

**Performance analysis of the system contents in a discrete-time non-preemptive priority queue with general service times**

Joris Walraevens, Bart Steyaert and Herwig Bruneel

SMACS Research Group  
Ghent University, Vakgroep TELIN (TW07V)  
Sint-Pietersnieuwstraat 41, B-9000 Gent, Belgium.  
Phone: 0032-9-2648902  
Fax: 0032-9-2644295  
E-mail: jw@telin.rug.ac.be

**Abstract**

We consider a discrete-time queueing system with non-preemptive priority scheduling. Two classes of traffic will be considered, i.e., high priority and low priority traffic, which both generate variable-length packets. We will derive an expression for the joint Probability Generating Function (pgf) of the steady-state system contents of the high and the low priority traffic. From these, some performance measures (such as the mean value of steady-state system contents and packet delay of high and low priority packets) will be derived. These will be used to illustrate the significance of priority scheduling. Our results can be used to analyse performance of buffers in voice/data networks.

**Keywords** discrete-time queueing, priority scheduling, generating functions

## 1 Introduction

In recent years, there has been much interest devoted to incorporating multimedia applications in IP networks. Different types of traffic need different QoS standards. For real-time applications, it is important that mean delay and delay-jitter are bounded, while for non real-time applications, the Loss Ratio (LR) is the restrictive quantity.

In general, one can distinguish two priority strategies, which will be referred to as Time Priority and Space Priority. Time priority schemes attempt to guarantee acceptable delay boundaries to delay-sensitive traffic (such as voice/video). This is achieved by giving it non-preemptive priority over non-delay-sensitive traffic, and/or by sharing access to the server among the various traffic classes in such a way so that each can meet its own specific delay requirements. Several types of Time priority (or scheduling) schemes (such as Weighted-Round-Robin (WRR), Weighted-Fair-Queueing(WFQ)) have been proposed and analyzed, each with their own specific algorithmic and computational complexity (see e.g. [6] and the references therein). On the other hand, Space Priority schemes attempt to minimize the packet loss of loss-sensitive traffic (such as data). Again, various types of Space Priority (or discarding) strategies (such as Push-Out Buffer (POB), Partial Buffer Sharing (PBS)) have been presented in the literature (see e.g. [15]), mainly in the context of ATM buffers. An overview of both types of priority schemes can be found in [1].

In this paper, we will focus on the effect of non-preemptive time priority scheduling. We assume that time-sensitive traffic has non-preemptive priority over time-insensitive traffic, i.e., when the server becomes idle, a packet of time-sensitive traffic, when available, will always be scheduled next, but newly arriving time-sensitive traffic can not interrupt transmission of a time-insensitive packet that has already commenced.

In literature, there have been a number of contributions with respect to non-preemptive priority scheduling. An overview of some basic non-preemptive priority queueing models can be found in Jaiswal [3], Takacs [10] and Takagi [11] and the references therein. Khamisy et al. [4], Laevens et al. [5], Takine et al. [13] and Walraevens et al. [16] have studied discrete-time non-preemptive priority queues with deterministic service times equal to one slot. Khamisy [4] analyzes the system contents for the different classes, for a queue fed by a two-state Markov modulated arrival process. Laevens [5] analyses the system contents and cell delay in the case of a multiserver queue. In Takine [13], the system contents and the delay for Markov modulated high priority arrivals and geometrically distributed low priority arrivals are presented. Walraevens [16] studies the system contents and cell delay, in the special case of an output queueing switch with Bernoulli arrivals. Furthermore, non-preemptive priority queues have been considered by Rubin et al. [7], Stanford [8], Sugahara et al. [9] and Takine et al. [12, 14]. Rubin [7] studies the mean waiting time, for a queue fed by an i.i.d. arrival process. Stanford [8] analyses the interdeparture time distribution in a queue fed by a Poisson process. In Sugahara [9], a non-preemptive queue in continuous time is presented, with a Switched Poisson Process arrival process for the high priority packets. Finally, Takine [12, 14] studies a discrete-time MAP/G/1 queue, using matrix-analytic techniques.

In this paper, we analyse the system contents of high and low priority traffic in a discrete-time single-server buffer for a non-preemptive priority scheme and a per-slot i.i.d. number of arrivals. The transmission times of the packets generated by both types are assumed to be generally distributed. We will demonstrate that an analysis based on generating functions is extremely suitable for modelling this type of buffers with priority scheduling. From these generating functions, we can then easily calculate expressions for some interesting performance measures, such as the mean value of system contents and packet delay of both traffic types. These closed-form expressions require virtually no computational effort at all, and are well-suited for evaluating the impact of the various system parameters on the overall performance. This makes it possible to study the effect of priority scheduling and the impact of the *non-preemptive* priority scheduling on the high priority traffic.

The remainder of this paper is structured as follows. In the following section, we present the mathematical model. Before analysing the system contents at arbitrary slot boundaries, we first analyse the system contents at the beginning of special slots in section 3. In section 4 we will then analyse the steady-state system contents at arbitrary slot boundaries. In section 5, we calculate the

moments of the system contents, while we give some numerical examples in section 6. Finally, some conclusions are formulated in section 7.

## 2 Mathematical model

We consider a discrete-time single-server queueing system with infinite buffer space. Time is assumed to be slotted. There are 2 types of traffic arriving in the system, namely packets of class 1 and packets of class 2. We denote the number of arrivals of class  $j$  during slot  $k$  by  $a_{j,k}$  ( $j = 1, 2$ ). The bivariate random variables  $(a_{1,k}, a_{2,k})$  are assumed to be i.i.d. and are characterized by the joint probability mass function

$$a(m, n) \triangleq \text{Prob}[a_{1,k} = m, a_{2,k} = n],$$

and joint probability generating function (pgf)  $A(z_1, z_2)$ ,

$$A(z_1, z_2) \triangleq E[z_1^{a_{1,k}} z_2^{a_{2,k}}] = \sum_{m,n=0}^{\infty} a(m, n) z_1^m z_2^n.$$

Notice that the number of packet arrivals from different classes (within a slot) can be correlated. We define the pgf of the total number of arrivals during a slot by  $A_T(z) \triangleq E[z^{a_{1,k} + a_{2,k}}] = A(z, z)$ . Further, we define the marginal pgf's of the arrivals from class 1 and class 2 during a slot by  $A_1(z) \triangleq E[z^{a_{1,k}}] = A(z, 1)$  and  $A_2(z) \triangleq E[z^{a_{2,k}}] = A(1, z)$  respectively. We furthermore denote the arrival rate of class  $j$  ( $j = 1, 2$ ) by  $\lambda_j = A_j'(1)$  and the total arrival rate by  $\lambda_T \triangleq \lambda_1 + \lambda_2$ .

The service times of the class  $j$  packets are assumed to be i.i.d. and are characterized by the probability mass function

$$s_j(m) \triangleq \text{Prob}[\text{service of a class } j \text{ packet takes } m \text{ slots}], \quad m \geq 1,$$

and probability generating function  $S_j(z)$ ,

$$S_j(z) = \sum_{m=1}^{\infty} s_j(m) z^m,$$

with  $j = 1, 2$ . We furthermore denote the mean service time of a class  $j$  packet by  $\mu_j = S_j'(1)$ .

The system has one server that provides the transmission of packets. Class 1 packets are assumed to have non-preemptive priority over class 2 packets, and within one class the service discipline is FCFS. Due to the priority scheduling mechanism, it is as if class 1 packets (the high priority packets) are stored in front of class 2 packets (the low priority packets) in the queue. So, if there are any class 1 packets in the queue when the server becomes idle, the one with the longest elapsed waiting time will be served next. If, on the other hand, no class 1 packets are present in the queue at that moment, the class 2 packet with the longest elapsed waiting time, if any, will be served next. Since the priority scheduling is non-preemptive, service of a packet will not be interrupted by newly arriving packets.

Finally, we define the load offered by class  $j$  packets as  $\rho_j \triangleq \lambda_j \mu_j$  ( $j = 1, 2$ ). The total load is then given by  $\rho \triangleq \rho_1 + \rho_2$ .

## 3 System contents at the beginning of start slots

To be able to analyze the system contents at the beginning of arbitrary slots, we will first analyze the system contents at the beginning of so-called start slots, i.e., slots at the beginning of which a service of a packet (if available) can start. Note that every slot during which the system is empty, is also a start slot. We denote the system contents of class  $j$  packets at the beginning of the  $l$ -th start slot by  $n_{j,l}$  ( $j = 1, 2$ ). Their joint pgf is denoted by  $N_l(z_1, z_2)$ , i.e.,

$$N_l(z_1, z_2) \triangleq E[z_1^{n_{1,l}} z_2^{n_{2,l}}].$$

Clearly, the set  $\{(n_{1,l}, n_{2,l})\}$  forms a Markov chain, since the number of arrivals of both classes are i.i.d. from slot-to-slot and only random variables during start slots are involved. If  $s_l^*$  indicates the service time of the packet that enters service at the beginning of start slot  $l$  (which is - by definition - regular slot  $k$ ) the following system equations can be established:

1. If  $n_{1,l} = n_{2,l} = 0$ :

$$\begin{aligned} n_{1,l+1} &= a_{1,k}; \\ n_{2,l+1} &= a_{2,k}, \end{aligned}$$

i.e., the only packets present in the system at the beginning of start slot  $l+1$  are the packets that arrived during the previous slot, i.e., start slot  $l$ . If the buffer is empty at the beginning of slot  $l$ , we set  $s_l^* = 0$ .

2. If  $n_{1,l} = 0$  and  $n_{2,l} > 0$ :

$$\begin{aligned} n_{1,l+1} &= \sum_{i=0}^{s_l^*-1} a_{1,k+i}; \\ n_{2,l+1} &= n_{2,l} + \sum_{i=0}^{s_l^*-1} a_{2,k+i} - 1, \end{aligned}$$

i.e., the class 2 packet in service leaves the system just before start slot  $l+1$ . In this case,  $s_l^*$  is characterized by probability mass function  $s_2(m)$ , since a class 2 packet enters the server at the beginning of start slot  $l$ .

3. If  $n_{1,l} > 0$ :

$$\begin{aligned} n_{1,l+1} &= n_{1,l} + \sum_{i=0}^{s_l^*-1} a_{1,k+i} - 1; \\ n_{2,l+1} &= n_{2,l} + \sum_{i=0}^{s_l^*-1} a_{2,k+i}, \end{aligned}$$

i.e., the class 1 packet in service leaves the system just before start slot  $l+1$ . For  $n_{1,l} > 0$ ,  $s_l^*$  is characterized by probability mass function  $s_1(m)$ , since a class 1 packet enters the server at the beginning of start slot  $l$ .

Using these system equations, we can derive a relation between  $N_l(z_1, z_2)$  and  $N_{l+1}(z_1, z_2)$ . In the remainder, we define  $E[X\{Y\}]$  as  $E[X|Y]\text{Prob}[Y]$ . We proceed as follows, taking into account the statistical independence of the random variables  $s_l^*$ ,  $(n_{1,l}, n_{2,l})$  and  $(a_{1,k+i}, a_{2,k+i}), i \geq 0$ :

$$\begin{aligned} N_{l+1}(z_1, z_2) &\triangleq E[z_1^{n_{1,l+1}} z_2^{n_{2,l+1}}] \\ &= E[z_1^{a_{1,k}} z_2^{a_{2,k}} \{n_{1,l} = n_{2,l} = 0\}] \\ &\quad + E \left[ \frac{z_1^{\sum_{i=0}^{s_l^*-1} a_{1,k+i}} z_2^{n_{2,l} + \sum_{i=0}^{s_l^*-1} a_{2,k+i} - 1}}{z_2} \{n_{1,l} = 0, n_{2,l} > 0\} \right] \\ &\quad + E \left[ \frac{z_1^{n_{1,l} + \sum_{i=0}^{s_l^*-1} a_{1,k+i} - 1} z_2^{n_{2,l} + \sum_{i=0}^{s_l^*-1} a_{2,k+i}}}{z_2} \{n_{1,l} > 0\} \right] \\ &= A(z_1, z_2) \text{Prob}[n_{1,l} = n_{2,l} = 0] + \frac{S_2(A(z_1, z_2))}{z_2} E[z_2^{n_{2,l}} \{n_{1,l} = 0, n_{2,l} > 0\}] \end{aligned}$$

$$\begin{aligned}
& + \frac{S_1(A(z_1, z_2))}{z_1} \mathbb{E} [z_1^{n_{1,l}} z_2^{n_{2,l}} \{n_{1,l} > 0\}] \\
= & A(z_1, z_2) N_l(0, 0) + \frac{S_2(A(z_1, z_2))}{z_2} [N_l(0, z_2) - N_l(0, 0)] \\
& + \frac{S_1(A(z_1, z_2))}{z_1} [N_l(z_1, z_2) - N_l(0, z_2)].
\end{aligned} \tag{1}$$

We assume that the system is stable (implying that the equilibrium condition requires that  $\rho < 1$ ) and as a result  $N_l(z_1, z_2)$  and  $N_{l+1}(z_1, z_2)$  converge both to a common steady-state value:

$$N(z_1, z_2) \triangleq \lim_{l \rightarrow \infty} N_l(z_1, z_2).$$

By taking the  $l \rightarrow \infty$  limit of equation (1), we obtain:

$$\begin{aligned}
[z_1 - S_1(A(z_1, z_2))] N(z_1, z_2) = & z_1 \frac{z_2 A(z_1, z_2) - S_2(A(z_1, z_2))}{z_2} N(0, 0) \\
& + \frac{z_1 S_2(A(z_1, z_2)) - z_2 S_1(A(z_1, z_2))}{z_2} N(0, z_2).
\end{aligned} \tag{2}$$

It now remains for us to determine the unknown function  $N(0, z_2)$  and the unknown parameter  $N(0, 0)$ . This can be done in two steps. First, we notice that  $N(z_1, z_2)$  must be bounded for all values of  $z_1$  and  $z_2$  such that  $|z_1| \leq 1$  and  $|z_2| \leq 1$ . In particular, this should be true for  $z_1 = Y(z_2)$ , with  $Y(z_2) \triangleq S_1(A(Y(z_2), z_2))$  and  $|z_2| \leq 1$ , since it follows from Rouché's theorem that there is exactly one solution such that  $|Y(z_2)| \leq 1$  for all such  $z_2$ . Notice that  $Y(1)$  equals 1. The above implies that if we choose  $z_1 = Y(z_2)$  in equation (2), where  $|z_2| \leq 1$ , the left hand side of this equation vanishes. The same must then be true for the right hand side, yielding

$$N(0, z_2) = N(0, 0) \frac{z_2 A(Y(z_2), z_2) - S_2(A(Y(z_2), z_2))}{z_2 - S_2(A(Y(z_2), z_2))}. \tag{3}$$

Finally, in order to find an expression for  $N(0, 0)$ , we put  $z_1 = z_2 = 1$  and use de l'Hospital's rule in equation (2). Therefore, we need the first derivative of  $Y(z)$  for  $z = 1$  and this is given by

$$\begin{aligned}
Y'(1) &= \mu_1(\lambda_1 Y'(1) + \lambda_2) \\
&= \frac{\lambda_2 \mu_1}{1 - \rho_1}.
\end{aligned} \tag{4}$$

We then obtain  $N(0, 0)$ :

$$N(0, 0) = \frac{1 - \rho}{1 - \rho + \lambda_1 + \lambda_2}. \tag{5}$$

A fully determined expression for  $N(z_1, z_2)$  can now be derived by combining equations (2) and (3):

$$\begin{aligned}
N(z_1, z_2) = & N(0, 0) \left[ \frac{z_1(z_2 A(z_1, z_2) - S_2(A(z_1, z_2)))}{(z_1 - S_1(A(z_1, z_2)))(z_2 - S_2(A(Y(z_2), z_2)))} \right. \\
& + \frac{S_2(A(Y(z_2), z_2))(S_1(A(z_1, z_2)) - z_1 A(z_1, z_2))}{(z_1 - S_1(A(z_1, z_2)))(z_2 - S_2(A(Y(z_2), z_2)))} \\
& \left. + \frac{A(Y(z_2), z_2)(z_1 S_2(A(z_1, z_2)) - z_2 S_1(A(z_1, z_2)))}{(z_1 - S_1(A(z_1, z_2)))(z_2 - S_2(A(Y(z_2), z_2)))} \right],
\end{aligned} \tag{6}$$

with  $N(0, 0)$  given by equation (5).

#### 4 System contents at the beginning of arbitrary slots

In this section, we analyse the system contents at the beginning of arbitrary slots. The joint pgf of the system contents of both priority classes at the beginning of slot  $k$  is defined as:

$$U_k(z_1, z_2) \triangleq E[z_1^{u_1, k} z_2^{u_2, k}].$$

In order to derive an expression for  $U_k(z_1, z_2)$ , we have to know the status of the server during slot  $k$ . There are 3 possibilities: the server can be idle, a low priority or a high priority packet can be in service during slot  $k$ . This yields

$$\begin{aligned} U_k(z_1, z_2) &= E[z_1^{u_1, k} z_2^{u_2, k} \{\text{no service}\}] + E[z_1^{u_1, k} z_2^{u_2, k} \{\text{service class 2 packet}\}] \\ &\quad + E[z_1^{u_1, k} z_2^{u_2, k} \{\text{service class 1 packet}\}] \\ &= U_k(0, 0) E[z_1^{u_1, k} z_2^{u_2, k} | \text{no service}] + (1 - U_k(0, 0)) \\ &\quad \left\{ \frac{\rho_2}{\rho} E[z_1^{u_1, k} z_2^{u_2, k} | \text{service class 2 packet}] \right. \\ &\quad \left. + \frac{\rho_1}{\rho} E[z_1^{u_1, k} z_2^{u_2, k} | \text{service class 1 packet}] \right\}. \end{aligned} \quad (7)$$

The last transition is found as follows: the server is idle during a slot if and only if the system was empty at the beginning of the slot, i.e.,  $\text{Prob}[\text{no service}] = U_k(0, 0)$ . On the other hand, if the server is busy during slot  $k$ , a class  $j$  packet is being served with probability  $\rho_j / \rho_T$  ( $j = 1, 2$ ). If slot  $k$  is a start slot, we will assume that it is start slot  $l$ . If slot  $k$  is not a start slot on the other hand, the last start slot preceding slot  $k$  is start slot  $l$ . Equation (7) then becomes

$$\begin{aligned} U_k(z_1, z_2) &= U_k(0, 0) + (1 - U_k(0, 0)) \\ &\quad \left\{ \frac{\rho_2}{\rho} E[z_1^{u_1, k} z_2^{u_2, k} | n_{1, l} = 0, n_{2, l} > 0] + \frac{\rho_1}{\rho} E[z_1^{u_1, k} z_2^{u_2, k} | n_{1, l} > 0] \right\}. \end{aligned} \quad (8)$$

This can be understood as follows: the server is idle during slot  $k$  if there were no packets in the system at the beginning of slot  $k$ , a class 2 packet is being served during slot  $k$  if there were no class 1 packets and at least one class 2 packet in the system at the beginning of start slot  $l$  and a class 1 packet is in service during slot  $k$  if there was at least one class 1 packet in the system at the beginning of start slot  $l$ . We denote the elapsed service time of the packet in service (if any) during slot  $k$  by  $s_k^+$ . The system contents at the beginning of slot  $k$  is a superposition of the system contents at the beginning of start slot  $l$  and the arrivals during  $s_k^+$ , yielding

$$\begin{aligned} U_k(z_1, z_2) &= U_k(0, 0) + (1 - U_k(0, 0)) \\ &\quad \left\{ \frac{\rho_2}{\rho} E \left[ z_1^{\sum_{i=1}^k a_{1, k-i}} z_2^{\sum_{i=1}^k a_{2, k-i}} \mid n_{1, l} = 0, n_{2, l} > 0 \right] \right. \\ &\quad \left. + \frac{\rho_1}{\rho} E \left[ z_1^{n_{1, l} + \sum_{i=1}^k a_{1, k-i}} z_2^{n_{2, l} + \sum_{i=1}^k a_{2, k-i}} \mid n_{1, l} > 0 \right] \right\} \\ &= U_k(0, 0) + (1 - U_k(0, 0)) \left\{ \frac{\rho_2}{\rho} S_{2, k}^-(A(z_1, z_2)) \frac{N_l(0, z_2) - N_l(0, 0)}{N_l(0, 1) - N_l(0, 0)} \right. \\ &\quad \left. + \frac{\rho_1}{\rho} S_{1, k}^+(A(z_1, z_2)) \frac{N_l(z_1, z_2) - N_l(0, z_2)}{1 - N_l(0, 1)} \right\}. \end{aligned} \quad (9)$$

Hereby is  $S_{j, k}^+(z)$  ( $j = 1, 2$ ) defined as the pgf of the elapsed service time of the class  $j$  packet in service at the beginning of slot  $k$ .

We denote the steady-state version of  $U_k(z_1, z_2)$  by  $U(z_1, z_2)$ , i.e.,

$$U(z_1, z_2) = \lim_{k \rightarrow \infty} U_k(z_1, z_2).$$

It is shown in e.g. [2] that the steady-state version of  $S_{j,k}(z)$  yields

$$\lim_{k \rightarrow \infty} S_{j,k}^+(z) = \frac{S_j(z) - 1}{S_j'(1)(z - 1)}, \quad (10)$$

for  $j = 1, 2$ . It now remains for us to determine the unknown parameter  $U(0, 0)$ . Keeping in mind that if the server is idle during slot  $k$ , slot  $k$  is a start slot,  $U(0, 0)$  can easily be found as follows:

$$\begin{aligned} U(0, 0) &= \lim_{k \rightarrow \infty} \text{Prob}[u_{1,k} = u_{2,k} = 0] \\ &= \lim_{k,l \rightarrow \infty} \text{Prob}[n_{1,l} = n_{2,l} = 0 \text{ and slot } k \text{ is a start slot}] \\ &= \lim_{k,l \rightarrow \infty} \text{Prob}[n_{1,l} = n_{2,l} = 0 | \text{slot } k \text{ is a start slot}] \text{Prob}[\text{slot } k \text{ is a start slot}] \end{aligned}$$

There are three possibilities for slot  $k$  to be a start slot: the system is empty at the beginning of slot  $k$ , slot  $k$  is the first slot of the service time of a class 1 packet or slot  $k$  is the first slot of the service time of a class 2 packet.  $U(0, 0)$  then becomes

$$\begin{aligned} U(0, 0) &= N(0, 0) \left[ U(0, 0) + \frac{1 - U(0, 0) \rho_1}{\mu_1 \rho} + \frac{1 - U(0, 0) \rho_2}{\mu_2 \rho} \right] \\ &= 1 - \rho. \end{aligned} \quad (11)$$

Using equations (6) and (10) in the steady-state version of equation (9), we derive a fully determined version for  $U(z_1, z_2)$ :

$$\begin{aligned} U(z_1, z_2) &= U(0, 0) \left\{ \frac{S_1(A(z_1, z_2))(z_1 - 1)}{z_1 - S_1(A(z_1, z_2))} + \frac{A(Y(z_2), z_2) - 1}{A(z_1, z_2) - 1} \right. \\ &\quad \left[ \frac{z_1 S_2(A(z_1, z_2))(S_1(A(z_1, z_2)) - 1)}{(z_1 - S_1(A(z_1, z_2)))(z_2 - S_2(A(Y(z_2), z_2)))} \right. \\ &\quad + \frac{z_1 z_2 (S_2(A(z_1, z_2)) - S_1(A(z_1, z_2)))}{(z_1 - S_1(A(z_1, z_2)))(z_2 - S_2(A(Y(z_2), z_2)))} \\ &\quad \left. \left. + \frac{z_2 S_1(A(z_1, z_2))(1 - S_2(A(z_1, z_2)))}{(z_1 - S_1(A(z_1, z_2)))(z_2 - S_2(A(Y(z_2), z_2)))} \right] \right\}, \end{aligned} \quad (12)$$

with  $U(0, 0)$  given by (11). From the two-dimensional pgf  $U(z_1, z_2)$ , we can easily derive expressions for the pgf's of the system contents of high and low priority packets at the beginning of an arbitrary slot - denoted by  $U_1(z)$  and  $U_2(z)$  respectively - yielding

$$\begin{aligned} U_1(z) &\triangleq \lim_{k \rightarrow \infty} E[z^{u_{1,k}}] \\ &= U(z, 1) \\ &= (1 - \rho_1) \frac{S_1(A_1(z))(z - 1)}{z - S_1(A_1(z))} + \lambda_2 \frac{S_1(A_1(z))(z - 1)}{z - S_1(A_1(z))} \left\{ \frac{S_2(A_1(z)) - 1}{A_1(z) - 1} - \mu_2 \right\} \end{aligned} \quad (13)$$

$$\begin{aligned} U_2(z) &\triangleq \lim_{k \rightarrow \infty} E[z^{u_{2,k}}] \\ &= U(1, z) \\ &= (1 - \rho) \frac{S_2(A_2(z))(z - 1)}{z - S_2(A(Y(z), z))} \frac{A(Y(z), z) - 1}{A_2(z) - 1}. \end{aligned} \quad (14)$$

We can also derive expressions for the pgf of the total system contents at the beginning of an arbitrary slot - denoted by  $U_T(z)$  - yielding

$$U_T(z) \triangleq \lim_{k \rightarrow \infty} E[z^{u_{1,k} + u_{2,k}}]$$

$$\begin{aligned}
&= U(z, z) \\
&= (1 - \rho) \left[ \frac{S_1(A_T(z))(z-1)}{z - S_1(A_T(z))} + \frac{A(Y(z), z) - 1}{A_T(z) - 1} \right. \\
&\quad \left. \frac{z(z-1)(S_2(A_T(z)) - S_1(A_T(z)))}{(z - S_1(A_T(z)))(z - S_2(A_T(z)))} \right]. \tag{15}
\end{aligned}$$

In the special case that  $S_1(z) = S_2(z) (= S(z))$ , i.e., when the distributions of the service times of high and low priority packets are the same,  $U_T(z)$  becomes

$$U_T(z) = (1 - \rho) \frac{S(A_T(z))(z-1)}{z - S(A_T(z))}. \tag{16}$$

This is the expression of the pgf of the system contents in a single-class GI-G-1 queue with FIFO scheduling. Indeed, if all packets have the same service distribution, the scheduling does not influence the total system contents.

## 5 Calculation of moments

The function  $Y(z)$  can only be explicitly found in case of some simple arrival and service processes. Its derivatives for  $z = 1$ , necessary to calculate the moments of the system contents, on the contrary, can be calculated in closed-form. For example,  $Y'(1)$  is given by equation (4). Let us define  $\lambda_{ij}$ ,  $\lambda_{TT}$  and  $\mu_{jj}$  as

$$\begin{aligned}
\lambda_{ij} &\triangleq \left. \frac{\partial^2 A(z_1, z_2)}{\partial z_i \partial z_j} \right|_{z_1=z_2=1} \\
\lambda_{TT} &\triangleq \left. \frac{d^2 A_T(z)}{dz^2} \right|_{z=1} \\
\mu_{jj} &\triangleq \left. \frac{d^2 S_j(z)}{dz^2} \right|_{z=1},
\end{aligned}$$

with  $i, j = 1, 2$ . Now we can calculate the mean values of the system contents of both classes by taking the first derivatives of the respective pgf's for  $z = 1$ . We find

$$E[u_1] = \rho_1 + \frac{1}{2} \frac{\mu_1 \lambda_{11}}{1 - \rho_1} + \frac{\lambda_1 \lambda_{1\mu_{11}} + \lambda_2 \mu_{22}}{2(1 - \rho_1)}, \tag{17}$$

for the mean value of the system contents of class 1 packets and

$$E[u_2] = \rho_2 + \frac{1}{2} \frac{\mu_2^2 \lambda_{2\lambda_{11}}}{(1 - \rho)(1 - \rho_1)} + \frac{1}{2} \frac{2\mu_1 \lambda_{12} + \mu_2 \lambda_{22}}{1 - \rho} + \frac{\lambda_2 \lambda_{1\mu_{11}} + \lambda_2 \mu_{22}}{2(1 - \rho)(1 - \rho_1)}, \tag{18}$$

for the mean value of the system contents of class 2 packets. Furthermore, the mean total system contents can be found by taking the first derivative of  $U_T(z)$  for  $z = 1$ , yielding

$$\begin{aligned}
E[u_T] &= \rho + \frac{1}{2} \frac{\mu_1 \lambda_{TT}}{1 - \mu_1 \lambda_T} + \frac{1}{2} \frac{\lambda_2^2 (\mu_2 - \mu_1) \mu_1^2 \lambda_{11}}{(1 - \rho)(1 - \mu_1 \lambda_T)(1 - \rho_1)} + \frac{\lambda_2 (\mu_2 - \mu_1) \mu_1 \lambda_{12}}{(1 - \rho)(1 - \mu_1 \lambda_T)} \\
&\quad + \frac{1}{2} \frac{(\mu_2 - \mu_1)(1 - \rho_1) \lambda_{22}}{(1 - \rho)(1 - \mu_1 \lambda_T)} + \frac{1}{2} \frac{(\lambda_1 \mu_{11} + \lambda_2 \mu_{22})(\lambda_T - \lambda_1 \rho)}{(1 - \rho)(1 - \rho_1)}, \tag{19}
\end{aligned}$$

for the mean total system contents. It is easily verified that equations (17) - (19) satisfy  $E[u_T] = E[u_1] + E[u_2]$ .

In a similar way, expressions for the variance can be calculated by taking the appropriate derivatives of the respective generating functions. By using Little's law, the mean packet delay of a class  $j$  packet, denoted by  $E[d_j]$ , i.e., the mean time a class  $j$  packet stays in the system, can be calculated as well, i.e.,  $E[d_j] = E[u_j] / \lambda_j$  ( $j = 1, 2$ ).

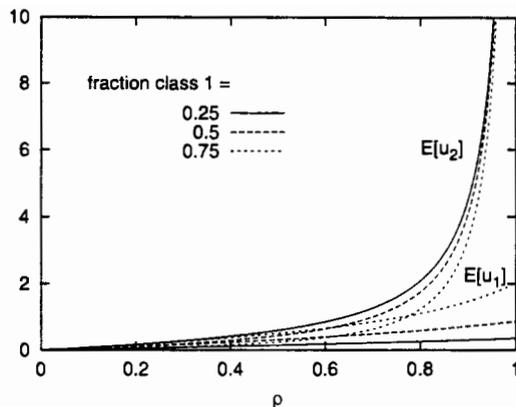


Figure 1: Mean system contents versus the total load

## 6 Numerical examples

In this section, we present some numerical examples. We assume the traffic of the two classes to be arriving according to a two-dimensional binomial process. The joint pgf  $A(z_1, z_2)$  is given by:

$$A(z_1, z_2) = \left(1 - \frac{\lambda_1}{N}(1 - z_1) - \frac{\lambda_2}{N}(1 - z_2)\right)^N.$$

The arrival rate of class  $j$  traffic is thus given by  $\lambda_j$  ( $j = 1, 2$ ). This arrival process occurs for instance at an output queue of a  $N \times N$  switch fed by a Bernoulli process at the  $N$  inlets (see [16]). Notice also that if  $N \rightarrow \infty$ , the arrival process is a superposition of two Poisson streams. In the remainder of this section, we assume that  $N = 16$ . We will furthermore assume deterministic service times for both classes.

In Figure 1 and 2, the mean value and the variance of the system contents of class 1 and class 2 packets is shown as a function of the total load  $\rho$ , when  $\mu_1 = \mu_2 = 2$ . The fraction of the arrival rate of class 1 traffic is 0.25, 0.5 and 0.75 respectively of the total arrival rate. One can easily see the influence of priority scheduling: the mean, as well as the variance of the number of class 1 packets in the system is severely reduced by the non-preemptive priority scheduling; the opposite holds for class 2 packets. In addition, it also becomes apparent that increasing the fraction of high priority packets in the overall mix increases the amount of the high priority packets while decreasing the amount of low priority packets in the buffer.

In Figure 3, the mean packet delay of class 1 and class 2 packets is shown - found by using Little's law - as a function of the total load  $\rho$ , when  $\mu_1 = \mu_2 = 2$  and the fraction of the arrival rate of class 1 traffic is again 0.25, 0.5 and 0.75 respectively of the total arrival rate. In order to compare with FIFO scheduling, we have also shown the mean value of the packet delay in that case. Since, in this example, the service times of the class 1 and class 2 packets are equal, the packet delay is then of course the same for class 1 and class 2 packets, and can thus be calculated as if there is only one class of packets arriving according to an arrival process with pgf  $A(z, z)$ . This has already been analyzed, e.g., in [2]. One can observe the influence of priority scheduling: mean delay of class 1 packets reduces significantly. The price to pay is of course a larger mean delay for class 2 packets. If this kind of traffic is not delay-sensitive, as assumed, this is not a problem. Also, the smaller the fraction of high priority packets in the overall traffic mix, the lower the mean packet delay of both classes will be.

Figure 4 shows the mean system contents of class 1 and class 2 packets as a function of the mean service time of the class 1 packets, when  $\lambda_T = 0.2$ ,  $\mu_2 = 2$  and the fraction of the arrival rate of class

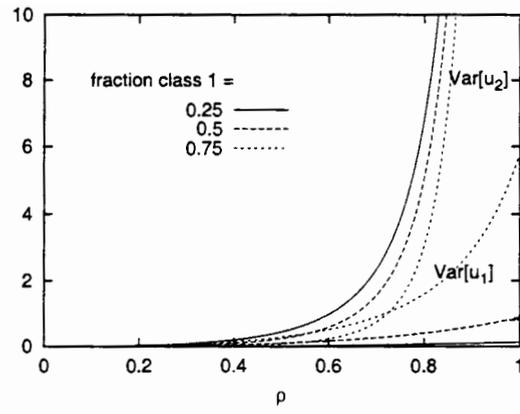


Figure 2: Variance of the system contents versus the total load

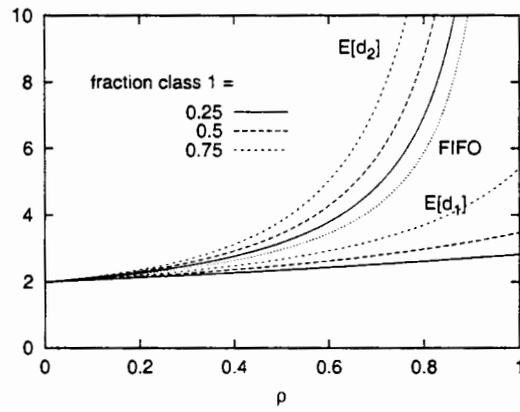


Figure 3: Mean packet delay versus the total load

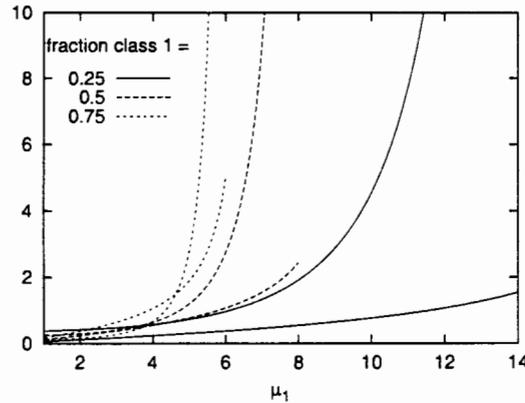


Figure 4: Mean system contents versus the mean service time of class 1 packets

1 traffic is 0.25, 0.5 and 0.75 respectively of the total arrival rate. In order to have a stable system, i.e., the total load has to be less than 1,  $\mu_1$  must be strictly less than 14, 8 or 6 when the fraction of the arrival rate of class 1 traffic is 0.25, 0.5 or 0.75 respectively. For these values the mean system contents of the low priority packets becomes infinite (as shown in the figure). Because of the priority scheduling, it can be seen that the service time of the high priority packets has a large influence on the system contents of both classes, and this influence is more pronounced when the fraction of high priority packets is larger.

Figure 5 shows the mean system contents of class 1 and class 2 packets as a function of the mean service time of the class 2 packets, when  $\lambda_T = 0.2$ ,  $\mu_1 = 2$  and the fraction of the arrival rate of class 1 traffic is 0.25, 0.5 and 0.75 respectively of the total arrival rate. In order to have a stable system, i.e., the total load has to be less than 1,  $\mu_2$  must be strictly less than 6, 8 or 14 when the fraction of the arrival rate of class 1 traffic is 0.25, 0.5 or 0.75 respectively. Again, for these values the mean system contents of the low priority packets becomes infinite. It can be seen that the service time of the low priority packets has a large influence on the mean system contents of low priority packets, while the influence on the mean high priority system contents is not too large. Nevertheless, the low priority traffic has an impact on the characteristics of the high priority packets, since the priority scheduling is non-preemptive.

Figure 6 shows the mean value of the system contents of class 1 packets as a function of the total load, when  $\lambda_1 = 0.25$ ,  $\mu_1 = 2$  and  $\mu_2 = 1, 2, 4, 8, 16$ . This figure shows the influence of the *non-preemptive* priority scheduling. When the service time of a class 2 packet is assumed to be deterministically 1 slot, i.e.,  $\mu_2 = 1$ , the preemptive priority scheduling has the same effect as the non-preemptive priority scheduling. If  $\mu_2 > 1$ , the non-preemptive priority has worse performance than the preemptive priority scheduling in terms of mean system contents for class 1 packets. Furthermore, for a given value of the low priority packet length, the mean high priority system contents increases linearly to the total load  $\rho$ .

## 7 Conclusion

In this paper, we analyzed the system contents in a queueing system with non-preemptive HOL priority scheduling. A generating-functions-approach was adopted, which led to closed-form expressions of performance measures, such as mean of system contents and packet delay of both classes, that are easy to evaluate. The model included possible correlation between the number of arrivals of the two classes during a slot and general service times for packets of both classes. Therefore, the results could

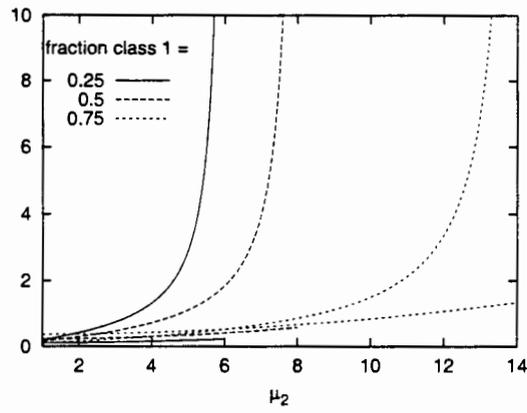


Figure 5: Mean system contents versus the mean service time of class 2 packets

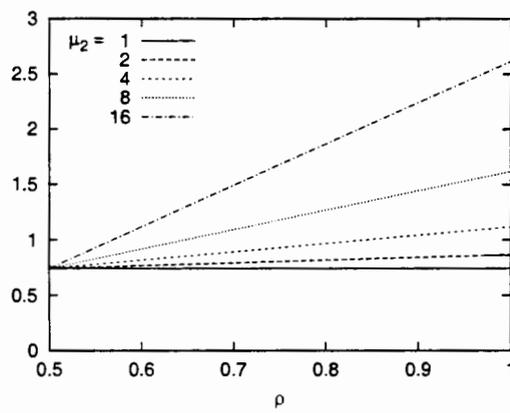


Figure 6: Mean packet delay of class 1 packets when the service time of class 2 packets equals 1, 2, 4, 8, 16

be used to analyse performance of buffers in an IP context.

## Acknowledgements

The authors would like to thank the anonymous referee for valuable comments.

## References

- [1] J.J. Bae and T. Suda, *Survey of traffic control schemes and protocols in ATM networks*, Proceedings of the IEEE 79(2), pp. 170-189, 1991.
- [2] H. Bruneel and B.G. Kim, *Discrete-time models for communication systems including ATM*, Kluwer Academic Publishers, Boston, 1993.
- [3] N.K. Jaiswal, *Priority queues*, Academic Press, New York, 1968.
- [4] A. Khamisy and M. Sidi, *Discrete-time priority queues with two-state Markov modulated arrivals*, Stochastic Models 8(2), pp. 337-357, 1992.
- [5] K. Laevens and H. Bruneel, *Discrete-time multiserver queues with priorities*, Performance Evaluation 33(4), pp. 249-275, 1998.
- [6] K. Liu, D.W. Petr, V.S. Frost, H. Zhu, C. Braun and W.L. Edwards, *Design and analysis of a bandwidth management framework for ATM-based broadband ISDN*, IEEE Communications Magazine, pp. 138-145, 1997.
- [7] I. Rubin and Z. Tsai, *Message delay analysis of multiclass priority TDMA, FDMA, and discrete-time queueing systems*, IEEE Transactions on Information Theory 35(3), pp. 637-647, 1989.
- [8] D.A. Stanford, *Interdeparture-time distributions in the non-preemptive priority  $\sum M_i/G_i/1$  queue*, Performance Evaluation 12, pp. 43-60, 1991.
- [9] A. Sugahara, T. Takine, Y. Takahashi and T. Hasegawa, *Analysis of a nonpreemptive priority queue with SPP arrivals of high class*, Performance Evaluation 21, pp. 215-238, 1995.
- [10] L. Takacs, *Priority queues*, Operations Research 12, pp. 63-74, 1964.
- [11] H. Takagi, *Queueing analysis A foundation of Performance Evaluation Volume 1: Vacation and priority systems*, North-Holland, 1991.
- [12] T. Takine, Y. Matsumoto, T. Suda and T. Hasegawa, *Mean waiting times in nonpreemptive priority queues with Markovian arrival and i.i.d. service processes*, Performance Evaluation 20, pp. 131-149, 1994.
- [13] T. Takine, B. Sengupta and T. Hasegawa, *An analysis of a discrete-time queue for broadband ISDN with priorities among traffic classes*, IEEE Transactions on Communications 42 (2-4), pp. 1837-1845, 1994.
- [14] T. Takine, *A nonpreemptive priority MAP/G/1 queue with two classes of customers*, Journal of the Operations Research Society of Japan 39(2), pp. 266-290, 1996.
- [15] P. Van Mieghem, B. Steyaert and G.H. Petit, *Performance of cell loss priority management schemes in a single server queue*, International Journal of Communication Systems 10, pp. 161-180, 1997.
- [16] J. Walraevens and H. Bruneel, *HOL priority in an ATM output queueing switch*. Proceedings of the seventh IFIP workshop on performance modelling and evaluation of ATM/IP networks, Antwerp. 28-30 June, 1999.