

Effective Bandwidth in High-Speed Digital Networks

Cheng-Shang Chang, *Senior Member, IEEE*, and Joy A. Thomas, *Member, IEEE*

Abstract—The theory of large deviations provides a simple unified basis for statistical mechanics, information theory and queuing theory. The objective of this paper is to use large deviation theory and the Laplace method of integration to provide an simple intuitive overview of the recently developed theory of effective bandwidth for high-speed digital networks, especially ATM networks. This includes 1) identification of the appropriate energy function, entropy function and effective bandwidth function of a source, 2) the calculus of the effective bandwidth functions, 3) bandwidth allocation and buffer management, 4) traffic descriptors, and 5) envelope processes and conjugate processes for fast simulation and bounds.

I. INTRODUCTION

THE next generation of communication networks will carry different classes of traffic, e.g., voice, video, fax and data, over the same network. Since different classes of traffic usually require different grade-of-service (GOS), an open and challenging problem for the network designer is to design schemes that integrate these different classes of traffic efficiently. One of the most interesting approaches in dealing with this problem is the recently developed theory of effective bandwidth. The effective bandwidth of a time varying source is the minimum amount of bandwidth required to satisfy its GOS, and the theory of effective bandwidth provides a method to compute (or approximate) the effective bandwidth. The concept of effective bandwidth for high-speed digital networks was first proposed independently in [36], [40], [34], where the concept was tested out for i.i.d. sources and ON-OFF sources. The general framework of the theory, including the computation of the effective bandwidth for Markov processes and other general processes and the associated calculus, was carried out in [6], [43], [10], [31], [56], [35]. (For detailed historical remarks, see e.g. [56]). Further development of the theory for traffic regulation, admission control and other applications can be found in [7], [42], [47], [24], [22], [28], among many others. The main objective of this paper is to provide an intuitive, and hopefully understandable, overview of this theory that, in some ways, parallels the development of statistical mechanics, and to briefly survey some of the applications of this theory. (We will also refer to the original papers for readers interested in formal proofs.)

Manuscript received September 30, 1994; revised April 1, 1995. This work was supported in part by the National Science Council, Taiwan, R.O.C., under Grants NSC 83-0208-M007-091 and NSC 83-0404-E007-052.

C.-S. Chang is with the Department of Electrical Engineering, National Tsing Hua University, Hsinchu 30043, Taiwan, R.O.C.

J. A. Thomas is with IBM Research Division, T. J. Watson Research Center, Yorktown Heights, NY 10598 USA.

IEEE Log Number 9412658.

Our approach in some ways is parallel to the development of statistical mechanics from classical mechanics. Boltzmann (see e.g., [60]) assumed that the probability density that a particle at time t is at position x with speed v is $f(x, v, t)$. Using this assumption, he introduced the H function to provide the link between the laws of mechanics (the equations of motions) and the laws of thermodynamics. The H function was later identified as the entropy function, which in turn inspired Shannon [52] in his derivation of information theory.

In Section II, we view a source like a particle in statistical mechanics and start from defining the “equations of motions” in a network. We assume a source that behaves as a constant rate fluid with rate α for a period of time t with probability $f(\alpha; t)$. Using this, the tail distribution of the queue length in a network can be represented by an integral of $f(\alpha; t)$ over a certain set. With the help of the recently developed theory of large deviations (see e.g., [5], [25], [27], [29], [55], [53]), the probability density function $f(\alpha; t)$ (and the density function $f(x, v, t)$ in Boltzmann’s framework [29]) can be shown to have the form of Gibb’s distribution, i.e., $\exp(-t\Lambda^*(\alpha))$, where $\Lambda^*(\alpha)$ is obtained from the Legendre transform of a function $\Lambda(\theta)$ that can be derived from the Gärtner-Ellis limit of a source [33], [29]. The functions $\Lambda(\theta)$ and $\Lambda^*(\alpha)$ are then called the “energy” function and the “entropy” function of a source, respectively.

By solving for the dominant exponent in the integral of $f(\alpha; t)$, we obtain an approximation to the queue length distribution. This corresponds to finding the minimum action path in classical mechanics. In particular, for a queue with capacity c subject to a source with the energy function $\Lambda(\theta)$, the distribution of its queue length has an exponential tail with rate θ^* , where θ^* is the unique solution of $\Lambda(\theta)/\theta = c$. In view of this equality, the function $\Lambda(\theta)/\theta$ is thus called the effective bandwidth function of the source.

In Section III, we introduce the calculus for effective bandwidth functions. Instead of following the formal arguments in [6], [8], [56], [13], we derive the calculus heuristically from the corresponding energy functions and entropy functions. The calculus includes the following network operations: 1) multiplexing independent arrivals, 2) output from a switch or a link with time varying capacity, and 3) demultiplexing or routing. Using these three rules, the effective bandwidth function in an intree network with routing can be derived inductively.

Using the calculus, we explore possible applications of the theory of effective bandwidth to bandwidth allocation and buffer management in Section IV. We consider two types of bandwidth allocation: independent bandwidth allocation and

dynamic bandwidth allocation. In the former, the bandwidth allocated for a time varying source is independent of the source. We show that the optimal *independent* bandwidth allocation sequence for a given average bandwidth is the constant bandwidth allocation. In the latter, we allow the bandwidth allocation to depend on the buffer occupancy, i.e., there is a bandwidth allocation function that depends on the buffer occupancy. For instance, if the buffer occupancy is high, then we may allocate more bandwidth. The performance of such an allocation scheme is derived by approximating the bandwidth allocation function as a piecewise linear function and using the calculus developed for the theory of effective bandwidth. These results are then used to compare two different approaches for implementing multiple GOS: *simple priority* and *cut-off threshold*. The simple priority scheme performs better at the cost of increasing implementation complexity.

In Section V, we address the problem of obtaining an estimate of the effective bandwidth function of a source. Via Taylor's expansions, the effective bandwidth function can be approximated by the four parameters: the average rate, the asymptotic variance, the peak rate and the average burst duration. These four parameters corresponds to a two-state Markov fluid whose energy function, entropy function and effective bandwidth function are known.

Though the theory of effective bandwidth yields approximations for the performance of local nodes in a network, better results including accurate estimates via fast simulations and tight bounds, could be obtained by examining more closely how buffers build up. In Section VI, we follow the approach used by Gibbs in statistical mechanics. The distribution of the space and momentum coordinates of a particle shifts from the uniform distribution to the Boltzmann distribution once the average energy is specified. Similarly, we look for the most likely distribution of a source given that the buffer builds up. This idea was previously used for fast simulations in [17], [49], [51], [32], [41], among others. In this paper, we consider possible candidates that are called the envelope processes and conjugate processes in [10]. The envelope process of a source satisfies a sample path criterion for the likelihood ratio between the source and the envelope process. The sample path criterion also establishes the connection between the entropy function in this paper and the relative entropy rate (the Kullback-Leibler distance) defined in information theory.

We conclude this paper in Section VII by discussing some issues that need to be examined more carefully before the theory of effective bandwidth is implemented. These issues include 1) traffic characterization for various real-time multimedia traffic sources such as video and voice, 2) admission control, and 3) traffic monitoring and regulation.

II. EFFECTIVE BANDWIDTHS

In general, we consider a high-speed digital network, e.g., an ATM network, that consists of a set of switches and links that operate in discrete time. The capacity of a link (or a switch) is the maximum number of cells that can be handled per unit of time. To investigate the performance of a link or a switch, a

discrete-time queueing model is adopted. Let $a(t)$ and $q(t)$ be the number of cells arriving at time t and the number of cells in the queue at time t respectively. Assume that the buffer size is infinite and that the capacity is c . Under a work-conserving policy, i.e., a policy that does not allow idling when there are cells in the queue, the queue is governed by the following recursive equation

$$q(t+1) = (q(t) + a(t+1) - c)^+ \quad (1)$$

where $(x)^+ \triangleq \max(0, x)$. For a high-speed ATM network, the unit of time is small, e.g., in the order of microseconds or smaller. We can view the arrival process as a continuous fluid with a rate process $a(t)$ and approximate the discrete-time equation in (1) by a continuous-time counterpart as follows

$$\dot{q}(t) = \begin{cases} a(t) - c, & \text{if } q(t) > 0 \\ (a(t) - c)^+, & \text{if } q(t) = 0. \end{cases} \quad (2)$$

The equation in (2) will be referred as the "equation of motion" for the queue length process in this paper. In particular, if the arrival process is a constant rate process with rate α , e.g., ATM adaptation layer (AAL) type 1, then one has from the equation of motion that $q(t) = (\alpha - c)^+ t$ if $q(0) = 0$. For this case, the analysis is easy. If the arrival process is a variable rate process, e.g., AAL type 2/3, we use the Gärtner-Ellis theorem to characterize the probability that such a source will behave like a constant rate source of rate α for a time t . This characterization, though unreasonable at first sight, yields a very simple and intuitive explanation of the theory of effective bandwidth. An discussed in the introduction, a similar assumption was made by Boltzmann in 1868 and 1872 [60] to provide an explanation of thermodynamic results. In particular, his H theorem showed that entropy is increasing in time. The probabilistic nature of the assumption was criticized by physicists at his time, but it provides an explanation of how the time-reversible laws of classical mechanics could lead to the time-irreversible laws of thermodynamics.

Analogous to Boltzmann's H theorem, the density function $f(\alpha; t)$ can be derived by the large deviation principle [30]. It is known (see e.g., Gärtner [33], Ellis [29]) if

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log E e^{\theta A(0, t)} = \Lambda(\theta) \quad (3)$$

for all $\theta \in \mathbb{R}$, and $\Lambda(\theta)$ is differentiable, then

$$f(\alpha; t) \approx e^{-t\Lambda^*(\alpha)} \quad (4)$$

where $\Lambda^*(\alpha)$ is obtained from the Legendre transform of $\Lambda(\theta)$, i.e.

$$\Lambda^*(\alpha) = \sup_{\theta} [\theta\alpha - \Lambda(\theta)]. \quad (5)$$

The approximation in (4) is up to logarithmic equivalence, i.e.

$$\lim_{t \rightarrow \infty} \frac{\log \Pr(A(0, t) \approx \alpha t)}{t} = \Lambda^*(\alpha). \quad (6)$$

The Legendre transform is also known as the convex transform. In view of the definition of the energy function in (3), $\Lambda(\theta)$ is convex. Also, from the differentiability of $\Lambda(\theta)$, one can solve the optimization problem $\sup_{\theta}[\theta\alpha - \Lambda(\theta)]$ and derive

$$\Lambda^*(\Lambda'(\theta)) = \theta\Lambda'(\theta) - \Lambda(\theta). \quad (7)$$

From the theory of the Legendre transform [29], its inverse transform has the same form. In fact, if $\Lambda^*(\alpha)$ is also differentiable, then

$$\Lambda(\Lambda^{*\prime}(\alpha)) = \alpha\Lambda^{*\prime}(\alpha) - \Lambda^*(\alpha). \quad (8)$$

Thus, one can derive these two functions from each other. These two functions are called a Legendre transform pair or convex conjugates. An instructive case occurs when $\Lambda'(\theta)$ is a continuous increasing function tending to ∞ as $\theta \rightarrow \infty$ and tending to $-\infty$ when $\theta \rightarrow -\infty$. In this case, $(\Lambda')^{-1}$ is a well-defined function, and we can write

$$\Lambda(\theta) = \int_0^{\theta} \Lambda'(x) dx \quad (9)$$

$$\Lambda^*(\alpha) = \int_0^{\alpha} (\Lambda')^{-1}(y) dy. \quad (10)$$

Drawing the graph of $\Lambda'(\theta)$, it is easy to see that that $\theta\alpha \leq \Lambda(\theta) + \Lambda^*(\alpha)$, with equality if and only if $\alpha = \Lambda'(\theta)$. From this it follows that $\Lambda^*(\alpha) = \sup_{\theta}[\theta\alpha - \Lambda(\theta)]$ and that $\Lambda(\theta) = \sup_{\alpha}[\theta\alpha - \Lambda^*(\alpha)]$. Extensions of these results to all convex functions and to the multidimensional case can be found in [30].

Following the meaning of the functions $\Lambda(\theta)$ and $\Lambda^*(\alpha)$ in statistical mechanics [30], we will call $\Lambda(\theta)$ and $\Lambda^*(\alpha)$ the “energy” function and the “entropy” function¹ of an arrival process (or source). In view of the Legendre transform and the fact that $\Lambda(0) = 0$, the entropy function $\Lambda^*(\alpha)$ is nonnegative, and strictly convex, and it has a global minimum at $\alpha^0 = \Lambda'(0)$ such that $\Lambda^*(\alpha^0) = 0$. Intuitively, α^0 is the most likely rate of the arrival process, namely the average rate. From (4), the probability that the actual rate deviates from α^0 is exponentially small.

Once $f(\alpha; t)$ is derived, we can compute the probability that there are at least x cells in the buffer using the equation of motion in (2). Thus

$$\begin{aligned} \Pr(q(t) \geq x) &= \int_{(\alpha-c)^+ t \geq x} f(\alpha; t) d\alpha \\ &\approx \int_{(\alpha-c)^+ t \geq x} e^{-t\Lambda^*(\alpha)} d\alpha. \end{aligned} \quad (11)$$

For two functions $f(x)$, $g(x)$, we will say that $f(x) \approx g(x)$ if $\lim_{x \rightarrow \infty} \frac{1}{x} \log \frac{f(x)}{g(x)} = 0$, i.e., the two functions are

¹The entropy function defined here should not be confused with the entropy rate of a discrete valued process [18] or the entropy rate of a point process [21]. In Section VI, we show that the entropy function $\Lambda^*(\alpha)$ is the relative entropy rate between a process $x(t)$ and the given process $a(t)$, where the process $x(t)$ is the closest (in relative entropy) stationary process to $a(t)$ with average value greater than α .

asymptotically equal to the first order in the exponent. Now simplifying (11) yields that

$$\begin{aligned} \sup_t \Pr(q(t) \geq x) &\approx \int_{\alpha} e^{-x \inf_{\alpha} \frac{\Lambda^*(\alpha)}{(\alpha-c)^+}} d\alpha \\ &\approx \exp\left(-x \inf_{\alpha} \frac{\Lambda^*(\alpha)}{(\alpha-c)^+}\right) \end{aligned} \quad (12)$$

where the approximation is obtained by the dominant exponent in the integration (an argument that originated with Laplace [38]). If we further assume that the steady state queue length, $q(\infty)$, exists (and we will throughout this paper), then $\sup_t \Pr(q(t) \geq x) = \Pr(q(\infty) \geq x)$ [45], [4]. Thus

$$\Pr(q(\infty) \geq x) \approx \exp\left(-x \inf_{\alpha} \frac{\Lambda^*(\alpha)}{(\alpha-c)^+}\right). \quad (13)$$

It remains to solve the minimization problem

$$\inf_{\alpha > c} \left[\frac{\Lambda^*(\alpha)}{\alpha - c} \right].$$

For the minimization problem, note from (5) that for all θ

$$\frac{\Lambda^*(\alpha)}{\alpha - c} \geq \frac{\theta\alpha - \Lambda(\theta)}{\alpha - c}. \quad (14)$$

If θ^* is a solution of $\Lambda(\theta)/\theta = c$, then

$$\inf_{\alpha > c} \left[\frac{\Lambda^*(\alpha)}{\alpha - c} \right] \geq \theta^*. \quad (15)$$

Also, if the solution of $\Lambda(\theta)/\theta = c$ is unique, then it follows from (7) that

$$\Lambda^*(\Lambda'(\theta^*)) = \theta^*\Lambda'(\theta^*) - \Lambda(\theta^*) = \theta^*(\Lambda'(\theta^*) - c). \quad (16)$$

Thus, the bound in (15) is achievable when $\alpha = \Lambda'(\theta^*)$ and we have

$$\Pr(q(\infty) \geq x) \approx e^{-\theta^* x} \quad (17)$$

if θ^* is the unique solution of $\Lambda(\theta)/\theta = c$. (For formal proofs, see [6], [56].) Since c is the capacity of a link or a switch, the function

$$a^*(\theta) = \frac{\Lambda(\theta)}{\theta} \quad (18)$$

is called the effective bandwidth function of the arrival process subject to the condition that the tail distribution of the queue length has the decay rate θ .

Since $\Lambda(\theta)$ is convex, the effective bandwidth function $a^*(\theta)$ is increasing (or nondecreasing) in θ . The function $a^*(\theta)$ converges to its average rate as $\theta \rightarrow 0$, and to its peak rate as $\theta \rightarrow \infty$. If the arrival process $a(t)$ is a finite state Markov arrival process, then $a^*(\theta)$ can be easily obtained from (3) to be the logarithm of the spectral radius of the a function of the associated transition kernel [6], [43], [31]. Typical plots for the energy function, the entropy function and the effective bandwidth function of a two state ON-OFF model can be found in [7].

III. THE CALCULUS OF EFFECTIVE BANDWIDTH

In this section, we derive a calculus for the effective bandwidth functions. These rules include 1) multiplexing independent arrivals, 2) output from a switch or a link with time varying capacity, and 3) demultiplexing or routing. We outline the derivation of these rules using the Laplace method of integration—more formal proofs can be found in [6], [13].

A. Multiplexing Independent Arrivals

Consider two independent arrival processes with rate processes $a_1(t)$ and $a_2(t)$, respectively. Let $a(t)$ be the rate process after multiplexing $a_1(t)$ and $a_2(t)$, i.e.

$$a(t) = a_1(t) + a_2(t). \quad (19)$$

Suppose that the effective bandwidth function associated with $a_i(t)$ is $a_i^*(\theta)$, $i = 1$ and 2 . Then the effective bandwidth function of $a(t)$, denoted by $a^*(\theta)$, is $a_1^*(\theta) + a_2^*(\theta)$.

This result could be explained either by the energy functions or by the entropy functions. Let $\Lambda_i(\theta)$ and $\Lambda_i^*(\alpha)$, $i = 1$ and 2 , be the energy function and the entropy function of $a_i(t)$. In view of the definition of the energy function in (3), it follows that the energy function of $a(t)$, denoted by $\Lambda(\theta)$, is the sum of the energy functions of the arrival processes being multiplexed, i.e.

$$\Lambda(\theta) = \Lambda_1(\theta) + \Lambda_2(\theta). \quad (20)$$

From the definition of the effective bandwidth function in (18), one has

$$a^*(\theta) = a_1^*(\theta) + a_2^*(\theta). \quad (21)$$

To see this from the entropy function, note that for $a(t)$ to behave as a constant rate fluid with rate α for a period time t , the arrival processes $a_1(t)$ and $a_2(t)$ must behave as constant rate fluids with rates α_1 and α_2 for a period of time t , and $\alpha_1 + \alpha_2 = \alpha$. Thus, the probability that $a(t)$ behaves as a constant rate fluid with rate α for a period time t is approximately

$$\int_{\alpha_1 + \alpha_2 = \alpha} e^{-t\Lambda_1^*(\alpha_1)} e^{-t\Lambda_2^*(\alpha_2)} d\alpha_1 d\alpha_2 \approx e^{-t \inf_{\alpha_1} [\Lambda_1^*(\alpha_1) + \Lambda_2^*(\alpha - \alpha_1)]} \quad (22)$$

where we once again choose the dominant exponent in the integral [38]. Thus, the entropy function of $a(t)$, denoted by $\Lambda^*(\alpha)$, can be derived as follows

$$\Lambda^*(\alpha) = \inf_{\alpha_1} [\Lambda_1^*(\alpha_1) + \Lambda_2^*(\alpha - \alpha_1)]. \quad (23)$$

Applying the Legendre transform yields the desired result for the energy function.

B. Output from a Switch or a Link with Time Varying Capacity

One variant of the rule for multiplexing is the problem with time varying capacity, where the capacity is independent of the input. Denote by $c(t)$ the capacity at time t (the service rate). Let $\Lambda_c(\theta)$ and $\Lambda_c^*(\alpha)$ be the corresponding energy function

and the entropy function. Now the equation of motion in (2) can be modified as follows

$$\dot{q}(t) = \begin{cases} a(t) - c(t), & \text{if } q(t) > 0 \\ (a(t) - c(t))^+, & \text{if } q(t) = 0. \end{cases} \quad (24)$$

In view of (24), this is equivalent to multiplexing $a(t)$ and $-c(t)$ and passing it to a server of with capacity 0. From (3), the energy function of $-c(t)$ is $\Lambda_c(-\theta)$. Denote by $\Lambda_a(\theta)$ (resp. $\Lambda_a^*(\alpha)$) the energy (resp. entropy) function of $a(t)$. Using the rule for multiplexing in (20) and (17) for capacity $c = 0$, one has

$$\Pr(q(\infty) \geq x) \approx e^{-\theta^* x} \quad (25)$$

if θ^* is the unique solution of $\Lambda_a(\theta) + \Lambda_c(-\theta) = 0$.

Now assume that the buffer in such a switch or a link is in its steady state at time 0. Let $b(t)$ be the (stationary) rate process of the output from such a switch or a link. Denote by $\Lambda_b(\theta)$ and $\Lambda_b^*(\alpha)$ the energy function and the entropy function of $b(t)$. If θ^* is the unique solution of $\Lambda_a(\theta) + \Lambda_c(-\theta) = 0$, then it was derived in [13] that

$$\Lambda_b^*(\alpha) = \theta^* \alpha - \sup_{\gamma \leq \alpha} [\theta^* \gamma - \Lambda_a^*(\gamma)] + \inf_{\alpha \leq \delta} \Lambda_c^*(\delta), \quad \alpha \geq \alpha^0 \quad (26)$$

$$\Lambda_b(\theta) = \sup_{\alpha \geq \alpha^0} [\theta \alpha - \Lambda_b^*(\alpha)], \quad \theta \geq 0 \quad (27)$$

where α^0 is the global minimum of $\Lambda_a^*(\alpha)$. We note that (26) and (27) are “one-sided” transforms. They only provide the information for the excursion that has a larger rate than the average rate α^0 . This information is enough for most applications in this paper since we are interested in how buffers build up.

Unlike the rule for multiplexing, it is not easy to derive this rule directly by the energy functions. We will derive this from the entropy functions instead. Note that for $b(t)$ to behave as a constant rate fluid with rate α for a period time t , the capacity process $c(t)$ must behave as a constant rate fluid with rate δ for some $\delta \geq \alpha$. The sum of the arrivals in $(0, t]$ and those already in the buffer at time 0 must equal αt . Thus, the arrival process must behave as a constant rate fluid with rate γ for some $\gamma < \alpha$ and the buffer must have $(\alpha - \gamma)t$ cells at time 0. Note from (25), the probability that there are $(\alpha - \gamma)t$ cells in the (stationary) buffer at time 0 is approximately $e^{-\theta^*(\alpha - \gamma)t}$. From these observations, the probability that $b(t)$ behaves as a constant rate fluid with rate α for a period time t is approximately

$$\int_{\delta \geq \alpha, \gamma \leq \alpha} e^{-t\Lambda_a^*(\gamma)} e^{-\theta^*(\alpha - \gamma)t} e^{-t\Lambda_c^*(\delta)} d\gamma d\delta \approx e^{-t \inf_{\delta \geq \alpha, \gamma \leq \alpha} [\Lambda_a^*(\gamma) - \theta^* \gamma + \theta^* \alpha + \Lambda_c^*(\delta)]} \quad (28)$$

Simplifying (28) yields the entropy function in (26). The energy function in (26) is then obtained by the Legendre transform. From the derivation above, we note that sometimes it is necessary to build up the buffer first in order to have a large excursion of the output. Thus, the entropy function for

the output from an empty buffer at time 0 is in general not the same as that from a stationary buffer.

As in [13], there are some cases that the energy function of the output can be obtained in closed form.

Example 3.1: (Constant capacity) In the case $c(t) = c$ for all t

$$\Lambda_b^*(\alpha) = \begin{cases} \Lambda_a^*(\alpha), & \text{if } \alpha \leq c \\ \infty, & \text{otherwise} \end{cases} \quad (29)$$

and

$$\Lambda_b(\theta) = \begin{cases} \Lambda_a(\theta), & \text{if } \theta \leq \tilde{\theta} \\ \theta c - \tilde{\theta} c + \Lambda_a(\tilde{\theta}), & \text{otherwise} \end{cases} \quad (30)$$

where $\tilde{\theta}$ is a solution of $\Lambda_a'(\theta) = c$. This was previously obtained in [10], [7], [8].

Example 3.2: (Constant arrival rate) In the case $a(t) = \alpha^0$ for all t

$$\Lambda_b^*(\alpha) = \begin{cases} \theta^*(\alpha - \alpha^0), & \text{if } \alpha^0 \leq \alpha \leq \Lambda_c'(0) \\ \theta^*(\alpha - \alpha^0) + \Lambda_c^*(\alpha), & \text{if } \alpha \geq \Lambda_c'(0) \end{cases} \quad (31)$$

where θ^* is the unique solution of $\alpha^0 = \frac{\Lambda_c(-\theta)}{-\theta}$. Also

$$\Lambda_b(\theta) = \begin{cases} \alpha^0 \theta, & \text{if } 0 \leq \theta \leq \theta^* \\ \alpha^0 \theta^* + \Lambda_c(\theta - \theta^*), & \text{if } \theta > \theta^*. \end{cases} \quad (32)$$

C. Demultiplexing or Routing

We now derive the rule for demultiplexing or routing. When an arrival process $a(t)$ passes through a demultiplexer (or router), it splits into several processes. Consider a particular output process from the demultiplexer and denote it by $b(t)$. We mark the cells in the arrival process that are routed to $b(t)$. To be precise, let $p(n) = 1$ if the n^{th} cell is routed to $b(t)$ and $p(n) = 0$ otherwise. We call the process $\{p(n), n \geq 1\}$ the routing process. We assume that the routing process is independent of the arrival process. Suppose that the effective bandwidth function associated with $a(t)$ and $p(n)$ are $a^*(\theta)$ and $p^*(\theta)$ respectively. Then the effective bandwidth function of $b(t)$, denoted by $b^*(\theta)$, is $p^*(\theta)a^*(\theta p^*(\theta))$.

To see this from the energy functions, let $A(0, t)$ be the number of cells arriving in $(0, t]$ and $B(0, t)$ be the number of cells routed to $b(t)$ in $(0, t]$. Also, let $P(0, n)$ be the number of cells routed to $b(t)$ up to the n^{th} cell. Clearly, one has $B(0, t) = P(0, A(0, t))$. In view of (3), the relation for the energy functions can be derived as follows

$$\Lambda_b(\theta) = \Lambda_a(\Lambda_p(\theta)) \quad (33)$$

where $\Lambda_a(\theta)$, $\Lambda_b(\theta)$ and $\Lambda_p(\theta)$ are the energy functions of $a(t)$, $b(t)$ and $p(n)$.

One can also derive this from the entropy functions. For $b(t)$ to behave as a constant rate fluid with rate α for a period time t , the arrival process $a(t)$ must behave as a constant rate fluid with rate γ for some γ and the routing process must behave as a constant rate fluid with rate δ for the period of γt cells such that $\gamma \delta = \alpha$. Denote by $\Lambda_a^*(\alpha)$, $\Lambda_b^*(\alpha)$ and $\Lambda_p^*(\alpha)$ the entropy functions of $a(t)$, $b(t)$ and $p(n)$. Then the probability that $b(t)$

behaves as a constant rate fluid with rate α for a period time t is approximately

$$\int_{\gamma \delta = \alpha} e^{-t \Lambda_a^*(\gamma)} e^{-\gamma t \Lambda_p^*(\delta)} d\gamma d\delta \approx e^{-t \inf_{\gamma} [\Lambda_a^*(\gamma) + \gamma \Lambda_p^*(\alpha/\gamma)]}. \quad (34)$$

Thus

$$\Lambda_b^*(\alpha) = \inf_{\gamma} [\Lambda_a^*(\gamma) + \gamma \Lambda_p^*(\alpha/\gamma)]. \quad (35)$$

The energy function in (33) is then obtained by the Legendre transform.

D. Applications to Networks

In the previous subsections, we have developed a calculus of effective bandwidth that can be used to analyze cascades of queues and intree networks fed by independent sources. In such a network, one can apply these rules inductively from the leaves to the root to derive the effective bandwidth function at each local node. Each of the nodes can be viewed as filter that transforms its input process to its output process. The performance at a local node is then obtained by (17).

Although the theory provides an elegant solution to the problem of analyzing such networks, a few caveats are in order here. First, the results are asymptotic, and though they provide the asymptotic rate of decay for the buffer overflow probability, the bounds might be very loose for finite buffer sizes. Also, the requirement of independence of the inputs to a queue makes it difficult to analyze networks with loops. In this case, though, one can use ideas of stochastic dominance to prove the stability of networks of queues [12]. And the analysis for queues with time varying capacity assumes that the capacity is independent of the input process; this therefore does not allow us to directly analyze different scheduling algorithms which allot capacity to a source depending on the input process. In the next section, we will describe some simple approaches to handle buffer allocation and bandwidth management using this theory.

One variant of intree networks is a tandem network in which cross traffic has priority. By subtracting out the bandwidth allocated to the time-varying cross traffic, this problem can be transformed into a tandem network with time varying capacities. However, a potential problem occurs when the same virtual circuits share many links along a path, at which point the independence assumption for the analysis of the time-varying capacity breaks down. Under such circumstances, one can use deterministic bounds on the peak rate and the variance of the source and use the calculus developed by [19], [20], [44] to bound the delay. Alternatively, the approach based on exponentially bounded tails of the traffic process [61] could be used to provide bounds on the delay for networks with loops.

IV. BANDWIDTH ALLOCATION AND BUFFER MANAGEMENT

In this section, we discuss how one can apply the theory of effective bandwidth to bandwidth allocation and buffer management. We discuss independent and dynamic bandwidth allocation schemes and also the partitioning of the capacity between different sources. Using these results, we compare

two different approaches for implementing multiple GOS in Section IV-D: simple priority and cut-off threshold.

A. Independent Bandwidth Allocation

As in (24), a sequence $\{c(t), t \geq 0\}$ is called an independent bandwidth allocation sequence if $\{c(t), t \geq 0\}$ is independent of the arrival process. Consider an arrival process $a(t)$ with the energy function $\Lambda_a(\theta)$. If we serve $a(t)$ with an independent bandwidth allocation sequence $c(t)$ with the energy function $\Lambda_c(\theta)$, then it follows from (25) that

$$\Pr(q(\infty) \geq x) \approx e^{-\theta^* x} \quad (36)$$

if θ^* is the unique solution of

$$\Lambda_a(\theta) + \Lambda_c(-\theta) = 0. \quad (37)$$

Let $\Lambda_c^*(\alpha)$ be the entropy function of $c(t)$ and α_c^0 be the global minimum of $\Lambda_c^*(\alpha)$, i.e., $\Lambda_c^*(\alpha_c^0) = 0$. We argue that one can do better (or at least as good) if the constant rate bandwidth allocation sequence $c(t) = \alpha_c^0$ for all t is used. In view of the Legendre transform

$$\Lambda_c(-\theta^*) \geq -\theta^* \alpha - \Lambda_c^*(\alpha) \quad (38)$$

for all α . In conjunction with (38) and (37)

$$a^*(\theta^*) = \frac{\Lambda_a^*(\theta^*)}{\theta^*} \leq \alpha_c^0 + \frac{\Lambda_c^*(\alpha_c^0)}{\theta^*} = \alpha_c^0. \quad (39)$$

Since $a^*(\theta)$ is increasing, we can have a larger decay rate if the constant rate bandwidth allocation sequence $c(t) = \alpha_c^0$ is used. Thus, the best *independent* bandwidth allocation sequences are the constant rate sequences.

B. Dynamic Bandwidth Allocation And Input Regulation

In this section, we allow the bandwidth allocation sequence in Section IV-A to depend on the arrival process. Specifically, we call a function $\mu(\gamma)$, $0 \leq \gamma \leq 1$, a dynamic bandwidth allocation function if $\mu(\gamma)$ amount of bandwidth is allocated when the buffer occupancy is γx . Since we have learned from Section IV-A that constant rate allocation sequences are optimal among independent allocation sequences, adding randomness to $\mu(\gamma)$ does not improve the performance of the tail distribution of queue length. Thus, it suffices to consider $\mu(\gamma)$ as a deterministic function.

For simplicity, let us start from a piecewise linear function with $\mu(\gamma) = c_1$ for $\gamma \leq \gamma_1$ and $\mu(\gamma) = c_2$ for $\gamma > \gamma_1$. Consider an arrival process $a(t)$ with the energy function $\Lambda_a(\theta)$, the entropy function $\Lambda_a^*(\alpha)$ and the effective bandwidth function $a^*(\theta)$. For the buffer to exceed x at time t , the buffer must exceed $\gamma_1 x$ at some time $t_1 \leq t$. It then takes $t - t_1$ to build up another $(1 - \gamma_1)x$ cells. Thus, the problem can be separated into two parts: the first part is for the buffer to build up $\gamma_1 x$ cells with capacity c_1 and the second part is for the buffer to build up $(1 - \gamma_1)x$ cells with capacity c_2 . Only when these two events happen, the buffer can exceed x . Using (3) for these two parts, one has

$$\Pr(q(\infty) \geq x) \approx e^{-x(\gamma_1 \theta_1^* + (1-\gamma_1)\theta_2^*)} \quad (40)$$

where θ_i^* , $i = 1$, and 2, are the unique solutions of $a^*(\theta) = c_i$. (Though our derivation is heuristic, it could be made rigorous using scaling and the method in [9].)

In general, we can rewrite (40) as follows

$$\Pr(q(\infty) \geq x) \approx e^{-x \int_0^1 \theta_\gamma^* d\gamma} \quad (41)$$

where θ_γ^* is the unique solution of $a^*(\theta) = \mu(\gamma)$. In view of (41), one might define the "effective bandwidth" for this problem as $a^*(\bar{\theta})$, where

$$\bar{\theta} = \int_0^1 \theta_\gamma^* d\gamma. \quad (42)$$

A variant of (41) is the input regulation, in which the effective bandwidth function also depends on the buffer size. To be precise, let $a_\gamma^*(\theta)$ be the effective bandwidth function when the buffer occupancy is γx , and θ_γ^* be the unique solution of $a_\gamma^*(\theta) = \mu(\gamma)$. Following the same argument, one can show that (41) still holds. A typical scenario for the input regulation is as follows: when the buffer occupancy exceeds a certain fraction of the total amount of buffer, a bandwidth renegotiation request is sent to the access point where a source enters a network. The regulator at the access point then adjusts its operation to meet the request. For other related regulator and control problems, we refer to [57], [58], [59] and references therein.

C. Capacity Partition And Buffer Sharing

Consider the problem of superposition of two independent arrival processes $a_1(t)$ and $a_2(t)$. Suppose the capacity is c . We argue that the performance of partitioning the capacity according to the effective bandwidths is the same as that of complete sharing (up to logarithmic equivalence). From the rule for multiplexing and (17), it follows that

$$\Pr(q(\infty) \geq x) \approx e^{-\theta^* x} \quad (43)$$

if θ^* is the unique solution of $a_1^*(\theta) + a_2^*(\theta) = c$. If we partition the capacity c in the way that $a_1(t)$ is assigned $a_1^*(\theta^*)$ amount of bandwidth and $a_2(t)$ is assigned $a_2^*(\theta^*)$ amount of bandwidth, then one has

$$\Pr(q_1(\infty) \geq x) \approx \Pr(q_2(\infty) \geq x) \approx e^{-\theta^* x} \quad (44)$$

where $q_1(\infty)$ and $q_2(\infty)$ are the number of cells of $a_1(t)$ and $a_2(t)$ in the steady state. Denote by $q_p(\infty)$ the total number of cells in the buffer in the steady state under the partition. Thus, if the buffer is still shared by these two arrival processes

$$\Pr(q_p(\infty) \geq x) \approx \int_{x_1} \Pr(q_1(\infty) \geq x - x_1) \cdot \Pr(q_2(\infty) = x_1) dx_1 \approx e^{-\theta^* x}. \quad (45)$$

In view of (43) and (45), the performance of partitioning the capacity according to the effective bandwidth functions and that of complete sharing are asymptotically identical. The intuition behind this is that congestion occurs when the input rate of $a_1(t)$ (resp. $a_2(t)$) is larger than $a_1^*(\theta^*)$ (resp. $a_2^*(\theta^*)$). In that case, the whole capacity is fully utilized and partitioning the capacity according to the effective bandwidth functions behaves as complete sharing when congestion occurs. However,

we note that if there is only a finite amount of buffer and one further partitions the buffer for each arrival process, then the performance is much worse than that of complete sharing. For example, if we assign x_1 to $a_1(t)$ and $x - x_1$ to $a_2(t)$. Then the probability of a cell loss of $a_1(t)$ (resp. $a_2(t)$) is approximately $e^{-\theta^* x_1}$ (resp. $e^{-\theta^*(x-x_1)}$), which is much larger than $e^{-\theta^* x}$ in the case with buffer sharing.

D. Implementing Multiple GOS

Suppose that the two arrival processes of the previous subsection require different loss probabilities. Assume $a_1(t)$ requires the loss probability in the order of $e^{-\theta_1^* x}$ for some $\theta_1^* > \theta^*$, where θ^* is the unique solution of $a_1^*(\theta) + a_2^*(\theta) = c$. Let θ_2^* be the unique solution of $a_2^*(\theta) = c$. We assume that $\theta_1^* < \theta_2^*$ so that the GOS for $a_1(t)$ is implementable. Consider the following two implementations.

- 1) *Simple priority*: Assign the arrival process $a_1(t)$ the priority to use capacity and buffer. Thus, the loss probability of $a_1(t)$ is approximately $e^{-\theta_1^* x}$, and the GOS of $a_1(t)$ is satisfied. Since we still have buffer sharing, a cell loss for $a_2(t)$ occurs when the total number of cells from both arrival processes exceeds x . This probability is once again approximately $e^{-\theta^* x}$.
- 2) *Cut-off threshold*: Let $\gamma_1 = (\theta_2^* - \theta_1^*) / (\theta_2^* - \theta^*)$. Before the buffer exceeds $\gamma_1 x$, we accept the cells from both processes. After the buffer exceeds $\gamma_1 x$, we only accept $a_1(t)$. Clearly, the loss probability of $a_2(t)$ is approximately $e^{-\theta^* \gamma_1 x}$. Using (41), the loss probability of $a_1(t)$ is approximately $e^{-\theta_1^* x}$.

In comparison of these two schemes, the simple priority scheme performs better for both arrival processes. This is at the cost of the complexity of implementing the priority scheme. When the buffer is full and there is a cell arrival of $a_1(t)$, one needs to look for a cell of $a_2(t)$ in the buffer and discard it. In an ATM network, this corresponds to checking the cell loss priority (CLP) bits of the cells in the buffer. On the contrary, in the second scheme all the cells from $a_2(t)$ are rejected when the buffer exceeds $\gamma_1 x$.

V. TRAFFIC DESCRIPTORS

Up to this point, we only know one can derive the energy function from (3) and then use that to derive the effective bandwidth function and the entropy function. This is fine if there is already a good mathematical model for an arrival process. If there does not exist a mathematical model for an arrival process, how can one obtain the effective bandwidth function? In view of (3), direct estimation seems difficult due to the exponential term in the expectation. In [7], a four parameter traffic descriptor was proposed to approximate the effective bandwidth function. The basic idea is to approximate the effective bandwidth function by Taylor's expansions (to the first order) at both $\theta = 0$ and $\theta = \infty$, i.e.

$$a^*(\theta) = \eta_1 + \eta_2 \theta + O(\theta^2), \quad a^*(\theta) = \eta_3 - \eta_4 \frac{1}{\theta} + O\left(\frac{1}{\theta^2}\right). \quad (46)$$

The first parameter η_1 turns to be the average rate, and $2\eta_2$ is the asymptotic variance. The third parameter η_3 is the peak rate, and $1/\eta_4$ is the average burst duration. Though these four parameters are intuitively clear, it is not easy to measure them efficiently, especially η_2 that requires measuring the second moment. More research and experiments along this line are needed.

Suppose these four parameters are available for an arrival process. One can then map these four-parameter traffic descriptors to a (continuous time) two-state Markov fluid [2] that is also characterized by four parameters: $\gamma_1, \gamma_2, \lambda_2$ and λ_1 , where γ_1 (resp. γ_2) is the transition rate from state 1 to state 2 (resp. from state 2 to 1) and λ_1 (resp. λ_2) is the rate of the fluid at state 1 (resp. 2). The mapping can be found in [7]. The energy function for such a Markov fluid is shown at the bottom of this page. The effective bandwidth function is simply $\Lambda_a(\theta)/\theta$ and the entropy function can be obtained by taking the Legendre transform.

The use of a simple parametric model allows one to precalculate acceptance regions for different qualities of service. In this case, the network controller would look at the superposition of the input parameters and immediately see whether the required grade of service conditions are met. Practical network control algorithms would involve some variation on this procedure.

We note that the three parameters, the average rate, the peak rate and the average burst duration, corresponds to a two-state Markov fluid with $\lambda_1 = 0$. This was previously discussed in [36]. For the results in heavy traffic ($\theta \rightarrow 0$), only the first two parameters are needed (see [54], [3], [16] and references therein).

VI. ENVELOPE PROCESSES AND CONJUGATE PROCESSES

One might wonder if a mathematical model for an arrival process is available, could better results than the approximation in (17) be derived using the theory of effective bandwidth? While the theory of effective bandwidth provides the asymptotic rate of buffer overflow probability, a more accurate estimate is needed for practical situations. The theory can be extended using two approaches: fast simulations and bounds. The approach adopted here follows the development in [10], [9], [11]. For other related work in fast simulation, see e.g., [17], [49], [51], [32], [41].

As discussed in Section II, when the buffer builds up, the arrival process behaves as a constant rate fluid with a rate larger than the capacity. If the arrival process has the entropy function $\Lambda_a(\theta)$ and the capacity is c , the effective rate is $\Lambda'(\theta^*)$, where θ^* is the unique solution of $\Lambda_a(\theta)/\theta = c$ (for formal proofs, see e.g., [1], [26], [8]). The method follows the procedure developed by Gibbs in statistical mechanics [60] or the conditional limit theorem in information theory [18]. We will look for the most likely distribution of an arrival process given that the buffer builds up. It turns out the the appropriate class of processes are the θ -envelope and the

$$\Lambda_a(\theta) = \frac{1}{2} \left(-(\gamma_1 + \gamma_2) + (\lambda_1 + \lambda_2)\theta + \sqrt{(\gamma_1 + \gamma_2 - (\lambda_1 + \lambda_2)\theta)^2 - 4(\lambda_1 \lambda_2 \theta^2 - \gamma_1 \lambda_2 \theta - \gamma_2 \lambda_1 \theta)} \right). \quad (47)$$

θ -conjugate processes [10], [9], [11], which are defined in terms of bounds on the likelihood ratio between the process and its envelope process [10]. Examples of processes that have conjugate processes include Markov arrival processes and autoregressive processes. The conjugate process of a Markov arrival process is also a Markov arrival process with an exponentially twisted probability transition function [10], [51].

Some intuition about the θ -conjugate process can be gained by considering the relative entropy rate [18] between processes. Let $D(\{x_1(t)\}||\{x_2(t)\})$ denote the relative entropy rate between the two processes $\{x_1(t), t \geq 1\}$ and $\{x_2(t), t \geq 1\}$, i.e.

$$D(\{x_1(t)\}||\{x_2(t)\}) = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{\mathbf{x}} \Pr((x_1(1), \dots, x_1(t)) = \mathbf{x}) \cdot \log \left(\frac{\Pr((x_1(1), \dots, x_1(t)) = \mathbf{x})}{\Pr((x_2(1), \dots, x_2(t)) = \mathbf{x})} \right) \quad (48)$$

where $\mathbf{x} = (x_1, \dots, x_t)$ (provided that the limit in (48) exists). It can then be shown [10] that if the effective bandwidth of the process is increasing and differentiable, and if $\tilde{a}(t; \theta)$ is the θ -conjugate process for $a(t)$, then it is the closest process in the sense of relative entropy among the class of stationary processes $x(t)$ with $Ex(t) \geq \Lambda'_a(\theta)$, i.e.

$$D(\{x(t)\}||\{a(t)\}) \geq D(\{\tilde{a}(t; \theta)\}||\{a(t)\}) = \Lambda_a^*(\Lambda'_a(\theta)). \quad (49)$$

This establishes the connection between the entropy function in this paper and the relative entropy rate defined in information theory.

Now consider the queue in (1) starting from an empty buffer and the event \mathcal{E}^x that the buffer exceeds x before it returns to empty. To obtain an estimate for the probability of this event, one can run N independent simulations and form the estimate $\frac{1}{N} \sum_{n=1}^N \mathbf{1}_{\mathcal{E}^x}(\omega_n)$, where ω_n is the sample path of the n^{th} simulation. Suppose that the average rate is smaller than the capacity and the buffer is large, then this event is rare and one requires many runs to achieve a reasonably good estimate of this probability. To speed up the occurrences of the event, one can run the N simulations with another process $\tilde{a}(t)$ instead and form the estimate $\frac{1}{N} \sum_{n=1}^N L(\omega_n) \mathbf{1}_{\mathcal{E}^x}(\omega_n)$, where $L(\omega_n)$ is the likelihood ratio between the two processes. This technique is known as importance sampling in the literature [37].

If the arrival process has a family of θ -conjugate processes, then the fastest speed up for a node with capacity c is achieved by simulating using the θ^* -conjugate process $\tilde{a}(t; \theta^*)$, where θ^* is the solution of $a^*(\theta) = c$. This change of measure is shown to be asymptotically optimal in the sense that no other importance sampling technique can do better when the buffer is large [10]. This methodology has been recently applied to Markov fluid sources in [46].

The notion of envelope processes can also be used to obtain stronger bounds [11]. Using the bounding technique developed in [12] and the martingale properties of the likelihood ratio, the bounds can be extended tointree networks with routing in [11]. Also, a similar bound for Markov processes in a single queue can be found in the recent paper [39].

VII. CONCLUSIONS

The theory of large deviations underlies the derivation of results in statistical mechanics [29] and can be used to derive results in information theory [5], [25]. Parallel to the development of statistical mechanics, we have provided an overview of the recently developed theory of effective bandwidth for high-speed digital networks. We started from defining the equations of motion in a network and then identified the appropriate energy function, entropy function and effective bandwidth function of a source. This led to a calculus for effective bandwidth functions. The calculus included the network operations for multiplexing, input-output relation through a node, and demultiplexing. We then applied the calculus to bandwidth allocation and buffer management. Two schemes for implementing multiple GOS, simple priority and cut-off threshold, were discussed. We also proposed the four parameters, the average rate, the asymptotic variance, the peak rate and the average burst duration, as the traffic descriptors for approximating the effective bandwidth functions. To obtain better results than approximations, we addressed the notions of envelope processes and conjugate processes that could be used for fast simulations and bounds.

So far, we have discussed some theoretical issues and some practical issues of the theory of effective bandwidth. However, at the current stage, it is not clear how the theory will be implemented. As discussed in Section IV-D, there are many options in implementing the theory for real networks. Also, as pointed out in [15], one might squeeze more sources into networks if accurate mathematical models for sources are available. Thus, in order to apply the theory, we think the following issues should be examined more carefully.

- *Traffic Characterization for Various Real-Time Multimedia Traffic Sources Such as Video and Voice:* Though we have suggested the four parameters for approximating the effective bandwidth function in Section V, it is not the only approach. Accurate mathematical modeling could be obtained by other methods, e.g., hidden Markov models [50].
- *Admission Control:* Using the theory of effective bandwidth, time varying sources in a scheme with shared buffers might be viewed as constant rate sources with the rates being their effective bandwidths. This shields the network layer (and higher layers) from the impact of the evolution to high-speed networks, e.g., ATM networks. Thus, the change to existing networks can be minimized. However, the problem becomes more complex in a scheme with distributed buffers, e.g., a wireless environment. For such a scheme, sophisticated scheduling methods are required, e.g., group randomly addressed polling [14] and generalized processor sharing [48], [23].
- *Traffic Monitoring and Regulation:* At the access point, the network needs to ensure the conformity of the effective bandwidth function of each source. For this, one needs to develop the concept of "filtering." The theory of effective bandwidth suggests that one can implement a filter either in the "energy" domain or in

the "entropy" domain. Research along this line can be found in [7], [42], [24].

The theory of effective bandwidth continues to be an active area of research and will continue to influence the design, simulation and analysis of future high-speed networks.

REFERENCES

- [1] V. Anantharam, "How large delays build up in a GI/G/1 queue," *Queue. Syst.*, vol. 5, pp. 345-368, 1988.
- [2] D. Anick, D. Mitra, and M. M. Sondhi, "Stochastic theory of a data handling system with multiple sources," *Bell Syst. Tech. J.*, vol. 61, pp. 1871-1894, 1982.
- [3] J. Abate and W. Whitt, "A heavy-traffic expansion for asymptotic decay rates of tail probabilities in multi-channel queues," preprint.
- [4] F. Baccelli and P. Bremaud, *Elements of Queueing Theory*. New York: Springer-Verlag, 1994.
- [5] J. Bucklew, *Large Deviation Techniques in Decision, Simulation and Estimation*. New York: Wiley, 1990.
- [6] C. S. Chang, "Stability, queue length and delay of deterministic and stochastic queueing networks," *IEEE Trans. Automat. Contr.*, vol. 39, pp. 913-931, 1994.
- [7] ———, "Approximations of ATM networks: effective bandwidths and traffic descriptors," IBM RC 18954, 1993.
- [8] ———, "Sample path large deviations andintree networks," to appear in *Queueing Systems*.
- [9] C. S. Chang, P. Heidelberger, and P. Shahabuddin, "Fast simulation of packet loss rates in a shared buffer communications switch," IBM RC 19250, 1993.
- [10] C. S. Chang, P. Heidelberger, S. Juneja, and P. Shahabuddin, "Effective bandwidth and fast simulation of ATMintree networks," *Perform. Eval.*, vol. 20, pp. 45-66, 1994.
- [11] C. S. Chang and J. Cheng, "Computable exponential bounds forintree networks with routing," *IEEE INFOCOM'95*, Boston, pp. 197-204.
- [12] C. S. Chang, J. A. Thomas, and S. -H. Kiang, "On the stability of open networks: an unified approach by stochastic dominance," *Queue. Syst.*, vol. 15, pp. 239-260, 1994.
- [13] C. S. Chang and T. Zajic, "Effective bandwidths of departure processes from queues with time varying capacities," in *Proc. IEEE INFOCOM'95*, Boston, pp. 1001-1009.
- [14] K. C. Chen and C. H. Lee, "Group randomly addressed polling for multi-access wireless data networks," *IEEE Int. Conf. Commun.*, New Orleans, 1994, pp. 913-917.
- [15] G. L. Choudhury, D. M. Lucantoni, and W. Whitt, "Squeezing the most out of ATM," preprint, 1993.
- [16] G. L. Choudhury and W. Whitt, "Heavy-traffic asymptotic expansions for the asymptotic decay rates in the BMAP/GI/1 queue," preprint.
- [17] M. Cottrell, J. -C. Fort, and G. Malgouyres, "Large deviations and rare events in the study of stochastic algorithms," *IEEE Trans. Automat. Contr.*, vol. 28, pp. 907-920, 1983.
- [18] T. M. Cover and J. A. Thomas, *Elements of Inform. Theory*. New York: Wiley, 1991.
- [19] R. L. Cruz, "A calculus for network delay, Part I: Network elements in isolation," *IEEE Trans. Inform. Theory*, vol. 37, pp. 114-131, 1991.
- [20] ———, "A calculus for network delay, Part II: Network analysis," *IEEE Trans. Inform. Theory*, vol. 37, pp. 132-141, 1991.
- [21] D. J. Daley and D. Vere-Jones, "An introduction to the theory of point processes," New York: Springer-Verlag, 1988.
- [22] G. de Veciana, C. Courcoubetis and J. Walrand, "Decoupling bandwidths for networks: A decomposition approach to resource management," submitted to *IEEE/ACM Trans. Networking*, 1993.
- [23] G. de Veciana and G. Kesidis, "Bandwidth allocation for multiple quality of service using generalized processor sharing," SCC Tech. Rep. 94-01, Univ. Texas at Austin, 1994.
- [24] G. de Veciana and J. Walrand, "Effective Bandwidths: Call admission, traffic policing & filtering for ATM networks," submitted to *IEEE/ACM Trans. Networking*, 1993.
- [25] A. Dembo and O. Zeitouni. *Large Deviations Techniques and Applications*. Boston, MA: Jones and Barlett Publishers, 1992.
- [26] A. Dembo and T. Zajic, "Large deviations from empirical mean and measure to partial sums processes," preprint, 1993.
- [27] P. Dupuis and R. S. Ellis, "A weak convergence approach to the theory of large deviations," LCDS Rep. #93-6, Brown Univ., RI 1993.
- [28] N. G. Duffield, "Exponential bounds for queues with Markovian arrivals," preprint, 1993.
- [29] R. S. Ellis, "Large deviations for a general class of random vectors," *Annu. Probab.*, vol. 12, pp. 1-12, 1984.
- [30] ———, *Entropy, Large Deviations, and Statistical Mechanics*. New York: Springer-Verlag, 1985.
- [31] A. I. Elwalid and D. Mitra, "Effective bandwidth of general Markovian traffic sources and admission control of high speed networks," *IEEE/ACM Trans. Networking*, vol. 1, pp. 329-343, 1993.
- [32] M. R. Frater, T. M. Lenon, and B. D. O. Anderson, "Optimally efficient estimation of the statistics of rare events in queueing networks," *IEEE Trans. Automat. Contr.*, vol. 36, pp. 1395-1405, 1991.
- [33] J. Gärtner, "On large deviations from invariant measure," *Theory Prob. Appl.*, vol. 22, pp. 24-39, 1977.
- [34] R. J. Gibbens and P. J. Hunt, "Effective bandwidths for the multi-type UAS channel," *Queue. Syst.*, vol. 9, pp. 17-28, 1991.
- [35] P. W. Glynn and W. Whitt, "Logarithmic asymptotics for steady-state tail probabilities in a single-server queue," *J. Appl. Prob.*, vol. 31A, pp. 131-156, 1994.
- [36] R. Guérin, H. Ahmadi and M. Naghshineh, "Equivalent capacity and its application to bandwidth allocation in high-speed networks," *IEEE J. Select. Areas Commun.*, vol. 9, pp. 968-981, 1991.
- [37] P. Heidelberger, "Fast simulation of rare events in queueing and reliability models," Yorktown Heights, New York, IBM Res. Rep. RC 19028, 1993.
- [38] P. Henrici, *Applied and Computational Complex Analysis*. New York: Wiley, 1977, vol. 2.
- [39] Z. Liu, P. Nain, and D. Towsley, "Exponential bounds with application to call admission," preprint, 1994.
- [40] F. P. Kelly, "Effective bandwidths at multi-class queues," *Queue. Syst.*, vol. 9, pp. 5-16, 1991.
- [41] G. Kesidis and J. Walrand, "Quick simulation of ATM buffers," in *Proc. IEEE CDC'92 Conf.*, Tucson, Arizona, 1992, pp. 1018-1019.
- [42] ———, "Traffic policing and enforcement of effective bandwidth constraints in ATM networks," preprint, 1993.
- [43] G. Kesidis, J. Walrand, and C. S. Chang, "Effective bandwidths for multiclass Markov fluids and other ATM sources," *IEEE/ACM Trans. Networking*, vol. 1, pp. 424-428, 1993.
- [44] J. F. Kurose, "On computing per-session performance bounds in high-speed multi-hop computer networks," *SIGMETRICS and PERFORMANCE'92*.
- [45] R. M. Loynes, "The stability of a queue with nonindependent interarrival and service times," in *Proc. Camb. Phil. Soc.*, vol. 58, pp. 497-520, 1962.
- [46] M. Mandjes and A. Ridder, "Finding the conjugate of Markov fluid processes," preprint, 1994.
- [47] R. Nagarajan, J. Kurose and D. Towsley, "Local allocation of end-to-end quality-of-service in high speed networks," preprint.
- [48] A. K. Parekh and R. G. Gallager, "A generalized processor sharing approach to flow control in integrated service networks: The single-node case," *IEEE/ACM Trans. Networking*, vol. 1, pp. 344-357, 1993.
- [49] S. Parekh and J. Walrand, "A quick simulation method for excessive backlogs in networks of queues," *IEEE Trans. Aut. Contr.*, vol. 34, pp. 54-66, 1989.
- [50] L. R. Rabiner and B. H. Juang, "An introduction to hidden Markov models," *IEEE ASSP Mag.*, pp. 4-16, Jan, 1986.
- [51] J. S. Sadowsky, "Large deviations and efficient simulation of excessive backlogs in a GI/G/m queue," *IEEE Trans. Automat. Contr.* vol. 36, pp. 1383-1394, 1991.
- [52] C. E. Shannon, "A mathematical theory of communication," *Bell System Tech. J.*, vol. 27, (pt. 1) pp. 379-423, (pt. 2) pp. 623-656, 1948.
- [53] A. Shwartz and A. Weiss, *Large Deviations for Performance Analysis: Queues, Communication, and Computing*, 1994.
- [54] K. Sohraby, "On the asymptotic behavior of heterogeneous statistical multiplexer with applications," *INFOCOM'92*, Florence, Italy.
- [55] D. W. Stroock, *An Introduction to the Theory of Large Deviations*. New York: Springer-Verlag, 1984.
- [56] W. Whitt, "Tail probability with statistical multiplexing and effective bandwidths in multi-class queues," *Telecommun. Syst.*, vol. 2, pp. 71-107, 1993.
- [57] P. Whittle, "A risk-sensitive maximum principle," *Syst. Contr. Lett.*, vol. 15, pp. 183-192, 1990.
- [58] ———, "A risk-sensitive maximum principle: the case of imperfect state observation," *IEEE Trans. Automat. Contr.*, vol. 36, pp. 793-801, 1991.
- [59] ———, "Likelihood and cost as path integrals," *J. R. Statist. Soc., B*, vol. 53, pp. 505-538, 1991.
- [60] T. -Y. Wu, *Lectures on Non-equilibrium Thermodynamics. Kinetic Theory of Gases and Statistical Mechanics*, 1993.
- [61] O. Yaron and M. Sidi, "Performance and stability of communication networks via robust exponential bounds," *IEEE/ACM Trans. Networking*, vol. 1, pp. 372-385, 1993.



Cheng-Shang Chang (S'85-M'89-SM'93) received the B.S. degree from the National Taiwan University, Taipei, in 1983, and the M.S. and Ph.D degrees from Columbia University, New York, NY, in 1986 and 1989, respectively.

From 1989 to 1993, he was a Research Staff Member at the IBM T. J. Watson Research Center, Yorktown Heights, NY. In 1993, he joined the Department of Electrical Engineering at National Tsing Hua University, Taiwan, R.O.C., as Associate Professor. His current research interests

are concerned with queueing theory, stochastic scheduling, and performance evaluation of telecommunication networks and parallel processing systems.

Dr. Chang received an IBM Outstanding Innovation Award in 1992, and he is currently an Associate Editor for Operations Research.



Joy A. Thomas (S'83-M'90) received the B. Tech. degree in electrical engineering from the Indian Institute of Technology, Madras, in 1984, and the M.S. and Ph.D degrees in electrical engineering from Stanford University, Stanford, CA, in 1986 and 1990, respectively.

He is currently a Research Staff Member at the IBM T. J. Watson Research Center, Hawthorne, NY. He is the coauthor (with Prof. Thomas Cover) of the textbook, *Elements of Information Theory*, (Wiley, 1991). His current research interests are in data

compression, and the connections between information theory and queueing theory.

Dr. Thomas received the Indian National Talent Search Scholarship for 1978-1984, the IEEE Chareles LeGeyt Fortescue Fellowship for 1984-1985, and the IBM Graduate Fellowship for 1987-1990.