# Optimal Design and Control of Finite-Population Queueing Systems

Chao Deng

A dissertation submitted to the faculty of the University of North Carolina at Chapel Hill in partial fulfillment of the requirements for the degree of Doctor of Philosophy in the Department of Statistics and Operations Research.

Chapel Hill

2012

Approved by:

Nilay Tanık Argon

Vidyadhar G. Kulkarni

Shankar Bhamidi

J. Scott Provan

Serhan Ziya

# ABSTRACT

CHAO DENG: Optimal Design and Control of Finite-Population Queueing Systems
(Under the direction of Professor Nilay Tanık Argon and Professor Vidyadhar G. Kulkarni)

We consider a service system with a finite population of customers (or jobs) and a service resource with finite capacity. We model this finite-population queueing system by a closed queueing network with two stages. The first stage, which represents the arrivals of customers for service, consists of an automated station with ample capacity. The second stage, which represents the service for customers, consists of multiple service stations which share the finite service resource. We consider both discrete and continuous service resources. We are interested in static or dynamic allocation of the service resource to the service stations in the second stage in order to optimize a given system measure. Specifically, a static allocation refers to a design problem, while a dynamic allocation refers to a control problem. In this thesis, we study both.

For control problems, we specify a parallel-series structure for service stations. We first consider dynamically allocating a single flexible server under both preemptive and non-preemptive policies. We characterize the optimal policies of dynamically scheduling this single server in order to maximize the long-run average throughput of the system. In the special case of a series system, we show that the optimal policy is a sequential policy where each customer is served by the single server sequentially from the first station until the last one. For a parallel system, we show that there exists an optimal policy which gives the highest priority to the station that has the largest service rate. We also propose an index policy heuristic for the general parallel-series system and compare its performance as opposed to the optimal policy by a numerical study. Finally, we study dynamically allocating a finite amount of continuous service resource for the parallel system.

For design problems, we consider allocating a finite amount of service resource which is continuously divisible and can be used at any of the service stations. Suppose that

service times at a service station are exponentially distributed and their mean is a strictly increasing and concave function of the allocated service resource. We characterize the optimal allocation of the continuous resource in order to maximize the long-run average throughput of the system. We first show that the system throughput is non-decreasing in the number of customers. Then, we study the optimization problem in three cases depending on the population size of customers in the system. First, when there is a single customer, we show that the optimal allocation is given by a set of optimality equations. Secondly, when the number of customers approaches infinity, we show that the optimal allocation approaches to a limit. Finally, for any finite number of customers, we show that the system throughput is bounded up by a limit. Moreover, under a certain condition, we show that the system throughput function is Schur-concave.

**Keywords**: finite-population queueing systems, dynamic control, static design.

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

# LIST OF FIGURES

# Chapter 1

# INTRODUCTION

A finite-population queueing system, or sometimes called a finite-source queueing system, is a queueing system in which requests for service are generated by a finite number of customers (sources) and the requests are handled by a single or multiple server(s). The service times of the requests generated by the customers are random variables. It is assumed that the server can handle only one request at a time and uses a specified service discipline. A customer can be idle, waiting for service, or in service. An idle customer generates a request for service after a random amount of time independent of the states of the other customers. Once the service is completed, the customer becomes idle again, and the process repeats. In this thesis, we consider identical customers whose idle periods follow an independent and identically distributed (i.i.d.) exponential distribution.

Two widely known applications of finite-population queueing systems are machine interference problems (or alternatively machine repairman problems) and computer-communication systems. In the machine interference problems, the customers are machines. A machine operates for a random period of time, then breaks down and requests service from workers. While the worker resource is scarce, the service facility needs to make two decisions: which worker will serve each machine and in what order the machines will be served. There is a rich literature studying the machine interference problems. For a complete review of research in this area, see the survey paper by Haque and Armstrong [7].

Many computer and communication systems can be modeled as finite-population queueing systems. For example, in a computer-communication system, $N$ terminals request to use a computer (server) to process transactions. Each of the terminals takes a random time to generate a request for the computer. The computer works on each of the transactions and responds to the user at the terminal once the transaction is completed. The throughput rate at which transactions are processed, or equivalently, generated in steady-state is one of the most important performance measures showing the system's processing power.

See Sztrik [30] for a complete review of finite-population queueing systems applications and bibliography of related papers.

Finite-population queueing models can be also useful in developing effective policies in healthcare. Green [5] provides an overview of using queueing analysis to improve service in healthcare. For example, the finite-population queueing systems can be applied to the nurse staffing problems in hospital wards. For the hospitals with high demands and constrained resources, it is reasonable to assume that the number of patients staying is constant. Patients stay in a bed for a random period of time, then request service from nurses. For a given objective, for example, to maximize the steady-state processing rate of service requests, the hospital managers need to determine service policies which optimally assign the nurses to each of the requests.



Figure 1.1: A closed queueing system with two stages of stations.

In this thesis, we are interested in optimal allocation of the resource capacity to stations in a finite-population queueing system. We consider two management paradigms: dynamic and static. Under the dynamic paradigm, we are allowed to change resource allocation whenever the system changes its state. For the dynamic problems, we consider only non-idling policies, i.e., a server cannot be idle if there is any service request waiting. Under the static paradigm, allocation decisions are made before the system starts to operate. When the system is under operation, we are not allowed to adjust resource allocation any more.

We also consider two types of service resource: continuous and discrete. For the first

2

case, we assume that a finite amount of continuous service capacity can be used at any of the service stations in the second stage. For the discrete case, we assume that there is a single flexible server who is able to work at any of the service stations.

We model the finite-population service system by a closed queueing network with $K+1$ $(2 \leq K < \infty)$ stations (labeled as station $0, 1, 2, \ldots, K$) and a fixed population size of $B$ $(1 \leq B < \infty)$. These $K+1$ stations form the two stages of service, as shown in Figure 1.1, where $B$ jobs circulate between the two stages. The first stage consists of the automated station, which is labeled as station $0$, where $B$ servers are dedicated to this station. A customer receives service by one of the servers at station $0$ immediately after she arrives at this automated station. Assume that service times at station $0$ are independent and identically distributed with an exponential distribution and mean $1/\mu_0$. The second stage consists of the remaining $K$ service stations, i.e., stations $1, 2, \ldots, K$. We refer to these $K$ stations in the second stage as service stations and to station $0$ as an automated station. Suppose that a fixed amount of resource is available to be used by service stations. The decision is how to allocate the resource (statically or dynamically) to each of the service stations in order to optimize a given system measure.

In this thesis, we mainly consider the objective of maximizing the long-run average reward (throughput) of the system. We also conduct a brief study on dynamic control of finite-population queueing systems under the objective of waiting cost minimization.

The organization of this thesis is given as follows. In Chapter 2, we provide a literature review on optimal control and design for finite-population queueing system and closed queueing systems. In Chapter 3 and 4, we study the dynamic control problems and the static design problems, respectively, under the reward (throughput) maximization objective. In Chapter 5, we present our work on dynamic control problems with the objective of waiting cost minimization. Finally, we provide future research directions in Chapter 6.

# Chapter 2

# LITERATURE REVIEW

Optimal design and control of queueing networks have been an important research subject of Operations Research for decades. In this chapter, we review the literature that is most relevant to our work. We focus on the papers that study design and control problems for closed queueing networks where finite-population service systems fit best. We also briefly review the literature studying open queueing networks whose methodologies on design and control problems are relevant to our work. In the review of each paper, we indicate the following main features of the model considered by that paper:

1. Resource type (continuous, discrete single-server, or discrete multiple-server);

2. Objective function (e.g., throughput maximization or waiting cost minimization);

3. Policy constraints (e.g., preemption or non-preemption);

4. Structure of the queueing network (e.g., parallel or series);

5. The Results.

## 2.1 Dynamic Control

### 2.1.1 Finite-Population and Closed Queueing Systems

There is a rich literature that studies finite-population queueing systems and one of its major applications, namely, the machine interference problem. A research survey on the machine interference problem up to 1985 is provided by Stecke and Aronson [27]. Haque and Armstrong [7] extensively extend this survey on this area until 2007. For a complete literature review on finite-population queueing systems up to 2001, see Sztrik [30]. We here review only the work that is most relevant to this study.

Palesano and Chandra [21] study a machine interference problem with multiple types of failures. A single server is available to work on all types of failures. A machine stays functioning for a random time, then breaks down and requests service from the repairman. The type of failures is random. This paper studies the system performance under different service priorities by a numerical study and compares them. They do not prove any optimality results. They find that the mean number of machines waiting for repair increases when the failure types with higher mean service times are given a high priority. This observation is consistent with our work, in which we prove that the optimal policy that maximizes the long-run system throughput rate gives priority to the stations with higher service rates.

Both Iravani and Kolfal [10] and Iravani, Krishnamurthy, and Chao [12] study a single-server machine repairman problem with multiple classes of machines. The single server is available to serve all machine classes. Iravani and Kolfal [10] consider preemptive policies. Cost is incurred when a machine is waiting for service. The authors observe that in a finite-population queueing system, ignoring customers' arrival rate and applying the $c\mu$ rule is not always optimal to minimize the long-run average cost of the system. They find the conditions under which static-priority rules, e.g. the $c\mu$ rule, are optimal independent of customer arrival rate and customer population size. Iravani, Krishnamurthy, and Chao [12] consider the non-preemptive case. Cost is incurred when a machine is down. The authors investigate the dynamic assignment policies that minimize the total average customer waiting cost. The authors show that the optimal service policy may never serve some classes of machines. For those classes that are served, the paper shows that a static priority policy is optimal, and derives sufficient conditions that determines the optimal priority sequence.

Iravani and Krishnamurthy [11] study a machine repairman problem with partially cross-trained servers, i.e., each server is able to repair a set of machines. Cost is incurred when a machine is down and waiting for repair. The objective is to obtain the optimal policy that dynamically allocates servers to minimize the long-run average cost. The paper shows that static machine priority rules are effective in minimizing the waiting cost rate.

Many papers study production systems by modeling them as closed serial queueing networks. Koole and Righter [15] study a tandem manufacturing system with multiple flexible servers. The tandem stations can be divided into several non-overlapping sets of adjacent stations. Each server is able to work on one set of stations. The paper looks

for optimal policies of dynamically assigning the server to work within his set of stations. The objective is to maximize the departure process stochastically. They consider two cases where either preemption and idling are permitted or preemption and idling are not allowed. For both cases, they show that the optimal policy assigns each server to work at his last nonempty station.

Hopp, Iravani, Shou and Lien [9] study a manufacturing system with a mix of manual and automated equipment. The system operates under a constant work-in-process (CONWIP) protocol, and is staffed by a single cross-trained worker. The system is modeled by a tandem queueing network with three stations. The first station is an automated station with automatic processing times but requiring a manual loading time. The other two stations are manual stations requiring manual processing times. The single flexible server is able to work at all stations. The paper investigates the optimal control policy to maximize the average throughput rate. They show that the optimal control policy is a static priority policy.

All papers reviewed in this section study dynamic control problems with a single server or multiple servers. In other words, they consider discrete service resources in their models. In this study, we consider both continuous and discrete service resources for our problems. For the continuous resources, we assume that there is a fixed job processing capacity that can be divided continuously among the $K$ service stations in the second stage. For the discrete resources, we assume that a single server or multiple servers are available to be allocated among the $K$ service stations.

### 2.1.2 Open Queueing Systems

There is a rich literature that studies design and control problems for open queueing networks. In this thesis, we only review those papers that are most relevant to our work. It is important to point out that our focus is not open queueing networks but finite-population and closed queueing networks.

Klimov [13] is a pioneering study on service priorities of open queueing networks. Klimov [13] studies a dynamic control problem for an open queueing network with a finite number of nodes and a single server. Jobs arrive according to a Poisson process at all nodes, and service times are generally distributed at each node. Cost is incurred when a job is waiting

in a queue. Interruption of service is not allowed. The paper proves that a priority index policy is optimal in order to minimize the long-run average waiting cost. In a subsequent paper, Klimov [14] provides a simple and efficient algorithm to compute such a priority index for queueing networks with a forest structure.

Harrison [8] is another pioneering work on studying service priorities of open queueing networks. Harrison [8] considers a single-server queue with multiple customer classes. He assumes independent Poisson arrival processes. Service times have general distribution which depend only on customer classes. Cost is incurred when a customer is waiting in the system. A reward is gained when a customer is served. The objective is to maximize the discounted total profit over an infinite planning horizon. A priority index policy called modified static policy is shown to be optimal.

Van Oyen et al. [33] study a serial production system with flexible workers. Service times are generally distributed and depend only on stations. The paper considers both collaborative (servers are able to work together on one job) and non-collaborative (servers are not allowed to work together on one job) cases. Under the collaborative case, they show that the so-called expedite policy is optimal to minimize the cycle time for each job. The expedite policy places all the servers successively on a given job. Under the non-collaborative case, no optimal policy is found. However, they propose a so-called pick-and-run policy and demonstrate that it is near-optimal. The paper also extends some their insights to a capacity-constrained environment with a constant work-in-process protocol.

Andradottir et al. [2] consider dynamic control problems for multiple-server queueing systems. Their objective is to find optimal dynamic allocation policies to maximize the long-run average throughput. Travel times of servers between stations are negligible. They show that all non-idling policies are optimal when service rates depend only on either servers or stations. For a special two-station tandem queueing system with two flexible servers and finite number of buffers between the two stations, the paper shows that the optimal policy assigns one server to each station unless the first station is starved or the second station is blocked. Andradottir and Ayhan [1] later extend this result to the case with three servers.

Yankovic and Green [35] explore the appropriate nurse-to-patient levels to minimize the probability that a patient's service request is delayed. They use a two-dimensional open queueing system rather than finite-population queueing system to model the hospital

system. However, this nurse staffing problem can also be another application for finite-population queueing systems. We consider a finite-population queueing system with four service stations representing patients' admission, patients' stay in beds, patients' caring requests and patients' discharge. While this paper seeks an appropriate nurse-to-patient level for hospitals, we are interested in obtaining priority policies that optimally assign the nurses to each of the patient requests in order to achieve a given objective. Finite-population queueing systems can be used to model the hospital systems if we assume that the number of patients staying in the hospitals is constant.

## 2.2  Static Design

In this section, we review articles that study static workload or server allocation problems of closed queueing systems. Stecke and Morin [28] investigate optimality of balancing workloads in closed queueing systems. They consider a central server closed queueing network where stations are parallel to each other. They are interested in obtaining optimal policies for allocating workload in order to maximize the system throughput. The paper proves that the throughput of this system is a quasi-concave function of workloads, and shows that a balanced allocation of workloads maximizes the expected throughput of the system.

Stecke [26] studies the non-concavity property of throughput function in closed queueing networks. For a general-structured closed queueing network with multiple customers, she shows that the throughput function is not concave in workload. When the closed queueing network includes two single server stations, the paper proves that the throughput function is concave when there are two customers, and the throughput function is quasi-concave when there are more than two customers.

Yao [36] considers a closed queueing network with single-server stations and exponential service times. He investigates the concavity property of the long-run average throughput of the system. He proves that the system throughput as a function of loading is Schur-concave. As a consequence, the paper shows that, when the total loading of the system is a constant, the balanced (or equal) loading maximizes the system throughput based on the majorization property.

Shanthikumar and Yao [25] study the static allocation problem of a multiple-server closed queueing network. Their objective is to find optimal policies for allocating servers

to maximize throughput. The paper shows that the throughput of the system has a monotonicity property, which means that an optimal policy allocates more servers to a station with a higher workload. They provide a search algorithm to obtain an optimal policy within a small number of allocations satisfying the monotonicity property.

Lee, Srinivasan, and Yano [17] consider the problem of allocating workload in a closed queueing network with multi-server stations. Their objective is to maximize the long-run average throughput of the system. The paper assumes that the system throughput is product-form and there is a single class of customers. The paper proves that the throughput function is quasi-concave in workload for a single-server closed queueing network and a multi-server closed queueing network with two customers. For the general model, the authors develop two heuristic algorithms to search the optimal workload allocation.

These papers study static allocation problems of closed queueing networks where all stations are included for allocation decisions. In this study, we consider a two-stage queueing system where the first stage is a special automated station. Allocation decisions are made only for the service stations in the second stage. Furthermore, except for Shanthikumar and Yao's paper [25], the other papers consider allocating the workload in the system. In this thesis, we consider allocating a fixed amount of service resource, and we study both discrete and continuous resources.

# Chapter 3

# OPTIMAL DYNAMIC CONTROL OF FINITE-POPULATION QUEUEING SYSTEMS: REWARD MAXIMIZATION

In this chapter, we study dynamic control problems, i.e., we are allowed to change resource allocation when the system changes its state. We consider both preemptive and non-preemptive policies. Let $\Pi_P$ and $\Pi_{NP}$ denote the set of preemptive policies and non-preemptive policies, respectively. For preemptive policies, the server is allowed to make service decisions whenever a service at the automated station or the service stations is completed. Under non-preemptive policies, the server is allowed to switch to work at other stations only when it completes service. In this chapter, we only consider non-idling policies, i.e., the server (or service capacity) is not allowed to be idle whenever job(s) are available at the service stations.

We formulate the dynamic control problem in Section 3.1. In Sections 3.2 and 3.3, we consider the case where the service resource is discrete and a single server is available to work at $K$ service stations. We study preemptive and non-preemptive policies in Section 3.2 and 3.3, respectively. In Section 3.4, we consider the problem with continuous resource constraint.

## 3.1 Model Formulation

We consider a special case of the general finite-population queueing system, called a *parallel-series* system, as shown in Figure 3.1. The second stage in this system consists of $M$ parallel branches, where the $m$-th branch consists of $i_m$ service stations in series. We use an $M$-dimensional vector $(i_1, i_2, \ldots, i_M)$ to denote the parallel-series system which has $M$ service types (branches) and $i_m$ tasks on its $m$-th branch $(m = 1, 2, \ldots, M)$, satisfying $\sum_{m=1}^{M} i_m = K$. For example, a system with $M = 1$ and $i_1 = K$ represents a tandem queueing network with $K$ service stations, while a system with $i_1 = i_2 = \cdots = i_M = 1$

represents a parallel queueing network with $M$ service stations. Such a parallel-series system arises where the service needs can be classified into $M$ types. The $m$-th service type consists of a series of $i_m$ tasks, each taking a random amount of time.



Figure 3.1: A closed queueing system with parallel-series service stations.

We denote the $i$th node in the $m$-th branch as node $(i, m)$. A customer stays in node 0 for a random amount of time and then moves to node $(1, m)$ with probability $p_m > 0$ $(1 \leq m \leq M)$, where $\sum_{m=1}^{M} p_m = 1$. A customer stays in node $(i, m)$ until she receives service from the server, then moves to node $(i + 1, m)$ if $i < i_m$, and to node 0 if $i = i_m$. This process repeats forever. Let $S_{j,m}$ represent the random service time performed by the server in node $(j, m)$ $(1 \leq j \leq i_m,\ 1 \leq m \leq M)$. We assume that all service times are independent of each other.

Let $D_0^{\pi}(t)$ and $D_{(i,m)}^{\pi}(t)$ denote the number of service completions in node 0 and in node $(i, m)$ $(1 \leq i \leq i_m,\ 1 \leq m \leq M)$, respectively during $[0, t]$ under policy $\pi$, where $\pi$ is either in $\Pi_P$ or $\Pi_{NP}$. Suppose that a finite reward $R_{(i,m)}$ is gained when service is completed in node $(i, m)$ $(1 \leq i \leq i_m,\ 1 \leq m \leq M)$. We define

$$R^{\pi} \equiv \liminf_{t \to \infty} \sum_{m=1}^{M} \sum_{i=1}^{i_m} \frac{R_{(i,m)} D_{(i,m)}^{\pi}(t)}{t}, \qquad (3.1.1)$$

which denotes the long-run average reward of the system under policy $\pi$.

11

Let $TH_0^\pi$ denote the long-run average throughput of station 0 under policy $\pi$, i.e.,

$$TH_0^\pi \equiv \liminf_{t\to\infty} \frac{D_0^\pi(t)}{t}. \qquad (3.1.2)$$

Throughout the paper, we will refer to $TH_0^\pi$ as the system throughput. We first show that maximizing the long-run average reward of the system is equivalent to maximizing the long-run average throughput of the system.

**Theorem 3.1.1.** *For the parallel-series system, maximizing the long-run average reward of the system is equivalent to maximizing the long-run average throughput of the system.*

***Proof of Theorem 3.1.1.*** We define $TH_{(i,m)}^\pi$ as the long-run average throughput from node $(i,m)$, i.e.,

$$TH_{(i,m)}^\pi \equiv \liminf_{t\to\infty} \frac{D_{(i,m)}^\pi(t)}{t}.$$

We first show that the long-run average throughput at any two consecutive nodes are equal, i.e., $TH_{(i,m)}^\pi = TH_{(i+1,m)}^\pi$ ($1 \le i \le i_m - 1$, $1 \le m \le M$). Let $C_{(i,m)}^\pi(t)$ denote the number of customers in node $(i,m)$ at time $t$ ($1 \le i \le i_m$ and $1 \le m \le M$). Then, we have

$$C_{(i+1,m)}^\pi(t) = C_{(i+1,m)}^\pi(0) + D_{(i,m)}^\pi(t) - D_{(i+1,m)}^\pi(t), \quad \text{for } 1 \le i \le i_m - 1 \text{ and } 1 \le m \le M,$$

which implies

$$\liminf_{t\to\infty} \frac{D_{(i,m)}^\pi(t)}{t} = \liminf_{t\to\infty} \frac{D_{(i+1,m)}^\pi(t)}{t} + \liminf_{t\to\infty} \frac{C_{(i+1,m)}^\pi(t)}{t} - \liminf_{t\to\infty} \frac{C_{(i+1,m)}^\pi(0)}{t}. \qquad (3.1.3)$$

Since $C_{(i+1,m)}^\pi(t) \le B < \infty$ for all $t \ge 0$, the last two terms in (3.1.3) are equal to 0. Hence,

$$TH_{(i,m)}^\pi = TH_{(i+1,m)}^\pi, \quad \text{for } 1 \le i \le i_m - 1 \text{ and } 1 \le m \le M. \qquad (3.1.4)$$

Now, let $D_{0,m}^\pi(t)$ be the number of customers who request an $m$-th type of service after leaving node 0 during $[0,t]$ under policy $\pi$, and hence $\sum_{m=1}^M D_{0,m}^\pi(t) = D_0^\pi(t)$. A similar argument as above leads to

$$\liminf_{t\to\infty} \frac{D_{0,m}^\pi(t)}{t} = \liminf_{t\to\infty} \frac{D_{(1,m)}^\pi(t)}{t}.$$

By the law of large numbers, we know that

$$\liminf_{t\to\infty} \frac{D_{0,m}^\pi(t)}{D_0^\pi(t)} = p_m,$$

and hence using (3.1) we have

$$TH^\pi_{(1,m)} = p_m TH^\pi_0. \tag{3.1.5}$$

By (3.1.4) and (3.1.5), we can show that

$$
\begin{aligned}
R^\pi &= \liminf_{t \to \infty} \sum_{m=1}^{M} \sum_{i=1}^{i_m} \frac{R_{(i,m)} D^\pi_{(i,m)}(t)}{t} \\
&= \sum_{m=1}^{M} \sum_{i=1}^{i_m} R_{(i,m)} TH^\pi_{(i,m)} \\
&= \sum_{m=1}^{M} TH^\pi_{(1,m)} \sum_{i=1}^{i_m} R_{(i,m)} \\
&= TH^\pi_0 \sum_{m=1}^{M} p_m \sum_{i=1}^{i_m} R_{(i,m)}.
\end{aligned}
$$

Since $\sum_{m=1}^{M} p_m \sum_{i=1}^{i_m} R_{(i,m)}$ is a constant for given $R_{(i,m)}$ and $p_m$ $(1 \le i \le i_m, 1 \le m \le M)$, maximizing $R^\pi$ is equivalent to maximizing $TH^\pi_0$. $\qquad\square$

In the following discussions, our objective is to solve the following two optimization problems in order to maximize the long-run average throughput of the system:

$$\max_{\pi \in \Pi_P} \quad TH^\pi_0$$

and

$$\max_{\pi \in \Pi_{NP}} \quad TH^\pi_0.$$

### 3.2  Discrete Resource Constraint with a Single Server: Preemption

In this section, we consider allocating a single flexible server under preemptive policies. Assume that service times at station $k$ are exponentially distributed with rate $\mu_k$ where $0 < \mu_k < \infty$.

#### 3.2.1  Series System

We first study a series system as shown in Figure 3.2. We formulate the optimization problem as a Markov decision process. Let $\boldsymbol{n} = (n_1, n_2, \ldots, n_K)$ denote the state of the system, where $n_k \ge 0$ represents the number of jobs at station $k$ $(1 \le k \le K, 0 \le \sum_{j=1}^{K} n_j \le B)$. Let $n = \sum_{k=1}^{K} n_k$ denote the total number of jobs at the service stations in state $\boldsymbol{n}$. Let

13

$I_n \subseteq \{1, 2, \ldots, K\}$ denote the set of service stations that are non-empty in state $n$. Also, let $e^i$ denote a $K$-dimensional row vector with all elements 0 except where the $i$th element is equal to 1, and denote by $\mathbf{0}$ a $K$-dimensional row vector with all elements equal to 0. Define $V(n)$ as the bias of state $n$, and $g$ as the long-run average throughput of the system.



Figure 3.2: A closed queueing system with an automated station and $K$ tandem service stations.

Because both the state space and the action space are finite and the transition matrix consists of a single recurrent class for every deterministic stationary policy, the MDP under study is unichain. Hence, we know that there exists a stationary average optimal policy and hence $g$ exists (see, e.g., Theorem 8.4.5 in Puterman [22]). Define $\Lambda = B\mu_0 + \sum_{k=1}^{K} \mu_k$ as the uniformization constant. Without loss of generality, we assume that $\Lambda = 1$. Then, the optimality equation can be expressed as follows:

$$g + V(n) = (B - n)\mu_0 V(n + e^1) + n\mu_0 V(n) + f(n), \quad \text{for } 0 \leq n \leq B, \quad (3.2.1)$$

where

$$f(n) = \sum_{k=1}^{K} \mu_k V(n) + \begin{cases} 0, & \text{if } n = \mathbf{0} \\ \max_{i \in I_n} \Big\{ \mu_i [V(n - e^i + e^{i+1}) - V(n)] \mathbb{1}_{\{i \neq K\}}, \\ \qquad \mu_K [V(n - e^K) + 1 - V(n)] \mathbb{1}_{\{i = K\}} \Big\}, & \text{otherwise,} \end{cases}$$

where $\mathbb{1}_{\{A\}}$ is an indicator function with value of 1 if $A$ holds or value of 0 otherwise. Here, we use the fact that the throughput from each station in a tandem line is the same (see the proof of Theorem 3.1.1). We provide a complete characterization of the optimal policy in Theorem 3.2.1.

**Theorem 3.2.1.** *The policy that gives priority to the non-empty station with the largest index maximizes the long-run average throughput of the system within the set of all preemptive policies* $\Pi_P$.

**Proof of Theorem 3.2.1.** In order to prove Theorem 3.2.1, we first show that the result holds for a similar finite horizon problem given by (3.2.2) with $m$ periods for all $m \geq 0$. Let $\boldsymbol{N}_k$ denote the state of the system at period $k$ and $d_k(N_k)$ the decision rule at period $k$ in state $N_k$ under policy $\pi$. Let $r(\boldsymbol{N}, d)$ denote the gained throughput when the system is in state $\boldsymbol{N}$ and the action $d$ is taken. We define $V_m(\pi, \boldsymbol{n})$ as the $m$-period expected throughput under policy $\pi$ when the initial state is $\boldsymbol{n}$, i.e.,

$$V_m(\pi, \boldsymbol{n}) \equiv E\left[\sum_{k=0}^{m-1} r(\boldsymbol{N}_k, d_k(\boldsymbol{N}_k))\right].$$

Then, the optimal $m$-period expected total throughput is

$$V_m^*(\boldsymbol{n}) \equiv \sup_{\pi \in \Pi_P} V_m(\pi, \boldsymbol{n}). \tag{3.2.2}$$

We let $g(\pi, \boldsymbol{n})$ be the long-run average throughput under policy $\pi$, given that the initial state of the system is $\boldsymbol{n}$, i.e.,

$$g(\pi, \boldsymbol{n}) \equiv \liminf_{m \to \infty} \frac{1}{m} V_m(\pi, \boldsymbol{n}).$$

Let $\mu \equiv \sum_{k=1}^{K} \mu_k$. Then, the optimality equation for the finite-period problem can be expressed as follows. For all $m \geq 0$,

$$V_{m+1}(\boldsymbol{n}) = (B - n)\mu_0 V_m(\boldsymbol{n} + \boldsymbol{e}^1) + n\mu_0 V_m(\boldsymbol{n}) + f_m(\boldsymbol{n}), \quad \text{for } 0 \leq n \leq B, \tag{3.2.3}$$

where

$$f_m(\boldsymbol{n}) = \mu V_m(\boldsymbol{n}) + \begin{cases} 0, & \text{if } \boldsymbol{n} = \boldsymbol{0} \\ \max_{i \in \boldsymbol{I}_n} \Big\{ [\mu_i V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1}) - \mu_i V_m(\boldsymbol{n})]\mathbb{1}_{\{i \neq K\}}, \\ \qquad [\mu_K(V_m(\boldsymbol{n} - \boldsymbol{e}^K) + 1) - \mu_K V_m(\boldsymbol{n})]\mathbb{1}_{\{i=K\}}\Big\}, & \text{otherwise.} \end{cases}$$

We assume that $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$. For system state $\boldsymbol{n}$, where $2 \leq n \leq B$, let $l_1(\boldsymbol{n}) \equiv \max\{k : k \in \boldsymbol{I}_n\}$ and $l_2(\boldsymbol{n}) \equiv \max\{k : k \in \boldsymbol{I}_{\boldsymbol{n}-\boldsymbol{e}^{l_1(\boldsymbol{n})}}\}$. We will show that, for all $m \geq 0$, $2 \leq n \leq B$, and $j \in \boldsymbol{I}_{\boldsymbol{n}-\boldsymbol{e}^{l_1(\boldsymbol{n})}}$,

i. if $l_1(\boldsymbol{n}) < K$,

$$\mu_{l_1(\boldsymbol{n})} V_m(\boldsymbol{n} - \boldsymbol{e}^{l_1(\boldsymbol{n})} + \boldsymbol{e}^{l_1(\boldsymbol{n})+1}) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_{l_1(\boldsymbol{n})})V_m(\boldsymbol{n}) \geq 0;$$

$$\tag{3.2.4}$$

ii. if $l_1(\boldsymbol{n}) = K$,

$$\mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_K)V_m(\boldsymbol{n}) \geq -\mu_K; \qquad (3.2.5)$$

iii.

$$V_m(\boldsymbol{n} + \boldsymbol{e}^1) - V_m(\boldsymbol{n}) \geq 0. \qquad (3.2.6)$$

We will use induction on $m$. Because $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$, the inequalities automatically hold at period 0. Assume that the inequalities hold at period $m$. We will show that they also hold at period $m + 1$. In the remainder of this proof, we will let $i = l_1(\boldsymbol{n})$ for ease of notation whenever it does not cause any ambiguity.

Proof of (3.2.4): We will consider two cases.

(a) Suppose that $l_1(\boldsymbol{n}) < K - 1$. Using equation (3.2.3), we have

$$\mu_i V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1}) - \mu_j V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_i)V_{m+1}(\boldsymbol{n})$$

$$= (B - n)\mu_0[\mu_i V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1} + \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} + \boldsymbol{e}^1) + (\mu_j - \mu_i)V_m(\boldsymbol{n} + \boldsymbol{e}^1)]$$

$$+ n\mu_0[\mu_i V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1}) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_i)V_m(\boldsymbol{n})]$$

$$+ \mu_{i+1}\mu_i V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+2}) - \mu_i \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1})$$

$$\qquad + \mu_i(\mu_j - \mu_i)V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1})$$

$$+ (\mu - \mu_{i+1})\mu_i V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1}) - (\mu - \mu_i)\mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1})$$

$$\qquad + (\mu - \mu_i)(\mu_j - \mu_i)V_m(\boldsymbol{n})$$

$$= (B - n)\mu_0[\mu_i V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1} + \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} + \boldsymbol{e}^1) + (\mu_j - \mu_i)V_m(\boldsymbol{n} + \boldsymbol{e}^1)]$$

$$+ \mu_i[\mu_{i+1}V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+2}) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1})$$

$$\qquad + (\mu_j - \mu_{i+1})V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1})]$$

$$+ (n\mu_0 + \mu - \mu_i)[\mu_i V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1}) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_i)V_m(\boldsymbol{n})],$$

which is non-negative by the inductive hypothesis for (3.2.4) at period $m$, the facts that $l_1(\boldsymbol{n} + \boldsymbol{e}^1) = l_1(\boldsymbol{n})$ and $l_1(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^{i+1}) = i + 1$, and the assumption that $B\mu_0 + \mu = 1$.

(b) Suppose that $l_1(\boldsymbol{n}) = K - 1$. Using equation (3.2.3), we have

$$\mu_{K-1}V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K) - \mu_j V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_{K-1})V_{m+1}(\boldsymbol{n})$$

$$= (B - n)\mu_0[\mu_{K-1}V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K + \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} + \boldsymbol{e}^1)$$

$$+ (\mu_j - \mu_{K-1})V_m(\boldsymbol{n} + \boldsymbol{e}^1)]$$

$$+ n\mu_0[\mu_{K-1}V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_{K-1})V_m(\boldsymbol{n})]$$

$$+ \mu_K\mu_{K-1}[V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1}) + 1] - \mu_{K-1}\mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K)$$

$$+ \mu_{K-1}(\mu_j - \mu_{K-1})V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K)$$

$$+ (\mu - \mu_K)\mu_{K-1}V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K) - (\mu - \mu_{K-1})\mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1})$$

$$+ (\mu - \mu_{K-1})(\mu_j - \mu_{K-1})V_m(\boldsymbol{n})$$

$$= (B - n)\mu_0[\mu_{K-1}V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K + \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} + \boldsymbol{e}^1)$$

$$+ (\mu_j - \mu_{K-1})V_m(\boldsymbol{n} + \boldsymbol{e}^1)]$$

$$+ \mu_{K-1}[\mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1}) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K)$$

$$+ (\mu_j - \mu_K)V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K) + \mu_K]$$

$$+ (n\mu_0 + \mu - \mu_{K-1})[\mu_{K-1}V_m(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1})$$

$$+ (\mu_j - \mu_{K-1})V_m(\boldsymbol{n})],$$

which is non-negative by the inductive hypothesis for (3.2.4) and (3.2.5) at period $m$, the facts that $l_1(\boldsymbol{n} + \boldsymbol{e}^1) = K - 1$ and $l_1(\boldsymbol{n} - \boldsymbol{e}^{K-1} + \boldsymbol{e}^K) = K$, and the assumption that $B\mu_0 + \mu = 1$.

Proof of (3.2.5). We will consider two cases:

(a) Suppose that $l_2(\boldsymbol{n}) = K$. Using equation (3.2.3), we have

$$\mu_K V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^K) - \mu_j V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_K)V_{m+1}(\boldsymbol{n})$$

$$=(B - n + 1)\mu_0 \mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K + \boldsymbol{e}^1)$$

$$+ (B - n)\mu_0[-\mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} + \boldsymbol{e}^1)) + (\mu_j - \mu_K)V_m(\boldsymbol{n} + \boldsymbol{e}^1)]$$

$$+ (n - 1)\mu_0 \mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K) + n\mu_0[-\mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_K)V_m(\boldsymbol{n})]$$

$$+ \mu_K \mu_K - \mu_j \mu_K + (\mu_j - \mu_K)\mu_K$$

$$+ \mu_K[\mu_K V_m(\boldsymbol{n} - 2\boldsymbol{e}^K) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} - \boldsymbol{e}^K) + (\mu_j - \mu_K)V_m(\boldsymbol{n} - \boldsymbol{e}^K)]$$

$$+ (\mu - \mu_K)[\mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_K)V_m(\boldsymbol{n})]$$

$$=(B - n)\mu_0[\mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K + \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} + \boldsymbol{e}^1) + (\mu_j - \mu_K)V_m(\boldsymbol{n} + \boldsymbol{e}^1)]$$

$$+ \mu_0[\mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K + \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_K)V_m(\boldsymbol{n})]$$

$$+ \mu_K[\mu_K V_m(\boldsymbol{n} - 2\boldsymbol{e}^K) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} - \boldsymbol{e}^K) + (\mu_j - \mu_K)V_m(\boldsymbol{n} - \boldsymbol{e}^K)]$$

$$+ ((n - 1)\mu_0 + \mu - \mu_K)[\mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_K)V_m(\boldsymbol{n})],$$

which is greater than or equal to $-\mu_K$ by the inductive hypothesis for (3.2.5) and (3.2.6) at period $m$, the facts that $l_1(\boldsymbol{n} + \boldsymbol{e}^1) = K$ and $l_1(\boldsymbol{n} - \boldsymbol{e}^K) = K$, and the assumption that $B\mu_0 + \mu=1$.

(b) Suppose that $l_2(\boldsymbol{n}) < K$. Using equation (3.2.3), we have

$$\mu_K V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^K) - \mu_j V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_K)V_{m+1}(\boldsymbol{n})$$

$$=(B - n + 1)\mu_0 \mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K + \boldsymbol{e}^1)$$

$$+ (B - n)\mu_0[-\mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} + \boldsymbol{e}^1)) + (\mu_j - \mu_K)V_m(\boldsymbol{n} + \boldsymbol{e}^1)]$$

$$+ (n - 1)\mu_0 \mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K) + n\mu_0[-\mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1}) + (\mu_j - \mu_K)V_m(\boldsymbol{n})]$$

$$- \mu_j \mu_K + (\mu_j - \mu_K)\mu_K$$

$$+ \mu_K \mu_{l_2(\boldsymbol{n})} V_m(\boldsymbol{n} - \boldsymbol{e}^K - \boldsymbol{e}^{l_2(\boldsymbol{n})} + \boldsymbol{e}^{l_2(\boldsymbol{n})+1}) - \mu_j \mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1} - \boldsymbol{e}^K)$$

$$+ (\mu_j - \mu_K)\mu_K V_m(\boldsymbol{n} - \boldsymbol{e}^K)$$

$$+ \mu_K(\mu - \mu_{l_2(\boldsymbol{n})})V_m(\boldsymbol{n} - \boldsymbol{e}^K) - \mu_j(\mu - \mu_K)V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^{j+1})$$

$$+ (\mu_j - \mu_K)(\mu - \mu_K)V_m(\boldsymbol{n})$$

18

$$=(B-n)\mu_0[\mu_K V_m(\boldsymbol{n}-\boldsymbol{e}^K+\boldsymbol{e}^1)-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j+\boldsymbol{e}^{j+1}+\boldsymbol{e}^1)+(\mu_j-\mu_K)V_m(\boldsymbol{n}+\boldsymbol{e}^1)]$$
$$+\mu_0[\mu_K V_m(\boldsymbol{n}-\boldsymbol{e}^K+\boldsymbol{e}^1)-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j+\boldsymbol{e}^{j+1})+(\mu_j-\mu_K)V_m(\boldsymbol{n})]$$
$$+(n-1)\mu_0[\mu_K V_m(\boldsymbol{n}-\boldsymbol{e}^K)-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j+\boldsymbol{e}^{j+1})+(\mu_j-\mu_K)V_m(\boldsymbol{n})]$$
$$+\mu_K[\mu_{l_2(\boldsymbol{n})}V_m(\boldsymbol{n}-\boldsymbol{e}^K-\boldsymbol{e}^{l_2(\boldsymbol{n})}+\boldsymbol{e}^{l_2(\boldsymbol{n})+1})-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j+\boldsymbol{e}^{j+1}-\boldsymbol{e}^K)$$
$$+(\mu_j-\mu_{l_2(\boldsymbol{n})})V_m(\boldsymbol{n}-\boldsymbol{e}^K)-\mu_K]$$
$$+(\mu-\mu_K)[\mu_K V_m(\boldsymbol{n}-\boldsymbol{e}^K)-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j+\boldsymbol{e}^{j+1})+(\mu_j-\mu_K)V_m(\boldsymbol{n})],$$

which is greater than or equal to $-\mu_K$ by the inductive hypothesis for (3.2.4), (3.2.5) and (3.2.6) at period $m$, the facts that $l_1(\boldsymbol{n}+\boldsymbol{e}^1)=K$ and $l_1(\boldsymbol{n}-\boldsymbol{e}^K)=l_2(\boldsymbol{n})<K$, and the assumption that $B\mu_0+\mu=1$.

Proof of (3.2.6). We will consider two cases:

(a) Suppose that $l_1(\boldsymbol{n})<K$. Using equation (3.2.3), we have

$$V_{m+1}(\boldsymbol{n}+\boldsymbol{e}^1)-V_{m+1}(\boldsymbol{n})$$
$$=(B-n-1)\mu_0 V_m(\boldsymbol{n}+2\boldsymbol{e}^1)-(B-n)\mu_0 V_m(\boldsymbol{n}+\boldsymbol{e}^1)$$
$$+(n+1)\mu_0 V_m(\boldsymbol{n}+\boldsymbol{e}^1)-n\mu_0 V_m(\boldsymbol{n})$$
$$+\mu_{l_1(\boldsymbol{n})}[V_m(\boldsymbol{n}+\boldsymbol{e}^1-\boldsymbol{e}^{l_1(\boldsymbol{n})}+\boldsymbol{e}^{l_1(\boldsymbol{n})+1})-V_m(\boldsymbol{n}-\boldsymbol{e}^{l_1(\boldsymbol{n})}+\boldsymbol{e}^{l_1(\boldsymbol{n})+1})]$$
$$+(\mu-\mu_{l_1(\boldsymbol{n})})[V_m(\boldsymbol{n}+\boldsymbol{e}^1)-V_m(\boldsymbol{n})]$$
$$=(B-n-1)\mu_0[V_m(\boldsymbol{n}+2\boldsymbol{e}^1)-V_m(\boldsymbol{n}+\boldsymbol{e}^1)]$$
$$+\mu_{l_1(\boldsymbol{n})}[V_m(\boldsymbol{n}+\boldsymbol{e}^1-\boldsymbol{e}^{l_1(\boldsymbol{n})}+\boldsymbol{e}^{l_1(\boldsymbol{n})+1})-V_m(\boldsymbol{n}-\boldsymbol{e}^{l_1(\boldsymbol{n})}+\boldsymbol{e}^{l_1(\boldsymbol{n})+1})]$$
$$+(n\mu_0+\mu-\mu_{l_1(\boldsymbol{n})})[V_m(\boldsymbol{n}+\boldsymbol{e}^1)-V_m(\boldsymbol{n})],$$

which is non-negative by the inductive hypothesis for (3.2.6) at period $m$ and the assumption that $B\mu_0+\mu=1$.

(b) Suppose that $l_1(\boldsymbol{n}) = K$. Using equation (3.2.3), we have

$$V_{m+1}(\boldsymbol{n} + \boldsymbol{e}^1) - V_{m+1}(\boldsymbol{n})$$

$$=(B - n - 1)\mu_0 V_m(\boldsymbol{n} + 2\boldsymbol{e}^1) - (B - n)\mu_0 V_m(\boldsymbol{n} + \boldsymbol{e}^1)$$

$$+ (n + 1)\mu_0 V_m(\boldsymbol{n} + \boldsymbol{e}^1) - n\mu_0 V_m(\boldsymbol{n})$$

$$+ \mu_K[V_m(\boldsymbol{n} + \boldsymbol{e}^1 - \boldsymbol{e}^K) - V_m(\boldsymbol{n} - \boldsymbol{e}^K)]$$

$$+ (\mu - \mu_K)[V_m(\boldsymbol{n} + \boldsymbol{e}^1) - V_m(\boldsymbol{n})]$$

$$=(B - n - 1)\mu_0[V_m(\boldsymbol{n} + 2\boldsymbol{e}^1) - V_m(\boldsymbol{n} + \boldsymbol{e}^1)]$$

$$+ \mu_K[V_m(\boldsymbol{n} + \boldsymbol{e}^1 - \boldsymbol{e}^K) - V_m(\boldsymbol{n} - \boldsymbol{e}^K)]$$

$$+ (n\mu_0 + \mu - \mu_K)[V_m(\boldsymbol{n} + \boldsymbol{e}^1) - V_m(\boldsymbol{n})],$$

which is non-negative by the inductive hypothesis for (3.2.6) at period $m$ and the assumption that $B\mu_0 + \mu = 1$.

Let $\pi^*$ be the policy that gives priority to the non-empty station with the largest index in the series system. By (3.2.4) and (3.2.5), we have

$$V_m(\pi^*, \boldsymbol{n}) \geq V_m(\pi, \boldsymbol{n}) \tag{3.2.7}$$

for all $\pi \in \Pi_P$ and for all $m \geq 0$. Dividing both sides of (3.2.7) by $m$ and taking limits as $m$ approaches infinity the long-run average throughput result follows, i.e.,

$$g(\pi^*, \boldsymbol{n}) \geq g(\pi, \boldsymbol{n})$$

for all $\pi \in \Pi_P$. Hence, policy $\pi^*$ maximizes the long-run average throughput of the system.

$\square$

**Remarks.** Theorem 3.2.1 shows that we should put the server to work at the station which is as close to the entry into station 0 as possible when a job is available. The intuition is that the earlier a job goes back to the automated station, the earlier this job leaves station 0 to request for service, which increases the utilization of the server and the throughput of the system as well. Note that the optimal policy eventually becomes a policy under which the server picks a job from the queue in front of station 1 and completes the service of this job at all service stations $1, 2, \ldots, K$ before it starts working on another job waiting in front of station 1. We call this policy a sequential policy.

In the set of preemptive policies, the optimal policy for the series system shown in Theorem 3.2.1 is actually non-preemptive. The server works on a job from the first station to last station sequentially without being interrupted by new arrivals from station 0. Hence, the sequential policy is also optimal within $\Pi_{NP}$ under the Markovian case.

### 3.2.2 Parallel System

Next, we study a parallel system as shown in Figure 3.3. We formulate the optimization problem as a Markov decision process (MDP). We use the same notation defined in Section



Figure 3.3: A closed queueing system with an automated station and $K$ parallel service stations.

3.2.1 unless otherwise stated. Following the same argument, we know that there exists a stationary average optimal policy and hence $g$ exists. Without loss of generality, we assume that $\Lambda = 1$. Then, the optimality equation can be expressed as follows:

$$g + V(\boldsymbol{n}) = (B - n)\mu_0 \sum_{k=1}^{K} p_k V(\boldsymbol{n} + \boldsymbol{e}^k) + n\mu_0 V(\boldsymbol{n}) + f(\boldsymbol{n}), \quad \text{for } 0 \leq n \leq B, \quad (3.2.8)$$

where

$$f(\boldsymbol{n}) = \sum_{k=1}^{K} \mu_k V(\boldsymbol{n}) + \begin{cases} 0 & \text{if } \boldsymbol{n} = \boldsymbol{0}, \\ \max_{i \in \boldsymbol{I_n}} \{\mu_i(V(\boldsymbol{n} - \boldsymbol{e}^i) + 1) - \mu_i V(\boldsymbol{n})\} & \text{otherwise.} \end{cases}$$

We provide a partial characterization of the optimal policy in Theorem 3.2.2.

**Theorem 3.2.2.** *Suppose that there exists a station $i$ for which $\mu_i \geq \mu_j$, for all $j = 1, 2, \ldots, K$ and $j \neq i$. Then, there exists an optimal policy which gives the highest priority to station $i$ within the set of all preemptive policies $\Pi_P$.*

***Proof of Theorem 3.2.2.*** In order to prove Theorem 3.2.2 we show that the result holds for the $m$-period expected total throughput problem defined by (3.2.2) for all $m \geq 0$. Then, the optimality equation for the finite-period problem can be expressed as follows. For all $m \geq 0$,

$$V_{m+1}(\boldsymbol{n}) = (B-n)\mu_0 \sum_{k=1}^{K} p_k V_m(\boldsymbol{n} + \boldsymbol{e}^k) + n\mu_0 V_m(\boldsymbol{n}) + f_m(\boldsymbol{n}), \quad \text{for } 0 \leq n \leq B, \quad (3.2.9)$$

where

$$f_m(\boldsymbol{n}) = \sum_{k=1}^{K} \mu_k V_m(\boldsymbol{n}) + \begin{cases} 0 & \text{if } \boldsymbol{n} = \boldsymbol{0}, \\ \max_{i \in \boldsymbol{I_n}} \{\mu_i(V_m(\boldsymbol{n} - \boldsymbol{e}^i) + 1) - \mu_i V_m(\boldsymbol{n})\} & \text{otherwise.} \end{cases}$$

We assume that $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$.

Without loss of generality, we assume that $\mu_1 \geq \mu_j$, for all $j = 2, 3, \ldots, K$. For system state $\boldsymbol{n}$, where $2 \leq n \leq B$, let $l_1(\boldsymbol{n}) \equiv \min\{k : k \in \boldsymbol{I_n}\}$ and $l_2(\boldsymbol{n}) \equiv \min\{k : k \in \boldsymbol{I_{n-e^{l_1(n)}}}\}$. Let $\boldsymbol{n}^a = \boldsymbol{n} + \boldsymbol{e}^a$ ($1 \leq a \leq K$). We will show that, for all $m \geq 0$, $2 \leq n \leq B$, $n_1 \geq 1$, and $1 \leq a \leq K$,

$$\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n}) \geq \mu_j - \mu_1, \quad (3.2.10)$$

$$\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n}) \geq \mu_j - \mu_1, \quad (3.2.11)$$

where $j \neq 1$ and $j \in \boldsymbol{I_n}$. We will use induction on $m$. Since $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$, (3.2.10) and (3.2.11) automatically hold at period 0. Assume that inequalities (3.2.10) and (3.2.11) hold at period $m$. We will show that they also hold at period $m + 1$.

Proof of (3.2.10): We will consider two cases.

(a) Suppose that $l_2(\boldsymbol{n}) = 1$. Using equation (3.2.9), we have

$$\mu_1 V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_{m+1}(\boldsymbol{n})$$

$$=(B - n + 1)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k))]$$

$$+ (B - n)\mu_0 \sum_k p_k(\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)$$

$$+ (n - 1)\mu_0[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j)] + n\mu_0[(\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1\mu_1 - \mu_j\mu_1 + (\mu_j - \mu_1)\mu_1 + \mu_1[\mu_1 V_m(\boldsymbol{n} - 2\boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^1)$$

$$+ (\mu_j - \mu_1)V_m(\boldsymbol{n} - \boldsymbol{e}^1)]$$

$$+ (\mu - \mu_1)[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$=(B - n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k) + (\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ \mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k)) + (\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1[\mu_1 V_m(\boldsymbol{n} - 2\boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_1)V_m(\boldsymbol{n} - \boldsymbol{e}^1)]$$

$$+ ((n - 1)\mu_0 + \mu - \mu_1)[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})],$$

which is greater than or equal to $\mu_j - \mu_1$ by the inductive hypothesis for (3.2.10) and (3.2.11) at period $m$, the fact that $\mu_1 \geq \mu_j$ for all $j \in \boldsymbol{I_n}$, and the assumption that $B\mu_0 + \mu = 1$.

(b) Suppose that $l_2(\boldsymbol{n}) > 1$. Using equation (3.2.9), we have

$$\mu_1 V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_{m+1}(\boldsymbol{n})$$

$$=(B - n + 1)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k))]$$

$$+ (B - n)\mu_0 \sum_k p_k(\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)$$

$$+ (n - 1)\mu_0[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j)] + n\mu_0(\mu_j - \mu_1)V_m(\boldsymbol{n})$$

$$+ \mu_1\mu_{l_2(\boldsymbol{n})} - \mu_j\mu_1 + (\mu_j - \mu_1)\mu_1$$

$$+ \mu_1\mu_{l_2(\boldsymbol{n})}V_m(\boldsymbol{n} - \boldsymbol{e}^1 - \boldsymbol{e}^{l_2(\boldsymbol{n})}) - \mu_j\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_1)\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1)$$

$$+ \mu_1(\mu - \mu_{l_2(\boldsymbol{n})})V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j(\mu - \mu_1)V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)(\mu - \mu_1)V_m(\boldsymbol{n})$$

$$=(B-n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}-\boldsymbol{e}^1+\boldsymbol{e}^k)-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j+\boldsymbol{e}^k)+(\mu_j-\mu_1)V_m(\boldsymbol{n}+\boldsymbol{e}^k)]$$

$$+\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}-\boldsymbol{e}^1+\boldsymbol{e}^k)-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j+\boldsymbol{e}^k))+(\mu_j-\mu_1)V_m(\boldsymbol{n})]$$

$$+\mu_1[\mu_{l_2(\boldsymbol{n})}V_m(\boldsymbol{n}-\boldsymbol{e}^1-\boldsymbol{e}^{l_2(\boldsymbol{n})})-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j-\boldsymbol{e}^1)+(\mu_j-\mu_{l_2(\boldsymbol{n})})V_m(\boldsymbol{n}-\boldsymbol{e}^1)$$

$$+\mu_{l_2(\boldsymbol{n})}-\mu_1]$$

$$+((n-1)\mu_0+\mu-\mu_1)[\mu_1 V_m(\boldsymbol{n}-\boldsymbol{e}^1)-\mu_j V_m(\boldsymbol{n}-\boldsymbol{e}^j)+(\mu_j-\mu_1)V_m(\boldsymbol{n})],$$

which is greater than or equal to $\mu_j-\mu_1$ by the inductive hypothesis for (3.2.10) and (3.2.11) at period $m$ and $l_2(\boldsymbol{n})=l_1(\boldsymbol{n}-\boldsymbol{e}^1)$, the fact that $\mu_1 \geq \mu_j$ for all $j \in \boldsymbol{I_n}$, and the assumption that $B\mu_0+\mu=1$.

Proof of (3.2.11). We will consider two cases:

(a) Suppose that $l_2(\boldsymbol{n}^a)=1$. Using equation (3.2.9), we have

$$\mu_1 V_{m+1}(\boldsymbol{n}^a-\boldsymbol{e}^1)-\mu_j V_{m+1}(\boldsymbol{n}^a-\boldsymbol{e}^j)+(\mu_j-\mu_1)V_{m+1}(\boldsymbol{n})$$

$$=(B-n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}^a-\boldsymbol{e}^1+\boldsymbol{e}^k)-\mu_j V_m(\boldsymbol{n}^a-\boldsymbol{e}^j+\boldsymbol{e}^k))]$$

$$+(B-n)\mu_0 \sum_k p_k[(\mu_j-\mu_1)V_m(\boldsymbol{n}+\boldsymbol{e}^k)]$$

$$+n\mu_0[\mu_1 V_m(\boldsymbol{n}^a-\boldsymbol{e}^1)-\mu_j V_m(\boldsymbol{n}^a-\boldsymbol{e}^j)]$$

$$+n\mu_0[(\mu_j-\mu_1)V_m(\boldsymbol{n})]$$

$$+\mu_1\mu_1-\mu_j\mu_1+(\mu_j-\mu_1)\mu_1$$

$$+\mu_1\mu_1 V_m(\boldsymbol{n}^a-2\boldsymbol{e}^1)-\mu_j\mu_i V_m(\boldsymbol{n}^a-\boldsymbol{e}^j-\boldsymbol{e}^1)+(\mu_j-\mu_1)\mu_1 V_m(\boldsymbol{n}-\boldsymbol{e}^1)$$

$$+\mu_1(\mu-\mu_1)V_m(\boldsymbol{n}^a-\boldsymbol{e}^1)-\mu_j(\mu-\mu_i)V_m(\boldsymbol{n}^a-\boldsymbol{e}^j)+(\mu_j-\mu_1)(\mu-\mu_1)V_m(\boldsymbol{n})$$

$$=(B-n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}^a-\boldsymbol{e}^1+\boldsymbol{e}^k)-\mu_j V_m(\boldsymbol{n}^a-\boldsymbol{e}^j+\boldsymbol{e}^k)+(\mu_j-\mu_1)V_m(\boldsymbol{n}+\boldsymbol{e}^k)]$$

$$+n\mu_0[\mu_1 V_m(\boldsymbol{n}^a-\boldsymbol{e}^1)-\mu_j V_m(\boldsymbol{n}^a-\boldsymbol{e}^j)+(\mu_j-\mu_1)V_m(\boldsymbol{n})]$$

$$+\mu_1[\mu_1 V_m(\boldsymbol{n}^a-2\boldsymbol{e}^1)-\mu_j V_m(\boldsymbol{n}^a-\boldsymbol{e}^j-\boldsymbol{e}^1)+(\mu_j-\mu_1)V_m(\boldsymbol{n}-\boldsymbol{e}^1)]$$

$$+(\mu-\mu_1)[\mu_1 V_m(\boldsymbol{n}^a-\boldsymbol{e}^1)-\mu_j V_m(\boldsymbol{n}^a-\boldsymbol{e}^j)+(\mu_j-\mu_1)V_m(\boldsymbol{n})],$$

which is greater than or equal to $\mu_j-\mu_1$ by the inductive hypothesis for (3.2.11) at period $m$, the fact that $\mu_1 \geq \mu_j$ for all $j \in \boldsymbol{I_n}$, and the assumption that $B\mu_0+\mu=1$.

24

(b) Suppose that $l_2(\boldsymbol{n}^a) > 1$. Using equation (3.2.9), we have

$$\mu_1 V_{m+1}(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_{m+1}(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_{m+1}(\boldsymbol{n})$$

$$=(B-n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j + \boldsymbol{e}^k))]$$

$$+ (B-n)\mu_0 \sum_k p_k[(\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ n\mu_0[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j)] + n\mu_0[(\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1 \mu_{l_2(\boldsymbol{n}^a)} - \mu_j \mu_1 + (\mu_j - \mu_1)\mu_1$$

$$+ \mu_1 \mu_{l_2(\boldsymbol{n}^a)} V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 - \boldsymbol{e}^{l_2(\boldsymbol{n}^a)}) - \mu_j \mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_1)\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1)$$

$$+ \mu_1(\mu - \mu_{l_2(\boldsymbol{n}^a)})V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j(\mu - \mu_1)V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)(\mu - \mu_1)V_m(\boldsymbol{n})$$

$$=(B-n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j + \boldsymbol{e}^k) + (\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ n\mu_0[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1[\mu_{l_2(\boldsymbol{n}^a)} V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 - \boldsymbol{e}^{l_2(\boldsymbol{n}^a)}) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_{l_2(\boldsymbol{n}^a)})V_m(\boldsymbol{n} - \boldsymbol{e}^1)$$

$$+ \mu_{l_2(\boldsymbol{n}^a)} - \mu_1]$$

$$+ (\mu - \mu_1)[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})],$$

which is greater than or equal to $\mu_j - \mu_1$ by the inductive hypothesis for (3.2.10) and (3.2.11) at period $m$, the fact that $\mu_1 \geq \mu_j$ for all $j \in \boldsymbol{I_n}$, and the assumption that $B\mu_0 + \mu = 1$.

Hence, we show that jobs at station 1 should be served ahead of jobs at station $j$ ($j = 2, 3, \ldots, K$) if jobs are available at both stations. In other words, we should give the highest priority to station 1 in order to maximize the long-run average throughput of the system. □

**Remarks.** Theorem 3.2.2 says that we should give the highest priority to the service station which has the fastest service rate among all service stations. This follows the same intuition as in the series model where the optimal policy pushes jobs towards the automated station as early as possible. Even though Theorem 3.2.2 only gives a partial characterization of the optimal policy for $K > 2$, we believe that a complete characterization should follow a similar intuition as we state in the following conjecture.

25

**Conjecture 3.2.1.** *The policy that gives priority to the non-empty station with the largest service rate maximizes the long-run average throughput of the system within the set of all preemptive policies $\Pi_P$.*

A numerical study is conducted in Section 3.3.4 to support Conjecture 3.2.1.

### 3.3 Discrete Resource Constraint with a Single Server: Non-Preemption

In this section, we consider allocating a single flexible server under non-preemptive policies. In Sections 3.3.1, 3.3.2, and 3.3.3, we study a series system, a parallel system, and a two-branch three-station system, respectively. We characterize (either completely or partially) optimal policies for each case. In Section 3.3.4, we propose an heuristic policy and provide numerical results to compare the performance of the heuristic as opposed to the optimal policy.

#### 3.3.1 Series System

We first consider the series system as shown in Figure 3.2. Note that unlike in Section 3.2, we here do not make any distributional assumptions on the service times at the service stations $1, 2, \ldots, K$. We provide a complete characterization of the optimal policy in Theorem 3.3.1.

**Theorem 3.3.1.** *The optimal policy gives priority to the station which has the largest index.*

*Proof of Theorem 3.3.1.* We first show that the optimal policy gives the the highest priority to station $K$. Suppose policy $\pi$ is a policy under which there exists at least one decision epoch where the highest priority is not given to station $K$. Specifically, let $\varepsilon$ be the first time policy $\pi$ does not give priority to station $K$ even if there is a job in station $K$. Suppose $\{j_1, \ldots, j_m\}$ $(j_1, \ldots, j_m \neq K)$ gives the sequence of stations that the server visits after time $\varepsilon$ before it visits station $K$ under policy $\pi$. We will next construct a new policy $\gamma$ which serves station $K$ right before it serves the last job at station $j_m$.

Let $\tau$ be the time under policy $\pi$ that the server starts to work on a job at station $j_m$ with service time $S_j$ right before it moves to station $K$. Then after completing serving this job, the server immediately switches to station $K$ to serve a job there with service time $S_K$. We construct the new policy $\gamma$ as follows: $\gamma$ follows $\pi$ during $[0, \tau)$ and then serves a job at

station $K$ with service time $S'_K$, switches to station $j_m$ and serves a job with service time $S'_j$.

We directly couple the service times of all the jobs taken into service during $[0, \infty)$, which yields $S'_K = S_K$ and $S'_j = S_j$. Let $S_0$ be the service time at station 0 for the job entering station 0 at time $\tau + S_j + S_K$ under policy $\pi$ and for the job entering station 0 at time $\tau + S'_K$ under policy $\gamma$. $\gamma$ follows $\pi$ during $[\tau + S_j + S_K, \infty)$. This is possible because same or more jobs are available to policy $\gamma$ compared to policy $\pi$ after time $\tau + S_j + S_K$. The system states for the two policies will become identical after time $\tau + S_j + S_K + S_0$. The problem can be analyzed in three intervals as follows:

$$D_0^\gamma(t) = \begin{cases} D_0^\pi(t), & 0 \leq t < \tau + S_K + S_0, \\ D_0^\pi(t) + 1, & \tau + S_K + S_0 \leq t < \tau + S_j + S_K + S_0, \\ D_0^\pi(t), & t \geq \tau + S_j + S_K + S_0. \end{cases}$$

Therefore,

$$D_0^\gamma(t) - D_0^\pi(t) \geq 0, \text{ for all } t \geq 0,$$

and

$$TH_0^\gamma - TH_0^\pi = \liminf_{t \to \infty} \frac{1}{t}[D_0^\gamma(t) - D_0^\pi(t)] \geq 0.$$

We have shown that a job at station $K$ should be served ahead of a job at station $j_m$ if jobs are available at both stations. Following the same argument iteratively, we can show that a job at station $K$ should be served ahead of the sequence of stations $\{j_1, \ldots, j_{m-1}\}$. In other words, the optimal policy which maximizes the throughput of the system gives the highest priority to station $K$.

Next, we use an inductive argument to show that the optimal policy is a sequential policy. Suppose that the optimal policy gives priority to station $K$, station $K - 1$, ..., station $l + 1$, then the other stations, in the given order. We will show that the optimal policy gives $(K + 1 - l)$-th priority to station $l$.

Suppose policy $\pi$ is a policy which does not give $(K + 1 - l)$-th priority to station $l$. Specifically, let $\varepsilon$ be the first time policy $\pi$ does not give priority to a job in station $l$, when no jobs are available at stations $l+1, \ldots, K$. Suppose $\{j_1, \ldots, j_m\}$ $(j_1, \ldots, j_m < l)$ gives the sequence of stations that the server visits after time $\varepsilon$ before it visits station $l$ under policy

27

$\pi$. We will next construct a new policy $\gamma$ which serves a job at station $l$ before it serves the last job at station $j_m$.

Let $\tau$ be the time under policy $\pi$ that the server starts to work on the last job at station $j_m$ with service time $S_j$ before it moves to station $l$. Then after completing serving this job, the server immediately switches to station $l$ to serve the station $l$ job with service time $S_l$, then serves this job at stations $l+1, l+2, \ldots, K$ in the given order with service times $S_{l+1}, S_{l+2}, \ldots, S_K$, respectively. (Note that we do not need to consider any other actions for policy $\pi$ due to the inductive argument on the priority order of stations $l+1, l+2, \ldots, K$.) We construct the new policy $\gamma$ as follows: $\gamma$ follows $\pi$ during $[0, \tau)$ and then serves the next job at station $l$ with service time $S'_l$, then serves this job at stations $l+1, l+2, \ldots, K$ in the given order with service times $S'_{l+1}, S'_{l+2}, \ldots, S'_K$, respectively, then switches to station $j_m$ and serves a job with service time $S'_j$.

We directly couple the service times of all the jobs taken into service during $[0, \infty)$, which yields $S'_l = S_l$, $S'_{l+1} = S_{l+1}$, $\ldots$, $S'_K = S_K$, and $S'_j = S_j$. Let $S_0$ be the service time at station 0 for the job entering station 0 at time $\tau + S_j + S_l + S_{l+1} + \ldots + S_K$ under policy $\pi$ and for the job entering station 0 at time $\tau + S'_l + S'_{l+1} + \ldots + S'_K$ under policy $\gamma$. $\gamma$ follows $\pi$ during $[\tau + S_j + \sum_{k=l}^{K} S_k, \infty)$. This is possible because same or more jobs are available to policy $\gamma$ compared to policy $\pi$ after time $\tau + S_j + \sum_{k=l}^{K} S_k$. The system states for the two policies will become identical after time $\tau + S_j + \sum_{k=l}^{K} S_k + S_0$.

The problem can be analyzed in three intervals as below.

$$
D_0^\gamma(t) = \begin{cases} D_0^\pi(t), & 0 \le t < \tau + \sum_{k=l}^{K} S_k + S_0, \\ D_0^\pi(t) + 1, & \tau + \sum_{k=l}^{K} S_k + S_0 \le t < \tau + S_j + \sum_{k=l}^{K} S_k + S_0, \\ D_0^\pi(t), & t \ge \tau + S_j + \sum_{k=l}^{K} S_k + S_0. \end{cases}
$$

Therefore,

$$
D_0^\gamma(t) - D_0^\pi(t) \ge 0, \text{ for all } t \ge 0,
$$

and

$$
TH_0^\gamma - TH_0^\pi = \liminf_{t \to \infty} \frac{1}{t} [D_0^\gamma(t) - D_0^\pi(t)] \ge 0.
$$

We have shown that a job at station $l$ should be served ahead of a job at station $j_m$, where $j_m < l$, if jobs are available at both stations. Following the same argument iteratively, we can show that a job at station $l$ should be served ahead of the sequence of stations

$\{j_1, \ldots, j_{m-1}\}$. In other words, the optimal policy which maximizes the throughput of the system gives $(K + 1 - l)$-th priority to station $l$. Hence, an optimal policy that maximizes the long-run average throughput of the system is a sequential policy, which gives priority to station $K$, station $K - 1$, ..., and station 1 in the given order. $\qquad\square$

**Remarks.** Theorem 3.3.1 shows that for the series system, the optimal policy under non-preemptive policies suggests the same priority sequence among service stations as under preemptive policies. The optimal policy under non-preemption is also a sequential policy. In Section 3.2.1, we showed that the sequential policy is optimal within $\Pi_{np}$ under Markovian case. In this section, we show that the sequential policy is also optimal within $\Pi_{np}$ even when service times at the service stations are not exponentially distributed.

### 3.3.2   Parallel System

Next, we consider the parallel system shown in Figure 3.3. For ease of notation, we now define $X_k$ to be the random variable denoting the i.i.d. service time at station $k$ (instead of $S_{k,1}$) for $k = 1, 2, \ldots, K$. We give a partial characterization of the optimal policy for this system in Theorem 3.3.2.

We first define a specific stochastic order that we will use frequently in this section. Let $X$ and $Y$ be two continuous [or discrete] random variables with densities [or probability mass functions] $f(t)$ and $g(t)$, respectively, so that

$$\frac{g(t)}{f(t)} \text{ increases in } t \text{ over the union of the supports of } X \text{ and } Y,$$

or, equivalently,

$$f(x)g(y) \geq f(y)g(x) \quad \text{for all } x \leq y.$$

In this case, $X$ is said to be smaller than $Y$ in the likelihood ratio order (denoted by $X \leq_{lr} Y$). For more on stochastic orders, see, e.g., Shaked and Shanthikumar [24].

**Theorem 3.3.2.** *Suppose that there exists a station $i$ for which $X_i \leq_{lr} X_j$, for all $j = 1, 2, \ldots, K$ and $j \neq i$. Then, there exists an optimal policy which gives the highest priority to station $i$.*

**Proof of Theorem 3.3.2.** Without loss of generality, we assume that $i = 1$ and $X_1 \leq_{lr} X_j$, where $j \neq 1$. Suppose policy $\pi$ is a policy under which there exists at least one decision

29

epoch where the highest priority is not given to station 1. Specifically, let $\varepsilon$ be the first time policy $\pi$ does not give priority to station 1 even if there is a job in station 1. Suppose $\{j_1, \ldots, j_m\}$ $(j_1, \ldots, j_m > 1)$ gives the sequence of stations that the server visits after time $\varepsilon$ before it visits station 1 under policy $\pi$. We will next construct a new policy $\gamma$ which serves a job at station 1 right before it serves the last job at station $j_m$.

Let $\tau$ be the time under policy $\pi$ that the server starts working on a job at station $j_m$ with service time $S_j$ right before it moves to station 1. Then after completing serving this job, the server immediately switches to station 1 to serve a job there with service time $S_1$. We construct the new policy $\gamma$ as follows: $\gamma$ follows $\pi$ during $[0, \tau)$ and then serves a job at station 1 with service time $S_1'$, switches to station $j_m$ and serves a job with service time $S_j'$.

We directly couple the service times of all the jobs taken into service during $[0, \tau)$ and we cross couple $S_1$, $S_j$, $S_1'$, $S_j'$ as follows. We first generate the minimum and maximum of $S_1$ and $S_j$, namely $m$ and $M$, respectively, condition on their values, and use these values in both systems. Let $p = P(S_1 = M | m, M) = P(S_j = m | m, M)$ and $q = P(S_1 = m | m, M) = P(S_j = M | m, M) = 1 - p$. By Lemma 13.D.1(i) of Shaked and Shanthikumar [23], $p < q$. Thus, we can let

   i. $S_1 = m$, $S_j = M$, $S_1' = M$, $S_j' = m$, with probability $p$,

   ii. $S_1 = M$, $S_j = m$, $S_1' = m$, $S_j' = M$, with probability $p$,

   iii. $S_1 = m$, $S_j = M$, $S_1' = m$, $S_j' = M$, with probability $1 - 2p$.

The coupling yields $S_1' \leq S_j$ (and $S_1 \leq S_j'$) in all three cases. See Figure 3.4 for a visual description of this coupling. In the first two cases all the arrival times to station 0 under policies $\pi$ and $\gamma$ are identical. The system states for the two policies are identical after time $\tau + m + M$. Hence, policy $\gamma$ can follow policy $\pi$ thereafter. By directly coupling the service times of all jobs taken into service after time $\tau + m + M$, we find that $D_0^\gamma(t) = D_0^\pi(t)$ for all $t \geq 0$.

In the third case, let $S_0$ be the service time at station 0 for the job entering station 0 at time $\tau + S_j$ under policy $\pi$ and for the job entering station 0 at time $\tau + S_1'$ under policy $\gamma$. Consider the following two sub-cases.

Figure 3.4: Sample path couplings for the parallel system.

1. $0 < S_0 \leq m$. We directly couple the service times of all jobs taken into service after $\tau + m + M$. The system states for the two policies will become identical after time $\tau + m + M$. Hence, policy $\gamma$ can follow policy $\pi$ after $\tau + m + M$.

2. $S_0 > m$. We directly couple the service times of all jobs taken into service after $\tau + m + M$. The system states for the two policies will become identical after time $\tau + M + S_0$. However, policy $\gamma$ can follow $\pi$ after $\tau + m + M$ because same jobs or more will be available to policy $\gamma$ compared to policy $\pi$.

In both sub-cases, the problem can be analyzed in three intervals as below.

$$
D_0^\gamma(t) = \begin{cases} D_0^\pi(t), & 0 \leq t < \tau + m + S_0, \\ D_0^\pi(t) + 1, & \tau + m + S_0 \leq t < \tau + M + S_0, \\ D_0^\pi(t), & t \geq \tau + M + S_0. \end{cases}
$$

Therefore,

$$
D_0^\gamma(t) - D_0^\pi(t) \geq 0, \text{ for all } t \geq 0,
$$

and

$$
TH_0^\gamma - TH_0^\pi = \liminf_{t \to \infty} \frac{1}{t}[D_0^\gamma(t) - D_0^\pi(t)] \geq 0.
$$

31

We have shown that jobs at station 1 should be served ahead of the jobs at station $j_m$ if jobs are available at both stations. Following the same argument iteratively, we can show that a job at station 1 should be served ahead of the sequence of stations $\{j_1, \ldots, j_{m-1}\}$. In other words, if $X_1 \leq_{lr} X_j, j \neq 1$, then jobs in station 1 should be prioritized in order to maximize the long-run average throughput of the system. $\qquad\square$

**Remarks.** Theorem 3.3.2 says that we should give the highest priority to the service station which has the shortest service times in likelihood ratio ordering among all service stations. This follows the same intuition as in the series model where the optimal policy pushes jobs towards the automated station as early as possible. Note that there is no ordering condition required on the service times for the series system as opposed to the parallel system. In Section 3.2.2, under preemptive policies and Markovian assumption, we provide a complete characterization of the optimal policy. Under non-preemptive policies, even though Theorem 3.3.2 only gives a partial characterization of the optimal policy for $K > 2$, we believe that a complete characterization should follow a similar intuition as we state in the following conjecture.

**Conjecture 3.3.1.** *If the service times at the service stations follow a likelihood ratio ordering, such as $X_1 \leq_{lr} X_2 \leq_{lr} \ldots \leq_{lr} X_K$, then there exists an optimal policy which maximizes the long-run average throughput and gives priority to the non-empty station with the smallest index at any decision epoch.*

A numerical study is conducted in Section 3.3.4 to support Conjecture 3.3.1.

### 3.3.3 Two-Branch Three-Station System

In this subsection, we study a system with two branches and three service stations as shown in Figure 3.5. Stations 2 and 3 are in series and parallel to station 1 as a whole. A job after being served at station 0 will join station 1 with probability $p_1$ or station 2 with probability $p_2$. This queueing system is motivated by the nurse staffing problem studied in Yankovic and Green [35]. We consider this closed queueing system as a hospital ward, where a patient seeks admission (station 3), then stays at a bed (station 0), then requests nursing service (station 1) and returns back to his/her bed after receiving service (this process may repeat for several times), and at last seeks discharge (station 2). We assume that a new patient

comes to the ward immediately after a patient is discharged. We define $X_i$ to be the random variable denoting the i.i.d. service time at station $i$ for $i = 1, 2, 3$. We provide a partial



Figure 3.5: A closed queueing system with an automated station and three service stations.

characterization of the optimal policy that maximizes the long-run average throughput of the system in Theorem 3.3.3.

**Theorem 3.3.3.** *If $X_1[X_3] \leq_{lr} X_3[X_1]$, then there exists an optimal policy that gives the highest priority to station 1[3].*

Theorem 3.3.3 partially characterizes the optimal policy: the faster station between the two stations which directly connect to the entry of the automated station should be prioritized, if the service times at these two stations are ordered according to likelihood ratio ordering.

***Proof of Theorem 3.3.3.*** We show that if $X_1 \leq_{lr} X_3$, then there exists an optimal policy that gives the highest priority to station 1. The proof that an optimal policy gives the highest priority to station 3 if $X_3 \leq_{lr} X_1$ is similar.

Suppose policy $\pi$ is a policy under which there exists at least one decision epoch where the highest priority is not given to station 1. Specifically, let $\varepsilon$ be the first time policy $\pi$ does not give priority to station 1 even if there is a job in station 1. Suppose $\{j_1, \ldots, j_m\}$ $(j_1, \ldots, j_m \in \{2, 3\})$ gives the sequence of stations that the server visits after time $\varepsilon$ before it visits station 1 under policy $\pi$. We will next construct a new policy $\gamma$ which serves station 1 right before it serves the last job at station $j_m$.

Case 1. $(j_m = 2)$

Let $\tau$ be the time under policy $\pi$ that the server starts to work on a job at station 2 with service time $S_2$ right before it moves to station 1. Then after completing serving this job, the server immediately switches to station 1 to serve a job there with service time $S_1$. We construct the new policy $\gamma$ as follows: $\gamma$ follows $\pi$ during $[0, \tau)$ and then serves a job at station 1 with service time $S_1'$, switches to station 2, and serves a job with service time $S_2'$.

We directly couple the service times of all jobs taken into service during $[0, \infty)$, which yields $S_1' = S_1$ and $S_2' = S_2$. Let $S_0$ be the service time at station 0 for the job entering station 0 at time $\tau + S_2 + S_1$ under policy $\pi$ and for the job entering station 0 at time $\tau + S_1'$ under policy $\gamma$. $\gamma$ follows $\pi$ during $[\tau + S_1 + S_2, \infty)$. This is possible because same jobs or more are available to policy $\gamma$ compared to policy $\pi$. The system states for the two policies will become identical after time $\tau + S_2 + S_1 + S_0$. The problem can be analyzed in three intervals as follows:

$$
D_0^\gamma(t) = \begin{cases}
D_0^\pi(t), & 0 \le t < \tau + S_1' + S_0, \\
D_0^\pi(t) + 1, & \tau + S_1' + S_0 \le t < \tau + S_2 + S_1 + S_0, \\
D_0^\pi(t), & t \ge \tau + S_2 + S_1 + S_0.
\end{cases}
$$

Therefore,

$$
TH_0^\gamma(t) - TH_0^\pi(t) = \liminf_{t \to \infty} \frac{1}{t}[D_0^\gamma(t) - D_0^\pi(t)] \ge 0, \text{ for all } t \ge 0.
$$

Case 2. $(j_m = 3)$

Let $\tau$ be the time under policy $\pi$ that the server starts to work on a job at station 3 with service time $S_3$ before it moves to station 1. Then after completing serving this job, the server immediately switches to station 1 to serve a job there with service time $S_1$. We construct a new policy $\gamma$ as follows: $\gamma$ follows $\pi$ during $[0, \tau)$ and then serves a job at station 1 with service time $S_1'$, switches to station 3 and serves a job with service time $S_3'$.

We directly couple the service times of all jobs taken into service during $[0, \tau)$ and we cross couple $S_1, S_3, S_1', S_3'$ as in the proof of Theorem 3.3.2 by first generating the minimum and maximum of $S_1$ and $S_3$, namely $m$ and $M$, respectively, conditioning on their values and using these values in both policies. We need to consider three couplings:

i. $S_1 = m$, $S_3 = M$, $S_1' = M$, $S_3' = m$,

ii. $S_1 = M$, $S_3 = m$, $S_1' = m$, $S_3' = M$,

iii. $S_1 = m$, $S_3 = M$, $S_1' = m$, $S_3' = M$,

all of which yield $S_1' \leq S_3$ (and $S_1 \leq S_3'$). In the first two cases all arrival times to station 0 for policies $\pi$ and $\gamma$ are identical. The system states for the two policies are identical after time $\tau + m + M$. Hence, policy $\gamma$ can follow policy $\pi$ thereafter. By directly coupling the service times of all jobs taken into service after $\tau + m + M$, we find that the throughput of the system is the same for policy $\pi$ and $\gamma$.

In the third case, let $S_0$ be the service time at station 0 for the job entering station 0 at time $\tau + S_3$ under policy $\pi$ and for the job entering station 0 at time $\tau + S_1'$ under policy $\gamma$. Consider the following two sub-cases:

1. $0 < S_0 \leq m$. We directly couple the service times of all jobs taken into service after $\tau + m + M$. The system states for the two policies will become identical after time $\tau + m + M$. Hence, policy $\gamma$ can follow policy $\pi$ after $\tau + m + M$.

2. $S_0 > m$. We directly couple the service times of all jobs taken into service after $\tau + m + M$. The system states for the two policies will become identical after time $\tau + M + S_0$. However, policy $\gamma$ can follow $\pi$ after $\tau + m + M$ because same jobs or more are available to policy $\gamma$ compared to policy $\pi$.

In both sub-cases, the problem can be analyzed in three intervals as below.

$$
D_0^\gamma(t) = \begin{cases} D_0^\pi(t), & 0 \leq t < \tau + m + S_0, \\ D_0^\pi(t) + 1, & \tau + m + S_0 \leq t < \tau + M + S_0, \\ D_0^\pi(t), & t \geq \tau + M + S_0. \end{cases}
$$

Therefore,
$$
TH_0^\gamma - TH_0^\pi = \liminf_{t \to \infty} \frac{1}{t}[D_0^\gamma(t) - D_0^\pi(t)] \geq 0, \text{ for all } t \geq 0.
$$

We have shown that a job at station 1 should be served ahead of a job at station 3 if jobs are available at both stations.

If we follow the same argument iteratively for the two cases, we can show that a job at station 1 should be served ahead of the sequence of jobs $\{j_1, \ldots, j_{m-1}\}$. In other words, the optimal policy which maximizes the throughput of the system gives the highest priority to station 1. $\qquad\square$

Furthermore, we provide a corollary to Theorem 3.3.3 when $X_3 \leq_{lr} X_1$.

**Corollary 3.3.1.** *Suppose that $X_3 \leq_{lr} X_1$. If $X_1 \leq_{lr} X_2 + X_3$, then there exists an optimal policy that gives priority to station 1 over station 2. If $X_1 \geq_{lr} X_2 + X_3$, then there exists an optimal policy which gives priority to station 2 over station 1.*

*Proof.* Since $X_3 \leq_{lr} X_1$, Theorem 3.3.3 tells that we need to only consider policies that give the highest priority to station 3 whenever a job visits that station. This implies that we need to only consider policies that serve stations 2 and 3 sequentially, i.e., we can think of stations 2 and 3 as single station with service time $X_2 + X_3$. Then, the result follows from Theorem 3.3.2. □

Corollary 3.3.1 gives conditions under which two priority policies are optimal: 1) If $X_3 \leq_{lr} X_1 \leq_{lr} X_2 + X_3$, then the priority order is station 3, station 1 and station 2; 2) If $X_3 \leq_{lr} X_1$ and $X_2 + X_3 \leq_{lr} X_1$, then the priority order is station 3, station 2 and station 1.

### 3.3.4 An Heuristic Policy and Numerical Results

Finally, we propose an index policy, namely the *shortest expected remaining service time heuristic*, for the general parallel-series system under non-preemptive policies. We compare the performance of this heuristic along with the performance of the optimal policy by means of a numerical study.

Define

$$S_m^i = \sum_{j=i}^{i_m} S_{j,m}, \text{ for } i = 1, 2, \ldots, i_m \text{ and } m = 1, 2, \ldots, M,$$

which denotes the remaining service time to complete type $m$ service starting from its $i$th task. We describe the index policy as follows. Suppose at decision epoch $t$, $N(t)$ customers are in need of attention from the server. Let $(j_n, l_n)$ be the station where the $n$-th customer resides at time $t$, $n = 1, 2, \ldots, N(t)$. The heuristic policy ranks these $N(t)$ jobs in increasing order of $E[S_{j_n}^{l_n}]$ and serve the first of them. Since we consider non-preemptive policies, the server is allowed to make a decision at service completion epochs.

We conduct a numerical study for the parallel system with three service stations and the two-branch three-station system. The objective of the numerical study is to examine the performance of the heuristic policy as opposed to the optimal policy. We would like to

study many different scenarios with a wide range of system parameters. More specifically, we generate the service time of each service station from an exponential distribution. We fix the service rate of the automated station $\mu_0 = 1$ and vary the service rates of service stations in several combinations. We have considered two subsets of experiments depending on the number of jobs circuiting in the system $B$, namely, 5 and 10.

For the parallel model with three service stations, we let the service rate $\mu_i$ ($i = 1, 2, 3$) take values from the set $\{0.5, 0.75, 1, 2, 5\}$. We consider four cases for $[p_1, p_2, p_3]$, namely $[0.2, 0.4, 0.4]$, $[0.4, 0.3, 0.3]$, $[0.6, 0.2, 0.2]$ and $[0.8, 0.1, 0.1]$. There are 1,000 scenarios in total.

We use the method of policy iteration to obtain the optimal policy for each scenario. The numerical results show that the index policy is optimal for all 1,000 scenarios within the set of the non-preemptive policies. In Sections 3.2.2 and 3.3.2, we are able to prove a partial characterization of the optimal policy for the parallel system under preemptive and non-preemptive policies, respectively. In Conjecture 3.2.1 and Conjecture 3.3.1, we characterize the complete optimal policy which gives priority to the non-empty station with the largest service rate and the shortest service times in likelihood ratio ordering, respectively. The numerical results for the index policy are consistent with the partial results in Theorem 3.2.2 and Theorem 3.3.2 and support the conjectures.

For the two-branch three-station system, we let $\mu_1$ take values from the set $\{0.5, 0.75, 1, 2, 5\}$, $\mu_2$ take values from the set $\{0.5\mu_1, 1.5\mu_1, 3\mu_1\}$, and set $\mu_3 = \mu_2$. We consider three cases for $[p_1, p_2]$, namely $[0.25, 0.75]$, $[0.5, 0.5]$, and $[0.75, 0.25]$. There are 90 scenarios in total. For each scenario, we computed the percentage deviation (P.D.) of the performance of the index policy heuristic from that of the optimal policy as well as the optimal throughput ($TH^*$). These results are presented in Table 3.1. From Table 3.1, it can be seen that the index policy is not always optimal when the sum of the expected service times for the branch with two service stations $E[S_2] + E[S_3]$ is less than the expected service time of the other single-station branch $E[S_1]$. For example, consider the case where $(\mu_1, \mu_2, \mu_3) = (1, 3, 3)$, we have $E[S_2] + E[S_3] = \frac{2}{3}$ which is less than $E[S_1] = 1$. We can observe from the table that the index policy is not optimal (positive P.D.) for all scenarios when we vary $B$ and $[p_1, p_2]$. The numerical results show that the index policy is optimal for 69 scenarios out of 90 scenarios. Over all 90 scenarios, the average percentage deviation is 0.00072%, and the maximum deviation is 0.01273%. Overall, the shortest expected remaining service time

Table 3.1: Performance of the index policy for exponential service times under non-preemptive policies (in terms of the percentage deviation (P.D.) from the optimal performance).

| $(\mu_1, \mu_2, \mu_3)$ | B=5 | | | | | | B=10 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | p = (0.25,0.75) | | p = (0.5,0.5) | | p = (0.75,0.25) | | p = (0.25,0.75) | | p = (0.5,0.5) | | p = (0.75,0.25) | |
| | P.D. | TH* | P.D. | TH* | P.D. | TH* | P.D. | TH* | P.D. | TH* | P.D. | TH* |
| (0.5,0.25,0.25) | 0 | 0.1538 | 0 | 0.2 | 0 | 0.2857 | 0 | 0.1538 | 0 | 0.2 | 0 | 0.2857 |
| (0.5,0.75,0.75) | 0 | 0.4 | 0 | 0.4286 | 0 | 0.4615 | 0 | 0.4 | 0 | 0.4286 | 0 | 0.4615 |
| (0.5,1.5,1.5) | 6.97E-05 | 0.665 | 6.59E-05 | 0.5999 | 5.76E-05 | 0.5454 | 1.67E-10 | 0.667 | 1.66E-10 | 0.6 | 1.63E-10 | 0.5455 |
| (0.75,0.375,0.375) | 0 | 0.2308 | 0 | 0.3 | 0 | 0.4286 | 0 | 0.2308 | 0 | 0.3 | 0 | 0.4286 |
| (0.75,1.125,1.125) | 0 | 0.6 | 0 | 0.6428 | 0 | 0.6922 | 0 | 0.6 | 0 | 0.6429 | 0 | 0.6923 |
| (0.75,2.25,2.25) | 1.27E-04 | 0.9988 | 1.13E-04 | 0.899 | 8.98E-05 | 0.8174 | 7.86E-09 | 1 | 7.29E-09 | 0.9 | 6.80E-09 | 0.8182 |
| (1,0.5,0.5) | 0 | 0.3077 | 0 | 0.4 | 0 | 0.5714 | 0 | 0.3077 | 0 | 0.4 | 0 | 0.5714 |
| (1,1.5,1.5) | 0 | 0.7998 | 0 | 0.8568 | 0 | 0.9224 | 0 | 0.8 | 0 | 0.8571 | 0 | 0.9231 |
| (1,3,3) | 4.39E-05 | 1.3267 | 3.74E-05 | 1.195 | 3.33E-05 | 1.0871 | 9.89E-08 | 1.3333 | 8.64E-08 | 1.2 | 7.65E-08 | 1.0909 |
| (2,1,1) | 0 | 0.6154 | 0 | 0.7998 | 0 | 1.1406 | 0 | 0.6154 | 0 | 0.8 | 0 | 1.1429 |
| (2,3,3) | 0 | 1.5856 | 0 | 1.6915 | 0 | 1.8077 | 0 | 1.6 | 0 | 1.7143 | 0 | 1.8462 |
| (2,6,6) | 0 | 2.4844 | 0 | 2.2695 | 0 | 2.0853 | 5.76E-06 | 2.665 | 4.91E-06 | 2.3999 | 3.26E-06 | 2.1817 |
| (5,2.5,2.5) | 0 | 1.5271 | 0 | 1.9506 | 0 | 2.6028 | 0 | 1.5385 | 0 | 2 | 0 | 2.8571 |
| (5,7.5,7.5) | 0 | 3.2723 | 0 | 3.3787 | 0 | 3.4809 | 0 | 3.9928 | 0 | 4.2714 | 0 | 4.5847 |
| (5,15,15) | 0 | 4.0073 | 0 | 3.8606 | 0 | 3.7166 | 0 | 6.32 | 0 | 5.7866 | 0 | 5.317 |

heuristic is either optimal or near-optimal.

It is important to point out that for the two-branch three-station system, the likelihood-ratio ordering on total remaining service times does not hold within this numerical study with exponential service times. For example, again consider the case where $(\mu_1, \mu_2, \mu_3) = (1, 3, 3)$, even though $E[S_2] + E[S_3] < E[S_1]$, $S_2 + S_3 <_{lr} S_1$ does not hold. Hence, our numerical study does not rule out the optimality of an index policy when the total remaining service times can be ordered according to likelihood-ratio ordering.

## 3.4 Continuous Resource

In this section, we consider the parallel system shown in Figure 3.3 with a continuous service resource. Suppose that the total available service rate is fixed and denoted by $\mu$. All service stations share this fixed resource and satisfy the constraint $\sum_{k=1}^{K} \mu_k = \mu$, where $\mu_k$ is the service capacity allocated to station $k$. We assume that the amount of intrinsic work required at all stations are i.i.d. exponentials. Hence, when service capacity $\mu_k$ is allocated to station $k$, service times at station $k$ will be exponentially distributed with mean $1/\mu_k$ $(k = 1, 2, \ldots, K)$. We only consider preemptive policies. We are interested in dynamically allocating the total service rate to the service stations (i.e., determine $\mu_k$'s) over time in

order to maximize the long-run average throughput of the system.

We formulate this problem as a Markov decision process. We use the same notation defined in Section 3.2.1 unless otherwise stated. Because $p_k > 0$ for $k = 1, 2, \ldots, K$, again the transition matrix of the system consists of a single recurrent class for every deterministic stationary policy. Hence, the MDP is recurrent and $g$ exists. Define $\Lambda = n\mu_0 + \mu$ as the uniformization constant. Without loss of generality, we assume that $\Lambda = 1$. Then, the optimality equation can be expressed as follows. For $1 \leq n \leq B$,

$$g + V(\boldsymbol{n}) = (B - n)\mu_0 \sum_{i=1}^{K} p_i V(\boldsymbol{n} + \boldsymbol{e}^i) + n\mu_0 V(\boldsymbol{n}) + f(\boldsymbol{n}),$$

where

$$f(\boldsymbol{n}) = \begin{cases} \mu V(\boldsymbol{0}) & \text{if } \boldsymbol{n} = \boldsymbol{0}, \\ \max_{\boldsymbol{\mu} \in \mathcal{M}(\boldsymbol{n})} \{ \sum_{i=1}^{K} \mu_i (1 + V(\boldsymbol{n} - \boldsymbol{e}^i)) \} & \text{otherwise}, \end{cases}$$

$\boldsymbol{\mu}$ is the allocation vector, and $\mathcal{M}(\boldsymbol{n}) = \{(\mu_1, \mu_2, \ldots, \mu_K) : \sum_{i=1}^{K} \mu_i = \mu, \mu_i \geq 0 \text{ for } i \in \boldsymbol{I_n} \text{ and } \mu_i = 0 \text{ for } i \notin \boldsymbol{I_n}\}$. We characterize the optimal policy in Theorem 3.4.1.

**Theorem 3.4.1.** *Any non-idling policy maximizes the long-run average throughput of the system.*

***Proof of Theorem 3.4.1.*** Following the same argument as in the proof of Theorem 3.2.1, we first show that the result holds for a finite horizon problem with $m$ periods for all $m \geq 0$. We still use $V_m(\pi, \boldsymbol{n})$ to denote the $m$-period expected throughput under policy $\pi$ when the initial state is $\boldsymbol{n}$.

We will show that, for all $m \geq 0$,

$$V_m(\boldsymbol{n} - \boldsymbol{e}^i) - V_m(\boldsymbol{n} - \boldsymbol{e}^j) = 0, \tag{3.4.1}$$

where $i, j \in \boldsymbol{I_n}$ and $i \neq j$. We assume that $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$.

We will use induction on $m$. Because $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$, then (3.4.1) automatically holds at period 0. Assume that (3.4.1) holds at period $m$. We will show that it also holds

at period $m + 1$. For $2 \leq n \leq B$, we have

$$V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^i) - V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j)$$

$$=(B - n + 1)\mu_0 \sum_{k=1}^{K} p_k[V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^k) - V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k)]$$

$$+ (n - 1)\mu_0[V_m(\boldsymbol{n} - \boldsymbol{e}^i) - V_m(\boldsymbol{n} - \boldsymbol{e}^j)]$$

$$+ \mu[V_m(\boldsymbol{n} - \boldsymbol{e}^i - \boldsymbol{e}^j) - V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^i)]$$

The right-hand side of the equation is equal to 0 by the inductive hypothesis for (3.4.1) at period $m$. $\qquad\qquad\square$

In Theorem 3.4.1, a non-idling policy means that the entire service capacity is allocated when there is at least one job at service stations. The intuition behind the theorem is that any non-idling policy utilizes all the available resource whenever there are job(s) at service stations. The departure rate from service stations to the automated station is always maximized for all non-idling policies.

### 3.5    Conclusion

In this chapter, we study the optimization problems of dynamically allocating a finite amount of service capacity among service stations for a finite-population service system. We focus on assigning a single flexible server, and briefly study a continuous resource problem in a parallel system. We consider a parallel-series system where the service stations are connected in the parallel-series structure. The objective is to maximize the long-run average reward of the system. We show that for the parallel-series system, maximizing the long-run average reward is equivalent to maximizing the long-run average throughput.

For the discrete resource problems, we investigate the optimal assignment policy for a single server in series, parallel and a two-branch three-station models. For the series model, we show that the optimal policy is a sequential policy under both preemptive and non-preemptive policies. For the parallel model, the optimal policy gives the highest priority to the fastest station. For the two-branch three-station model, we provide a partial characterization of the optimal policy under non-preemptive policies. For each model, the optimal policy tries to push jobs back to the automated station as early as possible so that these

jobs will leave the automated station early and hence reduce the idling time of the server, which would increase the system throughput.

We also propose an index policy which gives priority to the non-empty station with the shortest expected remaining service time. The numerical results for the parallel model with three service stations show that the index policy is optimal for all the scenarios in our numerical study. This supports our conjecture on the complete characterization of the optimal policy for the parallel model. The numerical results for the two-branch three-station model show that the index policy gives a small percentage deviation on the performance from the optimal policy, either on average or under the worst case.

For the continuous resource problem, we consider a parallel system. We show that any non-idling policy maximizes the long-run average throughput of the system under Markovian case.

# Chapter 4

# OPTIMAL STATIC DESIGN OF FINITE-POPULATION QUEUEING SYSTEMS: THROUGHPUT MAXIMIZATION

In this chapter, we study a static design problem for the closed queueing system under study, where allocation decisions are made before the system starts to operate. When the system is under operation, we are not allowed to adjust the allocation of resources. A finite amount of service resource $U$ can be used at any of the $K$ service stations in the second stage. Suppose that $U$ is fixed and continuously divisible. Let $\boldsymbol{u} \equiv [u_1, u_2, \ldots, u_K]$ be the allocation vector in which its $k$-th element $u_k$ is a decision variable denoting the units of service resource allocated to station $k$ for $k = 1, \ldots, K$, where $\sum_{k=1}^{K} u_k = U$. Let $\boldsymbol{U}$ denote the set of allocation vectors, i.e., $\boldsymbol{U} = \{\boldsymbol{u} : \sum_{k=1}^{K} u_k = U\}$. Service times at station $k$ is exponentially distributed with mean $1/\mu_k(u_k)$, where $\mu_k(u_k)$ is a function of $u_k$. We assume that $\mu_k(u_k)$ is strictly increasing in $u_k$ and continuous on $[0, \infty)$, and $\mu_k(0) = 0$, for $k = 1, \ldots, K$. We want to characterize the optimal allocation $\boldsymbol{u}$ among the $K$ service stations in order to maximize the long-run average throughput of the system.

## 4.1 Model Formulation

We model the finite-population service system by a closed queueing network with $K + 1$ $(2 \leq K < \infty)$ stations and $B$ $(1 \leq B < \infty)$ customers as shown in Figure 1.1. Assume that service times at station 0 are i.i.d. exponential random variables with mean $1/\mu_0$. The second stage consists of the remaining $K$ service stations, where each of the $K$ service stations is served by a single dedicated server.

We next introduce additional notation and recapitulate some known and relevant results for closed Jackson networks. Let $X_i(t)$ be the number of customers at station $i$ at time $t$ for $i = 0, 1, \ldots, K$ and $t \geq 0$. Then, the state of the system at time $t$ is $X(t) = [X_0(t), X_1(t), \ldots, X_K(t)]$, and $\{X(t), t \geq 0\}$ is a continuous-time Markov Chain (CTMC)

representation of the closed Jackson network. The state space is $S = \{ \boldsymbol{n} = (n_0, n_1, \ldots, n_K) : \sum_{i=0}^{K} n_i = B \text{ and } n_i = 0, 1, \ldots, B \text{ for } i = 1, 2, \ldots, K \}$. Let $r_{ij}$ denote the customer routing probability from station $i$ to station $j$. Suppose that the routing probability matrix $R = [r_{ij}]$ is irreducible, so that $\{ X(t), t \geq 0 \}$ is an irreducible and positive recurrent CTMC with limiting distribution $P(n_0, n_1, \ldots, n_K) = \lim_{t \to \infty} P[X_0(t) = n_0, X_1(t) = n_1, \ldots, X_K(t) = n_K]$. Define $v_i$ as the visiting ratio to station $i$ for $i = 0, 1, \ldots, K$, satisfying the following equations:

$$\sum_{i=0}^{K} v_i r_{ij} = v_j, \quad j = 0, 1, \ldots, K.$$

We know that the stationary probability can be expressed as follows:

$$P(n_0, n_1, \ldots, n_K) = \frac{1}{C(B)} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} \left( \frac{\mu_k(u_k)}{v_k} \right)^{-n_k},$$

where $C(B)$ is the normalizing constant so that $\sum_{\boldsymbol{n} \in S} P(\boldsymbol{n}) = 1$. For more on closed queueing networks, see, e.g., Gross and Harris [6]. Hence, we have

$$
\begin{aligned}
C(B) &= \sum_{\boldsymbol{n} \in S} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} \left( \frac{\mu_k(u_k)}{v_k} \right)^{-n_k} \\
&= \sum_{n_0=0}^{B} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} \sum_{n_1 + \ldots + n_K = B - n_0} \prod_{k=1}^{K} \left( \frac{\mu_k(u_k)}{v_k} \right)^{-n_k} \quad (4.1.1) \\
&= \sum_{n_0=0}^{B} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} G(B - n_0),
\end{aligned}
$$

where $G(\cdot)$ is defined as

$$G(n) = \sum_{n_1 + \ldots + n_K = n} \prod_{k=1}^{K} \left( \frac{\mu_k(u_k)}{v_k} \right)^{-n_k}.$$

Let $TH(B, \boldsymbol{u})$ denote the long-run average throughput of the system, which is a function of $B$ and $\boldsymbol{u}$. $TH(B, \boldsymbol{u})$ is defined as the long-run average throughput at station 0, and can be computed as follows:

$$
\begin{aligned}
TH(B, \boldsymbol{u}) &= \mu_0 \sum_{\boldsymbol{n} \in S} n_0 P(n_0, n_1, \ldots, n_K) \\
&= \mu_0 \frac{\sum_{\boldsymbol{n} \in S} \frac{n_0 (\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} \left( \frac{\mu_k(u_k)}{v_k} \right)^{-n_k}}{\sum_{\boldsymbol{n} \in S} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} \left( \frac{\mu_k(u_k)}{v_k} \right)^{-n_k}}. \quad (4.1.2)
\end{aligned}
$$

43

With some algebra, we can show that

$$TH(B, \boldsymbol{u}) = v_0 \frac{C(B-1)}{C(B)}.$$

In the following discussion, we use $TH(B)$ to denote $TH(B, \boldsymbol{u})$ whenever it does not cause ambiguity. Let $TH^*(B)$ denote the maximum long-run average throughput of the system when $B$ customers circulate in the system, i.e.,

$$TH^*(B) \equiv \max_{\boldsymbol{u} \in U} \quad TH(B, \boldsymbol{u}).$$

In Lemma 4.1.1, we show that $TH(B, \boldsymbol{u})$ is non-decreasing in $B$.

**Lemma 4.1.1.** *$TH(B, \boldsymbol{u})$ is non-decreasing in $B$, i.e., for any positive integer $B$, we have*

$$\frac{C(B-1)}{C(B)} \geq \frac{C(B-2)}{C(B-1)} \geq \frac{C(B-3)}{C(B-2)} \geq \cdots \geq \frac{C(0)}{C(1)}.$$

***Proof of Lemma 4.1.1.*** The expected number of customers at station $0$ is equal to $\frac{v_0}{\mu_0} \frac{C(B-1)}{C(B)}$, while the probability of having positive number of jobs at station $i$ is that $P(n_i \geq 1) = \frac{v_i}{\mu_i} \frac{C(B-1)}{C(B)}$ for $i = 1, 2, \ldots, K$. It suffices to show that as $B$ increases, the number of customers at station $i$ increases stochastically. This proof follows the proof of Lemma 1 in Yao [36]. □

Next, we show that the long-run average throughput of the system is bounded above, and provide a set of equations to calculate the upper bound. Let $\boldsymbol{u}^{**} = [u_1^{**}, u_2^{**}, \ldots, u_K^{**}]$ be the solution to the following equations:

$$\frac{v_0}{v_1} \mu_1(u_1^{**}) = \frac{v_0}{v_2} \mu_2(u_2^{**}) = \ldots = \frac{v_0}{v_K} \mu_K(u_K^{**}). \tag{4.1.3}$$

This system of equations has a unique solution if $\mu_k(\cdot)$ $(k = 0, 1, \ldots, K)$ is strictly increasing.

**Theorem 4.1.1.** *The long-run average throughput of the system is bounded above by $\frac{v_0}{v_k} \mu_k(u_k^{**})$ where $u_k^{**}$ satisfies equations (4.1.3).*

**_Proof of Theorem 4.1.1_**. The throughput of the system is given by

$$
\begin{aligned}
TH(B) =& \mu_0 \sum_{n_0+n_1+\ldots+n_K=B} n_0 P(n_0, n_1, \ldots, n_K) \\
=& \mu_0 \frac{\sum_{n_0+n_1+\ldots+n_K=B} \frac{n_0(\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} (\frac{\mu_k(u_k)}{v_k})^{-n_k}}{\sum_{n_0+n_1+\ldots+n_K=B} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} (\frac{\mu_k(u_k)}{v_k})^{-n_k}} \\
=& v_0 \frac{\sum_{n_0+n_1+\ldots+n_K=B, n_0>0} \frac{(\mu_0/v_0)^{-(n_0-1)}}{(n_0-1)!} \prod_{k=1}^{K} (\frac{\mu_k(u_k)}{v_k})^{-n_k}}{\sum_{n_0+n_1+\ldots+n_K=B} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} (\frac{\mu_k(u_k)}{v_k})^{-n_k}} \\
=& v_0 \frac{\sum_{n_0+n_1+\ldots+n_K=B-1} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} (\frac{\mu_k(u_k)}{v_k})^{-n_k}}{\sum_{n_0+n_1+\ldots+n_K=B} \frac{(\mu_0/v_0)^{-n_0}}{n_0!} \prod_{k=1}^{K} (\frac{\mu_k(u_k)}{v_k})^{-n_k}} \\
=& v_0 \frac{C(B-1)}{C(B)},
\end{aligned}
$$

and the long-run average throughput at station $k$ for $k = 1, 2, \ldots, K$ is known as

$$
TH_k(B) = v_k \frac{C(B-1)}{C(B)}.
$$

Then we have

$$
TH(B) = \frac{v_0}{v_k} TH_k(B).
$$

$TH_k(B)$ is increasing in $B$ and bounded by $\mu_k(u_k)$ for $k = 1, 2, \ldots, K$. Then, the throughput of the system is bounded by

$$
TH(B) \leq \min\{\frac{v_0}{v_1}\mu_1(u_1), \frac{v_0}{v_2}\mu_2(u_2), \ldots, \frac{v_0}{v_K}\mu_K(u_K)\}.,
$$

Thus, the maximum throughput of the system bounded by

$$
TH^*(B) \leq \max_{\boldsymbol{u} \in \boldsymbol{U}} \min\{\frac{v_0}{v_1}\mu_1(u_1), \frac{v_0}{v_2}\mu_2(u_2), \ldots, \frac{v_0}{v_K}\mu_K(u_K)\}.
$$

This max-min problem can be written as

$$
\begin{aligned}
\max_{u_1,\ldots,u_K,z} \quad & z \\
s.t. \quad & z - f_k(u_k) \leq 0, k = 1, 2, \ldots, K, \\
& \sum_{k=1}^{K} u_k \leq U, \\
& u_k \geq 0, k = 1, 2, \ldots, K,
\end{aligned}
$$

where $f_k(u_k) = \frac{v_0}{v_k}\mu_k(u_k)$ for $k = 1, 2, \ldots, K$. We first show that the given solution of $u_k^{**}$s satisfy the Karush-Kuhn-Tucker condition. Under the given solution, the first and second set of constraints are binding.

$$1 - \sum_{k=1}^{K} \lambda_k \leq 0$$

$$\lambda_k f_k'(u_k) - \lambda_U \leq 0, k = 1, \ldots, K$$

$$\lambda_k(z - f_k(u_k)) = 0, k = 1, \ldots, K$$

$$\lambda_U(U - \sum_{k=1}^{K} u_k) = 0$$

$$[1 - \sum_{k=1}^{K} \lambda_k]z = 0$$

$$[\lambda_k f_k'(u_k) - \lambda_U]u_k = 0, k = 1, \ldots, K$$

$$\lambda_k \geq 0, k = 1, \ldots, K$$

$$\lambda_U \geq 0$$

These exist $\lambda_k$s and $\lambda_U$ satisfying

$$1 - \sum_{k=1}^{K} \lambda_k = 0$$

$$\lambda_k f_k'(u_k) - \lambda_U = 0, k = 1, \ldots, K$$

If not all $f_k(u_k)$s are equal, i.e., there is at least one constraint in the first or second constraint sets is not binding, then there do not exist $\lambda_k$s and $\lambda_U$ satisfying the Karush-Kuhn-Tucker condition.

Since the function $\mu_k(u_k)$ $(k = 1, \ldots, K)$ is monotonously increasing, the given solution of $u_k^{**}$s is the unique solution to the equations and therefore it is the optimal solution to the max-min problem. □

## 4.2  Characterization of Optimal Capacity Allocation

We study the optimal allocation problem in three cases depending on the number of customers circulating in the system. First, we consider a simple case when a single customer circulates in the system ($B = 1$). Secondly, we study the problem when the number of cus-

tomers in the system approaches infinity $(B \to \infty)$. Finally, we allow $B$ to be any positive integer number.

### 4.2.1  A Single Job

First, we study the problem when only one customer circulates in the system, i.e., $B = 1$. Let $\mu_i'(u_i)$ denote the first-order derivative of $\mu_i(u_i)$ over $u_i$ for $k = 1, 2, \ldots, K$. The main result is given in Theorem 4.2.1

**Theorem 4.2.1.** *If $\mu_k(\cdot)$ is increasing and concave for $k = 1, 2, \ldots, K$, then the optimal allocation vector $\boldsymbol{u}^* = [u_1^*, u_2^*, \ldots, u_K^*]$ satisfies the following equations:*

$$\frac{v_1 \mu_1'(u_1^*)}{\mu_1(u_1^*)^2} = \frac{v_2 \mu_2'(u_2^*)}{\mu_2(u_2^*)^2} = \ldots = \frac{v_K \mu_K'(u_K^*)}{\mu_K(u_K^*)^2}. \tag{4.2.1}$$

**Proof of Theorem 4.2.1.** Following equations (4.1.1) and (4.1.2), we have

$$C(1) = \frac{v_0}{\mu_0} + \frac{v_1}{\mu_1(u_1)} + \frac{v_2}{\mu_2(u_2)} + \ldots + \frac{v_K}{\mu_K(u_K)},$$

and

$$TH(1, \boldsymbol{u}) = \frac{v_0}{\frac{v_0}{\mu_0} + \frac{v_1}{\mu_1(u_1)} + \frac{v_2}{\mu_2(u_2)} + \ldots + \frac{v_K}{\mu_K(u_K)}}. \tag{4.2.2}$$

Notice that maximizing (4.2.2) is equivalent to the following problem:

$$\min_{u_1, \ldots, u_K} \quad \frac{v_1}{\mu_1(u_1)} + \frac{v_2}{\mu_2(u_2)} + \ldots + \frac{v_K}{\mu_K(u_K)}$$

$$s.t. \quad \sum_{k=1}^{K} u_k = U.$$

By applying the Lagrange Multiplier method, the optimization problem reduces to:

$$\min_{\lambda, u_1, \ldots, u_K} \quad Z = \frac{v_1}{\mu_1(u_1)} + \frac{v_2}{\mu_2(u_2)} + \ldots + \frac{v_K}{\mu_K(u_K)} + \lambda(\sum_{k=1}^{K} u_k - U), \tag{4.2.3}$$

where $\lambda$ is the Lagrange multiplier. The second-order partial derivatives of function $Z$ is calculated by

$$\frac{\partial^2 Z}{\partial u_i^2} = 2 v_i \mu_i(u_i)^{-3} [\mu_i'(u_i)]^2 - v_i \mu_i(u_i)^{-2} \mu_i''(u_i), \quad \frac{\partial^2 Z}{\partial u_i \partial u_j} = 0 (i \neq j), \quad \frac{\partial^2 Z}{\partial u_i \partial \lambda} = 1, \quad \frac{\partial^2 Z}{\partial \lambda^2} = 0.$$

If $\mu_i''(u_i) \leq 0$ or $\mu_i(u_i)$ is increasing and concave for $i = 1, 2, \ldots, K$, we know that the Hessian matrix of $Z$ is positive definite, and then the minimizer of the program is given by

$$\frac{v_1 \mu_1'(u_1)}{\mu_1(u_1)^2} = \frac{v_2 \mu_2'(u_2)}{\mu_2(u_2)^2} = \ldots = \frac{v_K \mu_K'(u_K)}{\mu_K(u_K)^2}, \quad \sum_{k=1}^{K} u_k = U.$$

$\square$

Theorem 4.2.1 yields a closed-form solution for $\boldsymbol{u^*}$ when $\mu_k(\cdot)$ $(k = 1, 2, \ldots, K)$ is linear and strictly increasing, as shown in Corollary 4.2.1.

**Corollary 4.2.1.** *If $\mu_k(u_k) = a_k u_k$, where $a_k > 0$ for $k = 1, \ldots, K$, then the optimal allocation vector $\boldsymbol{u^*} = [u_1^*, u_2^*, \ldots, u_K^*]$ is given by*

$$u_k^* = \frac{\left(\frac{v_k}{a_k}\right)^{\frac{1}{2}}}{\sum_{i=1}^{K}\left(\frac{v_i}{a_i}\right)^{\frac{1}{2}}} U, \quad k = 1, 2, \ldots, K.$$

*__Proof of Corollary 4.2.1.__* Following (4.2.3), we have the unconstrained program:

$$\min_{\lambda, u_1, \ldots, u_K} \quad Z = \frac{v_1}{a_1 u_1} + \frac{v_2}{a_2 u_2} + \ldots + \frac{v_K}{a_K u_K} + \lambda\left(\sum_{k=1}^{K} u_k - U\right)$$

The second-order partial derivatives of function $Z$ is calculated by

$$\frac{\partial^2 Z}{\partial u_i^2} = \frac{2v_i}{a_i u_i^3}, \quad \frac{\partial^2 Z}{\partial u_i \partial u_j} = 0 (i \neq j), \quad \frac{\partial^2 Z}{\partial u_i \partial \lambda} = 1, \quad \frac{\partial^2 Z}{\partial \lambda^2} = 0.$$

We see that the Hessian matrix of $Z$ is positive definite. Hence, by solving the first-order derivative equations, we obtain the minimum of the program:

$$u_k^* = \frac{\left(\frac{v_k}{a_k}\right)^{\frac{1}{2}}}{\sum_{i=1}^{K}\left(\frac{v_i}{a_i}\right)^{\frac{1}{2}}} U.$$

$\square$

**Remark.** The optimal allocation $\boldsymbol{u^*} = [u_1^*, u_2^*, \ldots, u_K^*]$ is proportional to the square-root of the ratio of visiting ratio $v_i$ and linear coefficient $a_i$. The optimal allocation to a station increases as its visiting ratio increases while keeping others intact. Similarly, the optimal allocation to a station increases as its linear coefficient decreases while keeping others intact.

48

**Example 4.2.1.** Consider a series system as shown in Figure 3.2. A customer after being served at station 0 receives service at station $1, 2, \ldots, K$, in the given order, and then returns to station 0 after being served at station $K$. Notice that all $v_k$s, the visiting ratio, are equal for this series system. Then the optimal allocation is given by the solution to following equations:

$$\frac{\mu_1'(u_1)}{\mu_1(u_1)^2} = \frac{\mu_2'(u_2)}{\mu_2(u_2)^2} = \ldots = \frac{\mu_K'(u_K)}{\mu_K(u_K)^2}, \quad \sum_{k=1}^{K} u_k = U.$$

If $\mu_k(u_k) = a_k u_k$, where $a_k > 0$ for $k = 1, \ldots, K$, then the optimal allocation reduces to

$$u_k^* = \frac{(\frac{1}{a_k})^{\frac{1}{2}}}{\sum_{i=1}^{K} (\frac{1}{a_i})^{\frac{1}{2}}} U.$$

**Example 4.2.2.** Consider a parallel system as shown in Figure 3.3. A customer after being served at station 0 joins station $k$ ($k = 1, 2, \ldots, K$) with probability $p_k > 0$ ($\sum_{k=1}^{K} p_k = 1$). After a customer finishes its service at station $k$ ($k = 1, 2, \ldots, K$), it returns to station 0. Notice that $v_k = p_k$ for $k = 1, 2, \ldots, K$ if we set $v_0 = 1$. Then the optimal allocation is given by the solution to the following equations:

$$\frac{p_1 \mu_1'(u_1)}{\mu_1(u_1)^2} = \frac{p_2 \mu_2'(u_2)}{\mu_2(u_2)^2} = \ldots = \frac{p_K \mu_K'(u_K)}{\mu_K(u_K)^2}, \quad \sum_{k=1}^{K} u_k = U.$$

If $\mu_k(u_k) = a_k u_k$, where $a_k > 0$ for $k = 1, \ldots, K$, then the optimal allocation is given by

$$u_k^* = \frac{(\frac{p_k}{a_k})^{\frac{1}{2}}}{\sum_{i=1}^{K} (\frac{p_i}{a_i})^{\frac{1}{2}}} U.$$

*4.2.2 When Population Size Approaches Infinity*

Next, we consider characterizing the optimal allocation when the number of customers circulating in the system approaches infinity.

**Theorem 4.2.2.** *As the number of customers ($B$) increases to infinity, the optimal allocation vector approaches $\boldsymbol{u}_k^{**}$ which satisfies (4.1.3), and the optimal long-run average throughput of the system increases to $\frac{v_0}{v_k} \mu_k(u_k^{**})$. Moreover, there exists a unique solution $\boldsymbol{u}_k^{**}$ if $\mu_k(u_k)$ ($k = 1, 2, \ldots, K$) is strictly increasing.*

**Proof of Theorem 4.2.2.** Define $\rho = \frac{v_0}{v_1}\mu_1(u_1^{**}) = \ldots = \frac{v_0}{v_K}\mu_K(u_K^{**})$. Then, the long-run average throughput of the system for the given $\boldsymbol{u}^{**}$ which satisfies (4.1.3) is given by

$$
\begin{aligned}
TH(B, \boldsymbol{u}^{**}) =&\, \mu_0 \frac{\sum_{n_0+n_1+\ldots+n_K=B} \frac{n_0 \mu_0^{-n_0}}{n_0!} \prod_{k=1}^K \rho^{-n_k}}{\sum_{n_0+n_1+\ldots+n_K=B} \frac{\mu_0^{-n_0}}{n_0!} \prod_{k=1}^K \rho^{-n_k}} \\
=&\, \mu_0 \frac{\sum_{n_0=1}^{B} \sum_{n_1+\ldots+n_K=B-n_0} \frac{\mu_0^{-n_0}}{(n_0-1)!} \rho^{-(n_1+\ldots+n_K)}}{\sum_{n_0=0}^{B} \sum_{n_1+\ldots+n_K=B-n_0} \frac{\mu_0^{-n_0}}{n_0!} \rho^{-(n_1+\ldots+n_K)}} \\
=&\, \mu_0 \frac{\sum_{n_0=0}^{B-1} \sum_{n_1+\ldots+n_K=B-1-n_0} \frac{\mu_0^{-B-1-n_0}}{n_0!} \rho^{-(B-1-n_0)}}{\sum_{n_0=0}^{B} \sum_{n_1+\ldots+n_K=B-n_0} \frac{\mu_0^{-n_0}}{n_0!} \rho^{-(B-n_0)}} \\
=&\, \rho \frac{\sum_{n_0=0}^{B-1} \binom{B-n_0+K-2}{K-1} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0}}{\sum_{n_0=0}^{B} \binom{B-n_0+K-1}{K-1} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0}} \\
=&\, \rho \frac{\sum_{n_0=0}^{B-1} \frac{(B-n_0+K-2)!}{(K-1)!(B-n_0-1)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0}}{\sum_{n_0=0}^{B} \frac{(B-n_0+K-1)!}{(K-1)!(B-n_0)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0}}.
\end{aligned}
$$

Define

$$
NR = \sum_{n_0=0}^{B-1} \frac{(B - n_0 + K - 2)!}{(K - 1)!(B - n_0 - 1)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0},
$$

$$
DR = \sum_{n_0=0}^{B} \frac{(B - n_0 + K - 1)!}{(K - 1)!(B - n_0)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0}.
$$

Then, $TH(B, \boldsymbol{u}^{**})/\rho$ can be written as follows.

$$
\frac{TH(B, \boldsymbol{u}^{**})}{\rho} = \frac{NR}{DR} = 1 - \frac{DR - NR}{DR} = 1 - \frac{\Delta}{DR},
$$

where $\Delta = DR - NR$, and can be calculated by

$$\Delta = \sum_{n_0=0}^{B} \frac{(B-n_0+K-1)!}{(K-1)!(B-n_0)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0} - \sum_{n_0=0}^{B-1} \frac{(B-n_0+K-2)!}{(K-1)!(B-n_0-1)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0}$$

$$= \sum_{n_0=0}^{B-1} \left[ \frac{(B-n_0+K-1)!}{(K-1)!(B-n_0)!} - \frac{(B-n_0+K-2)!}{(K-1)!(B-n_0-1)!} \right] \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0} + \frac{\mu_0^{-B}}{B!} \rho^B$$

$$= \sum_{n_0=0}^{B-1} \frac{[(B-n_0+K-1)-(B-n_0)](B-n_0+K-2)!}{(K-1)!(B-n_0)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0} + \frac{\mu_0^{-B}}{B!} \rho^B$$

$$= \sum_{n_0=0}^{B-1} \frac{(B-n_0+K-2)!}{(K-2)!(B-n_0)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0} + \frac{\mu_0^{-B}}{B!} \rho^B$$

$$= \sum_{n_0=0}^{B} \frac{(B-n_0+K-2)!}{(K-2)!(B-n_0)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0}.$$

Define

$$f(B, n_0, K) = \frac{(B-n_0+K-1)!}{(K-1)!(B-n_0)!} \frac{\mu_0^{-n_0}}{n_0!} \rho^{n_0},$$

$$g(B, n_0, K) = \frac{(B/2-n_0+K-1)!}{(K-1)!(B/2-n_0)!} \frac{\mu_0^{-B/2-n_0}}{(B/2+n_0)!} \rho^{B/2+n_0}.$$

When $B$ is an even number, we have

$$\Delta = f(B, 0, K-1) + \sum_{n_0=1}^{B/2} f(B, n_0, K-1) + \sum_{n_0=B/2+1}^{B} f(B, n_0, K-1)$$

$$= f(B, 0, K-1) + \sum_{n_0=1}^{B/2} [f(B, n_0, K-1) + g(B, n_0, K-1)],$$

$$DR = f(B, 0, K) + \sum_{n_0=1}^{B/2} f(B, n_0, K) + \sum_{n_0=B/2+1}^{B} f(B, n_0, K)$$

$$= f(B, 0, K) + \sum_{n_0=1}^{B/2} [f(B, n_0, K) + g(B, n_0, K)].$$

We first show that $g(B, n_0, K)/f(B, n_0, K)$ goes to 0 as $B$ goes to infinity for $n_0 = 1, 2, \ldots, B$ and any positive integer $K$.

$$\lim_{B \to \infty} \frac{g(B, n_0, K)}{f(B, n_0, K)}$$

$$= \lim_{B \to \infty} \frac{(B/2 - n_0 + K - 1)!(B - n_0)!n_0!(\mu/\mu_0)^{B/2}}{(B - n_0 + K - 1)!(B/2 - n_0)!(B/2 + n_0)!}$$

$$= \lim_{B \to \infty} \frac{(B/2 - n_0 + 1) \cdots (B/2 - n_0 + K - 1)}{(B - n_0 + 1) \cdots (B - n_0 + K - 1)} \frac{(\mu/\mu_0)^{B/2}}{(n_0 + 1) \cdots (B/2 + n_0)}$$

$$= (\frac{1}{2})^{K-1} \cdot 0 = 0.$$

We next show that $[f(B, n_0, K - 1) + g(B, n_0, K - 1)]/[f(B, n_0, K) + g(B, n_0, K)]$ goes to

0 as $B$ goes to infinity for $n_0 = 1, 2, \ldots, B$ and any positive integer $K > 2$. We have shown

that for every real number $\epsilon_1 > 0$, there exists an even number $B_1$ such that for any $B > B_1$,

$\frac{g(B, n_0, K)}{f(B, n_0, K)} < \epsilon_1$. Then for every real number $\epsilon_0 > 0$, there exists an even number $B_0 > B_1$

satisfying $\frac{K-1}{B_0/2 + K - 1} < \frac{\epsilon_0}{1 + \epsilon_1}$ such that for any $B > B_0$,

$$\frac{f(B, n_0, K - 1) + g(B, n_0, K - 1)}{f(B, n_0, K) + g(B, n_0, K)} < \frac{f(B, n_0, K - 1)(1 + \epsilon_1)}{f(B, n_0, K)}$$

$$= \frac{K - 1}{B - n_0 + K - 1}(1 + \epsilon_1)$$

$$< \frac{\epsilon_0}{1 + \epsilon_1} \cdot (1 + \epsilon_1) = \epsilon_0$$

for $n_0 = 1, 2, \ldots, B/2$, and it's easy to show that

$$\frac{f(B, 0, K - 1)}{f(B, 0, K)} = \frac{K - 1}{B + K - 1} < \frac{\epsilon_0}{1 + \epsilon_1} < \epsilon_0,$$

when $n_0 = 0$. Therefore, we can show that for every real number $\epsilon_0 > 0$, there exists an

even number $B_0 > B_1$ satisfying $\frac{K-1}{B_0/2 + K - 1} < \frac{\epsilon_0}{1 + \epsilon_1}$ such that for any $B > B_0$,

$$\frac{\Delta}{DR} = \frac{f(B, 0, K - 1) + \sum_{n_0=1}^{B/2}[f(B, n_0, K - 1) + g(B, n_0, K - 1)]}{f(B, 0, K) + \sum_{n_0=1}^{B/2}[f(B, n_0, K) + g(B, n_0, K)]} < \epsilon_0.$$

Hence, we have

$$\lim_{B \to \infty} \frac{\Delta}{DR} = 0,$$

$$\lim_{B \to \infty} \frac{TH}{\rho} = 1.$$

In Theorem 4.1.1, we show that $\rho$ is the upper-bound on the throughput of the system.

Hence the allocation of $u_k^{**}$ at station $k$ for $k = 1, 2, \ldots, K$ maximizes the long-run through-

put of the system when the number of customers circulating in the system goes to infin-

ity. $\qquad \square$

Theorem 4.2.2 yields a closed-form solution of $\boldsymbol{u}^{**}$ to (4.1.3) when $\mu_k(u_k)$ $(k = 1, \ldots, K)$ is linear and strictly increasing, as shown in Corollary 4.2.2.

**Corollary 4.2.2.** *If $\mu_k(u_k) = a_k u_k$ where $a_k > 0$ for $k = 1, \ldots, K$, the optimal allocation vector $\boldsymbol{u}^{**} = [u_1^{**}, u_2^{**}, \ldots, u_K^{**}]$ satisfying (4.1.3) is given by*

$$u_k^{**} = \frac{\frac{v_k}{a_k}}{\sum_{i=1}^{K} \frac{v_i}{a_i}} U, \quad k = 1, 2, \ldots, K.$$

*And the optimal long-run average throughput of the system is equal to $\frac{v_0 U}{\sum_{i=1}^{K} \frac{v_i}{a_i}}$.*

**Remark.** The optimal allocation $\boldsymbol{u}^{**} = [u_1^{**}, u_2^{**}, \ldots, u_K^{**}]$ is proportional to the ratio of visiting ratio $v_i$ and linear coefficient $a_i$, as opposed to the square-root of the ratio for the case where a single server circulates in the system.

**Example 4.2.3.** Consider the series system as shown in Figure 3.2. Then the optimal allocation are the solutions to the following equations:

$$\mu_1(u_1) = \mu_2(u_2) = \ldots = \mu_K(u_K), \quad \sum_{j=1}^{K} u_k = U.$$

If $\mu_k(u_k) = a_k u_k$ where $a_k > 0$ for $k = 1, \ldots, K$, then the optimal allocation is given by

$$u_k^{**} = \frac{\frac{1}{a_k}}{\sum_{i=1}^{K} \left(\frac{1}{a_i}\right)} U.$$

**Example 4.2.4.** Consider the parallel system as shown in Figure 3.3. Then the optimal allocation are the solutions to the following equations:

$$\frac{\mu_1(u_1)}{p_1} = \frac{\mu_2(u_2)}{p_2} = \ldots = \frac{\mu_K(u_K)}{p_K}, \quad \sum_{j=1}^{K} u_k = U.$$

If $\mu_k(u_k) = a_k u_k$ where $a_k > 0$ for $k = 1, \ldots, K$, then the optimal allocation is given by

$$u_k^{**} = \frac{\frac{p_k}{a_k}}{\sum_{i=1}^{K} \left(\frac{p_i}{a_i}\right)} U.$$

*4.2.3 Finite Population Size $(1 < B < \infty)$*

In this section, we study the optimal allocation problem when $B$ can be any positive integer number. In Theorem 4.2.3, we show that the maximum long-run average throughput of the system is non-decreasing in the number of customers circulating in the system.

**Theorem 4.2.3.** *The optimal long-run average throughput of the system is non-decreasing in $B$, i.e., for any positive integer $B$, we have*

$$TH^*(1) \leq TH^*(2) \leq \cdots \leq TH^*(B).$$

**Proof of Theorem 4.2.3.** Let $\boldsymbol{u^*}_B$ be the optimal allocation when $B$ customers circulate in the system. By Lemma 4.1.1, for any positive integer $B$, we have

$$TH^*(B) = TH(B, \boldsymbol{u^*}_B) \leq TH(B+1, \boldsymbol{u^*}_B).$$

Also, we know that

$$TH(B+1, \boldsymbol{u^*}_B) \leq TH^*(B+1).$$

Hence, the result $TH^*(B) \leq TH^*(B+1)$ follows. $\qquad\square$

**Remark.** Yao [36] studies the properties of the throughput function for closed queueing networks and shows that the optimal loading and server-assignment policy is balanced. In our model, two major differences from Yao [36] are (1) we do not make allocation decisions over the automated station (station 0); (2) we consider allocating constant service resource $U$ rather than constant loading.

For the remaining of this section, we assume that $\frac{\mu_1(u)}{v_1} = \frac{\mu_2(u)}{v_2} = \ldots = \frac{\mu_K(u)}{v_K}$ for $0 \leq u \leq U$. Denote $\rho(u) = \frac{\mu_k(u)}{v_k}$ $(k = 1, 2, \ldots, K)$. The main result is given in Theorem 4.2.4.

**Lemma 4.2.1.**
$$C'_{u_k}(B) = -\rho'(u_k) \sum_{n=0}^{B-1} C(B-1-n)\rho(u_k)^{-n-2}.$$

**Proof of Lemma 4.2.1.** $C(B)$ can be written as follows:

$$C(B) = \sum_{n_0=0}^{B} \frac{\mu_0^{-n_0}}{n_0!} \sum_{n=0}^{B-n_0} G_k(B - n_0 - n)\rho(u_k)^{-n}.$$

Therefore,

$$
\begin{aligned}
C'_{u_k}(B) &= \sum_{n_0=0}^{B-1} \frac{\mu_0^{-n_0}}{n_0!} \sum_{n=1}^{B-n_0} (-n) G_k(B-n_0-n) \rho(u_k)^{-n-1} \rho'(u_k) \\
&= \sum_{n_0=0}^{B-1} \frac{\mu_0^{-n_0}}{n_0!} \sum_{n=0}^{B-n_0-1} (-n-1) G_k(B-n_0-n-1) \rho(u_k)^{-n-2} \rho'(u_k) \\
&= \rho(u_k)^{-1} \sum_{n_0=0}^{B-2} \frac{\mu_0^{-n_0}}{n_0!} \sum_{n=1}^{B-n_0-1} (-n) G_k(B-n_0-n-1) \rho(u_k)^{-n-1} \rho'(u_k) \\
&\quad - \mu_i^{-2} \rho'(u_k) \sum_{n_0=0}^{B-1} \frac{\mu_0^{-n_0}}{n_0!} \sum_{n=0}^{B-n_0-1} G_k(B-n_0-n-1) \rho(u_k)^{-n} \\
&= \rho(u_k)^{-1} C'_{u_k}(B-1) - \rho(u_k)^{-2} \rho'(u_k) C(B-1).
\end{aligned}
$$

The equation then follows from recursion on $B$. $\qquad\square$

**Theorem 4.2.4.** *If $\frac{\rho'(u)}{\rho(u)^2}$ is non-increasing in $u$, then $TH(B, \boldsymbol{\mu})$ is a Schur-concave function of $\boldsymbol{u}$.*

***Proof of Theorem 4.2.4.***

$$
\begin{aligned}
\frac{\partial}{\partial u_1} TH(B, \boldsymbol{u}) &= C^{-2}(B)[C'_{u_1}(B-1)C(B) - C(B-1)C'_{u_1}(B)] \\
&= C^{-2}(B)\Big[ -C(B)\rho'(u_1) \sum_{n=0}^{B-2} C(B-2-n)\rho(u_1)^{-n-2} \\
&\quad + C(B-1)\rho'(u_1) \sum_{n=0}^{B-1} C(B-1-n)\rho(u_1)^{-n-2} \Big] \\
&= C^{-2}(B)\Big\{ \sum_{n=0}^{B-2} \frac{\rho'(u_1)}{\rho(u_1)^{n+2}} [-C(B)C(B-n-2) + C(B-1)C(B-n-1)] \\
&\quad + \frac{\rho'(u_1)}{\rho(u_1)^{B+1}} C(B-1) \Big\}.
\end{aligned}
$$

By symmetry, we immediately have

$$
\begin{aligned}
(\frac{\partial}{\partial u_1} - \frac{\partial}{\partial u_2})TH(B, \boldsymbol{u}) &= C^{-2}(B)\Big\{ \sum_{n=0}^{B-2} (\frac{\rho'(u_1)}{\rho(u_1)^{n+2}} - \frac{\rho'(u_2)}{\rho(u_2)^{n+2}})[C(B-1)C(B-n-1) \\
&\quad - C(B)C(B-n-2)] + (\frac{\rho'(u_1)}{\rho(u_1)^{B+1}} - \frac{\rho'(u_2)}{\rho(u_2)^{B+1}})C(B-1) \Big\}.
\end{aligned}
$$

The quantity in the '[ ]' of the above expression is nonnegative. Therefore we have $(u_1 - u_2)(\frac{\partial}{\partial u_1} - \frac{\partial}{\partial u_2})TH(B, \boldsymbol{u}) \le 0$. $\qquad\square$

**Remarks.** If $\mu_k(u_k)$ $(k = 1, 2, \ldots, K)$ is increasing and concave, then $\rho(u)$ is increasing and concave so that the above condition automatically holds.

**Corollary 4.2.3.** *If $\rho(u)$ is increasing and concave, the optimal allocation assigns equal service resource to each of the service stations. That is, $TH(B, \boldsymbol{u}) \leq TH(B, \boldsymbol{u^*})$ for all $u$, where $u^*$ is a vector with equal elements, i.e., $u_k = U/K$ for all $k$.*

***Proof of Corollary 4.2.3.*** Since $\boldsymbol{u^*} \leq_m \boldsymbol{u}$ for all $u$ with $\sum_{k=1}^{K} = U$.  □

**Example 4.2.5.** Consider the series system as shown in Figure 3.2. Assume that $\mu_1(u) = \mu_2(u) = \ldots = \mu_K(u)$. We know that all $v_i$s $(i = 0, 1, \ldots, K)$ are equal for the series system, then we have $\frac{\mu_1(u)}{v_1} = \frac{\mu_2(u)}{v_2} = \ldots = \frac{\mu_K(u)}{v_K}$. By Corollary 4.2.3, the optimal allocation is given by $\boldsymbol{u^*} = [U/K, U/K, \ldots, U/K]$.

## 4.3 Conclusion

In this chapter, we consider allocating a finite amount of service resource which is continuously divisible and can be used by any of the service stations. Service times at a service station are exponentially distributed and their mean is a strictly increasing and concave function of the allocated service resource. We first show that system throughput is non-decreasing in the number of customers. Then, we study the optimization problem in three cases depending on the number of customers circulating in the system.

First, when there is a single customer in the system, we show that the optimal allocation is given by a set of optimization equations. In a special case when the service rate function is linear, we show that the optimal allocation of service resource to a station is proportional to the square-root of the ratio of its visiting ratio and its linear coefficient. Secondly, when the number of customers in the system increases to infinity, we show that the optimal allocation approaches to a limit which is given by a set of equations provided in Theorem 4.2.2. In a special case when the service rate function is linear, we show that the optimal allocation of service resource to a station is proportional to the ratio of its visiting ratio and its linear coefficient. Finally, for any positive number of customers in the system, we show that the system throughput as a function of service resource is Schur-concave when a certain condition is satisfied.

# Chapter 5

# OPTIMAL DYNAMIC CONTROL OF FINITE-POPULATION QUEUEING SYSTEMS: WAITING COST MINIMIZATION

In Chapter 3 and 4, we study the optimization problems with the objective of maximizing the long-run average throughput of the system. In this chapter, we study two cost minimization problems on optimal control.

## 5.1  A Parallel System with A Single Server

In this section, we consider the closed queueing system with $K$ parallel service stations as shown in Figure 3.3. A single server is able to work at each of the $K$ service stations. Let $\mu_k$ denote the rate of exponential service times at station $k$ for $k = 1, 2, \ldots, K$. Preemption is allowed. Let $c_k$ denote the cost incurred by a job waiting at station $k$ for unit time for $k = 1, 2, \ldots, K$. We consider non-idling policies, i.e., the server is not allowed to be idle whenever there is available job(s) at service stations. We formulate this problem as a Markov decision process. We use the same notation defined in Section 3.2.1 unless otherwise stated. Define $V(\boldsymbol{n})$ as the bias of state $\boldsymbol{n}$, and $h$ as the long-run average waiting cost of the system. Since every state is accessible from another state, the transition matrix consists of a single recurrent class for every deterministic stationary policy. Hence, the MDP is recurrent and $h$ exists. Define $\Lambda = B\mu_0 + \sum_{k=1}^{K} \mu_k$ as the uniformization constant. Without loss of generality, we assume that $\Lambda = 1$. Then, the optimality equation can be expressed as follows. For $0 \leq n \leq B$,

$$h + V(\boldsymbol{n}) = \sum_{k=1}^{K} c_k n_k + (B-n)\mu_0 \sum_{k=1}^{K} p_k V(\boldsymbol{n} + \boldsymbol{e}^k) + n\mu_0 V(\boldsymbol{n}) + f(\boldsymbol{n}), \qquad (5.1.1)$$

where

$$f(\boldsymbol{n}) = \sum_{k=1}^{K} \mu_k V(\boldsymbol{n}) + \begin{cases} 0 & \text{if } \boldsymbol{n} = \boldsymbol{0}, \\ \min_{i \in \boldsymbol{I_n}} \{\mu_i V(\boldsymbol{n} - \boldsymbol{e}^i) - \mu_i V(\boldsymbol{n})\} & \text{otherwise.} \end{cases}$$

We provide a partial characterization of the optimal policy in Theorem 5.1.1.

**Theorem 5.1.1.** *Suppose that there exists a station $i$ for which $c_i \mu_i \geq c_j \mu_j$ and $\mu_i(c_i - c_a) \geq \mu_j(c_j - c_a)$ for all $j = 1, 2, \ldots, K$ and $j \neq i$, and $1 \leq a \leq K$. Then, there exists an optimal policy which gives the highest priority to station $i$ within the set of all preemptive policies $\Pi_P$.*

It is important to point out that Iravani and Kolfal [10] study a similar problem. They consider a finite-population queueing system with multiple service stations and a single flexible server, under preemptive policies, where cost is incurred when a customer is waiting for service. They consider multiple classes of customers, and each class of customers generates one type of service requests. The authors investigate that applying the $c\mu$ rule in their finite-population queueing system is not always optimal to minimize the long-run average cost of the system (see Van Mieghem [32] for a brief review of literature on the $c\mu$ rule). They find the conditions under which static-priority rules are optimal independent of customer arrival rate and customer population size. In contrast, in our model, we consider identical customers, and a customer requests for one of the $K$ types of service with a certain probability. We show that a stronger condition than the simple $c\mu$ rule is required so that a static-priority policy is optimal.

***Proof of Theorem 5.1.1.*** In order to prove Theorem 5.1.1 we show that the result holds for the $m$-period expected total waiting cost problem defined by (5.1.2) for all $m \geq 0$. Let $N_k$ denote the state of the system at period $k$ and $d_k(N_k)$ the decision rule at period $k$ in state $N_k$ under policy $\pi$. Let $c(\boldsymbol{N}, d)$ denote the waiting cost incurred when the system is in state $\boldsymbol{N}$ and the action $d$ is taken. We define $V_m(\pi, \boldsymbol{n})$ as the $m$-period expected waiting cost under policy $\pi$ when the initial state is $\boldsymbol{n}$, i.e.,

$$V_m(\pi, \boldsymbol{n}) \equiv E\left[\sum_{k=0}^{m-1} c(\boldsymbol{N}_k, d_k(\boldsymbol{N}_k))\right].$$

Then, the optimal $m$-period expected waiting cost is

$$V_m^*(\boldsymbol{n}) \equiv \inf_{\pi \in \Pi_P} V_m(\pi, \boldsymbol{n}). \tag{5.1.2}$$

We let $h(\pi, \boldsymbol{n})$ be the long-run average waiting cost under policy $\pi$, given that the initial

state of the system is $\boldsymbol{n}$, i.e.,

$$h(\pi, \boldsymbol{n}) \equiv \liminf_{m \to \infty} \frac{1}{m} V_m(\pi, \boldsymbol{n}).$$

Let $\mu \equiv \sum_{k=1}^{K} \mu_k$. Then, the optimality equation for the finite-period problem can be expressed as follows. For all $m \geq 0$,

$$V_{m+1}(\boldsymbol{n}) = \sum_{k=1}^{K} c_k n_k + (B - n)\mu_0 \sum_{k=1}^{K} p_k V_m(\boldsymbol{n} + \boldsymbol{e}^k) + n\mu_0 V_m(\boldsymbol{n}) + f_m(\boldsymbol{n}), \quad \text{for } 0 \leq n \leq B,$$

$$(5.1.3)$$

where

$$f_m(\boldsymbol{n}) = \mu V_m(\boldsymbol{n}) + \begin{cases} 0 & \text{if } \boldsymbol{n} = \boldsymbol{0}, \\ \min_{i \in \boldsymbol{I_n}} \{\mu_i V_m(\boldsymbol{n} - \boldsymbol{e}^i) - \mu_i V_m(\boldsymbol{n})\} & \text{otherwise.} \end{cases}$$

We assume that $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$.

Without loss of generality, we assume that $c_1\mu_1 \geq c_j\mu_j$ and $\mu_1(c_1 - c_a) \geq \mu_j(c_j - c_a)$ for all $j = 2, 3, \ldots, K$, and $1 \leq a \leq K$. For system state $\boldsymbol{n}$, where $2 \leq n \leq B$, let $l_1(\boldsymbol{n}) \equiv \min\{k : k \in \boldsymbol{I_n}\}$ and $l_2(\boldsymbol{n}) \equiv \min\{k : k \in \boldsymbol{I_{n-e^{l_1(n)}}}\}$. Let $\boldsymbol{n}^a = \boldsymbol{n} + \boldsymbol{e}^a$ $(1 \leq a \leq K)$. We will show that, for all $m \geq 0$, $2 \leq n \leq B$, $n_1 \geq 1$, and $1 \leq a \leq K$

$$\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n}) \leq 0, \quad (5.1.4)$$

$$\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n}) \leq 0. \quad (5.1.5)$$

where $j \neq 1$ and $j \in \boldsymbol{I_n}$. We will use induction on $m$. Since $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$, then the inequalities automatically hold at period 0. Assume that inequalities (5.1.4) and (5.1.5) hold at period $m$. We will show that they also hold at period $m + 1$.

Proof of (5.1.4): We will consider two cases.

(a) Suppose that $l_2(\boldsymbol{n}) = 1$. Using equation (5.1.3), we have

$$\mu_1 V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_{m+1}(\boldsymbol{n})$$

$$= - c_1\mu_1 + c_j\mu_j + (B - n + 1)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k))]$$

$$+ (B - n)\mu_0 \sum_k p_k[(\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ (n - 1)\mu_0[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j)] + n\mu_0[(\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1[\mu_1 V_m(\boldsymbol{n} - 2\boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_1)V_m(\boldsymbol{n} - \boldsymbol{e}^1)]$$

$$+ (\mu - \mu_1)[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$= - c_1\mu_1 + c_j\mu_j$$

$$+ (B - n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k) + (\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ \mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k)) + (\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1[\mu_1 V_m(\boldsymbol{n} - 2\boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_1)V_m(\boldsymbol{n} - \boldsymbol{e}^1)]$$

$$+ ((n - 1)\mu_0 + \mu - \mu_1)[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})],$$

which is non-positive by the inductive hypothesis for (5.1.4) and (5.1.5) at period $m$, the condition that $c_1\mu_1 \geq c_j\mu_j$ for all $j = 2, 3, \ldots, K$, and the assumption that $B\mu_0 + \mu = 1$.

(b) Suppose that $l_2(\boldsymbol{n}) > 1$. Using equation (5.1.3), we have

$$\mu_1 V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_{m+1}(\boldsymbol{n})$$

$$= - c_1\mu_1 + c_j\mu_j$$

$$+ (B - n + 1)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k))]$$

$$+ (B - n)\mu_0 \sum_k p_k[(\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ (n - 1)\mu_0[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j)] + n\mu_0[(\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1\mu_{l_2(\boldsymbol{n})}V_m(\boldsymbol{n} - \boldsymbol{e}^1 - \boldsymbol{e}^{l_2(\boldsymbol{n})}) - \mu_j\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_1)\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1)$$

$$+ \mu_1(\mu - \mu_{l_2(\boldsymbol{n})})V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j(\mu - \mu_1)V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)(\mu - \mu_1)V_m(\boldsymbol{n})$$

60

$$= -c_1\mu_1 + c_j\mu_j$$

$$+ (B-n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k) + (\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ \mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k)) + (\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1[\mu_{l_2(\boldsymbol{n})} V_m(\boldsymbol{n} - \boldsymbol{e}^1 - \boldsymbol{e}^{l_2(\boldsymbol{n})}) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_{l_2(\boldsymbol{n})})V_m(\boldsymbol{n} - \boldsymbol{e}^1)]$$

$$+ ((n-1)\mu_0 + \mu - \mu_1)[\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n} - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})],$$

which is non-positive by the inductive hypothesis for (5.1.4) and (5.1.5) at period $m$, the fact that $l_2(\boldsymbol{n}) = l_1(\boldsymbol{n} - \boldsymbol{e}^1)$, the condition that $c_1\mu_1 \geq c_j\mu_j$ for all $j = 2, 3, \ldots, K$, and the assumption that $B\mu_0 + \mu = 1$.

Proof of (5.1.5). We will consider two cases:

(a) Suppose that $l_2(\boldsymbol{n}^a) = 1$. Using equation (5.1.3), we have

$$\mu_1 V_{m+1}(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_{m+1}(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_{m+1}(\boldsymbol{n})$$

$$= -\mu_1(c_1 - c_a) + \mu_j(c_j - c_a)$$

$$+ (B-n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j + \boldsymbol{e}^k))]$$

$$+ (B-n)\mu_0 \sum_k p_k[(\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ n\mu_0[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j)] + n\mu_0[(\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1\mu_i V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 - \boldsymbol{e}^i) - \mu_j\mu_i V_m(\boldsymbol{n}^a - \boldsymbol{e}^j - \boldsymbol{e}^i) + (\mu_j - \mu_1)\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1)$$

$$+ \mu_1(\mu - \mu_i)V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j(\mu - \mu_i)V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)(\mu - \mu_1)V_m(\boldsymbol{n})$$

$$= -\mu_1(c_1 - c_a) + \mu_j(c_j - c_a)$$

$$+ (B-n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j + \boldsymbol{e}^k) + (\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ n\mu_0[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_i[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 - \boldsymbol{e}^i) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j - \boldsymbol{e}^i) + (\mu_j - \mu_1)V_m(\boldsymbol{n} - \boldsymbol{e}^1)]$$

$$+ (\mu - \mu_i)[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})],$$

which is non-positive by the inductive hypothesis for (5.1.4) and (5.1.5) at period $m$, the condition that $\mu_1(c_1 - c_a) \geq \mu_j(c_j - c_a)$ for all $j = 2, 3, \ldots, K$, and $1 \leq a \leq K.$, and the assumption that $B\mu_0 + \mu = 1$.

(b) Suppose that $a \geq i$ and $l_2(\boldsymbol{n}^a) > 1$. Using equation (5.1.3), we have

$$\mu_1 V_{m+1}(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_{m+1}(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_{m+1}(\boldsymbol{n})$$

$$= -\mu_1(c_1 - c_a) + \mu_j(c_j - c_a)$$

$$+ (B - n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j + \boldsymbol{e}^k))]$$

$$+ (B - n)\mu_0 \sum_k p_k[(\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ n\mu_0[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j)] + n\mu_0[(\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1 \mu_{l_2(\boldsymbol{n}^a)} V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 - \boldsymbol{e}^{l_2(\boldsymbol{n}^a)}) - \mu_j \mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_1)\mu_1 V_m(\boldsymbol{n} - \boldsymbol{e}^1)$$

$$+ \mu_1(\mu - \mu_{l_2(\boldsymbol{n}^a)})V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j(\mu - \mu_1)V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)(\mu - \mu_1)V_m(\boldsymbol{n})$$

$$= -\mu_1(c_1 - c_a) + \mu_j(c_j - c_a)$$

$$+ (B - n)\mu_0 \sum_k p_k[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 + \boldsymbol{e}^k) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j + \boldsymbol{e}^k) + (\mu_j - \mu_1)V_m(\boldsymbol{n} + \boldsymbol{e}^k)]$$

$$+ n\mu_0[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})]$$

$$+ \mu_1[\mu_{l_2(\boldsymbol{n}^a)} V_m(\boldsymbol{n}^a - \boldsymbol{e}^1 - \boldsymbol{e}^{l_2(\boldsymbol{n}^a)}) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j - \boldsymbol{e}^1) + (\mu_j - \mu_{l_2(\boldsymbol{n}^a)})V_m(\boldsymbol{n} - \boldsymbol{e}^1)$$

$$+ (\mu - \mu_1)[\mu_1 V_m(\boldsymbol{n}^a - \boldsymbol{e}^1) - \mu_j V_m(\boldsymbol{n}^a - \boldsymbol{e}^j) + (\mu_j - \mu_1)V_m(\boldsymbol{n})],$$

which is non-positive by the inductive hypothesis for (5.1.4) and (5.1.5) at period $m$, the condition that $\mu_1(c_1 - c_a) \geq \mu_j(c_j - c_a)$ for all $j = 2, 3, \ldots, K$, and $1 \leq a \leq K$, and the assumption that $B\mu_0 + \mu = 1$.

Hence, we show that jobs at station 1 should be served ahead of jobs at station $j$ ($j = 2, 3, \ldots, K$) if jobs are available at both stations. In other words, we should give the highest priority to station 1 in order to minimize the long-run average waiting cost of the system. $\qquad \square$

**Remarks.** Theorem 5.1.1 partially characterizes the optimal policy for the waiting cost minimization problem. It shows that for this finite-population queueing system, a stronger condition than the simple $c\mu$ rule is required to characterize the optimal static-priority policy. We conducted a brief numerical study for this system when $K = 2$. The numerical result shows that $c_1\mu_1 \geq c_2\mu_2$ on its own is not a sufficient condition for the optimal policy that gives priority to station 1. For example, when we set the values for the parameters

as $[\mu_1, \mu_2] = [1.3, 1]$, $[c_1, c_2] = [1, 1.1]$ and $[p_1, p_2] = [0.6, 0.4]$, and allow the number of customers in the system to be 3, 4 or 5, the optimal policy gives priority to station 2 when there are jobs available in both service stations. In this scenario, we have $c_1\mu_1 > c_2\mu_2$, but station 1 is not prioritized. Hence, for our finite-population queueing model, a stronger condition than the simple $c\mu$ rule is required to be sufficient.

It is important to point out that when the system reduces to $K = 2$, Theorem 5.1.1 provides a complete characterization of the optimal policy, and the required sufficient condition reduces to $c_1\mu_1 \geq c_2\mu_2$ and $c_1 \geq c_2$.

## 5.2  A Parallel System with Continuous Resource

In this section, we still consider the parallel system with $K \geq 2$ single-server service as shown in Figure 3.3. However, we will consider continuous service resource rather than discrete server(s). Suppose that there is a finite amount of continuous service capacity $\mu$ which can be used at any of the $K$ service stations. Preemption is allowed. We are interested in optimal policy of dynamically allocating the service capacity to each of the service stations in order to minimize the long-run average waiting cost. All service stations share this fixed resource and satisfy the constraint $\sum_{k=1}^{K} \mu_k = \mu$, where $\mu_k$ is the service capacity allocated to station $k$. We assume that the amount of intrinsic work required at all stations are i.i.d. exponentials with mean 1. Hence, when service capacity $\mu_k$ is allocated to station $k$, service times at station $k$ will be exponentially distributed with mean $1/\mu_k$ ($k = 1, 2, \ldots, K$). Let $c_k$ denote the cost incurred when a customer is waiting in station $k$ ($k = 1, 2, \ldots, K$) per unit time. We consider non-idling policies, i.e., the service capacity is completely allocated whenever there is available job(s) at service stations.

We formulate this problem as a Markov decision process. We use the same notation defined in Section 5.1 unless otherwise stated. Define $V(\boldsymbol{n})$ as the bias of state $\boldsymbol{n}$, and $h$ as the long-run average waiting cost of the system. Following the same argument, we know that there exists a stationary average optimal policy and hence $h$ exists. Define $\Lambda = B\mu_0 + \sum_{k=1}^{K} \mu_k$ as the uniformization constant. Without loss of generality, we assume that $\Lambda = 1$. Then, the optimality equation can be expressed as follows. For $0 \leq n \leq B$,

$$h + V(\boldsymbol{n}) = \sum_{i=1}^{K} c_i n_i + (B - n)\mu_0 \sum_{i=1}^{K} p_i V(\boldsymbol{n} + \boldsymbol{e}^i) + n\mu_0 V(\boldsymbol{n}) + f(\boldsymbol{n}),$$

63

where

$$f(\boldsymbol{n}) = \begin{cases} \mu V(\boldsymbol{0}) & \text{if } \boldsymbol{n} = \boldsymbol{0}, \\ \min_{\boldsymbol{\mu} \in \mathcal{M}(\boldsymbol{n})} \{ \sum_{i=1}^{K} \mu_i V(\boldsymbol{n} - \boldsymbol{e}^i) \} & \text{otherwise,} \end{cases}$$

$\boldsymbol{\mu}$ is the allocation vector, and $\mathcal{M}(\boldsymbol{n}) = \{ (\mu_1, \mu_2, \ldots, \mu_K) : \sum_{i=1}^{K} \mu_i = \mu, \mu_i \geq 0 \text{ for } i \in \boldsymbol{I_n} \text{ and } \mu_i = 0 \text{ for } i \notin \boldsymbol{I_n} \}$. We first partially characterize optimal policies in Theorem 5.2.1.

**Theorem 5.2.1.** *There exists an optimal policy that assigns all service capacity to a single station.*

***Proof of Theorem 5.2.1.*** The minimum in the optimality equations is

$$\min_{\boldsymbol{\mu} \in \mathcal{M}(\boldsymbol{n})} \sum_{i=1}^{K} \{ \mu_i V(\boldsymbol{n} - \boldsymbol{e}^i) \}$$

Because the terms in the minimum operator are linear in $\mu_i$ ($1 \leq i \leq K$) for each state $\boldsymbol{n}$, the minimum must occur at the extreme points (i.e., $\mu_i = 0$ or $\mu$ for $1 \leq i \leq K$). Since we consider non-idling policies, the result follows. $\qquad \square$

Theorem 5.2.1 says that the search of an optimal policy can be narrowed down to "bang-bang" policies only. For bang-bang policies, we need to only consider discrete service rates either at its minimum or maximum feasible levels for each non-empty service station. Then, function $f(\boldsymbol{n})$ can be expressed as follows:

$$f(\boldsymbol{n}) = \begin{cases} \mu V(\boldsymbol{0}) & \text{if } \boldsymbol{n} = \boldsymbol{0}, \\ \mu \min_{i \in \boldsymbol{I_n}} \{ V(\boldsymbol{n} - \boldsymbol{e}^i) \} & \text{otherwise.} \end{cases}$$

We characterize the optimal policy in Theorem 5.2.2. Note that the optimization problem is equivalent to (5.1.3) with a discrete service resource where $\mu_i = \mu$ for $K = 2$.

**Theorem 5.2.2.** *The optimal policy which minimizes the long-run average waiting cost assigns all of the service rate to the service station which has the largest value of $c_k$ among all non-empty service stations.*

***Proof of Theorem 5.2.2.*** For state $\boldsymbol{n}$ where $n \geq 2$, suppose $c_{i_r}$'s are ordered as $c_{i_1} \geq c_{i_2} \geq \ldots \geq c_{i_R}$ for all $i_r \in \boldsymbol{I_n}$ and $R$ is the number of elements of the set $\boldsymbol{I_n}$. In order to prove Theorem 5.2.2, we show that the result holds for the $m$-period expected waiting cost

problem for all $m \geq 0$. Then, the optimality equation for the finite-period problem can be expressed as follows: For all $m \geq 0$ and $0 \leq n \leq B$,

$$V_{m+1}(\boldsymbol{n}) = \sum_{i=1}^{K} c_i n_i + (B-n)\mu_0 \sum_{i=1}^{K} p_i V_m(\boldsymbol{n} + \boldsymbol{e}^i) + n\mu_0 V_m(\boldsymbol{n}) + f_m(\boldsymbol{n}), \qquad (5.2.1)$$

where

$$f_m(\boldsymbol{n}) = \begin{cases} \mu V_m(\boldsymbol{0}) & \text{if } \boldsymbol{n} = \boldsymbol{0}, \\ \mu \min_{i \in \boldsymbol{I_n}} \{V_m(\boldsymbol{n} - \boldsymbol{e}^i)\} & \text{otherwise.} \end{cases}$$

We assume that $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$.

We will show that, for all $m \geq 0$,

$$V_m(\boldsymbol{n} - \boldsymbol{e}^i) - V_m(\boldsymbol{n} - \boldsymbol{e}^j) \leq 0, \qquad (5.2.2)$$

for $i = i_1$ and $j \in \{i_2, \ldots, i_R\}$, and $n_i, n_j \geq 1$. We will use induction on $m$. Since $V_0(\boldsymbol{n}) = 0$ for all $\boldsymbol{n}$, (5.2.2) automatically hold at period 0. Assume that (5.2.2) hold at period $m$. We will show that they also hold at period $m + 1$.

Proof of (5.2.2). We will consider two cases:

1. Suppose that $n_i > 1$. Using equation (5.2.1), we have

$$\begin{aligned} &V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^i) - V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j) \\ &= -c_i + c_j \\ &\quad + (B - n + 1)\mu_0 \sum_{k=1}^{K} p_k [V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^k) - V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k)] \\ &\quad + (n-1)\mu_0 [V_m(\boldsymbol{n} - \boldsymbol{e}^i) - V_m(\boldsymbol{n} - \boldsymbol{e}^j)] \\ &\quad + \mu [V_m(\boldsymbol{n} - 2\boldsymbol{e}^i) - V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^i)], \end{aligned}$$

which is non-positive by the inductive hypothesis for (5.2.2) at period $m$, the condition that $c_i \leq c_j$, and the assumption that $\Lambda = 1$.

2. Suppose that $n_i = 1$. Using equation (5.2.1), we have

$$V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^i) - V_{m+1}(\boldsymbol{n} - \boldsymbol{e}^j)$$

$$= - c_i + c_j$$

$$+ (B - n + 1)\mu_0 \sum_{k=1}^{K} p_k [V_m(\boldsymbol{n} - \boldsymbol{e}^i + \boldsymbol{e}^k) - V_m(\boldsymbol{n} - \boldsymbol{e}^j + \boldsymbol{e}^k)]$$

$$+ (n - 1)\mu_0 [V_m(\boldsymbol{n} - \boldsymbol{e}^i) - V_m(\boldsymbol{n} - \boldsymbol{e}^j)]$$

$$+ \mu [V_m(\boldsymbol{n} - \boldsymbol{e}^i - \boldsymbol{e}^{i_2}) - V_m(\boldsymbol{n} - \boldsymbol{e}^j - \boldsymbol{e}^i)]. \tag{5.2.3}$$

Consider two cases where $j = i_2$ or not. If $j = i_2$, then the last term of (5.2.3) is zero; otherwise, it is non-positive by the inductive hypothesis for state $\boldsymbol{n} - \boldsymbol{e}^i$ at period $m$. Hence the right-hand side of equation (5.2.3) is non-positive by the inductive argument at period $m$, the condition that $c_i \leq c_j$, and the assumption that $\Lambda = 1$.

Let $\pi^*$ be the policy that gives priority to the non-empty station with the largest value of $c_k$. By (5.2.2), we have

$$V_m(\pi^*, \boldsymbol{n}) \leq V_m(\pi, \boldsymbol{n}) \tag{5.2.4}$$

for all $\pi \in \Pi_P$ and for all $m$. Dividing both sides of (5.2.4) by $m$ and taking limits as $m$ approaches infinity the long-run average waiting cost result follows, i.e.,

$$h(\pi^*, \boldsymbol{n}) \leq h(\pi, \boldsymbol{n})$$

for all $\pi \in \Pi_P$. $\qquad\square$

**Remark.** Theorem 5.2.2 shows that the service facility needs to put all of the resource to the non-empty station that has the largest waiting cost rate in order to minimize the long-run average waiting cost.

# Chapter 6

# FUTURE RESEARCH

In this chapter, we briefly discuss some possible extensions to our study for future research. We first propose some possibilities for dynamic control problems.

**Multiple servers:** Consider the parallel-series queueing network which is described in Chapter 3, as shown in Figure 3.1. In our work, we consider a single flexible server to be able to work at any of the service stations. To extend this, we may consider allocating multiple servers. Suppose that each server is allowed to work at any of the $K$ service stations. Service times at a service station follow a given distribution which depends only on the service station. We may consider collaborative or non-collaborative servers.

**Emergency v.s. non-emergency:** Consider a system with an emergency station and non-emergency stations. For example, in a hospital system, the emergency room is usually given the top priority over other clinical units. We may model this system by a parallel system as shown in Figure 3.3. Suppose that station 1 is an emergency station which has absolute priority over other non-emergency stations (preemptive priority). A single flexible server is assigned to work at the emergency and non-emergency stations. Future work could investigate how to prioritize non-emergency stations to maximize the long-run average throughput of the system.

**Collaborative station:** Consider a system with stations which require more than one servers to process a request. We may use the parallel system with two service stations as shown in Figure **??** to model it. Suppose there are two flexible servers working at station 1 and 2. The service operation at station 1 requires two servers, while the service operation at station 2 needs only one server. Service times are random variables whose distributions depend only on the service station. Suppose preemption is not allowed. Future work could seek the optimal policy to maximize the long-run average throughput of the system.

Next, we provide some possible extensions to our static design problems.

**Multiple servers:** Consider the closed queueing network studied in Chapter 4. In our work, we consider allocating a fixed amount of continuous service resource to the $K$ service stations. Another consideration can be allocating discrete service resource, i.e., multiple servers. Suppose that each server is able to work at any of the $K$ service stations. Service times at a service station follow a given distribution which depends only on the service station. We may consider either collaborative or non-collaborative servers.

**Waiting cost minimization:** In our work, we consider maximizing the long-run average throughput of the system. We may consider waiting cost setup for the design problems. Suppose cost is incurred when customers are waiting in the service stations at the second stage, either waiting in the queue or being served. The cost rate depends only on the service stations. Assume that a fixed amount of service resource is available to be allocated to each of the $K$ service stations. We are interested in looking for optimal allocation in order to minimize the long-run average total waiting cost of the system. We expect that the optimal allocation would assign more service resource to a service station which has a higher cost rate.

# BIBLIOGRAPHY

[1] S. Andradottir and H. Ayhan. Throughput maximization for tandem lines with two stations and flexible servers. *OPERATIONS RESEARCH*, 53(3):516–531, 2005.

[2] S. Andradottir, H. Ayhan, and D.G. Down. Server assignment policies for maximizing the steady-state throughput of finite queueing systems. *MANAGEMENT SCIENCE*, 47(10):1421–1439, 2001.

[3] G. Bolch, S. Greiner, H. de Meer, and K.S. Trivedi. *Queueing networks and Markov chains : modeling and performance evaluation with computer science applications*. Wiley-Interscience, Hoboken, N.J., 2nd edition, 2006.

[4] J.P. Buzen. Computational algorithms for closed queueing networks with exponential servers. *Communications of the ACM*, 16(9):527–531, 1973.

[5] L. Green. Queueing analysis in healthcare. In *Patient flow: reducing delay in healthcare delivery*, chapter 10. Springer, 2006.

[6] D. Gross and C.M. Harris. *Fundamentals of queueing theory*. Wiley, New York, 3rd edition, 1998.

[7] L. Haque and M.J. Armstrong. A survey of the machine interference problem. *EUROPEAN JOURNAL OF OPERATIONAL RESEARCH*, 179(2):469–482, 2007.

[8] J.M. Harrison. Dynamic scheduling of a multiclass queue - discount optimality. *OPERATIONS RESEARCH*, 23(2):270–282, 1975.

[9] W.J. Hopp, S.M.R. Iravani, B.Y. Shou, and R. Lien. Design and control of agile automated conwip production lines. *NAVAL RESEARCH LOGISTICS*, 56(1):42–56, 2009.

[10] S.M.R. Iravani and B. Kolfal. When does the $c\mu$ rule apply to finite-population queueing systems? *OPERATIONS RESEARCH LETTERS*, 33(3):301–304, 2005.

[11] S.M.R. Iravani and V. Krishnamurthy. Workforce agility in repair and maintenance environments. *MANUFACTURING SERVICE OPERATIONS MANAGEMENT*, 9(2):168–184, 2007.

[12] S.M.R. Iravani, V. Krishnamurthy, and G.H. Chao. Optimal server scheduling in non-preemptive finite-population queueing systems. *QUEUEING SYSTEMS*, 55(2):95–105, 2007.

[13] G.P. Klimov. Time-sharing service systems .1. *THEORY OF PROBABILITY AND ITS APPLICATIONS*, 19(3):532–551, 1974.

[14] G.P. Klimov. Time-sharing service systems .2. *THEORY OF PROBABILITY AND ITS APPLICATIONS*, 23(2):314–321, 1978.

[15] G. Koole and R. Righter. Optimal control of tandem reentrant queues. *QUEUEING SYSTEMS*, 28(4):337–347, 1998.

[16] V.G. Kulkarni. *Modeling and analysis of stochastic systems*. CRC Press, Boca Raton, 2nd edition, 2010.

[17] H.F. Lee, M.M. Srinivasan, and C.A. Yano. Characteristics of optimal workload allocation for closed queuing-networks. *PERFORMANCE EVALUATION*, 12(4):255–268, 1991.

[18] J. Medhi. *Stochastic models in queueing theory*. Academic Press, Amsterdam; Boston, 2nd edition, 2003.

[19] M.F. Neuts. *Matrix-geometric solutions in stochastic models : an algorithmic approach*. Johns Hopkins University Press, Baltimore, 1981.

[20] C.H. Ng and B.H. Soong. *Queueing modelling fundamentals with applications in communication networks*. Wiley, Chichester, England; Hoboken, NJ, 2nd edition, 2008.

[21] J. Palesano and J. Chandra. A machine interference problem with multiple types of failures. *INTERNATIONAL JOURNAL OF PRODUCTION RESEARCH*, 24(3):567–582, 1986.

[22] M.L. Puterman. *Markov decision processes : discrete stochastic dynamic programming*. John Wiley & Sons, New York, 1994.

[23] M. Shaked and J.G. Shanthikumar. *Stochastic orders and their applications*. Academic Press, Boston, 1994.

[24] M. Shaked and J.G. Shanthikumar. *Stochastic orders*. Springer, New York, 2007.

[25] J.G. Shanthikumar and D.D. Yao. On server allocation in multiple center manufacturing systems. *OPERATIONS RESEARCH*, 36(2):333–342, 1988.

[26] K.E. Stecke. On the nonconcavity of throughput in certain closed queuing-networks. *PERFORMANCE EVALUATION*, 6(4):293–305, 1986.

[27] K.E. Stecke and J.E. Aronson. Review of operator machine interference models. *INTERNATIONAL JOURNAL OF PRODUCTION RESEARCH*, 23(1):129–151, 1985.

[28] K.E. Stecke and T.L. Morin. The optimality of balancing workloads in certain types of flexible manufacturing systems. *EUROPEAN JOURNAL OF OPERATIONAL RESEARCH*, 20(1):68–82, 1985.

[29] S. Stidham Jr. *Optimal design for queueing systems*. Chapman & Hall/CRC/Taylor & Francis, Boca Raton, FL, 2009.

[30] J. Sztrik. Finite source queuing systems and their applications: a bibliography. Technical report, University of Debrecen.

[31] H.C. Tijms. *Stochastic modelling and analysis : a computational approach*. Wiley, Chichester; New York, 1986.

[32] J.A. Van Mieghem. Dynamic scheduling with convex delay costs: The generalized c$\mu$ rule. *The Annals of Applied Probability*, 5(3):809–833, 1995.

[33] M.P. Van Oyen, E.G.S. Gel, and W.J. Hopp. Performance opportunity for workforce agility in collaborative and noncollaborative work systems. *IIE TRANSACTIONS*, 33(9):761–777, 2001.

[34] J. Walrand. *An introduction to queueing networks*. Prentice Hall, Englewood Cliffs, N.J, 1988.

[35] N. Yankovic and L.V. Green. Identifying good nursing levels: A queuing approach. *OPERATIONS RESEARCH*, 59(4):942–955, 2011.

[36] D.D. Yao. Some properties of the throughput function of closed networks of queues. *Operations Research Letters*, 3(6):313–317, 1985.