



Article A Multi-Server Heterogeneous Queuing-Inventory System with Class-Dependent Inventory Access

Karumbathil Rasmi ^{1,†}^(D), Machuveettil Joseph Jacob ^{1,*,†}^(D), Alexander S. Rumyantsev ^{2,3,†}^(D) and Achyutha Krishnamoorthy ^{4,†}

- ¹ Department of Mathematics, National Institute of Technology Calicut, Kozhikode 673601, Kerala, India; rasmi_p180125ma@nitc.ac.in
- ² Institute of Applied Mathematical Research, Karelian Research Centre, Russian Academy of Sciences, 185910 Petrozavodsk, Russia; ar0@krc.karelia.ru
- ³ Department of Applied Mathematics and Cybernetics, Petrozavodsk State University, 185035 Petrozavodsk, Russia
- ⁴ Centre for Research in Mathematics, Department of Mathematics, CMS College, Kottayam 686001, Kerala, India; krishnamoorthy@cmscollege.ac.in
- * Correspondence: mjj@nitc.ac.in
- + These authors contributed equally to this work.

Abstract: In this paper, we consider a queuing inventory system with heterogeneous customers of K types arriving according to a marked Markovian arrival process. Each class of customers differs by nature of the service they seek and different priorities are assigned for each class resulting in different levels of inventory admitted to exhaust for customers of each class. A single service node is provided for each class with exponential services having class-dependent service rates. All classes of customers are served from a single source of inventory replenished according to (s, S) policy with exponentially distributed lead time. Stability condition and steady state probabilities are obtained by matrix-analytic method. Some important performance measures are also derived. Inventory recycle time was analyzed in detail. Useful cost function and numerical illustrations are also given. The optimization problem is interesting and can be solved in similar real scenario.

Keywords: MMAP[K]; heterogeneous inventory access; multi-server system; queuing-inventory; matrix analytic model

1. Introduction

Queuing-inventory systems are in the focus of recent research due to practical applicability in many fields including social, biological and technical systems. Access to a finite consumable and refillable resource is a natural way of modeling interaction with retail shop customers, office visitors, hospital patients, and even packages in the telecommunication network. In such systems, various sophisticated models arise, including the models with random demand and/or number of customers served [1,2], random order grouping [3] or duplicate ordering from several facilities [4]. In operations research, queuing-inventory system is a natural way to model load leveraging techniques, such as the leaky bucket congestion avoidance scheme used in a wide range of systems, from large scale routers [5] to electric vehicle charging stations [6].

Queuing inventory systems with positive service time were first investigated in Reference [7], followed by the work of Reference [8], in which an optimal quantity of inventory to be ordered to minimize the cost rate was obtained. In queuing inventory framework with positive service time, customers' queue is formed even when some inventory is available. We point the reader to a detailed survey of inventory systems with positive service time given in Reference [9], which includes classical, retrial, and production inventory.

Many retailers and banks find it helpful to partition the customers into different categories (classes) according to specific characteristics and adopt an inventory management



Citation: Rasmi, K.; Jacob, M.J.; Rumyantsev, A.; Krishnamoorthy, A. A Multi-Server Heterogeneous Queuing-Inventory System with Class-Dependent Inventory Access. *Mathematics* **2021**, *9*, 1037. https:// doi.org/10.3390/math9091037

Academic Editor: Vladimir V. Rykov

Received: 15 April 2021 Accepted: 28 April 2021 Published: 3 May 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). policy based on this differentiation strategy. In particular, long-term customers can be treated as high priority as compared to walk-in customers. Given a limited resource, the low priority customers may have to wait while some amount of resource is still available, reserved for customers with higher priority. The reservation may be made by imposing typical inventory levels for each class of customers according to their priority level. When the inventory comes below the level, customers from corresponding classes may have to wait until the inventory replenishment. Such a critical level policy is introduced and studied in Reference [10]. Different classes experience congestion with a single inventory; thus, customer sojourn times in a system are intrinsically correlated.

In most cases priorities are accompanied with heterogeneity of customer classes, e.g., in terms of service time distribution either in single-server [11,12] or in multi-server case [13–15]. As such, the arrival process also becomes heterogeneous, and the promising candidate is the so-called Marked Markovian Arrival Process (MMAP) used, e.g., in the works [11,14,16].

MMAP[K] is a generalization of Markov arrival processes (MAP) which have been studied and used extensively in queuing theory. MAP was introduced in Reference [17] to model non-Markovian point processes. While MAP is a useful tool to model point processes with one class of customers, MMAP[K] introduced by Neuts (see Reference [18]) is useful when multiple types of customers are present, while the model remains analytically tractable [19]. The basic characteristics of MMAP, such as peakedness of the arrival process, the first passage time to the arrival of an item of a specific type, and the behavior of the MMAP during that first passage, are analyzed in Reference [18].

In this paper, we analyze a multi-server queuing-inventory system with *K* classes of customers served from a single inventory which is managed according to the (s, S) with a positive lead time. Following Reference [11,14,16], we use MMAP[K] to model the arrival process. The servers are class-dependent, each server dedicated to one specific class of customers. Only the highest priority customers are allowed to wait in an infinite buffer. All other class customers can wait in respective finite buffers. The service for each class of customers is carried out with different exponential rates, and the inventory item is consumed at the end of service. A class-specific boundary level in the inventory is defined, causing customers of this specific class to wait for inventory replenishment when this boundary is down crossed. To the best of our knowledge, this model is new.

The structure of the paper is as follows. In Section 2, we give a detailed description of the model. The example of K = 2 is also given for better illustration of the model. In Section 3, an intuitive stability criterion for the system has been derived. Section 4 analyzes the steady state of the system and expresses a few important performance measures. In Section 5, a detailed analysis of the inventory recycle time has been carried out. Numerical illustrations are provided in Section 6, in which an optimization problem of practical importance has also been stated.

2. Model Description

We consider a multi-server queuing inventory system with heterogeneous customers. A *K*-server station provides service to *K* classes of customers. Arrivals occur according to the Marked Markovian Arrival Process (MMAP) driven by an irreducible continuous time Markov chain (CTMC) $\{Z(t)\}_{t\geq 0}$ with finite state space \mathcal{W} . Let $|\mathcal{W}| = W$. The sojourn time in each state $z \in \mathcal{W}$ is distributed exponentially with rate σ_z and d.f.

$$F_{z\hat{z}}(x) = 1 - e^{-\sigma_z x}.$$

At sojourn time expiration epoch, with probability $P^{(k)}(z, \hat{z})$ the process Z(t) moves from the state z to \hat{z} , generating a class k customer arrival, $z, \hat{z} \in W$, k = 0, 1, ..., K(conventionally, no customer arrival is generated if k = 0). Note that, for any $z \in W$,

$$\sum_{k=0}^{K} \sum_{\hat{z} \in \mathcal{W}} P^{(k)}(z, \hat{z}) = 1.$$
 (1)

Thus, the MMAP is characterized by a set of K + 1 *transition rate* square matrices D_0, \ldots, D_K of order *W* defined as follows:

$$(\mathbf{D}_k)_{z\hat{z}} = \begin{cases} \sigma_z P^{(k)}(z,\hat{z}), & z \neq \hat{z}, \\ -\sigma_z \mathbf{1}_{k=0}, & z = \hat{z}, \end{cases} \quad z, \hat{z} \in \mathcal{W}, \quad k = 0, 1, \dots, K,$$
(2)

where $1_{k=0}$ is the indicator of a non-random condition k = 0. The matrices D_k constitute a *generator* matrix D of the Markov process $\{Z(t)\}_{t>0}$ such that

$$\boldsymbol{D} = \sum_{k=0}^{K} \boldsymbol{D}_k$$

Recall that De = 0, which also follows from (1) and (2). Hereinafter, e(0, respectively) is the vector of ones (zeroes), and, if necessary, we designate the row vector with a transpose sign, e.g., 0'. As such, denoting by θ the *stochastic* invariant vector of D (i.e., $\theta D = 0$, that is, θ is the distribution of the corresponding Markov jump process governed by D), the class k arrival rate, λ_k , equals

$$\lambda_k = \theta D_k \mathbf{e}$$

We would like to note that the MMAP process being a versatile Markov process, is suitable for a wide range of applications due to variety of modeling features including correlated arrivals [20]. However, capturing sophisticated features, such as Long-Range Dependence (slow autocorrelation decay of the process that complicates application of the standard methods of performance estimation), may require infinite number of sources [21]. As such, it is necessary to balance the size *W* of the MMAP state space with practical capabilities of the model, including computational/storage capacity. Useful examples of MMAP processes may be found in Reference [19].

If the server *k* is busy, then the arriving class *k* customer can wait in a buffer space of capacity l_k , where $l_1 = \infty$ and $l_k < \infty$, k = 2, ..., K. The class *k* customer, finding the respective buffer completely filled on arrival, leaves the system forever.

The inventory is divided into K segments,

$$s_0 = 0 < s = s_1 < s_2 < \dots < s_K = S.$$

Class *k* customer spends an exponentially distributed, with rate μ_k , time at the (class-specific) server k = 1, ..., K. At *service completion* epoch, one item from a single common inventory is consumed by the customer. However, class *k* customers are served only if the inventory level exceeds s_{k-1} , k = 1, ..., K. Thus, *s* is not only inventory replenishment boundary, but also the critical level, at or below which only the class 1 customers (highest priority) are served. If the inventory level at service completion epoch is insufficient, the corresponding customer of class *k* repeats an independent service time with the same rate, μ_k , *until a service completion coincides with a sufficient inventory level*.

The inventory is replenished under (s, S) policy in an exponentially distributed *lead* time with rate γ . That is, when the inventory hits level s, a request for replenishment is issued, and the inventory level S is restored after an independent exponentially distributed random time with rate γ . Finally, we note that the arrival process, service times and replenishment times are *independent*. The structure of the system is presented on Figure 1.

Let $N_k(t)$ be the number of class k customers in the system, k = 1, ..., K, I(t) be the inventory level, $0 \le I(t) \le S$, and $Z(t) \in W$ be the phase of the MMAP, at time $t \ge 0$. Then, the considered system can be modeled by the regular irreducible CTMC

$$\zeta(t) = \{N_1(t), \dots, N_K(t), I(t), Z(t)\}, t \ge 0,$$

with state space $(n_1, n_2, ..., n_K, i, z), n_k \in \{0, ..., l_k\}, k = 1, ..., K, i \in \{0, ..., S\}$ and $z \in W$.



Figure 1. Structure of the system.

Fix the process $\zeta(t)$ at some state $(n_1, n_2, ..., n_k, i, z)$. Due to independence of the components of $\zeta(t)$, the transitions are now possible only to 2K + 2 states enumerated below:

- $(n_1, n_2, ..., n_j + 1_{n_j < l_j}, ..., n_K, i, \hat{z})$ with rate $(D_j)_{z\hat{z}}, j = 1, ..., K$ (arrival of class *j* customer);
- $(n_1, n_2, ..., n_j 1, ..., n_k, i 1, z)$ with rate μ_j , if $n_j > 0$ and $i > s_{j-1}, j = 1, ..., K$ (departure of class *j* customer);
- $(n_1, \ldots, n_K, i, \hat{z})$ with rate $(D_0)_{z\hat{z}}$ (MMAP phase switch);
- $(n_1, \ldots, n_K, S, \hat{z})$ with rate γ , if $i \leq s$ (inventory replenishment).

This gives a very special structure of the infinitesimal generator matrix of the considered process $\{\zeta(t)\}_{t\geq 0}$. First, since the component $N_1(t)$ corresponds to the unbounded buffer, while other components of $\zeta(t)$ are finite, and due to the fact that the transitions of $N_1(t)$ are *skip-free* in both directions (i.e., it may be incremented/decremented by at most one), the process $\zeta(t)$ is the so-called Quasi-Birth-Death (QBD) process with *level* $N_1(t)$ and *phase* $(N_2(t), \ldots, N_K(t), I(t), Z(t))$. Lexicographically ordered state space allows one to write the generator matrix in block-tridiagonal form

$$Q = \begin{vmatrix} A_1^{(0)} & A_0 & \mathbb{O} & \mathbb{O} & \dots \\ A_2 & A_1 & A_0 & \mathbb{O} & \dots \\ \mathbb{O} & A_2 & A_1 & A_0 & \dots \\ \vdots & \ddots & \ddots & \ddots & \ddots \end{vmatrix},$$
(3)

where the matrix \mathbb{O} is a zero block of corresponding dimension (we give the dimension explicitly if and when necessary).

Define an integer-valued function of two arguments, *i*, *j*, such that $2 \le i \le j \le K$:

$$\alpha(i,j) = \prod_{k=i}^{j} (l_k + 1),$$
(4)

and define $\alpha(i,j) = 0$ otherwise. Then, the blocks A_i , i = 0, 1, 2, and $A_1^{(0)}$ are square matrices having size

$$\alpha(2, K)(S+1)W.$$

The matrix A_0 consists of the rates of transitions corresponding to the arrival of class 1 customer, A_2 keeps the transition rates corresponding to departure of a class 1 customer, while A_1 is related to 2*K* remaining possible transitions, such that the *level* remains unaffected. It is quite straightforward to define A_0 , since, upon arrival of a class 1 customer, only the MMAP phase is switched; thus,

$$A_0 = I_{\alpha(2,K)(S+1)} \otimes D_1, \tag{5}$$

where \otimes is the Kronecker product and *I* is the identity matrix of corresponding dimension.

We will need the following notation: hereafter $\mathbf{e}_{i:j}^k$ is the (column) vector of dimension $k \ge 0$ with *i*th to *j*th components equal 1, and zero otherwise, $i \le j \le k$ (conventionally we take $\mathbf{e}_{i:i}^k \equiv 1$ if k = 0). To shorten the notation, we use for any $i \le k$

- $\mathbf{e}_{i}^{k} \equiv \mathbf{e}_{i,i}^{k}$ (a single non-zero component at *i*th row, dimension *k*);
- $\mathbf{e}_k \equiv \mathbf{e}_{k:k}^k$ (last non-zero component, dimension k);
- $\mathbf{e}^k \equiv \mathbf{e}_{1\cdot k}^k$ (all non-zero components of dimension *k*).

To define A_2 and for further use, we need to construct some auxiliary matrices. Define for j = 1, ..., K the following square matrix of order S + 1:

$$\boldsymbol{L}_{j} = \begin{bmatrix} \boldsymbol{0}' & \boldsymbol{0} \\ \operatorname{diag}(\boldsymbol{e}_{s_{j-1}+1,S}^{S}) & \boldsymbol{0} \end{bmatrix},$$
(6)

where the matrix L_j has non-zero entries below main diagonal only for rows $s_{j-1} + 2, ..., S + 1$, and the zero vectors are of size *S*. The block matrix L_j corresponds to possible transitions of the inventory for class *j* customer, and is indexed from 0 to *S*. In particular, for class 1 customers the matrix L_1 has the lower diagonal of ones (there is no constraint for the inventory to be decremented). Since, upon departure of class 1 customer, the *level*, as well as the inventory, are decremented by one, the matrix A_2 has the following form:

$$A_2 = I_{\alpha(2,K)} \otimes \mu_1 L_1 \otimes I_W. \tag{7}$$

Now, to define A_1 , which contains the transition rates related to arrivals and departures of class *j* customers, $j \ge 2$, inventory replenishment and MMAP phase change, we need auxiliary matrices of an increment/decrement in the corresponding (phase) component:

$$\boldsymbol{N}_{j}^{(+)} = \begin{bmatrix} \boldsymbol{0} & \boldsymbol{I}_{l_{j}} \\ \boldsymbol{0} & \boldsymbol{e}_{l_{j}}^{\prime} \end{bmatrix}, \quad \boldsymbol{N}_{j}^{(-)} = \begin{bmatrix} \boldsymbol{0}^{\prime} & \boldsymbol{0} \\ \boldsymbol{I}_{l_{j}} & \boldsymbol{0} \end{bmatrix}.$$
(8)

Note that $N_j^{(-)}$ is a square matrix of order $l_j + 1$ having lower diagonal of ones, while $N_i^{(+)}$ is semi-upper diagonal. This asymmetry will be explained below.

Using these constructions, it is rather straightforward to define the transition rates of all possible transitions constituting the matrix A_1 :

$$A_1 = A_1^{(a)} + A_1^{(b)} + A_1^{(c)} - \Delta, \text{ where}$$
(9)

$$\boldsymbol{A}_{1}^{(a)} = \sum_{j=2}^{K} \boldsymbol{I}_{\alpha(2,j-1)} \otimes \left(\boldsymbol{N}_{j}^{(+)} \otimes \boldsymbol{I}_{\alpha(j+1,K)} \otimes \boldsymbol{I}_{S+1} \otimes \boldsymbol{D}_{j} \right)$$
(10)

$$+ N_{j}^{(-)} \otimes I_{\alpha(j+1,K)} \otimes \mu_{j} L_{j} \otimes I_{W} \bigg), \qquad (11)$$

$$\boldsymbol{A}_{1}^{(b)} = \boldsymbol{I}_{\alpha(2,K)} \otimes \gamma \boldsymbol{e}_{1:s+1}^{S+1} \boldsymbol{e}_{S+1}' \otimes \boldsymbol{I}_{W},$$
(12)

$$A_1^{(c)} = I_{\alpha(2,K)} \otimes I_{S+1} \otimes D_0.$$
(13)

Conventionally, in (10) and (11), we define the zero-size identity matrix $I_0 = 1$. Note that (10) corresponds to arrival of class *j* customer, (11) is the corresponding class departure, (12) is a replenishment (where in fact the corresponding vector product gives a matrix with only last column being non-zero), while (13) is the MMAP phase change. It is worth noting that asymmetry in $N_j^{(+)}$ (non-zero last row), defined in (8), is used in (10) to indicate an arrival of a class *j* customer that is lost, while the MMAP phase is changed according to D_j . The matrix Δ is a diagonal matrix that guarantees ($A_0 + A_1 + A_2$) $\mathbf{e} = \mathbf{0}$; thus,

$$\mathbf{\Delta} = \text{diag}\Big[\Big(A_0 + A_1^{(a)} + A_1^{(b)} + A_1^{(c)} + A_2\Big)\mathbf{e}\Big].$$
(14)

Straightforward algebra allows to obtain Δ , the diagonal matrix of dimension $\alpha(2, K)(S+1)W$, from (14) in a closed form as follows:

$$\mathbf{\Delta} = ext{diag} \Big[\mathbf{\delta}^{(a)} + \mathbf{\delta}^{(b)} + \mathbf{\delta}^{(d)} \Big]$$
 ,

where $\delta^{(a)}$ is the vector of transition rates due to departure of customers of classes 2, ..., *K*, $\delta^{(d)}$ contains transition rates due to a class 1 customer departure, and $\delta^{(b)}$ is the vector of transition rates due to replenishment given as follows:

$$\boldsymbol{\delta}^{(a)} = \sum_{j=2}^{K} \mu_j \mathbf{e}^{\alpha(2,j-1)} \otimes \mathbf{e}^{l_j+1}_{2:l_j+1} \otimes \mathbf{e}^{\alpha(j+1,K)} \otimes \mathbf{e}^{S+1}_{s_{j-1}+2:S+1} \otimes \mathbf{e}^{W}, \tag{15}$$

$$\boldsymbol{\delta}^{(b)} = \gamma \mathbf{e}^{\alpha(2,K)} \otimes \mathbf{e}_{1:s+1}^{S+1} \otimes \mathbf{e}^{W}, \tag{16}$$

$$\boldsymbol{\delta}^{(d)} = \mu_1 \mathbf{e}^{\boldsymbol{\alpha}(2,K)} \otimes \mathbf{e}_{2:S+1}^{S+1} \otimes \mathbf{e}^W.$$
(17)

It remains to define the block $A_1^{(0)}$ corresponding to possible transitions of the model from within the states having zero class 1 customers. Note that such a matrix is very similar to the matrix A_1 , and the only difference is in the diagonal balancing matrix (14). Indeed, from the condition $(A_1^{(0)} + A_0)\mathbf{e} = \mathbf{0}$, define

$$\Delta_0 = \text{diag}\Big[\Big(A_0 + A_1^{(a)} + A_1^{(b)} + A_1^{(c)}\Big)\mathbf{e}\Big].$$
(18)

The matrix $A_1^{(0)}$ is then defined as follows:

$$A_1^{(0)} = A_1^{(a)} + A_1^{(b)} + A_1^{(c)} - \Delta_0,$$
(19)

where it follows from (14) and (18) that $\Delta_0 = \text{diag} \left[\delta^{(a)} + \delta^{(b)} \right]$.

We note that the definitions of subblocks may be rewritten in a more compact form using Kronecker sums. Moreover, by defining unbounded analogs of the matrices given in (8), the generator matrix itself may be rewritten similarly. However, we skip this possibility to keep parsimony of the notation.

To illustrate (3), we consider the case K = 2. Recall that class 1 customers have priority over class 2 customers. Class 1 customers are allowed to wait in an infinite buffer, whereas class 2 customers cannot enter the system if there are l_2 class 2 customers already in the system. Server 1 serves class 1 customers with rate μ_1 and server 2 serves class 2 customers with rate μ_2 . Customers are served with an inventory from a common source running according to (s, S) policy with exponential lead time (rate γ). Even if the server 2 is free, class 2 customers will be served only if the inventory level is at least s + 1.

We consider the CTMC $\zeta_t = \{N_1(t), N_2(t), I(t), Z(t)\}, t \ge 0$ with state space (n_1, n_2, i, z) , where $n_1 \ge 0, 0 \le n_2 \le l_2, 0 \le i \le S, z \in W$. Since K = 2, it follows from (4) that $\alpha(2, K) = l_2 + 1$. Corresponding to each *level* n_1 , there will be $(l_2 + 1)(S + 1)W$ *phase states*. The infinitesimal generator matrix of ζ_t is given by (3), where the non zero blocks $A_1^{(0)}, A_0, A_1$, and A_2 are of size $(l_2 + 1)(S + 1)W$ and have the following forms:

$$egin{aligned} &A_0 = I_{(l_2+1)(S+1)} \otimes D_1, \ &A_2 = I_{l_2+1} \otimes \mu_1 L_1 \otimes I_N. \end{aligned}$$

where L_1 is a matrix of order S + 1 with a non-zero lower diagonal, defined as

$$\boldsymbol{L}_1 = \begin{bmatrix} \boldsymbol{0}' & \boldsymbol{0} \\ \boldsymbol{I}_S & \boldsymbol{0} \end{bmatrix}$$

The subblocks A_1 , $A_1^{(0)}$ also have the block-tridiagonal structure with $l_2 + 1$ blocks over the main diagonal indexed by the number of class 2 customers:

$$A_{1}^{(0)} = \begin{bmatrix} A_{1}^{(2)} & A_{0}^{(1)} & \mathbb{O} & \mathbb{O} & \dots & \mathbb{O} \\ A_{2}^{(1)} & A_{1}^{(1)} & A_{0}^{(1)} & \mathbb{O} & \dots & \mathbb{O} \\ \mathbb{O} & A_{2}^{(1)} & A_{1}^{(1)} & A_{0}^{(1)} & \ddots & \mathbb{O} \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \dots & \dots & A_{1}^{(1)} & A_{0}^{(1)} \end{bmatrix}, \quad A_{1} = \begin{bmatrix} A_{1}^{(4)} & A_{0}^{(1)} & \mathbb{O} & \mathbb{O} & \dots & \mathbb{O} \\ A_{2}^{(1)} & A_{1}^{(3)} & A_{0}^{(1)} & \mathbb{O} & \dots & \mathbb{O} \\ \mathbb{O} & A_{2}^{(1)} & A_{1}^{(3)} & A_{0}^{(1)} & \ddots & \mathbb{O} \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \dots & \dots & A_{1}^{(1)} & A_{0}^{(1)} \end{bmatrix}, \quad A_{1} = \begin{bmatrix} A_{1}^{(4)} & A_{0}^{(1)} & \mathbb{O} & \dots & \mathbb{O} \\ A_{2}^{(1)} & A_{1}^{(3)} & A_{0}^{(1)} & \mathbb{O} & \dots & \mathbb{O} \\ \mathbb{O} & A_{2}^{(1)} & A_{1}^{(3)} & A_{0}^{(1)} & \ddots & \mathbb{O} \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \dots & \dots & A_{1}^{(3)} & A_{0}^{(1)} \\ \vdots & \dots & \dots & A_{2}^{(1)} & A_{1}^{(3)} \end{bmatrix}.$$

Submatrices $A_0^{(1)}, A_2^{(1)}, A_1^{(i)}, i = 1, ..., 4$, of order (S + 1)W, are given below.

$$egin{aligned} &A_0^{(1)} = I_{S+1} \otimes D_2, \ &A_2^{(1)} = \mu_2 L_2 \otimes I_W, \end{aligned}$$

where L_2 has non-zero entries in the lower diagonal from s + 2nd row onward,

	Γ0	0	• • •	0	•••	0	0	
$L_2 =$: 0	: 0	 	: 1	 	: 0	: 0	
		: 0	:	: 0	·	: 1	:	

It remains to define $A_1^{(i)}$, i = 1, ..., 4. These matrices have semi-block-diagonal structure with main diagonal and last column containing S + 1, possibly non-zero, blocks, indexed by the number of inventory items available, each block being a square matrix of order W. Indeed, since the departures of class 1 customers are not possible from *level* 0,

$$\boldsymbol{A}_{1}^{(2)} = \boldsymbol{I}_{S+1} \otimes \boldsymbol{D}_{0} - \begin{bmatrix} \boldsymbol{\gamma} \boldsymbol{I}_{W} & \boldsymbol{\mathbb{O}} & \dots & \boldsymbol{\mathbb{O}} & \dots & -\boldsymbol{\gamma} \boldsymbol{I}_{W} \\ \vdots & \ddots & \dots & \dots & \vdots \\ \boldsymbol{\mathbb{O}} & \dots & \boldsymbol{\gamma} \boldsymbol{I}_{W} & \boldsymbol{\mathbb{O}} & \dots & -\boldsymbol{\gamma} \boldsymbol{I}_{W} \\ \boldsymbol{\mathbb{O}} & \dots & \boldsymbol{\mathbb{O}} & \boldsymbol{\mathbb{O}} & \dots & \boldsymbol{\mathbb{O}} \\ \vdots & \dots & \dots & \dots & \vdots \\ \boldsymbol{\mathbb{O}} & \dots & \dots & \boldsymbol{\mathbb{O}} & \dots & \boldsymbol{\mathbb{O}} \end{bmatrix},$$

where, from s + 1st row onward, the diagonal elements are D_0 (only MMAP phase change is possible). Similarly, since $A_1^{(1)}$ corresponds to states with positive number of class 2 customers,

$$A_{1}^{(1)} = A_{1}^{(2)} - \begin{bmatrix} 0 & \dots & 0 & \dots & 0 \\ \vdots & \dots & \vdots & \dots & \vdots \\ 0 & \dots & \mu_{2}I_{W} & \dots & 0 \\ \vdots & \dots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & \mu_{2}I_{W} \end{bmatrix}$$

Since A_1 corresponds to positive levels of the process ζ_t , i.e., the number of class 1 customers is positive,

$$A_1^{(4)} = A_1^{(2)} - I_{S+1}^{(0)} \otimes \mu_1 I_W, \quad A_1^{(3)} = A_1^{(1)} - I_{S+1}^{(0)} \otimes \mu_1 I_W,$$

where $I_{S+1}^{(0)}$ is a square matrix of order S + 1 consisting of a zero column, zero row, and identity matrix as follows:

$$\boldsymbol{I}_{S+1}^{(0)} = \begin{bmatrix} \boldsymbol{0} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{I}_S \end{bmatrix}$$

Indeed, the matrices $A_1^{(4)}$ and $A_1^{(3)}$ correspond to positive number of class 1 customers; thus, their diagonals include the service rates of class 1 customers, except the case of the empty inventory. As a final note, the last diagonal blocks in matrices, $\hat{A}_1^{(1)}$ and $\hat{A}_1^{(3)}$, correspond to arrivals of class 2 customers that are lost; hence,

$$\widehat{A}_{1}^{(i)} = A_{1}^{(i)} + A_{0}^{(1)}, \quad i = 1, 3.$$

3. Stability Condition

The necessary and sufficient condition for existence of the non-zero steady-state probability is the specific version of Foster ergodicity condition known as the Neuts ergodicity criterion [22],

$$\pi A_0 \mathbf{e} < \pi A_2 \mathbf{e},\tag{20}$$

where the stochastic vector π is the solution of the system

$$\pi A = \mathbf{0},\tag{21}$$

and

$$A = A_0 + A_1 + A_2.$$

Note that, due to the properties of A_i , i = 0, 1, 2, the matrix A is a generator matrix of a finite state space CTMC giving the projection of the phase transition at high levels; thus, the vector π may be considered as the steady-state probability of the phase at high levels [22]. It now follows from (5), (7), and (9) that A has block-tridiagonal structure indexed by the number of class-2 customers as follows:

A =	$\begin{bmatrix} \mathbf{R}_1^{(0)} \\ \mathbf{R}_2 \\ 0 \end{bmatrix}$	$egin{array}{c} m{R}_0 \ m{R}_1 \ m{R}_2 \end{array}$	$\begin{array}{c} 0 \\ \boldsymbol{R}_0 \\ \boldsymbol{R}_1 \end{array}$	$0\\0\\\boldsymbol{R}_0$	 	0 0 0	
	÷	÷	·	·	·	÷	
	0	0	•••	R_2	R_1	R_0	
	0	0		0	R_2	\widehat{R}_1	

The blocks $\mathbf{R}_1^{(0)}$, $\hat{\mathbf{R}}_1$, and \mathbf{R}_i , i = 0, 1, 2, are square matrices, and their structure follows from (5), (7), and (9). Indeed,

$$\begin{split} \mathbf{R}_{0} = & \mathbf{I}_{\alpha(3,K)} \otimes \mathbf{I}_{S+1} \otimes \mathbf{D}_{2}, \\ \mathbf{R}_{2} = & \mathbf{I}_{\alpha(3,K)} \otimes \mu_{2} \mathbf{L}_{2} \otimes \mathbf{I}_{W}, \\ \mathbf{R}_{1} = & \mathbf{I}_{\alpha(3,K)} \otimes \left(\left(\gamma \mathbf{e}_{1:s+1}^{S+1} \mathbf{e}_{S+1}' + \mu_{1} \mathbf{L}_{1} \right) \oplus (\mathbf{D}_{0} + \mathbf{D}_{1}) \right) \\ &+ \sum_{j=3}^{K} \mathbf{I}_{\alpha(3,j-1)} \otimes \left(\mathbf{N}_{j}^{(+)} \otimes \mathbf{I}_{\alpha(j+1,K)} \otimes \mathbf{I}_{S+1} \otimes \mathbf{D}_{j} \\ &+ \mathbf{N}_{j}^{(-)} \otimes \mathbf{I}_{\alpha(j+1,K)} \otimes \mu_{j} \mathbf{L}_{j} \otimes \mathbf{I}_{W} \right) - \hat{\mathbf{\Delta}}, \end{split}$$

where \oplus is the Kronecker sum defined as $A \oplus B = I \otimes B + A \otimes I$; $\widehat{\Delta}$ being a diagonal matrix that guarantees $(R_0 + R_1 + R_2)\mathbf{e} = \mathbf{0}$, and obvious convention $\sum_{j=3}^{K} = 0$ for K < 3 is used. It remains to note that, since $R_1^{(0)}$ corresponds to states with no class-2

customers, $\mathbf{R}_1^{(0)} = \mathbf{R}_1 + \text{diag}[\mathbf{R}_2\mathbf{e}]$, while the matrix $\hat{\mathbf{R}}_1 = \mathbf{R}_1 + \mathbf{R}_0$ includes the lost arrivals of class-2 customers. Recall that the losses of class 2 customers appear here due to the specific indexing of the blocks of matrix \mathbf{A} , and corresponding losses of classes 3, . . . , K are encoded in the components of matrix \mathbf{R}_1 . It is worthwhile to note that, despite the expected independence of the MMAP component evolution on the system state, the changes in the phase occur upon arrival of customers of classes 2, . . . , K, both in the corresponding queue size and in MMAP phase; similar simultaneous changes occur upon departure of such customers (both counter and inventory). In the case K = 2, the subblocks of A have the following clear structure:

$$\begin{aligned} \mathbf{R}_0 &= \mathbb{O} \oplus \mathbf{D}_2, \quad \mathbf{R}_2 = \mu_2 \mathbf{L}_2 \oplus \mathbb{O}, \\ \mathbf{R}_1 &= \left(\gamma \mathbf{e}_{1:s+1}^{S+1} \mathbf{e}_{S+1}' + \mu_1 \mathbf{L}_1\right) \oplus \left(\mathbf{D}_0 + \mathbf{D}_1\right) - \hat{\boldsymbol{\Delta}}, \\ \widehat{\mathbf{R}}_1 &= \left(\gamma \mathbf{e}_{1:s+1}^{S+1} \mathbf{e}_{S+1}' + \mu_1 \mathbf{L}_1\right) \oplus \mathbf{D} - \hat{\boldsymbol{\Delta}}. \end{aligned}$$

Note that the Kronecker sums in the phase transition rate matrices appear due to the multidimensional structure of the phase state space. These sums highlight the independent changes of one of the components of the phase vector, and \mathbb{O} is used if the corresponding component remains unchanged due to a transition.

To solve the system (21), the following numerically stable algorithm is suggested. Let the row vector π be presented in the form $\pi = (\pi_0, \pi_1, ..., \pi_{l_2})$. The vectors $\pi_m, m = 0, ..., l_2$, are considered to have the following form:

$$\pi_m = \pi_{m-1} \boldsymbol{U}_{m-1} = \pi_0 \prod_{k=1}^m \boldsymbol{U}_{k-1}, \quad m = 1, \dots, l.$$
 (22)

The matrices U_m , $m = 0, ..., l_2 - 1$, are obtained from the system (21) using (22), starting from the last column, i.e., calculated using the backward recursion,

$$\boldsymbol{U}_{m} = -\boldsymbol{R}_{0}(\boldsymbol{R}_{1} + \boldsymbol{U}_{m+1}\boldsymbol{R}_{2})^{-1}, \quad m = l_{2} - 2, \dots, 0,$$
 (23)

under the initial condition (following from the last column of the matrix A)

$$\boldsymbol{U}_{l_2-1} = -\boldsymbol{R}_0 \left(\widehat{\boldsymbol{R}}_1 \right)^{-1}.$$

Note that, since *A* is a generator matrix, \hat{R}_1 is diagonally dominant and, hence, invertible. Finally, the vector π_0 is the unique solution to the system (obtained from the first column of *A* and the fact that π is stochastic)

$$\begin{cases} \boldsymbol{\pi}_0 \left(\boldsymbol{R}_1^{(0)} + \boldsymbol{U}_0 \boldsymbol{R}_2 \right) = \boldsymbol{0}, \\ \boldsymbol{\pi}_0 \left(\boldsymbol{e} + \sum_{m=1}^{l_2} \prod_{k=1}^m \boldsymbol{U}_{k-1} \boldsymbol{e} \right) = 1. \end{cases}$$
(24)

Now, let us consider the l.h.s. of (20). The vector π is the steady-state probability vector of the finite state space CTMC defined by the matrix *A* and having the state space

$$\mathcal{E} = \{ (n_2, \dots, n_K, i, z), n_k \in \{0, \dots, l_k\}, k = 2, \dots, K, i \in \{0, \dots, S\}, z \in \mathcal{W} \}.$$
(25)

However, it is possible to shrink the state space \mathcal{E} into the following subsets defined for all $z \in \mathcal{W}$:

$$\mathcal{E}_z = \{(n_2,\ldots,n_K,i,z_0) \in \mathcal{E} : z_0 = z\}.$$

Now, we show that, for each state from \mathcal{E}_z , the transition rate to the set $\mathcal{E}_{\hat{z}}$ (i.e., the sum of corresponding transition rates to individual states) equals $(D)_{z\hat{z}}$. Indeed, to obtain this transition rate, the matrix A needs to be multiplied by the matrix that would sum up

the corresponding transition rates, that is, $e^{\alpha(2,K)(S+1)} \otimes I_W$. As such, using the fact that, for any matrices B_i , i = 1, ..., 4, the following equation holds good,

$$(B_1 \otimes A_2)(B_3 \otimes B_4) = B_1 B_3 \otimes B_2 B_4, \tag{26}$$

it can be obtained from (5) that

$$A_{0}(\mathbf{e}^{\alpha(2,K)(S+1)} \otimes \mathbf{I}_{W}) = (\mathbf{I}_{\alpha(2,K)(S+1)} \otimes \mathbf{D}_{1})(\mathbf{e}^{\alpha(2,K)(S+1)} \otimes \mathbf{I}_{W}) = \mathbf{e}^{\alpha(2,K)(S+1)} \otimes \mathbf{D}_{1}.$$
 (27)

It follows from (7) and (17) that

$$(\mathbf{A}_{2} - \operatorname{diag}(\boldsymbol{\delta}^{(d)}))(\mathbf{e}^{\alpha(2,K)(S+1)} \otimes \mathbf{I}_{W}) = (\mathbf{I}_{\alpha(2,K)} \otimes \mu_{1} \mathbf{L}_{1} \otimes \mathbf{I}_{W})(\mathbf{e}^{\alpha(2,K)} \otimes \mathbf{e}^{S+1} \otimes \mathbf{I}_{W}) - \mu_{1} \mathbf{e}^{\alpha(2,K)} \otimes \mathbf{e}^{S+1}_{2:S+1} \otimes \operatorname{diag}(\mathbf{e}^{W})) = \mathbf{0}.$$
(28)

Finally, it follows from (9), (15), (16), after some algebra, that

$$\left(A_{1}^{(a)} + A_{1}^{(b)} + A_{1}^{(c)} - \operatorname{diag}(\delta^{(a)} + \delta^{(b)}) \right) (\mathbf{e}^{\alpha(2,K)(S+1)} \otimes \mathbf{I}_{W})$$

$$= \mathbf{e}^{\alpha(2,K)(S+1)} \otimes \left(\mathbf{D}_{0} + \sum_{j=2}^{K} \mathbf{D}_{j} \right).$$
(29)

Then, it follows from (27)–(29) that

$$A(\mathbf{e}^{\alpha(2,K)(S+1)} \otimes \mathbf{I}_W) = \mathbf{e}^{\alpha(2,K)(S+1)} \otimes \mathbf{D}.$$

Thus, the subsets $\mathcal{E}_z, z \in \mathcal{W}$, are the states of a CTMC defined by the matrix D. Now, we note that

$$\pi A(\mathbf{e}^{\alpha(2,K)(S+1)}\otimes I_W)=\pi(\mathbf{e}^{\alpha(2,K)(S+1)}\otimes I)D;$$

hence, $\theta = \pi(\mathbf{e}^{\alpha(2,K)(S+1)} \otimes \mathbf{I})$. Finally, from (5), obtain

$$\pi A_0 \mathbf{1} = \boldsymbol{\theta} D_1 \mathbf{1}^{W} = \lambda_1$$

Now, taking into account (7), (22), (24), and using the property (26) for the unit vector **e**, the r.h.s. of the stability condition (20) becomes

$$\sum_{m=0}^{l_2} \pi_m \Big(\mathbf{I}_{\alpha(3,K)} \otimes \mu_1 \mathbf{L}_1 \otimes \mathbf{I}_W \Big) \mathbf{e} = \sum_{m=0}^{l_2} \pi_m \mathbf{e}^{\alpha(3,K)} \otimes \mu_1 \mathbf{e}_{2:S+1}^{S+1} \otimes \mathbf{e}^W.$$

Finally, (20) becomes

$$\rho := \lambda_1 \left[\sum_{m=0}^{l_2} \pi_m \mathbf{e}^{\alpha(3,K)} \otimes \mu_1 \mathbf{e}_{2:S+1}^{S+1} \otimes \mathbf{e}^W \right]^{-1} < 1.$$
(30)

Note that the stability condition (30) has a nice interpretation, since λ_1 is the upward drift, while the sum in brackets is the mean downward drift of class 1 customers at high levels obtained by aggregation of the components phase probabilities vector π corresponding to the states with positive departure probability of class 1 customers.

4. Steady State Performance

If the stability condition (30) holds, the stochastic vector of stationary probabilities, q, exists and is the solution of the steady-state equation involving the infinitesimal generator Q given in (3):

$$qQ = 0. (31)$$

Note that the components of the vector q are ordered lexicographically, that is, for $n_k \in \{0, ..., l_k\}, k = 1, ..., K, i \in \{0, ..., S\}$ and $z \in W$,

$$q(n_1,...,n_K,i,z) = \lim_{t\to\infty} P(N_1(t) = n_1,...,N_K(t) = n_K, I(t) = i, Z(t) = z).$$

Since the matrix Q defines a level-independent QBD, the soluton can be obtained in matrix-geometric form [22]. Indeed, splitting the vector q into finite vectors q_0, q_1, \ldots by the (value of the) first coordinate, we assume that

$$\boldsymbol{q}_i = \boldsymbol{q}_0 \boldsymbol{R}^i. \tag{32}$$

The matrix R is the so-called rate matrix, being the minimal non-negative solution of the matrix quadratic equation (which follows from (31) using the block-tridiagonal structure of Q and (32)),

$$R^2 A_2 + R A_1 + A_0 = \mathbf{0}. ag{33}$$

The boundary conditions for obtaining q_0 follow from the first block-column of Q and are [23]

$$q_0(A_1^{(0)} + RA_2) = \mathbf{0}, \tag{34}$$

$$q_0(I-R)^{-1}\mathbf{1} = 1. (35)$$

However, in order to avoid numerical difficulties, it is recommended to use the alternative matrix quadratic equation for the substochastic matrix G being the minimal non-negative solution of the system

$$A_0 G^2 + A_1 G + A_2 = \mathbf{0}, (36)$$

obtain the *R* matrix by the known relation [23]

$$R = A_0(-A_1 - A_0 G)^{-1},$$

and, finally, use (32) to obtain the vector q.

After obtaining the steady-state probability vector, it is straightforward to define the steady-state performance measures of interest. To do so, we use an auxiliary vector $j^{(n)}$ of size n + 1 containing the sequence (0, ..., n), i.e.,

$$\mathbf{j}^{(n)} = (0, \ldots, n)$$

Average number of class 1 customers in the system:

$$\mathbf{E}_{1} = \sum_{i=1}^{\infty} i q_{i} \mathbf{e}^{\alpha(2,K)(S+1)W} = q_{0} R (I-R)^{-2} \mathbf{e}^{\alpha(2,K)(S+1)W}.$$
(37)

Average number of class k = 2, ..., K customers in the system:

$$\mathbf{E}_{\mathbf{k}} = \sum_{i=0}^{\infty} q_i \mathbf{e}^{\alpha(2,k-1)} \otimes j^{(l_k)} \otimes \mathbf{e}^{\alpha(k+1,K)} \otimes \mathbf{e}^{(S+1)W} \\
= q_0 (I - \mathbf{R})^{-1} \mathbf{e}^{\alpha(2,k-1)} \otimes j^{(l_k)} \otimes \mathbf{e}^{\alpha(k+1,K)} \otimes \mathbf{e}^{(S+1)W} .$$
(38)

Average inventory size in the system:

$$\bar{\mathbf{I}} = \sum_{i=0}^{\infty} q_i \mathbf{e}^{\alpha(2,K)} \otimes j^{(S)} \otimes \mathbf{e}^W = q_0 (I - R)^{-1} \mathbf{e}^{\alpha(2,K)} \otimes j^{(S)} \otimes \mathbf{e}^W.$$
(39)

Replenishment rate, where $\delta^{(b)}$ is given in (16),

$$\bar{\mathbf{R}} = \sum_{i=0}^{\infty} q_i \delta^{(b)} = q_0 (I - R)^{-1} \delta^{(b)}.$$
(40)

Class *k* loss rate, $k = 2, \ldots, K$:

$$\mathbf{L}_{k} = \lambda_{k} \sum_{i=0}^{\infty} \boldsymbol{q}_{i} \mathbf{e}^{\alpha(2,k-1)} \otimes \mathbf{e}_{l_{k}+1} \otimes \mathbf{e}^{\alpha(k+1,K)} \otimes \mathbf{e}^{(S+1)W} = \lambda_{k} (\boldsymbol{I} - \boldsymbol{R})^{-1} \mathbf{e}^{\alpha(2,k-1)} \otimes \mathbf{e}_{l_{k}+1} \otimes \mathbf{e}^{\alpha(k+1,K)} \otimes \mathbf{e}^{(S+1)W}.$$
(41)

We note that the average number of queued customers can be obtained by the corresponding formula for the average number of customers in the system by replacing *i* to (i - 1). Finally, the corresponding sojourn times of class $k \ge 1$ customer are obtained by Little's law as follows:

$$\mathbf{E_{s_k}} = \frac{\mathbf{E_k}}{\hat{\lambda}_k},$$

where $\hat{\lambda}_k$ is the effective arrival rate of class *k*, i.e., $\hat{\lambda}_1 = \lambda_1$ and

$$\hat{\lambda}_k = \lambda_k - \mathbf{L}_k.$$

5. Analysis of Inventory Recycle Time

Starting with inventory level *S*, the time taken to hit the level *S* again is called inventory cycle time, say, Γ . The distribution of inventory cycle time depends on the number of customers of all classes in the system and the phase of the arrival process. However, we note that, since the replenishment happens only at the event when the inventory hits the level *s*, the inventory cycle is essentially the lead time plus the time it takes to reach *s* from *S* by the steps decreasing the inventory, made by the process $\zeta(t)$. It is clear that such a process is a time until absorption of the process, where the absorption happens at the level *s*. It is known that such a time can be modeled by a phase-type distribution (for details on this type of distributions, see Reference [24]), and below we obtain the parameters of such a distribution.

Indeed, consider the inventory level I(t) = S. If $N_k(t) \ge S - s_{k-1}$, then the class k customers present in the system are capable of consuming all the inventory allowed for such a class, either before stopping service for this specific class at the boundary level s_{k-1} , or before hitting the absorbing inventory state I(t) = s, and in such case the time to absorption does not depend on the arrivals of class k customers after time t. Otherwise, arrivals of class k customers can be tracked until the condition $N_k(t) \ge S - s_{k-1}$ is met (if this happens before absorption). Moreover, if the level $N_k(t)$ hits the value $S - s_{k-1}$ from below, the departures of class k customers need not be tracked, but instead, only the inventory decreasing process should take into account the corresponding rate of the class k customer service, μ_k . Thus, the time it takes I(t) to reach s (time to absorption) can be modeled by a finite state space absorbing CTMC, which is essentially a restriction of $\zeta(t)$ to the set

$$\{0, \ldots, S-s\} \times \{0, \ldots, S-s\} \times \{0, \ldots, S-s_2\} \times \ldots \{0, \ldots, S-s_{K-1}\} \times \{s+1, \ldots, S\} \times \mathcal{W}.$$

This means that the time to reach inventory level *s* has a phase-type distribution $PH(\beta, T)$, where β is the initial distribution of the restricted chain $(\hat{N}_1(t), \dots, \hat{N}_k(t), \hat{I}(t), Z(t))$ at time 0, and *T* is the transition matrix which follows from *Q* and the aforementioned restrictions. In particular, it is straightforward to define the transition rate matrix *T* as follows.

$$T = \begin{bmatrix} B_{1}^{(0)} & B_{0} & \mathbb{O} & \mathbb{O} & \dots & \mathbb{O} \\ B_{2} & B_{1} & B_{0} & \mathbb{O} & \dots & \mathbb{O} \\ \mathbb{O} & B_{2} & B_{1} & B_{0} & \dots & \mathbb{O} \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \mathbb{O} & \dots & \dots & \mathbb{O} & \widehat{B}_{1} \end{bmatrix}.$$
(42)

Note that *T* is a block-tridiagonal finite matrix, with S - s + 1 blocks on the main diagonal, where the first block corresponds to states with $N_1(t) = 0$, while the last block is for $N_1(t) = S - s$, respectively. The rows corresponding to the last block have zeroes except the block on the main diagonal, since upon reaching the level $N_1(t) = S - s$, neither the departures, nor the arrivals, of class 1 customers are tracked, but the rate μ_1 is taken into account in the value of the inventory decreasing rate. Moreover, the inventory states become in fact re-numbered so as to have S - s states starting from inventory level s + 1, numbered sequentially, that is, compared to the original inventory component I(t), it holds that

$$\hat{I}(t) = I(t) - s, \tag{43}$$

and the absorption happens when $\hat{I}(t)$ makes a transition to 0. The size of the blocks is

$$\hat{\alpha}(2,K)(S-s)W,$$

where, for $2 \le i \le j \le K$,

$$\hat{\alpha}(i,j) = \prod_{k=i}^{j} (S - s_{k-1} + 1), \tag{44}$$

and $\hat{\alpha}(i, j) = 0$ otherwise.

It is straightforward to see

$$\boldsymbol{B}_0 = \boldsymbol{I}_{\hat{\boldsymbol{\alpha}}(2,K)(S-s)} \otimes \boldsymbol{D}_1, \tag{45}$$

and

$$\mathbf{B}_2 = \mathbf{I}_{\hat{\alpha}(2,K)} \otimes \mu_1 \widehat{L}_1 \otimes \mathbf{I}_W, \tag{46}$$

where \hat{L}_j is obtained from the matrix L_j defined in (6) by reducing the latter to rows and columns from $s_{i-1} + 2$ to S + 1, i.e.,

$$L_j = J_{s_{j-1}+1,S+1}L_jJ'_{s_{j-1}+1,S+1},$$

where the matrix $J_{a,b}$ removes the first *a* rows, $1 \le a \le b$, that is,

$$\boldsymbol{J}_{a,b} = \begin{bmatrix} \mathbb{O}_{(b-a) \times a} & \boldsymbol{I}_{b-a} \end{bmatrix},$$

with an exception for \widehat{L}_1 defined as

$$\hat{L}_1 = J_{s+1,S+1} L_1 J'_{s+1,S+1}$$

Let us restrict the matrices $N_j^{(+)}$ and $N_j^{(-)}$ defined in (8) to the finite state space of the absorbing CTMC, and modify the last row of $N_j^{(-)}$ so as to force $\hat{N}_j(t)$ to stay at the level $S - s_{j-1}$ once reached, $j \ge 2$. To do so, denote

$$H_{a,b} = \begin{bmatrix} I_a & \mathbb{O}_{a \times (b-a)} \end{bmatrix}.$$

Multiplication by $H_{a,b}$ on the left leaves only the first *a* rows in the resulting matrix and removes the last b - a rows, while multiplication on the right by $H'_{a,b}$ removes the last b - a columns. Then, the matrices $\widehat{N}_{j}^{(+)}$ and $\widehat{N}_{j}^{(-)}$ are defined as follows:

$$\widehat{N}_{j}^{(+)} = H_{S-s_{j-1}+1,S+1}N_{j}^{(+)}H'_{S-s_{j-1}+1,S+1} + \operatorname{diag}(\mathbf{e}_{S-s_{j-1}+1}),$$
$$\widehat{N}_{j}^{(-)} = H_{S-s_{j-1}+1,S+1}N_{j}^{(-)}H'_{S-s_{j-1}+1,S+1}\operatorname{diag}(\mathbf{e}^{S-s_{j-1}+1} - \mathbf{e}_{S-s_{j-1}+1}) + \operatorname{diag}(\mathbf{e}_{S-s_{j-1}+1}),$$

where the *S* + 1-dimensional vector \mathbf{e}_{S+1} is zero vector except the last component equal to one, while \mathbf{e}^{S+1} is the *S* + 1-dimensional vector of ones. Now, we are ready to define the matrix B_1 as follows:

$$\begin{split} \boldsymbol{B}_{1} &= \sum_{j=2}^{K} \boldsymbol{I}_{\widehat{\alpha}(2,j-1)} \otimes \left(\widehat{\boldsymbol{N}}_{j}^{(+)} \otimes \boldsymbol{I}_{\widehat{\alpha}(j+1,K)} \otimes \boldsymbol{I}_{S-s} \otimes \boldsymbol{D}_{j} \right. \\ &+ \widehat{\boldsymbol{N}}_{j}^{(-)} \otimes \boldsymbol{I}_{\widehat{\alpha}(j+1,K)} \otimes \mu_{j} \widehat{\boldsymbol{L}}_{j} \otimes \boldsymbol{I}_{W} \right) + \boldsymbol{I}_{\widehat{\alpha}(2,K)} \otimes \boldsymbol{I}_{S-s} \otimes \boldsymbol{D}_{0} - \widetilde{\boldsymbol{\Delta}}. \end{split}$$

Note that the diagonal matrix $\hat{\Delta}$ holds the exit rates from states which do not lead to absorption, that is, transitions from inventory levels $\hat{I}(t) = 2, ..., S - s$ according to enumeration (43). As such, this matrix can be defined explicitly as follows:

$$\tilde{\boldsymbol{\Delta}} = \mu_1 \mathbf{e}^{\hat{\alpha}(2,K)} \otimes \mathbf{e}^{S-s}_{2:S-s} \otimes \mathbf{e}^W + \sum_{j=2}^K \mu_j \mathbf{e}^{\hat{\alpha}(2,j-1)} \otimes \mathbf{e}^{S-s_{j-1}+1}_{2:S-s_{j-1}+1} \otimes \mathbf{e}^{\hat{\alpha}(j+1,K)} \otimes \mathbf{e}^{S-s}_{s_{j-1}-s+2:S-s} \otimes \mathbf{e}^W.$$

Similarly to (19), the matrix $B_1^{(0)}$ differs from B_1 only on the diagonal (the states corresponding to $B_1^{(0)}$ have no class 1 customers), so that

$$\boldsymbol{B}_{1}^{(0)} = \boldsymbol{B}_{1} - \mu_{1} \mathbf{e}^{\widehat{\alpha}(2,K)} \otimes \mathbf{e}_{2:S-s}^{S-s} \otimes \mathbf{e}^{W}.$$

Finally, we need to define the matrix \hat{B}_1 corresponding to the boundary states. Since at the *level* (number of class 1 customers) S - s, neither arrivals nor departures of the class 1 customers are tracked,

$$\mathbf{B}_1 = \mathbf{B}_0 + \mathbf{B}_1 + \mathbf{B}_2$$

It remains to denote $t_0 = -Te$ as the corresponding absorption rate vector, and note that the initial state probability vector should be taken so as to have initial inventory equal to *S*, that is, P(I(0) = S) = 1.

To simplify comprehension, we outline the transitions possible for the chain

$$\{(\hat{N}_1(t),\ldots,\hat{N}_K(t),\hat{I}(t),Z(t))\},t\geq 0,\$$

according to the subgenerator matrix T and the absorption vector t_0 . For some $t \ge 0$, fix $(\hat{N}_1(t), \ldots, \hat{N}_K(t), \hat{I}(t), Z(t)) = (n_1, \ldots, n_K, i, z)$. Then, the state after transition is one of the following:

- $(n_1, \ldots, n_j + 1_{n_j < S s_{j-1}}, \ldots, n_K, i, z')$, with rate $(D_j)_{z,z'}$ (arrival of a class *j* customer),
- $(n_1, \ldots, n_j, \ldots, n_K, i, z')$, with rate $(D_0)_{z,z'}$ (change of the MMAP phase),
- $(n_1, \ldots, n_j 1_{n_j < S s_{j-1}}, \ldots, n_K, i 1, z')$, with rate μ_j , if $n_j > 0$ and $i > s_{j-1} s$ (departure of a class *j* customer, if allowed).

Note that the absorption happens if the transition is made from the inventory level i = 1 downward. In particular, this means that, since $s_2 > s$, the absorption happens only due to a transition caused by service completion of either class 1, or class 2 customer.

It is well known that time to absorption, say *X*, of an absorbing CTMC defined by a subgenerator *T* and initial state probability vector β has a phase-type distribution with mean [24]

$$\mathbf{E}\mathbf{X} = \boldsymbol{\beta}(-\boldsymbol{T})^{-1}\mathbf{e}.$$
(47)

Thus, to define the mean inventory cycle time, $E\Gamma$, it only remains to convert the steady-state probability vector q obtained in (32) into the initial state probability vector for the corresponding phase-type distribution defined by the matrix (42).

We summarize our findings in the following Lemma.

Lemma 1. Let Γ be the inventory cycle time. Then,

$$\mathrm{E}\Gamma = \frac{1}{\gamma} + \sum_{n=0}^{\infty} \widehat{q}_n (-T)^{-1} \mathbf{e}, \tag{48}$$

where the vector \hat{q} , split into finite subvectors by the first component, is defined from the vector q in a component-wise manner as follows: if $n_k < S - s_{k-1}, k = 1, ..., K$, then

$$\widehat{q}(n_1,\ldots,n_K,S,z) = q(n_1,\ldots,n_K,S,z),$$

and if for some $m \ge 1$ and indices k_1, \ldots, k_m , holds $n_{k_j} = S - s_{k_j-1}, j = 1, \ldots, m$, then

$$\widehat{\boldsymbol{q}}(n_1,\ldots,n_K,S,z) = \sum_{\widehat{n}_1,\ldots,\widehat{n}_K \in \mathcal{N}} \boldsymbol{q}(\widehat{n}_1,\ldots,\widehat{n}_K,S,z),$$

where

$$\mathcal{N} = \Big\{ (\hat{n}_1, \dots, \hat{n}_K) : \hat{n}_k = n_k, k \in \{1, \dots, m\} \setminus \{k_1, \dots, k_m\}; \hat{n}_{k_j} \ge S - s_{k_j-1}, j = 1, \dots, m \Big\}.$$

The vector \hat{q} is zero elsewhere.

It remains to note that the expectation is taken so as to summarize the resulting times by the appropriate steady-state probabilities involving the inventory level *S*, and the term $\frac{1}{\gamma}$ is added to emphasize the exponentially distributed lead time.

6. Numerical Illustration

6.1. Stability Condition Parametric Sensitivity

In this section we illustrate the sensitivity of the stability condition (30) on the parameters. In particular, we take K = 2, fix the MMAP arrival process, inventory size *S* and queue size for the second class customers, l_2 . We then vary the rates μ_1 , μ_2 and γ for several values of replenishment level *s*, *ceteris paribus*, and plot the resulting dependency. In the experiment we use the following default settings:

$$S = 100, s = 7, l_2 = 15, \mu_1 = 7, \mu_2 = 15, \gamma = 10.$$

The MMAP arrival process with W = 2 is driven by the following matrices:

$$\boldsymbol{D}_0 = \begin{bmatrix} -13 & 1\\ 2 & -14 \end{bmatrix}, \ \boldsymbol{D}_1 = \begin{bmatrix} 1 & 2\\ 3 & 4 \end{bmatrix}, \ \boldsymbol{D}_2 = \begin{bmatrix} 5 & 4\\ 3 & 2 \end{bmatrix}.$$
(49)

Thus, the arrival rate of class 1 customers may be calculated as $\lambda_1 \approx 4.867$.

In the first experiment we consider ρ defined in (30) to be the function of $\mu_2 \in [2, 15]$ for values s = 5, 6, 7. The resulting curves depict $\rho_s(\mu_2)$ versus μ_2 on Figure 2 (top). While in absolute values the variability of ρ is rather small, a non-linear dependency is clearly visible, and the dependency on s is motivated by the relatively higher impact of the second class customers on the system load for smaller values of s. Counter-intuitively though, the load *increases* (slightly) with an increasing service rate of the second class customers.

To investigate this interesting effect, we consider the probability $P(N_2 = i)$ for two boundary values, $i = 0, l_2$, that is, the probability that second class customer queue is empty or full. We plot the corresponding estimates within the setup of the first experiment on Figure 3. It can be observed that, for smaller values of μ_2 (that is, for larger mean service times of class 2 customers), the queue is mostly overloaded (with high probability the queue is full), which causes relatively high loss of class 2 customers. In contrast, for $\mu_2 > 7$, the probability of an empty class 2 queue overtakes the probability of a queue completely occupied, and the former is increasing, while the latter is decreasing with increasing μ_2 .



We note that, in the experiment, we used only s = 7 since, as numerical results show, the effect of *s* on empty/full probability of class 2 customers is negligible for s = 5, 6, 7.

Figure 2. Dependency of ρ on $\mu_2 \in [2, 15]$ (top, fixed $\gamma = 10$) and on on $\gamma \in [2, 10]$ (bottom, fixed $\mu_2 = 15$) for s = 5, 6, 7 and other parameters fixed: $S = 100, l_2 = 15, \mu_1 = 7$, MMAP arrival process parameters given in (49).



Figure 3. Dependency of the probabilities of an empty (full) class 2 queue on $\mu_2 \in [2, 15]$, fixed $\gamma = 10, s = 7, S = 100, l_2 = 15, \mu_1 = 7$, MMAP arrival process parameters given in (49).

In the second experiment, we fix $\mu_2 = 15$ and consider ρ defined in (30) to be the function of $\gamma \in [2, 10]$ for values s = 5, 6, 7. The resulting curves depict $\rho_s(\gamma)$ versus γ on Figure 2 (bottom). In this experiment, the non-linear dependency on γ follows the intuition: higher replenishment rate decreases the customer's queuing time and results in a lower system load.

Finally, we fix $\mu_2 = 15$, $\gamma = 10$ and consider ρ defined in (30) to be the function of $\mu_1 \in [7, 12]$. However, since the dependency on *s* is rather weak, we take more contrast values s = 1 and s = 50. The resulting curves depict $\rho_s(\mu_1)$ versus μ_1 on Figure 4. As expected, increasing μ_1 , and, hence, decreasing the service time of class 1 customers, causes a dramatic decrease of the system load. The additional load caused by a "lazy" replenishment at the level s = 1 is also visible.

6.2. Steady-State Performance Sensitivity

In this section, we illustrate the sensitivity of the performance measures (37)–(41) described in Section 4 to the management parameter γ , that is, the lead time intensity. To do so, we take K = 2 and slightly modify the parameters used in the previous section, so as to make computations more convenient. Namely, we take

$$S = 50, s = 5, l_2 = 10, \mu_1 = 7, \mu_2 = 15.$$

We use the same MMAP defined in (49). We vary $\gamma = 1, ..., 10$ and obtain the performance measures (37)–(41) for given γ , *ceteris paribus*. Finally, we depict the obtained numerical results.

Figure 5 (top) describes the dependency of the mean number of class 1 and class 2 customers given in (37) and (38), on the replenishment rate γ . A decreasing pattern with increasing γ is caused by decreasing load ρ given in (30); see Figure 5 (bottom).



Figure 4. Dependency of ρ on $\mu_1 \in [7, 12]$ for s = 1,50 and other parameters fixed: S = 100, $l_2 = 15$, $\mu_1 = 7$, $\mu_2 = 15$, $\gamma = 10$, MMAP arrival process parameters given in (49).



Figure 5. Cont.



Figure 5. Dependency of \mathbf{E}_i , i = 1, 2 (top) and ρ (bottom) on $\gamma = 1, ..., 10$ for S = 50, s = 5, $l_2 = 10$, $\mu_1 = 7$, $\mu_2 = 15$, MMAP arrival process parameters given in (49).

Figure 6 (top) describes the dependency of the mean inventory level $\mathbf{\overline{I}}$ defined in (39), while, in Figure 6 (bottom), the replenishment rate $\mathbf{\overline{R}}$ is depicted, for $\gamma = 1, ..., 10$. As expected, with increasing γ the inventory contents increases, along with the replenishment rate.



Figure 6. Cont.



Figure 6. Dependency of $\overline{\mathbf{I}}$ (top) and $\overline{\mathbf{R}}$ (bottom) on $\gamma = 1, ..., 10$ for S = 50, s = 5, $l_2 = 10$, $\mu_1 = 7$, $\mu_2 = 15$, MMAP arrival process parameters given in (49).

Finally, Figure 7 describes the dependency of the class 2 customer loss rate L_2 defined in (41), on $\gamma = 1, ..., 10$. As expected, with increasing γ the customer loss rate is decreasing.



Figure 7. Dependency of L₂ on $\gamma = 1, ..., 10$ for S = 50, s = 5, $l_2 = 10$, $\mu_1 = 7$, $\mu_2 = 15$, MMAP arrival process parameters given in (49).

6.3. Total Cost Optimization

Based on the above performance measures, we obtain expected total cost per unit time in the considered system as follows:

$$E_{TC} = c_I \bar{\mathbf{I}} + \sum_{k=2}^{K} c_{L_k} \mathbf{L}_k + c_R \bar{\mathbf{R}} + \sum_{k=1}^{K} c_{w_k} \mathbf{E}_{\mathbf{k}},$$
(50)

where the non-negative coefficients are defined on a per system | customer | item per time unit basis:

- *c*_{*I*} is the holding cost of inventory (per item per unit time),
- *c*_{*L*_{*k*}} is the cost due to a class *k* customer loss (per system per unit time),
- *c_R* is the reorder cost (per system per unit time),
- c_{w_k} is the waiting cost of class *k* customer (per customer per unit time), $k \le K$. Using the same 2-class system with same parameters as in Section 6.2, that is,

$$S = 50, s = 5, l_2 = 10, \mu_1 = 7, \mu_2 = 15,$$

and MMAP defined in (49), we perform a numerical exploration of sensitivity of the cost function E_{TC} given in (50) on the parameter c_I . This is motivated as follows. Since the per-customer measures are decreasing, while the per-system measures are increasing, there should be a trade-off between keeping a high inventory level and compromising some second-class customer losses. This is regulated by the inventory holding cost parameter c_I . For the experiment, we take $c_I = 1, 5, 10$ and vary $\gamma = 1, ..., 10$, as in the previous section. We depict the resulting curves on Figure 8 and observe clearly that the optimal (minimal) cost shifts to lower values of γ with increasing c_I , which is intuitive.



Figure 8. Total cost E_{TC} versus $\gamma = 1, ..., 10$ for various $c_I = 1, 5, 10$; the parameters are S = 50, $s = 5, l_2 = 10, \mu_1 = 7, \mu_2 = 15$, MMAP arrival process parameters given in (49).

7. Conclusions

We have analyzed a queuing inventory system with correlated arrivals of heterogeneous customers. The general case of *K* types of customers was investigated, where each class of customers is assigned some priority according to which the inventory access is managed. We illustrated the general system and a special case of K = 2. An intuitive stability condition was derived and the steady state probability vector was obtained. Key performance measures of the system were obtained. The inventory cycle time was analyzed and it was shown that the inventory cycle follows a phase type distribution.

The large scale models have computational difficulties due to a large state space (while the matrices are sparse), so we cannot rely on numerical investigation only. Thus, we performed numerical experiments only to highlight interesting features of the model, and an interesting cost optimization problem has been stated which could be studied in future for practical applications.

As future directions of investigation we would like to point out the following possibilities. First, it is relatively easy to incorporate customer impatience (which will cause additional flow of customers towards decreasing the *level* and the first *K* components of the *phase* in the corresponding QBD process) and balking upon arrival finding the server busy (which modifies the rates related to corresponding class customer arrivals). Phasetype service times and replenishment times can also be incorporated leaving the model mathematically tractable. At the same time, the inventory replenishment discipline can be modified into booking-type, where the inventory items are booked for specific customers upon arrival; thus, arriving customers are rejected if no items are available for booking. Comparison of this type of model with the one analyzed in the present paper is one of the promising directions for future research.

Thinking beyond simple extensions of the model, the retrial queues instead of classical queues can be incorporated, while the stability part most likely can be extended towards a regenerative input. Finally, it might be interesting to consider the inventory with common lifetime of items and priorities related to item "freshness".

It might be also interesting to consider the model in transient regime, according to possible applications in social systems. A few methods are available for this type of analysis of QBD processes [25–27]. However, this might require inverting Laplace transforms which might cause numerical instability. In this regard, we refer to a recent work, Reference [28], where a novel method of numerically effective Laplace transform inversion is presented.

Author Contributions: Conceptualization, A.K., K.R., and M.J.J.; formal analysis, A.K., K.R., M.J.J., and A.S.R.; Writing—original draft preparation, K.R.; simulation and visualization, K.R. and A.S.R.; writing—review and editing, A.K., A.S.R., and M.J.J.; supervision, project administration, A.K. All authors have read and agreed to the published version of the manuscript.

Funding: This work is supported by Ministry of Science and Higher Education of the Russian Federation Grant 075-15-2019-1621. The work of AR is partially supported by RFBR, projects 19-57-45022, 19-07-00303.

Acknowledgments: The authors would like to thank anonymous referees for review and detailed comments that helped to improve the paper.

Conflicts of Interest: The authors declare no conflict of interest.

References

- Chakravarthy, S.R.; Maity, A.; Gupta, U.C. An '(s, S)' inventory in a queueing system with batch service facility. *Ann. Oper. Res.* 2017, 258, 263–283. [CrossRef]
- Chakravarthy, S.R.; Rumyantsev, A. Analytical and simulation studies of queueing-inventory models with MAP demands in batches and positive phase type services. *Simul. Models Pract. Theory* 2020, 103, 102092. [CrossRef]
- 3. Altendorfer, K. Influence of lot size and planned lead time on service level and inventory for a single-stage production system with advance demand information and random required lead times. *Int. J. Prod. Econ.* **2015**, *170*, 478–488. [CrossRef]
- 4. Armony, M.; Plambeck, E.L. The impact of duplicate orders on demand estimation and capacity investment. *Manag. Sci.* 2005, 51, 1505–1518. [CrossRef]
- 5. Sohraby, K.; Sidi, M. On the performance of bursty and modulated sources subject to leaky bucket rate-based access control schemes. *IEEE Trans. Commun.* **1994**, *42*, 477–487. [CrossRef]
- Choi, D.I.; Lim, D.E. Analysis of the State-Dependent Queueing Model and Its Application to Battery Swapping and Charging Stations. Sustainability 2020, 12, 2343. [CrossRef]

- Sigman, K.; Simchi-Levi, D. Light traffic heuristic for anM/G/1 queue with limited inventory. *Ann. Oper. Res.* 1992, 40, 371–380.
 [CrossRef]
- 8. Berman, O.; Kaplan, E.H.; Shevishak, D.G. Deterministic approximations for inventory management at service facilities. *IIE Trans.* **1993**, *25*, 98–104. [CrossRef]
- 9. Krishnamoorthy, A.; Lakshmy, B.; Manikandan, R. A survey on inventory models with positive service time. *Opsearch* 2011, 48, 153–169. [CrossRef]
- 10. Enders, P.; Adan, I.; Scheller-Wolf, A.; van Houtum, G.J. Inventory rationing for a system with heterogeneous customer classes. *Flex. Serv. Manuf. J.* **2014**, *26*, 344–386. [CrossRef]
- 11. Avrachenkov, K.; Dudin, A.; Klimenok, V. Retrial queueing model MMAP/M 2/1 with two orbits. In Proceedings of the International Workshop on Multiple Access Communications, Barcelona, Spain, 13–14 September 2010; pp. 107–118.
- 12. He, Q.M.; Stanford, D.A. Distributions of the interdeparture times in FCFS and nonpreemptive priority MMAP [2]/G [2]/1 queues. *Perform. Eval.* **1999**, *38*, 85–103. [CrossRef]
- 13. Klimenok, V.; Dudin, A.; Vishnevsky, V. Priority Multi-Server Queueing System with Heterogeneous Customers. *Mathematics* **2020**, *8*, 1501. [CrossRef]
- 14. Dudin, S.; Kim, C.; Dudina, O. MMAP | M | N queueing system with impatient heterogeneous customers as a model of a contact center. *Comput. Oper. Res.* 2013, 40, 1790–1803. [CrossRef]
- 15. He, Q.M. Quasi-birth-and-death Markov processes with a tree structure and the MMAP [K]/PH [K]/N/LCFS non-preemptive queue. *Eur. J. Oper. Res.* 2000, 120, 641–656. [CrossRef]
- 16. He, Q.M. Queues with marked customers. Adv. Appl. Probab. 1996, 28, 567–587. [CrossRef]
- 17. Neuts, M.F. A versatile Markovian point process. J. Appl. Probab. 1979, 16, 764-779. [CrossRef]
- 18. He, Q.M.; Neuts, M.F. Markov chains with marked transitions. Stoch. Process. Their Appl. 1998, 74, 37-52. [CrossRef]
- 19. He, Q.M. The Versatility of MMAP[K] and the MMAP[K]/G[K]/1 Queue. Queueing Syst. 2001, 38, 397–418. [CrossRef]
- 20. Chakravarthy, S.R. Markovian Arrival Processes. In *Wiley Encyclopedia of Operations Research and Management Science;* American Cancer Society: Atlanta, GA, USA, 2011; [CrossRef]
- 21. Geerts, F.; Blondia, C. Superposition of Markov Sources and Long Range Dependence. In *Broadband Communications*; Springer: Boston, MA, USA, 1998; pp. 550–561. [CrossRef]
- 22. He, Q.M. Fundamentals of Matrix-Analytic Methods; Springer: New York, NY, USA, 2014.
- 23. Latouche, G.; Ramaswami, V. Introduction to Matrix Analytic Methods in Stochastic Modeling; ASA-SIAM: Philadelphia, PA, USA, 1999.
- 24. Bladt, M.; Nielsen, B.F. *Matrix-Exponential Distributions in Applied Probability*; Probability Theory and Stochastic Modelling; Springer: Boston, MA, USA, 2017; Volume 81. [CrossRef]
- 25. Lucantoni, D. Further transient analysis of the BMAP/G/1 Queue. Commun. Stat. Stoch. Models 1998, 14, 461–478. [CrossRef]
- 26. Kulkarni, L.A.; Li, S.Q. Transient behaviour of queueing systems with correlated traffic. *Commun. Stat. Stoch. Models* **1998**, 14, 933–978. [CrossRef]
- 27. Zhang, J.; Coyle, E.J. Transient analysis of quasi-birth-death processes. Commun. Stat. Stoch. Models 1989, 5, 459–496. [CrossRef]
- 28. Almousa, S.A.D.; Horváth, G.; Telek, M. Transient analysis of piecewise homogeneous QBD process. *Stoch. Models* **2021**, *37*, 59–84. [CrossRef]