

# The Cost of Quality in Networks of Aggregate Traffic

N. G. Duffield      S.H. Low  
 AT&T Labs–Research      University of Melbourne

*Abstract*—We relate burstiness (or indifference) curves for stochastic traffic flows to the quality they experience under the allocation of the resources of bandwidth and buffer space. We use this relation to explore the quality experienced by merged flows under various rules for allocating resources to them, including one motivated by the Controlled Load service specification. We show how cost structures on the resources can be used to encourage optimal use of shared resources amongst heterogeneous flows.

*Keywords*—Pricing, economies of scale, indifference curve, admission control, performance analysis, large deviations, stochastic networks

## I. INTRODUCTION

It is critical to provide quality of service (QoS) guarantee in high speed networks in an efficient manner in terms of its use of network resources. For traffic that can tolerate a certain delay, buffer capacity, in addition to transmission bandwidth, is also a scarce resource that should be carefully managed in resource allocation. For instance methods to tradeoff bandwidth and buffer in resource allocation are developed in Elwalid, Mitra and Wentworth [9], Low and Varaiya [17], [16], and Presti, Zhang, Kurose and Towsley [22].

In this paper we examine the QoS, as determined by buffer overflow probabilities, of heterogeneous aggregates of traffic flows with variable rate but whose statistical properties are stationary. As we shall presently describe, a potential problem that arises when sharing resources, is that heterogeneous characteristics and quality requirements can lead to suboptimal use of resources. However, we shall show that this need not be the case. Bandwidth and buffer are substitutable resources to guarantee the same quality and can be traded off in order to optimize network performance.

In this paper we shall find the circumstances under which it is possible for flows to specify, out of all resource requests compatible with a desired quality, that which guarantees quality in the shared resource. Roughly speaking, this choice ensures that the parameters of the traffic and resource which determine the frequency of buffer overflow are common across all flows. In some sense, the overflow-determining properties of the heterogeneous flows can be rendered homogeneous over subsets of flows. Furthermore, we show how the relative cost of network resources may be chosen in order to encourage flows to specify resources in the optimal manner, and comment on how the consequences of such choice propagate in a network.

The approach in this paper will be stochastic in that we deal with the statistical multiplexing of flows described by stochastic

processes. However, part of our framework is motivated by studies of the (deterministic) properties of individual sample paths. The form of the tradeoff between bandwidth and buffer is directly used to characterize traffic flows in [17]. An arrival process  $A = (A_{t,t'})_{t \leq t' \in \mathbf{R}}$  specifies the amount of work  $A_{t,t'}$  generated in intervals  $[t, t')$ . To  $A$  we can associate its convex decreasing **burstiness curve** (or indifference curve)  $b_A(c)$ . This is defined as the minimal buffer size  $b_A(c)$  to prevent overflow when  $A$  is served at a deterministic rate  $c$ :

$$b_A(c) = \sup_{-\infty < t < t' < 0} \{A_{t,t'} - c(t' - t)\}. \quad (1)$$

This characterization can be thought of as a generalization of the  $(\sigma, \rho)$  characterization of Cruz [5] if we treat  $\sigma$  as a function of  $\rho$ . In [17], several burst reducing servers whose output process  $A'$  is less bursty than their input process  $A$ , i.e.,  $b_{A'}(c) \leq b_A(c)$ , are identified and their behavior in tandem is studied. Based on these results, in [16], the burstiness of different traffic classes and at different servers is traded off in allocation of bandwidth and buffer in a network of burst reducing servers in order to optimize overall network performance. The deterministic fluid flow model in [16], [17] leads to simpler analysis, but is not suitable to study multiplexing with nonzero loss requirement. In this paper we develop a stochastic framework where such issues can be better examined; for a different approach, see, e.g., Kurose [15], Yaron and Sidi [28] and Chang [3].

The stochastic counterpart to the burstiness curve arises as follows. Use  $A$  to model a stochastic traffic flow with stationary increments. We will abbreviate the random quantity  $A_{-t,0}$  by  $A_t$ . Suppose the aggregate  $A^n$  of  $n$  independent copies of  $A$  feeds a buffer which is drained at a constant rate  $nc$  for some  $c$ . Let  $Q^n$  denote the random queue length at time zero. Following Borovkov [1], we can express this pathwise in terms of  $A^n$  as

$$Q^n = \sup_{t > 0} (A_t^n - cnt). \quad (2)$$

Under very general circumstances the tail probability for  $Q^n$  obeys the asymptotic

$$\lim_{n \rightarrow \infty} n^{-1} \log \mathbf{P}[Q^n > nb] = -I(c, b) \quad (3)$$

where the **shape function**  $I$  is determined by the statistical properties of  $A$ . This suggests using the following **shape function approximation** for tail probabilities:

$$\mathbf{P}[Q^n > nb] \approx e^{-nI(c,b)} \quad (4)$$

Generally,  $I(c, b)$  determines the *economies of scale* which are available through statistical multiplexing. For example, for a

Nick Duffield is at AT&T Laboratories, Room A175, 180 Park Avenue, Florham Park, NJ 07932-0971 USA. E-mail duffield@research.att.com

Steven Low is at the Dept. of Electrical & Electronic Engineering, The University of Melbourne, Parkville, Victoria 3052, Australia. E-mail s.low@ee.mu.oz.au

broad class of short-range dependent flows  $I(c, b) \approx \delta b + \nu$ , for large  $b$ , where  $\nu$  is positive for flows with the positive correlation property that its increments are associated. Thus, estimates of system capacity based on the approximation  $\mathbf{P}[Q^n > nb] \approx e^{-nI(c, b)} \approx e^{-n(\delta b + \nu)}$  will be less conservative than approximations based on the corresponding *effective bandwidth approximation*  $\mathbf{P}[Q^n > nb] \approx e^{-n\delta b}$  for such flows. Thus by use of (3) we aim to develop connection admission control (CAC) algorithms which can make better use of system resources than those based on the effective bandwidth approximation; e.g. see Kelly [12], and Gibbens and Hunt [10].

The relation (3) expresses the asymptotic behavior of  $Q^n$  as  $n$  increases. In view of (2) it is not surprising then that the asymptotic properties of  $Q^n$  for large  $n$  can be expressed in terms of those of  $A^n$ . Define the moment generating function of an individual flow by

$$\mu(t, \theta) = \log \mathbf{E}[e^{\theta A_t}]. \quad (5)$$

Recall the definition of the Legendre transform  $f^*$  of a real function  $f$ , i.e.,

$$f^*(x) = \sup_{\theta \in \mathbf{R}} \{x\theta - f(\theta)\}. \quad (6)$$

Fix some  $t > 0$ . Then the usual Chernoff bound for sums of independent random variables shows that

$$\mathbf{P}[A_t^n \geq n(ct + b)] \leq e^{-n\mu^*(t, ct+b)}. \quad (7)$$

Duffield [7] has established very general conditions under which that (3) holds with

$$I(c, b) = \inf_{t \geq 0} \mu^*(t, b + ct) = \inf_{t > 0} \sup_{\theta \in \mathbf{R}} \{(b + ct)\theta - \mu(t, \theta)\}, \quad (8)$$

where  $\mu^*$  denotes the Legendre transform of  $\mu$  w.r.t its second argument. (8) says that, asymptotically as  $n \rightarrow \infty$ , the probability  $\mathbf{P}[Q^n > nb]$  is equal to the least upper bound (7) as  $t$  is varied, and that this bound is tight. This is the familiar large deviation heuristic that ‘‘rare events occur in the most likely way’’. It is worth remarking that the conditions established in [7] allow for long range dependent traffic flows; see also Botvich and Duffield [2], Courcoubetis and Weber [4] and Simonian and Guibert [24] for results for short-range dependent traffic flows.

Actually the results from [7] are not limited to aggregates of identical flows, nor of independent ones. More generally one can consider a sequence of aggregations  $A^{(n)}$ , with MGF  $\mu^{(n)}$ . (An i.i.d. aggregation would have  $\mu^{(n)} = n\mu$  where  $\mu$  is the single flow MGF). For  $B = bn$  the shape function approximation (4) then becomes  $\mathbf{P}[Q^n > B] \approx e^{-I^{(n)}(B)}$  where  $I^{(n)}$  is derived from  $\mu^{(n)}$  using (8). We will use this approximation for aggregate flows in what follows. We will refer to a specific value of the shape function as a **quality**.

The stochastic counterpart of the burstiness curve derives from the shape function. Informally at least, we can define the inverse of  $I$  in its second component as  $b(c, \varepsilon) = I(c, \cdot)^{-1}(\varepsilon)$ . We can interpret  $b(c, \varepsilon)$  as the buffer  $b$  per flow required in order that the quality, as specified by the log-tail probability  $\log \mathbf{P}[Q^n > nb]$ , is  $\varepsilon$  when the service rate is  $c$  per flow. Thus  $b(c, \varepsilon)$  is the stochastic counterpart of the indifference or burstiness curve (1) in that  $\{(b, c) \mid b \geq b(c, \varepsilon)\}$  is the locus of points for which the quality  $I(c, b)$  is at least  $\varepsilon$ . In Section II we will

give the formal definition of the burstiness curve in terms of the shape function, and establish its basic properties. In Section III we show that the burstiness of a flow is non-increasing on passage through a FIFO queue. This can be seen as a stochastic development of the results from [16] on deterministic flows, and also of asymptotic results for departure processes in the effective bandwidth formalism; see de Veciana, Courcoubetis and Walrand [25] and O'Connell [21].

As well as understanding the propagation of burstiness through the network, we also want to be able to express the burstiness of an aggregate of flows in terms of the burstiness of its constituents. This is important in understanding how to allocate resources for the aggregate flows. In Section IV we give two rules for allocating bandwidth and buffer to aggregations:

**Sum rate, sum buffer (SRSB)** The aggregate is allocated the sum over constituents of the bandwidth and buffer capacity. This is motivated in part by the rules for merging Traffic Specifications (TSpec's) when flows are to be combined in the Controlled Load service specification [27]. In Section V we investigate the conditions under which the quality of the aggregate, as determined by buffer overflow probabilities, can be at least as great as the minimum quality of its constituent flows.

**Sum rate, maximum buffer (SRMB)** A more stringent allocation of resources to the aggregate is to sum the bandwidths but, in distinction with SRSB, use the maximum buffer allocation of the constituents. In Section VI we show that flows whose autocorrelation satisfies a certain positivity condition can be (homogeneously) aggregated in this manner without loss of quality.

Since the burstiness curves and shape functions can be regarded as inverses, it is not surprising that our results on combinations of shape functions can also be recast as results about combinations of burstiness curves. We shall obtain stochastic counterparts of the combination law for burstiness curves of deterministic traffic flows, namely that

$$b_{\sum_i A_i} \left( \sum_i c_i \right) \leq \sum_i b_{A_i}(c_i), \quad (9)$$

where  $\sum_i A_i$  denotes the traffic flow obtained by aggregating traffic flows  $A_i$ . This inequality follows straightforwardly from the definition (1) of  $b_A(c)$ .

In the following sections we find that the quality of the aggregate satisfies a certain optimality condition if the constituent flows are, in some sense, well-matched. Suppose the extremum in (8) is achieved at  $(t^*, \theta^*)$ . Assuming the requisite differentiability, (8) implies that

$$\frac{\partial I(c, b)}{\partial c} = t^* \theta^*, \quad \frac{\partial I(c, b)}{\partial b} = \theta^*. \quad (10)$$

Thus  $\theta^*$  characterizes the sensitivity of quality to changes in resources, while  $t^*$  gives the relative sensitivity of the flows. Put another way,

$$t^* = \frac{\partial I(c, b)}{\partial c} / \frac{\partial I(c, b)}{\partial b} = - \frac{\partial b}{\partial c} \Big|_I = - \frac{\partial b(c, \varepsilon)}{\partial c} \Big|_{I(c, b) = \varepsilon} \quad (11)$$

gives the (negative of the) slope of the burstiness curve corresponding to the quality  $I(c, b)$ . In view of the preceding paragraph,  $t^*$  also determines the time over which arrivals build up

to cause the queue to exceed a level  $b$  per flow. Thus  $t^*$  and  $\theta^*$  together determine the manner in which the resource capacities  $c$  and  $b$  are exceeded. The flows are well-matched if the characteristic  $t^*$  and  $\theta^*$  are identical across flows. In a sense, the overflow-determining properties of the heterogeneous flows are rendered homogeneous. Since  $t^* = -\frac{\partial b(c, \varepsilon)}{\partial c}$ , an arbitrary set of flows can match at least their values of  $t^*$  by choosing an appropriate point on the burstiness curve by which to specify their resource requirements  $(b, c)$ .

The incentive to do this can be provided in the costs attached to the resources  $b$  and  $c$ . In Section VII we show that if the ratio of the cost of service rate to that of buffer space is  $\gamma$ , then individual users minimize their costs by choosing their point on the burstiness curve for which  $t^* = \gamma$ . The minimal cost can be specified directly in term of the MGF  $\mu$  of the flow. In Section VIII we formulate cost based admission controls to aggregate channels. Resources are treated as a commodity, and an aggregator of flows is able to pass on the economies of scale of multiplexing as reduced costs. We show how, if a potential user of the aggregator knows their MGF they can tell whether the charge for a given quality reflects optimal use of aggregate resources, and whether the same quality at lower cost could be achieved without going through an aggregator. All proofs are omitted and will be provided elsewhere.

## II. SHAPE FUNCTION, BURSTINESS AND QUALITY

The precise form of (3) is the following

*Theorem 1:* Under Hypothesis 1 in [7]

$$\begin{aligned} -I(c, b) &\geq \limsup_{n \rightarrow \infty} n^{-1} \log \mathbf{P}[Q^n > nb] \\ &\geq \liminf_{n \rightarrow \infty} n^{-1} \log \mathbf{P}[Q^n > nb] \geq -I(c, b^+), \end{aligned} \quad (12)$$

where  $I(c, b^+)$  denotes the limit from the right in the second argument. A sufficient condition for the upper and lower bounds to coincide is that  $\mu^*(t, \cdot)$  be continuous on  $\mathbf{R}_+$ .

We do not reproduce the hypotheses in detail here, but will make some brief comments. For  $\mu^{(n)}$  the MGF of  $A^{(n)}$ , the general aggregate flow described in the introduction, the main hypothesis is that the limit  $\mu = n^{-1}\mu^{(n)}$  exists as  $n \rightarrow \infty$  and is essentially smooth (see [23] for terminology). By the Gärtner-Ellis theorem (see [6]) this ensures that for each  $t$ , the pair  $(A_t^{(n)}, n)$  satisfies a *Large Deviation Principle* with rate function  $\mu^*(t, \cdot)$ .

For each quality  $\varepsilon > 0$  define the  $\varepsilon$ -**burstiness curve** (or indifference curve)  $b(\cdot, \varepsilon)$  as the right-inverse of  $I$  in its second component, i.e.,

$$b(c, \varepsilon) = \inf\{b > 0 \mid I(c, b) \geq \varepsilon\} \quad (13)$$

Thus  $b(c, \varepsilon)$  is the minimum buffer allocation required for quality  $\varepsilon$  at service rate  $c$ . Let  $a$  denote the mean arrival rate of a flow, i.e.,  $\mathbf{E}[A_t] = at$ . For each  $t > 0$  set  $b(c, \varepsilon, t) = \mu^*(t, \cdot)^{-1}(\varepsilon) - ct$ , where  $\mu^*(t, \cdot)^{-1}$  is the left-inverse of  $\mu^*(t, \cdot)$  on the domain  $[at, \infty)$  on which it is non-decreasing:

$$\mu^*(t, \cdot)^{-1}(x) = \inf\{y > at \mid \mu^*(t, y) \geq x\}. \quad (14)$$

Sometimes we will write  $b_\varepsilon(c)$  for  $b(c, \varepsilon)$  and  $b_\varepsilon(c, t)$  for  $b(c, \varepsilon, t)$ .

*Theorem 2:* (i) For each  $\varepsilon > 0$ ,  $b(c, \varepsilon)$  is a convex, non-increasing function on  $\mathbf{R}^+$ .

(ii)  $b \geq b(c, \varepsilon)$  iff  $I(c, b) \geq \varepsilon$ .

(iii)  $b^* = b(c, \varepsilon) = \sup_{t>0} b(c, \varepsilon, t)$  is achieved at  $t^*$  iff  $\varepsilon = I(c, b^*) = \inf_{t>0} \mu^*(t, b^* + ct)$  is achieved at  $t^*$ .

## III. DEPARTURE PROCESSES AND NETWORK QUALITY

Here we show that, under technical conditions, quality propagates in a network in the sense that the burstiness of a traffic flow not increased by passage through a buffer.

We model the departure process by a fluid, so that arriving work is presented continuously to the output as it is processed. When the service rate is  $nc$ , an upper bound for the amount of work departing in the interval  $[t, t']$  is thus  $D_{t,t'} = \min[nc(t' - t), A_{t,t'}^n + Q_t^n]$  where  $Q_t^n = \sup_{t' < t} (A_{t',t}^n - nc(t - t'))$  is the length of the queue of arrived work unprocessed at time  $t$ . For the purposes of establishing the properties of the  $D^n$  as an upper bound for a potential arrival process at another queue it is convenient to work with

$$D_t^n = D_{-,t,0}^n = cnt + \min[0, \sup_{t' > t} (A_{t',t}^n - nct')]. \quad (15)$$

*Theorem 3:* Assume the departure  $D_t^n$  satisfies the hypothesis of Theorem 1 for some MGF. Let  $b, \tilde{b}$  be the burstiness curves of arrivals and departures. Then for all  $\varepsilon > 0$ ,

(i)  $\tilde{b}(\tilde{c}, \varepsilon) \leq b(\tilde{c}, \varepsilon)$  for all  $\tilde{c}$ .

(ii)  $\tilde{b}(\tilde{c}, \varepsilon) = 0$  for all  $\tilde{c} \geq c$ .

We now present a number of subsidiary results used in the proof of Theorem 3 that are of independent interest, and illustrate these results with an example application for fractional Brownian motion arrivals. The next theorem characterizes the large deviation behavior of the departure process at the level of rate functions.

*Theorem 4:* Under the hypotheses of Theorem 1, for  $x \geq a$

$$\begin{aligned} -J(c, t, x) &\geq \limsup_{n \rightarrow \infty} n^{-1} \log \mathbf{P}[D_t^n \geq nx] \\ &\geq \liminf_{n \rightarrow \infty} n^{-1} \log \mathbf{P}[D_t^n \geq nx] \geq -J(c, t, x^+) \end{aligned}$$

where

$$J(c, t, x) = \begin{cases} \inf_{t' \geq t} \mu^*(t', x + c(t' - t)) & \text{if } t \geq x/c \\ +\infty & \text{otherwise} \end{cases}$$

Now assume the departure process feeds a downstream queue whose service rate is  $n\tilde{c}$  for some  $\tilde{c} \in (a, c)$ . The stationary queue length at time 0 is  $\tilde{Q}^n = \sup_{t>0} (D_t^n - n\tilde{c}t)$ .

*Theorem 5:* Assume the process  $D_t^n$  satisfies the hypotheses of Theorem 1 for some MGF. Then

$$\limsup_{n \rightarrow \infty} n^{-1} \log \mathbf{P}[\tilde{Q}^n > nb] \leq -\tilde{I}(\tilde{c}, b) \quad (16)$$

where

$$\tilde{I}(\tilde{c}, b) = \inf_{t' \geq b/(c-\tilde{c})} \mu^*(t', b + \tilde{c}t') \geq I(\tilde{c}, b). \quad (17)$$

**Example 1.** Consider the model of Norros [20] where  $A_t = at + V_t$  where  $V_t$  is standard fractional Brownian motion with

Hurst parameter  $H \in (1/2, 1)$ . Since  $V_i$  is Gaussian with zero mean and variance  $t^{2H}$  we have

$$\mu(t, \theta) = at\theta + \theta^2 t^{2H}/2 \quad \text{and hence} \quad \mu^*(t, x) = \frac{(x - at)^2}{2t^{2H}}. \quad (18)$$

The variational calculation (8) gives

$$I(c, b) = \frac{1}{2} \left( \frac{b}{1-H} \right)^{2-2H} \left( \frac{c-a}{H} \right)^{2H}. \quad (19)$$

Then it follows from (17) that

$$\tilde{I}(\tilde{c}, b) = \begin{cases} 0 & \text{if } \tilde{c} \leq a \\ I(\tilde{c}, b) & \text{if } a < \tilde{c} \leq c(1-H) + aH \\ \mu^*\left(\frac{b}{c-\tilde{c}}, \frac{bc}{c-\tilde{c}}\right) & \text{if } c(1-H) + aH \leq \tilde{c} < c \\ +\infty & \text{if } c \leq \tilde{c}. \end{cases}$$

To see this note that  $t \mapsto \mu^*(t, b + ct)$  is convex and achieved its infimum at  $t^* = Hb/((1-H)(c-a))$ . The result then follows by comparing  $t^*$  with  $t_0 \equiv b/(c-\tilde{c})$ , the lower range of the infimum in (17). Now  $t^* \geq t_0$  if  $\tilde{c} \leq c(1-H) + aH$ , in which case  $\tilde{I}(\tilde{c}, b) = I(\tilde{c}, b)$ . Otherwise,  $\tilde{I}(\tilde{c}, b)$  is the value of  $\mu^*(t', b + \tilde{c}t')$  at  $t_0$ . In Figure 1 we have plotted  $I(\tilde{c}, b)$  and  $\tilde{I}(\tilde{c}, b)$  for  $a = 1, c = 2$ . The principle modification of  $I$  to  $\tilde{I}$  is that  $\tilde{I}$  increases to  $\infty$  as  $\tilde{c} \rightarrow c$ . The fact that the arrival rate to the downstream queue is limited to  $c$  means that large occupation of the downstream queue becomes increasingly unlikely as the downstream service rate  $\tilde{c}$  approaches  $c$ . The corresponding burstiness curves are

$$b(c, \varepsilon) = \begin{cases} \infty & \text{if } c \leq a \\ (1-H)(2\varepsilon(H/(c-a))^{2H})^{1/(2-2H)} & \text{if } c > a \end{cases}$$

for the input traffic flow. Downstream, the burstiness is

$$\tilde{b}(\tilde{c}, \varepsilon) = \begin{cases} b(\tilde{c}, \varepsilon) & \text{if } \tilde{c} \leq c(1-H) + aH \\ (c-\tilde{c})\left(\frac{2\varepsilon}{(c-a)^2}\right)^{\frac{1}{2-2H}} & \text{if } c(1-H) + aH \leq \tilde{c} < c \\ 0 & \text{if } c \leq \tilde{c} \end{cases}$$

We have plotted this for  $\varepsilon = 4$  in Figure 2.

#### IV. RESOURCE ALLOCATION FOR AGGREGATES

In this section and those following we develop criteria by which it can be determined whether aggregates of traffic flows, either homogeneous or heterogeneous, will enjoy a target quality. We use the shape function  $I(c, b)$  for given service rate  $c$  and buffer size  $b$  as the measure of quality. We will generally work with a finite set of traffic flows characterized by their MGF  $\mu_i$ , where  $i$  is an index running over the set of flows. Observe that  $\sum_i \mu_i$  is the MGF of an aggregation of independent flows of MGF  $\mu_i$ .

Let  $c_i$  and  $b_i$  be the resources allocated to flow  $i$ . We will examine the effect of resource combination rules where quality guarantees are only probabilistic in nature. Sometimes the target probability may not even be explicitly given. For example, in the proposed Controlled Load service specification [27], a flow is characterized by a Traffic Specification (TSpec) including leaky bucket parameters  $(C, B)$ , a rate  $C$  and a bucket size  $B$ . (Peak rate, minimum policed unit and maximum packet size

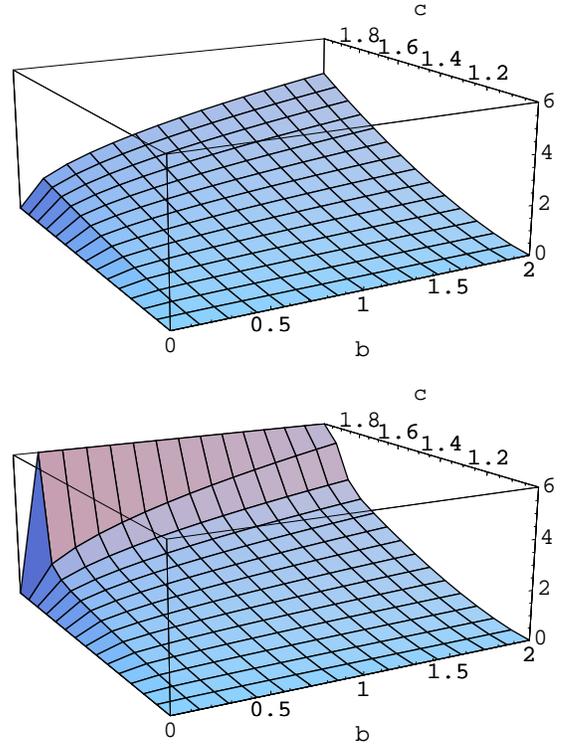


Fig. 1. Non-decreasing property of shape function under passage through FIFO queue. Fractional Brownian motion input of Example 1. Upper: inward shape function  $I(c, b)$  for service rate  $c$ . Lower: shape function  $\tilde{I}(\tilde{c}, b)$  at downstream queue for service rate  $\tilde{c} \in (a, c)$ , where  $a$  is mean arrival rate.

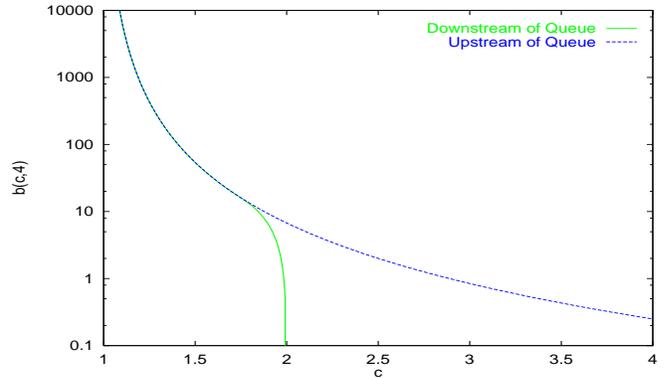


Fig. 2. Non-increasing property of burstiness curve under passage through FIFO queue. Fractional Brownian motion input of Example 1. Burstiness upstream of queue, and downstream of queue.

are also included, but these shall not concern us here). Packets conformant to a given TSpec should not experience queuing delay beyond that resulting from their own burstiness. Non-conformant packets will be sent by Best Effort service. Now, when traffic flows are to be combined at a network element, the merged TSpec of the combination (or more precisely, the  $(C, B)$  part of the combined TSpec which we consider here) is to be derived from that of its constituents by the following rule:

**Sum rate, sum buffer (SRSB)** The aggregate flow with MGF  $\mu = \sum_i \mu_i$  is allocated service rate  $c = \sum_i c_i$  and buffer  $b = \sum_i b_i$ .

The specific property we want to investigate is non-

conformance: if the constituent traffic flows are non-conformant to their Tspec with some probability  $\varepsilon$ , possibly zero, can the same be said when their aggregate shares resources under the SRSB rule. This we investigate in Section V following. In Section VI we briefly investigate this question for homogeneous aggregates under the following rule that allocates less resources:

**Sum rate, max buffer (SRMB)** The aggregate flow with MGF  $\mu = \sum_i \mu_i$  is allocated service rate  $c = \sum_i c_i$  and buffer  $b = \max_i b_i$ .

## V. SUMMED RATE SUMMED BUFFER ALLOCATIONS

The properties of SRSB for homogeneous flows follow simply from the definition of the shape function.

*Theorem 6:* The quality available to an  $n$ -fold homogeneous aggregate is  $n$  times that available to each component. In other words: for  $n \in \mathbb{N}$ ,  $I^{n\mu}(nc, nb) = nI^\mu(c, b)$ , and hence  $b^{n\mu}(nc, n\varepsilon) = nb^\mu(c, \varepsilon)$ .

For heterogeneous aggregate, the properties of the shape function are somewhat more complicated. Our first result is that quality is not impaired under SRSB.  $I^\mu$  will denote the shape function for  $\mu$ ,  $b^\mu$  its burstiness.

*Theorem 7:* (i) The quality of a heterogeneous aggregate under SRSB is at least as great as the minimum quality of its components, i.e.,

$$I^\mu\left(\sum_i c_i, \sum_i b_i\right) \geq \min_i I^{\mu_i}(c_i, b_i). \quad (20)$$

(ii) Hence  $b^\mu(\sum_i c_i, \varepsilon) \leq \sum_i b^{\mu_i}(c_i, \varepsilon)$  for all  $\varepsilon > 0$ .

Theorem 7(ii) can be viewed as a stochastic version of the aggregate inequality (9) for the deterministic burstiness curves (1). Recall that the deterministic burstiness of (1) and (9) shows that the burstiness curve of an aggregate, as defined in (1) to guarantee zero loss, is always smaller than the sum  $\sum_i b_{A_i}(c_i)$  of individual burstiness curves. Here we have shown that the same holds true for  $\varepsilon_i$ -burstiness of the aggregate and its components. Thus we have achieved one of the aims set out in the previous section.

However, we expect to be able to do better than Theorem 7 and obtain more than the minimal quality of a constituent flow in an aggregate. As we shall see now, this can happen only when the traffic flows are, in some sense, well matched.

*Theorem 8:* Let the infimum in (8) for  $I^\mu(\sum_i c_i, \sum_i b_i)$  be attained at  $t^*$ .

(i) Let the supremum in  $\mu_i^*(t^*, b_i + c_i t^*)$  be attained at the same point  $\theta$  for all  $i$ . Then

$$I^\mu\left(\sum_i c_i, \sum_i b_i\right) \geq \sum_i I^{\mu_i}(c_i, b_i) \quad (21)$$

with equality iff the infimum in  $I^{\mu_i}(c_i, b_i)$  is attained at  $t^*$  for all  $i$ . Hence  $b^\mu(\sum_i c_i, \sum_i \varepsilon_i) \leq \sum_i b^{\mu_i}(c_i, \varepsilon_i)$  where  $\varepsilon_i = I^{\mu_i}(c_i, b_i)$ .

(ii) Let the infimum in (8) for  $\varepsilon_i := I^{\mu_i}(c_i, b_i)$  be attained at the same point  $\hat{t}$ . Then

$$I^\mu\left(\sum_i c_i, \sum_i b_i\right) \leq \sum_i I^{\mu_i}(c_i, b_i) \quad (22)$$

with equality iff  $\hat{t} = t^*$  and the supremum in  $\mu_i^*(t^*, b_i + c_i t^*)$  is attained at the same  $\theta \forall i$ . Hence  $b^\mu(\sum_i c_i, \sum_i \varepsilon_i) \geq \sum_i b^{\mu_i}(c_i, \varepsilon_i)$

The importance of the conditions in the theorem is that without them, the effect on quality of adding new components to an aggregate is uncontrolled, even under the SRSB rule: it is possible to have  $I(\sum_i c_i, \sum_i b_i) < \max_i I^{\mu_i}(c_i, b_i)$ .

**Example 2.** Consider Gaussian flows  $A_i(t) = \lambda_i t + \sigma_i \sqrt{t} Z_i(t)$ ,  $i = 1, \dots, n$ , where  $Z_i(t)$  are white Gaussian and consider their aggregate  $\sum A_i(t)$ . Their shape functions are respectively

$$I^{\mu_i}(c_i, b_i) = \frac{2b_i(c_i - \lambda_i)}{\sigma_i^2}$$

$$I^\mu\left(\sum_i c_i, \sum_i b_i\right) = \frac{2\sum_i b_i \sum_i (c_i - \lambda_i)}{\sum_i \sigma_i^2}$$

When  $n = 2$ ,  $b_1 = c_1 - \lambda_1 = 1$ ,  $b_2 = c_2 - \lambda_2 = 2$ ,  $\sigma_1^2 = 2$ , and  $\sigma_2^2 = 1$ , we have

$$\max_i I^{\mu_i}(c_i, b_i) = I^{\mu_2}(c_2, b_2) = 8 > 6 = I^\mu(\sum_i c_i, \sum_i b_i)$$

## VI. SUMMED RATE MAXIMUM BUFFER FOR ASSOCIATED FLOWS

In the previous section we explored the extent to which increased quality could be obtained for an aggregate, as compared with its constituents, by SRSB allocation of the aggregate bandwidth and buffer  $(\sum_i c_i, \sum_i b_i)$ . Another approach to making use of the economies of scale is to offer less than the aggregate resources to the aggregate, but to have the aggregate quality no worse than that of the constituents. Theorem 7 showed that this is possible with SRSB allocation. In this section we show that this is possible for a wide class of traffic flows under SRMB allocation, i.e. of  $(\sum_i c_i, \max_i b_i)$ .

In part of what follows we shall consider traffic flows with the following correlation property. We say that an arrival process  $A$  is **associated** if all its increments over disjoint time intervals are associated, i.e., if for all increasing function  $f$  and  $g$  we have that

$$\mathbf{E}[f(A_{t_1, t_2})g(A_{t_3, t_4})] \geq \mathbf{E}[f(A_{t_1, t_2})]\mathbf{E}[g(A_{t_3, t_4})], \quad (23)$$

for disjoint time intervals  $[t_1, t_2)$  and  $[t_3, t_4)$ . In particular, taking  $f, g$  to be the identity we see that the increments are positively correlated when they are associated. Taking  $f(x) = g(x) = e^{\theta x}$  we see that for associated traffic flows,  $\mu(\cdot, \theta)$  is superadditive for  $\theta \geq 0$ , i.e.,

$$\mu(t+s, \theta) \geq \mu(t, \theta) + \mu(s, \theta), \quad \theta \geq 0. \quad (24)$$

As a simple consequence, for all  $n \in \mathbb{N}$   $\mu(nt, \theta) \geq n\mu(t, \theta)$  for  $t, \theta \geq 0$ . Then  $x \geq at$ :

$$\mu^*(t, x) \leq n\mu^*(t/n, x/n). \quad (25)$$

A simple example of a process with associated increments is a two-level Markov process in which successive arrivals are positively correlated, see Section 4 of Botvich and Duffield [2]. But generally, positive correlations are not a sufficient condition for association. We remark that one could establish (24) empirically for a given traffic flow.

*Theorem 9:* Under SRMB, the quality available to a homogeneous aggregate is at least as great as that of each component. In other words: letting  $\mu$  be the MGF of an associated

traffic flow, and  $n \in \mathbb{N}$ , then  $I^{n\mu}(nc, b) \geq I^\mu(b, c)$ , and hence  $b^{n\mu}(nc, \varepsilon) \leq b^\mu(c, \varepsilon)$ .

The assumption that a traffic flow has associated increments turns out to have strong ramifications for the form of the shape function  $I(c, b)$ , amongst which are that it is asymptotically linear in  $b$ . The restriction to associated traffic flows would appear to rule out certain classes of traffic, in particular long range dependent arrival flows; the asymptotic behavior of  $I$  is not in general linear for these, see Duffield [7], [8]. However, in order for Theorem 9 to hold, it is sufficient that (25) hold at the extremal  $t$  for  $I^{n\mu}$ . Thus only the details of the arrival process at the governing time-scale need be relevant.

## VII. COST STRUCTURES

In this section we introduce cost structures for bandwidth  $c$  and buffer  $b$ . Cost will act as a differentiator between points of the burstiness curve. A **cost structure** for a single set of resources  $(c, b)$  is a pair  $(\alpha, \beta) \in \mathbf{R}_+^2$ . Their interpretation is that

$$w = \alpha c + \beta b \quad (26)$$

is the cost of making the allocation  $(c, b)$ .

Let  $\mu$  be a MGF with  $\varepsilon$ -burstiness  $b_\varepsilon$ . Define the  $\varepsilon$ -**minimal cost**  $M(\gamma, \varepsilon)$  for all  $\gamma \geq 0$  by

$$M(\gamma, \varepsilon) = \inf_{c > 0} (b(c, \varepsilon) + \gamma c). \quad (27)$$

For each elementary cost structure  $(\alpha, \beta)$ ,  $\beta M(\beta^{-1}\alpha, \varepsilon)$  is the minimal cost at which traffic flow with an  $\varepsilon$ -burstiness curve  $b_\varepsilon$  can attain a quality  $\varepsilon$ .

At an informal level we can already see a striking relation between cost and quality.  $b_\varepsilon$  is convex, and if it is also differentiable then from (27) we see that the minimal cost is attained at that  $c$  for which  $b'_\varepsilon(c) = -\gamma$ . But as we saw by an informal argument in the introduction, (10–11) tell us that, with sufficient differentiability,  $b'_\varepsilon(c) = -t^*$ , where  $t^*$  achieves the infimum in (8). The conclusion from this would then be:

*For a given cost structure  $(\alpha, \beta)$ , a flow achieves minimal cost for a given quality  $\varepsilon$  when its resources  $(c, b)$  are such that  $t^* = \beta^{-1}\alpha$ . In this case the minimal cost is  $\alpha c + \beta b$ .*

We will now state more precisely sufficient conditions for the above to be true; differentiability of all functions concerned is not required. It will be convenient to use  $\mu_\varepsilon^{*-1}(t)$  as a shorthand for  $\mu^*(t, \cdot)^{-1}(\varepsilon)$ , and likewise  $M_\varepsilon$  for  $M(\cdot, \varepsilon)$ . We call an inverse of  $M(\gamma, \cdot)$ , namely

$$K(\gamma, x) = \sup\{\varepsilon > 0 \mid M(\gamma, \varepsilon) \leq x\} \quad (28)$$

the  **$x$ -maximal quality**. For the elementary cost structure  $(\alpha, \beta)$ ,  $K(\alpha/\beta, x/\beta)$  is the maximal quality available at cost  $x$ .

*Theorem 10:* For any  $\gamma > 0$  and  $x > 0$

$$\varepsilon = \mu^*(\gamma, M_\varepsilon(\gamma)), \quad \text{and} \quad K(\gamma, x) = \mu^*(\gamma, x) \quad (29)$$

provided any one of the following is true:

- (i)  $b_\varepsilon$  is differentiable
- (ii)  $M_\varepsilon$  is strictly concave;
- (iii)  $\mu_\varepsilon^{*-1}$  is concave;

(iv) for all  $b, c > 0$ ,  $I(b, c) = \inf_t \mu^*(t, b + ct)$  is uniquely attained.

In the following, by ‘cost’, we may mean  $\alpha c + \beta b$ ,  $(\alpha/\beta)c + b$ ,  $M(\alpha/\beta, \varepsilon)$ , or  $\beta M(\alpha/\beta, \varepsilon)$ ; the meaning should be clear from the context. If we assume that resources are available as required with charges according to the cost structure  $(\alpha, \beta)$ , then in order to allocate resources to a single flow, we must use the following **minimal cost resource allocation**:

- Determine the MGF  $\mu$ ; either by modelling, or measurement;
- Use Theorem 10 to find the minimal cost  $x$  for desired quality  $\varepsilon$
- Solve  $b'_\varepsilon(c) = \alpha/\beta$  to find resources  $(c^*, b^*)$  which realize this quality

If we assume that all resource allocation is done at minimal cost within the same cost structure  $(\alpha, \beta)$ , then all extremizing  $t^*$  in (8) is  $\alpha\beta^{-1}$  for all traffic flows, e.g., for all flows in an aggregate. Thus we obtain automatically the versions of previous theorems on quality for aggregates which presupposed equality of the various  $t^*$ .

*Theorem 11:* Let  $\mu = \sum_i \mu_i$ . Under any of the conditions (i) to (iv) of Theorem 10

(i) Maximal quality is subadditive over aggregates at constant total cost, i.e.,

$$K_\mu(\gamma, \sum_i x_i) \leq \sum_i K_{\mu_i}(\gamma, x_i). \quad (30)$$

Equality occurs when

$$\mu_i^{*'}(\gamma, x_i) = \mu^{*'}(\gamma, \sum_i x_i), \quad (31)$$

i.e., the values of  $\theta^*$  are identical across each  $\mu_i$  and their sum  $\mu$ . Here  $\mu^{*'}$  denotes the derivative of  $\mu^*$  w.r.t. its second argument.

(ii) Lower bound: if all elements in an aggregate have individual quality at least  $\varepsilon$  at optimal cost, then  $\varepsilon$  is a lower bound for the quality of the aggregate at the same total optimal cost, i.e.,

$$K_{\mu_i}(\gamma, y_i) \geq \varepsilon, \forall_i \implies K_\mu(\gamma, \sum_i y_i) \geq \varepsilon. \quad (32)$$

We saw in the introduction that the optimizing  $\theta$  in (8) can be interpreted as a quality sensitivity to resource. Cost structures allow us to interpret it as the sensitivity of maximal quality to optimal cost, since

$$\frac{\partial K(\gamma, x)}{\partial x} = \mu^{*'}(\gamma, x) = \theta^*. \quad (33)$$

Theorem 11(i) suggest that it might be possible to choose  $(c_i, b_i)$  such that the quality  $I^\mu(\sum_i c_i, \sum_i b_i)$  of the aggregate is at least  $\sum I^{\mu_i}(c_i, b_i)$ . To satisfy the sufficient condition in the theorem, however, would require coordinating the choice of  $(c_i, b_i)$  for all flows such that the supremum in  $\mu_i^*(t^*, b_i + c_i t^*)$  is attained at the same  $\theta$  for all  $i$ . Can a cost structure provide the incentive for this, i.e., does there exist an elementary cost structure  $(\alpha, \beta)$  such that, for all flows  $i$ , the minimizers  $(c_i, b_i)$  of resource costs

$$w_i(c_i, b_i) = \alpha c_i + \beta b_i$$

achieve given qualities  $\varepsilon_i$ , i.e.,  $I_i(c_i, b_i) = \varepsilon_i$ , and satisfy the condition of Theorem 8(i)? Such a cost structure does not unfortunately exist in general, as the next example shows. The reason is that quality requirement that  $I_i(c_i, b_i) = \varepsilon_i$  imposes an  $\varepsilon_i$ -burstiness curve  $b_i(c)$  on which  $(c_i, b_i)$  must lie. Minimizing the resource costs  $w_i(c_i, b_i)$  over a specified  $\varepsilon_i$ -burstiness curve yields allocations  $(c_i^*, b_i^*)$ , which may not satisfy the condition of Theorem 8(i).

**Example 3.** Consider the previous Gaussian flows from Example 2:  $A_i(t) = \lambda_i t + \sigma_i \sqrt{t} Z_i(t)$ ,  $i = 1, \dots, n$ , where  $Z_i(t)$  are white Gaussian. The shape function  $I_i(c_i, b_i)$  and its minimizing  $t_i$  and maximizing  $\theta_i$  are

$$I^{\mu_i}(c_i, b_i) = \frac{2b_i(c_i - \lambda_i)}{\sigma_i^2}, \quad t_i = \frac{b_i}{c_i - \lambda_i}, \quad \theta_i = \frac{2(c_i - \lambda_i)}{\sigma_i^2}.$$

The  $\varepsilon_i$ -burstiness curve of flow  $i$  is  $b_i(c) = \varepsilon_i \sigma_i^2 / 2(c_i - \lambda_i)$ . To satisfy the sufficient condition in Theorem 8(i) we must have  $\theta_i = \theta^*$  for all  $i$ . Hence the allocation

$$c_i^* = \lambda_i + \frac{1}{2} \theta^* \sigma_i^2, \quad b_i^* = b_i(c_i^*) = \frac{\varepsilon_i}{\theta^*}$$

achieves  $\varepsilon_i$  and  $I(C, B) \geq \sum I^{\mu_i}(c_i, b_i)$ . If  $(c_i^*, b_i^*)$  also minimizes  $w_i(c_i, b_i)$  then

$$b_i'(c_i^*) = -\frac{2\varepsilon_i}{\theta^{*2} \sigma_i^2} = \frac{\alpha}{\beta}$$

i.e.,  $\varepsilon_i / \sigma_i^2$  must be equal for all  $i$ , which may not be the case.

## VIII. ADMISSION CONTROL, COSTS AND INCENTIVES

We saw in the previous section that a single channel cannot in general provide additive quality across all potential flows. But we now show that when there are multiple channels with possibly different cost structures, it is possible to match flows to channels in order to promote matching of flow characteristics. Moreover, the cost structure provides a mechanism to encourage this which is tied closely to admission control. In what follows we assume that resources, buffer  $b$  and bandwidth  $c$  are available as required, subject to the charges for their use as specified by a cost structure.

We define **extended cost structure** as a triple  $(\alpha, \beta, \theta) \in \mathbf{R}_+^3$ . The role of  $\alpha$  and  $\beta$  is as in the cost structure defined previously. As the notation suggests,  $\theta$  will play the role of the second argument of  $\mu$ . By specifying an extended cost structure we are, in effect, announcing the extremal variables  $t = \alpha/\beta$  and  $\theta$  in (8) at which optimal quality will be obtained. The following theorem is an easy consequence of Theorem 10 and the definition of the Legendre transform.

*Theorem 12:* Consider a channel with extended cost structure  $(\alpha, \beta, \theta)$ , and set  $\gamma = \alpha\beta^{-1}$ . A flow with MGF  $\mu$  which satisfies any of the conditions of Theorem 10 can obtain quality  $\varepsilon$  at cost  $x$  if

$$\mu(\gamma, \theta) \leq x\theta - \varepsilon. \quad (34)$$

Moreover, the cost is minimal if  $\mu'(\gamma, \theta) = x$ .

When  $\mu = \sum_i \mu_i$  is an aggregation of multiple flows, then we obtain from (34) the following **cost based admission control rule**. Under the assumptions of Theorem 12:

*Flows with MGF  $\mu_i$  can obtain joint quality  $\varepsilon$  at total cost  $x$  if  $\sum_i \mu_i(\gamma, \theta) \leq x\theta - \varepsilon$ .*

The  $(\gamma\theta)^{-1} \mu_i(\gamma, \theta)$  are familiar as effective bandwidths as propounded by Kelly [13]. By choosing a value of  $\theta$  we obtain an admission control rule which is additive in the flows through their effective bandwidths, a desirable property. In some sense the value of  $\theta$  is arbitrary; choosing one value or another is simply specifying a condition for admission control, although this does not exclude the possibility that one choice of  $\theta$  may lead to more efficient use of resources than another. This aspect is familiar, for the same reason, from the admission control schemes of Gibbens and Kelly [11] based on a simple characterization of the one-time marginal distribution of the arrival process.

However, the arbitrariness can be resolved somewhat in the presence of different channels, with different extended cost structures. If  $\gamma, \theta$  and the actual charge  $y_i$  to flows  $\mu_i$  are known, then it can be determined to what extent  $y_i$  differs from that required for quality to be maximal, i.e.,  $\mu_i'(\gamma, \theta)$ . This provides the (owner of a) flow with a cost criterion by which to evaluate the service provided by different extended cost structures. With sufficiently many extended cost structures available, a flow could choose one for which the charge is closest to the optimal cost. Such behavior will lead to flows on individual channels becoming matched, in the sense that they share equal or closer values of  $\theta$ , and consequently their aggregate quality being maximized for the given total cost.

It can also be determined to what extent a flow is benefiting from the potential economies of scale of multiplexing. Let us suppose for example, that a given channel is available in “wholesale” and “retail” versions. These are distinct channels with extended cost structures  $(\alpha, \beta, \theta)$  and  $(\eta\alpha, \eta\beta, \theta)$  respectively, where  $\eta > 1$ . The channels have the same cost ratio  $\gamma = \alpha/\beta$  and sensitivity  $\theta$ , but the resources are a factor  $\eta$  more expensive for retail channel.

Suppose a flow is quoted a cost  $x$  to gain admittance to the wholesale channel at aggregate quality  $\varepsilon$ . But according to Theorem 10, the retail channel could give quality  $\varepsilon$  at cost  $\eta\mu^*(\gamma, \cdot)^{-1}(\varepsilon)$ . If this is less than  $x$  then the optimal cost of the retail channel is less for the same quality.

**Example 4.** Consider the flow of Example 1: a flow  $\mu$  with arrival process  $A_t = at + V_t$  with  $V_t$  fractional Brownian motion. Suppose the target quality  $\varepsilon$  is offered at the wholesale channel at cost  $y$ . Using (18) then we can tell that quality is not maximal at that cost if

$$y \geq \mu'(\gamma, \theta) = a\gamma + \theta\gamma^{2H}. \quad (35)$$

However, the retail channel will not be cheaper unless

$$y \geq \eta\mu^*(\gamma, \cdot)^{-1}(\varepsilon) = \eta(a\gamma + \sqrt{2\varepsilon\gamma^{2H}}). \quad (36)$$

Note that the inference on optimality is drawn without explicit reference to the other flows in the wholesale channel.

## IX. CONCLUSION

It is well known that a flow that is carried lowers the value of the network for every other flow, and that usage-based pricing is a way to internalize this externality, aligning social and individual optimality, see e.g., MacKie-Mason and Varian [18]. Usage-

based pricing is thus a useful tool to induce more efficient sharing of network resources. It may have several components such as an access charge, a connection charge per unit time, a volume charge, and a resource charge. For instance, based on effective bandwidth, Kelly [14] proposes a charging scheme with three components that encourages users to accurately estimate their mean rate, thus allowing more effective sharing of bandwidth. Distributed iterative schemes for setting usage-based prices in order to optimize resource allocation have been proposed, e.g., by Low and Varaiya [16], Murphy, Murphy and Posner [19] and de Veciana and Baldick [26].

The cost structure used in this paper is resource based where the price for a flow depends only on the resources buffer  $b$  and bandwidth  $c$  which are allocated, though the use of fixed and/or per unit time costs for these resources is envisaged. Although the determination of resource requirements via the MGF is quite sophisticated, the form of the cost itself is not, being just a linear combination of resources. The resources  $b$  and  $c$  can be specified when a connection is set up, e.g. as part of the TSpec used by Controlled Load. Hence once the cost structure is determined, the prices for a given resource allocation become readily available. Finally since bandwidth and buffer are separately priced, their prices can convey to the traffic flows the relative availability of bandwidth and buffer, allowing individual flows to pick bandwidth and buffer allocations, based on their burstiness curves, that incur the least cost. When different extended cost structures are available the resource costs, through their ratio, can be used to encourage well-matchedness among flows that are aggregated, leading to better joint quality.

## REFERENCES

- [1] A.A. Borovkov. *Stochastic Processes in Queueing Theory*, Springer, Berlin, 1976.
- [2] D.D. Botvich and N.G. Duffield, Large deviations, the shape of the loss curve, and economies of scale in large multiplexers, *Queueing Systems*, 20:293–320, 1995.
- [3] C.S. Chang . Stability, queue length and delay of deterministic and stochastic queueing networks. *IEEE Trans. on Automatic Control*. 39:913–931, 1994.
- [4] C. Courcoubetis and R. Weber. Buffer overflow asymptotics for a switch handling many traffic sources. *J. Appl. Prob.* vol. 33, pp. 886–903, 1996.
- [5] R. L. Cruz. A calculus for network delay Part I: network elements in isolation. *IEEE Transactions on Information Theory*, 37(1):114–131, January 1991.
- [6] A. Dembo and O. Zeitouni, *Large Deviation Techniques and Applications*. Jones and Bartlett, Boston-London, 1993.
- [7] N. G. Duffield, Economies of scale in queues with sources having power-law large deviation scalings, *J. Appl. Prob.* vol. 33, pp. 840–857, 1996.
- [8] N. G. Duffield, Queueing at large resources driven by long-tailed  $M/G/\infty$ -modulated processes, *Queueing Systems*, to appear.
- [9] A. Elwalid, D. Mitra, and R. Wentworth. A new approach for allocating buffer and bandwidth to heterogeneous, regulated traffic in an atm node. *IEEE Journal on Selected Area in Communications*, 13(6):1115–1127, August 1995.
- [10] R.J. Gibbens and P.J. Hunt. Effective Bandwidths for the multi-type UAS channel *Queueing Systems*, 9:17–28, 1991.
- [11] R.J. Gibbens and F.P. Kelly. Measurement-based connection admission control. *Proceedings ITC-15*, Washington, DC, USA, 22–27 June, 1997.
- [12] F.P. Kelly. Effective bandwidths at multi-type queues. *Queueing Systems* 9:5–16, 1991.
- [13] F.P. Kelly. Notes on effective bandwidths. in: *Stochastic Networks, Theory and Applications*, Eds. F.P. Kelly, S. Zachary and I. Ziedens, Royal Statistical Society Lecture Notes Series, vol. 4, pp.141–168, 1996.
- [14] F. P. Kelly, Charging and Accounting for Bursty Connections, In: *Internet Economics*, Eds. L. W. McKnight and J. P. Bailey, MIT Press, 1996.
- [15] J. Kurose. On computing per-session performance bounds in high-speed multi-hop computer networks. *Performance Evaluation Review*, 20(1):128–139, June 1992.
- [16] S. H. Low and P. P. Varaiya. A new approach to service provisioning in ATM networks. *IEEE/ACM Transactions on Networking*, 1(5):547–553, October 1993. For an updated version see <http://www.ee.mu.oz.au/staff/slow/research.html>.
- [17] S. H. Low and P. P. Varaiya. Burst reducing servers in ATM networks. *Queueing Systems*, 20:61–84, 1995.
- [18] J.K. MacKie-Mason and H.R. Varian, Pricing Congestible Network Resources, *IEEE Journal on Selected Areas in Communications*, 13:1141–1149, 1995.
- [19] J. Murphy, L. Murphy and E. C. Posner, Distributed Pricing for Embedded ATM Networks, In: *Proceedings of the 14th International Teletraffic Congress*, Ed. J. Labetoulle and J. W. Roberts, Elsevier Science, 1994
- [20] I. Norros, A storage model with self-similar input. *Queueing Systems*, 16:387–396, 1994.
- [21] N. O’Connell, Queue lengths and departures at single-server resources. in: *Stochastic Networks, Theory and Applications*, Eds. F.P. Kelly, S. Zachary and I. Ziedens, Royal Statistical Society Lecture Notes Series, vol. 4, pp.91–104, 1996.
- [22] F. L. Presti, Z. Zhang, J. Kurose, and D. Towsley. Source time scale and optimal buffer/bandwidth trade-off for regulated traffic in an atm node. In *Proceedings of Infocom’97*, April 1997.
- [23] R.T. Rockafellar, *Convex Analysis*. Princeton University Press, Princeton, 1970.
- [24] A. Simonian and J. Guibert, Large deviations approximation for fluid queues fed by a large number of on-off sources. *Proceedings of ITC 14, Antibes, 1994* pp. 1013–1022.
- [25] G. de Veciana, C. Courcoubetis and J. Walrand, Decoupling Bandwidths: A Decomposition Approach to Resource Management in Networks, *Proceedings, IEEE Infocom*, Toronto Canada, June 12–16 1994.
- [26] G. de Veciana and R. Baldick, Pricing multi-service networks, Technical Report, University of Texas, Austin, July 1994.
- [27] J. Wroclawski. Specification of the Controlled-Load Network Element Service, Internet Draft, 1995. <ftp://ftp.ietf.org/internet-drafts/draft-ietf-intserv-ctrl-load-svc-05.txt>
- [28] O. Yaron and M. Sidi. Performance and stability of communication networks via robust exponential bounds. *IEEE/ACM Transactions on Networking*, 1(3):372–385, June 1993.