

Probing strategies for distributed admission control in large and small scale systems

Peter Key

Microsoft Research Limited
7 JJ Thomson Avenue, Cambridge, UK
Email: peterkey@microsoft.com

Laurent Massoulié

Microsoft Research Limited
7 JJ Thomson Avenue, Cambridge, UK
Email: lmassoul@microsoft.com

Abstract—The aim of this article is to propose and analyse measurement-based admission control schemes. We distinguish between large-scale and small-scale systems, where scale is measured in the number of concurrent applications that can run simultaneously. For large scale systems, we show that simple end-user probing strategies, based on ECN-type feedback provided by the network, achieve a good utilisation/quality trade-off. We explicitly take account of feedback delay, and use limiting results for assessing performance. We illustrate the benefits of using ECN-type feedback rather than relying on loss. For small-scale systems, the previous strategies are no longer adequate and we propose alternative, more gradual probing strategies.

Methods Keywords: Stochastic Processes/Queuing Theory, Control theory, System Design

I. INTRODUCTION

Real-time flows have quality of service requirements that translate to requiring some minimal level of resource allocation, for example some minimal bandwidth. For packet-based networks such as the Internet, the minimal quality of service level for many real-time flows also typically involves requiring small packet-loss rates and small queuing delays to minimise latency. Although adaptive codes and Forward Error Correction can be used to help counteract variations in bandwidth and quality, they do not solve issues of latency or varying bandwidth for all applications. We are motivated by audio and video applications that have an interactive element, or tight latency bound, where the minimal requirement is necessary.

Admission control is the obvious way to ensure such flows get an acceptable level of performance from the network, assuring the quality of service of existing flows by refusing admission to others. Traditionally some form of signalling mechanism is used, for example RSVP, which makes reservations along a path. However there are scalability questions associated with this approach, which has led several researchers [1], [2], [3], [4], [5] to consider *end-point admission control*, or *distributed admission control*, where the end-system probes the network at some rate, receives some information back from the network in terms of packet-loss or ECN marks [6], and bases the entry decision on this information. The end-systems may be hosts, or gateways. The advantage of this approach is that no state has to be maintained in routers, all they have to do is drop or mark packets appropriately.

Breslau et al [2] discuss some of the architectural issues with this approach. They major on loss as the primary feedback mechanism, with the admission controlled traffic in a separate DiffServ service class, giving soft-guarantees that fit into the Controlled Load framework. Kelly [5] looks at a looser framework, an Integrated Services network, with ECN marks as the feedback mechanism, and hence router support for ECN is the main network requirement. We do not limit ourselves to any particular architecture, although our model is similar to Kelly's, and shall not go further into the implementation issues. Our model also includes loss based feedback as a special case, which requires no ECN support.

Many voice and audio applications use UDP to transfer data and for the associated control channel, hence to use ECN marks in the ways that we want requires that this IP level information is accessible to UDP. The feedback mechanism could be done using an extension to RTCP for example (the RTP control channel) as the feedback channel.

Our primary focus is on how the probing should be done. For simplicity, we only consider admission controlled traffic. For example, such traffic could be in its own traffic class. Much of the previous literature has focussed on simulations as a way of investigating probing schemes. Most researchers assume that since packet loss is a critical performance metric, the probing phase should attempt to estimate this loss, [2], [5], which implies a long probing phase, whereas we consider much simpler probing schemes which make a relatively quick decision based on feedback from a small number (perhaps 1) of probe packets; the cumulative decisions of many connections then keep the loss rate down. This is similar to [1], although there the effects of feedback delay were neglected. Previous researchers have suggested probing at the rate which you want to achieve, or to use something like slow-start ramp-up to achieve this probing rate. We attempt to determine an optimal probing rate. Indirectly, we are also able to illustrate the benefits of using ECN-type feedback rather than packet loss; essentially when loss is the feedback signal we are unable to cope with a very high connection arrival rate without damaging the system, which is not the case when 'early-warning' signals (such as ECN) are used.

There is a natural system scale measured by the number of applications that can run simultaneously. A small-scale system may represent a home network, where there is a small number

of simultaneous real-time audio/video applications that may compete for some resource. Conversely, an example of a large scale system is the Internet, with much larger capacities at bottleneck links. We attempt to look at both scales.

The outline of the paper is as follows: in Section II, we look at large scale systems. First we outline the model of a single resource, and then show how it is possible to use fluid limits and diffusion scaling to analyse fluctuations about this limit. The closest related work is that of [7], who look at adaptation rather than admission control. We use the diffusion approximation to the system with feedback delays to analyse the performance of different probing schemes, in terms of performance seen by connections and in terms of stability. In Section III we discuss some consequences of the analysis, as well as comparing the results with simulations. It turns out that there are several key parameters that affect the performance of the system, such as the derivative of the marking function and whether we mark before we lose packets, which we discuss. Stability of the system is affected by the relative size of the round-trip time compared to the flow life-time. We are able to show that choosing the key parameters appropriately enables us to design robust probing strategies for stable systems. In Section IV we turn to small scale systems, where the analysis and probing strategies for large systems do not apply. Finally, in Section V we conclude.

II. LARGE SCALE SYSTEMS

A. Model

Consider a single resource, with capacity N . Request arrivals occur at the instants of a Poisson process with intensity $\nu^{(N)} := N\nu$. Requests last for an exponentially distributed random time, with parameter μ .

New connections initially enter a *probing* state. When probing, connections send packets at some rate r_p , and these packets are ACKed. Upon receipt of an unmarked ACK, a probing connection moves to the *active* state with probability Π_a , and continues probing with probability $1 - \Pi_a$. Upon receipt of a marked ACK (or inference of loss for the corresponding packet), a connection drops out with probability Π_d , and keeps on probing with probability $1 - \Pi_d$. While active, connections do not react to marks, and send packets at constant rate r_a . ACKs are received after some fixed round trip time delay τ . As a special case, packets may be dropped (lost) rather than marked, which conceptually we treat *as if* a marked ACK is received. For ease of description, we continue to use the term ‘marked packets’ even when these correspond to dropped packets.

For tractability we assume that when active or probing, connections send packets at Poisson time instants, with respective rates r_a and r_p . Let $X_a^{(N)}(t)$ and $X_p^{(N)}(t)$, respectively, denote the number of active probing connections at time t . Thus packets are generated at a Poisson rate of $r_a X_a^{(N)}(t) + r_p X_p^{(N)}(t)$ at time t . Assume that the resource marks (alternatively, drops) packets independently of one another, with probability $f(y/N)$ when faced with a Poisson process of packet arrivals, with

intensity y (see [8] for a discussion of marking schemes that may achieve this). Notice that the function f is implicitly defined by the marking or packet dropping behaviour at a resource, for example drop-tail marking, or other marking schemes induce a particular f .

If we ignore the round-trip delay τ , then under these assumptions the state variable $X^{(N)}(t) := (X_p^{(N)}(t), X_a^{(N)}(t))$ is a Markov process, with transition rates

$$\begin{cases} (x_p, x_a) \rightarrow (x_p, x_a - 1) : & \mu x_a \\ (x_p, x_a) \rightarrow (x_p + 1, x_a) : & N\nu \\ (x_p, x_a) \rightarrow (x_p - 1, x_a + 1) : & q_a [1 - f(y/N)] x_p \\ (x_p, x_a) \rightarrow (x_p - 1, x_a) : & q_d f(y/N) x_p, \end{cases} \quad (1)$$

where $y := r_a x_a + r_p x_p$, $q_d = r_p \Pi_d$ and $q_a = r_p \Pi_a$, in other words q_a, q_d are the rates at which probing connections convert to active connections or drop out respectively. A similar model has been analysed in [1], where the authors also consider the case of multiple bottlenecks, but assume zero round-trip delays. For non-zero round-trip time τ , the same transition rates apply, but at time t one should use $y(t - \tau) = r_a x_a(t - \tau) + r_p x_p(t - \tau)$ rather than $y(t)$ in the argument of function f , and $x_p(t - \tau)$ rather than $x_p(t)$. Indeed, probing connections liable to change their rates at time t will do so because they were already probing at time $t - \tau$, and do so on the basis of feedback information delayed by τ .

Another system we shall consider is that where candidate connections are provided with a feedback signal, equal to 0 with probability $1 - f(y/N)$, and chose to enter when it equals 0 and to leave otherwise. In that system no probing traffic is created and the state of the system is simply the number of active connections, $X_a(t)$. The transitions for this system are

$$\begin{cases} x_a \rightarrow x_a - 1 : & \mu x_a \\ x_a \rightarrow x_a + 1 : & N\nu [1 - f(y/N)]. \end{cases} \quad (2)$$

Again, in the presence of round-trip delays, one should take $y = y(t - \tau) = r_a x_a(t - \tau)$ into the argument of f . In the sequel we refer to this system as operating with ‘free probing’. This model is appropriate for reflecting the behaviour of candidate traffic sources that send a single packet through the network to receive feedback, in the case where the additional traffic due to such probe packets can be neglected.

B. Fluid limits

We now let N go to infinity. Using results of Hunt and Kurtz [9] or [10], it can be shown that under suitable assumptions on the initial conditions $X^{(N)}(0)$, the rescaled process $\{N^{-1}X^{(N)}(t)\}$ converges to the process $\{n(t)\}$, where $n(t) = (n_p(t), n_a(t))^T$ satisfies the delay-differential equations

$$\begin{cases} \dot{n}_a(t) = q_a(1 - f_{t-\tau})n_p(t - \tau) - \mu n_a(t), \\ \dot{n}_p(t) = \nu - q_a(1 - f_{t-\tau})n_p(t - \tau) - q_d f_{t-\tau} n_p(t - \tau). \end{cases}$$

In the above, we have introduced the notation $f_s = f(r_p n_p(s) + r_a n_a(s))$. Setting these derivatives to zero one obtains the

following fixed point equations for the eventual limiting points (\bar{n}_p, \bar{n}_a) of these dynamics:

$$\begin{cases} \mu \bar{n}_a &= q_a [1 - f(r_p \bar{n}_p + r_a \bar{n}_a)] \bar{n}_p, \\ \nu &= q_a [1 - f(r_p \bar{n}_p + r_a \bar{n}_a)] \bar{n}_p + q_d f(r_p \bar{n}_p + r_a \bar{n}_a) \bar{n}_p. \end{cases} \quad (3)$$

For concreteness, assume now that

$$f(x) = \min(1, (1 + f' \cdot (x - c))^+), \quad (4)$$

so that the sensitivity of f is f' in the range where it is not constant, and f reaches its maximum value 1 at c . c is related to the capacity of the resource. Note that if $c = 1$, then we quench arrivals when the capacity limit is reached, whereas if $c < 1$, we start to stop accepting connections before this point.

Assume for the moment that $q_a = q_d = q$, which is saying that the probability that an unmarked packet causes a probing connection to become active is the same as the probability that a marked packet causes a connection to drop out. We let $\rho := r_a \nu / \mu$ denote the (normalised) potential load on the system. Assuming $\rho > c$, there exists a unique fixed point to the above equations, given by

$$\begin{cases} \bar{n}_p &= \frac{\nu}{q}, \\ r_a \bar{n}_a &= \frac{\rho f'}{1 + \rho f'} (c - r_p \bar{n}_p). \end{cases} \quad (5)$$

This can be interpreted as follows: the steeper the marking function (i.e., the larger f'), the larger the useful utilisation $r_a \bar{n}_a$. Also, the smaller the per-connection probe traffic r_p/q , the larger the useful utilisation. We now discuss how these conclusions should be modified when taking into account the dynamics of the system.

We assume existence and uniqueness of the solution \bar{n} . Let $\bar{f} := f(r_p \bar{n}_p + r_a \bar{n}_a)$, and $\bar{f}' := f'(r_p \bar{n}_p + r_a \bar{n}_a)$. The linearised dynamics for small perturbations $m(t) = (m_p(t), m_a(t))^T := n(t) - \bar{n}$ are

$$\dot{m}(t) = -Pm(t) - Qm(t - \tau), \quad (6)$$

where the matrices P and Q are given by

$$\begin{aligned} P &= \begin{pmatrix} 0 & 0 \\ 0 & \mu \end{pmatrix} \\ Q &= \begin{pmatrix} q_a(1 - \bar{f}) + q_d \bar{f} & 0 \\ -q_a(1 - \bar{f}) & 0 \end{pmatrix} + \\ &\quad \bar{n}_p \bar{f}' \begin{pmatrix} (q_d - q_a)r_p & (q_d - q_a)r_a \\ q_a r_p & q_a r_a \end{pmatrix}. \end{aligned} \quad (7)$$

The main stability properties for the linearised system (6) are summarised below; their proof is in the Appendix.

Theorem 1: In the absence of propagation delays (i.e., when $\tau = 0$), the system (6) is stable when $q_d \leq q_a$. For positive τ , stability holds if and only if any solution s to the characteristic equation

$$\text{Det}(sI + P + e^{-s\tau}Q) = 0 \quad (8)$$

is such that its real part $\Re(s)$ is negative. In the special case where $q_a = q_d =: q$, sufficient conditions for stability are given

by

$$q\tau < \frac{\pi}{2}, \quad (9)$$

$$(\nu \bar{f}' r_a \tau)^2 < (\mu\tau)^2 + \alpha^2, \quad (10)$$

where α is the solution to the equation $\alpha = -\mu\tau \tan(\alpha)$ that lies in the interval $(0, \pi)$.

Remark 1: Recently, several authors ([11], [8], [12]) have analysed the stability properties of congestion control schemes with an emphasis on the impact of feedback delays. These contrast with the present work not only because we focus on admission control rather than on congestion control, but also because we are mainly concerned with flow level dynamics, whereas the references above consider rate adaptation for a fixed population of users. On the other hand, these go beyond the single bottleneck case.

Note that q is the rate at which probing connections will either leave the system or move to the active state. Thus, q is also interpreted as the reciprocal of the probing time. Condition (9) requires a probing time that is at least proportional to the round-trip delay¹. As the quantity α in (10) is always larger than $\pi/2$, a sufficient condition for (10) to hold is given by

$$(\nu \bar{f}' r_a \tau)^2 < (\mu\tau)^2 + \left(\frac{\pi}{2}\right)^2. \quad (11)$$

Taking for instance the piecewise affine marking function (4), one sees that the system's stability is not affected by the probing volume r_p/q used by each connection. Thus the system's performance is increased by reducing this probe traffic r_p/q , as this allows one to accept more flows into the system.

C. Diffusion scaling: fluctuations around the fluid limits

In this section we investigate the ‘‘second-order’’ properties of the system under consideration. Recall that $X^{(N)}(t)$ is the state variable, expected to be close to $N\bar{n}$. Introduce the notation

$$Z(t) := (Z_p(t), Z_a(t))^T = \frac{1}{\sqrt{N}} (X(t) - N\bar{n}).$$

We argue that this (approximately) satisfies the following stochastic delay differential equation

$$dZ(t) = -PZ(t)dt - QZ(t - \tau)dt + dW(t), \quad (12)$$

where

$$\begin{aligned} W(t) &= \sqrt{\nu} \begin{pmatrix} 1 \\ 0 \end{pmatrix} B_1(t) + \sqrt{\mu \bar{n}_a} \begin{pmatrix} 0 \\ 1 \end{pmatrix} B_2(t) \\ &+ \sqrt{q_a \bar{n}_p (1 - \bar{f})} \begin{pmatrix} -1 \\ 1 \end{pmatrix} B_3(t) + \sqrt{q_d \bar{n}_p \bar{f}} \begin{pmatrix} -1 \\ 0 \end{pmatrix} B_4(t), \end{aligned} \quad (13)$$

and B_i , $i = 1, \dots, 4$ are independent scalar Brownian motions. This is a two-dimensional Ornstein-Uhlenbeck process, with

¹In the more general case where $q_a \neq q_d$, we also find that the two parameters q_a and q_d have to be chosen proportional to τ^{-1} . The corresponding stability conditions are less easily expressed and we do not provide them here.

a delayed input term. Indeed, the perturbations $Z(t)$ (approximately) satisfy the linearised equations of the previous subsection, with additional input noise terms which correspond to the randomness in the original Markov process². The result below characterises the stationary covariance of process Z .

Theorem 2: The process Z as described by (12) admits the steady state covariance matrix

$$\Sigma := \mathbf{E}[ZZ^T] = \frac{1}{2\pi} \int_{-\infty}^{+\infty} \left\{ [i\omega I + P + e^{-i\omega\tau} Q]^{-1} M \right. \\ \left. [-i\omega I + P^T + e^{i\omega\tau} Q^T]^{-1} \right\} d\omega, \quad (14)$$

where

$$M = \begin{pmatrix} v + q_a \bar{n}_p \bar{f} + q_a \bar{n}_p (1 - \bar{f}) & -q_a \bar{n}_p (1 - \bar{f}) \\ -q_a \bar{n}_p (1 - \bar{f}) & \mu \bar{n}_a + q_a \bar{n}_p (1 - \bar{f}) \end{pmatrix}.$$

The proof is in the Appendix. We now discuss two cases of particular interest.

a) variance analysis with zero delays: In the undelayed case $\tau = 0$, the integral in (14) admits an explicit solution. In the special case where $q_a = q_d = q$, the corresponding closed form expression for Σ in turn yields the following formula for the variance of the normalised total utilisation $r_a Z_a + r_p Z_p$:

$$\text{Var}(r_a Z_a + r_p Z_p) = \frac{r_a^2 v (1 - \bar{f})}{\mu + \bar{f} r_a v} \\ + \left(\frac{r_p}{q} \right) v \frac{r_p \mu^2 + q \left(-(1 - \bar{f}) \bar{f}' r_a^2 v + r_p (\mu + v r_a \bar{f}') \right)}{(\mu + \bar{f} r_a v) (q + \mu + \bar{f} r_a v)}. \quad (15)$$

Note that the first term does not depend directly on the two parameters r_p, q that characterise the probing phase (it depends indirectly on the ratio r_p/q , which affects the value of \bar{f}). The second term will be zero if both r_p and q are set to zero, while the ratio r_p/q is held fixed. This limiting situation corresponds to probing at an infinitesimal rate, which still sees the true marking probability, and reacting accordingly. The first term in (15) does in fact correspond to the variance of the quantity $r_a Z_a$ where Z_a is the (approximation to the) perturbation term $(X_a(t) - N \bar{n}_a) / \sqrt{N}$ for the system with free probing discussed previously.

But note that the second term in (15) may be negative, hence it is possible to *decrease* the variance by having a non-zero probing rate. Recall that $q = \Pi r_p$. To ease exposition, and without loss of generality, we now set $\mu = 1$, corresponding to a time rescaling, which requires that we replace Π by Π/μ . Define γ by

$$\gamma = v r_a \bar{f}'. \quad (16)$$

Then the variance can be written as

$$\frac{v r_a^2 (1 - \bar{f})}{1 + \gamma} + \frac{v r_p (1 - (1 - \bar{f}) \Pi \gamma r_a + \Pi (1 + \gamma) r_p)}{\Pi (1 + \gamma) (1 + \gamma + \Pi r_p)}. \quad (17)$$

²A rigorous treatment would consist in showing that as N goes to infinity, the corresponding perturbation process Z converges weakly to the Ornstein-Uhlenbeck process given here. Such a proof might for instance rely on the methods in Ethier and Kurtz [13], but is beyond the scope of the present paper.

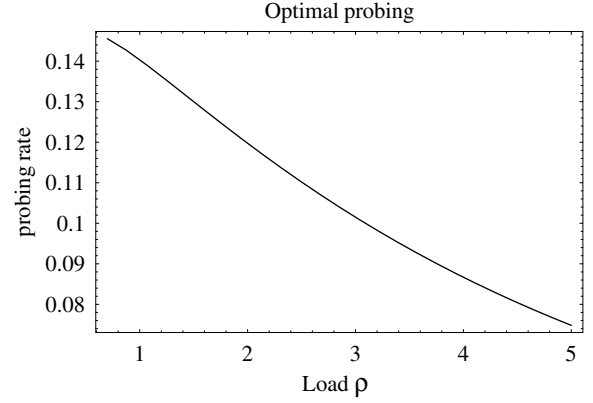


Fig. 1. Optimal probing rate r_p as a function of offered load. Parameter values: $\mu = 0.01$, $c = 1, N = 100, r_a = 1, f' = 10$.

The second term will be negative if γ is large enough, that is if

$$\gamma > \frac{1 + \Pi r_p}{(1 - \bar{f}) \Pi r_a - \Pi r_p}. \quad (18)$$

Now the fixed point equation can be written as

$$\bar{n}_a = v \left(1 - f \left(r_a \bar{n}_a + \frac{v}{\Pi} \right) \right) \quad (19)$$

hence the utilisation increases with Π , while the variance is relatively insensitive to Π and decreasing in Π , hence we want to chose Π as large as possible (1 in the original units, or $1/\mu$ in the rescaled units here). Provided the condition (18) is satisfied, then for Π fixed, the variance is minimised by choosing the probing rate to satisfy,

$$r_p = \frac{-1 - \gamma + \sqrt{\gamma(2 + (1 - \bar{f}) \Pi r_a + \gamma)}}{\Pi}. \quad (20)$$

As γ becomes large, this tends to the limit

$$r_p = \frac{1}{2} r_a (1 - \bar{f}). \quad (21)$$

The limiting value (21) has an interesting interpretation, namely that the optimal probing rate is to probe at half the rate at which active connections are marked. Note that it is always best to probe at less than half the active rate.

Figure 1 shows how the optimal rate decreases with increasing load ρ on the system, where we have taken a system of capacity $N = 100$ with $\mu = 0.01$.

b) variance analysis with free probing and non-zero delays: We now provide an exact variance analysis for the non-zero delay case, for 'free' probing. The diffusion equation satisfied by $Z_a(t)$ can either be derived from (12) by replacing \bar{n}_p by v/q and then letting r_p, q go to zero, or alternatively be derived directly. The resulting equation is

$$dZ_a(t) = -v \bar{f}' r_a Z_a(t - \tau) dt - \mu Z_a(t) dt \\ + \sqrt{2v(1 - \bar{f})} dB(t), \quad (22)$$

where $B(t)$ is a standard Wiener process. Rewrite equation (22) as

$$dZ_a(t) = -\gamma Z_a(t - \tau)dt - \mu Z_a(t)dt + bdB(t) \quad (23)$$

where γ is defined in (16) and $b = \sqrt{2v(1 - \bar{f})}$. We then have the following

Theorem 3: The process Z_a as described by (22) admits the steady state variance

$$\text{Var}(Z_a) = v(1 - \bar{f}) \cdot \frac{\cos\left(\frac{\tau}{2}\sqrt{\gamma^2 - \mu^2}\right) + \sqrt{\frac{\gamma + \mu}{\gamma - \mu}} \sin\left(\frac{\tau}{2}\sqrt{\gamma^2 - \mu^2}\right)}{(\gamma + \mu) \cos\left(\frac{\tau}{2}\sqrt{\gamma^2 - \mu^2}\right) - \sqrt{\gamma^2 - \mu^2} \sin\left(\frac{\tau}{2}\sqrt{\gamma^2 - \mu^2}\right)} \quad (24)$$

when $\gamma > \mu$, and

$$\text{Var}(Z_a) = v(1 - \bar{f}) \cdot \frac{\cosh\left(\frac{\tau}{2}\sqrt{\mu^2 - \gamma^2}\right) + \sqrt{\frac{\gamma + \mu}{\mu - \gamma}} \sinh\left(\frac{\tau}{2}\sqrt{\mu^2 - \gamma^2}\right)}{(\gamma + \mu) \cosh\left(\frac{\tau}{2}\sqrt{\mu^2 - \gamma^2}\right) + \sqrt{\mu^2 - \gamma^2} \sinh\left(\frac{\tau}{2}\sqrt{\mu^2 - \gamma^2}\right)} \quad (25)$$

otherwise.

The proof is in the Appendix.

c) Minimal variance for central controller: We would like to know how efficient the performance of the scheme with free probing is. To this end, we now evaluate the minimal variance achievable by a central controller, whose aim is to bring the utilisation as close as possible to a target value. The set-up is as follows: the controller chooses at each time t the rate of accepted new arrivals into the system, $A(t)$, on the basis of the history of the number $X(s)$ in the system for all $s \leq t - \tau$, and of the previous acceptance rates $A(s)$, $s < t$. The durations of accepted calls are random, exponentially distributed with parameter μ , and the controller has no information at time t about the departures that have taken place during $[t - \tau, t]$. We assume that the target utilisation, N , is large, and consider the deviations $Z(t) := (X(t) - N)/\sqrt{N}$. We formalise this as follows:

$$dZ(t) = a(t)dt - \mu Z(t)dt + \sqrt{\mu}dB(t),$$

where the controlled (perturbation to the) rate of arrivals $a(t)$ is $\mathcal{F}_{t-\tau}$ -adapted, and the filtration \mathcal{F}_t keeps track of the history of the process Z . One has the following explicit expression

$$Z(t) = \int_{-\infty}^t e^{-\mu(t-s)}[a(s)ds + \sqrt{\mu}dB(s)].$$

The conditional variance formula ensures that

$$\text{Var}(Z(t)) \geq \mathbf{E}(\text{Var}(Z(t) | \mathcal{F}_{t-\tau})).$$

Because $a(s)$ is $\mathcal{F}_{s-\tau}$ -adapted, the right-hand side of this expression reads

$$\text{Var}\left(\int_{t-\tau}^t e^{-\mu(t-s)}dB(s)\right) = \int_{t-\tau}^t \mu e^{-2\mu(t-s)}ds.$$

We thus obtain the following lower bound for the optimal variance:

$$\text{Var}(Z(t)) \geq \frac{1 - e^{-2\mu\tau}}{2}. \quad (26)$$

This lower bound can effectively be achieved by a central controller, as long as it is able to pool from a large reservoir of candidate connections (which holds when the offered load is larger than the capacity): indeed it suffices to define the control $a(t)$ so as to set to zero the non-negative term we have neglected when using the conditional variance formula.

This minimal variance (26) is typically one order of magnitude smaller than the one achieved by free probing. One thing to note is that this depends only on the quantity $T := 1/(\mu\tau)$, which is the ratio of holding time to round trip time. As we shall see in the next section, this ratio is also critical for the performance of the schemes we consider, although none achieves variances close to the optimal (26).

III. NUMERICAL RESULTS AND DISCUSSION

We now explore some of the consequences of the analysis, and compare with simulation. One striking fact that emerges from the analysis is the crucial effect the derivative of the marking function, f' , has on performance. On the one hand, increasing f' increases the efficiency of the system, both by increasing the utilisation, as is immediate from (5), and by decreasing the variance; on the other hand the f' may compromise stability, since the larger the f' the smaller the offered load before stability breaks down. We return to this point later.

First we consider the system with negligible delays. Equation (21) says that we should probe at not more than half the active rate, in fact at exactly half the rate at which an active connection would be marked. To assess the performance of different probing schemes, we look at the probability the offered load exceeds the capacity limit, N . Recall that connections are generating packets at certain rates, hence when the load (total utilisation) exceeds this limit packets are lost, hence we use the performance measure

$$P_{\text{loss}} = \Pr[r_a X_a + r_p X_p > N] = \Pr\left[r_a \bar{n}_a + r_p \bar{n}_p + \frac{1}{\sqrt{N}}(r_a Z_a + r_p Z_p) > 1\right] \quad (27)$$

which is the proportion of connections that lose packets. This relates to the proportion of connections that see an unacceptable quality of service. Using the large system limit, for a system of size N we approximate this by the probability that a normal distribution with mean $r_a \bar{n}_a + r_p \bar{n}_p$ and variance $\text{Var}(r_a Z_a + r_p Z_p)/N$ exceeds 1.

Figure 2 compares utilisation against performance for three schemes, where we put $c = 1$, in other words where the point at which we mark all packets is also the point at which we start to lose packets. This corresponds to drop-tail marking, or the case where there is no marking, only packet loss. Here the utilisation is the normalised goodput, $r_a \bar{n}_a$. The three schemes are

- Free Probing,

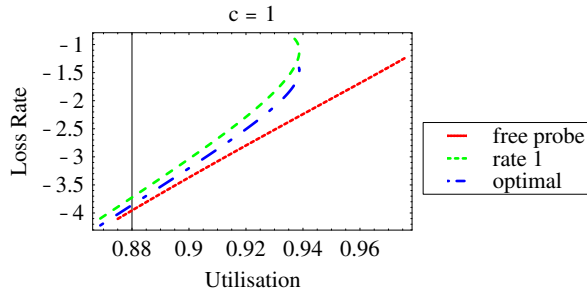


Fig. 2. Utilisation (goodput) and $\log_{10} P_{loss}$, based on analysis. Parameter values: $\mu = 0.01$, $c = 1$, $N = 100$, $r_a = 1$, $f' = 10$.

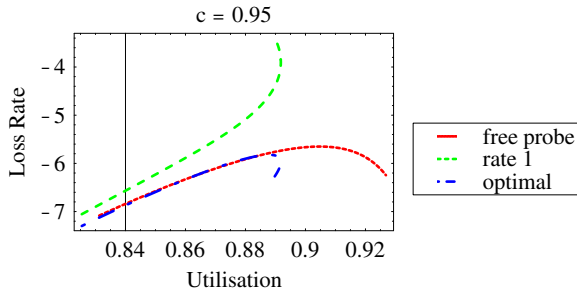


Fig. 3. Utilisation (goodput) and $\log_{10} P_{loss}$, based on analysis. Parameter values: $\mu = 0.011$, $c = 0.95$, $N = 100$, $r_a = 1$, $f' = 10$.

- ‘Rate 1 probing’: probing at active rate (i.e. $r_p = r_a$)
- ‘Optimal probing’: probing at the rate which minimises the variance of the load, given by (20).

We have taken f' to be 10, and set $\Pi/\mu = 100$, which could correspond to a mean holding time of say 100 (seconds), with $\Pi = 1$. We have also taken a moderately sized system, with $N = 100$, and the graphs are obtained using the analytic formulae of the previous section as the offered load ρ is varied between 0.7 and 4. The results show that free probing performs the best, with a log-linear relationship between goodput and performance. For optimal and rate 1 probing the performance degrades badly above a critical offered load.

Contrast this with Figure 3, where c is less than 1, namely $c = 0.95$, implying that we mark all packets before we start losing packets. Now both free probing and ‘optimal’ probing have room to work, and the performance of the system is bounded, in other words the performance of existing flows is not damaged excessively by probing traffic, whereas probing at the active rate leads to a performance collapse above a critical offered load. This illustrates dramatically the benefits of sacrificing a small amount of utilisation (say 5%) to gain control of the system. In real terms there is no loss of utilisation, since trying to run at full utilisation is likely to produce unacceptable performance for most applications.

We now introduce delay. In Figure 4 the mean and the standard deviation of the utilisation with free probing are reported, based on both simulations and analysis. We observe a good match between the predicted mean and that achieved

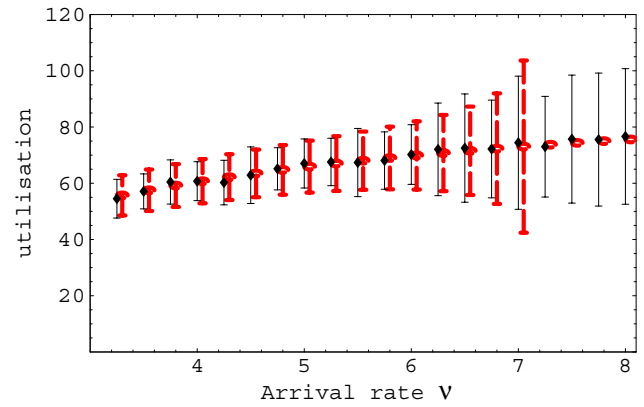


Fig. 4. Average utilisation and corresponding standard deviation, based on analysis (dotted lines) and simulations (plain lines). Parameter values: $\mu = .025$, $\tau = 30$, $c = 100$, $r_p = 0$, $r_a = 1$, $f' = 0.43$.

in simulations. The match between the predicted and observed standard deviations is also good in a wide range of values of v where the stability condition (10) is satisfied. When condition (10) fails, the Ornstein-Uhlenbeck system no longer admits a stationary regime, and thus the analysis cannot provide predictions for the standard deviation. When this condition fails, utilisation in the real system displays oscillations, the amplitude of which is essentially kept bounded by non-linearity of the marking function f .

Figures 5 and 6 illustrate how delays may compromise the behaviour of probing schemes. We use the same three schemes as before, where the optimal probing refers to the optimal rate if the delay is zero. Recall the definition $T = 1/(\mu\tau)$, hence T is the number of round trip-times per holding time. If T is moderate, 100 or smaller, implying the feedback delay is large, for example, corresponding to a 1 second round trip time and 100 second mean holding time, then optimal rate probing may not do much better than rate 1 probing. With a smaller delay, $T = 1000$ ($\mu\tau = 0.001$), the story becomes the same as without delays, i.e. almost identical to Figures 2 and 3. Hence, whether it is safe to use a non-zero probing rate depends critically on the ratio of the mean holding time to the round trip time.

Figure 7 shows a simulation of optimal probing for a system of size 100 with large delays, which tells the same kind of story as the analytic graphs. Notice that we need a large amount of headroom to preserve quality of service bounds, which here are represented by the maximum (total) utilisation. In other words here we need to have $c < 0.8$, approximately, to ensure a good performance for large loads in this case, i.e. 20% spare capacity.

For comparison, consider the experiments of Breslau et al. [2]. They consider a 10Mbps link with bursty connections having an average rate of 128k, a mean holding time of 300 seconds and a round time comprising a propagation delay of 20ms and a queuing delay bounded by 20ms, with loads between 1.2 and 4. Under our scaling, this corresponds to a capacity N of 80, and a T value of $T = 300/0.04 = 7500$. This

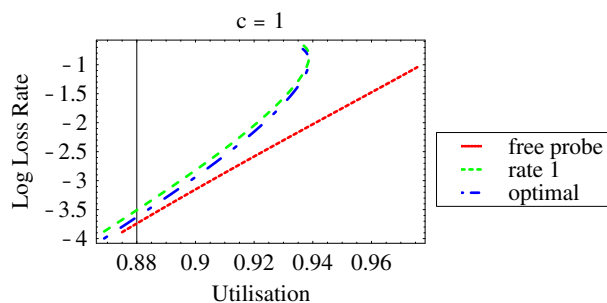


Fig. 5. Utilisation (goodput) and $\log_{10} P_{loss}$. Parameter values: $\mu = 0.01$, $c = 1$, $N = 100$, $r_a = 1$, $f' = 10$, $\tau = 1$, $T = 100$.

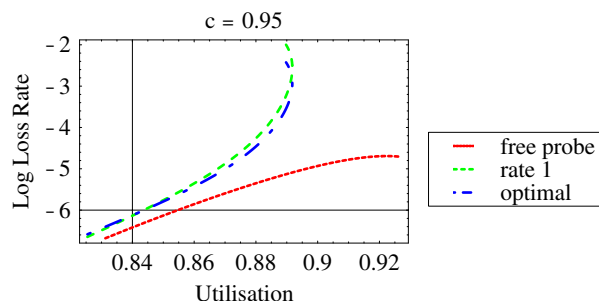


Fig. 6. Utilisation (goodput) and $\log_{10} P_{loss}$. Parameter values: $\mu = 0.01$, $c = 0.95$, $N = 100$, $r_a = 1$, $f' = 10$, $\tau = 1$, $T = 100$

means that the impact of delayed measurements is negligible, stability is not an issue, and Figure 6 can be compared with their results. In contrast to the schemes of their paper, we are probing for a very short time. Notice that we get very good performance with either optimal rate probing or ‘free-probing’. Their paper also considers a virtual queue, with $c = 0.9$ rather than $c = 0.95$, which is one reason why their utilisations are lower. They also effectively have a much higher value of f' than we do (an order of magnitude higher).

The analytical condition (11) seems to be a reliable predictor of when oscillations will emerge. Rewriting the condition as

$$\rho f' < \sqrt{1 + \left(\frac{\pi}{2\mu\tau}\right)^2} \quad (28)$$

gives the sufficient condition

$$\rho f' < \frac{\pi}{2\mu\tau} = \frac{\pi T}{2}. \quad (29)$$

There are two ways to interpret the condition (29): for a given T and f' we can regard it as a condition describing the limits on the offered load before the system destabilises, or we can seek to tune the parameter f' to ensure stability. Note that for long-lived flows such as streaming media, the system will self-stabilise under typical operating conditions. For example, assuming such flows last 100 seconds and have a 500msec

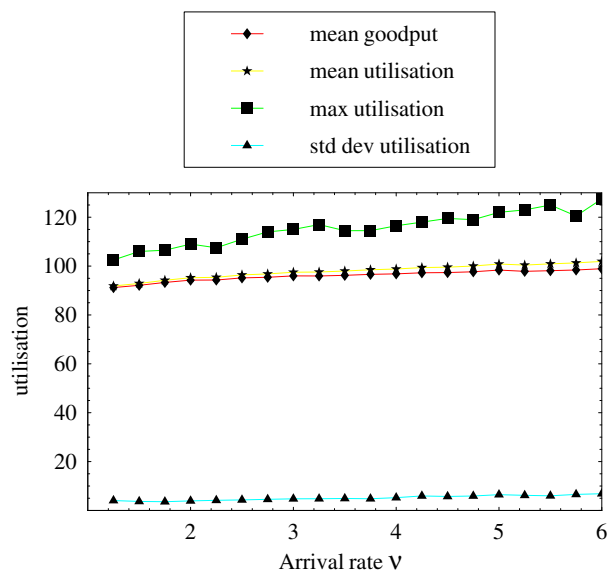


Fig. 7. Utilisation (from simulations) with optimal probing. Parameter values: $\mu = .001$, $\tau = 5$, $N = 100$, $r_p = opt$, $r_a = 1$, $f' = 10$.

round trip time implies that $T = 200$; if we set $f' = 10$, then ρ can go as large as 30 (extreme overload!) and still the system will be stable.

As a specific example of how f' can be chosen, consider the Virtual Queue marking mechanism, described in [14], [1], [15]. If the real queue (output port of a router) has a service rate N , then the virtual queue runs at lower rate cN , and marks packets if the virtual buffer capacity exceeds some threshold B . No real scheduling is done in the virtual queue, its sole purpose is to mark packets. Under our scaling, for this queue f' is approximately B , hence if we chose B to be 10 packets, (taking the average packet size) then we have the system modelled here.

Adapting f' could correspond to adapting the buffer size in the virtual queue. An alternative approach is to adapt both parameters f' and c in (4) to tune average utilisation and fluctuations. Srikant [16] considered an adaptive virtual capacity algorithm; f' is a harder quantity to tune.

IV. SMALL SCALE SYSTEMS

In small scale systems, simple strategies such as free probing might no longer provide a satisfactory goodput/quality of service trade-off. This is illustrated in Figure 8, where the mean and maximal utilisation of a system with virtual capacity cN equal to 10. We see that, when the arrival rate v is large, the maximum utilisation can be as large as twice the average utilisation. Another observation we make is that standard deviation is not an appropriate performance descriptor in such small scale systems. Indeed, the standard deviations corresponding to the simulations reported in Figure 8 are all between 1 and 2, which suggests better performance than we actually observe.

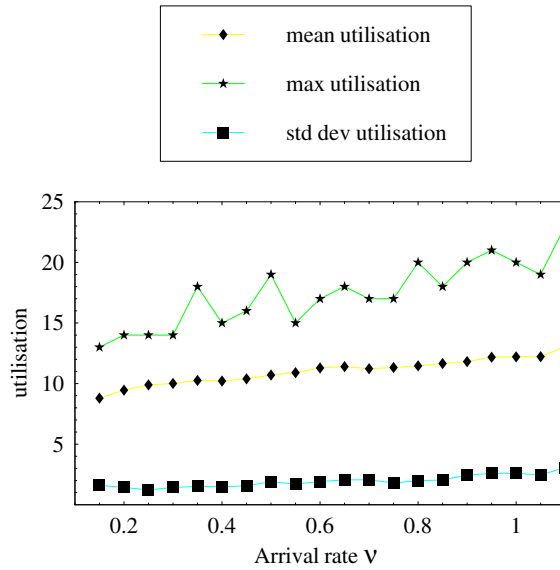


Fig. 8. Maximal and average utilisation (from simulations) with free probing. Parameter values: $\mu = .001$, $\tau = 5$, $N = 10$, $r_p = 0$, $r_a = 1$, $f' = 10$.

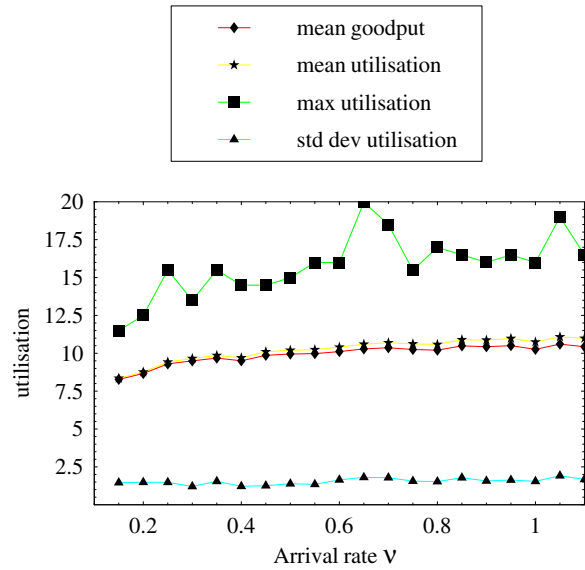


Fig. 9. Utilisation and goodput for “optimal” probing. Parameter values: $\mu = .001$, $\tau = 5$, $N = 10$, $r_p = 0$, $r_a = 1$, $f' = 10$.

Probing at non-zero rates provides a solution to reduce the variance. The question is then of designing probing strategies which allow to keep the overflow probability small, while maintaining the probing traffic small in comparison to the useful traffic carried by the resource. The rule (21) derived in the previous section provides a natural candidate. Figure 9 illustrates the corresponding performance. As we see, this is only marginally better than what free probing achieves.

Yet another option consists of using more than one probing stage. In Figure 10, users go through a first phase at which they probe at rate $1/4$, and then have to go through a second probing phase, during which they probe at rate $1/2$, before eventually entering the system.

As is seen on Figure 10, the added complexity of using several probing phases can be beneficial. More generally, one might consider the following family of probing strategies. Connections can probe at several rates, r_1, \dots, r_{K-1} , while the rate when active is still r_a . For instance, one might let

$$r_i = r_a z^{K-i},$$

where z is less than 1, and z^{-1} is the inflation ratio between probing rates $i-1$ and i . We shall also assume that a connection probing at rate r_{i-1} will either abandon, or move to the next probing level r_i (if $i = K$, the connection will in fact become active). Again, we assume that the decision between the two options is based on ECN-type binary feedback marks received from the congested resource, and that the probability of a feedback mark generated at time t being equal to 1 is exactly $f(X(t)/N)$, where N is the bottleneck capacity, and $X(t)$ is the rate submitted to the bottleneck at time t . For each feedback signal received while probing at rate r_i , a connection will choose to react to this information with probability Π_i .

Denoting by $X_i(t)$ the number of connections in probing state i at time t , we thus have

$$X(t) = \sum_{i=1}^K r_i X_i(t).$$

In the previous equation we let $X_K(t)$ denote the number of active connections at time t , and $r_K = r_a$. For an arrival rate $v^{(N)} = Nv$, exponentially distributed sojourn times with parameter μ , and in the absence of delay in the feedback loop, $(X_i(t))_{1 \leq i \leq K}$ is a Markov process with transition rates

$$\begin{cases} x_1 \rightarrow x_1 + 1 : & Nv \\ (x_i, x_{i+1}) \rightarrow (x_i - 1, x_{i+1} + 1) : & x_i r_i \Pi_i [1 - f(x)] \\ & (1 \leq i < K - 1) \\ x_i \rightarrow x_i - 1 : & x_i r_i \Pi_i f(x) \\ & (i < K) \\ x_K \rightarrow x_K - 1 : & \mu x_K \end{cases}.$$

Although the analysis techniques used in the previous section are based on the assumption that the system’s capacity N is large — which we no longer assume here — they might still provide insight into the performance of systems with small N . Mimicking the preceding derivations, we arrive at the following fixed point equations for the normalised quantities $\bar{n}_i = \mathbf{E}[X_i]/N$:

$$\begin{cases} \bar{n}_i = (\Pi_i r_i)^{-1} v (1 - \bar{f})^{i-1} & (1 \leq i \leq K - 1) \\ \bar{n}_K = (v/\mu) (1 - \bar{f})^{K-1} \\ \bar{f} = f(\sum_{i=1}^K r_i \bar{n}_i). \end{cases}$$

These might be used to infer the average resource utilisation, and the average goodput, that is the average traffic due to active rather than probing connections. We do not try to reproduce the variance analysis in the previous section, not only because it becomes so intricate with several probing levels, but also

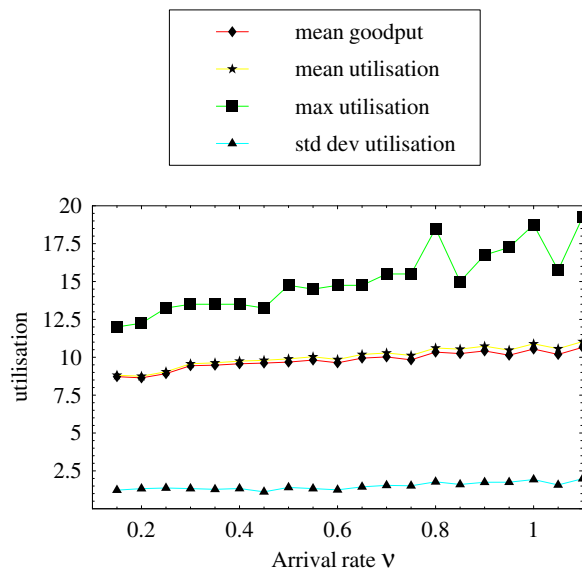


Fig. 10. Utilisation and goodput for two-level probing. Parameter values: $\mu = .001$, $\tau = 5$, $N = 10$, $r_p = 0$, $r_a = 1$, $f' = 10$. Probing rates: $1/4$ and $1/2$.

because variances might be too crude to describe the system's performance.

V. CONCLUDING REMARKS

A number of authors have proposed various probing mechanisms for end-point admission control, yet little analysis has been done to enable different schemes to be assessed. One of the aims of this work was to look at analytic techniques as a vehicle to gain insight into designing good strategies that can then be tested by simulation or experiment. A concern with any probing strategy is that the act of probing can damage existing connections, and we want to use strategies that provide good information while having minimal impact on existing connections.

We have analysed a number of probing schemes for distributed admission control, and analysed the effect of the probing traffic on the system performance. The analytic techniques take account of the delayed information available to flows attempting to join the system, which has enabled us to look at stability. We use limiting results for large systems, which give rise to fluid limits; deviations from the fluid limits are described by delayed diffusion approximations, which can be explicitly solved.

The large system results seem accurate, even for moderate size systems, and enable a number of practical conclusions to be drawn. First, the parameter $T = 1/(\mu\tau)$, the ratio of mean holding time to the round trip time, equivalently the number of round times per holding time, is a critical parameter for stability of the system. If this is large (round trip times are relatively small) then the influence of delay on the system is minimal, and performance is similar to undelayed systems; whereas if T is small (large relative RTTs), the delayed feedback can compromise the stability of the systems. We have

concentrated on strategies that are relatively light-weight, and that react rather quickly, for example making a decision on the basis of a single marked packet. Indeed, within the restricted set of probing strategies we consider, this strategy is better than one which uses several packets to make the decision. This is in contrast to much previous work which assumes that relatively long probing time is required. For the typical streaming or interactive applications of today, T is likely to be large enough that stability is guaranteed for our schemes. For emergent applications which may have short call holding times, this condition may change.

Two other critical parameters that affect the performance of the system are the derivative of the marking function, f' and the relationship between the point at which all packets are marked, c in normalised units, and the point at which packets are dropped. The larger f' , the better the utilisation, and the better the performance, but the tighter the stability criterion. If $c = 1$, corresponding to tail marking or packet loss feedback for example, then any non-zero probing scheme can become unstable under moderate loading, the effect of probing traffic pushing the system into overload. Even free probing cannot bound the performance for high load. In contrast, if we set c to be smaller than capacity, say $c = 0.95$, which could be achieved by using a virtual queue running at 95% of the real queues' rate, then the system can be stabilised. Indeed with free probing, the system is self-adjusting: the performance is bounded regardless of the load. This illustrates the benefits of giving up a small amount of potential utilisation in return for controllability and better performance, and illustrates the benefits of using something like ECN for admission control.

The 'free-probing' scheme appears the best way to probe in large systems, which receives delayed feedback and has negligible impact on the network. This could be implemented by using a very small packet to do the probing. For stable (large) systems, the limiting optimal non-zero probing rate is exactly half the rate at which active connections are marked, hence we should never probe at more than half the active rate, and if there is some knowledge of the current marking rate, then we can use this to approximate the optimal probing rate. Such a scheme can induce a smaller variance than free probing, which is offset by the extra traffic non-zero probing generates, which means that free-probing and optimal rate probing behave similarly if the system is stable.

The choice of probing strategies for small systems is less clear: the large system models are not appropriate, and indeed the variance of the load can be a poor predictor of performance. For small systems free probing is not ideal, and there is an incentive to probe at a gradually increasing rate, and we have some evidence that probing in several phases is beneficial.

The analysis has relied on simplifying assumptions for tractability, such as exponential holding times, and homogeneous connections. It is possible to relax some of the assumptions at the cost of a rather more complex analysis, though we believe that our models and conclusions for large systems are robust. If we assume marking functions that are

conditionally independent across resources given the load, then for certain probing strategies our methodology can be extended to network models.

APPENDIX

A. Proof of Theorem 1

When $\tau = 0$, the linear system (6) is stable if and only if the eigenvalues of $P + Q$ all have a positive real part. This in turn holds if and only if the trace and the determinant of $P + Q$ are positive. In view of the expressions

$$\begin{aligned} \text{Trace}(P + Q) &= q_a(1 - \bar{f}) + q_d \bar{f} + \mu + \\ &\quad \bar{n}_p \bar{f}' [(q_d - q_a)r_p + q_a r_a], \\ \text{Det}(P + Q) &= \mu [q_a(1 - \bar{f}) + q_d \bar{f} + \\ &\quad \bar{n}_p \bar{f}' (q_d - q_a)r_p] + \bar{n}_p \bar{f}' q_d r_a q_a, \end{aligned}$$

it follows that these conditions are always met provided $q_d \geq q_a$.

The system (6) is of retarded type (see Bellman and Cooke [17] for a definition and Theorem 12.12, p. 418 in [17] for a proof); Corollary 6.1, p. 190 in [17] can thus be applied, yielding the corresponding necessary and sufficient condition for stability (8).

Nyquist's criterion guarantees that stability holds for some value of τ if for all $t \in [0, \tau]$, the characteristic equation (8) has no root $i\omega$ on the imaginary axis. In the special case where $q_d = q_a =: q$, the matrices P and Q are both lower triangular, so that the eigenvalues of the matrix in (8) are simply its diagonal elements. Thus stability follows if for all $t \in [0, \tau]$, and all $\omega \in \mathbb{R}$, the two conditions

$$\begin{aligned} -i\omega &\neq e^{-i\omega t} [q_a(1 - \bar{f}) + q_d \bar{f}] = e^{-i\omega t} q \\ -i\omega &\neq \mu + e^{-i\omega t} \bar{n}_p \bar{f}' q r_a \end{aligned}$$

hold. The latter condition simplifies to

$$i\omega + \mu + e^{-i\omega t} \bar{f}' v r_a \neq 0,$$

since $\bar{n}_p = v/q$.

Hayes' lemma [18] (see also Theorem 13.8 in Bellman and Cooke [17]) can then be applied to yield the equivalent conditions (9)-(10), thus concluding the proof of Theorem 1. \square

B. Proof of Theorem 2

Using Laplace transform techniques, it can be seen that a stationary solution to (12) is given by

$$Z(t) = \int_{-\infty}^t H(t-s) dW(s),$$

where the matrix-valued function $s \rightarrow H(s)$ is characterised by its Laplace transform:

$$\hat{H}(z) := \int_{\mathbb{R}} H(s) e^{-zs} ds = (zI + P + e^{-z\tau} Q)^{-1}.$$

Equation (14) then follows from an application of Plancherel-Parseval's isometry formula. \square

C. Proof of Theorem 3

We adapt the proof of Ott [19] who considered the case $\mu = 0$, to calculate the steady state variance of this process. Consider the covariance process,

$$C(t) \stackrel{\text{def}}{=} \mathbf{E}[Z_a(s)Z_a(s+t)], \quad (30)$$

which satisfies $C(t) = C(-t)$. It follows directly from (23) that, for $t > 0$,

$$\frac{d}{dt} C(t) = -\gamma C(t - \tau) - \mu C(t). \quad (31)$$

By considering $d(Z_a^2(t))$, using Ito's formula and taking expectations it follows that

$$\gamma C(\tau) + \mu C(0) = \frac{b^2}{2} \quad (32)$$

giving a boundary condition. Now for $0 < t < \tau$,

$$\frac{d}{dt} C(t) = -\gamma C(t - \tau) - \mu C(t) = -\gamma C(\tau - t) - \mu C(t),$$

hence for $|t| \leq \tau/2$,

$$\frac{d}{dt} C\left(\frac{\tau}{2} + t\right) = -\gamma C\left(\frac{\tau}{2} - t\right) - \mu C\left(\frac{\tau}{2} + t\right). \quad (33)$$

Now putting

$$C\left(\frac{\tau}{2} + t\right) = \sum_{j=0}^{\infty} a_j t^j$$

and equating coefficients in equation (33), gives

$$\begin{aligned} C\left(\frac{\tau}{2} + t\right) &= a_0 \left\{ \cos\left(t\sqrt{\gamma^2 - \mu^2}\right) \right. \\ &\quad \left. - \sqrt{\frac{\gamma + \mu}{\gamma - \mu}} \sin\left(t\sqrt{\gamma^2 - \mu^2}\right) \right\} \quad (34) \end{aligned}$$

for $\gamma \geq \mu$, and

$$\begin{aligned} C\left(\frac{\tau}{2} + t\right) &= a_0 \left\{ \cosh\left(t\sqrt{\mu^2 - \gamma^2}\right) \right. \\ &\quad \left. - \sqrt{\frac{\gamma + \mu}{\mu - \gamma}} \sinh\left(t\sqrt{\mu^2 - \gamma^2}\right) \right\} \quad (35) \end{aligned}$$

for $\gamma < \mu$. Hence using the boundary condition (32) to find a_0 , gives the expressions (24)-(25) for $C(0)$. \square

REFERENCES

- [1] F. Kelly, P. Key, and S. Zachary, "Distributed admission control," *IEEE Journal on Selected Areas in Communications*, vol. 18, pp. 2617-2628, 2000.
- [2] L. Breslau, E. W. Knightly, S. Shenker, I. Stoica, and H. Zhang, "Endpoint admission control: Architectural issues and performance," in *SIGCOMM 2000*, 2000, pp. 57-69.
- [3] V. Elek, G. Karlsson, and R. Rönngren, "Admission control based on end-to-end measurements," in *INFOCOM 2000*. IEEE, 2000, www.comnet.technion.ac.il/infocom2000.
- [4] G. Bianchi, A. Capone, and C. Petrioli, "Throughput analysis of end-to-end measurement based admission control in IP," in *INFOCOM 2000*. IEEE, 2000, www.comnet.technion.ac.il/infocom2000.
- [5] T. Kelly, "An ECN probe-based connection acceptance control," *Computer Communication Review*, vol. 31, no. 3, July 2001.
- [6] K. K. Ramakrishnan, S. Floyd, and D. Black, "The addition of Explicit Congestion Notification (ECN) to IP," September 2001.

- [7] A. Bain and P. B. Key, "Modelling the performance of distributed admission control for adaptive applications," *Performance Evaluation Review*, December 2001.
- [8] L. Massoulié, "Stability of distributed congestion control with heterogeneous feedback delays," *IEEE Transactions on Automatic Control*, vol. 47, pp. 895–902, 2002.
- [9] P. Hunt and T. Kurtz, "Large loss networks," *Stochastic Processes and their Applications*, vol. 53, pp. 363–378, 1991.
- [10] A. Mandelbaum, W. A. Massey, and M. I. Reiman, "Strong approximations for Markovian service networks," *Queueing Systems – Theory and Applications*, vol. 30, pp. 149–201, 1998.
- [11] R. Johari and D. Tan, "End-to-end congestion control for the internet: delays and stability," *IEEE/ACM Trans. Networking*, vol. 9, pp. 818–832, 2001.
- [12] G. Vinnicombe, "On the stability of networks operating tcp-like congestion control," in *IFAC conference*, 2002, <http://www-control.eng.cam.ac.uk/gv/internet/index.html>.
- [13] S. Ethier and T. Kurtz, *Markov processes: characterization and convergence*. John Wiley, New York, 1986.
- [14] R. Gibbens and F. Kelly, "Resource pricing and the evolution of congestion control," *Automatica*, 1999.
- [15] R. J. Gibbens, P. B. Key, and S. R. E. Turner, "Properties of the Virtual Queue marking algorithm," in *17th UK Teletraffic Symposium*. IEE, 2001.
- [16] S. Kunniyur and R. Srikant, "Analysis and design of an Adaptive Virtual Queue (AVQ) algorithm for active queue management," in *Proceedings of SIGCOMM 2001*, San Diego, California, USA, 8 2001.
- [17] R. Bellman and K. Cooke, *Differential-Difference Equations*. Academic Press: New York, 1963.
- [18] N. Hayes, "Roots of the transcendental equation associated with a certain differential-difference equation," *Journal of the London Mathematical Society*, vol. 25, pp. 226–232, 1950.
- [19] T. Ott, "On the Ornstein-Uhlenbeck process with delayed feedback," available at http://web.njit.edu/~ott/Papers/Del_LO_U.ps.