



# LUND UNIVERSITY

## Analysis and design of admission control systems in web-server systems

Robertsson, Anders; Wittenmark, Björn; Kihl, Maria

*Published in:*  
Proceedings of the 2003 American Control Conference. Vol. 1

*DOI:*  
[10.1109/ACC.2003.1238947](https://doi.org/10.1109/ACC.2003.1238947)

2003

[Link to publication](#)

*Citation for published version (APA):*  
Robertsson, A., Wittenmark, B., & Kihl, M. (2003). Analysis and design of admission control systems in web-server systems. In *Proceedings of the 2003 American Control Conference. Vol. 1* (pp. 254-259). IEEE - Institute of Electrical and Electronics Engineers Inc.. <https://doi.org/10.1109/ACC.2003.1238947>

*Total number of authors:*  
3

### General rights

Unless other specific re-use rights are stated the following general rights apply:  
Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Read more about Creative commons licenses: <https://creativecommons.org/licenses/>

### Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

LUND UNIVERSITY

PO Box 117  
221 00 Lund  
+46 46-222 00 00



# Analysis and design of admission control in web-server systems

Anders Robertsson

Department of Automatic Control  
Lund Institute of Technology  
Box 118, SE-221 00 Lund, Sweden  
e-mail: {andersro, bjorn}@control.lth.se

Björn Wittenmark

Maria Kihl  
Department of Communication Systems  
Lund Institute of Technology  
Box 118, SE-221 00 Lund, Sweden  
e-mail: maria@telecom.lth.se

## Abstract

All service control nodes, for example web sites or Mobile Switching Centers, can be modeled as server systems with one or more servers processing the incoming requests. To avoid overload at the service control node different types of admission control mechanisms are usually implemented. In this paper we discuss, from a control point of view, the modeling of a service control-node. The queue is assumed to be an  $M/G/1$ -system and is modeled by a nonlinear flow model and a simplified discrete-time model is used in the analysis and design of the system.

An admission control system based on a PI-controller combined with an anti-reset windup feature is developed. The stability of the admission control system is analyzed and the stability of the closed loop system is proved. The results from the simplified queue model are verified through discrete-event simulations of the system.

## 1 Introduction

Today, many Internet users experience performance problems, such as long response times, when visiting different web sites. A web site consists of one or more web servers that process the incoming HTTP requests sent by the users. When the arrival rate of new requests increases above the maximum service rate, the queues build up and the response times increase. A user experiencing long response times when down-loading a page, will most probably abandon the site, leading to profit loss if the site is commercial,

The performance problems occurring due to overload can be solved by implementing admission control mechanisms in the sites. An admission control mechanism rejects some requests whenever the arrival rate is too high and thereby maintains an acceptable load in the system. Without admission control, the throughput will decrease rapidly when the arrival rate reaches the maximum service rate for the site. This because many requests will experience long response times and thereby be dropped.

Like all computer systems, web servers may be analyzed with the help of queuing theory. [6] show that a web server may be modeled as an  $M/G/1/K$ -system in case the arrival process is Poissonian. However, queuing theory will not be enough when investigating admission control mechanisms. The queuing theoretic methods usually assume stable systems working in the non-overloaded region. Further, there are no mathematical tools in queuing theory for developing optimal admission control mechanisms. Instead, control theory may be used when investigating control mechanisms for queuing systems. However, queuing systems are both stochastic and nonlinear, which introduces a number of modeling problems. These problems may be solved by finding mathematical approximations that mimics the behavior of the system. Usually, a fluid flow model is used in which the arrival process is seen as a continuous incoming flow.

Very few papers have investigated server systems using nonlinear control theory. In [1] and [2] a web server was modeled as a static gain to find optimal controller parameters for a PI-controller. A scheduling algorithm for an Apache web server was designed using system identification methods and linear control theory in [10]. In [9] a discrete-time queuing system with geometrically distributed inter-arrival times and service times was analyzed. An admission control algorithm was developed using optimal control theory. In [14] it is argued that deterministic models cannot be used when analyzing queuing systems. Also, an admission control scheme for an  $M/G/1$ -system is developed using optimal control theory. Recent results for admission control based on linearized models are reported in [7, 12].

This paper models and analyzes a general web-server system, that we model as an  $M/G/1$ -system. We use nonlinear control theory to investigate an admission control mechanism based on a PI-controller. In the investigation we use three different models describing the server system. The most accurate model is the *discrete-event model*. This is a simulation model based on discrete events where the arrival process is a Poisson process and where the service process is a stochastic process with a general distribution, i.e. we have an  $M/G/1$ -system. The second model is a *nonlinear flow model* for  $M/G/1$ -systems. This nonlinear fluid flow approximation is further described in Section 3. The third

model is the *discrete-time model*, which is a discrete-time approximation of the discrete-event model. The discrete-time model, described in Section 4, is used to investigate different designs for admission control instead of using the more complex and time-consuming discrete-event model. The final designs are tested on the discrete-event model.

The main objective of the paper is to show how nonlinear control theory may be used when designing admission control mechanisms for server systems. Further, the paper shows that it is beneficial to use several models when making the design and analysis of the admission controller.

## 2 Problem formulation

A good admission control mechanism will improve the performance of a server system during overload. The main idea is to only admit a certain amount of the requests coming into the server. The general structure of the admission control system is shown in Figure 1. The system consists of four parts: the server system, the controller, the gate, and the monitor. The server system is modeled as a simple queue of M/G/1 type, where the service rate is modeled as a general distribution with mean value  $\bar{\mu}$  and variance of the service rate distribution  $\sigma_\mu^2$ . The arrival process is modeled as a Poisson process with mean arrival rate  $\bar{\lambda}$ . The mean arrival rate will, in practice, vary with time and has to be estimated, but it is in this paper assumed to be constant and known, i.e. we have a stochastic arrival process with known statistics. The queue length in the server is the output of the system. The objective of the controller is to allow as many requests as possible and at the same time be sure that the queue length is below the maximum allowed queue length  $x_{max}$ . The reference value for the queue length is  $x_{ref}$ .

The output of the controller  $u$  is the desired admission rate. Using the control variable  $u$  the gate rejects those requests that cannot be admitted. The gate uses a throttling mechanism, which in the paper is assumed to be *percent blocking*. In this mechanism a certain fraction of the requests are admitted, see [5]. However, the actual admission rate to the system is  $\bar{u}$  and can never be larger than the arrival rate, i.e.

$$\bar{u}(t) = f_{sat}(u(t)) = \begin{cases} 0 & u(t) < 0 \\ u(t) & 0 \leq u(t) \leq \lambda(t) \\ \lambda(t) & u(t) > \lambda(t) \end{cases} \quad (1)$$

This corresponds to the saturation block in Figure 2.

The controller has to work in discrete time. The sampling interval is denoted  $h$  and the control signal is assumed to be constant over time intervals of length  $h$ . The control signal  $u(kh)$  will be calculated based on measurements of queue length  $x(kh)$ . In mathematical terms we want to design the controller and select the reference value  $x_{ref}$  to minimize

$$J_1 = \sum_{k=1}^N (\lambda(kh) - \bar{u}(kh)) \quad (2)$$

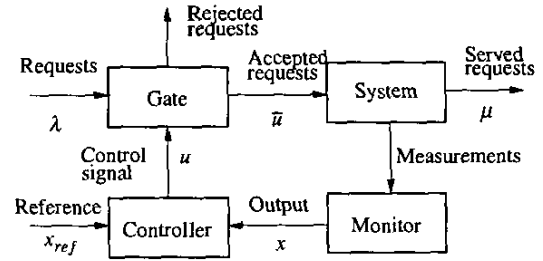


Figure 1: An admission control mechanism.

with respect to  $x_{ref}$  and the controller parameters, subject to the constraint  $x(kh) \leq x_{max}$  over the control horizon  $N$ . Another measure of the performance of the closed-loop system is the variance of the queue length around the desired value  $x_{ref}$ , i.e., to use

$$J_2 = \sum_{k=0}^N (x(kh) - x_{ref})^2 \quad (3)$$

To analyze a queuing system it is necessary to have a mathematical model that mimics the behavior of the discrete-event system with stochastic arrival and service times. This model is called the *discrete-event model*. A simulation model of this M/G/1-system is implemented in C and is used for evaluations of the designs.

## 3 Nonlinear flow models for queues

Since queuing systems have a stochastic behavior it is difficult to find equations that are simple enough to use in the analysis. A nonlinear flow model was first developed by [3] and was further investigated in [11], [13], [15], and [16]. In the references they show that the steady-state behavior of the single server queue is described by the *nonlinear flow model*

$$\frac{dx}{dt} = \bar{u} - \bar{\mu} G(x(t)) \quad (4)$$

For an M/M/1-system we have

$$G(x(t)) = \frac{x(t)}{x(t) + 1} \quad (5)$$

and for an M/G/1-system the expression becomes

$$G(x(t)) = \frac{x(t) + 1 - \sqrt{x^2(t) + 2C^2x(t) + 1}}{1 - C^2} \quad (6)$$

where  $C$  is the variation coefficient, i.e. the ratio of the standard deviation over the mean for the service time distribution (i.e.  $C^2 = 1$  for a M/M/1-system and (6) approaches (5) when  $C^2 \rightarrow 1$ ). The flow model (4) has been verified against the discrete-event model and it has been shown in the cited references that it is correct in terms of average number of customers and service utilization during steady state.

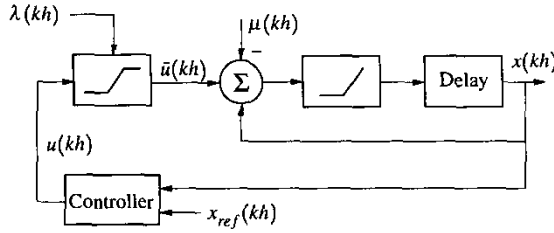


Figure 2: The discrete-time queue model.

#### 4 Discrete-time queuing model

In the design and analysis we will also use the simple process model shown in Figure 2. The model is called the *discrete-time model* and is a simple flow or liquid model in discrete-time, where  $\lambda(kh)$  is the number of requests over the time interval  $[(k-1)h, kh]$  and  $\mu(kh)$  is the number of served requests during the same interval. The queue length at time  $kh$  is  $x(kh)$ . To get an accurate model it is also necessary to limit the queue length such that  $x(kh) \geq 0$ . This is an essential limitation for low load on the system. For heavy loads there will always be jobs waiting for service.

In the analysis of the discrete-time model  $\lambda(kh)$  and  $\mu(kh)$  are stochastic processes that are integrated stochastic processes over one sampling period. If the incoming service rate is a Poisson process with arrival rate  $\lambda_c$  then  $\lambda(kh)$  is also a Poisson process with arrival rate  $\bar{\lambda} = \lambda_c h$ .  $\mu(kh)$  is a stochastic process with a distribution obtained from the underlying service rate distribution. In this paper it is assumed that  $\lambda(kh)$  and  $\mu(kh)$  are independent from between sampling instants and uncorrelated to each other.

The discrete-time model is given by

$$x(kh+h) = f(x(kh) + \bar{u}(kh) - \mu(kh)) \quad (7)$$

where the limit function

$$f(z) = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}$$

assures that  $x(kh+h) \geq 0$ . When the limit function is disregarded then the queue is a discrete-time integrator.

The discrete-time model is an averaging model in the sense that we are not considering the specific timing of the different events, arrivals, or departures from the queue. We assume that the sampling period is sufficiently long to guarantee that the "quantization effects" around the sampling times are negligible.

We will use the steady-state properties of the flow model to verify the relevance of the discrete-time model. The number of requests arriving during an interval  $h$  is Poisson distributed with mean  $\bar{\lambda}h$ . The service process is generated such that it corresponds to  $C^2 = 3.74$ , which implies that we use (6) in (4). The discrete-time model has been simulated for different loads  $\rho = \bar{\lambda}/\bar{\mu}$  and the average queue length

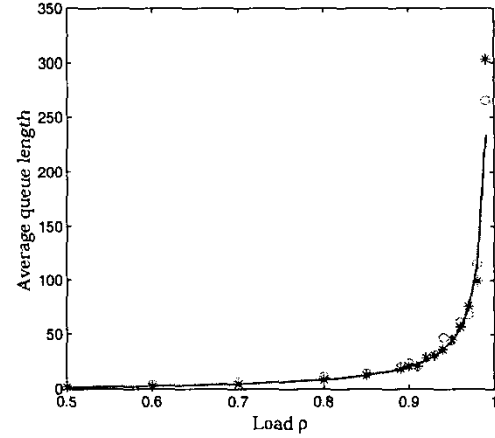


Figure 3: The average queue length of the discrete-time model when simulating an M/G/1-system with Poisson distributed arrivals and a service rate corresponding to  $C^2 = 3.74$  in (6). The different load cases for  $\rho$  are shown with asterisks. The steady-state queue length obtained from (4) and (5) is shown as a full line. The discrete-event simulation results are marked with rings.

has been estimated for  $5 \cdot 10^4$  sampling intervals of length 1 second. Figure 3 shows that there is a good agreement in the steady-state behavior of the discrete-time, the non-linear flow, and the discrete event models. Based on these simulations we regard the discrete-time model to be accurate, at least for steady-state analysis. For transient properties we have to evaluate the discrete-time model against the discrete-event model.

#### 5 Design of admission controller

The choice of controlled variable or output is an important issue when developing an admission control scheme. First, the output must be easy to measure. Second, the value of the output must accurately show the status of the controlled system. Finally, the output must in some way relate to the Quality of Service (QoS) demands on the system.

Traditionally, server utilization or queue lengths have been the output variables mostly used in admission control schemes. For many service control nodes the main objective of the control scheme is to protect a single server system from overload. As long as the server utilization is below a certain level, the response times are low. Therefore, server utilization and queue length are two appropriate variables that are easy to measure. In this paper we use the queue length as the main measurement of the status of the system. Control of utilization and server load is discussed in [8], where also other arrival processes than the Poisson distribution are considered.

A PI-controller will now be designed to investigate different types of admission controllers. Using (4) and (6) and

linearizing the system around  $x_{ref}$  and  $\bar{u}^0 = \bar{\mu}G(x_{ref})$  we get the linearized system

$$\frac{d\Delta x}{dt} = -a\Delta x + \Delta \bar{u}$$

where

$$a = \frac{\bar{\mu}}{1-C^2} \left( 1 - \frac{x_{ref} + C^2}{\sqrt{x_{ref}^2 + 2C^2x_{ref} + 1}} \right)$$

Notice that  $a \in (0, \bar{\mu}]$  for  $x_{ref} \geq 0$  and when  $x_{ref}$  is large we get  $a \approx 0$ , i.e. the queue can be approximated by an integrator. Compare the discrete-time model.

Let the controller be a PI-controller with

$$\Delta \bar{u}(t) = K_c \left( e(t) + \frac{1}{T_i} \int e(\tau) d\tau \right)$$

where  $e(t) = x_{ref}(t) - x(t)$ . The closed loop system is in Laplace transfer form given by

$$\Delta X(s) = \frac{K_c(s + 1/T_i)}{s^2 + (a + K_c)s + K_c/T_i} \Delta X_{ref}(s) \quad (8)$$

Assume that the desired characteristic equation is

$$s^2 + a_{m1}s + a_{m2} = 0 \quad (9)$$

The values of the controller parameters that gives this are

$$K_c = a_{m1} - a \quad T_i = \frac{a_{m1} - a}{a_{m2}} \quad (10)$$

Depending on the desired speed of the response of the closed loop system we can determine the coefficients in (9) and from (10) obtain the controller parameters.

When implementing the PI-controller we will use a discrete-time implementation. The saturation function given by (1) depends on the number of arriving requests that are really arriving during the sampling period. When there is a saturation in the process and when an integrating controller is used it is necessary to introduce anti-reset windup in the controller, see [4]. Since  $\lambda(kh)$  is a stochastic process we will instead use  $\bar{\lambda}$  in the anti-reset windup device. The stability of the closed loop system with the saturation is investigated in Section 6.

In heavy load situations we can simplify the discrete-time model (7) to

$$x(kh+h) = x(kh) + \bar{u}(kh) - \mu(kh) \quad (11)$$

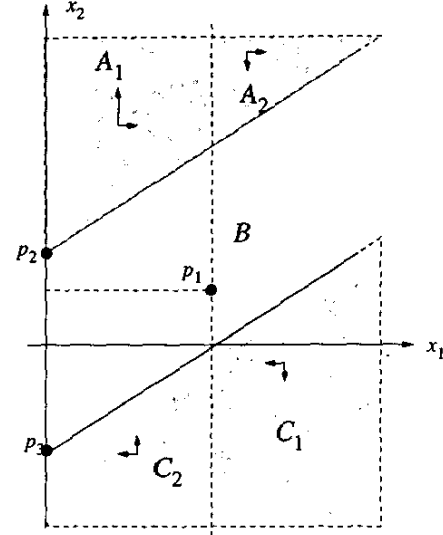
Introduce

$$\mu(kh) = \bar{\mu} + \Delta\mu(kh)$$

where  $\Delta\mu(kh)$  is an independent identically distributed stochastic process with zero mean. The minimum variance controller, see [4], around  $x_{ref}$  is given by the controller

$$\bar{u}(kh) = x_{ref} - x(kh) + \bar{\mu} \quad (12)$$

The controller (12) can be interpreted as a proportional controller with  $K_c = 1$  and a bias term  $\bar{\mu}$ . The bias term plays the same role as the integrator in the PI-controller.



**Figure 4:** Phase-plane for the controlled system. The equilibrium  $p_1 = (x_{ref}, \bar{\mu}G(x_{ref}))$ , belongs to region B for  $\bar{\mu} < \bar{\lambda}$  and  $x_{ref} > 0$ .

## 6 Stability of PI-controlled M/G/1-systems

This section contains a deterministic stability analysis of the PI-controlled nonlinear flow model for an M/G/1-system. The stability region for the controlled system is determined. This means that we evaluate for which values of the controller parameters  $K_c$  and  $T_i$  the controlled system is stable.

Let the controlled system be described by ( $K_c > 0, T_i > 0$ )

$$\dot{x}_1 = -\bar{\mu}G(x_1) + f_{sat}(K_c e + x_2)$$

$$\dot{x}_2 = \frac{K_c}{T_i} e$$

$$e = x_{ref} - x_1$$

$$G(x_1) = \frac{x_1 + 1 - \sqrt{x_1^2 + 2C^2x_1 + 1}}{1 - C^2}, \quad C \neq 1$$

The states  $x_1$  and  $x_2$  correspond to the queue length and the integrator state in the controller, respectively. The function  $f_{sat}(z)$  corresponds to having a controlled non-negative flow upper limited by the arrival rate  $\bar{\lambda}$ , see (1).  $K_c$  and  $T_i$  are positive controller parameters to be designed. Notice that the nonlinearity  $G \in [0, 1]$  and that  $\frac{dG}{dx} \in [0, 1]$  and is monotonically decreasing.

Depending on the value of  $K_c e + x_2$ , the phase plane can be divided into three regions: A, B and C, see Figure 4.

- Region A =  $\{(x_1, x_2) | K_c e + x_2 > \bar{\lambda}\}$ 

$$\begin{aligned}\dot{x}_1 &= -\bar{\mu}G(x_1) + \bar{\lambda} \\ \dot{x}_2 &= \frac{K_c}{T_i}(x_{ref} - x_1)\end{aligned}\quad (13)$$

- Region B =  $\{(x_1, x_2) | 0 \leq K_c e + x_2 \leq \bar{\lambda}\}$ 

$$\begin{aligned}\dot{x}_1 &= -\bar{\mu}G(x_1) + K_c(x_{ref} - x_1) + x_2 \\ \dot{x}_2 &= \frac{K_c}{T_i}(x_{ref} - x_1)\end{aligned}$$

- Region C =  $\{(x_1, x_2) | K_c e + x_2 < 0\}$ 

$$\begin{aligned}\dot{x}_1 &= -\bar{\mu}G(x_1) \\ \dot{x}_2 &= \frac{K_c}{T_i}(x_{ref} - x_1)\end{aligned}\quad (14)$$

Furthermore, consider the following subregions, depicted in Figure 4

$$\begin{aligned}A_1 &= \{(x_1, x_2) \in A | x_1 < x_{ref}\}, & A_2 &= \{(x_1, x_2) \in A | x_1 \geq x_{ref}\} \\ C_1 &= \{(x_1, x_2) \in C | x_1 > x_{ref}\}, & C_2 &= \{(x_1, x_2) \in C | x_1 \leq x_{ref}\}\end{aligned}$$

### 6.1 Stability of equilibrium

The system will have the equilibrium, see Section 5,

$$(x_1^o, x_2^o) = (x_{ref}, \bar{\mu}G(x_{ref}))$$

which will belong to region B provided that

$$\bar{\mu} < \bar{\lambda} \text{ and } x_{ref} > 0$$

The local stability property of the equilibrium is given by the linearization around the equilibrium point with the Jacobian matrix  $J_A$ , where

$$J_A = \frac{df}{dx|_{(x_1^o, x_2^o)}} = \begin{bmatrix} -K_c - \bar{\mu} \frac{dG}{dx_1}|_{x_1=x_{ref}} & 1 \\ -K_c/T_i & 0 \end{bmatrix}$$

The characteristic polynomial of the system matrix will be given by the denominator of (8) with  $a \geq 0$ .

For positive parameters  $K_c$  and  $T_i$ , the equilibrium will thus be locally asymptotically stable, since the roots of the characteristic polynomial will be in the left half-plane.

It can be shown that for  $x_{ref} > 0$  there are no other equilibria for this system. It remains to show that all trajectories starting outside the unsaturated region will enter this region and thereby globally converge to the equilibrium point.

### 6.2 Transition between regions

From (13) and (14) we see that the  $x_1$ -dynamics in region A and C are independent of the  $x_2$ -dynamics. Furthermore, the value of  $x_2$  is decreasing in all regions A, B, and C when

$x_1 > x_{ref}$ . For  $\bar{\lambda} > \bar{\mu}$ , any trajectory starting in region  $A_1$  will have  $\dot{x}_1 > 0$  and  $\dot{x}_2 > 0$ , so the trajectory will thus leave region  $A_1$ , either into region B or into region  $A_2$ . Once in region  $A_2$  we have that  $\dot{x}_1 > 0$  and  $\dot{x}_2 < 0$  so the trajectory will eventually leave region A into B.

In the same way, for any trajectory starting in region  $C_1$  we have that  $(\dot{x}_1 < 0, \dot{x}_2 < 0)$  so the trajectory will either leave into region B or into  $C_2$ . Once in region  $C_2$  we have that  $(\dot{x}_1 < 0, \dot{x}_2 > 0)$  if  $x_1 > 0$  and  $(\dot{x}_1 = 0, \dot{x}_2 > 0)$  if  $x_1 = 0$ , so it will eventually enter region B.

If there exists a trajectory starting in a point on the boundary, say, for instance,  $(x_1, x_2) = (0, 0)$ , which will reach the equilibrium we can conclude that no other trajectory can cross that since the system is autonomous. Thus, there can be no limit cycle and we have asymptotic stability for all trajectories starting in  $x_1 \geq 0$  (that is, for all feasible initial queue lengths).

## 7 Numerical investigation

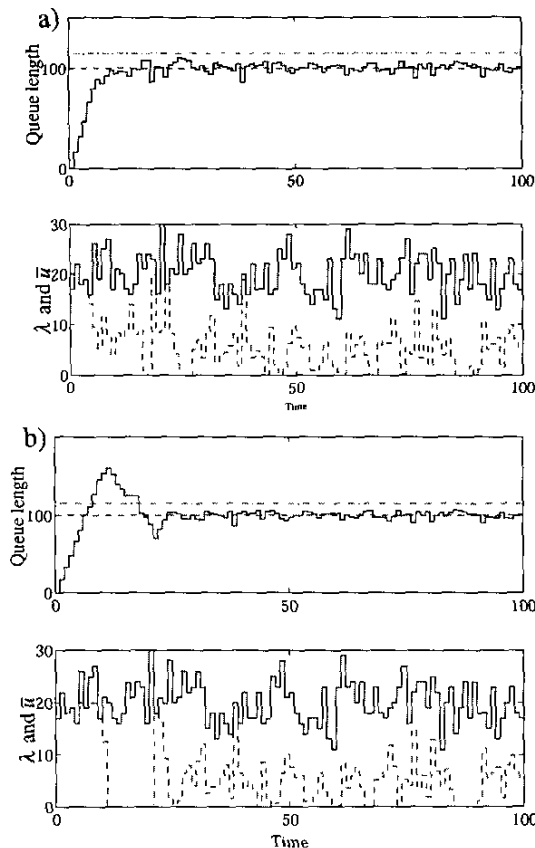
The numerical investigations have shown that the anti-reset windup mechanism is important for the transient behavior of the queue-length control at start-up.

The loss function (2) is after a short transient practically independent of  $x_{ref}$  and the controller parameters as long as the closed-loop system is stable. The reason is that when the queue has settled around the desired reference value then the mean value of (3) will depend on difference of the mean values of the arrival and service rates. A discrete-time PI-controller designed such that the sampled data system has two poles in  $z = 0.5$  gives only marginally larger loss  $J_2$  than when the minimum variance controller (12) is used.

The transient responses with and without antireset windup compensation are shown in Figure 5 and align well with the transient for the discrete-event model (not included). The distributions of the controlled queue lengths for the discrete-time and discrete-event models around a fixed reference are shown in Figure 6. The simulations show a good match between the models.

## 8 Conclusions

In the paper we have shown different ways of modeling a web-server system. The admission control problem is here solved with a simple PI-controller with an anti-reset windup modification and the stability of the closed loop system is investigated. The simulations have shown that the anti-reset windup is important when implementing the admission controller. Further research is needed on how to model the arrival and service processes and how to measure the quality of service of the controlled system.



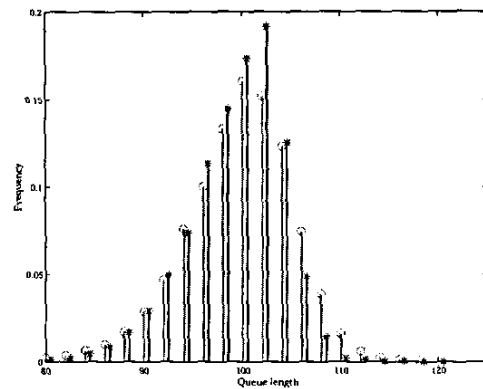
**Figure 5:** Transient behavior of queue-length control for the discrete-time model. The queue length and the arrival- (full) and control (dashed) signals for a) PI-control with antireset windup compensation b) PI-control without antireset windup compensation.

### Acknowledgments

The work in this paper has partially been supported by the Swedish Research Council through the Multi Project Grant 621-2001-3020 and contract 621-2001-3053.

### References

- [1] T.F. Abdelzaher and C. Lu. Modeling and performance control of Internet servers. In *Proc. of the 39th IEEE Conference on Decision and Control*, pages 2234–2239, 2000.
- [2] T.F. Abdelzaher, K.G. Shin, and N. Bhatti. Performance guarantees for web server end-systems: A control theoretic approach. *IEEE Transactions on Parallel and Distributed Systems*, 13(1):80–96, Jan 2002.
- [3] C. E. Agnew. Dynamic modeling and control of congestion-prone systems. *Operations Research*, 24(3):400–419, May–June 1976.
- [4] Karl Johan Åström and Björn Wittenmark. *Computer-Controlled Systems*. Prentice Hall, 1997.
- [5] A. Berger. Comparison of call gapping and percent blocking for overload control in distributed switching systems and



**Figure 6:** Histogram of the average queue length for the PI-controlled discrete-time system ('\*') and for the discrete-event system ('o'), respectively. A PI-controller without anti-windup was used.

telecommunications networks. *IEEE Trans. on Communications*, 39:407–414, 1991.

- [6] J. Cao, M. Andersson, C. Nyberg, and M. Kihl. Web server performance modelling using an M/G/1/K\*PS queue. In *Proc. of International Conference on Telecommunications*, 2003.
- [7] N. Gandhi, D.M. Tilbury, Y. Diao, J. Hellerstein, and S. . MIMO control of an Apache web server — modeling and controller design. In *Proc. of the American Control Conference 2002*, pages 4922–4927, 2002.
- [8] M. Kihl, A. Robertsson, and B. Wittenmark. Analysis of admission control mechanisms using non-linear control theory. Accepted to IEEE Int Symp on Computer Communications 2003.
- [9] J. Kuri and A. Kumar. Optimal control of arrivals to queues with delayed queue length information. *IEEE Transactions on Automatic Control*, 40(8):1444–1450, Aug 1995.
- [10] C. Lu, T.F. Abdelzaher, J.A. Stankovic, and S.H. Son. A feedback control approach for guaranteeing relative delays in web servers. In *Proc. of the 7th IEEE Real-Time Technology and Applications Symposium*, pages 51–62, 2001.
- [11] A. Pitsillides, J. Lambert, and D. Tipper. A multilevel optimal control approach to dynamic bandwidth allocation in broadband ISDN. *Telecommunication Systems*, 4:71–96, 1995.
- [12] L. Sha, X. Liu, Y. Lu, and T. Abdelzaher. Queueing model based network server performance control. In *Proc of the Real-Time Systems Symposium 2002*, pages 81–90, 2002.
- [13] S. Sharma and D. Tipper. Approximate models for the study of nonstationary queues and their applications to communication networks. In *Proc. of IEEE International Conference on Communications*, pages 352–358, 1993.
- [14] S. Stidham Jr. Optimal control of admission to a queueing system. *IEEE Transactions on Automatic Control*, 30(8):705–713, 1985.
- [15] D. Tipper and M. K. Sundareshan. Numerical models for modeling computer networks under nonstationary conditions. *IEEE Journal on Selected Areas in Communications*, 8(9):1682–1695, Dec 1990.
- [16] W. Wang, D. Tipper, and S. Banerjee. A simple approximation for modeling nonstationary queues. In *Proc. of IEEE Info-com '96*, pages 255–262, 1996.