# Estimation of Available Bandwidth for an Unidentified Queueing System

Seung Yeob Nam, Sunggon Kim, and Dan Keun Sung

Department of Electrical Engineering and Computer Science, KAIST,

Daejeon, 305-701, KOREA

Korea Advanced Institute of Science and Technology Technical Report KAIST-CNR-06-01

June 14, 2006

### Abstract

This paper is concerned with estimating the available bandwidth of a single-server queueing system whose service rate and the input traffic load are not known in advance, called an *unidentified queueing* system in this paper. In order to estimate the available bandwidth, we propose a probing method called a *minimal-backlogging method* and propose two statistics. The first statistic is based on the delay of each probing packet and the second statistic is based on the amount of probing packets served in a specific time interval. We first show that an M/G/1 queueing system is stable when probing packets are sent to the system according to the minimal-backlogging method. We also show that the available bandwidth can be estimated by using either of the two statistics if the probing packets are sent to the queueing system by the minimal-backlogging method. We also estimate the available bandwidth for a local server that is connected to the probing source node with non-zero delay as an application of the theory developed for a single-server queue. We evaluate the accuracy of the proposed available bandwidth estimation scheme numerically under a Poisson and a self-similar traffic loads.

*Keywords*: M/G/1 queue; minimal backlogging; probing; available bandwidth; G/G/1 queue; unidentified queueing system;

## I. INTRODUCTION

Estimation of the residual processing capacity called the available bandwidth for a local server such as a web server or a router is one of important issues to service providers or network operators. If a web server admits every service request without limitation, the throughput of the server and the quality of service (QoS) provided for customers can be significantly degraded [1]. The same situation can be expected for a local router in case of overload. Thus, in order to protect severe degradation in throughput and to improve QoS, it is necessary to monitor and manage the available bandwidth of the local server. The configuration information of the local server such as service rate may not be easily accessible to a monitoring entity. Even though such information is available, the service rate of each class may change over time if the local server



Figure 1. Unidentified Queueing System

serves multiple classes. Thus, in this paper, we assume that the service rate of the server is not known in advance.

In this paper, we investigate how to estimate the available bandwidth of a queueing system with an unknown service rate. Fig. 1 shows a queueing system of interest. C and  $\lambda$  denote the service rate and the arrival rate of packets except probing packets, respectively. Let L denote the average size of packets except probing packets. Then, for the queueing system, available bandwidth  $C_a$  is defined as

$$C_a = C(1-\rho),$$

where  $\rho = \lambda L/C$ . If all of the parameters C,  $\lambda$  and L representing a queueing system are unknown, this system is said to be *unidentified* in this paper. We propose a new method to estimate the available bandwidth  $C(1 - \rho)$  of an unidentified queueing system.

Sharma and Mazumdar [2] considered a similar problem. They investigated the problem of estimating the traffic intensity of a local node by sending a probing traffic stream. Thus, a queue receives two streams of traffic: one is a probing stream of local user and the other is the data traffic stream obtained by superposition of all the data traffic passing through the node. Their result for the estimation of available bandwidth is valid under the assumption that the total input load of the probing and non-probing streams is less than 1. However, this assumption is not reasonable. If we do not know the available bandwidth, then it is difficult to find the appropriate probing rate which makes the total traffic load less than 1. They also assume that the service time of probing traffic stream or cross traffic stream is known in advance. For the model of Fig. 1, this is equivalent to assuming that C is known in advance. On the contrary, we propose a method to estimate the available bandwidth under an assumption of unknown C.

Alouf et al. [3] developed inference models based on finite capacity single server queues for estimating the buffer size and the intensity of cross traffic. They consider two inference models based on M/M/1/K and M/D/1/K queues. They also assume that the service rate  $\mu$ is known and the estimation performance is usually evaluated for congested queues since their relatively good estimation scheme is based on loss probability. However, congestion may not occur frequently and congesting a queue for estimation purpose may significantly degrade the QoS of cross traffic.

Recently, estimation problems of available bandwidth on an end-to-end internet path have received a lot of attentions and many techniques have been proposed [4–11]. Most of them are based either on the probe gap model (PGM), which exploits the property that the time gap between two successive probe packets is closely related with the amount of cross traffic at the bottleneck node, or the probe rate model (PRM), which exploits the concept of self-induced congestion [11].

In this paper, we propose a probing method based on the concept of a minimal backlogging and develop a theory to estimate the available bandwidth of a single-server queueing system. The concept of minimal backlogging was introduced by Knightly [12, 13] in order to define *available service* between a specific node pair in communication networks. The available service is a useful concept to understand the service capability of a network path. In order to investigate the residual service capability of a queueing system, we define the available service differently from that defined in [12] and [13]. We investigate the limiting behavior of the available service in detail and find that the available service is closely related with the available bandwidth in the limiting case, in other words, the available service normalized with time converges to the available bandwidth. Thus, we can estimate the available bandwidth of an unidentified queueing system by sending minimally backlogging probing packets and monitoring the probing packets.

Assuming that it is possible to send minimally backlogging probing packets, we propose two estimation schemes. The first scheme is to estimate the available bandwidth by measuring the delay of each probing packet, and the second scheme is to estimate the available bandwidth by measuring the total amount of probing packets served during a specific time period. The first estimation scheme is analyzed for an M/G/1 queueing system. Furthermore, both schemes can

be used to estimate the available bandwidth of a G/G/1 queueing system.

The rest of this paper is organized as follows. In Section 2, we propose a probing method called a minimal-backlogging method and investigate the stability of an M/G/1 queueing system when the minimal-backlogging method is used. We also analyze the effect of the minimal-backlogging method on the delay of non-probing traffic. In Section 3, we propose a statistic based on the delay of each probing packet to estimate the available bandwidth of the M/G/1 queueing system. In Section 4, we propose another statistic based on the amount of probing packets served in a specific time interval to estimate the available bandwidth of a G/G/1 queueing system. In Section 5, we consider estimation of available bandwidth for a local server that is separated from a probing source by a fixed delay as an application of the theory developed for a single-server queue. In Section 6, we evaluate the performance of the proposed available bandwidth estimation schemes numerically under a Poisson and a self-similar traffic loads. Finally, conclusions are presented in Section 7.

# II. MINIMAL-BACKLOGGING METHOD

We consider an M/G/1 queueing system with a First-Come-First-Served (FCFS) service policy.  $\lambda$  denotes the arrival rate of packets and L is the average packet size. Suppose that the service time of a packet is given by the packet size divided by the service rate C of the system. Let G be the service time distribution of the packets and let S be a random variable corresponding to G. Then, the traffic load to the system is  $\rho = \lambda E[S]$ , which has the same value as  $\lambda L/C$ . We assume that  $\rho < 1$  for the stability of the system. To consider the problem generally, we assume that  $G_p$ , the service time distribution of probing packets, may be different from G. We let  $S_p$  denote a random variable corresponding to  $G_p$ . We define some terminologies as follows:

*Definition 1:* A session is a sequence of packets sent to a queueing system by a user. A session is said to be in a backlogging state if there is at least one packet belonging to the session in the queueing system.

*Definition 2:* Suppose that probing packets are sent to a queueing system so that there exists one and only one probing packet in the system. This probing method is called a minimal-backlogging method.

If we send a new probing packet to a queueing system just at the departure time of the previous probing packet, then there exists one and only one probing packet in the system. Let  $X_i$ , i = 1, 2, ... be the number of non-probing packets in the system seen by the *i*-th probing packet on arrival. Suppose that we start the probing for the M/G/1 queueing system in a stationary state. Then,  $X_1$ , the number of packets in the system seen by the first probing packet, is equal to the stationary queue length in number of packets, whose moment generating function is given in [14] as

$$\Pi(z) = \frac{(1-\rho)(1-z)\tilde{G}[\lambda(1-z)]}{\tilde{G}[\lambda(1-z)] - z},$$
(1)

where  $\tilde{G}(s) = \int_0^\infty e^{-sx} dG(x)$  is the Laplace transform of G.

Clearly,  $X_{i+1}$  is the number of packets arriving during the total service time of the  $X_i$  packets and the *i*-th probing packet. Let  $N_k^i$  be the number of non-probing packets arriving during the service time of the *k*-th non-probing packet among the  $X_i$  packets and let  $N_p^i$  be the number of non-probing packets arriving during the service time of the *i*-th probing packet. Since the arrival process of non-probing packets is a Poisson process,  $N_k^i$  depends only on the service time of the *k*-th packet. Thus, for all *i* and *k*,  $N_k^i$ 's are independent and identically distributed. By the same reason, for all *i*,  $N_p^i$  are also independent and identically distributed. Now, we obtain the following relation:

$$X_{i+1} = \sum_{k=1}^{X_i} N_k + N_p,$$
(2)

where for all k,  $N_k$  is a random variable with the same distribution as  $N_1^1$  and  $N_p$  with the same distribution as  $N_p^1$ , and each random variable is independent of the others. For simplicity, we will use N instead of  $N_1^1$ .

The probing based on the minimal-backlogging method keeps the queueing server continuously busy. Thus, the probing may make the queueing system unstable. Theorem 1 answers this question.

*Theorem 1:* Let  $X_i$  be the number of packets in the system upon arrival of the *i*-th probing packet. Then,  $\{X_i, i = 1, 2, ...\}$  is an aperiodic and irreducible Markov Chain and it is positive recurrent.

*Proof.* By Eqn. (2), we can see that  $\{X_i, i = 1, 2, ...\}$  is a Markov chain. Since  $N_k, k = 1, 2, ...$  and  $N_p$  can have any nonnegative integers with a positive probability,  $\{X_i, i = 1, 2, ...\}$ 

i) 
$$|E[X_{i+1} - X_i | X_i = n]| < \infty, n = 0, 1, 2, ...$$
  
ii)  $\limsup_{n \to \infty} E[X_{i+1} - X_i | X_i = n] < 0.$ 

By conditioning on  $X_i$  in Eqn. (2), we have

$$E[X_{i+1}|X_i = n] = nE[N] + E[N_p].$$
(3)

Since N is the number of Poisson arrivals during a random time of mean E[S], it can be easily shown that  $E[N] = \lambda E[S]$ . By the similar reason,  $E[N_p] = \lambda E[S_p]$ . Then, Eqn. (3) is rewritten as

$$E[X_{i+1}|X_i = n] = n\rho + \lambda E[S_p].$$
(4)

By subtracting n from the both sides of the above equation, we have

$$E[X_{i+1} - X_i | X_i = n] = n(\rho - 1) + \lambda E[S_p].$$

Thus, for any n,  $E[X_{i+1} - X_i | X_i = n]$  is finite. From the assumption that  $\rho < 1$ , it follows that  $\lim_{n\to\infty} E[X_{i+1} - X_i | X_i = n] = -\infty.$ 

By taking expectation on  $X_i$  in Eqn. (4), we derive

$$E[X_{i+1}] = \lambda E[S_p] + \rho E[X_i], \quad i = 1, 2, \dots$$

The solution of the above recurrence relation is given by

$$E[X_i] = \frac{\lambda E[S_p]}{1 - \rho} + \rho^{i-1} \left( E[X_1] - \frac{\lambda E[S_p]}{1 - \rho} \right), \quad i = 1, 2, \dots,$$
(5)

where  $E[X_1]$  has a value of  $\lambda^2 E[S^2]/[2(1-\rho)] + \rho$ , the expected queue length of a stationary M/G/1 queueing system.

Let  $W_i$  be the waiting time of the *i*-th probing packet. By conditioning on  $X_i$ , we derive the Laplace transform of  $W_i$  as follows:

$$E[e^{-sW_i}] = \sum_{n=0}^{\infty} E[e^{-sW_i} | X_i = n] \operatorname{Pr}\{X_i = n\}$$
  
$$= \sum_{n=0}^{\infty} \tilde{G}_p(s)\tilde{G}(s)^n \operatorname{Pr}\{X_i = n\}$$
  
$$= \tilde{G}_p(s)\Pi_i(\tilde{G}(s)).$$
  
(6)

where  $\tilde{G}_p$  is the Laplace transform of  $G_p$  and  $\Pi_i(z)$  is the moment-generating function of  $X_i$ . Differentiating the above equation and substituting s = 0, we obtain

$$E[W_i] = E[S_p] + E[S]E[X_i], \quad i = 1, 2, \dots$$
(7)

From Theorem 1 we can see that the embedded Markov chain  $\{X_i\}$  has a limiting distribution. To extend this result to the queue length process of an M/G/1 queueing system probed by the minimal-backlogging method, we obtain the following theorem:

Theorem 2: Suppose that we start probing an M/G/1 queueing system according to the minimal-backlogging method. Let  $\{X(t), t \in [0, \infty)\}$  be the queue length process of the queueing system. Then,  $\{X(t)\}$  is a stable process, i.e.  $\{X(t)\}$  converges to a stationary process. *Proof.* We assume that the first probing packet is sent to the queueing system at time 0 without loss of generality. Consider the epochs  $\{\tau_1, \tau_2, \tau_3, \ldots\}$  such that there is no non-probing packet upon arrival of probing packets. Then,  $\{X(t)\}$  is a regenerative process with regeneration points of  $\{\tau_1, \tau_2, \tau_3, \ldots\}$ . In order to show that  $\{X(t)\}$  is stable, it is sufficient to show that the expectation of the length of a regeneration cycle is finite [14, Theorem 17 of Chapter 2].

Let  $a_t^p$  be the number of probing packets arriving until time t, i.e.,

$$a_t^p = 1 + \max\{n | \sum_{i=1}^n W_i \le t\}.$$
(8)

Let  $Z(t) = X_{a_t^p}$ , where  $\{X_n\}$  is the Markov chain defined in Theorem 1. Since the sojourn time of Z(t) in state k is the total sum of service times of the number of k non-probing packets and a probing packet, the sojourn time only depends on k. This implies that  $\{Z(t)\}$  is a semi-Markov process with embedded Markov chain  $\{X_n\}$ . Let  $\mu_k$  be the expectation of the sojourn time of Z(t) in state k, and  $\pi_k$  be the stationary distribution of  $\{X_n\}$ . Then,

$$\sum_{k=0}^{\infty} \pi_k \mu_k = \sum_{k=0}^{\infty} \pi_k (kE[S] + E[S_p])$$
$$= E[S]E[X_{\infty}] + E[S_p].$$

Since  $E[X_{\infty}]$  is finite by Eqn. (5),  $\sum_{k=0}^{\infty} \pi_k \mu_k$  is also finite. By [14, Theorem 9 of Chapter 4], we can see that  $\{Z(t)\}$  is positive recurrent. Thus, the expectation of  $\tau_{i+1} - \tau_i$  is finite. Now, we have shown that  $\{X(t)\}$  is a stable process.

Thus, we can know that an M/G/1 queueing system is stable when it is probed by the minimal-backlogging method. Now we investigate the impact of the minimal-backlogging method upon the performance of the non-probing traffic. Especially, we derive the average number of non-probing packets under probing analytically and compare it with the case of no probing. Let  $X'_i$  be the number of non-probing packets seen by the *i*-th departing probing packet. If probing packets are sent to the queueing system according to the minimal-backlogging method, the (i + 1)-th probing packet arrives at the queueing system upon departure of the *i*-th probing packet. Thus,  $X_{i+1} = X'_i$  for  $i = 1, 2, \cdots$ . Let  $X'_{i,j}$  be the number of non-probing packets in the system seen by the j-th departing non-probing packet  $(j \leq X'_i)$  among  $X'_i$  non-probing packets observed by the *i*-th departing probing packet. Let  $N'_{i,j}$  be the number of non-probing packets arriving during the service time of the j-th departing non-probing packet among  $X'_i$  nonprobing packets. Then, when the *j*-th non-probing packet departs from the system, only  $X'_i - j$ non-probing packets remain among  $X'_i$  non-probing packets and  $N'_{i,1} + \cdots + N'_{i,j}$  non-probing packets additionally arrive. Thus, we have  $X'_{i,j} = X'_i - j + \sum_{k=1}^j N'_{i,k}$ . Then, the following theorem gives an analytical result for the average number of non-probing packets in the system under probing.

*Theorem 3:* Let L(n) be the average number of non-probing packets observed by  $\sum_{i=1}^{n} X'_{i}$  departing non-probing packets. Then,

$$\lim_{n \to \infty} L(n) = \rho + \frac{\lambda^2 E[S^2]}{2(1-\rho)} + \frac{\lambda E[S_p^2]}{2E[S_p]} + \frac{2\lambda^2 E[S_p]E[S]}{2(1-\rho)}.$$
(9)

*Proof.* The proof is given in Appendix A.

If we let  $L = \lim_{n\to\infty} L(n)$ , then L is also equal to the average number of non-probing packets in the system due to Burke's theorem [16, p.7] and PASTA property. The summation of the first two terms of Eqn. (9) is equal to  $E[X_1]$ , i.e. the expected queue length of a stationary M/G/1queueing system. Thus, the last two terms correspond to the excess delays induced by probing packets. The excess delay is finite for  $\rho < 1$  and the delays of non-probing packets and probing packets are also numerically evaluated in Section 6.

#### III. ESTIMATION BASED ON DELAY

In this section, we investigate how to estimate the available bandwidth of an M/G/1 queueing system by measuring the delay of each probing packet sent according to the minimalbacklogging method.

Theorem 4: Let  $W_i$  be the waiting time of the *i*-th probing packet. If we fix the size of the probing packets to a constant of  $L_p$  and let  $\overline{W}_n = (W_1 + W_2 + \ldots + W_n)/n$ , then

$$\lim_{n \to \infty} E\left[\frac{\bar{W}_n}{L_p}\right] = [(1-\rho)C]^{-1}.$$

Proof. It follows from Eqn. (7) that

$$E\left[\sum_{i=1}^{n} W_i\right] = nE[S_p] + E[S]\sum_{i=1}^{n} E[X_i].$$

Since  $\lim_{i\to\infty} E[X_i] = \lambda E[S_p]/(1-\rho)$  by Eqn. (5), we obtain

$$\lim_{n \to \infty} E\left[\frac{\sum_{i=1}^{n} W_i}{n}\right] = E[S_p] + E[S]E[X_{\infty}]$$
$$= \frac{E[S_p]}{1 - \rho}.$$

Since the size of the probing packets is fixed to  $L_p$ ,  $S_p$  is equal to  $L_p/C$ , which completes the proof.

Theorem 4 says that  $\overline{W}_n/L_p$  can be a candidate for an estimator of the available bandwidth. By the following theorem and corollary, we can observe that  $\overline{W}_n/L_p$  is a good candidate.

Theorem 5: Let  $W_i$  be the waiting time of the *i*-th probing packet and let  $\overline{W}_n = (W_1 + W_2 + \dots + W_n)/n$ . Then, the variance of  $\overline{W}_n$  converges to zero with order of 1/n, moreover, for a constant *c* not depending on *n*,

$$Var[\bar{W}_n] \le \frac{c}{n}.$$

*Proof.* The proof is given in [17, Appendix].

Corollary 6: Let  $W_i$  be the waiting time of the *i*-th probing packet. If we fix the size of the probing packets to a constant of  $L_p$  and let  $\overline{W}_n = (W_1 + W_2 + \ldots + W_n)/n$ , then

$$\lim_{n \to \infty} E\left[ \left| \frac{\bar{W}_n}{L_p} - [C(1-\rho)]^{-1} \right|^2 \right] = 0.$$

*Proof.* Let  $Z_n = \overline{W}_n / L_p$ . Then, by Minkowski's inequality, we can obtain

$$E\left[\left|Z_n - [C(1-\rho)]^{-1}\right|^2\right]^{\frac{1}{2}} \le E\left[\left|Z_n - E[Z_n]\right|^2\right]^{\frac{1}{2}} + \left|E[Z_n] - [C(1-\rho)]^{-1}\right|.$$

By Theorems 4 and 5, the right hand side of the above inequality converges to zero. This completes the proof.  $\Box$ 

We can obtain the following relation by Chebychev's inequality:

$$\Pr\left(\left|\frac{\bar{W}_n}{L_p} - \frac{1}{C(1-\rho)}\right| > \varepsilon\right) \le \frac{E\left[\left|\bar{W}_n/L_p - 1/(C(1-\rho))\right|^2\right]}{\varepsilon^2}.$$

Since the right hand term of the above inequality goes to zero by Corollary 6,  $\bar{W}_n/L_p$  is a consistent estimator [18] of  $[C(1-\rho)]^{-1}$ .

# IV. ESTIMATION BASED ON THE AMOUNT OF PACKETS

In Section 3, we proposed a statistic to estimate the available bandwidth of an unidentified queueing system when the arrival process of non-probing packets is a Poisson process. We can estimate the available bandwidth by measuring the delay of each probing packet. In this section, we propose another statistic to estimate the available bandwidth of a queueing system when the arrival process of non-probing packets is a general process. The available bandwidth can be estimated by measuring the total amount of minimally backlogging probing packets that are served during a specific time period. We define the concept of *Available Service*, which is defined in a different way from that in [12, 13].

Definition 3: The available service  $\hat{Y}_{[s,t]}$  for a queueing system is the amount of probing packets ets served in interval [s,t] when probing packets are sent to the queueing system according to the minimal-backlogging method.

Before we investigate the characteristics of the available service analytically, we briefly explain why the term of *available service* is used for  $\hat{Y}_{[s,t]}$ . In case that the minimal-backlogging method is not used, an *idle period*, i.e. a time interval when the server is not busy, can exist if the load of non-probing packets are less than 1. In case that the probing packets are sent to the queueing system according to the minimal-backlogging method, there always exists at least one probing packet in the queueing system, and thus, there is no idle period during the probing time. If there is no non-probing packet in the system, probing packets will be served continuously until a new non-probing packet arrives. Thus, we can know that the amount of probing packets served in a given time interval will be at least the maximum amount of service that the server can additionally support while serving all arriving non-probing packets according to an FCFS policy. On the other hand, the available service defined in [12, 13] represents the maximum amount of service that the server can do in a given time interval.

The size of each probing packet is fixed to a constant of  $L_p$  in this section. We assume that the first probing packet is sent to the system at time 0 without loss of generality. For simplicity, we will use  $\hat{Y}_t$  instead of  $\hat{Y}_{[0,t]}$ . Then, the available service  $\hat{Y}_t$  is expressed as

$$\hat{Y}_t = L_p \cdot \max\{n | \sum_{i=1}^n W_i \le t\}.$$

Let  $Q_t$  denote the amount of packets in the queueing system at time t. Then,

$$Q_t = L_p + \sum_{k=1}^{X_t^n} L_k,$$

where  $X_t^n$  is the number of non-probing packets in the system at time t and  $L_k$  is the size of the k-th non-probing packet in the system. Let  $A_t$  be the amount of packets arriving during [0, t] and let  $Y_t$  be the amount of packets served during [0, t]. Note that  $A_t$  consists of probing packets,  $A_t^p$ , and non-probing packets,  $A_t^n$ . Then,

$$Q_t = Q_0 + A_t - Y_t = Q_0 + A_t^n + A_t^p - Y_t.$$
(10)

The following lemma and theorem say that  $\hat{Y}_t/t$  converges to  $C(1-\rho)$  in  $L^q$ .

*Lemma 7:* Let  $a_t^p$  be the number of probing packets arriving until time t. If probing packets are sent according to the minimal-backlogging method, then  $\lim_{t\to\infty} a_t^p = \infty$  almost surely (a.s.). *Proof.* Eqn. (10) is rewritten as

$$A_t^p = Q_t - Q_0 + Y_t - A_t^n.$$

Since  $Q_t \ge 0$ , we have  $A_t^p \ge Y_t - A_t^n - Q_0$ . Thus,

$$\liminf_{t \to \infty} \frac{A_t^p}{t} \ge \liminf_{t \to \infty} \frac{Y_t - A_t^n - Q_0}{t}.$$
(11)

By the assumption that the input load of non-probing packets is  $\rho$ ,  $\lim_{t\to\infty} A_t^n/t = \rho C \ a.s.$ Since the server is continuously busy during the period of probing,  $\lim_{t\to\infty} Y_t/t = C \ a.s.$  Thus, it follows from Eqn. (11) that

$$\liminf_{t \to \infty} \frac{A_t^p}{t} \ge (1 - \rho)C, \quad a.s.$$
(12)

Since  $A_t^p = L_p a_t^p$ ,  $\liminf_{t\to\infty} a_t^p/t \ge (1-\rho)C/L_p a.s.$  Then,  $\liminf_{t\to\infty} a_t^p = \infty a.s.$  because  $\rho < 1.$ 

*Theorem 8:* Let  $\hat{Y}_t$  be the available service for a G/G/1 queueing system. The size of each probing packet is fixed to a constant of  $L_p$ . Then, for  $0 < q < \infty$ ,

$$\lim_{t \to \infty} E\left[ \left| \frac{\hat{Y}_t}{t} - C(1-\rho) \right|^q \right] = 0.$$
  
7. Theorem 3.71.

Proof. The proof is given in [17, Theorem 3.7

Since  $\Pr(|\hat{Y}_t/t - C(1-\rho)| > \varepsilon) \le E[|\hat{Y}_t/t - C(1-\rho)|^2]/\varepsilon^2$  by Chebychev's inequality and the right hand term of the inequality goes to zero by Theorem 8,  $\hat{Y}_t/t$  is a consistent estimator of  $C(1-\rho)$ .

We need to note that the statistic based on delay of probing packets are closely related with the statistic based on the amount of packets. We consider the case that probing packets are sent to the queueing system according to the minimal-backlogging method. The first statistic  $\overline{W}_n/L_p$  is a consistent estimator of  $1/((1 - \rho)C)$ . Thus,  $L_p/\overline{W}_n$  is an estimator of the available bandwidth and it can be rewritten as  $nL_p/\sum_{i=1}^n W_i$ . If we consider time t until the service completion time of the *n*-th probing packet, then  $\hat{Y}_t = nL_p$  and thus, both statistics agree with each other exactly. Thus, though we showed that the statistic based on the amount of probing packets is a consistent estimator of the available bandwidth for a G/G/1 queueing system, the statistic based on packet delay can also be used to estimate the available bandwidth of a G/G/1 queueing system since the two statistics are identical as the time goes to infinity.

## V. APPLICATION TO ESTIMATION OF AVAILABLE BANDWIDTH OF A LOCAL SERVER

Thus far, we considered a problem of estimating the available bandwidth of a queueing system which is directly accessible with no access delay. However, in real situation, an unidentified queueing system may be physically separated from the probing site such that the access time



Figure 2. A measurement setup for estimation of the available bandwidth at a local server

delay is not zero due to a propagation delay of pre-processing time. Thus, we consider the problem of estimating the available bandwidth when there exists access time delay between the target queueing system, also called a local server, and a probing site. Extending the approach developed in the previous sections, we propose a scheme to estimate the available bandwidth of a local server based on the minimal backlogging concept.

Fig. 2 illustrates a measurement process for estimation of the available bandwidth of a local server. The application or machine at a measurement point A sends probing packets to the local server and receives feedback information. Probing packets sent from Node A arrive at the local server after a delay of  $D_f$ . Probing packets served by the local server return to Node A after a delay of  $D_b$ . We assume that both delays are fixed and known to the monitoring node A in advance.

Due to the delays of  $D_f$  and  $D_b$ , it is not easy to send probing packets while maintaining one and only one probing packet in the server. Thus, we attempt to maintain a minimal-backlogging condition with the following heuristic method. The proposed method is based on the idea that if probing packets are sent to the server according to the minimal-backlogging method, the inter-packet spacing between two consecutive probing packets is equal to the sojourn time of the former probing packet of the two. The proposed available bandwidth estimation method is described as follows:

- 1) The measurement node sends a probing packet to the local server and obtains the roundtrip delay  $d_0$  of the probing packet upon receiving the returning packet.
- 2) The measurement node sends the first probing packet  $p_1$  for estimation of the available

bandwidth after acquiring  $d_0$ .

- 3) Let  $p_j$  be the last probing packet that was sent toward the server and let  $v_j$  be the time when  $p_j$  was sent to the server. If the last round-trip delay value available to the measurement node is  $d_i$ , we estimate the sojourn time of  $p_j$  in the server as  $d_i D_f D_b$ , and thus, the next probing packet is sent at time  $v_j + (d_i D_f D_b)$ . Exceptionally, if the last probing packet  $p_j$  arrives before  $v_j + (d_i D_f D_b)$ , there is no probing packet in the path, especially in the server. Thus, the next probing packet is sent upon arrival of  $p_j$  in order to maintain at least one probing packet in the server.
- 4) The measurement node measures the available service  $\hat{Y}_t$  or the delay experienced by each probing packet in the server whenever a returning packet arrives, and estimates the available bandwidth by using either of the two statistics proposed in Sections 3 and 4.

Since the round-trip delay can be measured solely at Node A only if all probing packets are returned to Node A, there is no clock synchronization problem.

# VI. NUMERICAL RESULTS

In this section, we compare the performance of the two proposed statistics with that of the method proposed in [2] in terms of accuracy and the effect on the delay of non-probing packets. Sharma and Mazumdar proposed a method to estimate the utilization of the system by sending probing packets according to a Poisson process in Subsection 2.3 of their paper [2]. Then, the available bandwidth (AB) can be obtained if the service rate C of the system is known in advance. This estimation scheme is called a *Poisson probing scheme* in this paper, and this scheme is compared with our proposed scheme. In addition, we also evaluate the accuracy of the available bandwidth estimation scheme developed in Section 5 for a local server whose access time is non-zero from a probing source.

Fig. 2 shows a measurement node interconnected to an unidentified queueing system. We first consider the case of no access delay, i.e.  $D_f = D_b = 0$ . The measurement node directly connected to the queueing system sends probing packets to the queueing system by the minimal-backlogging method, i.e., the node sends a new probing packet upon arrival of the previous probing packet and calculates the values of two statistics. The measurement node bypasses every non-probing packet.

Two types of traffic patterns are used for non-probing packet traffic streams: Poisson and self-similar traffic. The traffic patterns of today's IP networks have been known to exhibit self-similarity and long-range dependence [19–21]. Neither of them can be modeled using conventional Markovian models. Thus, we use a multi-fractal model [22] to generate self-similar traffic. The Hurst parameter is 0.8. The sizes of both probing and non-probing packets are fixed to 500 bytes. The service rate (C) of the unidentified queueing system is 10 Mbps.

Figs. 3, 4 and 5 compare our proposed scheme with the Poisson probing scheme under a Poisson cross traffic load of 0.3, 0.5, and 0.7, respectively. The proposed estimation scheme is based on the minimal-backlogging method and the statistic of probing packet delay. Since both statistics based on the delay and the amount of probing packets yield an identical estimation result, we show only the result obtained from the statistic based on the delay of probing packets. The value of *Measured AB* is obtained in the queueing system by subtracting the service rate of non-probing packets from the service rate C when the probing traffic is not sent. The same traffic patterns are used for both estimation and measurement of the AB at the same load. We can observe that the estimation results of the minimal-backlogging method agree with the measured AB for all traffic loads. Furthermore, the estimation results are accurate even when the observation time duration is short. The reason is explained as follows. We know that the estimation result converges to an AB value of  $C(1 - \rho)$  when the observation time goes to infinity by Theorem 8. Let us consider a finite time interval [s, t] after start of probing. Then, the server is continuously busy for the interval [s, t] because there is at least one probing packet in the queueing system. When the server does not serve non-probing packets, the server surely serves probing packets. Thus, all unused capacity of the server is used by probing packets in any finite interval. If the probing traffic is greedy like TCP flows, then the throughput of non-probing packets may be degraded. However, since probing traffic tries to prevent from being greedy by maintaining only one probing packet in the queueing system, the AB is estimated reasonably in a finite time interval.

On the other hand, the Poisson probing scheme is inaccurate and sometimes does not converge when the probing rate is as low as 0.1 Mbps. As the probing rate increases, the accuracy improves and the convergence time decreases. However, if the probing rate increases so that



Figure 3. Comparison of the proposed scheme and the Poisson probing scheme under a Poisson traffic load of 0.3



Figure 4. Comparison of the proposed scheme and the Poisson probing scheme under a Poisson traffic load of 0.5

the aggregate input traffic load approaches to 1, the accuracy is not guaranteed as observed in Figs. 4 and 5. Especially, if the total input load exceeds 1, then the system becomes unstable and the Poisson probing scheme yields erroneous results.

Tables I, II, and III compare the performance of our proposed scheme and the Poisson probing scheme in terms of average delay under a Poisson traffic load of 0.3, 0.5, and 0.7, respectively. We can observe that the average delay of probing packets is lower than that of non-probing



Figure 5. Comparison of the proposed scheme and the Poisson probing scheme under a Poisson traffic load of 0.7

packets in case of the minimal-backlogging method. This is because the probing packets are sent to the queueing system in a special way to maintain only one probing packet in the queueing system. On the other hand, the average delay of non-probing packets is very close to that of probing packets in case of Poisson probing since both cross traffic and probing traffic follow Poisson processes. The Poisson probing scheme yields smaller delays of non-probing packets than for the minimal-backlogging method when the aggregate traffic load is rather low. However, when the aggregate traffic load exceeds 0.8 in Tables I and II or 0.9 in Table III, the Poisson probing scheme yields worse delay performance. From Figs. 3, 4 and 5, we can observe that the convergence time is rather long and thus, it is difficult to obtain an accurate estimate of available bandwidth in a short time by the Poisson probing scheme. Increasing probing rate improves convergence time and accuracy. However, high probing rates in the Poisson probing scheme may degrade the delay performance of non-probing packets, as observed in Tables I, II, and III. Thus, it is a challenging and difficult problem to find a good probing rate of the Poisson probing scheme without knowing the available bandwidth. Especially, if the aggregate input traffic load exceeds 1, the system becomes unstable and the Poisson probing scheme can not give the correct value of the available bandwidth. However, our proposed scheme based on the minimal-backlogging method does not suffer from such problems and the delay of non-probing

## TABLE I

Comparison of the proposed scheme and the Poisson probing scheme in terms of average delay (load = 0.3)

	Method	Total input load	non-probing	probing
			packet delay	packet delay
Minim	al-backlogging method	1.0	0.000856	0.000569
	probing rate = 0.1Mbps	0.31	0.000488	0.000482
Poisson	probing rate $= 0.5$ Mbps	0.35	0.000505	0.000504
probing	probing rate = 1.0Mbps	0.4	0.000531	0.000531
scheme	probing rate $= 2.0$ Mbps	0.6	0.000597	0.000598
	probing rate = 5.0Mbps	0.8	0.001203	0.001200
	probing rate = 7.0Mbps	1.0	0.048535	0.048348

#### TABLE II

Comparison of the proposed scheme and the Poisson probing scheme in terms of average delay (load =

#### 0.5)

	Method	Total input load	non-probing	probing
			packet delay	packet delay
Minim	al-backlogging method	1.0	0.001199	0.000800
	probing rate $= 0.1$ Mbps	0.51	0.000604	0.000599
Poisson	probing rate = 0.5Mbps	0.55	0.000644	0.000654
probing	probing rate = 1.0Mbps	0.6	0.000698	0.000702
scheme	probing rate = 2.0Mbps	0.7	0.000867	0.000865
	probing rate = 3.0Mbps	0.8	0.001206	0.001207
	probing rate = 5.0Mbps	1.0	0.109547	0.109673

packets are maintained stable for any load of cross traffic.

Fig. 6 compares our proposed scheme with the Poisson probing scheme under a self-similar traffic load. The sigma/mean ratio of self-similar traffic is approximately 0.5 for all traffic loads of 0.3, 0.5, and 0.7. The probing rate of the Poisson probing scheme is fixed to 1.0Mbps. First, we can observe that the traffic is even burstier than the case of Poisson traffic. The AB estimation results of the minimal-backlogging method agree well with the measured AB for all traffic loads and even for short duration of observation time. However, it takes several seconds to obtain a converged value of available bandwidth by the Poisson probing scheme.

Thus far, we have evaluated the performance of the proposed estimation scheme based on the proposed statistics and the minimal-backlogging method for an unidentified queueing system that is directly connected to the probing source without delay. We now evaluate the accuracy of the available bandwidth estimation algorithm proposed in Section 5 for a local server that is

### TABLE III

Comparison of the proposed scheme and the Poisson probing scheme in terms of average delay (load = 0.7)

	Method	Total input load	non-probing	probing
			packet delay	packet delay
Minim	al-backlogging method	1.0	0.002032	0.001344
	probing rate = 0.1Mbps	0.71	0.000902	0.000907
Poisson	probing rate = $0.5$ Mbps	0.75	0.001016	0.001018
probing	probing rate = 1.0Mbps	0.8	0.001221	0.001210
scheme	probing rate = 2.0Mbps	0.9	0.002274	0.002288
	probing rate = 3.0Mbps	1.0	0.142773	0.142999



Figure 6. Comparison of the proposed scheme and the Poisson probing scheme under a self-similar traffic load

accessible with delay from the probing source, i.e.  $D_f \ge 0$  and  $D_b \ge 0$  in Fig. 2. The local server is a queueing system with an FCFS polity and the service rate C is fixed to 10 Mbps. The value of the forward delay  $D_f$  is assumed to be the same as that of feedback delay  $D_b$ and is increased from 0 to 0.5 msec in the simulation. The sizes of both probing packets and non-probing packets are fixed to 500 bytes.

Fig. 7 shows the accuracy of the proposed AB estimation algorithm for various values of fixed delay when the self-similar traffic loads of 0.3, 0.5, and 0.7 are offered. The observation time duration is 50 seconds. In this case, we used the statistic based on the amount of probing packets. We can observe that the accuracy degrades as the value of fixed delay  $(D_f + D_b)$  increases. The reason is that long response time makes it difficult to maintain the minimal



Figure 7. Comparison of the estimated AB and the measured AB for various values of fixed delay  $(D_f + D_b)$ 

backlogging condition for the local server. Especially, as considered in the third stage of the AB estimation procedure, if the last probing packet sent arrives before the next probing packet is sent, the local server remains in the probing-packet-free state for at least  $(D_f + D_b)$ . In other words, the next probing packet arrives late at the local server  $(D_f + D_b)$ , compared with the case that the probing packets are sent ideally according to the minimal-backlogging method. Thus, the amount of probing packets sent to the local server in a given time is always less than that of the ideal case due to  $(D_f + D_b)$ . Thus, the estimation result of the proposed method is conservative if  $(D_f + D_b)$  is significantly large. However, if the value of  $(D_f + D_b)$  is not much larger than the queueing delay at a local server or router, the proposed estimation method can work reliably. Similar tendency is observed for Poisson traffic loads.

# VII. CONCLUSIONS

A new estimation method of the available bandwidth for an unidentified queueing system is proposed using a minimal backlogging concept. Two statistics are also proposed to estimate the available bandwidth: the first one is based on the delay of each probing packet and the second one is based on the amount of probing packets served during a specific time period. If the probing packets are sent to the queueing system according to the minimal-backlogging method, the available bandwidth of the system can be estimated by either of two statistics. If the load of input traffic for an M/G/1 queueing system is less than 1, the queueing system is still stable when the minimal-backlogging method is used. The first statistic is a consistent estimator of the reciprocal of the available bandwidth and the mean square error converges to zero. The second statistic is a consistent estimator of the available bandwidth with a mean square error converging to zero. Both statistics can be used to estimate the available bandwidth of a G/G/1queueing system. Though the two statistics are unbiased estimators of the available bandwidth or its reciprocal in case of an infinite probing time duration, since infinite probing time can not be realized, we evaluated the performance of two statistics by simulation and observed that two statistics agree well with the measured available bandwidth even for a finite probing time.

We also proposed a scheme to estimate the available bandwidth of a local server that is separated from the probing source by a fixed delay by exploiting the theory for a single-server queue. The proposed scheme yields an accurate estimation result for various traffic loads when the fixed delay is relatively small compared with the queueing delay at the local server.

## APPENDIX

# A. Proof of Theorem 3

L(n) can be expressed as

$$L(n) = \frac{\sum_{i=1}^{n} \sum_{j=1}^{X'_{i}} X'_{i,j}}{\sum_{i=1}^{n} X'_{i}}$$
  
=  $\frac{\frac{1}{2} \sum_{i=1}^{n} (X'_{i})^{2} + \sum_{i=1}^{n} \sum_{j=1}^{X'_{i}} (X'_{i} - j + 1) N'_{i,j}}{\sum_{i=1}^{n} X'_{i}} - \frac{1}{2}$ 

Since  $\{X_i, i \ge 0\}$  is an ergodic Markov chain,  $\{X'_i, i \ge 0\}$  is also an ergodic Markov chain. Since the possible values of  $X'_i$ 's are non-negative,  $\{(X'_i)^2, i \ge 0\}$  is also an ergodic Markov chain. Thus, we have

$$\frac{\sum_{i=1}^{n} X'_{i}}{n} \to E[X_{\infty}], \text{ with prob. 1 and } \frac{\sum_{i=1}^{n} (X'_{i})^{2}}{n} \to E[X_{\infty}^{2}] \text{ with prob. 1}$$

Now, we consider the problem of evaluating  $\sum_{i=1}^{n} \sum_{j=1}^{X'_{i}} (X'_{i} - j + 1)N'_{i,j}/n$ . Since  $N'_{i,j}$ 's are i.i.d,  $\sum_{j=1}^{X'_{i}} (X'_{i} - j + 1)N'_{i,j}$  and  $\sum_{j=1}^{X'_{j}} jN'_{i,j}$  have the same distribution. Thus,

$$\lim_{n \to \infty} \frac{\sum_{i=1}^{n} \sum_{j=0}^{X'_{i}} (X'_{i} - j + 1) N'_{i,j}}{n} = \lim_{n \to \infty} \frac{\sum_{i=1}^{n} \sum_{j=0}^{X'_{i}} j N'_{i,j}}{n}.$$

By using the indicator function  $I(X'_i \ge j)$ , we have

$$\frac{\sum_{i=1}^{n} \sum_{j=1}^{X'_{i}} j N'_{i,j}}{n} = \frac{\sum_{j=1}^{\infty} j \sum_{i=1}^{n} I(X'_{i} \ge j) N'_{i,j}}{n}.$$
 (A.1)

To apply the renewal rates theorem for Markov-chain [14, p.164], we assign the reward  $R_j(k) = I(k \leq j)N'_{i,j}$  to the visit to the state k of  $X'_i$ . Note that  $N'_{i,j}$ 's are i.i.d and  $E[N'_{i,j}] = \lambda E[S]$ . Then, we have for each  $j \geq 0$ ,

$$\lim_{n \to \infty} \frac{\sum_{i=1}^{n} I(X'_i \ge j) N'_{i,j}}{n} = \sum_{k=0}^{\infty} \pi_k E[R_j(k)]$$
$$= \lambda E[S] \sum_{k=j}^{\infty} \pi_k$$
(A.2)

under the condition that for any initial state s,

$$E\left[\sum_{i=1}^{T_s} I(X'_i \ge j) N'_{i,j}\right] < \infty, \tag{A.3}$$

where  $T_s$  is the first return time to s of  $\{X'_i, i \ge 1\}$  starting from the initial state s. Clearly,  $E[\sum_{i=1}^{T_s} I(X'_i \ge j)N'_{i,j}] \le E[\sum_{i=1}^{T_s} N'_{i,j}]$ . Since  $T_s$  is the stopping time for the random sequence  $\{N'_{1,j}, N'_{2,j}, \cdots\}$  and  $N'_{i,j}$  are i.i.d, we obtain from the Wald's identity [14, p.97]

$$E[\sum_{i=1}^{T_s} N'_{i,j}] = E[T_s] \cdot \lambda E[S]$$

Since  $\{X'_i, i \ge 1\}$  is positive recurrent,  $E[T_s] < \infty$ . This implies Eqn. (A.3). From Eqns. (A.1) and (A.2), we have

$$\lim_{n \to \infty} \frac{\sum_{i=1}^{n} \sum_{j=1}^{X'_{i}} j N'_{i,j}}{n} = \sum_{j=0}^{\infty} j \cdot \lambda E[S] \sum_{k=j}^{\infty} \pi_{k}$$
$$= \lambda E[S] \sum_{k=0}^{\infty} \pi_{k} \sum_{j=0}^{k} j$$
$$= \lambda E[S] \frac{E[X_{\infty}^{2}] + E[X_{\infty}]}{2}.$$

Thus, we can obtain

$$\lim_{n \to \infty} L(n) = \frac{E[X_{\infty}^2]/2 + \lambda E[S](E[X_{\infty}^2] + E[X_{\infty}])/2}{E[X_{\infty}]} - \frac{1}{2}$$

$$= \frac{1}{2}(1+\rho)\frac{E[X_{\infty}^2]}{E[X_{\infty}]} - \frac{1}{2}(1-\rho).$$
(A.4)

Eqn. (5) implies  $E[X_{\infty}] = \lambda E[S_p]/(1 - \rho)$ . From Eqn. (2), we can obtain the following relation between  $\Pi_i(z)$  and  $\Pi_{i+1}(z)$ :

$$\Pi_{i+1}(z) = \tilde{G}_p(\lambda(1-z))\Pi_i(\tilde{G}(\lambda(1-z))).$$

By differentiating twice the above equation with respect to z and substituting z = 1 into the equation, we can obtain a recurrence relation between  $E[X_i^2]$  and  $E[X_{i+1}^2]$ . From the recurrence relation, we have

$$E[X_{\infty}^2] = \frac{\lambda^2 E[S_p^2]}{1 - \rho^2} + \frac{(1 - \rho^2 + 2\lambda^2 E[S_p] E[S] + \lambda^2 E[S^2])\lambda E[S_p]}{(1 - \rho)(1 - \rho^2)}.$$
 (A.5)

Combining Eqns. (A.4) and (A.5) yields

$$\lim_{n \to \infty} L(n) = \rho + \frac{\lambda^2 E[S^2]}{2(1-\rho)} + \frac{\lambda E[S_p^2]}{2E[S_p]} + \frac{2\lambda^2 E[S_p]E[S]}{2(1-\rho)}.$$

#### REFERENCES

- [1] L. Cherkasova and P. Phaal, "Session-based admission control: a mechanism for peak load management of commercial web sites," *IEEE Transactions on Computers*, vol. 51, no. 6, pp. 669-685, June 2002.
- [2] V. Sharma and R. Mazumdar, "Estimating traffic parameters in queueing systems with local information," *Performance evaluation*, vol. 32, no. 3, pp. 217-230, Apr. 1998.
- [3] S. Alouf, P. Nain, and D. Towsley, "Inferring network characteristics via moment-based estimators," *Proc. of IEEE Infocom* '01, pp. 1045-1054, Apr. 2001.
- [4] R. L. Carter and M. E. Crovella, "Measuring bottleneck link speed in packet-switched networks," *Performance Evaluation*, vol. 27-28, pp. 297-318, 1996.
- [5] C. Dovrolis, P. Ramanathan, and D. Moore, "What do packet dispersion techniques measure?," Proc. of IEEE Infocom '01, pp. 905-914, Apr. 2001.
- [6] M. Jain and C. Dovrolis, "End-to-end available bandwidth: measurement methodology, dynamics, and relation with TCP throughput," *Proc. of ACM SIGCOMM '02*, pp. 295-308, Aug. 2002.
- [7] B. Melander, M. Bjorkman, and P. Gunningberg, "A new end-to-end probing and analysis method for estimating bandwidth bottlenecks," *Proc. of IEEE Globecom* '00, pp. 415-421, Nov. 2000.
- [8] V. Ribeiro, M. Coates, R. Riedi, S. Sarvotham, B. Hendricks, and R. Baraniuk, "Multifractal cross-traffic estimation," *Proc. of ITC Specialist Seminar on IP Traffic Measurement, Modeling, and Management*, Sep. 2000.
- [9] V. Ribeiro *et al*, "pathChirp: efficient available bandwidth estimation for network paths," *Proc. Passive and Active Measurements Workshop*, Apr. 2003.
- [10] N. Hu and P. Steenkiste, "Evaluation and characterization of available bandwidth probing techniques," *IEEE J. Select. Areas Commun.*, vol. 21, no. 6, pp. 879-894, Aug. 2003.
- [11] J. Strauss, D. Katabi, and F. Kaashoek, "A measurement study of available bandwidth estimation tools," *Proc. The Internet Measurements Conference*, Florida, Oct. 2003.
- [12] C. Cetinkaya, V. Kanodia, and E. W. Knightly, "Scalable services via egress admission control," *IEEE Transactions on Multimedia*, vol. 3, no. 1, pp. 69-81, Mar. 2001.
- [13] J. Qiu and E. W. Knightly, "Inter-class resource sharing using statistical service envelopes," Proc. of IEEE INFOCOM '99, New York, pp. 1404-1411, Mar. 1999.
- [14] R. W. Wolff, Stochastic Modeling and the Theory of Queues, Prentice Hall, 1988.
- [15] A. G. Pakes, "Some conditions for ergodicity and recurrence of Markov chains," Oper. Res., vol. 17, pp. 1058-1061, 1969.
- [16] H. Takagi, Queueing analysis : a foundation of performance evaluation, vol. 1, Elsevier, North-holland, 1991.
- [17] S. Y. Nam, "Available bandwidth estimation and measurement-based admission control in IP networks," *Ph.D. thesis*, Dept. EECS, KAIST, Korea, 2004.
- [18] R. V. Hogg and A. T. Craig, Introduction to Mathematical Statistics, 5th ed., Prentice Hall, 1995.
- [19] M. E. Crovella and A. Bestavros, "Self-similarity in world wide web traffic: evidence and possible causes," *IEEE/ACM Trans. Networking*, vol. 5, pp. 835-846, 1997.
- [20] W. E. Leland, M. S. Taqqu, W. Willinger, and D. V. Wilson, "On the self-similar nature of Ethernet traffic," *IEEE/ACM Trans. Networking*, vol. 2, pp. 1-15, 1994.
- [21] V. Paxson and S. Floyd, "Wide-area traffic: the failure of Poisson modeling," *IEEE/ACM Trans. Networking*, vol. 3, pp. 226-244, 1995.
- [22] R. H. Riedi, M. S. Course, V. J. Ribeiro, and R. G. Baranuik, "A multifractal wavelet model with application to network traffic," *IEEE Trans. Inform. Theory*, vol. 45, pp. 992-1018, 1999.
- [23] K. L. Chung, A Course in Probability Theory, 2nd ed., Academic Press, 1974.