

# Measurement-based Computation of Equivalent Bandwidth Under Performance Constraints

Mario Marchese, *Senior Member, IEEE*, Maurizio Mongelli, *Member, IEEE*

**Abstract** – *Equivalent bandwidth (EqB)* is defined as the minimum service rate to be provided to a traffic buffer to guarantee a certain degree of *Quality of Service (QoS)* in terms of objective parameters (packet loss, delay, jitter). EqB techniques are usually obtained analytically for homogeneous traffic trunks, with respect to a single QoS constraint. Modern network solutions often imply the aggregation of service classes with different QoS constraints, thus generating heterogeneous trunks from the point of view of both traffic sources and QoS requirements. This situation leads to the need to develop new EqB techniques so to match heterogeneity. This paper proposes a novel measurement-based EqB technique that computes the minimum required bandwidth to be allocated to a buffer, which conveys heterogeneous traffic (both concerning the traffic sources and the QoS requirements), without using closed-form expressions. The effectiveness of the algorithm is checked through simulations.

**Keywords** – *QoS, Equivalent bandwidth, Measurement control.*

## I. INTRODUCTION

**E**quivalent bandwidth (EqB) is defined as the minimum service rate to be provided to a traffic buffer to guarantee a certain degree of *Quality of Service (QoS)* in terms of objective parameters (packet loss, delay, jitter). EqB techniques are usually obtained analytically for homogeneous traffic trunks, with respect to a single QoS constraint. Modern network solutions (e.g., satellite or WiMAX Service Access Points, edge routers of DiffServ domains) often imply the aggregation of service classes with different QoS constraints, thus generating heterogeneous trunks from the point of view of both traffic sources and QoS requirements. This situation leads to the need to develop new EqB techniques so to match heterogeneity.

## II. ASSUMPTIONS

There are  $N$  traffic classes.  $a_i(t)$  is the input rate process of the  $i$ -th traffic class and  $a(t)$  the aggregate process of all  $a_i(t)$ ,  $i = 1, \dots, N$ . Traffic is conveyed towards a single buffer, modeled through a *Stochastic Fluid Model* [1]. Each  $a_i(t)$ ,  $i = 1, \dots, N$  is supposed ergodic for now, so that a single realization is representative of the entire process. This assumption will be relaxed later. There is no knowledge of  $a_i(t)$  processes, as well as of the aggregate process  $a(t)$ . Additionally, aggregation may involve also buffering and encapsulation operations as typically done in real network nodes, so making more complex the process  $a(t)$ . The only information about  $a_i(t)$  and  $a(t)$  may be got measures. The

service rate of the buffer is  $R(t)$ .  $l_i(R(t), t)$  is the overflow rate process of the  $i$ -th traffic class, measured in [bps]. The average value of the overflow rate is defined in (1). No analytical expression for  $l_i(R(t), t)$  is supposed available. Information about its behavior is got by measures. The entire system model is reported in Fig. 1. SLA (*Service Level Agreement*) for each traffic class,  $i = 1, \dots, N$ , is composed of a *Packet Loss Probability* threshold ( $PLP_i^*$ ). It means that the amount of feasible loss rate must be limited by the process  $l_i^*(t) = PLP_i^* \cdot a_i(t)$ , measured in [bps], whose average value is contained in (2).

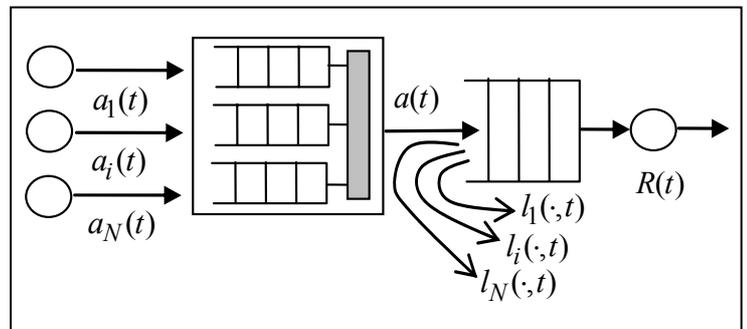


Fig. 1. System model.

$$\bar{l}_i = \lim_{\tau \rightarrow \infty} \frac{1}{\tau} \int l_i(t) dt, \quad i = 1, \dots, N \quad (1)$$

$$\bar{l}_i^* = \lim_{\tau \rightarrow \infty} \frac{1}{\tau} \int l_i^*(t) dt; \quad l_i^*(t) = PLP_i^* \cdot a_i(t); \quad i = 1, \dots, N \quad (2)$$

## III. PROBLEM DEFINITION

The aim is to provide a minimum buffer service rate so that the average loss rate  $\bar{l}_i$  equalizes  $\bar{l}_i^*$ ,  $\forall i = 1, \dots, N$ . It corresponds to define the optimization problem in (3), identified as *Generalized Equivalent Bandwidth (GEqB)*.

$$R^* = \arg \min_R I_{\Delta}(\cdot, R), \quad I_{\Delta}(\cdot, R) = \text{Max}_i \left[ \bar{l}_i - \bar{l}_i^* \right]^2 \quad (3)$$

The  $\text{Max}[\cdot]$  operator allows capturing the traffic class more “greedy” in terms of bandwidth requirements.

#### IV. PROBLEM SOLUTION

Being the involved stochastic processes unknown, GEqB problem is solved by taking measures over a given  $(k+1)$ -th observation horizon (OH),  $T_{k+1} = [t_k, t_{k+1}]$ ,  $k = 0, 1, 2, \dots$ , and performing a sequence of bandwidth reallocations  $R(t_k)$ ,  $k = 0, 1, 2, \dots$ , each  $T_k$ , based on the gradient method. The overflow rate  $l_i(R(t_k), t)$  and the feasible loss rate  $\hat{l}_i^*(t)$  processes are averaged over each OH, giving origin to the quantities  $\hat{l}_i^{(k+1)}$  in (4) and  $\hat{l}_i^{*,(k+1)}$  in (5). Being used to solve the GEqB problem,  $\hat{l}_i^{(k+1)}$  and  $\hat{l}_i^{*,(k+1)}$  must be representative of the average values  $\bar{l}_i$  and  $\bar{l}_i^*$ ,  $\forall i = 1, \dots, N$  and  $\forall k$ .

$$\hat{l}_i^{(k+1)} = \frac{1}{T_{k+1}} \int_{T_{k+1}} l_i(t) dt; \quad i = 1, \dots, N; k = 0, 1, 2, \dots \quad (4)$$

$$\hat{l}_i^{*,(k+1)} = \frac{1}{T_{k+1}} \int_{T_{k+1}} \hat{l}_i^*(t) dt; \quad \hat{l}_i^*(t) = PLP_i^* \cdot a_i(t); \quad i = 1, \dots, N; k = 0, 1, \dots \quad (5)$$

Bandwidth allocation at instant  $T_{k+1}$  is ruled by the algorithm reported in the frame below and called *Gradient-based Generalized Equivalent Bandwidth* (G<sup>2</sup>EqB) algorithm.  $step_k$  is the gradient stepsize. Condition **a)** means that the allocated bandwidth needs to be increased (i.e.,  $\hat{l}_i^{(k+1)} - \hat{l}_i^{*,(k+1)}$  for some  $i = 1, \dots, N$ ). The  $Max[\cdot]$  operator in **a)** allows exploiting the largest bandwidth need deriving from the currently most ‘‘suffering’’ traffic class. The state of sufferance is measured by cost sensitivity  $2 \cdot \frac{\partial \hat{l}_i(R)}{\partial R} \Big|_{R=R(t_k)} [\hat{l}_i^{*,(k+1)} - \hat{l}_i^{(k+1)}]$ . Condition **b)** states the opposite; the  $min[\cdot]$  operator leads to a bandwidth reduction with respect to the traffic class whose sensitivity distance from the performance constraint (i.e.,  $2 \cdot \frac{\partial \hat{l}_i(R)}{\partial R} \Big|_{R=R(t_k)} [\hat{l}_i^{(k+1)} - \hat{l}_i^{*,(k+1)}]$ ) is the smallest one. In other words, when bandwidth is increased (condition **a)**), largest bandwidth steps are followed; when bandwidth is decreased (condition **b)**), smallest steps are used.

**G<sup>2</sup>EqB algorithm**

**a)** if  $\hat{l}_i^{(k+1)} - \hat{l}_i^{*,(k+1)} \geq 0$  for at least one  $i$

$$\Delta_i(t_{k+1}) = \begin{cases} 2 \cdot \frac{\partial \hat{l}_i(R)}{\partial R} \Big|_{R=R(t_k)} [\hat{l}_i^{*,(k+1)} - \hat{l}_i^{(k+1)}], & \text{if } \hat{l}_i^{(k+1)} - \hat{l}_i^{*,(k+1)} \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\Delta(t_{k+1}) = \text{Max}_i |\Delta_i(t_{k+1})|, \quad R(t_{k+1}) = R(t_k) + step_k \cdot \Delta(t_{k+1})$$

**b)** if  $\hat{l}_i^{(k+1)} - \hat{l}_i^{*,(k+1)} < 0, \forall i$ ;  $\Delta_i(t_{k+1}) = 2 \cdot \frac{\partial \hat{l}_i(R)}{\partial R} \Big|_{R=R(t_k)} [\hat{l}_i^{(k+1)} - \hat{l}_i^{*,(k+1)}]$

$$\Delta(t_{k+1}) = \text{min}_i |\Delta_i(t_{k+1})|, \quad R(t_{k+1}) = R(t_k) - step_k \cdot \Delta(t_{k+1})$$

Derivatives  $\frac{\partial \hat{l}_i(R)}{\partial R}$  represent the sensitivity of the packet loss to infinitesimal variations of the rate serving the buffer. Intuitively they depend on the speed with which the system passes from an empty to a full state. They can be obtained [1] by observing the buffer state evolution within each OH, which is divided into  $N_k$  busy periods (where the buffer is not empty) identified by the variable  $bp$ .  $\frac{\partial \hat{l}_i(R)}{\partial R} \Big|_{R=R(t_k)}$  is approximated as in (6).  $[{}^i at_{T_{k+1}}^{bp}(R(t_k)) - {}^i ll_{T_{k+1}}^{bp}(R(t_k))]$  is the contribution to information loss of the  $i$ -th traffic class for

the busy period  $bp$  within  $T_{k+1}$ ,  $k = 0, 1, 2, \dots$ .  ${}^i at_{T_{k+1}}^{bp}$  is the arrival time of the first packet of service class  $i$  within the busy period  $bp$ .  ${}^i ll_{T_{k+1}}^{bp}$  is the time when the last loss of class  $i$  occurs during  $bp$ . If there is one traffic class, (6) is equality as proved in [2]. It is an approximation, introduced here, in case of aggregation.

$$\frac{\partial \hat{l}_i(R)}{\partial R} \Big|_{R=R(t_k)} \cong -\frac{1}{T_{k+1}} \sum_{bp=1}^{N_{T_{k+1}}} [{}^i at_{T_{k+1}}^{bp}(R(t_k)) - {}^i ll_{T_{k+1}}^{bp}(R(t_k))] \quad (6)$$

## V. ALGORITHM CONVERGENCE

Technical conditions for G<sup>2</sup>EqB convergence to global optimum ( $R(t_k) \xrightarrow{k \rightarrow \infty} R^*$ ) are: **1)** Ergodic stochastic processes, **2)** decreasing behavior of  $step_k$ , **3)** boundness of gradient within the control domain  $R(t) \in \mathfrak{R}^+$ ,  $\forall t$ , and **4)** non-existence of local optima.

**1)** and **2)** are assumptions. **3)** the lengths of buffer busy periods are bounded by OH size; measured loss rate at the end of each OH cannot be infinite. **4)** the loss rate of a traffic queue can be reasonably assumed to be continuous, differentiable, with a negative derivative in the service rate, so the cost function is also continuous, differentiable with a unique minimum.

The length of OH is very important because, on one hand, it must be long enough to assure that  $\hat{l}_i^{(k+1)}$  and  $\hat{l}_i^{*,(k+1)}$  are representative of the average values  $\bar{l}_i$  and  $\bar{l}_i^*$ ,  $\forall i=1, \dots, N$  and  $\forall k$ , but, on the other hand, it must be short to assure quick convergence. OH length is important also to get fast reactions to traffic variations.

In this context, the assumption of process ergodicity may be relaxed and limited to the time the sequence  $R(t_k)$ ,  $k=0,1,\dots$  needs to converge to  $R^*$ . When  $a(t)$  changes its statistical behavior, a new GEqB problem is solved by starting G<sup>2</sup>EqB again.

## VI. PERFORMANCE ANALYSIS AND DISCUSSION

### A. Regular EqB

Due to the complexity of the overall input rate process  $a(t)$ , equivalent bandwidth approaches which use mathematical descriptors may be hardly applied. The approach in [3] (called EqB in the following) is applicable in this context and used as a comparison. The EqB here is also the one of [4], where drives *Call Admission Control* decisions in similar heterogeneous conditions. It is chosen as performance comparison because it is the only applicable closed-form expression in heterogeneous traffic conditions.

Quantities  $m_a(t_{k+1})$  and  $\sigma_a(t_{k+1})$  are defined as the measured *mean* and *standard deviation* of  $a(t)$  over the  $(k+1)$ -th OH. Bandwidth is assigned in  $t_{k+1}$  as in (7).  $PLP_{EqB}^*$  is the upper bound on the allowed PLP and is defined as the most stringent PLP requirement out of  $N$  SLAs.

$$R(t_{k+1}) = m_a(t_{k+1}) + d \cdot \sigma_a(t_{k+1}), \quad d = \sqrt{-2 \ln(PLP_{EqB}^*) - \ln(2\pi)} \quad (7)$$

### B. G<sup>2</sup>EqB versus EqB: rate provision and convergence

VoIP SLA is considered. Each source is an on-off process. Mean on and off times durations are exponentially distributed with mean 1.008 s and 1.587 s, respectively. Peak bandwidth is 16 kbps. VoIP traffic enters an IP buffer whose length and service rate (set by the traffic peak bandwidth) guarantee no packet loss rate. IP traffic is encapsulated in ATM (via AAL5) so generating the process  $a(t)$  as output of the “Buffering and Encapsulation” box in Fig. 1.  $a(t)$  enters the ATM buffer (1600 bytes), where the loss rate in IP packets is measured.  $PLP_{VoIP}^*$  is set to  $2 \cdot 10^{-2}$ ; G<sup>2</sup>EqB OH to 30 s; Gradient stepsize to 6.0. EqB OH is either fixed to 30 s or tuned through the *Dominant Time Scale* (DTS) principle (see [4] and references therein), which computes a proper OH size to estimate the EqB statistics.

Fig. 2 and Fig. 3 show PLP and corresponding allocated bandwidth, respectively of G<sup>2</sup>EqB and EqB. The number of VoIP sources is increased of 10 from 70 to 110 each 3000 seconds. The step amounts of 10 connections to stress working conditions. Average PLP results are:  $4.40 \cdot 10^{-3}$  for G<sup>2</sup>EqB and  $1.18 \cdot 10^{-2}$  for EqB (OH size 30 s). EqB OH size DTS assures null packet loss. Average allocated bandwidths are: 0.842 Mbps for G<sup>2</sup>EqB, 0.867 Mbps for EqB OH size 30 s, and 1.31 Mbps for EqB OH size DTS. Even if the average PLP values seem to be satisfying for all the schemes, the simple observation of Fig. 2 suggests that: G<sup>2</sup>EqB reacts quickly to traffic changes also minimizing bandwidth oscillations; EqB OH size DTS always matches  $PLP_{VoIP}^*$  request but implies a relevant bandwidth waste; EqB OH size 30 s often fails to guarantee  $PLP_{VoIP}^*$  and introduces wide oscillations.

G<sup>2</sup>EqB allows tracking PLP threshold over time. Quantitative metrics may help the interpretation of this qualitative behavior. PLP standard deviation is  $7.4 \cdot 10^{-3}$  for G<sup>2</sup>EqB and  $1.33 \cdot 10^{-2}$  for EqB OH size 30 s. The percentage of the OH periods where PLP is over threshold is 5% for G<sup>2</sup>EqB and 18.6% for EqB OH size 30 s. The average difference value between measured PLP and  $PLP_{VoIP}^*$  selecting the OH periods where PLP is over threshold is  $4.22 \cdot 10^{-4}$  for G<sup>2</sup>EqB and  $2.77 \cdot 10^{-3}$  for EqB OH size 30 s.

### C. G<sup>2</sup>EqB versus EqB: heterogeneous traffic

A real video trace (“Jurassic park” from [5]) is added within the VoIP scenario. Peak and average rate are 1.418 and 0.280 Mbps. Video enters an IP buffer whose length and service rate (set to 1.418 Mbps) guarantee no loss. Both VoIP and video traffic at the exit of the IP buffers are encapsulated over DVB (packets of 188 bytes) and generate the process  $a(t)$ .  $PLP_{video}^*$  is set to  $5 \cdot 10^{-3}$ .

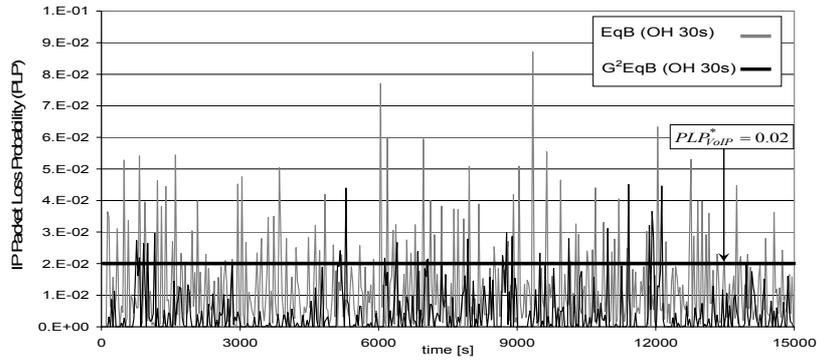


Figure 2. G<sup>2</sup>EqB and EqB: PLP.

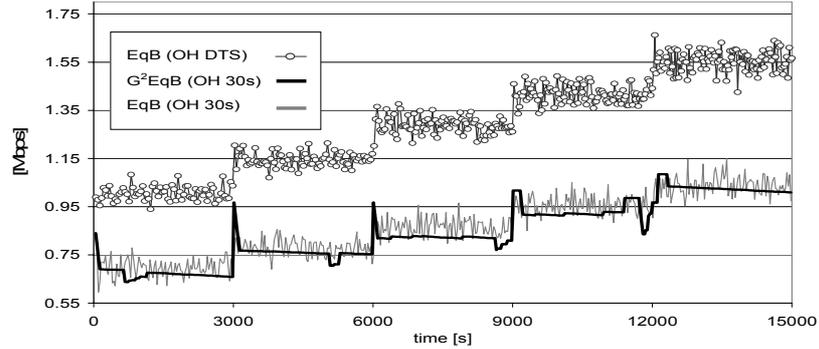


Figure 3. G<sup>2</sup>EqB versus EqB: bandwidth allocation.

Num. VoIP Calls	Buffer Length [bytes]	Allocated Bandwidth G <sup>2</sup> EqB/EqB	PLP video G <sup>2</sup> EqB/EqB	PLP VoIP G <sup>2</sup> EqB/EqB
30	9400	0.94 / 0.79	$2.62 \cdot 10^{-3}$ / $8.26 \cdot 10^{-2}$	$1.20 \cdot 10^{-4}$ / $5.58 \cdot 10^{-3}$
30	18800	0.84 / 0.79	$3.20 \cdot 10^{-3}$ / $2.96 \cdot 10^{-2}$	$2.22 \cdot 10^{-4}$ / $2.13 \cdot 10^{-3}$
30	28200	0.82 / 0.79	$1.78 \cdot 10^{-3}$ / $1.34 \cdot 10^{-2}$	$1.24 \cdot 10^{-4}$ / $9.14 \cdot 10^{-4}$
60	9400	1.47 / 1.26	$3.94 \cdot 10^{-3}$ / $8.26 \cdot 10^{-2}$	$2.26 \cdot 10^{-4}$ / $7.77 \cdot 10^{-3}$
60	18800	1.41 / 1.26	$2.23 \cdot 10^{-3}$ / $5.13 \cdot 10^{-2}$	$1.55 \cdot 10^{-4}$ / $3.96 \cdot 10^{-3}$
60	28200	1.39 / 1.26	$1.38 \cdot 10^{-3}$ / $3.0 \cdot 10^{-2}$	$1.04 \cdot 10^{-4}$ / $2.32 \cdot 10^{-3}$
90	9400	1.97 / 1.73	$4.13 \cdot 10^{-3}$ / $1.07 \cdot 10^{-1}$	$2.38 \cdot 10^{-4}$ / $8.11 \cdot 10^{-3}$
90	18800	1.87 / 1.73	$3.26 \cdot 10^{-3}$ / $5.33 \cdot 10^{-2}$	$1.94 \cdot 10^{-4}$ / $3.97 \cdot 10^{-3}$
90	28200	1.84 / 1.73	$1.73 \cdot 10^{-3}$ / $3.07 \cdot 10^{-2}$	$1.02 \cdot 10^{-4}$ / $2.35 \cdot 10^{-3}$

Table 1. G<sup>2</sup>EqB versus EqB: packet loss and bandwidth allocation.

Table 1 contains average measured PLP and allocated bandwidth in [Mbps] both for G<sup>2</sup>EqB and EqB. DVB buffer dimension is changed as well as the number of VoIP calls. OH is set to 3 minutes. Each single test simulates 107 overall minutes. Proper G<sup>2</sup>EqB gradient stepsizes are chosen for each single test. The average G<sup>2</sup>EqB PLP is always close to but below the threshold of the most restrictive requirement. G<sup>2</sup>EqB is adaptive to buffer length because its behavior depends only on loss measures.

## VII. CONCLUSIONS AND FUTURE WORK

A novel equivalent bandwidth algorithm to automatically adapt the rate assigned to a buffer and counteract time varying system conditions is introduced. It is based only on measures.

No closed-form expressions, no a-priori information about source statistical properties, and no assumptions about buffer dimension are imposed. To achieve quick on-line convergence, the algorithm requires: proper off line dimensioning of gradient descent stepsize and proper length of the observation horizon. Time-scale of traffic changes must be slower than the time required for convergence. Otherwise, the resulting allocations are suboptimal.

Future research may regard: **1)** multi-objective control, namely, joint control of non homogeneous (sometimes conflicting) performance metrics (loss versus delay, versus delay jitter of the packets); **2)** implementation and validation of the algorithm within a Linux-based testbed architecture [6]; **3)** tuning of algorithm parameters to reliably support QoS in heterogeneous conditions; **4)** performance analysis of the control algorithm with respect to different traffic categories, for example in the presence of congestion control (e.g., TCP).

## REFERENCES

- [1] C. G. Cassandras, G. Sun, C. G. Panayiotou, Y. Wardi, "Perturbation Analysis and Control of Two-Class Stochastic Fluid Models for Communication Networks," *IEEE Trans. Automat. Contr.*, vol. 48, n. 5, May 2003, pp. 23-32.
- [2] Y. Wardi, B. Melamed, C.G. Cassandras, C.G. Panayiotou, "Online IPA Gradient Estimators in Stochastic Continuous Fluid Models," *J. of Optimization Theory and Applic.*, vol. 115, no. 2, pp. 369-405, November 2002.
- [3] R. Guérin, H. Ahmadi, M. Naghshineh, "Equivalent Capacity and its Application to Bandwidth Allocation in High-Speed Networks," *IEEE J. Select. Areas Commun.*, vol. 9, no. 7, September 1991, pp. 968-981.
- [4] S. Georgoulas, P. Trimintzios, G. Pavlou, K. Ho, "Heterogeneous Real-time Traffic Admission Control in Differentiated Services Domains," Proc. IEEE Global Telecommunications Conference 2005 (*Globecom 2005*), St. Louis, MO, 28 Nov.-2 Dec. 2005, pp. 523-528.
- [5] <http://www-tnk.ee.tu-berlin.de/research/trace/trace.html>.
- [6] <http://www.tlc-networks.polito.it/projects/borabora/>.