and L. Zhang, "An Architectural Comparison of ST-II and RSVP," in *Proceedings of the Conference on Computer Communications (IEEE Infocom)*, (Toronto, Canada), June 1994.
  - [8] D. E. Stewart. *Meschach: Matrix Computations in C*. Centre for Mathematics and its Applications, School of Mathematical Sciences, Australian National University, Canberra, 1994. CMA Proceedings #32.
  - [9] M. F. Neuts. *Matrix-geometric solutions in stochastic models : an algorithmic approach*. Johns Hopkins series in the mathematical sciences. The Johns Hopkins University press, 1981.
  - [10] R. Núñez Queija and O.J. Boxma, "Analysis of a multi-server queueing model of ABR", manuscript, 1996.
  - [11] S. Tohme and U. Yechiali, *Research proposal for a French-Israeli scientific cooperation on information highways*, May 1996.
  - [12] L. Zhang, S. Deering, D. Estrin, S. Shenker, and D. Zappala, "RSVP: a new resource ReSerVation protocol," *IEEE Network*, vol. 7, pp. 8-18, Sept. 1993.



---

Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,

615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY

Unité de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex

Unité de recherche INRIA Rhône-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN

Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex

Unité de recherche INRIA Sophia Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA ANTIPOLIS Cedex

---

Éditeur

INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)

ISSN 0249-6399