



ELSEVIER

Performance Evaluation 27&28 (1996) 391–409



Analysis of communication systems with timed token protocols using the power-series algorithm

J.P.C. Blanc^{a,*}, L. Lenzini^{b,1}

^a *Department of Econometrics, Tilburg University, P.O. Box 90153, 5000 LE Tilburg, Netherlands*

^b *Department of Information Engineering, University of Pisa, via Diotisalvi, 2, I-56126 Pisa, Italy*

Abstract

The IEEE 802.4 and FDDI (Fibre Distributed Data Interface) standards are high speed MAC (Medium Access Control) protocols for LAN/MANs employing a timer-controlled token passing mechanism, the so-called Timed Token Protocol, to control station access to the shared media. MAC protocols belonging to the class of timed token protocols support synchronous and real-time (i.e., time-critical) applications, and provide priority among asynchronous (i.e., non time-critical) applications. During the last few years, a lot of research has focused on the study of timed token protocols to obtain performance measures such as throughputs or mean waiting times. The recent development of the Power-Series Algorithm (PSA) has opened new perspectives in the analysis of this class of protocols. This paper shows the versatility of the PSA technique to evaluate the station buffer occupancy and delay distributions of a very general model which can be used to represent the behavior of several LAN/MANs MAC protocols, among which the timed token MAC protocols. Specifically, the focus of the paper is on the solution of an almost exact model of the IEEE 802.4 MAC protocol. Since the model we propose and solve numerically by exploiting the PSA technique is an approximate model of the FDDI MAC protocol, the paper also reports on a comparison between performance measures obtained for this model and simulation results for the corresponding (exact) model of FDDI.

Keywords: Timed token protocol; Medium access control; IEEE 802.4 Token Bus; FDDI (Fibre Distributed Data Interface); Power-series algorithm

1. Introduction

The idea behind the class of timed token protocols, such as IEEE 802.4 and FDDI, is firstly to partition the services they can provide their users into two main classes, time-critical and non time-critical type of services; and secondly, to employ a token passing MAC protocol with a cycle-dependent timing mechanism which limits the amount of data (organized into frames) transmitted by a station for each class of service in a cycle. The non time-critical class of service may be further subdivided into

* Corresponding author. E-mail: blanc@kub.nl.

¹ E-mail: lenzini@iet.unipi.it.

subclasses according to a priority scheme which is normally optional. The elements of the IEEE 802.4 MAC protocol which are relevant to our analysis, and the main difference with respect to FDDI are described in Appendix A. The FDDI protocol is described in Appendix B.

The main difficulty in the analysis of a timed token MAC protocol is the high degree of complexity and interdependence of the various processes that describe the operations of the protocol itself. In fact, when a station has seized the token, time-critical (priority 6 in IEEE 802.4 or synchronous in FDDI) frames (if any) are always transmitted, whereas asynchronous frames are only transmitted if the token is early. This implies that there are interdependencies between the total service time given at one station, the service time required at subsequent stations, and the total cycle time. Therefore, exact analytically-tractable solutions for timed token protocols are very difficult to formulate. Simplifying assumptions thus have to be made in order to obtain analytically tractable solutions. There are many papers on FDDI (see, for example, [10,11,15,17,18,22,24,25]) and IEEE 802.4 (see, for example, [12,18,5,26]), providing bounds and mean performance figures (typically throughput and mean waiting time). Takagi [22] studies the effect of the token rotation timer on the delay-throughput performance of a symmetric single buffer system operating under a timed-token protocol. A gauss elimination method is used to calculate the finite limiting state distribution of the embedded Markov chain for an FDDI network connecting a moderate number of stations. In [24] marginal queue length distributions in FDDI networks with only asynchronous traffic are computed with an iterative scheme utilizing several approximations. A summary of the research work related to the FDDI performance evaluation can be found in [9].

The model

What we want to stress in the present paper is that the recent Power-Series Algorithm (PSA) [6–8] allows the numerical calculation of station buffer occupancy and delay distributions for detailed models of moderate size. In [4], numerical results are obtained with the PSA for a timed token protocol (FDDI) model in which the switchover times between stations are zero, interarrival and service times are exponential, and either each station implements the 1-limited service discipline (i.e., asynchronous traffic has an additional limitation beside the token rotation time restriction), or no station has a service limit (i.e., no synchronous traffic). Furthermore, in [4] the influence of the accumulated lateness is not considered, and the constraint on the actual cycle time has been replaced by a constraint on the mean total time needed by the various queues to serve (transmit) packets in the last cycle.

The model proposed in this paper removes several of the above limitations, and becomes almost exact for the IEEE 802.4 token bus MAC protocol. Specifically, in our model the number of frames allowed to be transmitted by a station during one cycle of the server may be restricted by any limit, the switchover times are different from zero, and the switchover and service times are described by Erlang distributed random variables, cf. e.g. [16]. Note that by using an Erlang distribution with a sufficient number of stages, we can, in principle, approximate as closely as we want both the deterministic switchover times between stations in a real LAN/MAN and constant service times.² On the basis of an example, it will be shown that accurate estimates for the case of constant service and switchover times can be obtained from the cases of exponentially and Erlang-2 distributed times. The first example in Section 5 also shows that replacing the constraints on the actual cycle times by constraints based on the actual numbers of served

² Constant service times are particularly important because the wide area ATM subnetworks, to which LAN/MANs will be interconnected, manage packets of fixed length (i.e., cells) rather than of variable lengths.

frames times the corresponding mean service times and on the mean switchover times has only a minor influence on the waiting time characteristics.

In our model a distinction is made between access class 6 queues (type ST queues, i.e., synchronous or time-constraint traffic) and the other queues (type AT queues, i.e., asynchronous or non time-constraint traffic). We allow an arbitrary distribution of type ST and type AT queues within the logical ring; furthermore, type ST and type AT queues can have arbitrary *hi_pri_token_hold_time* and TTRT (to be indicated briefly by *K* and *R* respectively in the rest of the paper) parameters respectively. Thus, there can be as many different *hi_pri_token_hold_time*'s as the number of type ST queues, and likewise for type AT queues. Stations generating frames of different access classes are modelled by multiple queues, one for each relevant access class.

A virtual substation may initiate a transmission of a non time-critical frame if the token hold timer has not reached the TTRT threshold. This might cause an additional delay in the release of the token, hereafter called overflow transmission (asynchronous overrun in FDDI), which is bounded by the time for the transmission of a frame of maximum length. According to the IEEE 802.4 standard, the on-going transmission shall nevertheless be completed. To make our model as general as possible and in order to make a comparison with FDDI (which allows for asynchronous overrun), we model this overflow transmission as well.

To conclude, in the timed token model we propose and solve throughout, we only neglect the accumulated lateness. Therefore, the proposed model is an exact model (to the extent that an Erlang distribution with a sufficient number of stages represents a constant distribution) for the IEEE 802.4 standard and an approximate (but fairly precise) model for FDDI.

Finally, in our model, time-critical and non time-critical frames are assumed to be generated by a Poisson process. While for non time-critical frames this choice is commonly made, for time-critical frames, at first glance, it may seem inadequate. However, when time-critical frames are generated by a Variable Bit Rate (VBR) video source [21] or by an aggregate number of voice sources [13], it has been shown that a Poisson distribution very well approximates the real sources. In principle, the PSA can also handle systems with Markovian arrival processes (MAPs), cf. [27], but this requires a still larger supplementary space than the MAC protocol already demands.

Organization of the paper

Notations for our model are introduced in Section 2. This section also contains a detailed description of the model, in particular of the timed token access control protocol, and some remarks on the stability of the system. In Section 3 the queue-length process for this model is transformed into a Markov process with the aid of some supplementary variables, and the balance equations for the stationary state probabilities are given. The recurrence relations of the PSA for this model are derived in Section 4. Several numerical examples are presented in Section 5, where results of our model are also compared with simulation results for systems with an FDDI protocol. Section 6 contains some concluding remarks.

2. The model: notations and assumptions

The communication system consists of *S* stations (queues) and a single token (server) which visits the stations in cyclic order. Frames arrive at queue *j* according to a Poisson process with rate λ_j , $j = 1, \dots, S$. The superposition of the arrival processes at the various queues is a Poisson process with

rate $\Lambda \doteq \sum_{j=1}^S \lambda_j$. Each queue may contain an unbounded number of frames. At each station frames are served in order of arrival. Service times of frames arriving at queue j are assumed to be Erlang distributed with Ψ_j exponential phases, each with rate μ_j , $j = 1, \dots, S$. The mean service time β_j for frames at queue j is then given by $\beta_j = \Psi_j / \mu_j$, $j = 1, \dots, S$. The load offered at queue j is defined as $\rho_j \doteq \lambda_j \beta_j$, $j = 1, \dots, S$, and $\rho \doteq \sum_{j=1}^S \rho_j$ will denote the total load offered to the system. The times the server needs for switching from queue $j - 1$ (queue 0 indicating queue S) to queue j are also assumed to be Erlang distributed, with Ω_j exponential phases, each with rate ν_j , $j = 1, \dots, S$. The mean switchover time δ_j from queue $j - 1$ to queue j is given by $\delta_j = \Omega_j / \nu_j$, $j = 1, \dots, S$. The mean total switchover time during one cycle of the server along the queues is denoted by $\Delta \doteq \sum_{j=1}^S \delta_j$. The target token rotation time, exclusive the total expected switchover time Δ , at queue j will be denoted by R_j , $j = 1, \dots, S$. Let v_j denote the actual number of frames served at queue j since the beginning of the last completed visit of the server to that queue, $j = 1, \dots, S$. This means that if the server is currently serving a frame at queue j , $j = 1, \dots, S$, then v_j is the sum of the number of frames served at that queue during the previous visit and those already served during the current visit. When the server leaves queue j , the value of v_j is reset to the number of services performed during the just completed visit to this queue, $j = 1, \dots, S$. For a given $\mathbf{v} \doteq (v_1, \dots, v_S)$, let

$$T(\mathbf{v}) = T(v_1, \dots, v_S) \doteq \sum_{j=1}^S v_j \beta_j \quad (2.1)$$

denote the (approximate) actual token rotation time, also exclusive the total expected switchover time Δ , based on constant service times. The maximum number of frames that the server is allowed to serve during a visit to queue j will be denoted by K_j , $j = 1, \dots, S$. When the server arrives at queue j , it will pass if that queue is empty or if the (expected) token rotation timer has expired, i.e. if $T(\mathbf{v}) \geq R_j$; otherwise, it will start servicing frames at that queue until either the queue becomes empty, or the maximal number of frames, K_j , has been served, or the rotation timer $T(\mathbf{v})$ which is augmented by β_j after each service completion exceeds the target R_j , $j = 1, \dots, S$. Note that this target will have no effect on the number of frames served during a visit of the server if $R_j > \sum_{h=1}^S K_h \beta_h + (K_j - 1) \beta_j$, $j = 1, \dots, S$. On the other hand, if $K_j \beta_j \geq R_j$ then the limit K_j has no effect because the maximal number of frames served during a visit of the server to queue j is equal to $\lceil R_j / \beta_j \rceil$, $j = 1, \dots, S$ ($\lceil x \rceil$ denotes the smallest integer larger than or equal to x). Let $\hat{K}_j \doteq \min\{K_j, \lceil R_j / \beta_j \rceil\}$ be the effective service limit at queue j , $j = 1, \dots, S$.

Since in our model frames are approximately of constant length, frames belonging to type ST queues are served (i.e., transmitted) according to a K -limited service discipline that limits the number of frames that can be served during the token visit; the value of the limit is generally station dependent. From a purely mathematical standpoint it is convenient to assume that the target token rotation time R is infinite for a type ST queue, and that the frame limit K is infinite for a type AT queue. Hence, we model a generic type ST queue j by K_j finite and R_j infinite, and we model a generic type AT queue by R_j finite and K_j infinite.

Because the number of frames served per cycle at a queue in the above described polling systems with token rotation time restrictions cannot be more than in polling systems without such restrictions, but with the same effective service limits \hat{K}_j , $j = 1, \dots, S$, a necessary condition for stability of the

former systems is, cf. e.g. [7],

$$\rho + \Delta \max_{j=1,\dots,S} \left\{ \lambda_j / \hat{K}_j \right\} < 1. \quad (2.2)$$

Further, it is shown in [19] that in the special case when the mean service times β_j , $j = 1, \dots, S$, are all equal to, say, β , and when the targets R_j , $j = 1, \dots, S$ are all equal to the same multiple of the mean service time, say, $R_j = R = M\beta$, and the service limits do not influence the system, i.e. $K_j = M = R/\beta$, for $j = 1, \dots, S$, the restriction on the rotation times implies that the condition

$$\rho + \frac{\Delta}{R} \max_{j=1,\dots,S} \{\rho + \rho_j\} < 1, \quad (2.3)$$

should hold in case of stability. In all other cases of our model, the condition for stability seems to be unknown (in [23] a generalization of (2.3) to systems with different targets R_j is given, but our experiments indicate that this condition is not always correct, cf. the comments on Table 5 in Section 5). Still we will assume stability throughout the paper. In numerical experiments with the PSA instability can be detected by the occurrence of negative state probabilities.

3. The queue-length process

The random variable N_j will indicate the number of frames present at queue j in steady state, $j = 1, \dots, S$. Beside the vector of random variables $\mathbf{N} \doteq (N_1, \dots, N_S)$ a number of supplementary variables are needed to obtain a Markov process. The supplementary variable U_j will indicate the number of services which have been performed during the last completed visit to queue j , $j = 1, \dots, S$. The range of values of the vector $\mathbf{U} \doteq (U_1, \dots, U_S)$ is the product set

$$\mathcal{K} \doteq \bigotimes_{j=1}^S \{0, 1, \dots, K_j\}. \quad (3.1)$$

The supplementary variable H will indicate the queue to which the server is switching or to which the server is attending. The supplementary variable Z will indicate the action of the server. More precisely, $Z = 0$ will indicate that the server is switching and $Z = \kappa$ will indicate that the server is serving the κ th frame during the current visit. The supplementary variable Φ will indicate the actual phase of the current switchover time or service time. We will assume that the Markov process $(\mathbf{N}, \mathbf{U}, H, Z, \Phi)$ is stable, and denote the stationary state probabilities of this process by $p(\mathbf{n}, \mathbf{u}, h, \kappa, \phi)$, $\mathbf{n} \in \mathbb{N}^S$, $\mathbf{u} \in \mathcal{K}$, $h = 1, \dots, S$, $\kappa = 1, \dots, K_h$, $\phi = 1, \dots, \Omega_h$ if $\kappa = 0$, $\phi = 1, \dots, \Psi_h$ if $\kappa > 0$. In order to formulate the balance equations for this stationary Markov process we will use the indicator function $I_{\{C\}}$ taking the values 0 (if C is false) or 1 (if C is true), and the unit vectors \mathbf{e}_j , $j = 1, \dots, S$, in \mathbb{N}^S . The balance equations for the probabilities of states in which the server is switching are, for $\mathbf{n} \in \mathbb{N}^S$, $\mathbf{u} \in \mathcal{K}$, $h = 1, \dots, S$, $\phi = 1, \dots, \Omega_h$,

$$\begin{aligned} [\Lambda + \nu_h] p(\mathbf{n}, \mathbf{u}, h, 0, \phi) &= \sum_{j=1}^S \lambda_j I_{\{n_j \geq 1\}} p(\mathbf{n} - \mathbf{e}_j, \mathbf{u}, h, 0, \phi) + \nu_h I_{\{\phi \geq 2\}} p(\mathbf{n}, \mathbf{u}, h, 0, \phi - 1) \\ &+ \mu_{h-1} I_{\{u_{h-1} \geq 1, \phi = 1\}} \sum_{\kappa=0}^{K_{h-1}} I_{\{n_{h-1}=0 \vee u_{h-1}=K_{h-1} \vee T(\mathbf{u} + \kappa \mathbf{e}_{h-1}) \geq R_{h-1}\}} \end{aligned}$$

$$\begin{aligned} & \times I_{\{T(\mathbf{u}+(\kappa-1)\mathbf{e}_{h-1}) < R_{h-1}\}} p(\mathbf{n} + \mathbf{e}_{h-1}, \mathbf{u} + (\kappa - u_{h-1})\mathbf{e}_{h-1}, h-1, \kappa, \Psi_{h-1}) \\ & + v_{h-1} I_{\{u_{h-1}=0, \phi=1\}} \sum_{\kappa=0}^{K_{h-1}} I_{\{n_{h-1}=0 \vee T(\mathbf{u}+\kappa\mathbf{e}_{h-1}) \geq R_{h-1}\}} p(\mathbf{n}, \mathbf{u} + \kappa\mathbf{e}_{h-1}, h-1, 0, \Omega_{h-1}). \end{aligned} \quad (3.2)$$

The first term at the right-hand side stands for transitions caused by an arrival of a frame at one of the queues. The second term stands for a phase transition in the Erlang distributed switchover time. The third term describes a transition from a last service at queue $h-1$ to a switch to queue h ; such a transition can only occur if $u_{h-1} \geq 1$, indicating that at least one service has been performed during the last visit to queue $h-1$, if the token rotation timer at queue $h-1$ had not expired at the instant when the server was ready to start the u_{h-1} th service, and if either queue $h-1$ became empty or the service limit of queue $h-1$ had been reached or the token rotation timer at queue $h-1$ had expired after this u_{h-1} th service. The fourth term describes a transition from a switch to queue $h-1$ to a switch to queue h ; such a transition can only occur if $u_{h-1} = 0$, indicating that no service has been performed during the last visit to queue $h-1$, and if either queue $h-1$ was empty or the token rotation timer at queue $h-1$ had expired at the instant when the server completed its switch to this queue. Note that the third and fourth terms only contribute if the phase of the current switchover time is $\phi = 1$.

The balance equations for the probabilities of states in which the server is serving frames are, for $\mathbf{n} \in \mathbb{N}^S$, $\mathbf{u} \in \mathcal{K}$, $h = 1, \dots, S$, $n_h \geq 1$, $\kappa = 1, \dots, K_h$, $T(\mathbf{u} + (\kappa-1)\mathbf{e}_h) < R_h$, $\phi = 1, \dots, \Psi_h$,

$$\begin{aligned} [\Lambda + \mu_h] p(\mathbf{n}, \mathbf{u}, h, \kappa, \phi) &= \sum_{j=1}^S \lambda_j I_{\{n_j \geq 1\}} p(\mathbf{n} - \mathbf{e}_j, \mathbf{u}, h, \kappa, \phi) \\ &+ \mu_h I_{\{\phi \geq 2\}} p(\mathbf{n}, \mathbf{u}, h, \kappa, \phi-1) + v_h I_{\{\kappa=1, \phi=1\}} p(\mathbf{n}, \mathbf{u}, h, 0, \Omega_h) \\ &+ \mu_h I_{\{\kappa \geq 2, \phi=1\}} p(\mathbf{n} + \mathbf{e}_h, \mathbf{u}, h, \kappa-1, \Psi_h). \end{aligned} \quad (3.3)$$

The first term at the right-hand side stands for transitions caused by an arrival of a frame at one of the queues. The second term stands for a phase transition in the Erlang distributed service time. The third term describes a transition from a switch to queue h to the first service at queue h ($\kappa = 1$). The fourth term describes a transition from one service at queue h to another service at queue h ($\kappa \geq 2$). The last two types of transitions can only occur if the timer had not expired before the (new) service started, i.e., if $T(\mathbf{u} + (\kappa-1)\mathbf{e}_h) < R_h$, and if the phase of the current service time is $\phi = 1$. It should be noted that for all $\mathbf{n} \in \mathbb{N}^S$, $\mathbf{u} \in \mathcal{K}$, $h = 1, \dots, S$, $\kappa = 1, \dots, K_h$, $\phi = 1, \dots, \Psi_h$,

$$p(\mathbf{n}, \mathbf{u}, h, \kappa, \phi) = 0, \quad \text{if } n_h = 0, \text{ or } T(\mathbf{u} + (\kappa-1)\mathbf{e}_h) \geq R_h, \quad (3.4)$$

because the server cannot be serving a frame at a queue which is empty or at which the token rotation timer had already expired when the server was ready to start a (new) service. Finally, it holds by the law of total probability that

$$\sum_{n_1=0}^{\infty} \cdots \sum_{n_S=0}^{\infty} \sum_{u_1=0}^{K_1} \cdots \sum_{u_S=0}^{K_S} \sum_{h=1}^S \left[\sum_{\phi=1}^{\Omega_h} p(\mathbf{n}, \mathbf{u}, h, 0, \phi) + \sum_{\kappa=1}^{K_h} \sum_{\phi=1}^{\Psi_h} p(\mathbf{n}, \mathbf{u}, h, \kappa, \phi) \right] = 1. \quad (3.5)$$

4. The power-series algorithm

Before the recurrence relations of the PSA for the present model are derived, we introduce the following bilinear mapping of the interval $[0,1]$ onto itself:

$$\theta = \frac{(1+G)\rho}{1+G\rho}, \quad \rho = \frac{\theta}{1+G-G\theta}. \quad (4.1)$$

This mapping is needed to enlarge the radius of convergence of the power-series expansions and to avoid numerical instabilities. The choice of the parameter G depends on the model on hand. For the present type of models values in the order of $G = 1.5$ give good results. Next, we introduce power-series expansions of the state probabilities as functions of θ :

$$p(\mathbf{n}, \mathbf{u}, h, \kappa, \phi) = \theta^{|\mathbf{n}|} \sum_{k=0}^{\infty} \theta^k b(k; \mathbf{n}, \mathbf{u}, h, \kappa, \phi). \quad (4.2)$$

Here, we use the notation $|\mathbf{n}| \doteq n_1 + \dots + n_S$. In order to obtain a parametrization of the model as a function of θ we write $\lambda_j = a_j \rho = a_j \theta / (1+G-G\theta)$, $j = 1, \dots, S$, and $\Lambda = A\rho = A\theta / (1+G-G\theta)$, cf. (4.1). These expressions and the expansions (4.2) are substituted into the equations (3.2), (3.3) and (3.5) for the state probabilities. By equating coefficients of corresponding powers of θ on both sides of these equations one obtains relations for the coefficients of the power-series expansions of the state probabilities. The recurrence relations for the coefficients of the probabilities of states in which the server is switching are, for $k = 0, 1, 2, \dots$, $\mathbf{n} \in \mathbb{N}^S$, $\mathbf{u} \in \mathcal{K}$, $h = 1, \dots, S$, $\phi = 1, \dots, \Omega_h$,

$$\begin{aligned} & (1+G)v_h b(k; \mathbf{n}, \mathbf{u}, h, 0, \phi) \\ &= \sum_{j=1}^S a_j I_{\{n_j \geq 1\}} b(k; \mathbf{n} - \mathbf{e}_j, \mathbf{u}, h, 0, \phi) + [Gv_h - A] I_{\{k \geq 1\}} b(k-1; \mathbf{n}, \mathbf{u}, h, 0, \phi) \\ & \quad + v_h I_{\{\phi \geq 2\}} [(1+G)b(k; \mathbf{n}, \mathbf{u}, h, 0, \phi-1) - G I_{\{k \geq 1\}} b(k-1; \mathbf{n}, \mathbf{u}, h, 0, \phi-1)] \\ & \quad + \mu_{h-1} I_{\{u_{h-1} \geq 1, \phi=1\}} \sum_{\kappa=0}^{K_{h-1}} I_{\{n_{h-1}=0 \vee u_{h-1}=K_{h-1} \vee T(\mathbf{u}+\kappa \mathbf{e}_{h-1}) \geq R_{h-1}\}} I_{\{T(\mathbf{u}+(\kappa-1)\mathbf{e}_{h-1}) < R_{h-1}\}} \\ & \quad \times [(1+G) I_{\{k \geq 1\}} b(k-1; \mathbf{n} + \mathbf{e}_{h-1}, \mathbf{u} + (\kappa - u_{h-1})\mathbf{e}_{h-1}, h-1, \kappa, \Psi_{h-1}) \\ & \quad - G I_{\{k \geq 2\}} b(k-2; \mathbf{n} + \mathbf{e}_{h-1}, \mathbf{u} + (\kappa - u_{h-1})\mathbf{e}_{h-1}, h-1, \kappa, \Psi_{h-1})] \\ & \quad + v_{h-1} I_{\{u_{h-1}=0, \phi=1\}} \sum_{\kappa=0}^{K_{h-1}} I_{\{n_{h-1}=0 \vee T(\mathbf{u}+\kappa \mathbf{e}_{h-1}) \geq R_{h-1}\}} [(1+G)b(k; \mathbf{n}, \mathbf{u} + \kappa \mathbf{e}_{h-1}, h-1, 0, \Omega_{h-1}) \\ & \quad - G I_{\{k \geq 1\}} b(k-1; \mathbf{n}, \mathbf{u} + \kappa \mathbf{e}_{h-1}, h-1, 0, \Omega_{h-1})]. \end{aligned} \quad (4.3)$$

The recurrence relations for the coefficients of the probabilities of states in which the server is serving frames are, for $k = 0, 1, 2, \dots$, $\mathbf{n} \in \mathbb{N}^S$, $\mathbf{u} \in \mathcal{K}$, $h = 1, \dots, S$, $n_h \geq 1$, $\kappa = 1, \dots, K_h$, $T(\mathbf{u} + (\kappa-1)\mathbf{e}_h) < R_h$, $\phi = 1, \dots, \Psi_h$,

$$\begin{aligned} & (1+G)\mu_h b(k; \mathbf{n}, \mathbf{u}, h, \kappa, \phi) \\ &= \sum_{j=1}^S a_j I_{\{n_j \geq 1\}} b(k; \mathbf{n} - \mathbf{e}_j, \mathbf{u}, h, \kappa, \phi) + [G\mu_h - A] I_{\{k \geq 1\}} b(k-1; \mathbf{n}, \mathbf{u}, h, \kappa, \phi) \end{aligned}$$

$$\begin{aligned}
& + \mu_h I_{\{\phi \geq 2\}} [(1 + G)b(k; \mathbf{n}, \mathbf{u}, h, \kappa, \phi - 1) - GI_{\{k \geq 1\}} b(k - 1; \mathbf{n}, \mathbf{u}, h, \kappa, \phi - 1)] \\
& + \nu_h I_{\{\kappa = 1, \phi = 1\}} [(1 + G)b(k; \mathbf{n}, \mathbf{u}, h, 0, \Omega_h) - GI_{\{k \geq 1\}} b(k - 1; \mathbf{n}, \mathbf{u}, h, 0, \Omega_h)] \\
& + \mu_h I_{\{\kappa \geq 2, \phi = 1\}} [(1 + G)I_{\{k \geq 1\}} b(k - 1; \mathbf{n} + \mathbf{e}_h, \mathbf{u}, h, \kappa - 1, \Psi_h) \\
& - GI_{\{k \geq 2\}} b(k - 2; \mathbf{n} + \mathbf{e}_h, \mathbf{u}, h, \kappa - 1, \Psi_h)].
\end{aligned} \tag{4.4}$$

The law of total probability implies: for $k = 0, 1, 2, \dots$,

$$\begin{aligned}
\sum_{0 \leq |\mathbf{n}| \leq k} \cdots \sum_{u_1=0}^{K_1} \cdots \sum_{u_S=0}^{K_S} \sum_{h=1}^S \left[\sum_{\phi=1}^{\Omega_h} b(k - |\mathbf{n}|; \mathbf{n}, \mathbf{u}, h, 0, \phi) \right. \\
\left. + \sum_{\kappa=1}^{K_h} \sum_{\phi=1}^{\Psi_h} b(k - |\mathbf{n}|; \mathbf{n}, \mathbf{u}, h, \kappa, \phi) \right] = I_{\{k=0\}}.
\end{aligned} \tag{4.5}$$

The relations (4.3) and (4.4) can be used to compute the coefficients of the power-series expansions of the state probabilities in a mainly recursive manner when a suitable ordering of the states is adopted, cf. [7,8]. The only term which may prevent recursive computation is the term with $b(k; \mathbf{n}, \mathbf{u} + \kappa \mathbf{e}_{h-1}, h - 1, 0, 1)$ in (4.3). This term is only relevant if $u_{h-1} = 0$. This suggests that the coefficients should be computed, for fixed k and \mathbf{n} , in decreasing order of u_j , $j = 1, \dots, S$. In this way, only the term with $\kappa = 0$ may cause a problem. It is readily verified that the only case in which the coefficients cannot be computed recursively is the case $\mathbf{n} = \mathbf{0}$ and $\mathbf{u} = \mathbf{0}$; this is the only situation in which the server can make a complete cycle along the queues without a change in the values of \mathbf{N} and \mathbf{U} . In the case $\mathbf{n} = \mathbf{0}$ and $\mathbf{u} = \mathbf{0}$ Eq. (4.3) reduces to: for $k = 0, 1, 2, \dots$, $h = 1, \dots, S$, $\phi = 1, \dots, \Omega_h$,

$$\begin{aligned}
& (1 + G)\nu_h b(k; \mathbf{0}, \mathbf{0}, h, 0, \phi) \\
& = [G\nu_h - A]I_{\{k \geq 1\}} b(k - 1; \mathbf{0}, \mathbf{0}, h, 0, \phi) + \nu_{h-1} I_{\{\phi = 1\}} \sum_{\kappa=0}^{K_{h-1}} [(1 + G)b(k; \mathbf{0}, \kappa \mathbf{e}_{h-1}, h - 1, 0, \Omega_{h-1}) \\
& - GI_{\{k \geq 1\}} b(k - 1; \mathbf{0}, \kappa \mathbf{e}_{h-1}, h - 1, 0, \Omega_{h-1})] \\
& + \nu_h I_{\{\phi \geq 2\}} [(1 + G)b(k; \mathbf{0}, \mathbf{0}, h, 0, \phi - 1) - GI_{\{k \geq 1\}} b(k - 1; \mathbf{0}, \mathbf{0}, h, 0, \phi - 1)].
\end{aligned} \tag{4.6}$$

This forms, for each fixed k , $k = 0, 1, 2, \dots$, a dependent set of equations for the coefficients $b(k; \mathbf{0}, \mathbf{0}, h, 0, \phi)$, $h = 1, \dots, S$, $\phi = 1, \dots, \Omega_h$. Note that these sets of equations have a similar structure as those which have been encountered in cyclic polling models without token rotation timers, cf. [6,7], although the general form of the recursions (4.3) and (4.4) is quite different from that of the recursions for the latter models. The sets of Eqs. (4.6) can be solved together with (4.5). For the case $k = 0$ — which corresponds to the case $\rho = 0$ — we note that for $h = 1, \dots, S$, $\phi = 1, \dots, \Omega_h$,

$$b(0; \mathbf{0}, \mathbf{u}, h, 0, \phi) = 0, \text{ if } \mathbf{u} \neq \mathbf{0}, \tag{4.7}$$

because the components of \mathbf{U} will all vanish if there are no arrivals. In this case, (4.5) reduces to

$$\sum_{h=1}^S \sum_{\phi=1}^{\Omega_h} b(0; \mathbf{0}, \mathbf{0}, h, 0, \phi) = 1. \tag{4.8}$$

Table 1

The maximal number of terms with a storage capacity of 5 000 000 coefficients

K	1	1	1	2	2	2	2	2	3	3	4
Ψ	1	2	4	1	1	2	2	4	1	2	1
Ω	1	2	4	1	2	1	2	4	1	2	1
$S = 2$	788	556	392	427	370	330	301	212	277	195	197
$S = 3$	82	64	50	46	42	39	36	28	31	24	22
$S = 4$	27	22	18	16	16	14	13	10	10	7	7

For $k = 1, 2, \dots$, relation (4.5) can be rewritten as

$$\begin{aligned} \sum_{h=1}^S \sum_{\phi=1}^{\Omega_h} b(k; \mathbf{0}, \mathbf{0}, h, 0, \phi) = & - \sum_{u_1=0}^{K_1} \cdots \sum_{u_S=0}^{K_S} I_{\{\mathbf{u} \neq \mathbf{0}\}} \sum_{h=1}^S \sum_{\phi=1}^{\Omega_h} b(k; \mathbf{0}, \mathbf{u}, h, 0, \phi) \\ & - \sum_{1 \leq |\mathbf{n}| \leq k} \cdots \sum_{u_1=0}^{K_1} \cdots \sum_{u_S=0}^{K_S} \sum_{h=1}^S \left[\sum_{\phi=1}^{\Omega_h} b(k - |\mathbf{n}|; \mathbf{n}, \mathbf{u}, h, 0, \phi) + \sum_{\kappa=1}^{K_h} \sum_{\phi=1}^{\Psi_h} b(k - |\mathbf{n}|; \mathbf{n}, \mathbf{u}, h, \kappa, \phi) \right]. \end{aligned} \quad (4.9)$$

Hence, for each k , $k = 0, 1, 2, \dots$, one set of linear equations of size $\sum_{h=1}^S \Omega_h$ — all with the same determinant — has to be solved to obtain the coefficients for states with $\mathbf{n} = \mathbf{0}$ and $\mathbf{u} = \mathbf{0}$; and all other coefficients can be computed recursively. See Appendix C for a detailed computation scheme.

The number of coefficients which have to be computed in order to determine the power-series expansions up to the M th power of θ (or ρ) is

$$\binom{M+S+1}{S+1} \prod_{h=1}^S [1 + K_h] \sum_{j=1}^S (\Omega_j + K_j \Psi_j). \quad (4.10)$$

Note that quite some coefficients may vanish, cf. e.g. (3.4), (4.7). Moreover, it is not necessary to keep all computed coefficients in memory until the end of the execution of the algorithm if the coefficients of the power-series expansions of the desired performance measures are updated when those of the state probabilities are computed, cf. [8]. In Table 1 the number of terms of the power-series expansions is listed that can be computed with a given storage capacity, for systems with the same service limit K , the same number of phases of the service time distributions Ψ and the same number of phases of the switchover time distributions Ω for all queues.

Finally, it should be noted that the convergence of the power series can be accelerated with the aid of the so-called ϵ -algorithm, cf. e.g. [7,8]. The accuracy of the results obtained with the PSA can be estimated by inspection of the series produced with the aid of the ϵ -algorithm. The relative errors in the data to be presented in the next section are estimated to be in the order of 0.1% or (much) less. The correctness of the implementation of the PSA has been carefully checked by comparison with simulation experiments.

Table 2

Three-queue model: the influence of the service and switchover time distributions

Ψ	Ω	$E\{W_1\}$	$E\{W_2\}$	$E\{W_3\}$	$E\{W\}$	$\sigma\{W_1\}$	$\sigma\{W_2\}$	$\sigma\{W_3\}$	$\sigma\{W\}$
1	1	1.57	10.27	29.05	10.62	1.73	12.44	36.64	22.40
1	2	1.57	10.24	28.98	10.59	1.73	12.42	36.82	22.46
1	4	1.56	10.23	28.95	10.58	1.73	12.41	36.54	22.33
2	1	1.22	7.98	22.96	8.35	1.25	9.53	28.72	17.56
2	2	1.22	7.95	22.89	8.32	1.25	9.50	28.67	17.53
2	4	1.21	7.94	22.86	8.31	1.25	9.49	28.63	17.50
4	1	1.05	6.88	19.90	7.22	1.00	8.11	24.74	15.14
4	2	1.04	6.85	19.82	7.19	1.00	8.09	24.65	15.09
4	4	1.04	6.84	19.79	7.18	1.00	8.08	24.63	15.07

Table 3

Three-queue model: estimating performance measures for constant times

Est.	Ψ, Ω	$E\{W_1\}$	$E\{W_2\}$	$E\{W_3\}$	$E\{W\}$	$\sigma\{W_1\}$	$\sigma\{W_2\}$	$\sigma\{W_3\}$	$\sigma\{W\}$
1,2	4	1.04	6.80	19.81	7.17	1.00	8.04	24.71	15.10
1,2	∞	0.86	5.64	16.73	6.02	0.76	6.57	20.74	12.67
2,4	∞	0.86	5.73	16.69	6.02	0.75	6.66	21.13	12.83
Sim.	∞	0.86	5.73	16.32	5.94	0.76	6.66	19.83	12.22

5. Examples

Once the moments of the joint queue length distribution have been computed those of the waiting time distributions can be determined in the usual manner for polling systems with Poisson arrival streams, cf. e.g. [7]. In the examples below W_j denotes the waiting time, without service time, at queue j , $j = 1, \dots, S$, and W denotes the waiting time, without service time, of an arbitrary frame. In all examples the mean switchover times between the queues are chosen to be equal, i.e. $\delta_j = \Delta/S$, $j = 1, \dots, S$. The influence of individual switchover times on performance measures is usually limited, cf. e.g. [6]. The most important characteristics of the switchover times are the first two moments of the total switchover time of the server during a cycle along the stations.

In the first example we study the effects of the number of phases of the Erlang service time distributions and switchover time distributions on the waiting time characteristics. To this end, we consider a three-queue system with mean service times $\beta_j = 1.0$, $j = 1, 2, 3$, and with total mean switchover time $\Delta = 0.15$. Queue 1 is a station with arrival rate $\lambda_1 = 0.4$ and with synchronous traffic. The service discipline at this station is $K_1 = 2$, $R_1 = \infty$. Queues 2 and 3 are stations with asynchronous traffic, and identical characteristics. The arrival rates are $\lambda_2 = \lambda_3 = 0.2$ and the service disciplines at these stations are $R_2 = R_3 = 1.0$, $K_2 = K_3 = \infty$. Due to the target rotation times $R_2 = R_3 = 1.0$ these stations cannot send more than one frame per cycle, so that the effective service limits are $\hat{K}_2 = \hat{K}_3 = 1$. The offered load to this system is readily seen to be $\rho = 0.8$. Table 2 shows the means and the standard deviations of the waiting times for this system. In each instance in this table all service time distributions

consist of the same number of phases, and all switchover time distributions consist of the same number of phases. These numbers of phases are indicated in the first two columns. It is seen that the influence of the number of phases of the switchover time distributions is only minor, but that the influence of the number of phases of the service time distributions is important. However, the latter seems to affect mainly the absolute values, and not so much the relative values, of the waiting time characteristics. Moreover, the waiting time characteristics turn out to be almost linear functions of the reciprocal of the number of phases of the Erlang service time distributions. If we compute for each performance measure the difference between its value in the case of exponential distributions (row 1 in Table 2) and that in the case of Erlang-2 distributions (row 5 in Table 2), and subtract the half of this difference from its value in the case of Erlang-2 distributions we obtain very good approximations for the values in the case of Erlang-4 distributions (the results are displayed in the first row of Table 3; compare these values with those in the last row in Table 2). If we subtract the full above mentioned differences from the values in the case of Erlang-2 distributions we obtain approximations for the case of constant service and switchover times. These results are given in the second row of Table 3. The third row of this table contains a similar estimate for the deterministic case, but obtained from linear extrapolation from the Erlang-2 and Erlang-4 cases (row 5 and row 9 in Table 2; the first column of Table 3 indicates the source of the estimation). The last row of this table contains simulation results for our model, but with constant service and switchover times. These estimates have been obtained from a simulation run of 8,000,000 time units (service times). The relative widths of the 95% confidence intervals are less than 1% for station 1, about 2% for station 2 and still in the order of 5% for station 3. Finally, note the large differences between the waiting time characteristics of queues 2 and 3 which have the same parameter values. These differences are due both to the different positions of these stations with respect to station 1 (the synchronous station) and to the fact that queues 2 and 3 manage asynchronous traffic. In any cycle, stations 2 and 3 can only transmit if station 1 does not have any frame to transmit, due to the values of R_2 and R_3 . Furthermore, station 3 can only transmit if station 2 does not have any frame to transmit. These observations explain the larger average waiting times experienced by frames at station 3 compared to that experienced by frames at station 2.

Consider next a system with $S = 2$ queues, Erlang E_4 service times and Erlang E_2 switchover times. The arrival rate at the first queue is four times as high as that at the second queue, i.e. $\lambda_1 = 4\lambda_2$. The mean service times are $\beta_1 = \beta_2 = 1.0$. Note that the token rotation timer $T(\mathbf{v})$ will only take values that are multiples of 1.0 for this model. Table 4 shows the means and the standard deviations of the waiting time distributions for this model as function of the target token rotation times R_1 and R_2 , while the service limits K_1 and K_2 are chosen such that they do not influence the performance of the system (i.e., $\hat{K}_j = R_j$, $j = 1, 2$). Both queues represent stations with asynchronous traffic, but possibly with different priorities. The first five entries of the table concern cases in which the stations have the same priority ($R_1 = R_2$). Note that increasing the TTRT leads to decreasing mean waiting times at station 1, while $R_1 = R_2 = 2.0$ yields a minimal mean waiting time at station 2. The larger value of $E\{W_1\}$ compared to that of $E\{W_2\}$ can be explained by the fact that $\lambda_1 = 4\lambda_2$. The other entries concern cases in which the stations have different priorities ($R_1 \neq R_2$). By properly adjusting the values of R_1 and R_2 it is possible to have $E\{W_1\} < E\{W_2\}$ although $\lambda_1 > \lambda_2$. See, e.g., the entries with $(R_1 = 2.0, R_2 = 1.0)$, $(R_1 = 3.0, R_2 = 1.0)$ and $(R_1 = 3.0, R_2 = 2.0)$. Also note that increasing the TTRT of some station may lead to smaller mean waiting times at other stations of which the TTRT is kept fixed; see for example the entries with $R_1 = 1.0$ fixed and R_2 increasing: here, $E\{W_1\}$ is decreasing.

Table 5 concerns similar quantities as Table 4, the difference being a larger total switchover time

Table 4

Two-queue model with $K_1 = K_2 = \infty$, $\rho = 0.8$ and $\Delta = 0.1$

R_1	R_2	$E\{W_1\}$	$E\{W_2\}$	$E\{W\}$	$\sigma\{W_1\}$	$\sigma\{W_2\}$	$\sigma\{W\}$
1.0	1.0	14.24	1.36	11.67	15.07	1.45	14.44
2.0	2.0	5.31	1.26	4.50	6.09	1.30	5.71
3.0	3.0	4.18	1.40	3.63	4.79	1.38	4.62
4.0	4.0	3.71	1.59	3.29	4.52	1.54	4.19
6.0	6.0	3.26	1.96	3.00	4.09	1.95	3.80
2.0	1.0	2.14	17.05	5.12	2.21	23.79	12.36
3.0	1.0	1.89	12.89	4.09	2.02	18.44	9.52
3.0	2.0	3.39	4.76	3.66	3.99	6.86	4.74
4.0	3.0	3.48	2.57	3.30	4.24	3.07	4.05
4.0	2.0	2.05	8.92	3.42	2.16	14.44	7.28
4.0	1.0	1.77	11.42	3.70	1.93	16.42	8.47
1.0	2.0	11.68	0.72	9.49	12.39	0.67	11.92
1.0	4.0	10.97	0.69	8.92	11.65	0.63	11.20

Table 5

Two-queue model with $K_1 = K_2 = \infty$, $\rho = 0.7$ and $\Delta = 0.4$

R_1	R_2	$E\{W_1\}$	$E\{W_2\}$	$E\{W\}$	$\sigma\{W_1\}$	$\sigma\{W_2\}$	$\sigma\{W\}$
2.0	2.0	18.74	1.72	15.34	19.62	1.60	18.84
3.0	3.0	5.63	1.65	4.83	6.41	1.51	5.99
4.0	4.0	3.82	1.73	3.40	4.56	1.61	4.22
6.0	6.0	2.70	1.96	2.55	3.37	1.96	3.15
3.0	1.0	2.41	36.02	9.13	2.44	43.71	23.82
3.0	2.0	4.58	6.79	5.03	5.22	9.34	6.33
4.0	3.0	3.59	2.82	3.44	4.30	3.25	4.12
4.0	2.0	2.24	10.62	3.92	2.30	16.52	8.37
4.0	1.0	1.96	20.22	5.61	2.06	25.59	13.70
2.0	3.0	14.24	1.20	11.63	14.99	0.99	14.40
2.0	4.0	12.84	1.15	10.50	13.56	0.94	13.00

and a smaller offered load. In this case, the values $R_1 = R_2 = 1.0$ and the values $R_1 = 2.0$, $R_2 = 1.0$ correspond to unstable systems. This example reveals a remarkable difference concerning stability between polling systems with and without rotation time restrictions. In the latter systems it is necessary to increase the service limit of the bottleneck station to prevent instability, cf. (2.2). Here, we observe a system which is not stable for $R_1 = 2.0$, $R_2 = 1.0$, but which is stable for $R_1 = 3.0$, $R_2 = 1.0$, while station 2 forms the bottleneck. Note that the foregoing observations are in contradiction with the stability condition given in [23], formula (13). The latter condition applied to the present example would give a stability condition $\rho < 1/1.48 \approx 0.676$ for both cases $R_1 = 2.0$, $R_2 = 1.0$, and $R_1 = 3.0$, $R_2 = 1.0$. Experiments with the PSA indicate that the former system is still stable for $\rho = 0.697$, while the latter is still stable for $\rho = 0.750$. These properties have been confirmed by simulation of completely deterministic systems.

Table 6

Symmetrical models with $\rho = 0.75$, $\Delta = 0.12$

K	R	$P\{L=0\}$	$E\{L\}$	$\sigma\{L\}$	$P\{N=0\}$	$E\{N\}$	$\sigma\{N\}$	$E\{W\}$	$\sigma\{W\}$
1	∞	0.203	3.58	3.94	0.470	1.19	1.72	3.77	5.24
1	3.0	0.198	3.72	4.09	0.466	1.24	1.80	3.96	5.59
1	2.0	0.183	4.11	4.49	0.444	1.37	1.96	4.48	6.22
1	1.0	0.129	6.10	6.41	0.356	2.03	2.76	7.13	9.39
2	∞	0.219	3.31	3.69	0.488	1.10	1.61	3.41	4.79
2	6.0	0.217	3.36	3.75	0.488	1.12	1.65	3.48	4.99
2	4.0	0.211	3.50	3.90	0.479	1.17	1.71	3.67	5.22
2	3.0	0.203	3.67	4.07	0.467	1.22	1.78	3.89	5.47
2	2.0	0.187	4.07	4.47	0.444	1.36	1.94	4.42	6.10
3	∞	0.223	3.23	3.61	0.493	1.08	1.58	3.30	4.65
3	10.0	0.223	3.24	3.64	0.493	1.08	1.45	3.32	3.94
3	8.0	0.221	3.28	3.68	0.493	1.09	1.46	3.37	3.97
3	6.0	0.219	3.34	3.46	0.489	1.11	1.64	3.46	4.94
3	4.0	0.212	3.49	3.90	0.478	1.16	1.71	3.65	5.20
3	3.0	0.204	3.66	4.07	0.467	1.22	1.76	3.88	5.39

Table 7

Four-queue models with offered load $\rho = 0.75$

K_1	R_2	$\lambda_1 = \lambda_2$			$\lambda_1 = 2\lambda_2$			$\lambda_1 = \frac{1}{2}\lambda_2$		
		$E\{W_1\}$	$E\{W_2\}$	$E\{W\}$	$E\{W_1\}$	$E\{W_2\}$	$E\{W\}$	$E\{W_1\}$	$E\{W_2\}$	$E\{W\}$
1	1.0	1.41	9.22	5.32	1.80	10.83	4.81	1.16	8.36	5.96
2	1.0	1.35	8.43	4.88	1.68	9.44	4.27	1.13	7.94	5.67
1	2.0	1.86	6.39	4.12	2.33	5.20	4.05	1.54	5.57	4.23
2	2.0	1.54	6.36	3.95	1.83	7.65	3.78	1.36	5.52	4.13

Table 6 concerns symmetrical systems with three stations, and exponential service and switchover times. The mean service times are $\beta_j = 1.0$, $j = 1, 2, 3$. This table captures the effects of the service disciplines (i.e., the K and R values) on the performance measures in the case in which each station has a service limit as well as a target rotation time. These limits and targets are the same for each station in this symmetrical system. In the table, N stands for the number of frames present in an individual queue, and $L \doteq \sum_{j=1}^S N_j$ indicates the total number of frames present in the system, both including the frame in service if any. From the table it can be observed that the mean waiting time decreases with increasing K when all stations are synchronous ($R = \infty$). Furthermore, when all stations are asynchronous the mean waiting time increases with decreasing R for any given K . This behaviour can easily be justified by taking into account the way traffic is managed by each station. Finally, note that the mean waiting time in the first row of this table can be computed exactly from the well-know pseudo-conservation law for polling systems with 1-limited service at each station. This yields $E\{W\} = 83/22 \approx 3.7727$.

The next examples concern systems with $S = 4$ queues and with exponential service and switchover time distributions. The mean service times are $\beta_j = 1.0$, $j = 1, \dots, 4$. The total mean switchover time

Table 8

Three-queue model with $\lambda_1 = 2\lambda_2 = 2\lambda_3$, and $\beta_1 = 1.0$, $\beta_2 = \beta_3 = 0.5$

ρ	$E\{W_1\}$	$E\{W_2\}$	$E\{W_3\}$	$E\{W\}$	$\sigma\{W_1\}$	$\sigma\{W_2\}$	$\sigma\{W_3\}$	$\sigma\{W\}$
0.1	0.16/0.16	0.15/0.15	0.15/0.15	0.15/0.16	0.29/0.31	0.25/0.27	0.25/0.27	0.27/0.29
0.3	0.39/0.41	0.31/0.33	0.32/0.33	0.36/0.37	0.67/0.71	0.49/0.54	0.51/0.54	0.59/0.63
0.4	0.59/0.62	0.43/0.46	0.44/0.47	0.51/0.54	0.93/0.99	0.61/0.71	0.66/0.72	0.80/0.87
0.5	0.90/0.96	0.57/0.65	0.61/0.66	0.75/0.80	1.32/1.40	0.76/0.94	0.86/0.96	1.11/1.20
0.6	1.44/1.55	0.76/0.91	0.84/0.93	1.12/1.24	1.95/2.11	0.95/1.28	1.14/1.31	1.60/1.78
0.7	2.58/2.81	1.02/1.35	1.19/1.32	1.84/2.07	3.23/3.49	1.20/1.79	1.53/1.74	2.59/2.86
0.8	6.27/7.61	1.41/2.12	1.77/2.06	3.93/4.86	7.18/8.67	1.57/2.61	2.11/2.46	5.75/6.97

Table 9

Three-queue model with $\lambda_1 = \lambda_2 = \lambda_3$, and $\beta_1 = \beta_2 = \beta_3 = 1.0$

ρ	$E\{W\}$	$\sigma\{W\}$	$P\{N>0\}$	$P\{N>1\}$	$P\{N>2\}$	$P\{N>3\}$	$P\{N>4\}$	$P\{N>5\}$	$P\{N>6\}$
0.1	0.16/0.17	0.32/0.36	0.0379	9.1e-4	1.8e-5	3.3e-7	5.9e-9	1.1e-10	2.0e-12
0.3	0.42/0.47	0.74/0.83	0.1289	0.0118	9.9e-4	8.5e-5	7.7e-6	7.2e-7	7.0e-8
0.4	0.62/0.71	1.03/1.19	0.1859	0.0261	0.0035	5.0e-4	7.5e-5	1.2e-5	1.8e-6
0.5	0.93/1.08	1.44/1.68	0.2549	0.0523	0.0108	0.0023	5.3e-4	1.2e-4	3.0e-5
0.6	1.43/1.72	2.09/2.55	0.3415	0.1000	0.0301	0.0096	0.0032	0.0011	3.6e-4
0.7	2.42/3.11	3.35/4.29	0.4560	0.1893	0.0819	0.0371	0.0170	0.0079	0.0036
0.8	5.23/8.27	6.87/10.6	0.5298	0.2628	0.1357	0.0730	0.0390	0.0208	0.0108

is $\Delta = 0.12$. Queues 1 and 3 are stations with the same characteristics and with synchronous traffic ($\lambda_1 = \lambda_3$, $R_1 = R_3 = \infty$, $K_1 = K_3$). Queues 2 and 4 are also stations with the same characteristics, but with asynchronous traffic ($\lambda_2 = \lambda_4$, $R_2 = R_4$, $K_2 = K_4 = \infty$). Table 7 shows the mean waiting times for this system for various service disciplines and for various proportions between the arrival rates of the synchronous and the asynchronous traffic. In all cases, the total offered load is the same, and equals $\rho = 0.75$. Due to the symmetry, $E\{W_1\} = E\{W_3\}$ and $E\{W_2\} = E\{W_4\}$ in all instances. Note that for all three considered ratios between the arrival rates the transition from $K_1 = K_3 = 1$ to $K_1 = K_3 = 2$, with $R_2 = R_4 = 1.0$ fixed, is advantageous for all stations. However, the transition from $K_1 = K_3 = 1$ to $K_1 = K_3 = 2$, with $R_2 = R_4 = 2.0$ fixed, is only advantageous for all stations in the cases $\lambda_1 = \lambda_2$ and $\lambda_1 = \frac{1}{2}\lambda_2$; in the case $\lambda_1 = 2\lambda_2$ this transition is not advantageous for stations 2 and 4. This can be justified by the fact that increasing the synchronous arrival rate over the asynchronous one penalizes the asynchronous stations since their transmissions are subject to time constraints. Increasing R_2 and R_4 with $K_1 = K_3$ fixed, is advantageous for stations 2 and 4, and disadvantageous for stations 1 and 3, in all examples.

The next examples concern systems with $S = 3$ queues and with Erlang E_4 service and switchover time distributions. Tables 8 and 9 concern systems with fixed service limits $K_1 = K_2 = K_3 = 1$, fixed target token rotation times $R_1 = R_2 = R_3 = 2.0$, and with $\Delta = 0.15$. They show waiting time characteristics as function of the offered load ρ , for two combinations of the arrival rates and the mean service times. For each performance measure two values are listed. The left values are the results of computations with the PSA for the model described in Section 2. The right values are simulation results

Table 10

Model with $\lambda_1 = \lambda_2 = \lambda_3$, $K_1 = 1$, $K_2 = K_3 = \infty$, $R_1 = \infty$, $R_2 = R_3 = 2.0$

ρ	$E\{W_1\}$	$E\{W_2\}$	$E\{W_3\}$	$E\{W\}$	$\sigma\{W_1\}$	$\sigma\{W_2\}$	$\sigma\{W_3\}$	$\sigma\{W\}$
0.1	0.16/0.16	0.16/0.17	0.16/0.18	0.16/0.17	0.32/0.29	0.30/0.35	0.31/0.37	0.31/0.34
0.3	0.39/0.35	0.40/0.48	0.41/0.51	0.40/0.45	0.65/0.53	0.72/0.90	0.76/0.96	0.71/0.82
0.4	0.54/0.46	0.60/0.78	0.64/0.79	0.60/0.67	0.81/0.62	1.02/1.39	1.11/1.38	0.99/1.20
0.5	0.73/0.58	0.92/1.22	1.00/1.29	0.88/1.03	0.98/0.70	1.46/2.01	1.65/2.08	1.40/1.75
0.6	0.96/0.72	1.45/2.08	1.67/2.15	1.36/1.65	1.17/0.79	2.16/3.18	2.58/3.17	2.08/2.72
0.7	1.25/0.88	2.48/3.99	3.11/4.15	2.28/3.00	1.38/0.85	3.46/5.52	4.49/5.73	3.46/4.86
0.8	1.63/1.06	5.09/12.9	7.88/11.2	4.85/8.37	1.64/0.92	6.54/15.8	10.4/13.6	7.63/13.2

Table 11

Model with $\lambda_1 = 2\lambda_2 = 2\lambda_3$, $K_1 = 1$, $K_2 = K_3 = \infty$, $R_1 = \infty$, $R_2 = R_3 = 2.0$

ρ	$E\{W_1\}$	$E\{W_2\}$	$E\{W_3\}$	$E\{W\}$	$\sigma\{W_1\}$	$\sigma\{W_2\}$	$\sigma\{W_3\}$	$\sigma\{W\}$
0.1	0.17/0.16	0.16/0.18	0.16/0.18	0.16/0.17	0.32/0.30	0.30/0.37	0.30/0.39	0.31/0.34
0.3	0.41/0.37	0.39/0.50	0.40/0.53	0.41/0.44	0.69/0.57	0.71/0.96	0.74/1.05	0.71/0.82
0.4	0.60/0.50	0.59/0.80	0.61/0.86	0.60/0.67	0.91/0.70	1.02/1.46	1.08/1.59	0.98/1.20
0.5	0.84/0.66	0.90/1.31	0.96/1.43	0.89/1.02	1.17/0.82	1.48/2.25	1.61/2.42	1.37/1.79
0.6	1.20/0.84	1.44/2.31	1.58/2.42	1.36/1.60	1.54/0.96	2.26/3.67	2.52/3.72	2.02/2.80
0.7	1.77/1.08	2.56/4.82	2.94/4.86	2.26/2.96	2.10/1.09	3.83/7.02	4.43/6.70	3.32/5.27
0.8	2.81/1.37	5.84/14.9	7.39/13.2	4.72/7.70	3.10/1.25	8.10/19.2	10.0/15.9	7.08/14.1

for a corresponding system with an FDDI protocol, cf. Appendix B. The simulations have been carried out with constant service and switchover times, and the accumulated lateness is dealt with as described in Appendix B. The relative widths of the 95% confidence intervals are in the order of 10% or less. In the completely symmetrical case (Table 9) also some excess probabilities for the number of frames present at a station are displayed for the PSA-model.

In the final examples, the model of the previous examples is considered, but now there is a service limit for queue 1 (synchronous traffic), while there are target token rotation times for queue 2 and 3 (asynchronous traffic). In both cases, $\beta_j = 1.0$, $j = 1, 2, 3$, $K_1 = 1$ and $R_2 = R_3 = 2.0$. In Table 10 the arrival rates are equal, while in Table 11, $\lambda_1 = 2\lambda_2 = 2\lambda_3$. Also in these tables computations with the PSA for the model of the token bus are compared with simulations for comparable systems with an FDDI protocol.

When only non time-critical traffic is managed by all the stations, performance figures of the token bus and FDDI are very close up to very high offered load (approx. 70%). On the other hand, when there is at least one station managing time-critical traffic the agreement is poor. The above facts are obviously due to the different ways the accumulated lateness is managed by the token bus and FDDI MAC protocols.

The computations with the PSA for the models considered in Tables 8–11 took each about 15 minutes of CPU time on a workstation to determine 40 terms of the power-series expansions. Note that only one run of the PSA is needed to produce performance measures for various values of the offered load ρ for a given system configuration.

6. Conclusions

In this paper we have proposed a general model for communication systems with timed token access protocols, including the IEEE 802.4 token bus protocol, and solved it with the aid of the Power-Series Algorithm. The model can readily be modified to include finite buffer spaces; this modification only requires the addition of a few indicator functions in the balance Eqs. (3.2), (3.3), and hence in the recurrence relations (4.3), (4.4). Comparison of performance measures computed with our model with those obtained by simulating systems with an FDDI protocol has revealed that these values diverge with increasing load. Further improvement of the model could be achieved by including the accumulated lateness, and also by using more general (Markovian) arrival processes. Both these extensions of the model go at the cost of larger supplementary spaces than indicated in (4.10). As an alternative to the approximation of the rotation timers by their expected values, cf. (2.1), one could approximate the rotation timers by Erlang distributed timers. An important subject for further research is the determination of the stability conditions for the proposed model.

Appendix A. IEEE 802.4 Token Bus

The IEEE 802.4 Token Bus standard, cf. [1], specifies a token passing protocol on a bus with an optional priority mechanism. Specifically, this MAC protocol identifies four priority classes, denoted access classes, termed 0, 2, 4, and 6, with 6 being the highest priority and 0 the lowest. Each access class acts as a virtual substation in that the token is passed, internally, from the highest access class downward, in the order 4, 2, 0. A time parameter, denoted `hi_pri_token_hold_time`, is assigned to class 6, whereas each of the other three classes is assigned a parameter, called “target” token rotation time (abbreviated as $TTRT_i$, $i = 4, 2, 0$). Hence, for the IEEE 802.4 standard, there can be at most three different values for the $TTRT$ parameters. Each station using the optional priority scheme will have three rotation timers for the three lower access classes, and each access class has its own queue for frames to be transmitted. When a station receives the token, it is guaranteed to transmit data frames of class 6 until either the station becomes empty, or a period of time equal to `hi_pri_token_hold_time` has elapsed, whichever comes first. For each of the other access classes, the corresponding virtual substation measures the time it takes the token to circulate around the logical ring. If the token returns in less than $TTRT$, then the substation transmits frames of that class until such frames are transmitted or $TTRT$ has elapsed, whichever comes first. Otherwise, if the token returns later than $TTRT$ the station cannot send frames of that priority on this pass of the token, and forwards the token immediately. Hence, priority 6 class supports the time-constraint service whereas priorities 4, 2, and 0 support non time-constraint services. Obviously, if the total transmission time of class 6 data frames in a token cycle exceeds all the $TTRTs$, then no lower class frames can be transmitted at all. The aim of the cycle-dependent timing mechanism is that, as the aggregate offered load of class 6 traffic decreases, lower classes are allowed to access the channel successively starting from the access class with the largest $TTRT$ down to the one with the smallest $TTRT$.

The access class service algorithm consists of loading the residual value (target token rotation time minus the contents of the token rotation timer for the corresponding access class) from the token rotation timer into a token hold timer, and resetting the same token rotation timer. The main difference between the IEEE 802.4 token bus and the FDDI MAC protocols is the management of negative residual values or accumulated latency as it is called in FDDI (see Appendix B); the latter MAC protocol takes this into

account, whereas the former MAC protocol does not. In other words, the IEEE 802.4 standard loses memory of this accumulated latency by resetting the proper token rotation timer, whereas FDDI keeps track of the accumulated latency (by setting to 1 the *Late_Ct* counter; see Appendix B) until it has been recovered because, for example, one station does not transmit time-critical frames in some cycle.

Appendix B. FDDI (Fibre Distributed Data Interface)

FDDI, standardized by the American National Standards Institute X3T9 committee (e.g., [2,3]) is based on a dual fiber optic ring. To provide guaranteed service to time-critical (synchronous) traffic, FDDI enforces a limitation on how much synchronous traffic each node can send per token received. Specifically, a Target Token Rotation Time (TTRT) is negotiated among stations during ring initialization and whenever a station captures the token it can transmit synchronous data up to a maximum duration of $T_{ST} \doteq (TTRT - \alpha)/S$, where S is the number of active stations while α is a constant term which takes into account the maximum ring latency, the maximum frame length, and the time it takes to transmit a token. Hence, priority 6 service plays in the Token Bus the same role as the synchronous service in FDDI. To compute the maximum time a station can transmit non time-critical traffic (asynchronous) data when it captures a token, two timers are used: the Token Rotation Timer (TRT) and the Token Holding Timer (THT). TRT measures the time between the receipt of two consecutive tokens while THT is used to limit the transmission of a station when a token is captured. If TRT reaches TTRT before the token returns to the station, a variable, named *Late_Ct*, is set to 1 and TRT is restarted. When the token arrives at a station with *Late_Ct*=1 the token is called a *late token*. Whenever a late token is captured, only synchronous transmissions are enabled, TRT is not restarted, and *Late_Ct* is set to 0. TRT is left running to count both the amount of time by which the token arrived late (accumulated lateness) plus the next rotation time of the token. On the other hand, if the token arrives before TRT reaches TTRT and *Late_Ct* is 0, the token is named an *early token*. Whenever an early token is captured, the current value of TRT is stored in the THT, TRT is reset to time the next rotation of the token and synchronous frames are transmitted for a time up to T_{ST} . After synchronous transmission, THT is enabled and asynchronous transmissions start (THT is disabled during Synchronous frame transmissions). The difference between TTRT and the content of THT is the maximum time available for asynchronous transmissions in this cycle. Any unused time remaining in THT at the end of asynchronous frame transmissions is lost; it cannot be retained until the next token arrives. A station may initiate a transmission of an asynchronous frame if the timer THT has not reached the TTRT threshold. This may cause an additional delay in the release of the token (hereafter called asynchronous overrun) since the transmission of an asynchronous frame is always completed. The asynchronous overrun is bounded by the time for the transmission of a frame of maximum length. Multiple levels of asynchronous priorities may be distinguished by a station. For each priority level n , a threshold value ($T_Pri(n)$) is defined. $T_Pri(n)$ are an ordered sequence of values in the range $[0, TTRT]$, higher priorities have higher T_Pri values and the highest priority has a threshold which is equal to TTRT. Asynchronous transmissions start from the highest priority. Asynchronous frames of priority n may only be transmitted if THT is less than $T_Pri(n)$. If multiple asynchronous priority levels are not implemented, all asynchronous frames have a threshold value which is equal to TTRT.

It has been formally proved [14,20] that FDDI guarantees upper bounds for mean and maximum cycle times, e.g., the average token rotation time does not exceed TTRT, and the maximum token rotation time does not exceed twice the TTRT.

Appendix C. The computation scheme

The computation scheme for calculating the coefficients of the power-series expansions of the state probabilities up to the M th power of θ reads:

```

for  $k = 0, \mathbf{n} = \mathbf{0}, \mathbf{u} = \mathbf{0}$  do solve set of Eqs. (4.6), (4.8) with  $k = 0$ .
for  $m = 1$  to  $M$  do
  for  $k = 0$  to  $m$  do
    for all  $\mathbf{n}$  with  $|\mathbf{n}| = m - k$  do
      for  $u_1 = K_1$  downto  $0$  do
        .....
      for  $u_s = K_s$  downto  $0$  do
        if  $\mathbf{n} \neq \mathbf{0}$  or  $\mathbf{u} \neq \mathbf{0}$  then
          if  $\mathbf{u} = \mathbf{0}$  then determine an  $h = h(\mathbf{n}, \mathbf{u})$  such that  $n_h \geq 1$ 
            else determine an  $h = h(\mathbf{n}, \mathbf{u})$  such that  $u_h \geq 1$ ,
          for  $h = h(\mathbf{n}, \mathbf{u}) + 1, \dots, s, 1, \dots, h(\mathbf{n}, \mathbf{u})$  do
            for  $\phi = \Omega_h$  downto  $1$  do compute  $b(k; \mathbf{n}, \mathbf{u}, h, 0, \phi)$  according to (4.3),
            for  $\kappa = 1$  to  $K_h$  do
              for  $\phi = \Psi_h$  downto  $1$  do
                compute  $b(k; \mathbf{n}, \mathbf{u}, h, \kappa, \phi)$  according to (4.4),
            else solve set of Eqs. (4.6), (4.9) with  $k = m$ .

```

References

- [1] ANSI/IEEE Standard 802.4-1985 Token-Passing Bus Access Method and Physical Layer Specifications, The Institute of Electrical and Electronic Engineers, Inc., New York (1985).
- [2] FDDI Token Ring Media Access Control, ANSI Standard X3.139 (1987).
- [3] FDDI Hybrid Ring Control, Draft Proposed American National Standard ANSI X3T9.5, Rev. 6, May 11 (1990).
- [4] E. Altman, Analysing timed-token ring protocols using the power-series algorithm, in: J. Labetoulle and J.W. Roberts (eds.), *The Fundamental Role of Teletraffic in the Evolution of Telecommunications Networks*, Elsevier, Amsterdam (1994) pp. 961–971.
- [5] S. Ayandeh, Performance evaluation of the IEEE 802.4 token bus protocol for distributed real-time applications, in: *Proc. IEEE INFOCOM*, San Francisco, Ca. (1990) pp. 1102–1110.
- [6] J.P.C. Blanc, The power-series algorithm applied to cyclic polling systems, *Commun. Statist.-Stochastic Models* **7** (1991) 527–545.
- [7] J.P.C. Blanc, Performance evaluation of polling systems by means of the power-series algorithm, *Annals Oper. Res.* **35** (1992) 155–186.
- [8] J.P.C. Blanc, Performance analysis and optimization with the power-series algorithm, in: L. Donatiello and R. Nelson (eds.), *Performance Evaluation of Computer and Communication Systems*, Lecture Notes in Computer Science **729** (Springer, Berlin, 1993) 53–80.
- [9] M. Conti, E. Gregori and L. Lenzini, Metropolitan Area Networks (MANs): protocols, modeling and performance evaluation, in: L. Donatiello and R. Nelson (eds.), *Performance Evaluation of Computer and Communication Systems*, Lecture Notes in Computer Science **729** (Springer, Berlin, 1993) pp. 81–120.
- [10] D. Dykeman and W. Bux, Analysis and tuning of the FDDI media access control protocol, *IEEE J. Selec. Areas in Commun.* **6** (1988) 997–1010.
- [11] R. M. Grow, A timed token protocol for local area networks, presented at Electro 82, *Token Access Protocols* (17/3) (May 1982).
- [12] A.P. Jayasumana, Throughput analysis of the IEEE 802.4 priority scheme, *IEEE Trans. on Commun.* **37** (1989) 565–571.

- [13] Y.C. Jenq, Approximations for packetized voice traffic in statistical multiplexer, in: *Proc. IEEE INFOCOM*, San Francisco, Ca. (1984) 256–259.
- [14] M.J. Johnson, Proof that timing requirements of the FDDI token ring protocol are satisfied, *IEEE Trans. Commun.* **35** (1987) 620–625.
- [15] R.O. LaMaire, An M/G/1 vacation model of an FDDI station, *IEEE J. Select. Areas Commun.* **9** (1991) 257–264.
- [16] K.K. Leung, D.M. Lucantoni, Two vacation models for token ring networks where service is controlled by timers, *Performance Evaluation* **20** (1994) 165–184.
- [17] K. Nakamura, T. Takine, Y. Takahashi and T. Hasegawa, An analysis of an asymmetric polling model with cycle-time constraint, in: *Proc. of NATO Adv. Res. Workshop*, Sophia Antipolis, France (June 1990).
- [18] J.W.M. Pang and F.A. Tobagi, Throughput analysis of a timer controlled token passing protocol under heavy load, *IEEE Trans. Communications* **37** (1989) 694–702.
- [19] I. Rubin and J.C.-H. Wu, Throughput and stability analysis of FDDI networks under nonsymmetric load, in: *Proc. IEEE GLOBECOM'93*, Houston (Nov. 29–Dec. 2, 1993) pp. 1154–1158.
- [20] K.C. Sevcik and M.J. Johnson, Cycle time properties of the FDDI token ring protocol, *IEEE Trans. Software Eng.* **13** (1987) 376–385.
- [21] P. Skelly, M. Schwartz and S. Dixit, A histogram-based model for video traffic behaviour in an ATM multiplexer, *IEEE/ACM Transactions on Networking* **1** (1993) 446–459.
- [22] H. Takagi, Effects of the target token rotation time on the performance of a timed-token protocol, in: *Proc. Performance '90*, Edinburgh (1990) 363–370.
- [23] M. Tangemann, Mean waiting time approximations for FDDI, in: J. Labetoulle and J.W. Roberts, eds., *The Fundamental Role of Teletraffic in the Evolution of Telecommunications Networks*, Elsevier, Amsterdam (1994) 973–984.
- [24] M. Tangemann and K. Sauer, Performance analysis of the timed token protocol of FDDI and FDDI-II, *IEEE J. Select. Areas Commun.* **9** (1991) 271–278.
- [25] J.M. Ullm, A timed token ring local area network and its performance characteristics, in: *Proc. 7th Conf. Local Comput. Networks*, Minneapolis (1982) 50–56.
- [26] A. Valenzano, P. Montuschi and L. Ciminiera, Some properties of timed token medium access protocols, *IEEE Trans. on Software Engineering* **16** (1990) 858–869.
- [27] W.B. Van den Hout and J.P.C. Blanc, Development and justification of the power-series algorithm for BMAP-systems, *Commun. Statist.-Stochastic Models* **11** (1995) 471–496.