# UC Irvine
## UC Irvine Previously Published Works

**Title**

Violation Probability in Processor-Sharing Queues

**Permalink**

https://escholarship.org/uc/item/8q09k6k9

**Journal**

IEEE Transactions on  Automatic Control, 53(8)

**ISSN**

**Author**

Jordan, Scott

**Publication Date**

2008-09-01

**DOI**

Peer reviewed

# Violation Probability in Processor-Sharing Queues

Na Chen and Scott Jordan

*Abstract*—**Processor-sharing queues are often used to model file transmission in networks. While sojourn time is a common performance metric in the queueing literature,** *average transmission rate* **is the more commonly discussed metric in the networking literature. Whereas much is known about sojourn times, there is little known about the average service rate experienced by jobs in processor-sharing queues. We focus here upon performance requirements in the form of an upper bound on the probability of failing to achieve a specified minimum transmission rate or a specified minimum average rate. For an M/G/1 processor-sharing queue, we give a closed-form expression for this violation probability. We derive closed-form expressions for the marginal service rate with respect to the violation probability and to the minimum transmission rate, and characterize when each is binding. We then consider the effect of using connection access control by modeling an M/G/1/K processor-sharing queue, and discuss the relationship between queue service rate, queue limit, violation probability, and blocking probability. Finally, we consider a two-class discriminatory processor-sharing queue, and discuss what combinations of class weighting and service rate can be used to achieve specified minimum rate violation probabilities for both classes.**

*Index Terms*—**Discriminatory processor-sharing (DPS), transmission control protocol (TCP).**

## I. INTRODUCTION

We are motivated here by the goal of providing performance guarantees for elastic data applications which require a specified transmission rate or higher at least a specified proportion of time. Intuition suggests that the cost of providing such a threshold violation probability guarantee should be increasing with the level of the guarantee, but this intuition has not yet been grounded with a theoretical basis. In this paper, we consider such guarantees in the context of processor-sharing queues and a simple discriminatory processor-sharing queue.

Processor-sharing (PS) queues were originally intended to model time-sharing computer systems, and nowadays are also widely used to analyze the call-level performance of bandwidth sharing in communication systems (cf. [1]). Each user or job represents transmission of a file, and the queue service rate (in bits/sec) represents the bandwidth, which is shared equally among multiple transmissions. A user starts transmission when it arrives and departs when the transmission has completed. The transmission rate of each user changes whenever the number of users in the system changes.

The processor-sharing service discipline is an appropriate model when the time scale of interest is call-level and all objects share bandwidth equally. The call-level time scale applies when the relevant performance metrics are measured over the typical length of a file transmission; if the relevant metrics are measured on a packet-level time scale, then the scheduler is usually modeled as one that swaps between jobs. The equal bandwidth assumption is often made when there exists a mechanism, e.g., the Transmission Control Protocol

(TCP), that attempts to equalize bandwidth between multiple streams over multiple round trip times.

There is a rich literature concerning processor-sharing queues and discriminatory processor-sharing (DPS) queues. The most important performance metric for such queues is sojourn time. When modeling file transmission, sojourn time corresponds to the time required to completely transmit the file, which is certainly of interest. For M/M/1-PS queues, Coffman [2] derived the Laplace transform of the waiting time distribution conditioned on the required service time and the number of users in the system seen on arrival. By removing the conditioning and inverting the Laplace transform, Morrison [3] obtained an integral representation for the complementary distribution of the sojourn time, which was further refined by Guillemin [4] via spectral theory to obtain the sojourn time distribution conditioned on the number of users in the system seen on arrival. For M/M/1/K-PS queues, Morrison [5] obtained an asymptotic approximation to the distribution of the waiting time. In heavy traffic, Morrison [6] also found the distribution of the response time conditioned on the required service time, and Knessl [7] constructed an asymptotic approximation to the sojourn time distribution.

For discriminatory processor-sharing queues, Fayolle [8] derived the expected sojourn time conditioned on the required service time. Rege and Sengupta decomposed the conditional sojourn time into constituent parts which can be obtained by solving a system of non-linear integral equations [9], and found linear simultaneous equations for moments of the stationary distribution [10]. Borst [11] derived the sojourn time asymptotics for a G/G/1-DPS queue with regularly varying service requirements and proved that the sojourn time of a given class has the same tail behavior as its service requirement.

However, the most common performance metric for data applications in the Internet is throughput, not sojourn time. Indeed, many Internet service providers advertise a speed of some type when selling residential broadband service, and there are many online speed tests that measure throughput on a broadband connection. Throughput is casually perceived as the rate at which a computer or network sends or receives data. However, a more precise definition of throughput is with respect to a time window, as the number of bits transmitted divided by the length of the time window. The time window is traditionally chosen to correspond to the time scale on which users judge performance, typically ranging from tenths of a second for highly interactive applications such as gaming, to seconds for moderately interactive applications such as web browsing, to minutes for non-interactive applications such as file downloads.

While there are many available results on sojourn time in PS queues, there is little literature on the throughput in such queues. Definitions of average rate as observed by users and of average rate as observed by the queue were introduced in [12]. Comparisons between the two, and the marginal costs associated with performance requirements in terms of average rate were provided in [13]. Other common metrics include slowdown (cf. [14]), mean slowdown (cf. [15]), and flow throughput (cf. [1]). We have found no literature addressing threshold violation probabilities of throughput in PS queues.

In Section II, we introduce definitions of minimum rate violation probability and minimum throughput violation probability. In Section III, we give a closed-form expression for the minimum rate violation probability, derive closed-form expressions for the marginal service rate with respect to the violation probability and to the rate threshold, and characterize when each is binding. In Section IV, we then consider the effect of using connection access control by modeling an M/M/1/K-PS queue, and discuss the relationship between queue service rate, queue size, violation probability, and blocking

probability. In Section V, we turn to a two-class DPS queue, and discuss what combinations of class weighting and service rate can be used to achieve specified violation probabilities for both classes.

We hope that these results may inspire researchers to consider violation probabilities of throughput when they are more closely related to application performance than traditional metrics such as sojourn time. The expressions for marginal service rate with respect to these metrics may be used to guide dimensioning algorithms, and to set prices based on the marginal cost of the required resources.

## II. VIOLATION PROBABILITIES

In this section, we present definitions of minimum rate violation probability seen by the queue and minimum throughput violation probability seen by users. Consider a processor-sharing queue. Let $n(u)$ and $r(u)$ denote the number of users in the system and the transmission rate per user at time $u$ respectively.

From the queue's perspective, the bandwidth is split among all users present in the system; therefore, the transmission rate per user changes immediately after a user arrives or departs. A violation probability as observed by the queue is defined as the proportion of time the (instantaneous) transmission rate of a user in the system drops below a specified value. We call this quantity the *minimum rate violation probability*

$$P^I(x) \equiv \lim_{t \to \infty} \frac{\int_0^t I_{\{r(u) < x, \, n(u) > 0\}} du}{\int_0^t I_{\{n(u) > 0\}} du} \qquad (1)$$

where $x$ is the minimum transmission rate, and $I_{\{.\}}$ is the indicator function. Under the PS discipline, $r(u)$ is inversely proportional to $n(u)$ for $n(u) > 0$; thus, in the steady state, $P^I(x)$ depends on the stationary distribution of the number of users.

From the users' perspective, the average transmission rate (throughput) of user $i$, $r_i$, is defined as its job length divided by the time required to finish this job. A violation probability as observed by users can be defined by the proportion of users with average transmission rate lower than a specified value. We call this quantity the *minimum throughput violation probability*

$$P^A(x) \equiv \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^n I_{\{r_i < x\}}. \qquad (2)$$

In the steady state, $P^A(x)$ relies on the distribution of average transmission rate, and therefore depends on the stationary distribution of the sojourn time conditioned on the job length.

## III. VIOLATION PROBABILITY IN AN M/G/1-PS QUEUE

In this section, we investigate the minimum rate violation probability and the minimum throughput violation probability in an M/G/1-PS queue. Starting by examining the expressions of $P^I(x)$ and $P^A(x)$, we then turn to performance requirements in the form of $P^I(x) \le p$ and in the form of $P^A(x) \le p$.

Users arrive as a Poisson process with rate $\lambda$ (jobs/sec), and job lengths are i.i.d. random variables with a general distribution with mean $l$ (bits). The bandwidth (in bits/sec) is denoted by $R$, the queue service rate (in jobs/sec) is given by $\mu = R/l$, and the offered load is denoted by $\rho \equiv \lambda/\mu = \lambda l/R$. We assume that the queue is ergodic, namely $\rho < 1$. Then the stationary distribution of the queue length is geometric, given by $\pi_n \equiv \Pr\{N = n\} = (1 - \rho)\rho^n = (1 - \lambda l/R)(\lambda l/R)^n, n \in \mathbb{Z}$[16].

*Theorem 1:* In an M/G/1-PS queue, the minimum rate violation probability for $x$ is given by $P^I(x) = \rho^{\lfloor R/x \rfloor} = (\lambda l/R)^{\lfloor R/x \rfloor}$, where $\lfloor . \rfloor$ is the floor function.
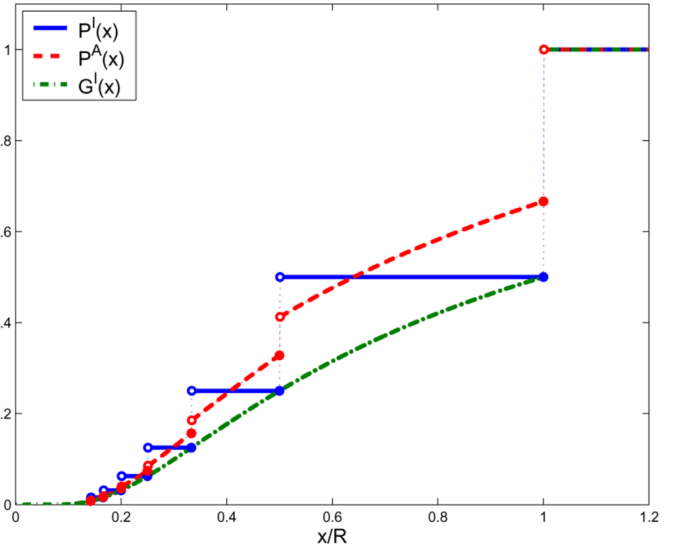


Fig. 1.   $P^I(x)$, $P^A(x)$ or $G^I(x)$ versus $x$ in an M/M/1-PS queue ($\rho = 0.5$).

*Proof:* Since the bandwidth is equally shared among all active users, the transmission rate per user, denoted by $X$, equals the bandwidth $R$ divided by the number of users $N$ for $N > 0$. It follows that $\Pr\{X = R/n\} = \pi_n/(1 - \pi_0) = (1 - \rho)\rho^{n-1}, n \in \mathbb{Z}^+$. Under the assumption of ergodicity, $P^I(x)$ can be expressed as the probability of a user receiving an transmission rate lower than $x$ during a busy period, i.e., $P^I(x) = \Pr\{X < x\} = \sum_{n=\lfloor R/x \rfloor + 1}^{\infty} \Pr\{X = R/n\} = \rho^{\lfloor R/x \rfloor}$. The theorem follows. ∎

The minimum throughput violation probability $P^A(x)$ depends on the stationary distribution of sojourn time conditioned on the job length. The unconditional distribution of sojourn time for an M/M/1-PS queue has been presented in [3] and [4]; however, the conditional distribution of sojourn time given the job length for an M/G/1-PS queue, even for an M/M/1-PS queue, is unknown. As a result, the minimum throughput violation probability is unavailable in closed form and must be obtained from simulation results.

Fig. 1 shows $P^I(x)$ and $P^A(x)$ versus $x$ (normalized by the bandwidth $R$) in an M/M/1-PS queue with load $\rho = 0.5$. For $x \in (R/n, R/(n+1)) \, \forall n \in \mathbb{Z}^+$, $P^I(x)$ remains constant at $P^I(R/n)$, whereas $P^A(x)$ is monotonically increasing with $x$ and crosses above $P^I(x)$. Both $P^I(x)$ and $P^A(x)$ are discontinuous at the points where $R$ is a multiple of $x$, and we observe that $P^A(x) > P^I(x)$ at all these points, i.e., the minimum throughput threshold is more likely to be violated than the same level of minimum rate threshold. (Given that $P^A(x)$ is monotonically increasing with $x$ in between these steps, this behavior seems natural, but we have not been able to analytically prove it.) When $x = R$, a simple calculation shows that $P^A(R) = 2\rho/(1+\rho) > P^I(R) = \rho$ for $0 < \rho < 1$.

We now consider performance guarantees on violation probability of the form $P^I(x) \le p$ and of the form $P^A(x) \le p$ with $p$ denoting the maximum allowed violation probability. We will examine the marginal cost of bandwidth with respect to $p$ and to the rate threshold $x$, and show that this cost is monotonically and convexly decreasing with $p$ in both minimum rate and minimum throughput cases, consistent with our intuition.

We start by considering the performance guarantee of the form $P^I(x) \le p$. As $P^I(x)$ is discontinuous and nondifferentiable at some points, to calculate the marginal cost, we use an approximation $G^I(x) \equiv \rho^{R/x} = (\lambda l/R)^{R/x}$, $x \le R$, of $P^I(x)$. The approximation $G^I(x)$ is continuous and differentiable except at $x = R$. For purposes
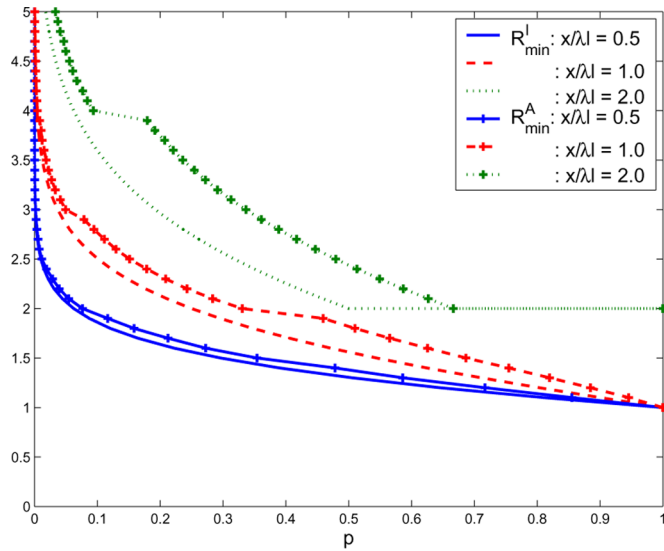
Fig. 2. $R^I_{min}$ or $R^A_{min}$ versus $p$ in an M/M/1-PS queue.



Fig. 3. $R^I_{min}$ or $R^A_{min}$ versus $x$ in an M/M/1-PS queue.

of discussion, we assume that the user arrival rate $\lambda$ and the average job length $l$ are fixed, but that the bandwidth $R$ can be chosen to satisfy the performance requirement. We expect that satisfaction of $G^I(x) \leq p$ thus requires that $R$ exceeds a related lower bound, denoted by $R^I_{min}$. The first results are expressions of the marginal bandwidth $R^I_{min}$ with respect to the maximum allowed violation probability $p$ and to the minimum required transmission rate $x$.

*Theorem 2:* In an M/G/1-PS queue

$$\frac{\partial R^I_{min}}{\partial p} = \begin{cases} \frac{\frac{x}{p}}{\ln\left(\frac{\lambda l}{R^I_{min}}\right) - 1}, & 0 < p < \min\left(\frac{\lambda l}{x}, 1\right) \\ 0, & \min\left(\frac{\lambda l}{x}, 1\right) < p < 1 \end{cases} \quad (3)$$

$$\frac{\partial R^I_{min}}{\partial x} = \begin{cases} \frac{\ln p}{\ln\left(\frac{\lambda l}{R^I_{min}}\right) - 1}, & \frac{x < \lambda l}{p} \\ 1, & \frac{x > \lambda l}{p}. \end{cases} \quad (4)$$

*Proof:* Taking the derivative of $G^I(x) = (\lambda l/R)^{R/x}$ with respect to $R$ and simplifying yields $\partial G^I(x)/\partial R = (\ln \rho - 1)\rho^{R/x}/x < 0$, $R \geq x, R > \lambda l$. Hence $G^I(x)$ decreases monotonically with $R$ for a fixed $x \leq R$. To satisfy $G^I(x) \leq p$, it is required that $R \geq R^I_{min}$. The minimum bandwidth $R^I_{min}$ is determined by the fixed point equation $p = (\lambda l/R^I_{min})^{R^I_{min}/x}$, $R^I_{min} \geq x$. When taking the derivative of $R^I_{min}$ with respect to $p$ or $x$, two cases apply. When $R^I_{min} = x$, the fixed point equation simplifies to $p = \lambda l/R^I_{min} = \lambda l/x$, and hence $\partial R^I_{min}/\partial p = 0$ and $\partial R^I_{min}/\partial x = 1$. However when $R^I_{min} < x$, taking first partial derivatives of the fixed point equation and solving for $\partial R^I_{min}/\partial p$ and $\partial R^I_{min}/\partial x$ gives the first case in (3)–(4). ∎

The expressions given in (3)–(4) can be used to guide dimensioning algorithms, as they relate the minimum required bandwidth $R^I_{min}$ to the traffic intensity $\lambda l$ and the violation probability requirement $p$ for a given threshold $x$. In addition, it is common in the research literature to suggest pricing of service based on the marginal cost of the required resources. In this context users would be charged a joint function of $x$ and $p$. If price is based on the marginal cost of resources required to provide the performance guarantee, then the price would be based on (3)–(4), the cost per unit $R$, and perhaps a fixed cost to ensure that revenues covers average costs.

$R^I_{min}$ versus $p$ is shown in Fig. 2, where each curve corresponds to a constant value of $x/\lambda l$. The required bandwidth $R^I_{min}$ decreases monotonically and convexly over $0 < p < \min(\lambda l/x, 1)$, and then
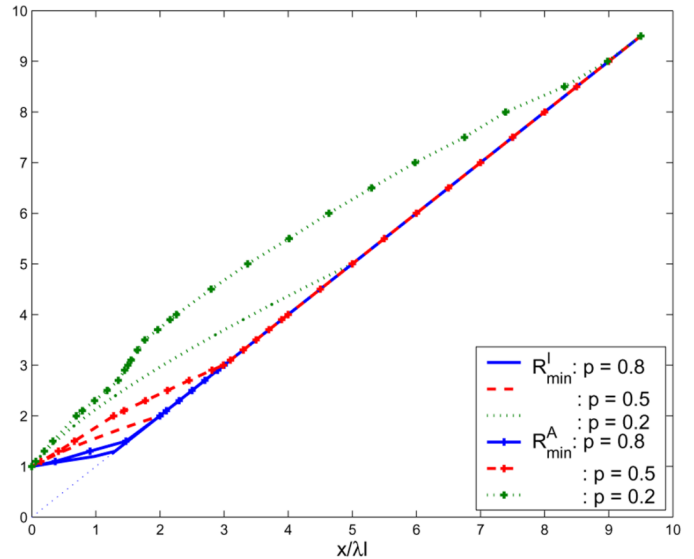
remains constant at $x$ for higher values of $p$. The convexity comes from the fact that the second partial of $R^I_{min}$ with respect to $p$ is positive. We use the term *p-limited* to denote the case in which $R^I_{min}$ decreases monotonically with $p$, and the term *x-limited* to denote the case in which $R^I_{min} = x$. In the $x$-limited case, the performance bound simplifies to a guarantee that the user obtains the full bandwidth with probability $p$ or higher. Given that the user is in the system, there are no other users with probability $\lambda l/R^I_{min}$, and hence the $x$-limited case only occurs when $p \geq \lambda l/x$.

$R^I_{min}$ versus $x$ is shown in Fig. 3, where each curve represents a constant value of $p$. The required bandwidth $R^I_{min}$ is concavely and then linearly increasing as $x$ increases. The concavely increasing portion, i.e., $x < \lambda l/p$, corresponds to the $p$-limited case discussed above; the linearly increasing portion, i.e., $x \geq \lambda l/p$, corresponds to the $x$-limited case, and $R^I_{min} = x$.

We now proceed to the performance bound on the minimum throughput violation probability. We expect that satisfaction of $P^A(x) \leq p$ requires that $R$ exceeds a related lower bound, denoted by $R^A_{min}$. The required bandwidth $R^A_{min}$ versus $p$, obtained via simulation, is plotted (curves with markers) in Fig. 2. The characterization of $R^A_{min}$ is similar to that of $R^I_{min}$, except that $R^A_{min}$ is discontinuous when it is a multiple of $x$. These discontinuous portions correspond to the jumps in $P^A(x)$ in Fig. 1. As before, there are two cases: $x$-limited and $p$-limited. In the $x$-limited case $R^A_{min} = x$, each user requires the full bandwidth throughout its transmission. A user fails to achieve that if there exists another user in the system during its transmission. It can be derived that this probability is $\frac{2\lambda l/x}{1+\lambda l/x}$. In the $p$-limited case, we observe that $R^A_{min}$ decreases with $p$ over $0 < p < \min(\frac{2\lambda l/x}{1+\lambda l/x}, 1)$.

$R^A_{min}$ versus $x$ for different values of $p$ is also plotted in Fig. 3. The required bandwidth $R^A_{min}$ is monotonically increasing with $x$. When $x < \lambda l(2-p)/p$, the system is $p$-limited and $R^A_{min}$ is increasing with $p$; when $x \geq \lambda l(2-p)/p$, the system becomes $x$-limited, and $R^A_{min} = x$.

## IV. VIOLATION PROBABILITY IN AN M/G/1/K-PS QUEUE

In the M/G/1-PS queue, the system provides a higher level of performance guarantees by increasing the bandwidth. However, in many practical situations, the available bandwidth is physically limited and can not be increased, or increasing bandwidth may not be the most efficient manner to satisfy the performance requirements. In this case,

connection access control (CAC) is commonly used to maintain acceptable performance for admitted jobs, at the cost of blocking some jobs. This approach can be modeled by an M/G/1/K-PS queue, where $K$ denotes the queue limit, i.e., the maximum number of users that can transmit simultaneously. CAC gives the network designer additional flexibility, by allowing for a tradeoff between the bandwidth and the blocking probability.

In this section, our goal is to explore the relationship between CAC and performance requirements on the violation probability. After presenting a closed-form expression for the minimum rate violation probability, we numerically investigate the relationship between the bandwidth, the queue limit, the violation probability and the blocking probability. In particular, we demonstrate the nature of binding constraints on the violation probability and on blocking probability.

The stationary distribution of queue length for an M/G/1/K-PS queue is given by $\pi_n = \Pr\{N = n\} = (1-\rho)\rho^n/(1-\rho^{K+1})$, $n = 0, 1, \ldots, K$, where $\rho = \lambda l/R$[16].

*Theorem 3:* In an M/G/1/K-PS queue, the minimum rate violation probability is given by

$$P^I(x, K) = \begin{cases} \frac{\rho^{\lfloor R/x \rfloor} - \rho^K}{1-\rho^K}, & x > R/K; \\ 0, & \text{otherwise.} \end{cases} \tag{5}$$

where the load $\rho = \lambda l/R$.

*Proof:* With the ergodicity assumption, $P^I(x)$ can be expressed as the probability of a user receiving a transmission rate below $x$ during a busy period, i.e., $P^I(x, K) = \Pr\{X < x\} = \sum_{n=\lfloor R/x \rfloor+1}^{K} \pi_n/(1-\pi_0)$. The theorem follows by substituting $\pi_n$ by the expression above. ∎

*Corollary 1:* For $R/K < x < R$, the minimum rate violation probability $P^I(x, K)$ increases with $K$, with an upper bound $\rho^{\lfloor R/x \rfloor}$, and decreases with $R$.

*Proof:* From (5), $P^I(x, K) = 1 - \frac{1-\rho^{\lfloor R/x \rfloor}}{1-\rho^K}$, $x > R/K$. As $K$ increases, $P^I(x, K)$ also increases; as $K$ approaches infinity, the queue reduces to an M/G/1-PS queue, and $P^I(x, K)$ reaches its maximum, $\rho^{\lfloor R/x \rfloor}$. When $R$ increases by $\Delta R$, it follows that $P^I(x, K; R + \Delta R) - P^I(x, K; R) < \frac{1-\rho^{\lfloor R/x \rfloor}}{1-\rho^K} - \frac{1-(\rho-\Delta\rho)^{\lfloor R/x \rfloor}}{1-(\rho-\Delta\rho)^K} < 0$, where $\Delta\rho = \frac{\lambda l \Delta R}{R(R+\Delta R)}$. The last inequality holds for $K > \lfloor R/x \rfloor$. The corollary follows. ∎

The blocking probability is given by $P^B = \pi_K = (1-\rho)\rho^K/(1-\rho^{K+1})$ where $\rho = \lambda l/R$. It is readily shown that $P^B$ decreases as $R$ increases or as $K$ increases.

We now consider a performance bound on the minimum rate violation probability in the form $P^I(x, K) \leq p_v$. The bandwidth is assumed fixed and the queue size is determined by the performance requirement. From (5), the maximum number of users $K_{max}^I$ is given by

$$K_{\max}^I = \begin{cases} \lfloor \frac{R}{x} \rfloor, & p_v < a; \\ \lfloor \frac{1}{\ln \rho} \ln \frac{p_v - \rho^{\lfloor R/x \rfloor}}{p_v - 1} \rfloor, & a \leq p_v < \rho^{\lfloor R/x \rfloor}; \\ +\infty, & \text{otherwise} \end{cases} \tag{6}$$

where $a = P^I(x, \lfloor R/x \rfloor + 1)$ is the minimum rate violation probability at $K = \lfloor R/x \rfloor + 1$.

The first case in (6) comes from the second case in (5). When $K \leq R/x$, $P^I(x, K) = 0$. We can view this as $p_v = 0$ implies that $K_{max}^I = \lfloor R/x \rfloor$. As $p_v$ increases, we remain in this case until $p_v = a$, where the first case in (5) and the second case in (6) apply. As $p_v$ increases, $K_{max}^I$ is monotonically and convexly increasing. When $p_v > \rho^{\lfloor R/x \rfloor}$, which is the upper bound of $P^I(x, K)$, $K$ can be any positive integer.

We next consider the joint selection of the bandwidth $R$ and the queue limit $K$ to satisfy performance requirements on the violation probability $P^I(x, K) \leq p_v$ and on the blocking probability $P^B \leq p_b$. We are interested in the feasible region of such pairs $(R, K)$.
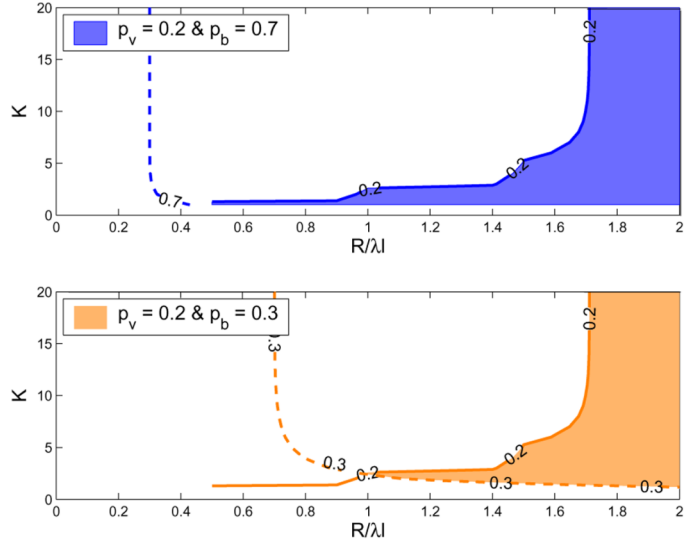


Fig. 4. Two cases of feasible region for $R$ and $K$ ($x/\lambda l = 0.5$).

*Theorem 4:* When $x \geq (1-p_b)\lambda l/p_b$, the feasible region of $(R, K)$ depends only on $p_v$ and can be determined by (6); otherwise, the feasible region of $(R, K)$ relies on both $p_v$ and $p_b$, which is the intersection of the set determined by (6) and the set of $\{(R, K) : K \geq b\}$ where $b = \ln[p_b/(1-\rho+\rho p_b)]/\ln \rho$.

*Proof:* Using the expression for $P^B$ above, it is readily shown that satisfaction of $P^B \leq p_b$ iff $K \geq b$. When $\rho < \frac{p_b}{1-p_b}$, $b < 1$; hence $K \geq b$ which implies that $P^B \leq p_b$. Since $R \geq x$, i.e., $\rho \leq \frac{\lambda l}{x}$, we have that if $\frac{\lambda l}{x} < \frac{p_b}{1-p_b}$, the blocking probability requirement can be always satisfied for $K \geq 1$ and $R \geq x$. Therefore, in this case the feasible set can be determined only by the violation probability requirement $p_v$. ∎

There exist two possible cases for the feasible region of $(R, K)$ given a pair of performance requirements $(p_v, p_b)$, as shaded areas shown in Fig. 4. The first case occurs when $x > (1-p_b)\lambda l/p_b$. In this case, only the constraint on $P^I(x, K)$ is binding; as a result, the minimum cost solutions lie on the left boundary of the shaded region, on the curve for which $P^I(x, K) = p_v$. The minimum bandwidth occurs when $K = 1$. In the second case, both constraints are binding; as a result, the minimum cost solutions lie on the left boundaries of the shaded region, on either or both curves. The minimum bandwidth occurs at the intersection of $P^I(x, K) = p_v$ and $P^B = p_b$. This minimum bandwidth increases along the upper boundary when a lower blocking probability is required and increases along the lower boundary when a lower violation probability is required.

## V. VIOLATION PROBABILITY IN AN M/M/1-DPS QUEUE

In this section, we consider the system's ability to serve multiple classes that may require different violation probabilities. This can be modeled by a discriminatory processor sharing (DPS) queue, which is a multi-class generalization of the PS queue. Under the DPS discipline, each class is associated with a weight, and the transmission rate of a user present in the system is controlled by this vector of weights. By varying the weights, the achievable performance can be varied over a wide range. Our goal here is to explore the relationship between the bandwidth, class weights and violation probabilities, and to explore how to set system parameters to meet performance requirements. We discuss these issues in a two-class system.

In a two-class M/M/1 DPS queue, the arrivals of each class independently constitute a Poisson process with rate $\lambda_1$ or $\lambda_2$, and job lengths are assumed to be i.i.d. exponentially distributed with mean
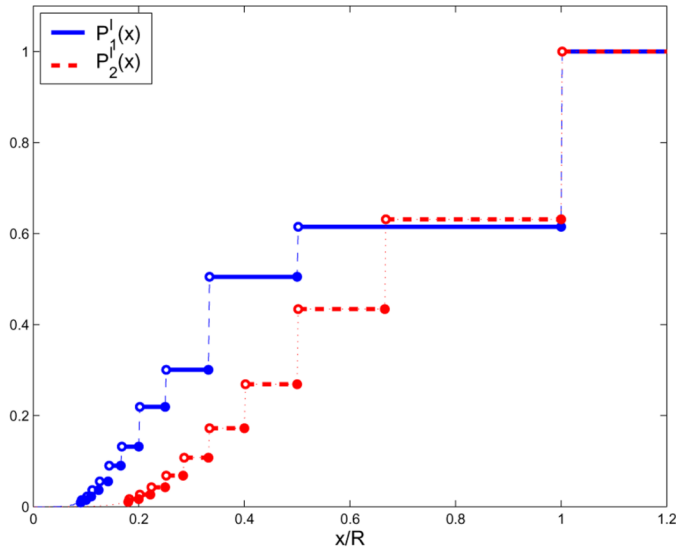
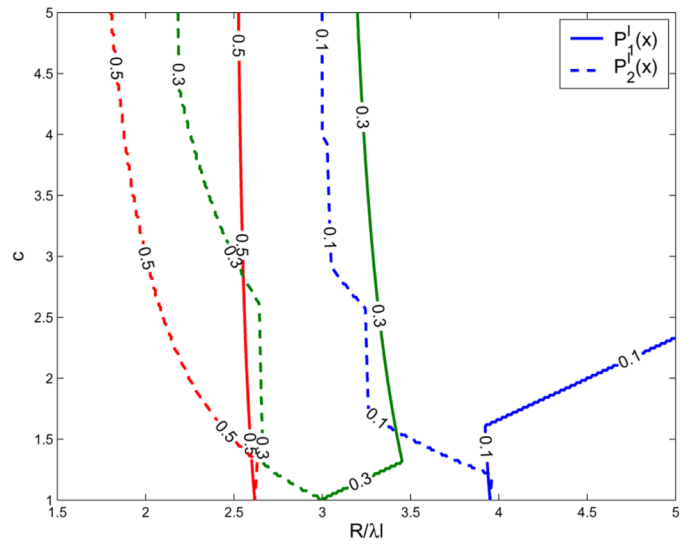Fig. 5. $P_1^I(x)$ or $P_2^I(x)$ versus $x$ in an M/M/1-DPS queue ($\rho = 0.5, c = 2$).



Fig. 6. $c$ versus $R$ for fixed $p_1$ or $p_2$ in an M/M/1-DPS queue ($\rho = 0.5, x = 1.5\lambda l$).

$l_1$ or $l_2$. The transmission rate of a user depends on the user numbers of both classes and on their assigned weights. Specifically, denote the user numbers by $(n_1, n_2)$ and the weights by $(w_1, w_2)$. The transmission rate of a class-$j$ user is given by $w_j R/(w_1 n_1 + w_2 n_2), j = 1, 2$. The queue is ergodic when the system is nonsaturated, i.e., when the total load $\rho = (\lambda_1 l_1 + \lambda_2 l_2)/R < 1$.

The minimum rate violation probability of class $j$ is given by

$$P_j^I(x) = \frac{\sum\sum_{\{(n_1,n_2)|n_j\neq 0\}} \pi_{n_1 n_2} I_{\{w_j R/(w_1 n_1 + w_2 n_2) < x\}}}{\sum\sum_{\{(n_1,n_2)|n_j\neq 0\}} \pi_{n_1 n_2}}, j = 1, 2.$$

Unfortunately, $(n_1(t), n_2(t))$ is an irreversible Markov chain and its joint stationary distribution is unavailable in closed form. Our following discussion is based on numerical computation.

Define $c \equiv w_2/w_1$. Without loss of generality, we assume that $c > 1$, i.e., a class-2 user can transmit at a higher rate than a class-1 user when both are present in the system. To focus on the effect of $c$, let $\lambda_1 = \lambda_2$ and $l_1 = l_2$. Fig. 5 plots $P_1^I(x)$ and $P_2^I(x)$ versus the rate threshold $x$ when $\rho = 0.5$ and $c = 2$. Since class 2 is given preferential treatment, intuition leads one to expect that $P_1^I(x) > P_2^I(x) \ \forall \ x$. However, our numerical computation results show that while $P_1^I(x) > P_2^I(x)$ for $x \leq \frac{c}{c+1}R$, the order is reversed for $x > \frac{c}{c+1}R$. To understand this, we investigate the effect of the increase in $c$ on the stationary distribution of queue length. For $x > \frac{c}{c+1}R$, the corresponding minimum rate violation probabilities are given by

$$P_1^I(x) = 1 - \frac{\pi_{10}}{\sum_{m=1}^{\infty}\sum_{n=0}^{\infty}\pi_{mn}}, \ P_2^I(x) = 1 - \frac{\pi_{01}}{\sum_{m=0}^{\infty}\sum_{n=1}^{\infty}\pi_{mn}}.$$

When $c = 1$, the system reduces to a two-class PS queue, and $P_1^I(x) = P_2^I(x) \ \forall \ x$. As $c$ increases, numerical results show that $P_1^I(x) > P_2^I(x)$ for $x > \frac{c}{c+1}R$ and $c > 1$. In contrast, if the minimum throughput violation probabilities are considered, $P_1^A(x) > P_2^A(x) \ \forall x \ \forall c$, since the conditional sojourn times are stochastically ordered according to the DPS weights [17].

We now consider the joint selection of the bandwidth $R$ and the weight ratio $c$ to satisfy the requirements on the minimum rate violation probability for both classes. Fig. 6 plots combinations of $R$ and $c$ that attain constant $P_1^I(x)$ and $P_2^I(x)$, given a fixed $x$, as solid and dashed curves, respectively. The minimum rate violation probability is denoted on each curve. Some curves are monotonic, e.g., $R/\lambda l$ is monotonically decreasing with $c$ at $P_1^I(x) = 0.5$. However, some curves display more
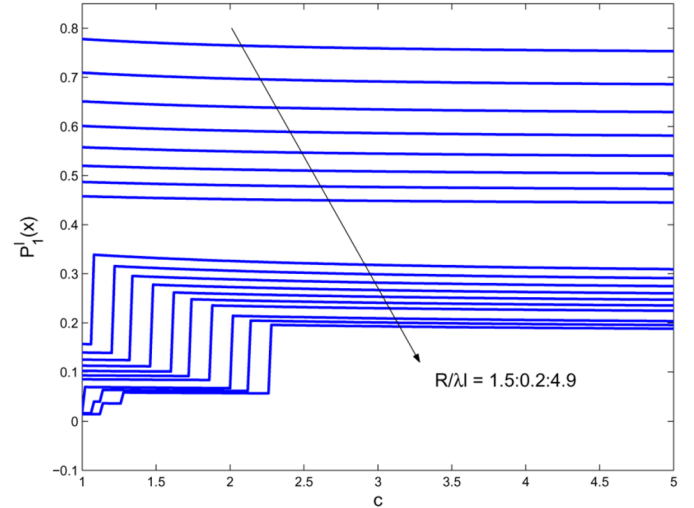


Fig. 7. $P_1^I(x)$ versus $c$ in an M/M/1-DPS queue ($\rho = 0.5, x = 1.5\lambda l$).

complex behavior, e.g., at $P_1^I(x) = 0.3 \ R/\lambda l$ initially increases with $c$ and then decreases with $c$. This complex behavior is caused by the variation of the minimum rate violation probability with $c$. Fig. 7 displays $P_1^I(x)$ versus $c$ at various values of $R/\lambda l$. Each curve consists of positive jumps and intervals of monotonically decreasing $P_1^I(x)$. To illustrate this behavior, consider a set $\{(n_1, n_2)|\frac{R}{n_1+n_2 c} < x\}$. If the set does not change membership, $P_1^I(x)$ decreases monotonically (and $P_2^I(x)$ increases monotonically) with $c$. However, when the set increases in size as $c$ increases, this causes a positive jump in $P_1^I(x)$ (and a negative jump in $P_2^I(x)$). The pairs of $(R, c)$ that satisfy $P_1^I(x) \leq p_1, P_2^I(x) \leq p_2$ are the intersection of the areas to the right of the pairs of corresponding contours.

## VI. CONCLUSION

We have focused on the threshold violation probability of the transmission rate or the average rate as a performance metric in processor-sharing or discriminatory processor-sharing queues. We introduced definitions of threshold violation probability as observed by users or by the queue. In an M/G/1-PS queue, we found the minimum bandwidth and marginal bandwidth needed for a given performance

guarantee on $P^I(x)$, and showed that under a constraint on the minimum rate violation probability, the required system bandwidth might be limited by either the probability $p$ or by the threshold of the rate requirement $x$, and gave conditions explaining when each case occurs. In an M/G/1/K-PS queue, we discussed the relationship between total transmission rate, queue limit, threshold violation probability, and blocking probability. We finally considered a two-class DPS, and discussed what combinations of class weighting and bandwidth can be used to achieve specified threshold violation probabilities of transmission rate for both classes.

We believe such results are useful in dimensioning processor-sharing queues when performance is measured by tail probabilities. In particular, we expect such results can be used within networking to design scheduling and connection access control policies for data services.

## REFERENCES

[1] J. W. Roberts, "A survey on statistical bandwidth sharing," *Computer Networks*, vol. 45, no. 3, pp. 319–332, 2004.
[2] J. E. G. Coffman, R. R. Muntz, and H. Trotter, "Waiting time distributions for processor-sharing systems," *J. ACM*, vol. 17, no. 1, pp. 123–130, 1970.
[3] J. A. Morrison, "Response-time distribution for a processor-sharing system," *SIAM J. Appl. Math.*, vol. 45, no. 1, pp. 152–167, 1985.
[4] F. Guillemin and J. Boyer, "Analysis of the M/M/1 queue with processor sharing via spectral theory," *Queueing Syst.*, vol. 39, pp. 377–397, 2001.
[5] J. A. Morrison, "Aysmptotic analysis of the waiting-time distribution for a large closed processor-sharing system," *SIAM J. Appl. Math.*, vol. 46, pp. 140–170, 1986.
[6] J. A. Morrison, "Conditioned response-time distribution for a large closed processor-sharing system in very heavy usage," *SIAM J. Appl. Math.*, vol. 47, pp. 1117–1129, 1987.
[7] C. Knessl, "On the sojourn time distribution in a finite capacity processor shared queue," *J. ACM*, vol. 40, no. 5, pp. 1238–1301, 1993.
[8] G. Fayolle, I. Mitrani, and R. Iasnogorodski, "Sharing a processor among many job classes," *J. ACM*, vol. 27, pp. 519–532, 1980.
[9] K. M. Rege and B. Sengupta, "A decomposition theorem and related results for discriminatory processor sharing queue," *Queueing Syst.*, vol. 18, pp. 333–351, 1994.
[10] K. M. Rege and B. Sengupta, "Queue-length distribution for the discriminatory processor-sharing queue," *Operat. Res.*, vol. 44, pp. 653–657, 1996.
[11] S. Borst, "User-level performance of channel-aware scheduling algorithms in wireless data networks," *IEEE/ACM Trans. Networking*, vol. 13, no. 3, pp. 636–647, Jun. 2005.
[12] A. A. Kherani and A. Kumar, "Stochastic models for throughput analysis of randomly arriving elastic flows in the Internet," in *Proc. INFOCOM'02*, 2002, vol. 2, pp. 1014–1023.
[13] N. Chen and S. Jordan, "Throughput in processor-sharing queues," *IEEE Trans. Automat. Control*, vol. 52, no. , pp. 299–305, 2007.
[14] M. Harchol-Balter, K. Sigman, and A. Wierman, "Asymptotic convergence of scheduling policies with respect to slowdown," *Proc. of IFIP Perform.*, pp. 241–256, 2002.
[15] M. Harchol-Balter and A. B. Downey, "Exploiting process lifetime distributions for dynamic load balancing," *ACM Trans. Comput. Syst.*, vol. 15, no. 3, pp. 253–285, 1997.
[16] L. Kleinrock, *Queueing Systems*. New York: Wiley, 1975.
[17] K. E. Avrachenkov, U. Ayesta, P. Brown, and R. Nunez-Queija, "Discriminatory processor sharing revisited," in *Proc. INFOCOM'05*, 2005, vol. 2, pp. 784–795.

# Improvement on Stability Analysis for Linear Systems Under State Saturation

Xiaofu Ji, Yukun Sun, and Tailiu Liu

*Abstract*—This note considers the problem of stability analysis for linear systems under state saturation. With the introduction of set coverage that gives less constraint on the free matrix $G$, a less conservative sufficient global asymptotic stability condition is obtained and the corresponding iterative linear matrix inequality algorithm is given. A numerical example is given to show the effectiveness of the proposed method.

*Index Terms*—Iterative linear matrix inequality, linear systems, state saturation.

## I. INTRODUCTION AND PROBLEM STATEMENT

In this note, we consider the following linear system under state saturation:

$$\dot{x} = h(Ax) \tag{1}$$

where $x \in \mathbb{D}^n := \{x = [x_1 \ x_2 \ \cdots \ x_n]^T \in \mathbb{R}^n : -1 \leq x_i \leq 1, i \in [1, n]\}$ is the state vector, $A = [a_{ij}] \in \mathbb{R}^{n \times n}$ is a constant matrix, and $h(\cdot)$ is a saturation function defined as

$$h(Ax) = \begin{bmatrix} h_1 \left( \sum_{j=1}^n a_{1j} x_j \right) \\ h_2 \left( \sum_{j=1}^n a_{2j} x_j \right) \\ \vdots \\ h_n \left( \sum_{j=1}^n a_{nj} x_j \right) \end{bmatrix} \tag{2}$$

with, for each $i \in [1, n]$

$$h_i \left( \sum_{j=1}^n a_{ij} x_j \right) = \begin{cases} 0, & \text{if } |x_j| = 1 \\ & \text{and } \left( \sum_{j=1}^n a_{ij} x_j \right) x_i > 0 \\ \sum_{j=1}^n a_{ij} x_j, & \text{otherwise.} \end{cases} \tag{3}$$

Global asymptotic stability of this system has been studied recently [1]–[3], and discrete counterpart [4], [5]. In [3], Fang *et al.* introduced a diagonally dominant matrix $G$ with negative diagonal element to confine the system state under saturation to a convex polyhedron, based on which, a less conservative stability condition was obtained. In fact, $G$