

**ANALYSIS OF DYNAMIC CONGESTION  
CONTROL PROTOCOLS---  
A FOKKER-PLANCK APPROXIMATION**

**by**

**Amarnath Mukherjee and John C. Strikwerda**

**Computer Sciences Technical Report #1003**

**February 1991**



**Analysis of  
Dynamic Congestion Control Protocols-  
A Fokker-Planck Approximation**

Amarnath Mukherjee  
*Computer and Information Science Department  
University of Pennsylvania  
Philadelphia, PA 19104*

John C. Strikwerda  
*Department of Computer Sciences  
University of Wisconsin - Madison  
Madison, WI 53706*

*February 1991*



## Abstract

We present an approximate analysis of a queue with dynamically changing input rates that are based on implicit or explicit feedback. This is motivated by recent proposals for adaptive congestion control algorithms [RaJa 88, Jac 88], where the sender's window size at the transport level is adjusted based on perceived congestion level of a bottleneck node. We develop an analysis *methodology* for a simplified system; yet it is powerful enough to answer the important questions regarding stability, convergence (or oscillations), fairness and the significant effect that delayed feedback plays on performance. Specifically, we find that, in the absence of feedback delay, the linear increase/exponential decrease algorithm of Jacobson and Ramakrishnan-Jain [Jac 88, RaJa 88] is *provably* stable and fair. Delayed feedback, on the other hand, introduces oscillations for *every* individual user as well as unfairness across those competing for the same resource. While the simulation study of Zhang [Zha 89] and the fluid-approximation study of Bolot and Shankar [BoSh 90] have observed the oscillations in cumulative queue length and measurements by Jacobson [Jac 88] have revealed some of the unfairness properties, the *reasons* for these have not been identified. We identify *quantitatively* the *cause* of these effects, vis-a-vis the system parameters and properties of the algorithm used.

The model presented is fairly general and can be applied to evaluate the performance of a wide range of feedback control schemes. It is an extension of the classical Fokker-Planck equation. Therefore, it addresses traffic variability (to some extent) that fluid approximation techniques do not address.

## 1. Introduction

We investigate the performance of congestion control protocols that dynamically change input rates based on feedback information received from the network. This is motivated by proposals for adaptive congestion control algorithms [Jac 88, RaJa 88,90], where the sender's window size at the transport layer is adjusted based on perceived congestion level of a bottleneck node.

Demers et. al. [DeKeSh 89] report a simulation study that compares the Jacobson and Ramakrishnan-Jain algorithms [Jac 88, RaJa 88,90] vis-a-vis scheduling disciplines used in intermediate gateways. Zhang [Zhang 89] compares the TCP protocol, which incorporates the Jacobson algorithm, to her Virtual Clock Protocol. In there, she reports some interesting (albeit undesirable) oscillatory properties of the Jacobson algorithm. She also observes that connections with larger number of hops receive a poorer share of an intermediate resource than those with a smaller number of hops. Jacobson had also reported this in his measurements [Jac 88]. Bolot and Shankar [BoSh 90] have recently studied the behavior of the Ramakrishnan-Jain algorithm using a fluid approximation model and they too observe the oscillatory characteristics. Recently, some interesting studies have been reported by Mitra and Seery [MiSe 90, Mit 90] and Shenker [She 90]. Mitra and Seery have developed a new feedback based dynamic window adjustment algorithm based on asymptotic analysis of queueing networks, while Shenker has studied some intrinsic properties of feedback based flow control.

In this study, we develop, from first principles, a Fokker-Planck equation for the evolution of the joint probability density function of queue length and arrival rate at

the bottleneck node. This approximates the *transient* behavior of a queue subjected to adaptive rate-control. We then seek answers to questions regarding *stability* (or oscillations) and *fairness* of a particular adaptive algorithm. We also investigate the effect of *delayed feedback* on performance.

We find that, in the absence of feedback delay, senders using the Jacobson-Ramakrishnan-Jain (or JRJ) Algorithm [Jac 88, RaJa 88,90] (or rather an equivalent rate-based algorithm) *converge* to an equilibrium. Further, this algorithm is *fair* in that all sources sharing a resource get an equal share of the resource if they use the same parameters for adjusting their rates. The exact share of the resource that different sources get when they use different parameters is also determined.

A delay in the feedback information introduces cyclic behavior. If different sources get the feedback information after *different* amounts of delay, then the algorithm may also be *unfair*, i.e., the sources may get unequal throughput. These results strengthen the observations in previous studies and also identify the underlying reasons. For instance, if the adaptive algorithm is linear-increase/exponential-decrease, then the oscillations are due to delayed feedback. However, if the adaptive algorithm is linear-increase/linear-decrease, then the oscillations could be due to both the algorithm itself and the delay in the feedback path. Also unfairness is *partly* due to the larger (feedback) delay suffered by the longer connections as compared to the shorter ones.

The rest of the paper is organized as follows. Section 2 presents the model. Section 3 motivates the analysis methodology. In Section 4, a Fokker-Planck approximation for the time dependent queue behavior is derived. Section 5 discusses the properties of the JRJ-algorithm when only one source is using the resource. Section 6 investigates the properties of the system with multiple sources. Section 7 re-investigates these properties in the presence of delayed feedback. Section 8 presents our conclusions.

## 2. Model

The model we have chosen is motivated by the Jacobson-Ramakrishnan-Jain Algorithm for window adjustment. In the JRJ algorithm, when congestion is detected (by implicit or explicit feedback), the window size is *decreased multiplicatively*. However, when there is no congestion, it is *increased linearly* — to probe for more bandwidth, i.e.,

$$w \leftarrow \begin{cases} w/d; & \text{if congested;} \\ w + a; & \text{if not congested.} \end{cases} \quad (1)$$

While this makes good intuitive sense, it is far from clear as to what *values* the parameters  $a$  and  $d$  should take. Further, it is not *provably* clear if the algorithm is *fair* or *stable*<sup>1</sup> and if so, *under what circumstances*.

To understand the behavior of dynamic congestion control algorithms, we study a queueing system with a time varying input *rate*. The latter is adjusted periodically based

---

<sup>1</sup> An algorithm is *fair* if everybody gets a 'fair' share of the resource (*Fair* share and *equal* share are synonymous if all the demands are equal). Stability, on the other hand, implies that the algorithm converges to a particular value.

on some feedback that the end-point receives about the state of the queue. We are interested in the time evolution of the queue length density function.

Let us assume that we are changing the arrival rate,  $\lambda(t)$ , based on the current queue length,  $Q(t)$ , at some bottleneck node. An example adaptive control algorithm could be

$$\frac{d\lambda}{dt} = \begin{cases} +C_0, & \text{if } Q(t) \leq \bar{q}, \\ -C_1\lambda, & \text{if } Q(t) > \bar{q} \end{cases} \quad (2)$$

where  $\bar{q}$  is some target queue length.  $C_0$  and  $C_1$  are positive constants.

Equation 2 models a *linear increase* in  $\lambda$  for  $Q(t) \leq \bar{q}$  and an *exponential decrease* in it for  $Q(t) > \bar{q}$ . It is therefore the rate-analogue of the dynamic window adjustment algorithm given by Equation 1. For purposes of generality however, we shall denote

$$\frac{d\lambda(\cdot)}{dt} = g(\cdot) \quad (3)$$

$g(\cdot)$  can be viewed as a generic rate-control algorithm.

In the following section, we motivate the methodology chosen. The method adopted can not only lead to a better theoretical understanding of a key problem but also be useful in solving other problems that might involve some form of feedback.

### 3. Methodology

To analyze the effect of Equation (2), Bolot and Shankar [BoSh 90] have used two separate differential equations, one for the queue length,  $Q(t)$ , and another for the arrival rate  $\lambda(t)$ .  $Q(t)$  depends on  $\lambda(t)$  as follows:

$$\frac{dQ(t)}{dt} = \lambda(t - d) - \mu$$

where  $\mu$  is the mean service rate.  $d\lambda(t)/dt$  is given by Equation 2. These are then coupled together, i.e.,  $\lambda(t)$  drives the differential equation for  $Q(t)$  and vice-versa. Their model assumes that  $Q(t)$  and  $\lambda(t)$  are both deterministic.

Suppose, however that  $Q(t)$  were a random variable and one attempts to characterize the time evolution of this process. Consider the classical Fokker-Planck, or diffusion, equation in one-dimension:

$$\frac{\partial f}{\partial t} + \frac{\partial}{\partial q} ((\lambda(\cdot) - \mu)f) = \frac{1}{2} \frac{\partial^2}{\partial q^2} (\sigma^2 f)$$

where  $f(t, q)$  is the probability density function of the queue length at time  $t$ ,  $(\lambda(t, q) - \mu)$  and  $\sigma^2(t, q)$  are the instantaneous mean and variance of the queue growth rate given the queue length is  $q$ . Now, if this equation is to be extended for the congestion-control problem at hand, how should the equation for  $\lambda$  be expressed so that the control part is properly reflected? It turns out that one cannot use a coupled set of equations (one for the density and one for the control) at all.

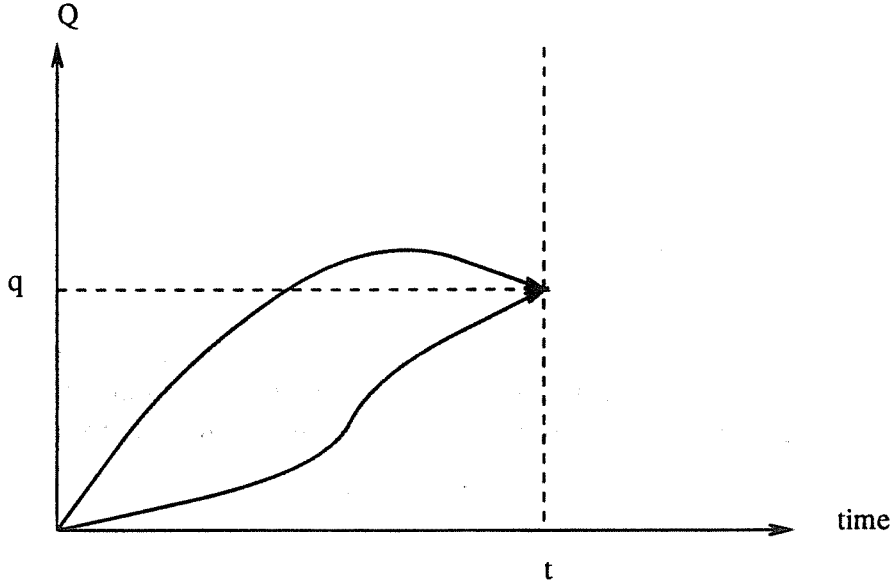


Figure 1: Queue length trajectory as a function of time.

To see this, suppose that  $Q(t)$  were a random variable and say, we were observing the process  $\{(Q(t), \lambda(t))\}$  as time progressed (see Figure 1). Given some initial values  $(Q(0), \lambda(0))$ , let the queue length at time  $t$  be  $Q(t) = q$ , for some  $q$ . At this point, the value of  $\lambda(t)$  is dependent on not just the current value of  $q$ , but also on the *sample path* of  $Q(s)$ ,  $0 \leq s \leq t$ . Intermediate values of the queue length affects  $\lambda$  because of Equation 2 and since the sample path of  $Q$  is random,  $\lambda(t)$  itself is a *random variable*. One cannot therefore couple the two equations.

We hence choose an alternate route. Let  $\mu$  be the average service rate of the queue and let  $\nu(t) = (\lambda(t) - \mu)$  be the instantaneous queue growth rate (with the convention that  $\nu(t) = 0$  if  $Q(t) = 0$  and  $\lambda(t) < \mu$ ). We define  $f(t, q, \nu)$  to be the joint probability density function of  $(Q(t), \nu(t))$ . Our goal is to understand the time dependent behavior of  $f(\cdot)$  based on  $g(\cdot)$  and the variabilities of  $Q(t)$  and  $\nu(t)$ . We address this in the next section.

#### 4. Fokker-Planck approximation for queue with feedback control

Suppose that at time  $t$ , the queue length and queue growth rate are given by  $Q(t) = \hat{q}$  and  $\nu(t) = \hat{\nu}$ . We want to express the density function  $f(t + \tau, q, \nu)$  in terms of  $f(t, \hat{q}, \hat{\nu})$ . We assume that variability in  $\nu$  is caused only by the random sample path of  $Q$  and there is no ‘intrinsic’ variability in  $\nu$ . Then, given  $Q(t + \tau) = q$ , and some small  $\tau$ ,

$$\nu(t + \tau) = \hat{\nu} + g(\cdot)\tau. \quad (4)$$

Let  $h(t + \tau, q, \nu | t, \hat{q}, \hat{\nu})$  be the conditional probability of the transition between  $(\hat{q}, \hat{\nu})$  and  $(q, \nu)$  in time  $(t, t + \tau)$ . Then by the law of total probability,

$$f(t + \tau, q, \nu) = \int \int f(t, \hat{q}, \hat{\nu}) h(t + \tau, q, \nu | t, \hat{q}, \hat{\nu}) d\hat{q} d\hat{\nu} \quad (5)$$



The integral over  $\hat{\nu}$  in Equation 5 is essentially a delta function which is zero for all values of  $\hat{\nu}$  except that satisfying Equation 4. We then have

$$f(t + \tau, q, \nu) = \int f(t, \hat{q}, \hat{\nu}) h(t + \tau, q, \nu | t, \hat{q}, \hat{\nu}) \frac{d\hat{\nu}}{d\nu} d\hat{q} \quad (6)$$

with the understanding that  $\hat{\nu}$  and  $\nu$  are related by Equation 4.

Now, let us further assume that the central limit theorem (approximately) holds for the conditional density function  $h(\cdot)$ , i.e.,

$$h(t + \tau, q, \nu | t, \hat{q}, \hat{\nu}) \approx \eta \left( \frac{q - \hat{q} - \hat{\nu}\tau}{\sigma/\sqrt{\tau}} \right) \quad (7)$$

where  $\sigma^2$  is the variance of  $Q$ . Validity of this assumption is key to the Fokker-Planck approximation that follows.<sup>2</sup>

Combining Equations 6 and 7 gives

$$f(t + \tau, q, \nu) = \int f(t, \hat{q}, \hat{\nu}) \eta \left( \frac{q - \hat{q} - \hat{\nu}\tau}{\sigma/\sqrt{\tau}} \right) \frac{d\hat{\nu}}{d\nu} d\hat{q} \quad (8)$$

To derive the differential equation of  $f(\cdot)$  with respect to time, we subtract  $f(t, q, \nu)$  from both sides, divide by  $\tau$  and let  $\tau \rightarrow 0$ . Using  $d\hat{\nu}/d\nu = 1 - g_\nu\tau$  from Equation 4, we then get<sup>3</sup>

$$\begin{aligned} f_t &= \lim_{\tau \rightarrow 0} \frac{1}{\tau} \int \{f(t, \hat{q}, \hat{\nu})(1 - g_\nu\tau) - f(t, q, \nu)\} \eta(\cdot) d\hat{q} \\ &= \lim_{\tau \rightarrow 0} \frac{1}{\tau} \int \{f(t, \hat{q}, \hat{\nu}) - f(t, q, \nu)\} \eta(\cdot) d\hat{q} - \int g_\nu(\cdot) f(t, q, \nu) \eta(\cdot) d\hat{q} + o(\tau) \\ &= \lim_{\tau \rightarrow 0} \frac{1}{\tau} \int \{f(t, \hat{q}, \hat{\nu}) - f(t, q, \nu)\} \eta(\cdot) d\hat{q} - g_\nu(\cdot) f(t, q, \nu) + o(\tau) \end{aligned} \quad (9)$$

Let

$$I = \lim_{\tau \rightarrow 0} \frac{1}{\tau} \int \{f(t, \hat{q}, \hat{\nu}) - f(t, q, \nu)\} \eta(\cdot) d\hat{q} \quad (10)$$

Adding (and subtracting)  $f(t, q, \hat{\nu})$  to (and from) the right hand side of this equation, we get

$$= \lim_{\tau \rightarrow 0} \frac{1}{\tau} \int \{f(t, \hat{q}, \hat{\nu}) - f(t, q, \hat{\nu})\} \eta(\cdot) d\hat{q} + \frac{1}{\tau} \int \{f(t, q, \hat{\nu}) - f(t, q, \nu)\} \eta(\cdot) d\hat{q} \quad (11)$$

The first integral in Equation 11 is similar to expressions arising in the derivation of the standard Fokker-Planck equation [New 68, New 71, Kle 76]:

$$-\hat{\nu} f_q(t, q, \hat{\nu}) + \frac{1}{2} \sigma^2 f_{qq}(t, q, \hat{\nu})$$

<sup>2</sup> higher order moments may be needed to express more burstiness in  $h$ .

<sup>3</sup> notation:  $f_t = \partial f / \partial t$ ,  $f_q = \partial f / \partial q$ ,  $f_{qq} = \partial^2 f / \partial q^2$  etc.

As  $\tau \rightarrow 0$ ,  $\hat{\nu} \rightarrow \nu$ , (see Equation 4). so this becomes

$$-\nu f_q(t, q, \nu) + \frac{1}{2}\sigma^2 f_{qq}(t, q, \nu) \quad (12)$$

The second integral is equal to

$$\begin{aligned} & \lim_{\tau \rightarrow 0} \frac{1}{\tau} \{f(t, q, \nu - g(\cdot)\tau) - f(t, q, \nu)\} \\ &= \lim_{\tau \rightarrow 0} \frac{1}{\tau} \{[f(t, q, \nu) - g(\cdot)\tau f_\nu(t, q, \nu) - f(t, q, \nu)] + o(\tau)\} \\ &= -g(\cdot) f_\nu(t, q, \nu) \end{aligned} \quad (13)$$

Combining Equations 9, 10, 11, 12 and 13. and noting that  $g_\nu f + g f_\nu = (gf)_\nu$ , we have

$$f_t + \nu f_q + (gf)_\nu = \frac{1}{2}\sigma^2 f_{qq} \quad (14)$$

Equation 14 describes the basic equation of motion for the density function  $f(\cdot)$ .

## 5. Properties of Algorithm 2

We now investigate the properties of Algorithm 2 in conjunction with Equation 14. For the purposes of an intuitive discussion, we suppress the  $\sigma^2$  term in Equation 14 and study a reduced system. We therefore have a hyperbolic partial differential equation whose properties can be explored by studying its *characteristics*. We will address the effects of  $\sigma^2$  being positive at the end of this section. The characteristics are the family of curves satisfying

$$\frac{dq}{dt} = \nu \quad \text{and} \quad \frac{d\nu}{dt} = g.$$

This is equivalent to

$$\frac{dq}{dt} = \lambda - \mu \quad \text{and} \quad \frac{d\lambda}{dt} = g.$$

Consider the  $q - \nu$  diagram of Figure 2. The  $x$ -axis represents the queue length,  $Q$ , and the  $y$ -axis represents the instantaneous queue growth rate,  $\nu$ . Two lines corresponding to  $Q = \bar{q}$  and  $\nu = 0$ , shown by dotted lines, divide the  $q - \nu$  plane into four quadrants. The behavior of Equation 14 is best described by considering each quadrant separately.

First consider Quadrant I in Figure 2. This corresponds to  $\nu > 0$  (i.e.,  $\lambda > \mu$ ) and  $Q < \bar{q}$ . Since  $\lambda > \mu$ , the *instantaneous* queue length at any point in this quadrant is increasing. The instantaneous  $\nu$  is also increasing because  $d\lambda/dt = C_0 > 0$ . The *resultant* direction of instantaneous motion (i.e., the characteristic) is increasing in both  $Q$  and  $\nu$  as shown in the figure. Notice that Equation 14 confirms this intuition: the coefficient of  $f_q$  which represents the  $Q$ -drift is  $\nu$  and this is positive in Quadrant I; the coefficient of  $f_\nu$  which represents the  $\nu$ -drift is  $g(\cdot) = +C_0$  which is positive as well. The *characteristic* is the resultant of these two drifts.

Next, consider Quadrant II. Here  $Q > \bar{q}$  and  $\nu > 0$  (i.e.,  $\lambda > \mu$ ). From Equation 14, the  $Q$ -drift is again positive since  $\nu > 0$ . However, the  $\nu$ -drift is now negative because

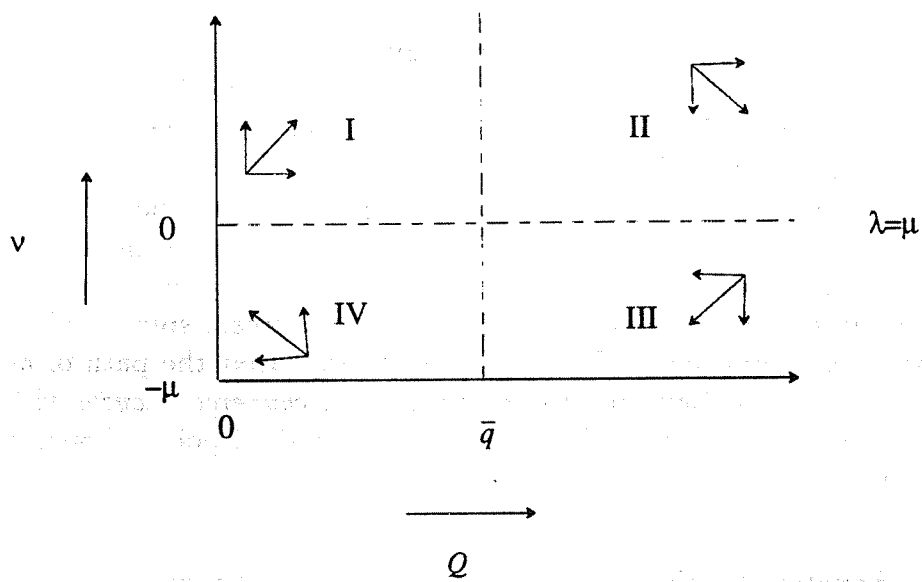


Figure 2: Characteristics and their directions.

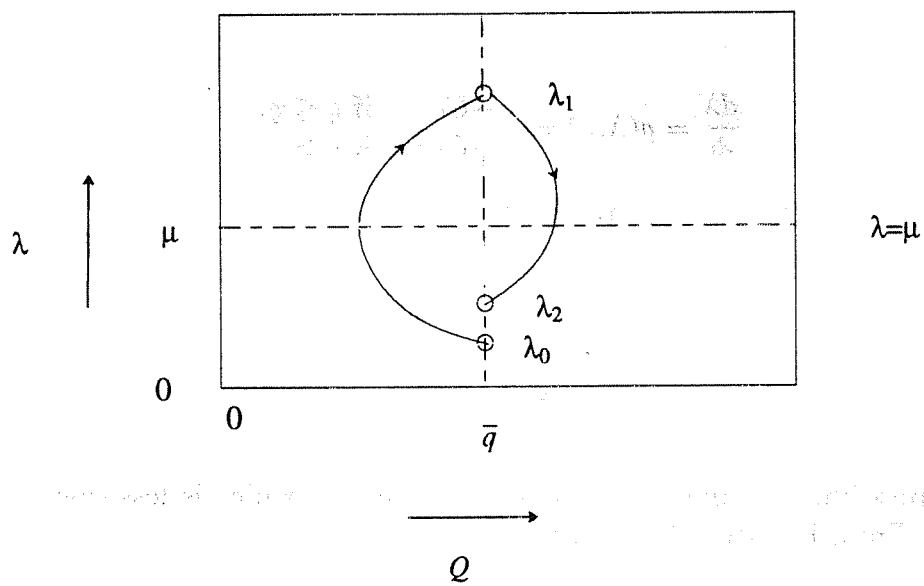


Figure 3: Covergent Spiral

$d\lambda/dt$  is  $-C_1\lambda$  for  $Q > \bar{q}$ . The characteristic, which is the resultant of these two drifts, is increasing in  $Q$  but decreasing in  $\nu$  as shown in Figure 2.

We can similarly check that in Quadrant III, both the  $Q$ -drift and the  $\nu$ -drift are negative while in Quadrant IV, the  $Q$ -drift is negative but the  $\nu$ -drift is positive. The directions of individual drifts and the characteristics are shown in the figure.

Now, suppose we were to trace the path of a 'particle' that obeys both Equation 14 and Equation 2. This path will follow the characteristic. Therefore, from the above argument, it is clear that the trajectory would either be a *cycle* or a *spiral*; the latter could be one that *converges inwards* or *diverges outward*. Further, a convergent spiral could home in to either a *limit point* or a *limit cycle*. Theorem 1 below says that the path of any particle obeying Equations 2 and 14 (ignoring the  $\sigma^2$  term) is a *convergent cycle* with the *limit point*  $Q = \bar{q}$  and  $\nu = 0$ . Notice that this is exactly the desired point of operation of the adaptive algorithm.

**Theorem 1:**

If  $\sigma^2 = 0$  in Equation 14, then Algorithm 2 converges in the limit. The *limit point* is  $q = \bar{q}, \lambda = \mu$ .

*Proof:*

We have

$$\frac{dq}{dt} = \lambda - \mu \quad (15)$$

and

$$\frac{d\lambda}{dt} = g(\lambda, q) = \begin{cases} +C_0, & \text{if } q \leq \bar{q}, \\ -C_1\lambda, & \text{if } q > \bar{q}. \end{cases} \quad (16)$$

Since  $\mu$ , the average service rate, is not changing with time,

$$\begin{aligned} \frac{d^2q}{dt^2} &= g(\lambda, q) = g(\mu + dq/dt, q) \\ &= \begin{cases} +C_0, & \text{if } q \leq \bar{q}, \\ -C_1\lambda, & \text{if } q > \bar{q}. \end{cases} \end{aligned} \quad (17)$$

Now, suppose that at time  $t = 0$ ,  $\lambda$  is some value  $\lambda_0$  which is less than  $\mu$  and  $q$  is  $\bar{q}$  (see Figure 3). From Equation 17, we have

$$\frac{d^2q}{dt^2} = C_0.$$

Its solution is

$$q = \frac{1}{2}C_0t^2 + (\lambda_0 - \mu)t + \bar{q} \quad (18)$$

After a certain time, say  $t_1$ , the characteristic hits  $q = \bar{q}$  line again. Let  $\lambda$  be  $\lambda_1$  now. For the moment, let us assume that the characteristic did not hit the  $q = 0$  boundary, so that Equation 18 is valid all the way up to  $t = t_1$ .

The two roots of Equation 18 with  $q = \bar{q}$  are  $t = 0$  and  $t = 2(\mu - \lambda_0)/C_0$ . The first one corresponds to the initial point. Therefore,

$$t_1 = \frac{2(\mu - \lambda_0)}{C_0} \quad (19)$$

Also, since  $\lambda = \mu + dq/dt$ , we have, from Equation 18 and 19,

$$\begin{aligned} \lambda_1 &= \mu + C_0 t_1 + (\lambda_0 - \mu) \\ &= 2\mu - \lambda_0 \end{aligned} \quad (20)$$

Notice that  $\lambda_1 - \mu$  is equal to  $\mu - \lambda_0$  which says that the *overshoot* above  $\mu$  is *exactly equal to  $\mu - \lambda_0$ , irrespective of the value of  $C_0$* . This is therefore an inherent property of the *linear* increase component of Algorithm 2.

Let us next evaluate the characteristic when  $q$  is greater than  $\bar{q}$ . We have

$$\frac{dq}{dt} = \lambda - \mu \quad (21)$$

and

$$\begin{aligned} \frac{d^2 q}{dt^2} &= g(\lambda, \mu) = -C_1(\mu + dq/dt) \\ \Rightarrow \frac{d^2 q}{dt^2} + C_1 \frac{dq}{dt} + C_1 \mu &= 0 \end{aligned} \quad (22)$$

Since at  $t = t_1$ ,  $q = \bar{q}$  and  $dq/dt = \lambda_1 - \mu$ , its solution is

$$q = -\mu(t - t_1) + \frac{\lambda_1}{C_1} \left(1 - e^{-C_1(t-t_1)}\right) + \bar{q} \quad (23)$$

Let the characteristic again hit the  $q = \bar{q}$  line at some later time  $t_2$  and let  $\lambda$  now be  $\lambda_2$ . Then from Equation 23, we have at time  $t_2$ ,

$$-\mu(t_2 - t_1) + \frac{\lambda_1}{C_1} \left(1 - e^{-C_1(t_2-t_1)}\right) = 0$$

Putting  $\alpha = C_1(t_2 - t_1)$ , we get

$$\mu\alpha = \lambda_1(1 - e^{-\alpha}). \quad (24)$$

Since  $dq/dt$  is equal to  $\lambda - \mu$ ,  $\lambda_2$  can be obtained by differentiating Equation 23. We get

$$\begin{aligned} \lambda_2 &= \lambda_1 e^{-C_1(t_2-t_1)} \\ &= \lambda_1 e^{-\alpha} \end{aligned} \quad (25)$$

Substituting the value of  $\lambda_1$  from Equation 20, we have

$$\lambda_2 = (2\mu - \lambda_0)e^{-\alpha} \quad (26)$$

Therefore

$$\frac{\lambda_2}{\lambda_0} = (2\frac{\mu}{\lambda_0} - 1)e^{-\alpha} \quad (27)$$

The question then is whether  $\lambda_2/\lambda_0$  is greater than 1, less than 1 or equal to 1. From Figure 3, we see that if  $\lambda_2/\lambda_0$  were greater than 1, we would have a converging spiral. We verify next that this is indeed the case.

Let  $\gamma = \mu/\lambda_1$  in Equation 24. Then, using Equation 20, we have

$$\begin{aligned} \gamma &= \frac{\mu}{\lambda_1} = \frac{\mu}{2\mu - \lambda_0} = \frac{1}{2 - \lambda_0/\mu} \\ \Rightarrow 2 - \frac{\lambda_0}{\mu} &= \gamma^{-1} \end{aligned} \quad (28)$$

Substituting into Equation 27, we get

$$\begin{aligned} \frac{\lambda_2}{\lambda_0} &= \left[ \frac{2}{2 - \gamma^{-1}} - 1 \right] e^{-\alpha} \\ &= \left[ \frac{2\gamma}{2\gamma - 1} - 1 \right] e^{-\alpha} \\ &= \left[ \frac{1}{2\gamma - 1} \right] e^{-\alpha} \end{aligned} \quad (29)$$

From Equation 24,  $\gamma$  is given by

$$\gamma = \frac{1 - e^{-\alpha}}{\alpha}$$

Therefore,

$$2\gamma - 1 = \frac{2(1 - e^{-\alpha})}{\alpha} - 1 = \frac{2 - \alpha - 2\alpha e^{-\alpha}}{\alpha} \quad (30)$$

and from 29 and 30,

$$\frac{\lambda_2}{\lambda_0} = \frac{\alpha e^{-\alpha}}{2 - \alpha - 2\alpha e^{-\alpha}} \quad (31)$$

Let us next define a function,  $h(\alpha)$ , such that

$$h(\alpha) = (2 - \alpha - 2e^{-\alpha}) - \alpha e^{-\alpha} \quad (32)$$

If  $h(\alpha)$  is less than 0, then from Equation 31,  $\lambda_2/\lambda_0$  is greater than 1. Notice that  $h(0)$  is 0 and

$$h'(\alpha) = -1 + e^{-\alpha} + \alpha e^{-\alpha}.$$

So,

$$h'(0) = 0.$$

Differentiating once again,

$$h''(\alpha) = -\alpha e^{-\alpha} < 0 \quad \text{for} \quad \alpha > 0$$

Therefore,

$$h'(\alpha) = \int_0^\alpha h''(\alpha) < 0$$

Similarly,

$$h(\alpha) = \int_0^\alpha h'(\alpha) < 0 \quad (33)$$

From Equations 31, 32 and 33, we have

$$\frac{\lambda_2}{\lambda_0} > 1 \quad (34)$$

which implies that *the spiral is convergent*.

So far, we have assumed that the characteristic starting at  $(\bar{q}, \lambda_0)$  never hits the  $q = 0$  boundary. In this case, we have established that we have a convergent spiral. To complete the proof, let us next consider the case when a characteristic hits the left boundary,  $q = 0$ .

Notice that this characteristic cannot hit the boundary for  $\lambda > \mu$ , because the  $q$ -drift which is positive for  $\lambda > \mu$ , will pull it to the right. Therefore, if it hits the  $q = 0$  boundary then  $\lambda \leq \mu$ . Suppose that for some initial  $(\bar{q}, \hat{\lambda}_0)$ , the characteristic barely *touches* the boundary. This point is  $(q = 0, \lambda = \mu)$ , as shown by arc 'a' in Figure 4. Since Equations 18, 19 and 20 hold for this characteristic, it will converge by the earlier argument. Any point corresponding to  $\lambda_0 < \hat{\lambda}_0$  first hits the  $q = 0$  boundary (as shown by arc e), then goes vertically up until  $\lambda = \mu$ , (arc f), and then follows the characteristic corresponding to  $\hat{\lambda}_0$ , (arcs b, c, d). This too, therefore, converges. The partial differential equation 14 is however, not quite valid in this range.

This completes the proof of Theorem 1. ■

**Corollary 1:** If both the increase and the decrease components are *linear*, then the system will never converge.

*Proof:*

We saw from Equation 20 that the amount of overshoot exactly equals the amount of undershoot during the linear increase phase *irrespective of the value of  $C_0$* . The same is true in the reverse direction for a linear decrease algorithm. Hence, the system moves in a non-convergent cycle. ■

We now address the changes that occur due to  $\sigma$  being nonzero and small. To do this, consider an initial state that is zero except for a small rectangle in which the function  $f$  is constant. Assume that this rectangle is to the left of the line  $q = \bar{q}$ . Let the rectangle be given by  $q_1 < q < q_2 < \bar{q}$  and  $\nu_1 < \nu < \nu_2$ . The main mass of the solution will proceed as it would under the influence of the characteristics, but with the additional change due to diffusion in the  $q$  direction. According to our analysis there is no diffusion in the  $\nu$  direction. Thus the solution to the left of  $q = \bar{q}$  will be sharply limited between the two lines  $\nu_1 + C_0 t$  and  $\nu_2 + C_0 t$ . As the solution encounters the line  $q = \bar{q}$ , it will change the direction of motion, and there will be a spreading of the solution in the  $\nu$  direction because of the different times that the different parts encounter the line  $q = \bar{q}$ . The main mass of the solution will follow the path given by the characteristics for small times.

For longer times, the convergence of the characteristics to the limit point, suggests that the probability distribution will converge to a limiting distribution. Most likely the limiting solution will be independent of the initial conditions. More study is required to resolve this speculation. This limiting distribution will be a smooth function, except perhaps at the line  $q = \bar{q}$  where  $g$  changes sign, due to the diffusion in  $q$  and the spreading in  $\nu$ . Note that the steady-state equation

$$\nu f_q + (gf)_\nu = \frac{1}{2}\sigma^2 f_{qq}$$

is locally of parabolic type (with  $\nu$  being the time-like variable and  $q$  being the space-like variable) and thus has infinitely differentiable solutions. The analysis of this equation is nontrivial since the coefficient of the time-like direction changes sign with  $q$ .

## 6. Multiple Sources

We have assumed so far that there is only a single source transmitting through a particular node. We next investigate the properties of the system with multiple sources. Specifically, we are interested in the *convergence* and *fairness* properties when multiple sources compete for a resource. There are two ‘feedback schemes’ that we consider; one where all the sources receive the (same) cumulative queue length information [RaJa 88, Jac 88] and another, where each source receives its own queue length information only.<sup>4</sup> In the latter case, fairness is guaranteed by the scheduler; the analysis of the previous sections then apply directly to each source: if there are  $n$  sources, we change  $\mu$  to  $\mu/n$  and apply Equations 2 and 12. The conclusion is that the system is both convergent and fair.

Next, let us consider the case when all sources receive the common queue length information. All of them adjust their rates according to Algorithm 2. If there are  $n$  sources, let  $(\lambda_1(t), \lambda_2(t), \dots, \lambda_n(t))$  denote their transmission rates at time  $t$ . Let  $\lambda(t) = \sum_{i=1}^n \lambda_i(t)$  be the cumulative transmission rate and let  $Q(t)$  be the cumulative queue length at time  $t$ . Then

$$\frac{d}{dt}\lambda(t) = \sum_{i=1}^n \frac{d}{dt}\lambda_i(t) = \begin{cases} +nC_0, & \text{if } Q(t) \leq \bar{q}, \\ -C_1 \sum_i \lambda_i(t) = -C_1 \lambda(t), & \text{if } Q(t) > \bar{q}. \end{cases} \quad (35)$$

This is the equivalent version of Equation 2 for multiple sources. Equations 12 and 35 completely specify the behavior of the system. From Theorem 1, this system of multiple sources *converges*. Notice that the increase rate is proportional to  $n$ , but the decrease rate is unchanged. Therefore, the length of the spiral trajectory (the path to convergence) is the same, but the time to traverse it is shortened (see Equations 18 and 19).

We next investigate if Algorithm 35 is *fair*. If it is, then the  $\lambda_i$ 's must be equal to each other in the limit.

### Theorem 2:

Algorithm 35 is *fair*.

---

<sup>4</sup> Possible with a Fair-Queue-like scheduling algorithm at the resource.



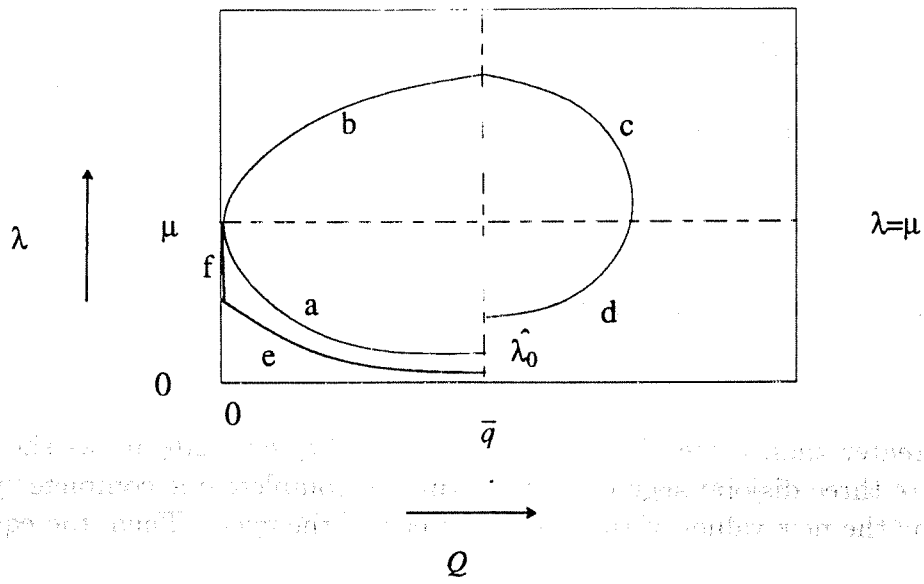


Figure 4: Converging spiral when characteristics touch the  $q=0$  boundary.

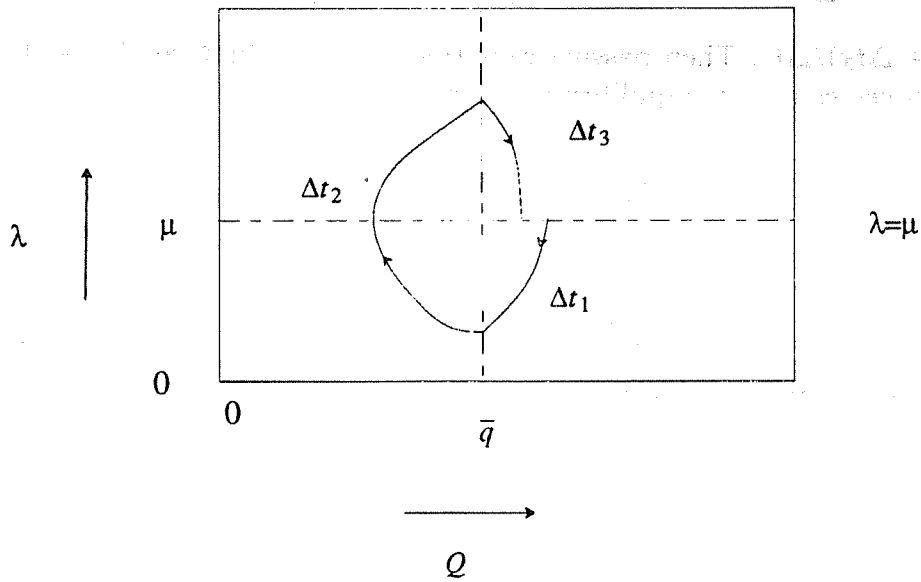


Figure 5: Meaning of  $\Delta t_1, \Delta t_2$  and  $\Delta t_3$ .

*Proof:*

For the purposes of this proof, let us assume that the different sources use different increase and decrease parameters.<sup>5</sup> Suppose there are  $n$  sources and let source  $i$  use an increase parameter  $C_{0,i}$  and a decrease parameter  $C_{1,i}$ . Let  $\lambda_1, \lambda_2, \dots, \lambda_n$  denote their respective transmission rates in the limit (notice that convergence is guaranteed by Theorem 1). Then

$$\lambda = \lambda_1 + \lambda_2 + \dots + \lambda_n = \mu \quad (36)$$

Suppose  $\lambda_1^0, \lambda_2^0, \dots, \lambda_n^0$  are the transmission rates at some time such that

$$\lambda_1^0 + \lambda_2^0 + \dots + \lambda_n^0 = \mu$$

but let  $q$  be greater than  $\bar{q}$  (see Figure 4). Let  $\Delta t_1, \Delta t_2$  and  $\Delta t_3$  be as shown in the figure. These are three disjoint segments of the time to complete one complete cycle.<sup>6</sup> Let  $\lambda_1^1, \lambda_2^1, \dots, \lambda_n^1$  be the new values of the  $\lambda_i$ 's at the end of the cycle. Then, the equation for  $\lambda_1^1$  is given by

$$\lambda_1^1 = (\lambda_1^0 e^{-C_{1,1}\Delta t_1} + C_{0,1}\Delta t_2) e^{-C_{1,1}\Delta t_3}$$

Other  $\lambda_i^1$ 's are similar. We then get.

$$\frac{\lambda_1^1 - \lambda_1^0}{\Delta t_2} = \frac{\lambda_1^0 (e^{-C_{1,1}(\Delta t_1 + \Delta t_3)} - 1)}{\Delta t_2} + C_{0,1} e^{-C_{1,1}\Delta t_3} \quad (37)$$

Let  $\gamma = (\Delta t_1 + \Delta t_3)/\Delta t_2$ . Then passing Equation 37 to the limit as  $\Delta t_2 \rightarrow 0$  which will occur as the processes tend to equilibrium, we get

$$\frac{d\lambda_1}{dt_2} = -C_{1,1}\gamma\lambda_1 + C_{0,1} \quad (38)$$

Similarly,

$$\frac{d\lambda_i}{dt_2} = -C_{1,i}\gamma\lambda_i + C_{0,i} \quad (39)$$

In the limit, when convergence occurs.

$$\frac{d\lambda_i}{dt_2} = 0,$$

so

$$\lambda_i = \frac{C_{0,i}}{\gamma C_{1,i}} \quad (40)$$

Since,  $\sum_i \lambda_i = \mu$ , we have

$$\frac{1}{\gamma} = \frac{\mu}{\sum_i C_{0,i}/C_{1,i}}$$

<sup>5</sup> this is therefore, a more general proof.

<sup>6</sup> i.e., when the process hits  $\lambda = \mu$  and  $q > \bar{q}$  again.

Therefore

$$\lambda_i = \frac{\mu}{\frac{C_{1,i}}{C_{0,i}} \sum_{j=1}^n \frac{C_{0,j}}{C_{1,j}}} \quad (41)$$

Thus, if the  $C_{0,i}$ 's and the  $C_{1,i}$ 's are equal, then  $\lambda_i = \mu/n$ , which implies complete fairness.

■

In real systems, this may be violated because the sources get the feedback information after different amounts of delay and due to finite queue capacity.

## 7. Effect of feedback delay

We next investigate the effect of feedback delay on the control algorithm. Figure 6 shows the mechanics of the system;  $r$  is the delay in obtaining the feedback information from the queue to the control point;  $d$  is the inertia in the forward direction in that it takes the control algorithm this much time to take effect after  $\lambda$  is changed. Let us, for the moment, assume that  $d$  is 0.

The control algorithm can now be precisely stated as:

$$\frac{d}{dt}\lambda(t) = \begin{cases} +C_0. & \text{if } Q(t-r) \leq \bar{q}, \\ -C_1\lambda(t). & \text{if } Q(t-r) > \bar{q} \end{cases} \quad (42)$$

The queueing system with delay in the feedback path is harder to analyze. If we ignore the  $\sigma^2$ -component, the resulting reduced system is readily tractable. The study of the Fokker-Planck equation with delay is the subject of on-going investigation; we shall address the issues involved at the end of this sub-section.

With the  $\sigma^2$ -component deleted, the model is similar to the Bolot and Shankar model [BoSh 90]. It is encouraging that the results obtained are similar. The main difference is that their model assumes  $\bar{q} = 0$ . If  $\bar{q}$  is greater than 0, some interesting properties of  $g(\cdot)$  are revealed, regarding convergence and fairness.

With feedback delay, Algorithm 42 does not converge to a point. To see this, suppose that at time  $t_0$ , the process is at the target equilibrium point  $Q(t_0) = \bar{q}$  and  $\lambda(t_0) = \mu$ . We shall show that it cannot remain here for any significant amount of time.

We need to consider two cases. First, let us say that the process arrived at this point from the left, i.e.,  $Q(t_0 - r) < \bar{q}$ . Then

$$d\lambda(t)/dt = C_0. \quad t \in (t_0, t_0 + r) \quad (43)$$

Therefore

$$\lambda(t_0 + r) = \lambda(t_0) + rC_0 = \mu + rC_0 > \mu \quad (44)$$

and

$$Q(t_0 + r) = \bar{q} + \frac{1}{2}C_0r^2 > \bar{q} \quad (45)$$

Figure 7 shows this pictorially (see Quadrant II). The process overshoots the equilibrium point because  $r > 0$ .

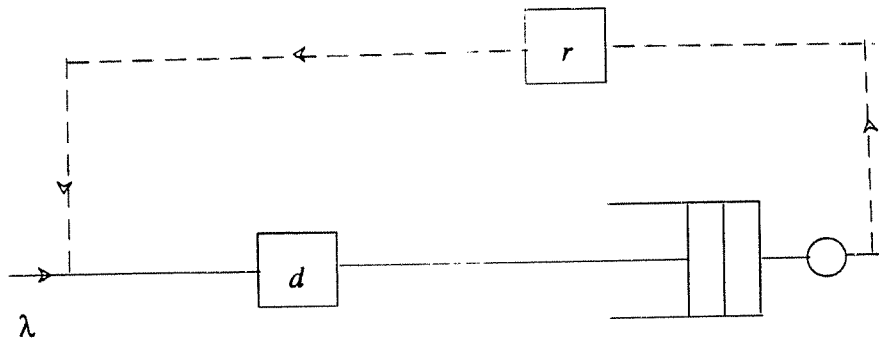


Figure 6: Delayed feedback.

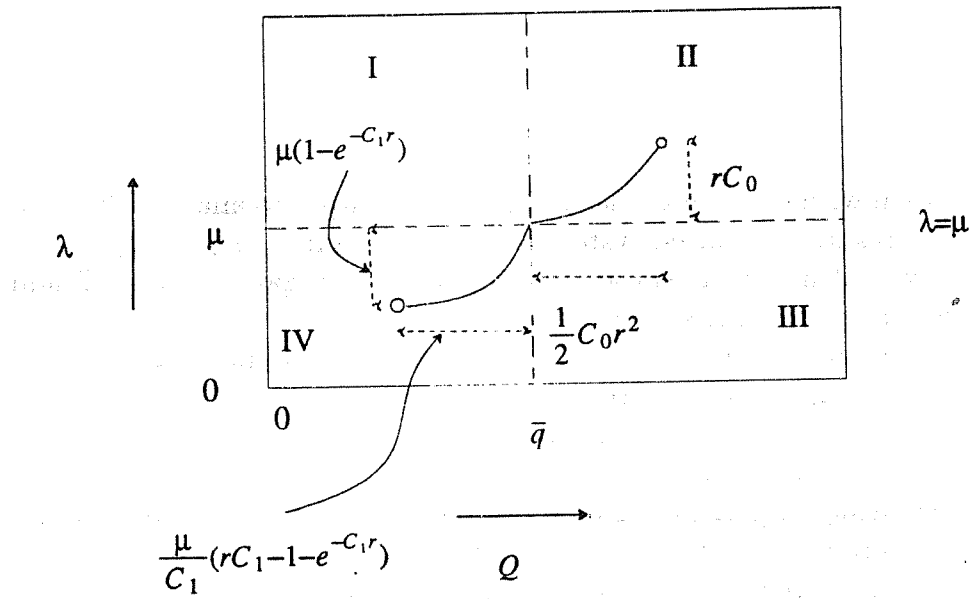


Figure 7: Consequence of delayed feedback.

Next, let us consider the case when the process arrives at  $(\bar{q}, \mu)$  from the right, i.e.,  $Q(t_0 - r) > \bar{q}$ . Then

$$d\lambda(t)/dt = -C_1\lambda(t), \quad t \in (t_0, t_0 + r) \quad (46)$$

Therefore

$$\lambda(t_0 + r) = \lambda(t_0)e^{-C_1r} = \mu e^{-C_1r} < \mu \quad (47)$$

and

$$Q(t_0 + r) = \bar{q} - \frac{\mu}{C_1} (rC_1 - 1 + e^{-C_1r}) < \bar{q} \quad (48)$$

Figure 7 shows this case too (Quadrant IV). The process, here, undershoots the equilibrium.

Notice that the overshoot and the undershoot are going to be larger than what is shown above because when  $Q(t_0) = \bar{q}$ , the value of  $\lambda$  will either be greater than  $\mu$  or less than  $\mu$  (depending on whether the process came from left or right respectively). Clearly the system cannot stabilize at  $(\bar{q}, \mu)$ . Further, at any other point in the  $q - \lambda$  space, the process is forced to be in motion. Therefore the system oscillates.

These oscillations cannot however, become unbounded. To see this, suppose the process is currently at  $(\lambda_0, \bar{q})$  and  $\lambda_0$  is large.  $g()$  is  $C_0$ .  $r$  time units later the control algorithm switches to the exponential decay phase. After some time, say  $\Delta t_1$ , the process hits the  $q = \bar{q}$  line again. Another  $r$  time units later it switches to the linear increase phase. Let this point be  $(\lambda_1, q_1)$ . Then

$$\lambda_1 = (\lambda_0 + C_0r)e^{-C_1(\Delta t_1+r)} \geq 0$$

During the linear increase phase, the process once again hits  $q = \bar{q}$  line (say, after time  $\Delta t_2$ ). Let the value of  $\lambda$  now be  $\lambda_2$ . Then

$$\lambda_2 = \lambda_1 + C_0\Delta t_2$$

Notice that  $\Delta t_2$  is bounded because  $\lambda_1 \geq 0$  and  $q_1 \geq 0$ . Hence,  $\lambda_2$  is bounded from above. This means that if  $\lambda_0$  is large, the diameter of the oscillation has to shrink in the next cycle. This, together with the fact that there can be no stable point, proves the following theorem.

### Theorem 3:

Feedback delay, as expressed by Equation 42, introduces oscillations. These oscillations converge to a limit cycle. ■

While we believe that this limit cycle is unique, we do not have a proof for it.

The diameter of the oscillatory cycle increases with the delay,  $r$ . If different sources experience different delays, they have different oscillatory cycles. This could lead to unfairness in resource usage.

Equations 44, 45, 47 and 48 point to an important difficulty with choosing parameters  $C_0$  and  $C_1$ . The oscillations are larger with higher values of these parameters. Thus, while larger values of  $C_0$  and  $C_1$  help to converge faster in the absence of delay (see Equation 18 for example), they cause larger oscillations in the presence of delay.

Next, let us consider the effect of the inertia  $d$ . We still have

$$\frac{d}{dt}\lambda(t) = \begin{cases} +C_0, & \text{if } Q(t-r) \leq \bar{q}, \\ -C_1\lambda(t), & \text{if } Q(t-r) > \bar{q} \end{cases} \quad (49)$$

However,  $d^2Q/dt^2$  is now given by

$$\frac{d^2}{dt^2}Q(t) = \frac{d}{dt}\lambda(t-d) = \begin{cases} +C_0, & \text{if } Q(t-r-d) \leq \bar{q}, \\ -C_1\lambda(t), & \text{if } Q(t-r-d) > \bar{q} \end{cases} \quad (50)$$

i.e., the queue length now lags  $r+d$  time units while  $\lambda$  still lags  $r$  time units. The oscillatory effect is now more severe, but qualitatively similar to the previous case, i.e., larger values of  $C_0$  and  $C_1$  cause larger oscillations.

The Fokker-Planck equation on delayed feedback looks like:

$$\frac{\partial f}{\partial t} + \nu \frac{\partial f}{\partial q} + \frac{\partial}{\partial \nu} (E[g(Q(t-r)|q, \nu)]f) = \frac{1}{2}\sigma^2 \frac{\partial^2 f}{\partial q^2} \quad (51)$$

The difference is in the appearance of the term  $E[g(Q(t-r)|q, \nu)]$  instead of  $g(q, \nu)$ . The former is the expected value of  $g()$  at time  $t-r$  given that  $Q(t) = q$  and  $\nu(t) = \nu$ , while the latter is the value of  $g()$  at the current time itself. Computing the value of  $E[g(\cdot)]$  turns out to be non-trivial and is the subject of ongoing investigation. One goal is to find a way to do so. A second goal is to find ways to ensure that  $E[g(\cdot)]$  maintains desirable properties of convergence. For example, if one were in Quadrant II in Figure 2, (i.e.,  $Q(t) > \bar{q}$  and  $\lambda(t) > \mu$ ), a desirable property of  $E[g(\cdot)]$  would be for it to be negative and proportional to  $\lambda - \mu$ . While this may in fact turn out to be a formidable problem, it is interesting to see that if the system were deterministic, it can easily be ensured.

In the presence of delayed feedback, one may separate random fluctuations into two categories — those that are short term and those that are medium term. By short-term fluctuations, we mean those which have a time constant smaller than the round-trip delay (or as it turns out two round-trip delays); the feedback mechanism is not useful for tracking this phenomenon. Feedback may however be used to track medium term fluctuations — those that have a larger time constant. To filter out short-term fluctuations, Ramakrishnan-Jain have used averaging of the feedback information over a period of time. Exponential averaging is another method that is often employed.

## 8. Summary and conclusions

We presented an approximate analysis of a queue with dynamically changing input rates based on implicit or explicit feedback. This was motivated by recent proposals for adaptive congestion control algorithms [RaJa 88, 90, Jac 88], where the sender's window size at the transport level was adjusted based on perceived congestion level of a bottleneck node. We developed an analysis methodology for a simplified system; yet it was powerful enough to answer the important questions regarding stability, convergence (or oscillations), fairness and the significant effect that delayed feedback plays on performance. Specifically, we found that, in the absence of feedback delay, the linear increase/exponential decrease