

A state-dependent polling model with k-limited service

Citation for published version (APA):

Winands, E. M. M., Adan, I. J. B. F., & Houtum, van, G. J. J. A. N. (2006). *A state-dependent polling model with k-limited service*. (BETA publicatie : working papers; Vol. 171). Technische Universiteit Eindhoven.

Document status and date:

Published: 01/01/2006

Document Version:

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

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

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

A state-dependent polling model with k -limited service

E.M.M. Winands^{1,2}, I.J.B.F. Adan¹ and G.J. van Houtum²

¹Department of Mathematics and Computer Science

²Department of Technology Management

Technische Universiteit Eindhoven

P.O. Box 513, 5600 MB Eindhoven, The Netherlands

{e.m.m.winands,i.j.b.f.adan,g.j.v.houtum}@tue.nl

June 26, 2006

Abstract

We consider a two-queue model with state-dependent setups, in which a single server alternately serves the two queues. The high-priority queue is served exhaustively, whereas the low-priority queue is served according to the k -limited strategy. A setup at a queue is incurred only if there are customers waiting at the polled queue. We obtain the transforms of the queue length and sojourn time distributions under the assumption of Poisson arrivals, generally distributed service times and generally distributed setup times. The interest for this model is fueled by an application in the field of logistics. It is shown how the results of this analysis can be applied in the evaluation of a stochastic two-item single-capacity production system. From these results we can conclude that significant cost reductions are possible by bounding the production runs of the low-priority item, which indicates the potential of the k -limited service discipline as priority rule in production environments.

Keywords: polling model, two queues, k -limited and exhaustive service policy, state-dependent setups.

1 Introduction

The present paper considers a two-queue *state-dependent* polling model, in which a setup is incurred for a queue only when it is non-empty. In this model, the single server serves the high-priority queue *exhaustively* and the low-priority queue according to the *k-limited* service strategy. Exhaustive means that the queue must be empty before the server moves on, while under the *k-limited* strategy the server continues working until either a predefined number of *k* customers is served or until the queue becomes empty, whichever occurs first.

The motivation for the present study is two-fold. The first one is application-oriented. Although in the past the *k-limited* strategy proved its merit in communication systems (see, e.g., [4, 5]), the specific application that attracted our attention is in the field of logistics. In particular, in many stochastic multi-product single-capacity make-to-stock production systems considerable setup times are incurred, i.e., the so-called *stochastic economic lot scheduling problem* (SELSP). For surveys on the SELSP, see Sox *et al.* [20] and Winands *et al.* [26]. The SELSP is a common problem in practice, e.g., in glass and paper production, injection molding, metal stamping and semi-continuous chemical processes, but also in bulk production of consumer products such as detergents and beers.

The presence of these setup times in combination with the stochastic environment are the key complicating factors of the SELSP. On the one hand, one aims for short cycle lengths, and thus frequent production opportunities for the various products, in order to be able to react to the stochasticity in the system. On the other hand, short cycle lengths will increase the setup frequency, which has a negative influence on the amount of capacity available for production. Consequently, this effect will hinder the timely fulfillment of demand. An important issue is the fact that in the SELSP - and in various other applications of polling systems - the objective function typically depends not only on the *mean* queue lengths, but on the *complete* marginal queue length distributions. The main interest of the present paper is, therefore, in the marginal queue length distributions in the aforementioned queueing model.

In the context of the SELSP, the *exhaustive* service discipline with *state-independent* setups, i.e., the machine sets up irrespective of whether the polled product has a positive shortfall, i.e. has outstanding production orders, has been studied by Federguen and Katalan [7, 8]. Their strategy has two major drawbacks. First of all, in the exhaustive policy one single product, for which a high demand arrives in a certain period of time, may occupy the machine for quite a while. The impacts of this phenomenon on the other products are stock outs, highly variable cycle lengths and high costs. The *k-limited* policy circumvents this drawback and offers the possibility to the manager to control both the setup frequencies and the production runs (and, thus, the cycle lengths). Secondly, Federguen and Katalan [7, 8] assume that the machine incurs a setup for a certain product even when there is no shortfall for this product, which is of course suboptimal, as argued by Sox *et al.* [20]. This observation makes the practical relevance of the inclusion of state-dependent setups in the studied queueing model evident.

From the application point of view, the present paper offers, albeit in an idealized mathematical setting, a preliminary exploration of an important managerial issue: What is the gain in performance of bounding production runs - at the expense of higher setup frequencies - in multi-item production-inventory systems. In order to answer this question analytically and not to be diverted by other effects, the present paper focusses on a basic occurrence of the SELSP. Ungainsayable, the two-queue model is, however, also of interest in its own right: Production applications, in which only two items have to be produced on one single

production facility, are certainly not inconceivable. In such settings, it is quite natural to apply a limited-service mechanism to provide different service to the different items in order to improve system performance.

A second motivation for the present work is the fact that we have a theoretical interest in the proposed queueing model. That is, the focus of attention is on the *exact* evaluation of this model in order to get better fundamental insights. However, so far, hardly any exact results for polling systems with the k -limited service policy have been obtained; even mean performance measures are, in general, not known. This can be explained by the fact that the k -limited strategy does not satisfy a well-known branching property for polling systems independently discovered by Fuhrmann [9] and Resing [18]. For general k , an exact evaluation for the queue length distribution is, therefore, only available for very few special two-queue cases (see Lee [14] and Ozawa [16, 17]). In these models one (low-priority) queue is served by the k -limited service strategy, whereas the other (high-priority) queue is served by the exhaustive policy. Furthermore, all of these papers make the restrictive assumption of *zero setup times*. In many applications such as the SELSP, however, the setup times may be substantial and the presence of these setup times may be crucial for the operation of the system.

A second feature of the proposed model seriously complicating the analysis is the presence of the state-dependent setups (the server sets up for a queue only when it is non-empty). Nearly all of the existing literature on polling systems makes the assumption of state-independent setups. Notable exceptions are the recent studies of Altman *et al.* [3], Günalay and Gupta [10], Gupta and Srinivasan [11] and Singh and Srinivasan [19], where exhaustive-type and gated-type service disciplines are explored in combination with state-dependent setups. The choice of modeling state-independent setups is generally not motivated by an application but by the tractability of the resulting analysis. The present work is the first study combining the limited service discipline and state-dependent setups. Finally, we refer to Takagi [22, 23, 24] and to Levy [15] for extensive general surveys on polling models and their applications.

Finally, we should bring the paper of Borst *et al.* [4] to the attention, in which approximate optimal values of the service limits with respect to a weighted sum of mean waiting times are obtained for general k -limited polling systems with state-independent setup times. Of particular interest to the present paper is the fact that they derive a (partially conjectured) rule stating that for optimal operation of these systems the queues with the highest priority must have their service limit set at infinity. In a two-queue system this would result into the priority rule of the present paper providing additional evidence for the significance of the present study.

In sum, the main contribution of the present paper is two-fold. Firstly, the model in [14] is generalized by including state-dependent setups. In particular, we obtain the transforms of the queue length and sojourn time distributions under the assumption of Poisson arrivals, generally distributed service times and generally distributed setup times. Secondly, we demonstrate how the results of the analysis can be applied in the evaluation of a stochastic two-item single-capacity production system. We observe significant cost reductions by application of the k -limited policy, compared to the standard exhaustive policies, in such settings indicating the potential of the k -limited service discipline as priority rule in production environments.

The rest of the present paper is organized as follows. In Section 2, we present the model description including the stability conditions and the balance equations. Section 3 derives the *probability generating functions* (PGFs) of the joint and marginal queue length distributions both at service completions epochs and at arbitrary instants. The penultimate section is

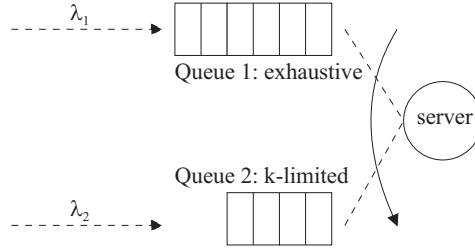


Figure 1: The model.

devoted to the application of the analysis in the field of the SELSP. Some concluding remarks are presented in Section 5.

2 Model description

We start this section with the description of the notation and assumptions used throughout the present paper. Then, Subsections 2.2 and 2.3 present the stability condition and the state description together with the corresponding balance equations, respectively.

2.1 Notation and assumptions

We study a two-queue model, in which a single server alternately serves the two queues. The high-priority queue 1 is served *exhaustively*, whereas the low-priority queue 2 is served according to the *k-limited* strategy (see Figure 1). The combination of these service policies obviously creates a preferential treatment of type-1 customers. Customers arrive at queue i according to a Poisson process with rate $\lambda_i > 0$. The service times at queue i are independent, identically distributed random variables with mean $\beta_i > 0$ and *Laplace Stieltjes Transform* (LST) $B_i(\cdot)$. When the server starts service at queue i , a setup time is incurred with mean $\tau_i \geq 0$ and LST $T_i(\cdot)$. These setup times are identically distributed random variables. The occupation rate ρ_i at queue i is defined by $\rho_i = \lambda_i \beta_i$ and the total occupation rate ρ is given by $\rho = \rho_1 + \rho_2 < 1$. In the next subsection, we give a stability condition for the system in terms of the total occupation rate and the service parameter $k \in \{1, 2, \dots\}$.

The setup times are assumed to be *state-dependent*, i.e., the server incurs a setup for a queue only when it is non-empty. When both queues are empty, the server stops working. He starts again upon arrival of a new customer and, then, he has to setup irrespective of the type of the last customer served before the idle time. Moreover, if the server has served k customers of the low-priority queue and the high-priority queue is empty, the server starts a new sequence up to k customers of this low-priority class after a new setup time of the low-priority queue. However, it is important to note that in this case no setup is incurred for the high-priority queue (which is the standard assumption in polling systems with state-independent setup times).

Finally, we define $S_i(z_1, z_2)$ and $R_i(z_1, z_2)$ as the PGFs of the number of type-1 and type-2 arrivals during a service time and a setup time at queue i , respectively. That is,

$$S_i(z_1, z_2) = B_i(\lambda_1(1 - z_1) + \lambda_2(1 - z_2)), \quad i = 1, 2, \quad (1)$$

$$R_i(z_1, z_2) = T_i(\lambda_1(1 - z_1) + \lambda_2(1 - z_2)), \quad i = 1, 2. \quad (2)$$

The quantities

$$r_1 = \frac{\lambda_1}{\lambda_1 + \lambda_2}, \quad \text{and} \quad r_2 = \frac{\lambda_2}{\lambda_1 + \lambda_2}, \quad (3)$$

denote the probabilities that the server switches to queue 1 and 2 after an idle period, respectively.

Remark 2.1. It is important to note that the analysis of the present paper also fully holds in the case of state-independent setups, *mutatis mutandis*, which are observed occasionally in telecommunication applications. However, we envision production systems as the main application for the present paper (see also Section 4) and, in such applications, state-independent setups are definitely no realistic assumption. In particular, the specific setup rule assumed here is quite common in, for example, chemical industry. After the cleaning of the machine and delivery of raw materials, the machine is able to produce a batch with a certain maximum size. At the moment this maximum batch has been produced, the machine has to be cleaned again and raw materials have to be replenished. \square

2.2 Stability conditions

In this subsection, we derive the stability conditions for the two queues by deploying (heuristic) arguments similar to those used by Ibe and Cheng [12], who study the stability conditions for a variety of polling systems without state-dependent setups.

Since the k -limited service policy bounds the time the server spends at queue 2, queue 2 will not affect the stability of queue 1. Hence, the stability condition of the latter simply reads

$$\rho_1 < 1. \quad (4)$$

If this condition is not satisfied, queue 2 is unstable as well: the server will serve queue 1 for an infinite amount of time and, in the meantime, the queue length of queue 2 grows to infinity.

For queue 2, we define a cycle time as the time between the start of two successive setups of the server at this queue. In a long period of time T , the number of customers arriving at queue 2 equals $\lambda_2 T$. Each customer brings, on average, β_1 units of work into the system. We approximately need $\frac{\lambda_2 T}{k}$ cycles to serve these customers in heavy traffic, which is the regime of interest for the derivation of the stability condition. In each cycle, the server always has to setup for queue 2 and has to setup for queue 1 in case queue 1 is non-empty at the instant that the server has served k customers at queue 2. The latter occurs with the probability q that the number of type-1 arrivals during a setup time plus k successive service times at queue 2 is not equal to zero, i.e.,

$$q = 1 - R_2(0, 1)(S_2(0, 1))^k. \quad (5)$$

Hence, the mean total setup time in each cycle equals $\tau_2 + q\tau_1$. Moreover, a fraction ρ_1 of time is used to serve customers at queue 1. Since queue 2 is stable if the server manages to serve all customers arrived in T in an amount of time less than T , we have

$$\rho_1 T + \rho_2 T + \frac{\lambda_2 T}{k}(\tau_2 + q\tau_1) < T, \quad (6)$$

which gives us, after division by T ,

$$\rho_1 + \rho_2 + \frac{\lambda_2}{k}(\tau_2 + q\tau_1) < 1. \quad (7)$$

When queue 2 satisfies (7), ρ_1 is also smaller than 1 and, thus, in this case both queues are stable. Throughout the present paper, we assume that (7) is satisfied, as we restrict ourselves to steady-state behavior.

2.3 State description and balance equations

We study the system at embedded epochs of service completions of customers. The state of the system $\mathbf{Q}(n)$ just after the n^{th} departure from the system can be described by the following three variables:

1. $Q_1(n)$: the number of customers in queue 1;
2. $Q_2(n)$: the number of customers in queue 2;
3. $C(n)$: equals zero when the n^{th} departure is a type-1 customer, while it equals the number of type-2 departures since the last setup when the n^{th} departure is a type-2 customer.

The associated stochastic process,

$$\mathbf{Q}(n) = \{(Q_1(n), Q_2(n), C(n)), n = 1, 2, \dots\}, \quad (8)$$

is an aperiodic and irreducible three-dimensional Markov chain. Let

$$\pi(q_1, q_2, c) = \lim_{n \rightarrow \infty} P[(Q_1(n), Q_2(n), C(n)) = (q_1, q_2, c)], \quad (9)$$

be the equilibrium state probability and define the corresponding generating functions for this Markov chain as follows

$$p_1(z_1, z_2) = \sum_{q_1=0}^{\infty} \sum_{q_2=0}^{\infty} \pi(q_1, q_2, 0) z_1^{q_1} z_2^{q_2}, \quad (10)$$

$$p_{2,j}(z_1, z_2) = \sum_{q_1=0}^{\infty} \sum_{q_2=0}^{\infty} \pi(q_1, q_2, j) z_1^{q_1} z_2^{q_2}, \quad j = 1, 2, \dots, k. \quad (11)$$

Now, the following set of $k+1$ balance equations for equally many unknowns $p_1(z_1, z_2)$ and $p_{2,j}(z_1, z_2)$, $j = 1, 2, \dots, k$, holds

$$p_1(z_1, z_2) = \frac{S_1(z_1, z_2)}{z_1} \left\{ p_1(z_1, z_2) - p_1(0, z_2) + \left[c_0 r_1 z_1 + \sum_{j=1}^{k-1} [p_{2,j}(z_1, 0) - p_{2,j}(0, 0)] + p_{2,k}(z_1, z_2) - p_{2,k}(0, z_2) \right] R_1(z_1, z_2) \right\}, \quad (12)$$

$$p_{2,1}(z_1, z_2) = \frac{S_2(z_1, z_2) R_2(z_1, z_2)}{z_2} \alpha(z_2), \quad (13)$$

$$p_{2,j}(z_1, z_2) = \frac{S_2(z_1, z_2)}{z_2} \left\{ p_{2,j-1}(z_1, z_2) - p_{2,j-1}(z_1, 0) \right\}, \quad j = 1, 2, \dots, k, \quad (14)$$

with

$$c_0 = p_1(0, 0) + \sum_{j=1}^k p_{2,j}(0, 0), \quad (15)$$

the probability that the system is left idle at a departure epoch and

$$\alpha(z_2) = c_0 r_2 z_2 + p_1(0, z_2) - p_1(0, 0) + p_{2,k}(0, z_2) - p_{2,k}(0, 0). \quad (16)$$

These balance equations are formulated by considering all the possible states at the previous departure epoch from which we can reach the current state. We explain (12), which describes the case that the current departure is a type-1 customer. First of all, the previous departure could be a type-1 departure which did not leave the system idle. This event corresponds to the term $p_1(z_1, z_2) - p_1(0, z_2)$. Secondly, the term $c_0 r_1 z_1$ represents the event that the previous departure left the system idle and that the first new arriving customer is of type 1. Thirdly, the term $\sum_{j=1}^{k-1} [p_{2,j}(z_1, 0) - p_{2,j}(0, 0)]$ represents the event that the last departure was a type-2 customer that was not the k^{th} in the sequence and that left queue 2, but not the complete system, idle. Finally, $p_{2,k}(z_1, z_2) - p_{2,k}(0, z_2)$ corresponds to the event that the last departure was a type-2 customer that was the k^{th} in the sequence and that queue 1 was not empty. The explanations of (13) and (14) are similar.

The function $\alpha(\cdot)$ is recognized as the generating function of the number of type-2 customers at moments a setup for queue 2 is initiated. Since at these specific points in time Q_1 and C are both equal to zero, the state description can be reduced to one single dimension represented by the function $\alpha(\cdot)$. The analysis of the next section is oriented towards relating the unknown generating functions $p_1(\cdot)$ and $p_{2,j}(\cdot, \cdot)$ to this function $\alpha(\cdot)$.

3 Exact analysis of queue lengths

In Subsection 3.1, we present the derivation of the generating functions of the joint queue length distributions at service completion epochs, which is a generalization of the method used by Lee [14] for the model without setup times. In Subsection 3.2, these results are used to derive expressions for the PGFs of the marginal queue size distributions at arbitrary instants.

3.1 Joint queue lengths at service completion epochs

To derive the generating functions of the joint queue length distributions at service completion epochs, we successively substitute (14) into itself and, then, into (13) which yields, for $j = 1, 2, \dots, k$,

$$p_{2,j}(z_1, z_2) = \frac{S_2^j(z_1, z_2) R_2(z_1, z_2) \alpha(z_2) - \sum_{l=1}^{j-1} z_2^{j-l} S_2^l(z_1, z_2) p_{2,j-l}(z_1, 0)}{z_2^j}. \quad (17)$$

Notice that (17) gives an expression of $p_{2,j}(\cdot, \cdot)$, $j = 1, 2, \dots, k$, as a function of the unknown functions $\alpha(\cdot)$ and $p_{2,l}(\cdot, 0)$, $l = 1, 2, \dots, j - 1$.

Now, we turn our attention to $p_1(\cdot, \cdot)$. Substituting (17) for $j = k$ into (12) and using (15) and (16) gives us, after some straightforward manipulations,

$$\begin{aligned} (z_1 - S_1(z_1, z_2))p_1(z_1, z_2) &= S_1(z_1, z_2) \left\{ c_0 r_1 z_1 R_1(z_1, z_2) + (R_1(z_1, z_2) - 1)p_1(0, z_2) + \right. \\ &\quad \left(\left(\frac{S_2(z_1, z_2)}{z_2} \right)^k R_2(z_1, z_2) \alpha(z_2) - \alpha(z_2) + c_0 r_2 z_2 + \right. \\ &\quad \left. \sum_{j=1}^{k-1} \left[1 - \left(\frac{S_2(z_1, z_2)}{z_2} \right)^j \right] p_{2,k-j}(z_1, 0) - c_0 \right) R_1(z_1, z_2) \left. \right\}. \end{aligned} \quad (18)$$

We eliminate $p_1(0, z_2)$ from the above equation by rewriting (16) as follows

$$\begin{aligned} p_1(0, z_2) &= \alpha(z_2) - c_0 r_2 z_2 + p_1(0, 0) - p_{2,k}(0, z_2) + p_{2,k}(0, 0) \\ &= \alpha(z_2) + c_0(1 - r_2 z_2) - \left(\frac{S_2(0, z_2)}{z_2} \right)^k R_2(0, z_2) \alpha(z_2) \\ &\quad - \sum_{j=1}^{k-1} \left[1 - \left(\frac{S_2(0, z_2)}{z_2} \right)^j \right] p_{2,k-j}(0, 0), \end{aligned} \quad (19)$$

which yields

$$\begin{aligned} (z_1 - S_1(z_1, z_2))p_1(z_1, z_2) &= S_1(z_1, z_2) \left\{ \left(\frac{\beta(z_1, z_2)}{z_2^k} - 1 \right) \alpha(z_2) + \right. \\ &\quad R_1(z_1, z_2) \sum_{j=1}^{k-1} \frac{\delta_j(z_1, z_2)}{z_2^k} p_{2,k-j}(z_1, 0) + \\ &\quad \left. D(z_1, z_2) - (R_1(z_1, z_2) - 1) \sum_{j=1}^{k-1} \frac{\delta_j(0, z_2)}{z_2^k} p_{2,k-j}(0, 0) \right\}, \end{aligned} \quad (20)$$

where

$$D(z_1, z_2) = c_0 \left[r_1 z_1 R_1(z_1, z_2) + r_2 z_2 - 1 \right], \quad (21)$$

$$\beta(z_1, z_2) = S_2^k(z_1, z_2) R_1(z_1, z_2) R_2(z_1, z_2) - S_2^k(0, z_2) (R_1(z_1, z_2) - 1) R_2(0, z_2), \quad (22)$$

$$\delta_j(z_1, z_2) = z_2^k - z_2^{k-j} S_2^j(z_1, z_2). \quad (23)$$

It is again important to notice that via (20), $p_1(\cdot, \cdot)$ is also expressed as a function of the unknown functions $\alpha(\cdot)$ and $p_{2,j}(\cdot, 0)$, $j = 1, 2, \dots, k-1$.

It is well-known that for each (fixed) $|z_2| \leq 1$ the term $z_1 - S_1(z_1, z_2)$ in (20) has exactly one zero $z_1 = \xi(z_2)$ with $|z_2| \leq 1$ if $\rho_1 < 1$. More specifically,

$$z_1 = \xi(z_2) = \gamma_1[\lambda_2(1 - z_2)], \quad (24)$$

where $\gamma_1(\cdot)$ is the LST of the busy period of a standard M/G/1 queue with arrival rate λ_1 and LST of the service time $B_1(\cdot)$ (see, e.g., Takács [21]). Thus, $\xi(\cdot)$ can be seen as the PGF of the distribution of the number of type-2 arrivals during such an M/G/1 busy period.

Remark 3.1. It is interesting to note that the function $\beta(\xi(z), z)$ is the PGF of the number of type-2 customers arriving in a cycle for queue 2 in which the maximum of k customers is served. \square

By analyticity of $p_1(z_1, z_2)$, the right-hand side of (20) should vanish when $z_1 = \xi(z_2)$. Hence,

$$\alpha(z) = \frac{D(\xi(z), z)z^k + R_1(\xi(z), z) \sum_{j=1}^{k-1} \delta_j(\xi(z), z)p_{2,k-j}(\xi(z), 0) - (R_1(\xi(z), z) - 1) \sum_{j=1}^{k-1} \delta_j(0, z)p_{2,k-j}(0, 0)}{z^k - \beta(\xi(z), z)}, \quad (25)$$

and $\alpha(z)$ is formulated as a function of the unknown functions $p_{2,j}(\cdot, 0)$, $j = 1, 2, \dots, k-1$.

To eliminate these unknown functions, we differentiate the numerator and denominator of (17) j times with respect to z_2 and, by L'Hospital's rule, we obtain the following recursion, for $j = 1, 2, \dots, k-1$,

$$p_{2,j}(x, 0) = \sum_{l=1}^j \frac{c_l \frac{d^{j-l}}{dy^{j-l}} [S_2^j(x, y) R_2(x, y)] \Big|_{y=0}}{(j-l)!} - \sum_{l=1}^{j-1} \frac{\frac{d^l}{dy^l} [S_2^l(x, y)] \Big|_{y=0}}{l!} p_{2,j-l}(x, 0), \quad (26)$$

where c_l , $l = 1, 2, \dots, k-1$, represent the probabilities that l type-2 customers are present at the start of a setup for this queue, i.e.,

$$c_l = \frac{\frac{d^l}{dy^l} [\alpha(y)] \Big|_{y=0}}{l!}, \quad l = 1, 2, \dots, k-1. \quad (27)$$

From this interpretation of c_l one can easily deduce a probabilistic interpretation of (26) as well, i.e., left and right hand side clearly represent the PGF of the number of customers in queue 1 when the server leaves queue 2 due to the fact that the latter queue is empty after j customers served, $j = 1, 2, \dots, k-1$. Of course, we could have derived (26) directly by using this probabilistic interpretation and, hence, avoid use of L'Hospital's rule.

By (26) we can write $p_{2,j}(\cdot, 0)$ as a function of the unknown probabilities c_j , $j = 0, 1, \dots, k-1$. Moreover, with the help of (17), (20) and (25) the generating functions $p_{2,j}(\cdot, \cdot)$, $p_1(\cdot)$ and $\alpha(\cdot)$ can be expressed in terms of these constants as well. The problem of finding these generating functions is, thus, reduced to finding the unknown probabilities c_j , which can be computed as follows. Given that (7) holds, the following theorem states that the denominator of (25) has exactly k zeros on or within the unit circle.

Theorem 3.2. *Under (7), it holds that $z^k = \beta(\xi(z), z)$ has k roots on or within the unit circle.*

Proof The derivative of $\beta(\xi(z), z)$ at $z = 1$ equals

$$\beta'(\xi(1), 1) = \frac{\lambda_2}{1 - \rho_1} (k\beta_2 + \tau_1 + \tau_2) - \frac{\lambda_2}{1 - \rho_1} S_2^k(0, z_2) R_2(0, z_2) \tau_1 = \frac{k\rho_2 + \lambda_2(q\tau_1 + \tau_2)}{1 - \rho_1}, \quad (28)$$

where q is defined by (5). Because we have assumed (7), we obtain $\beta'(\xi(1), 1) < k$ and the result follows from Theorem 3.2 of Adan *et al.* [2]. \square

Since $\alpha(z)$ is bounded in $|z| \leq 1$, the zeros in the numerator must be canceled by corresponding zeros in the denominator. One of the zeros equals one and leads to a trivial equation. However, the normalization condition provides an additional equation and, therefore, we have a set of k linear equations. By making the assumption that the k roots of $z^k = \beta(\xi(z), z)$ on or within the unit circle are all distinct (see Remark 3.4), this set of equations has a unique solution for c_j , $j = 0, 2, \dots, k-1$. This completes the determination of the generating functions of the queue length distributions at service completion epochs and, hence, in the remainder of the present paper we assume that these generating functions are known.

Remark 3.3. In Subsection 2.2 we have heuristically derived the stability condition for the model under consideration. However, the fact that the same condition shows up in Theorem 3.2 supports the validity of this condition. \square

Remark 3.4. If one or more roots of $z^k = \beta(\xi(z), z)$ on or within the unit circle coincide, our reasoning needs to be slightly modified. That is, for $\alpha(z)$ to be bounded in $|z| \leq 1$, the numerator of (25) should still have the same zeros as the denominator of (25) *and* with the same multiplicity. Additional equations can, therefore, be obtained by requiring that the derivative(s) of the numerator should also vanish where the denominator has a zero of higher multiplicity. \square

3.2 Marginal queue lengths at arbitrary instants

From the results of the previous subsection, we can obtain expressions for the PGF $q_i(\cdot)$ of the marginal queue size distributions of queue i at type- i departure epochs, i.e.,

$$q_1(z) = \frac{p_1(z, 1)}{r_1}, \quad \text{and} \quad q_2(z) = \frac{\sum_{j=1}^k p_{2,j}(1, z)}{r_2}. \quad (29)$$

By using a standard level crossing argument, in combination with PASTA, it can be shown that the marginal queue length distribution of queue i at type- i departure epochs and at arbitrary instants in time are the same. Hence, the PGFs for these marginal distributions are given by (29).

From (29) we can easily obtain the LST $W_i(\cdot)$ of the sojourn time distribution of a type- i customer. Since the number of type- i customers left behind by a tagged type- i customer equals the number of customers arrived during the sojourn time of this tagged customer, we have

$$W_i(z) = q_i\left(1 - \frac{z}{\lambda_i}\right), \quad i = 1, 2, \quad (30)$$

which is known as the distributional form of Little's law (see, e.g., Keilson and Servi [13]).

4 Application

In this section, we use the analysis of the present paper to evaluate a stochastic two-item single-capacity production system. After the description of this application in Subsection 4.1, we present a numerical illustration in Subsection 4.2.

4.1 Stochastic economic lot scheduling problem

Consider a system with one single production capacity for two products, in which there is an infinite stock space for each product and raw material is always available. Demands for the two products arrive according to stationary and mutually independent unit Poisson processes with rate $\lambda_i > 0$. Demand that cannot be satisfied directly from stock is backlogged until the product becomes available after production. The individual products are produced make-to-stock in a cyclical order with stochastic production times with mean $\beta_i > 0$ and LST $B_i(\cdot)$. Stochastic setup times with mean $\tau_i \geq 0$ and LST $T_i(\cdot)$ occur *before* the start of the

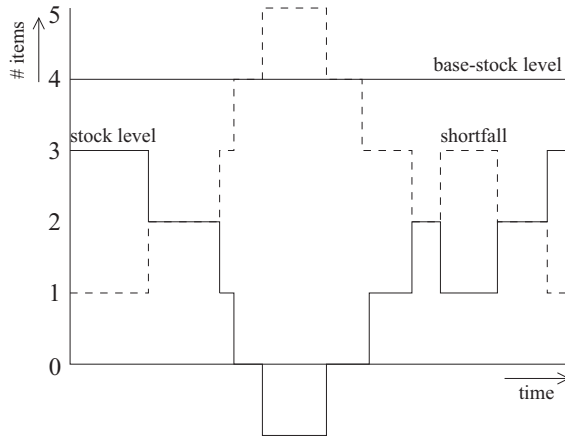


Figure 2: Example of the relation between base-stock level, stock level and shortfall.

production of a product. For an arbitrary number of products, this setting is often referred to as the *stochastic economic lot scheduling problem* (see Winands *et al.* [26], for a recent survey).

For product 1, a standard *base-stock* policy is implemented, i.e., when production is commenced, the machine will continue production until a pre-defined target stock level b_1 has been reached. For product 2, a *quantity-limited base-stock policy* is used. That is, when the machine starts production, it will continue production until either the base-stock level b_2 has been reached or a maximum number k of products has been produced. A production order is placed immediately after demand for the corresponding product has arrived. These production orders queue up at the production facility.

The number of outstanding production orders of product i at this facility is indicated as the *shortfall* of this product. The setup times are assumed to be state-dependent, i.e., no setup for a product is incurred when there is no *shortfall* (no outstanding production orders). When both products have no shortfall, the machine is turned off and is turned on again upon demand arrival. When a type-2 batch of size k has been produced and product 1 has no shortfall, a new type-2 batch up to k products is started after a new setup time.

For given values of the base-stocks, b_1 and b_2 , and the quantity limit, k , the steady-state stock level distribution I_i for product i is given by (see Figure 2)

$$I_i = b_i - L_i, \quad i = 1, 2, \quad (31)$$

where L_i denotes the steady-state shortfall of product i . Notice that the stock level becomes negative, when the shortfall of a certain product is larger than its base-stock level.

It is easily seen that the shortfall of a product is independent of the base-stock levels. Moreover, it is easily verified that the shortfall distribution of product i at the production facility is identical to the queue length distribution of queue i in the queueing model of the present paper. Hence, by the procedure presented in Section 3, in combination with (31), the steady-state stock level for both products can be computed. With these distributions, various performance measures of interest can be computed.

In the remainder of the paper, we consider the total expected costs C , i.e., the sum of holding and backlogging costs. For these individual cost components, we assume a simple

linear structure. That is, the cost rate function for product i as a function of the stock level x equals

$$c_i(x) = \begin{cases} h_i x, & x \geq 0, \\ -p_i x, & x < 0, \end{cases} \quad (32)$$

where h_i and p_i represent the holding and penalty cost coefficients for product i , respectively. With the help of (31), the total expected costs C can be written as follows

$$C = C_1 + C_2 = \mathbb{E}[c_1(I_1)] + \mathbb{E}[c_2(I_2)] = \mathbb{E}[c_1(b_1 - L_1)] + \mathbb{E}[c_2(b_2 - L_2)]. \quad (33)$$

For a given value of k , it can be shown that the optimal base-stock levels b_i^* are given by

$$b_i^* = \min\{n \in \mathbb{N}_0 | P[L_i \leq n] \geq \frac{p_i}{p_i + h_i}\}, \quad i = 1, 2, \quad (34)$$

which is recognized as the solution of a standard newsboy problem. Finally, we note that we do not consider the optimization of the quantity limit k here.

4.2 Numerical illustration

We now present some cases, which illustrate the value of the procedure of Section 3 in the evaluation of the described production system. Of course, a whole plethora of cases can be studied: different values of the quantity limit, choice of service time distributions and their parameters, choice of setup distributions and their parameters, different (ratios between) cost factors, etcetera. However, the aim of the present section is to present some illustrative cases which show the potential of the k -limited service policy in the context of the SELSP.

As described before, the PGFs derived in the previous section have to be finished off with a number of zeros, which are numerically computed by using the Chaudhry QROOT software package [6]. We, then, use the method presented by Abate and Whitt [1] to numerically invert these PGFs. Unfortunately, for large quantity limits, numerical problems have been encountered in the procedure. Therefore, we confine ourselves to cases with small limits.

In the numerical evaluation we also present results for the special case of $k = \infty$, which amounts to a two-product model with exhaustive base-stock policy for both products. In this case, we do not use the procedure of Section 3, but we have implemented a discrete event simulation. Each simulation run is sufficiently long such that the widths of the 95% confidence intervals of the performance measures of interest are smaller than 1%.

The examination starts with an initial case, which is subsequently being perturbed into 6 cases to study the effect of (1) (ratio between) the loads, (2) (ratio between) cost factors, (3) setup times.

Case 1 (*Initial case*). Suppose that product 1 is a product with high costs, whereas product 2 is of secondary importance compared to the first product. Table 1 shows the detailed specifications for these two products. Table 2 shows the costs C_i per product, the total costs C and the optimal base-stock levels b_i^* as a function of the quantity limit k . Firstly, we observe that the marginal costs for product 2 are decreasing in the quantity limit, whereas the marginal costs for product 1 increase with this limit. Secondly, the optimal value of the quantity limit with respect to total costs is equal to 2. For smaller limits the amount of capacity available for production is too low, while larger limits lead to more variable cycle lengths. Thirdly, the optimal base-stock levels for product 2 are non-increasing in the

Product information		
Parameter	Product 1	Product 2
Demand	Poisson(0.375)	Poisson(0.375)
Processing times	Exp(1.0)	Exp(1.0)
Setup times	Exp(0.25)	Exp(0.25)
Holding costs	10	1
Backlogging costs	90	9

Table 1: Product information for Case 1.

Output					
k	b_1^*	C_1	b_2^*	C_2	C
1	3	27.5	13	13.9	41.4
2	3	28.8	9	9.8	38.5
3	3	30.5	7	8.7	39.2
4	3	32.5	7	8.1	40.6
5	3	34.6	6	7.7	42.3
6	4	36.5	6	7.4	43.9

Table 2: Case 1.

quantity limit, while the optimal base-stock levels for product 1 are non-decreasing in this limit. Finally, it is interesting to compare this table with the total costs equalling 59.2 that would be incurred if a standard exhaustive policy were implemented for both products. Via a k -limited policy for product 2, we may, thus, save 35.0% compared to the latter policy clearly showing the advantage of the k -limited policy in a production environment.

Case 2 and 3 (*Effect of the load*). These cases are similar to the initial case, except that the demand rates are perturbed, i.e., in Case 2 the demand rates of product 1 and 2 equal 0.15 and 0.6, whereas in Case 3 product 1 and 2 have demand rates 0.6 and 0.15, respectively. Comparing the results in Tables 3 and 4 with the costs for exhaustive base-stock policies equalling 49.3 and 61.7, respectively, once more significant cost reductions are observed by application of the k -limited policy. Of course, the advantages of the k -limited policy are much more pronounced in case the low-priority product 2 has the highest demand rate.

Output					
k	b_1^*	C_1	b_2^*	C_2	C
1	1	13.8	24	24.9	38.7
2	1	15.4	12	12.5	27.8
3	1	17.1	10	10.6	27.7
4	2	18.6	9	9.8	28.4
5	2	19.2	8	9.3	28.5
6	2	19.9	8	9.0	28.8

Table 3: Case 2.

Output					
k	b_1^*	C_1	b_2^*	C_2	C
1	5	47.6	5	6.6	54.2
2	5	48.7	4	5.8	54.5
3	5	49.9	4	5.4	55.3
4	5	51.0	4	5.2	56.2
5	5	52.0	4	5.0	57.0
6	5	53.0	3	4.9	57.9

Table 4: Case 3.

Case 4 and 5 (*Effect of the cost factors*). In the fourth and fifth case, we consider systems similar to that of the initial case and perturb the cost factors. That is, the cost factors for product 2 remain unaltered, while the holding and penalty costs for product 1 are decreased to 5 and 45 for Case 4 and to 2 and 18 for Case 5, respectively. Although the same conclusions as in Case 1 can be drawn from Tables 5 and 6, it should be observed that the advantages of the k -limited discipline are a bit less pronounced in these cases due to the leveling of the costs among the products. In fact, application of exhaustive policies for both products, which is normally done, would lead to total costs of 32.3 and 16.1, respectively.

Output					
k	b_1^*	C_1	b_2^*	C_2	C
1	3	13.7	13	13.9	27.6
2	3	14.4	9	9.8	24.2
3	3	15.2	7	8.7	24.0
4	3	16.2	7	8.1	24.3
5	3	17.3	6	7.7	25.0
6	4	18.3	6	7.4	25.6

Table 5: Case 4.

Output					
k	b_1^*	C_1	b_2^*	C_2	C
1	3	5.5	13	13.9	19.4
2	3	5.8	9	9.8	15.5
3	3	6.1	7	8.7	14.8
4	3	6.5	7	8.1	14.6
5	3	6.9	6	7.7	14.6
6	4	7.3	6	7.4	14.7

Table 6: Case 5.

Case 6 and 7 (*Effect of the setup times*). In Cases 6 and 7 we examine what the effects are of the sizes of the setup times. We therefore study two cases similar to the initial case, but in which the mean setup times equal 0.5 and 0.1, respectively. See Tables 7 and 8 for the results. Notice that for Case 6 the system is not stable when the quantity limit is chosen to be equal to 1. In case we implement the exhaustive base-stock policy for both products, the total costs are given by 62.7 and 56.3. This leads to the intuitively appealing conclusion that the advantages of the k -limited service discipline slightly increase in case the setup times vanish, but we stress that even in the case of large setup times the k -limited strategy still shows its superiority over the exhaustive policy.

Output					
k	b_1^*	C_1	b_2^*	C_2	C
1	—	—	—	—	—
2	3	30.2	16	16.9	47.1
3	3	32.4	11	12.1	44.5
4	4	34.7	9	10.4	45.0
5	4	36.0	8	9.4	45.4
6	4	37.6	7	8.8	46.4

Table 7: Case 6.

Output					
k	b_1^*	C_1	b_2^*	C_2	C
1	2	26.7	8	8.8	35.4
2	3	28.3	7	7.9	36.2
3	3	29.8	7	7.5	37.2
4	3	31.5	7	7.2	38.6
5	3	33.3	6	7.0	40.3
6	3	35.3	5	6.7	42.0

Table 8: Case 7.

Conclusion. In the present section we have presented some cases, by using the results of the exact analysis of the present paper, which aim to answer the question: What is the gain in performance of bounding production runs - at the expense of higher setup frequencies - in multi-item production-inventory systems. The results of the present section lead to the conjectures that the widely used exhaustive policy is not the most effective strategy in (frequently encountered) asymmetric production situations as well as that it may be desirable that production runs of low-priority products are bounded in these environments. We touch upon the possibilities of a large-scale (approximate) study to support this statement in more realistic systems in the next section.

5 Conclusions

The present paper has presented an exact analysis of a two-queue state-dependent polling system with k -limited service extending the polling literature on both the k -limited service discipline and on state-dependent setups; containing the non-preemptive priority model and the model of Lee [14] as special cases. Moreover, the results of the analysis have been applied to a make-to-stock production setting, which provides us with theoretical evidence that the k -limited strategy leads to considerable cost reductions compared to widely used (standard) exhaustive policies.

As stated in the introduction, the k -limited policy violates the branching property for polling systems (see Fuhrmann [9] and Resing [18]) implying that extensions of the analysis of the present paper to more realistic systems are, in most likelihood, outside the borders of possibility and that one has to resort to approximations in these cases. Recently, Van Vuuren and Winands [25] developed an efficient and accurate approximate decomposition approach for k -limited polling systems under the assumption of generally distributed arrival, service and setup distributions.

It is our hope that the theoretical evidence on the potential of the k -limited discipline given by the present paper in combination with the efficient approximate algorithm for the k -limited strategy of [25] tempts people to a large-scale study comparing the k -limited and exhaustive strategies in the context of the (make-to-stock) production-inventory environments. Such a study would get us a second step nearer to answering the question raised in the introduction of the present paper: What is the gain in performance of bounding production runs - at the expense of higher setup frequencies - in multi-item production-inventory systems?

Acknowledgement

The authors would like to thank Onno Boxma for valuable comments.

References

- [1] Abate, J., Whitt, W., (1992). *Numerical inversion of probability generating functions* (Operations Research Letters, vol. 12, no. 4, pp. 245-251).
- [2] Adan, I.J.B.F., Leeuwaarden, J.S.H. van, Winands, E.M.M., (2006). *On the application of Rouché's theorem in queueing theory* (Operations Research Letters, vol. 34, no. 3, pp. 355-360).
- [3] Altman, E., Blanc, H., Khamisy, A., Yechiali, U., (1994). *Gated-type polling systems with walking and switch-in times* (Stochastic Models, vol. 10, pp. 741-764).
- [4] Borst, S.C., Boxma, O.J., Levy, H., (1995). *The use of service limits for efficient operation of multistation single-medium communication systems* (IEEE/ACM Transactions on Networking, vol. 3, no. 5, pp. 602-612).
- [5] Charzinski, J., Renger, T., Tangemann, M., (1994). *Simulative comparison of the waiting time distributions in cyclic polling systems with different service strategies* (Proceedings of the 14th International Teletraffic Congress, Antibes Juan-les-Pins, pp. 719-728).
- [6] Chaudhry, M.L., (1994). *QROOT Software Package* (A&A Publications, Kingston).
- [7] Federgruen, A., Katalan, Z., (1996). *The stochastic economic lot scheduling problem: cyclical base-stock policies with idle times* (Management Science, vol. 42, no. 6, pp. 783-796).

- [8] Federgruen, A., Katalan, Z., (1998). *Determining production schedules under base-stock policies in single facility multi-item production systems* (Operations Research, vol. 46, no. 6, pp. 883-898).
- [9] Fuhrmann, S.W., (1981). *Performance analysis of a class of cyclic schedules* (Bell Laboratories Technical Memorandum 81-59531-1).
- [10] Günalay, Y., Gupta, D., (1997). *A polling system with a patient server and state-dependent setup times* (IIE Transactions, vol. 29, pp. 469-480).
- [11] Gupta, D., Srinivasan, M.M., (1996). *Polling systems with state-dependent setup times* (Queueing Systems, vol. 22, pp. 403-423).
- [12] Ibe, O.C., Cheng, X., (1988). *Stability conditions for multiqueue systems with cyclic service* (IEEE Transactions Automatic Control, vol. 33, no. 1, pp. 102-103).
- [13] Keilson, J., Servi, L.D., (1990). *The distributional form of Little's law and the Fuhrmann-Cooper decomposition* (Operations Research Letters, vol. 9, no. 4, pp. 239-247).
- [14] Lee, D.-S., (1996). *A two-queue model with exhaustive and limited service disciplines* (Stochastic Models, vol. 12, no. 2, pp. 285-305).
- [15] Levy, H., Sidi, M., (1990). *Polling systems: applications, modeling and optimization* (IEEE Transactions on Communications, vol. COM-38, no. 10, pp. 1750-1760).
- [16] Ozawa, T., (1990). *Alternating service queues with mixed exhaustive and K-limited services* (Performance Evaluation, vol. 11, pp. 165-175).
- [17] Ozawa, T., (1997). *Waiting time distribution in a two-queue model with mixed exhaustive and gated-type K-limited services* (Proceedings of International Conference on the Performance and Management of Complex Communication Networks, Tsukuba, pp. 231-250).
- [18] Resing, J.A.C., (1993). *Polling systems and multitype branching processes* (Queueing Systems, vol. 13, pp. 409-426).
- [19] Singh, M.P., Srinivasan, M.M., (2002). *Exact analysis of the state dependent polling model* (Queueing Systems, vol. 41, pp. 371-399).
- [20] Sox, C.R., Jackson, P.L., Bowman, A., Muckstadt, J.A., (1999). *A review of the stochastic lot scheduling problem* (International Journal of Production Economics, vol. 62, pp. 181-200).
- [21] Takács, L., (1968). *Two queues attended by a single server* (Operations Research, vol. 16, no. 3, pp. 639-650).
- [22] Takagi, H., (1990). *Queueing analysis of polling models: an update* (In Stochastic Analysis of Computer and Communication Systems, H. Takagi (ed.), North-Holland, Amsterdam, pp. 267-318).
- [23] Takagi, H., (1997). *Queueing analysis of polling models: progress in 1990-1994* (In Frontiers in Queueing: Models, Methods and Problems, J.H. Dshalalow (ed.), CRC Press, Boca Raton, pp. 119-146).
- [24] Takagi, H., (2000). *Analysis and application of polling models* (In Performance Evaluation: Origins and Directions, G. Haring, C. Lindemann and M. Reiser (eds.), Lecture Notes in Computer Science, vol. 1769, Springer, Berlin, pp. 423-442).
- [25] Vuuren, M. van, Winands, E.M.M., (2006). *Iterative approximation of k-limited polling systems* (Report, Eindhoven University of Technology, Eindhoven).
- [26] Winands, E.M.M., Adan, I.J.B.F., Houtum, G.J. van, (2005). *The stochastic economic lot scheduling problem: a survey* (BETA WP-133, Beta Research School for Operations Management and Logistics, Eindhoven).