

# On Occupancy Based Randomized Load Balancing for Large Systems with General Distributions

by

Thirupathaiah Vasantam

A thesis  
presented to the University of Waterloo  
in fulfillment of the  
thesis requirement for the degree of  
Doctor of Philosophy  
in  
Electrical and Computer Engineering

Waterloo, Ontario, Canada, 2019

© Thirupathaiah Vasantam 2019

## Examining Committee Membership

The following served on the Examining Committee for this thesis. The decision of the Examining Committee is by majority vote.

External Examiner: Sem Borst  
Professor, Dept. of Mathematics and Computer Science  
Eindhoven University of Technology

Supervisor(s): Ravi R. Mazumdar  
Professor, Dept. of ECE, University of Waterloo

Internal Members: Andrew Heunis  
Professor, Dept. of ECE, University of Waterloo  
Patrick Mitran  
Associate Professor, Dept. of ECE, University of Waterloo

Internal-External Member: Hossein Abouee Mehrizi  
Associate Professor, Dept. of Management Sciences  
University of Waterloo

I hereby declare that I am the sole author of this thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

I understand that my thesis may be made electronically available to the public.

## Abstract

Multi-server architectures are ubiquitous in today's information infrastructure whether for supporting cloud services, web servers, or for distributed storage. The performance of multi-server systems is highly dependent on the load distribution. This is affected by the use of load balancing strategies. Since both latency and blocking are important features, it is most reasonable to route an incoming job to a server that is lightly loaded. Hence a good load balancing policy should be dependent on the states of servers. Since obtaining information about the remaining workload of servers for every arrival is very hard, it is preferable to design load balancing policies that depend on occupancy or the number of progressing jobs of servers. Furthermore, if the system has a large number of servers, it is not practical to use the occupancy information of all the servers to dispatch or route an arrival due to high communication cost. In large-scale systems that have tens of thousands of servers, the policies which use the occupancy information of only a finite number of randomly selected servers to dispatch an arrival result in lower implementation cost than the policies which use the occupancy information of all the servers. Such policies are referred to as occupancy based randomized load balancing policies.

Motivated by cloud computing systems and web-server farms, we study two types of models. In the first model, each server is an Erlang loss server, and this model is an abstraction of Infrastructure-as-a-Service (IaaS) clouds. The second model we consider is one with processor sharing servers that is an abstraction of web-server farms which serve requests in a round-robin manner with small time granularity. The performance criterion for web-servers is the response time or the latency for the request to be processed. In most prior works, the analysis of these models was restricted to the case of exponential job length distributions and in this dissertation we study the case of general job length distributions.

To analyze the impact of a load balancing policy, we need to develop models for the system's dynamics. In this dissertation, we show that one can construct useful Markovian models. For occupancy based randomized routing policies, due to complex interdependencies between servers, an exact analysis is mostly intractable. However, we show that the multi-server systems that have an occupancy based randomized load balancing policy are examples of weakly interacting particle systems. In these systems, servers are interacting particles whose states lie in an uncountable state space. We develop a mean-field analysis to understand a server's behavior as the number of servers becomes large. We

show that under certain assumptions, as the number of servers increases, the sequence of empirical measure-valued Markov processes which model the systems' dynamics converges to a deterministic measure-valued process referred to as the mean-field limit. We observe that the mean-field equations correspond to the dynamics of the distribution of a non-linear Markov process. A consequence of having the mean-field limit is that under minor and natural assumptions on the initial states of servers, any finite set of servers can be shown to be independent of each other as the number of servers goes to infinity. Furthermore, the mean-field limit approximates each server's distribution in the transient regime when the number of servers is large.

A salient feature of loss and processor sharing systems in the setting where their time evolution can be modeled by reversible Markov processes is that their stationary occupancy distribution is insensitive to the type of job length distribution; it depends only on the average job length but not on the type of the distribution. This property does not hold when the number of servers is finite in our context due to lack of reversibility. We show however that the fixed-point of the mean-field is insensitive to the job length distributions for all occupancy based randomized load balancing policies when the fixed-point is unique for job lengths that have exponential distributions. We also provide some deeper insights into the relationship between the mean-field and the distributions of servers and the empirical measure in the stationary regime.

Finally, we address the accuracy of mean-field approximations in the case of loss models. To do so we establish a functional central limit theorem under the assumption that the job lengths have exponential distributions. We show that a suitably scaled fluctuation of the stochastic empirical process around the mean-field converges to an Ornstein-Uhlenbeck process. Our analysis is also valid for the Halfin-Whitt regime in which servers are critically loaded. We then exploit the functional central limit theorem to quantify the error between the actual blocking probability of the system with a large number of servers and the blocking probability obtained from the fixed-point of the mean-field. In the Halfin-Whitt regime, the error is of the order inverse square root of the number of servers. On the other hand, for a light load regime, the error is smaller than the inverse square root of the number of servers.

## Acknowledgements

First and foremost, I would like to thank my supervisor Prof. Ravi R. Mazumdar for his valuable personal and professional guidance and support throughout my doctoral studies. I am grateful to him for his valuable advice about academic career. His enthusiasm and positive attitude towards research is inspirational.

In addition, I would like to extend my gratitude to Dr. Arpan Mukhopadhyay of University of Warwick, UK for his constant encouragement and guidance that helped me to grow as a researcher. His collaboration led to the results of Chapters 2 and 3. I thoroughly enjoyed our collaboration.

I thank my thesis committee members: Prof. Andrew Heunis, Prof. Patrick Mitran, Prof. Hossein Abouee Mehrizi, and Prof. Sem Borst for serving on my thesis committee and carefully reading my dissertation. Their valuable comments have helped me to improve the quality of the presentation of this dissertation.

I extend my utmost gratitude to my friends for enriching my time in the beautiful city of Waterloo.

Finally, thanks to my mother and brother who supported and believed me throughout my education. I am deeply indebted to them for all their efforts to support my family. I am also thankful to my wife Shireesha for her consistent support throughout my doctoral studies. I am grateful to my son Yuvan for coming into our life and bringing us joy.

## Dedication

This dissertation is dedicated to my family.

# Table of Contents

|  |           |
|--|-----------|
| List of Tables   | xii       |
| List of Figures  | xiii      |
| <b>1 Introduction</b>  | <b>1</b>  |
| 1.1 Models   | 4         |
| 1.1.1 The SQ( $d$ ) Policy for Erlang Loss Systems   | 5         |
| 1.1.2 Occupancy Based Randomized Routing for Processor Sharing Systems                                 | 13        |
| 1.1.3 A Functional Central Limit Theorem for Erlang Loss Systems Under the SQ( $d$ ) Policy            | 19        |
| 1.2 Common Notation  | 24        |
| 1.3 An Introduction to Mean-field Techniques   | 26        |
| 1.3.1 An Overview of the Mean-field Techniques   | 26        |
| 1.3.2 Mean-field Analysis of the SQ( $d$ ) Policy  | 31        |
| 1.4 Outline  | 35        |
| <b>2 Insensitivity of the Mean-field Limit of Loss Systems Under SQ(<math>d</math>) Load Balancing</b> | <b>36</b> |
| 2.1 System Model   | 37        |
| 2.2 Additional Notation and Terminology  | 38        |
| 2.3 System Dynamics  | 40        |
| 2.4 Summary of Main Results  | 44        |



|          |   |           |
|----------|---|-----------|
| 2.4.1    | An Overview of the Analysis . . . . .   | 45        |
| 2.4.2    | Transient Regime . . . . .  | 46        |
| 2.4.3    | Stationary Regime . . . . .   | 49        |
| 2.5      | Proof of Theorem 2.3 . . . . .  | 51        |
| 2.6      | Proof of Theorem 2.5 . . . . .  | 55        |
| 2.7      | Numerical Results . . . . .   | 56        |
| 2.8      | Proofs of Main Results . . . . .  | 59        |
| 2.8.1    | Proof Theorem 2.2 . . . . .   | 60        |
| 2.8.2    | Proof of Lemma 2.2 . . . . .  | 69        |
| 2.8.3    | Proof of Theorem 2.1 . . . . .  | 74        |
| 2.8.4    | Proof of Theorem 2.4 . . . . .  | 76        |
| 2.8.5    | Proof of Lemma 2.4 . . . . .  | 77        |
| 2.8.6    | Proof of Lemma 2.3 . . . . .  | 78        |
| 2.9      | Conclusions . . . . .   | 80        |
| <b>3</b> | <b>On Occupancy Based Randomized Load Balancing for Large Processor Sharing Systems</b> | <b>81</b> |
| 3.1      | System Model and Routing Policy . . . . .   | 82        |
| 3.1.1    | System Model . . . . .  | 82        |
| 3.1.2    | Routing Policy . . . . .  | 83        |
| 3.2      | Additional Notation and Terminology . . . . .   | 85        |
| 3.3      | System Dynamics . . . . .   | 87        |
| 3.4      | Main Results . . . . .  | 88        |
| 3.5      | Numerical Results . . . . .   | 96        |
| 3.6      | On the Stationary Regime under the SQ( $\sigma$ ) Policy . . . . .                      | 99        |
| 3.6.1    | On the Propagation of Chaos in the Stationary Regime . . . . .                          | 105       |
| 3.6.2    | Discussion on a Related Work . . . . .  | 108       |

|          |   |            |
|----------|---|------------|
| 3.7      | Proofs of Main Results . . . . .  | 109        |
| 3.7.1    | Proof of Lemma 3.1 . . . . .  | 110        |
| 3.7.2    | Proof of Theorem 3.3 . . . . .  | 110        |
| 3.7.3    | Proof of Theorem 3.1 . . . . .  | 111        |
| 3.7.4    | Proof of Theorem 3.2 . . . . .  | 118        |
| 3.7.5    | Evolution of $(h_t; i; t = 0)$ for $\mathcal{L}C_b(U)$ . . . . .  | 128        |
| 3.7.6    | Proof of Lemma 3.2 . . . . .  | 133        |
| 3.7.7    | Single Server System with State-dependent Arrival Rates . . . . .   | 138        |
| 3.7.8    | Proof of Theorem 3.4 . . . . .  | 139        |
| 3.7.9    | Proof of Theorem 3.5 . . . . .  | 144        |
| 3.8      | Conclusions . . . . .   | 145        |
| <b>4</b> | <b>A Functional Central Limit Theorem for Multi-Server Erlang Loss Systems Under SQ(<math>d</math>) Load Balancing Policy</b> . . . . . | <b>146</b> |
| 4.1      | System Model . . . . .  | 147        |
| 4.2      | Additional Notation and Preliminary Results . . . . .   | 148        |
| 4.2.1    | Additional Notation . . . . .   | 148        |
| 4.2.2    | Preliminary Results . . . . .   | 149        |
| 4.3      | Summary of Main Results . . . . .   | 153        |
| 4.3.1    | Transient Regime . . . . .  | 154        |
| 4.3.2    | Stationary Regime . . . . .   | 156        |
| 4.4      | Proofs of Main Results . . . . .  | 159        |
| 4.4.1    | Proof of Theorem 4.2 . . . . .  | 159        |
| 4.4.2    | Proof of Lemma 4.1 . . . . .  | 161        |
| 4.4.3    | Proof of Theorem 4.3 . . . . .  | 162        |
| 4.4.4    | Proof of Lemma 4.3 . . . . .  | 164        |
| 4.4.5    | Proof of Theorem 4.6 . . . . .  | 166        |

|          |   |            |
|----------|---|------------|
| 4.4.6    | Proof of Theorem 4.7 . . . . .                    | 167        |
| 4.4.7    | Proof of Lemma 4.2 . . . . .                      | 168        |
| 4.5      | Conclusions . . . . .                             | 169        |
| <b>5</b> | <b>Summary and Future Research</b>                | <b>170</b> |
| 5.1      | Future Research . . . . .                         | 172        |
|          | <b>References</b>                                 | <b>174</b> |
|          | <b>A Some Key Background Results</b>              | <b>182</b> |
|          | <b>APPENDICES</b>                                 | <b>182</b> |
| A.1      | Weak Convergence and Prohorov's Theorem . . . . . | 183        |
| A.2      | Riesz Markov Kakutani Theorem . . . . .           | 184        |
| A.3      | Jakubowski's Criteria . . . . .                   | 184        |
| A.4      | Gronwall's Inequality . . . . .                   | 185        |
| A.5      | Doob's Inequality . . . . .                       | 186        |

# List of Tables

|     |   |     |
|-----|---|-----|
| 3.1 | The average number of probed servers per arrival . . . . .  | 97  |
| 3.2 | $\#_{tv}(\Gamma^{(\mathcal{N})}; \Gamma^{(exp)})$ for different JLDs for $\rho = 0.7$ and $\mathcal{N} = 300$ . . . . . | 100 |
| 3.3 | $\#_{tv}(\Gamma^{(\mathcal{N})}; \Gamma^{(exp)})$ for different JLDs for $\rho = 0.8$ and $\mathcal{N} = 300$ . . . . . | 100 |

# List of Figures

|     |  |     |
|-----|--|-----|
| 1.1 | A multi-server system with three servers . . . . .                                     | 2   |
| 1.2 | A multi-server Erlang loss model with $C = 3$ . . . . .                                | 3   |
| 1.3 | A multi-server PS model . . . . .  | 4   |
| 2.1 | Convergence of the mean-field to the fixed-point . . . . .                             | 59  |
| 2.2 | Comparison of the average blocking probability under $SQ(d)$ with lower bound. . . . . | 60  |
| 3.1 | The average response time versus . . . . .   | 98  |
| 3.2 | $d_{tv}(x(t; u); \cdot)$ versus $t$ . . . . .  | 103 |
| 3.3 | $\#_{tv}(y(t; v); \cdot)$ versus $t$ . . . . .   | 103 |
| 3.4 | $\#_{tv}(y(t; v); \cdot)$ versus $t$ . . . . .   | 104 |
| 3.5 | $\#_{ds}(y(t; v); \cdot)$ versus $t$ . . . . .   | 104 |
| 3.6 | The average response time ( $E(T_s)$ ) versus . . . . .                                | 105 |

# Chapter 1

## Introduction

Computer systems such as web-server farms, high-traffic websites, and cloud computing systems often contain hundreds of thousands of servers for meeting the incoming demands [1–3]. All these systems are examples of large-scale multi-server systems. The main challenge is to distribute the incoming jobs (requests) that generate tasks to servers so that the system achieves good system performance. Multi-server systems use job dispatchers to dispatch an incoming job to one of the servers. A basic multi-server model with three servers and one job dispatcher is shown in Figure 1.1.

Load balancing in multi-server systems is the process of balancing the load on servers so that individual servers are not overburdened that can lead to performance degradation. The job dispatchers dispatch an incoming job according to a load balancing policy that can improve the system performance. The main focus of this dissertation is the design and analysis of load balancing policies for multi-server systems with one central job dispatcher.

When a job arrives at the system, the best method is to dispatch the arrival to the least loaded or best server in the entire system as the job would get processed more efficiently. Since the job sizes are unknown in practice, the estimation and exchange of the remaining workload information is impractical and hence the policies that require only the knowledge of the occupancy or number of jobs of servers are preferred. If an incoming job is dispatched to a server with the least occupancy among all the servers, then the load balancing policy is referred to as the Join-the-Shortest-Queue (JSQ) policy [4–6]. The implementation of this policy requires information about the occupancy of all the servers for every arrival. Since many cloud computing systems and web-server farms often contain hundreds of thousands of servers, the implementation of the JSQ policy for large-scale systems results

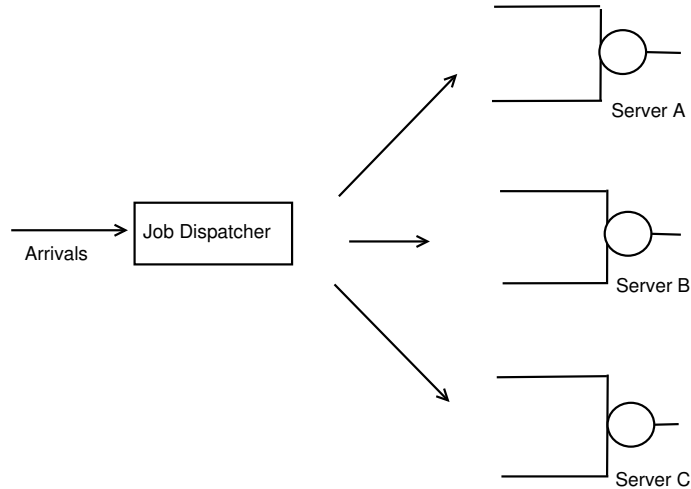


Figure 1.1: A multi-server system with three servers

in huge communication overhead and implementation cost. Hence, the JSQ policy is not preferred for large systems. An alternative type of policies that are easy to implement are randomized load balancing policies according to which the dispatcher first samples a few servers uniformly at random upon an arrival and one of them is chosen as the destination server for the arrival based on their occupancies. The criteria for selecting the destination server depends on the load balancing policy.

A popular randomized load balancing policy is the Shortest-Queue-Among  $d$  ( $SQ(d)$ ) or Power-of- $d$  policy according to which the dispatcher selects the server with the least occupancy among  $d$  randomly selected servers as the destination for an incoming job. This policy was introduced in [7] for multi-server systems with First-Come-First-Served (FCFS) servers when the job lengths have exponential distributions. The exact analysis of the system is not tractable due to complex interactions between servers as a result of the fact that the dispatcher uses the occupancy information of multiple servers to dispatch an arrival. Although servers are coupled, it was shown that the model is an example of a weakly interacting particle system. They used mean-field techniques to characterize the impact of the  $SQ(d)$  policy on the system performance. We give an introduction to mean-field techniques in Section 1.3.

The analysis of a randomized load balancing policy strongly depends on the type of servers that are present in the system. For example, servers could be FCFS, Erlang loss, or processor sharing (PS) servers. Each type of servers are used to study systems that serve jobs coming from different types of applications. For example, the FCFS model can be used to study supermarkets, hospitals, and path selection in networks. The loss models have

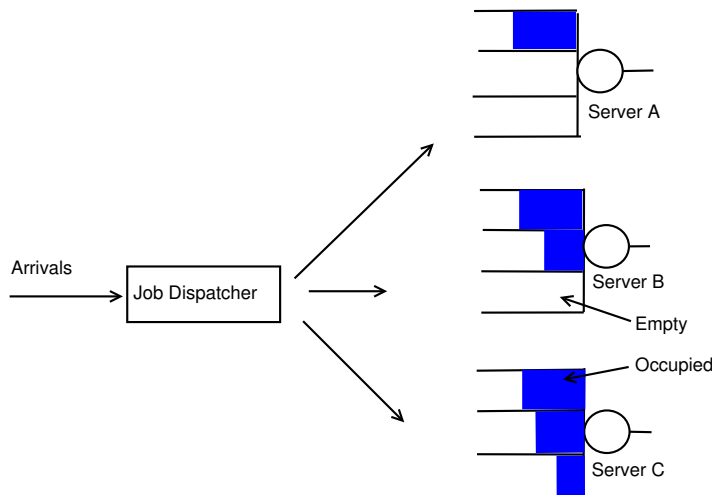


Figure 1.2: A multi-server Erlang loss model with  $C = 3$ .

applications in call centers and cloud computing systems. The PS servers have applications in web-server farms that provide service to delay-sensitive applications. Furthermore, the metric to measure the system performance also depends on the type of the system's servers. For FCFS, loss, and PS systems, the metrics of interest are the average delay, the average blocking probability, and the average response time (total execution time of a job) to measure the system performance, respectively. In this dissertation, we study occupancy based randomized load balancing for multi-server systems with Erlang loss servers and PS servers.

We now give a brief introduction to the architecture of the multi-server loss model that we study in this dissertation. Each loss server has capacity to serve only a finite number of jobs simultaneously, the maximum number of simultaneous jobs that a server can process is referred to as the server's capacity denoted by  $C$ . Furthermore, each job is processed at a constant unit rate irrespective of the occupancy of the server. A job arriving at a server will be accepted for service if its occupancy is less than the capacity  $C$ . Otherwise, the job is blocked from service and it is considered to be discarded from the system immediately. In Figure 1.2, we show a multi-server loss model with three servers each with capacity  $C = 3$ . In this figure, Server A, Server B, and Server C have occupancies one, two, and three, respectively. If an arrival is routed to Server C, then it will be blocked. On the other hand, if the destination for an arrival is either Server A or Server B, then it will be accepted. This example clearly shows the importance of a careful design of load balancing policies for efficient usage of resources so that the resulting average blocking probability is minimized.



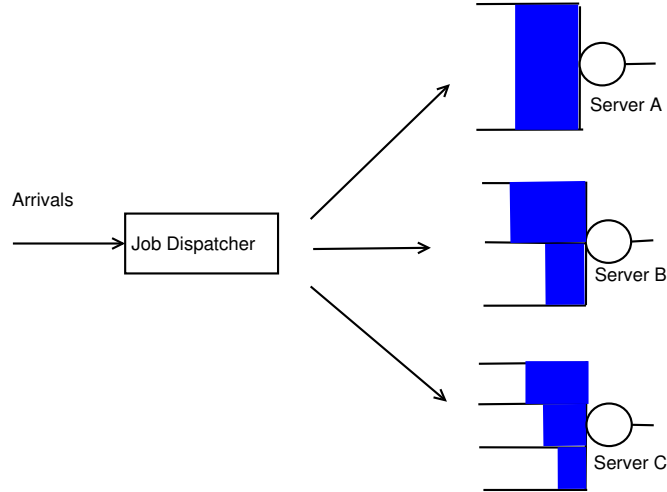


Figure 1.3: A multi-server PS model

Next, we give a brief introduction to the PS model. We assume that each server has capacity to serve jobs at a unit rate. A server with  $n$  jobs simultaneously processes all the  $n$  jobs at the rate of  $\frac{1}{n}$ . A job arriving at a server will be accepted for service and its processing begins immediately. In Figure 1.3, we show a PS model with three servers each having capacity to serve jobs at unit rate. The servers  $A$ ,  $B$ , and  $C$  have occupancies one, two, and three, respectively. Hence, the processing speeds of jobs at servers  $A$ ,  $B$ , and  $C$  are  $1$ ,  $\frac{1}{2}$ , and  $\frac{1}{3}$ , respectively. If an incoming job is routed to Server  $A$ , then it would get the highest possible instantaneous service rate in the system.

Our focus of interest in this dissertation is to study occupancy based randomized load balancing for the loss and PS models. In most prior works [8–12], the analysis was restricted to the case when job lengths have exponential distributions by using mean-field techniques. However, in many applications, the service time distributions are not exponential. For example, the service times follow log-normal distributions in call centers [13], and Gamma distributions in automatic teller machines (ATMs) [14] etc. Hence, it is of great interest to study the case of general job length distributions (JLDs). Our objective is to develop a mean-field analysis for this case.

## 1.1 Models

We now introduce the problems which we investigate in this dissertation. We study mean-field models in which servers are interacting particles, and the interactions between servers

take place when a job arrives at the system. The task of dispatching an incoming job to a server introduces dependence amongst the servers.

We study three problems. The first problem is on establishing the mean-field limit for large-scale multi-server loss models with general JLDs when the  $SQ(d)$  policy is used to dispatch the incoming jobs. This model is inspired by cloud computing systems which can be modeled as blocking models. We give an analysis of this problem in Chapter 2.

The second problem is on the mean-field analysis of randomized routing schemes for the PS model under the assumption of general JLDs. In this problem, we study not only the  $SQ(d)$  policy but also all occupancy based randomized routing schemes under a common framework. This model has applications in cloud computing systems and web-server farms that provide service to delay-sensitive applications. We provide an analysis of this problem in Chapter 3.

In Chapter 4, we investigate the third problem on the accuracy of mean-field approximations for multi-server loss models with the  $SQ(d)$  policy under the assumption of exponential JLDs. For this, we establish a suitable functional central limit theorem.

We next give a detailed description of all the three problems that we study in this dissertation. We also discuss the contributions of previous works and provide a summary of our contributions.

### 1.1.1 The $SQ(d)$ Policy for Erlang Loss Systems

Due to tremendous growth in the trend to externalize storage and computing resources, cloud computing systems maintain a large number of servers to provide an efficient service to jobs. In particular, in Infrastructure-as-a-service (IaaS) clouds such as Microsoft's Azure [2] and Amazon EC2 [1], the incoming job requests are mapped into virtual machines (VMs) that request resources such as processor power, I/O bandwidth, disk etc. from a server that is picked from a large set of distributed servers. When a job arrives, the incoming request is routed to one of the servers where it is accepted for service if the requested amount of resources are available. Otherwise, it is blocked from service. Then the job is discarded from the system immediately. The resources allocated to a job will be released once the service of a job ends. Every incoming job requests a certain amount of resources for usage from a server in the system and clearly, since servers have finite amount of resources, we assume that servers can process only a limited number of jobs simultaneously. Therefore,

the available resources in the system should be used as efficiently as possible. The system's performance is measured in terms of the average blocking probability of a job. To serve jobs efficiently, the service provider of a cloud computing system should design a mechanism to quickly identify a server that has enough resources for an incoming arrival. Since the job requests arrive randomly and their job lengths are random too, the main challenge is to quickly choose the server that is the least loaded or has the smallest number of jobs in the system as the destination server for an arrival.

Cloud computing systems are distributed systems in practice. Hence, they have job dispatchers to distribute incoming jobs to servers. The job dispatchers follow a predefined load balancing policy to distribute the incoming job requests to servers. We focus on the homogeneous model in which all the servers have the same amount of resources and also the incoming jobs request same amount of resources and have identical JLDs. Each accepted job is assumed to be processed at a constant unit rate. Furthermore, we assume that the system has a single central job dispatcher to distribute the incoming jobs to servers. A server in IaaS clouds with capacity  $C$  indicating the maximum number of jobs that can be served simultaneously at the server can be viewed as an Erlang loss server with capacity  $C$ . In particular, an IaaS cloud with  $N$  servers each with capacity  $C$  can be modeled by a multi-server Erlang loss system with  $N$  servers each with capacity  $C$ .

Since cloud computing systems often contain hundreds of thousands of servers, the JSQ policy for such large-scale systems suffers from high implementation cost. On the other hand, if an arrival is routed to a server that is picked uniformly at random, then the resulting blocking probability is much larger than the resulting blocking probability under the JSQ policy although the implementation cost is less. In the literature [8, 10, 15], the SQ( $d$ ) ( $d \geq 2$ ) policy was studied for loss models. This policy was shown to improve the system performance significantly over the random routing ( $d = 1$ ) policy even for the small value of  $d = 2$ . Although the SQ( $d$ ) policy performs worse than the JSQ policy, but it reduces the implementation cost drastically.

For systems with a finite number  $N$  of servers, characterizing the impact of the SQ( $d$ ) policy on each server thereby on the average blocking probability is a challenging task due to inter-dependencies between servers. However, under the assumption that the job arrival process is a Poisson process with rate  $N$ , due to weak interactions between servers, it was shown in previous works [8, 10, 15] by using mean-field techniques that a deterministic process referred to as the mean-field limit and its fixed-point approximate the distribution of a server's state of the system with large  $N$  in the transient and stationary regimes,

respectively. As a consequence, it can be shown that a function of the fixed-point of the mean-field approximates the average blocking probability.

## Objective and Key Challenges

Our objective is to develop the mean-field analysis to study the impact of the SQ( $d$ ) routing policy for a multi-server Erlang loss system with  $N$  servers each with capacity  $C$  under the assumption of general JLDs. Furthermore, we consider the case when jobs arrive according to a Poisson process with rate  $N$ . If a routed job gets accepted at its destination, then its processing begins immediately at a unit rate until its service is completed. This model was investigated in previous works under the assumption of exponential JLDs [8, 10]. Since in practice jobs have general JLDs, it is important to characterize the impact of the SQ( $d$ ) policy on each server and more importantly, on the average blocking probability. Motivated by this, we aim to obtain the mean-field limit and understand whether this can approximate the distribution of a server under the assumption of general JLDs. Moreover, we would like to investigate the impact of the routing policy on the average blocking probability when job lengths have general distributions.

To analyze complex stochastic systems, the starting point is often the mathematical modeling of the time evolution of the system as the time evolution of a Markov process. For the mean-field analysis, the first step is also to obtain a mathematical formulation of the time evolution of the system starting from an initial state. For exponential JLDs, due to the memoryless property of exponential distributions, it is sufficient to keep track of only the occupancy information of a server and tracking the age indicating the amount of time for which the job is in service or the residual service time is not required to construct a Markov process that models the system dynamics. In this simple case, each server's state is its occupancy. Due to symmetry of the system to servers' identities (as servers are homogeneous and the SQ( $d$ ) policy never uses the identities of servers), the process  $(\mathbf{X}^{(N)}(t); t \geq 0)$  where  $\mathbf{X}^{(N)}(t) = (\mathbf{X}_i^{(N)}(t); 0 \leq i \leq C)$  and  $\mathbf{X}_i^{(N)}(t)$  denotes the fraction of servers with at least  $i$  progressing jobs at time  $t$  out of  $N$  servers is a Markov process. Let  $\mathbf{S}_i^{(N)}(t)$  be the random variable denoting the state of the  $i^{\text{th}}$  server at time  $t$ , then

$$\mathbf{X}_i^{(N)}(t) = \frac{1}{N} \sum_{k=1}^N I_{\mathbf{S}_k^{(N)}(t) \geq i}.$$

The mean-field limit which is a deterministic process is the limit of the process  $(\mathbf{X}^{(N)}(t); t \geq 0)$  as  $N \rightarrow \infty$ . The proof of the existence of a mean-field limit mainly uses techniques on

the convergence of a sequence of Markov processes. In this case, the results in [8,10] imply that the following interchange holds,

$$\lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \mathbf{X}^{(N)}(t) = \lim_{t \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbf{X}^{(N)}(t) \quad (1.1)$$

where we consider the limit of a random element as its weak convergence limit. From (1.1), we obtain the equivalence between the stationary distribution of the limiting system given by the left-hand side and the globally stable fixed-point or equilibrium of the mean-field given by the right-hand side under the SQ( $d$ ) load balancing policy. The key is that the mean-field evolves according to a set of ordinary differential equations (ODEs) referred to as the mean-field equations (MFEs). We can study these equations to understand the behavior of the mean-field in stationary regime. Moreover, (1.1) can be used to show that any finite set of servers become statistically independent of each other [10] in the limit as  $N \rightarrow \infty$ . As a result, if  $\mathbf{x} = (x_n; 0 \leq n \leq c)$  is the unique fixed-point of the mean-field, then  $x_n$  is the probability that a server has at least  $n$  jobs and  $x_c$  is the blocking probability in the stationary regime as  $N \rightarrow \infty$ .

For general JLDs, due to the absence of the memoryless property, each progressing job's state must also be included in server's state. We can use either the age or residual service time of a job to track its status. In this dissertation, we use the ages of progressing jobs as their states. The state of a server must include the ages of all the jobs that are currently in progress at the server. At any time  $t$ , the state of a server for the case of general JLDs is of the form  $(n; x_1, \dots, x_n)$  where  $n$  denotes the number of progressing jobs and  $x_i$  denotes the age of the  $i^{\text{th}}$  progressing job at the server. Since the age of a job increases linearly with time and the state of a server lies in  $\mathbb{R}_+^n$  when it has  $n$  jobs, the classical methods used to study the exponential case cannot be used. We can overcome these challenges if we model the system evolution by a measure-valued Markov process, and then the mean-field limit can be established by using techniques on the convergence of a sequence of measure-valued Markov processes. We model the time evolution of the system as the time evolution of the measure-valued Markov process  $(\mu_t^{(N)}; t \geq 0)$  where

$$\mu_t^{(N)} = \frac{1}{N} \sum_{k=1}^N \delta_{\mathbf{s}_k^{(N)}(t)}$$

and  $\delta_a$  is the Dirac measure with unit mass at  $a$ . The result that  $(\mu_t^{(N)}; t \geq 0)$  is a Markov process follows from [16, Proposition 2.3.3]. Our main aim is to establish the limit of the process  $(\mu_t^{(N)}; t \geq 0)$  as  $N \rightarrow \infty$ . If the limit exists, then the limiting process is

referred to as the mean-field limit. We then plan to show that the obtained mean-field limit approximates the distribution of a server’s state in the transient regime.

To study the impact of the  $SQ(d)$  policy in the stationary regime, we first need to characterize the fixed-point of the mean-field, we then need to show that the stationary distribution of a server’s state can be approximated by the fixed-point of the mean-field if it is unique.

An important feature of Erlang loss systems for the case  $d = 1$  is that the system is insensitive. That is, the stationary distribution of a server’s occupancy is independent of the type of the JLD, but it depends only on the average job length. This implies the robustness of the system to variations in the JLDs. Hence, one can dimension the system by assuming the exponential JLD for the given average job length. Motivated by this, we would like to understand whether the stationary distribution of a server’s occupancy is insensitive or not in our model. In the system with  $N$  servers, the stationary distribution is not insensitive when  $d \geq 2$  due to violation of the Whittle balance condition [17] by the resulting routing probabilities which is a necessary condition to have insensitivity [18]. However, it was observed in simulation results of [10] that the stationary distribution of a server’s occupancy for different classes of JLDs coincide with that of the exponential case as  $N \rightarrow \infty$ . This implies that if the fixed-point of the mean-field coincides with the stationary distribution of a server’s state as  $N \rightarrow \infty$ , then the fixed-point of the mean-field should exhibit insensitivity. The proof of insensitivity of the fixed-point of the mean-field and the stationary distribution of a server’s occupancy when  $N \rightarrow \infty$  are still open problems. We investigate these issues by characterizing the fixed-point of the mean-field. A major challenge for our model is that the resulting MFEs are partial differential equations (PDEs). As a result, understanding the behavior of the mean-field as  $t \rightarrow \infty$  is much more challenging than the exponential case since in that case the MFEs are ODEs.

## Related Literature

The  $SQ(d)$  scheme was first introduced in [7] for multi-server systems with the First-Come-First-Served (FCFS) service discipline assuming  $d = 2$  and exponential service times. When the number of servers  $N$  is finite, analyzing the impact of the  $SQ(d)$  policy is a difficult task due to dependence amongst the servers introduced by the  $SQ(d)$  policy. However, they obtained an elegant tractable way of characterizing the stationary distributions in the limit as  $N \rightarrow \infty$  by using mean-field techniques. Their results were then extended for the

case of  $d > 2$  in [12] where it was argued that the case  $d = 2$  provides most of the gains and hence, the term ‘The power-of-2’ came to be used.

Multi-server loss models with randomized load balancing schemes were first studied in [15] when job lengths are exponentially distributed using a formal mean-field approach. However, the existence and uniqueness of the fixed-point of the mean-field were not shown. In [8, 10, 19], the existence and uniqueness of the fixed-point of the mean-field for the homogeneous loss model of [15] was addressed. For the heterogeneous loss model, the existence and uniqueness of the fixed-point of the mean-field were established in [8] under the assumption of asymptotic independence of servers in the stationary regime while [10] showed the asymptotic independence of servers and the interchange of limits (1.1). The propagation of chaos on path space implying that any finite set of random elements  $\{f(\mathbf{S}_i^{(N)}(t); t = 0)g_{1 \dots i N}\}$  become independent of each other as  $N \rightarrow \infty$  was established by [20, 21] in the context of alternate routing in circuit-switched networks. The notion of propagation of chaos is further explained in Definition 1.1.

The impact of the SQ( $d$ ) policy for the loss models similar to the one considered here was analyzed in [8, 10, 15, 19] under the assumption of exponential JLDs. They also considered the more general heterogeneous case with an appropriate modification of the SQ( $d$ ) policy to account for server and job heterogeneity. It was shown in [19] that the SQ( $d$ ) load balancing scheme yields almost optimal blocking performance in that the average blocking is very close to the theoretical lower bound on the minimum average blocking achievable by any load balancing policy.

It is well known that the stationary distributions of single server loss systems even with state-dependent Poisson arrival rates are insensitive to the service time distribution, i.e., they only depend on the mean of the service times [22]. Hence, it is essential to investigate whether the insensitivity property carries over to systems with randomized load balancing such as the SQ( $d$ ) policy. When  $N$  is finite, randomized load balancing policies result in the individual servers being coupled. It can be shown as in [18, 23] that when  $N$  is finite, the system is not insensitive since the SQ( $d$ ) policy does not satisfy the necessary condition of state-dependent arrival rates to be balanced. However, the insensitivity of the fixed-point or equilibrium of the mean-field was observed for the limiting case (i.e., when  $N \rightarrow \infty$ ) via simulations in [8, 10] but no proofs were provided. This requires a proof in our context because as  $N \rightarrow \infty$ , the resulting MFEs represent the time evolution of the distribution of a nonlinear Markov process that models a queueing system with a single server in which the arrival rates of jobs depend both on the instantaneous state and the

instantaneous distribution of the state of the server whereas in the classical models, the underlying Markov process is linear since the arrival rates of jobs are independent of the instantaneous distribution of the server's state. The main goal of the chapter is to address this issue.

For the case of general JLDs, Markovian modeling of the system requires us to keep track of the age or residual service time of each job that is in progress in the system. Therefore, the underlying space on which the Markov process is defined is not discrete, and hence the classical Markov chain techniques cannot be used. As a result, establishing the mean-field limit and characterizing the properties of its stationary behavior for general JLDs is a challenging task.

In the literature, some prior works discussed later in this section have studied queueing models in which either jobs or servers are interacting particles under the assumption of general JLDs by using mean-field techniques. Also, some works have investigated the randomized load balancing for multi-server systems with general JLDs but the analysis is not complete. There are still several open problems in the literature that need to be addressed. We now discuss contributions and drawbacks of some related works for the case of general JLDs below.

In [24], randomized load balancing schemes for queueing systems with general service time distributions were investigated when servers use the service disciplines FCFS, PS, and Last-In-First-Out (LIFO). The steady-state results were characterized by assuming the asymptotic independence of servers in the system when  $N \rightarrow \infty$ . However, the proof of asymptotic independence of servers remains an open problem and also, the mean-field limit and its fixed-point were not studied in any detail.

In [25], mean-field techniques were used to study a closed queueing network with  $M$  jobs and  $N$  queues under the FCFS service discipline. In this model, a job once served at a queue joins another queue chosen with probability  $\frac{1}{N}$  from  $N$  queues. The mean-field limit was established for the regime when  $\lim_{M, N \rightarrow \infty} \frac{M}{N} = \rho$ . However, the equilibrium behavior of the system was not studied. Recently, the SQ( $d$ ) setting in a system of  $N$  servers with FCFS service discipline where jobs arrive according to a time-inhomogeneous Poisson process and general i.i.d. service times were studied in [26]. They studied only the transient behavior of the mean-field. The steady-state analysis of the model was not investigated.

The mean-field and fluid analysis of queues are closely related, the former usually in



the space of probability measures and the latter on the space of finite measures. The fluid limit analysis of the FCFS and the PS queues with general service time distributions has been studied using a measure-valued processes approach developed by Dawson [16] in [27–31]. In this chapter, we use the age of a progressing job, indicating the amount of time elapsed since its arrival to construct a measure-valued Markov process that models the system dynamics. We use this framework to establish the mean-field limit by studying the limit of the empirical measure-valued process when the number of servers  $N \rightarrow \infty$ , as in [16, 29]. Our proof techniques closely follow those in [29] in which a fluid limit and a functional central limit theorem for a sequence of M/GI/1 systems were established by using measure-valued processes.

In [32], the FCFS model was studied with exponential distributions under the SQ( $d$ ) policy by using measure-valued processes. In the exponential case, the set of states of servers is  $Z_+$ . In [32], the law of large numbers on path space is established by studying the limit of the sequence of empirical measures with samples in  $M_1(D_{Z_+}([0; 1]))$  where  $D_{Z_+}([0; 1])$  is the space of right continuous with left limits (RCLL) functions in  $Z_+$  and  $M_1(D_{Z_+}([0; 1]))$  is the space of probability measures on  $D_{Z_+}([0; 1])$ . For our analysis, we study stochastic processes with sample paths in  $D_{M_1(U)}([0; 1])$  where  $U$  is the set of states of servers which is uncountable.

## Summary of Contributions

We obtained the following results for the two regimes of interest. Preliminary results were presented in [33] and the detailed results will appear in [34].

*For the Transient Regime:* We first establish the mean-field limit by using measure-valued processes. We then use this result to show the asymptotic independence of any finite set of servers at a given finite time  $t$  as in [26] under the exchangeability assumptions of initial servers' states. The MFEs represent dynamics of a classical single server system with Poisson job arrival process whose intensity depends both on the occupancy of the server and the distribution at any time  $t$ . From the initial form of the MFEs whose solution is a measure-valued process, we show that under certain assumptions, the mean-field satisfies a set of PDEs.

*For the Stationary Regime:* We show the existence of a unique fixed-point of the mean-

field. The fixed-point of the mean-field satisfies a product form given in (2.21) and the occupancy distribution obtained from the fixed-point of the mean-field coincides with the fixed-point in the exponential case having the same average job length. This concludes insensitivity of the fixed-point of the mean-field. We also provide numerical evidence suggesting that the fixed-point is globally asymptotically stable (GAS) as  $t \rightarrow \infty$  when job lengths have mixed-Erlang distributions.

### 1.1.2 Occupancy Based Randomized Routing for Processor Sharing Systems

Due to the emergence of the cloud computing paradigm and other applications, the server farms often contain large numbers of servers to process the incoming jobs. Unlike the model discussed in Section 1.1.1 which is a blocking model, we now consider a model in which servers provide service to delay-sensitive applications such as online-search (Google), Social networking (Facebook) etc. The main objective in this model is to minimize the average response time experienced by a job in the system. Furthermore, each server has an infinite buffer, and there is no blocking of jobs in the system. A server always accepts an incoming job.

The front end job dispatchers in web-server farms route an arriving job to one of the servers that provide minimal response times as the tasks in most cases are delay-sensitive. Therefore, the main challenge in these systems is to design low complexity load balancing algorithms that yield efficient usage of available resources, thereby good system performance. The system performance is measured in terms of the average response time experienced by a job in the system. In server farms, the resources are shared by processing requests in a round-robin manner with small time granularity. This model can be well approximated by a model with servers having the PS service discipline [4, 35, 36] in which the processing speed of a server is equally shared by all the progressing jobs. In practice, the service requirements of jobs are highly variable. Under the PS service discipline, as the processing speed of a server is shared by all the jobs, the short jobs will be processed quickly and the PS service discipline avoids the chance of getting waited for a long time by small jobs before their service begins. As a result, the average response time is minimized. Hence, the PS service discipline is preferred for delay-sensitive applications.

Traditionally, small-scale server farms such as Cisco Local Director, IBM Network Dispatcher, Microsoft Sharepoint use the classical Join-the-Shortest-Queue (JSQ) load balanc-

ing policy [4, 8]. It was shown in [4–6] that the JSQ policy is nearly optimal and further, it is robust to the JLDs since it is almost insensitive. For large scale systems that contain hundreds of thousands of servers, the JSQ policy requires information about the number of progressing jobs at all the servers. However, it was shown in [7, 9, 10, 12, 24, 37] that the SQ( $d$ ) policy can achieve almost the same gains at a much smaller sampling cost.

## Objective and Key Challenges

We consider a multi-server system with  $N$  servers which employ the PS service discipline. Each server has capacity to process jobs at a unit rate and hence, a server with  $n$  jobs processes every job at rate  $\frac{1}{n}$  because of the PS service discipline assumption. The job arrival process is a Poisson process with rate  $N$  and jobs have general JLDs with finite average job length. We study occupancy based randomized routing schemes which imply sampling of a finite number of servers upon an arrival and the dispatcher selects a server from this set as the destination server. Any occupancy based policy can be viewed as a mapping of occupancies of sampled servers to routing probabilities. For example, let us assume  $d$  servers are sampled upon an arrival and let  $n_i$  be the occupancy of the  $i^{\text{th}}$  sampled server. Furthermore, let  $(n_1; \dots; n_d)$  be the vector of occupancies of sampled servers. Then the routing decision according to a load balancing policy can be viewed as finding a probability vector  $(p_1; \dots; p_d)$  based on  $(n_1; \dots; n_d)$  where  $p_i$  denotes the probability with which the arrival should be routed to the  $i^{\text{th}}$  sampled server. Based on this logic, we aim to develop a mean-field analysis under a common framework to include all occupancy based routing policies which require sampling of a finite number of servers upon an arrival. The study of occupancy based routing policies for the loss model easily follows from the analysis developed for the PS model. This is because each server in the loss model can process at most  $C$  jobs simultaneously each with constant unit processing rate and hence, it can be seen that the analysis of the loss model follows easily from the analysis of the PS model.

In the literature, the complete mean-field analysis for the model of interest exists only for the SQ( $d$ ) policy [7, 11] and for some threshold based policies [12] under the assumption of exponential JLDs but the analysis for general JLDs remains an open problem. The mean-field analysis of the SQ( $d$ ) policy was developed in [7] under the assumption of Poisson arrival process of jobs with rate  $N$  and jobs have exponential JLDs for the case  $d = 2$ . Later in [12], the analysis was extended to the case  $d > 2$ . Unlike the Erlang loss model

considered in Section 1.1.1 that is stable iff the average job length is finite, if  $\rho < 1$  is the average job length, then the PS model requires  $\rho < 1$  for the stability of the system [38]. In [7, 12], under the assumption of  $\rho < 1$ , they used mean-field techniques to understand the impact of the SQ( $d$ ) policy on large-scale system behavior. In this approach, they first characterized the system's dynamics with time by studying the evolution of the empirical process  $(\mathbf{X}^{(N)}(t); t \geq 0)$  where  $\mathbf{X}^{(N)}(t) = (\mathbf{X}_i^{(N)}(t); i = 0)$  and  $\mathbf{X}_i^{(N)}(t)$  denotes the fraction of servers with at least  $i$  progressing jobs at time  $t$ . Then they showed that under certain initial conditions, the empirical process  $(\mathbf{X}^{(N)}(t); t \geq 0)$  converges to a deterministic process referred to as the mean-field. The mean-field is a unique solution of a set MFEs which are ODEs. Moreover, they showed that the interchange of limits (1.1) holds. As a result, the mean-field and its fixed-point can be used to approximate the transient and stationary distribution of a server when  $N$  is large. In fact, for  $d \geq 2$ , they found the fixed-point of the mean-field  $\mathbf{x} = (x_i; i = 0)$  given by  $x_i = \left(-\right)^{\frac{d^i - 1}{d - 1}}$ . Then the average delay of a job in the system when the number of servers  $N$  is large is nearly  $\sum_{i=0}^{\infty} x_i \left(-\right)^{\frac{d^i - 1}{d - 1}}$ . For the random routing case of  $d = 1$ , we have  $x_i = (-)^i$ . This clearly shows a significant improvement in the system performance since  $x_i$  decreases double exponentially with  $i$  for  $d \geq 2$  whereas  $x_i$  decreases exponentially with  $i$  for the case  $d = 1$ .

Since the JLDs are not exponential in practice, it is crucial to understand the impact of the SQ( $d$ ) policy for the PS model with general JLDs. In the literature, although there are some previous works [39], the complete analysis remains an open problem. There is another reason why the PS model is interesting: it is known to be insensitive if the arrival process to each individual server is a Poisson process, and so showing insensitivity of the mean-field under very general occupancy based randomized routing can then lead us to conclude that insensitivity is a generic property of the PS service discipline rather than the precise input to a queue that is controlled by the routing policy. Also, by using our analysis, one can obtain the mean-field limit immediately for the policies that fall into our framework and thereby, the considered policy's impact can be studied by using the obtained mean-field limit. It is known that the random routing SQ( $d$ ) with  $d = 1$  is insensitive. We need to understand whether this holds for the SQ( $d$ ) policy with  $d \geq 2$  and also, for other occupancy based routing policies. The system with finite  $N$  under the SQ( $d$ ) policy is not insensitive as the arrival rate of jobs do not satisfy the Whittle balance condition [17], which is a necessary condition [18, 23]. However, when  $N \rightarrow \infty$ , it was observed through simulations that the insensitivity property holds for the SQ( $d$ )

policy [11]. For other routing policies, insensitivity was not studied.

We now discuss why it is essential to study occupancy based randomized routing policies in addition to the SQ( $d$ ) policy. Upon an arrival, if the dispatcher samples a fixed number of  $d$  servers to decide the destination server, then the SQ( $d$ ) policy is optimal [5, 6]. Now let us consider the following example. For the case  $d = 2$ , if the first server sampled upon an arrival is found to be idle, then the dispatcher can stop further sampling of servers to decide the destination. For some systems, limiting the amount of information exchanged to dispatch an arrival is important. Hence, while designing a routing policy, we need to consider both its performance and complexity. Furthermore, a routing policy should provide the system performance that is robust to variations in the system congestion level. For example, when the dispatcher uses the SQ( $d$ ) policy, if the system is heavily loaded, then it could result in bad system performance when  $d$  is fixed at a small value. On the other hand, if the system is lightly loaded, then the dispatcher might unnecessarily sample more servers when  $d$  is fixed at a large value. Since variations in the system congestion level are common in practice, the routing policy should be able to adapt to the system congestion level. Our framework should also include policies that are robust to variations in the system traffic.

Our objective is to study the mean-field of the empirical distributions for policies that fall into our framework introduced in Section 3.1.2. This includes policies that are adaptive to the system congestion level and threshold based policies. We provide results for four occupancy based policies defined in Section 3.4 that include the SQ( $d$ ) policy, and the remaining three are generalized threshold based policies. Such policies have been studied in other contexts. For example, in [12], a threshold based policy was studied in which for an incoming arrival, the dispatcher randomly selects a server as the destination server if its occupancy is less than or equal to the chosen threshold value. Otherwise, the dispatcher samples another server, and the same procedure is repeated until the dispatcher samples  $d$  servers. If all the  $d$  sampled servers have occupancy higher than the threshold value, then the shortest queue is chosen as the destination. This policy has lower complexity than the SQ( $d$ ) policy at the cost of a small drop in the performance. We investigate this policy in this dissertation, and we show that the fixed-point of the mean-field is insensitive, which is also confirmed by our simulation results. In [12], insensitivity of the threshold based policy was not studied.

In [24], a  $d$ -adaptive algorithm was studied using simulations to overcome the bad performance of the SQ( $d$ ) policy for small values of  $d$  for the case of Power-law distributions

in First-In-First-Out (FIFO) models. According to this policy, if the occupancy of the first selected server is  $k$ , then the dispatcher selects further  $f(k)$  servers and then chooses the server with the least occupancy as the destination. They observed through simulations that for  $f(k) = k$ , the  $d$ -adaptive algorithm samples with high probability three servers per arrival providing the system performance that is similar to the SQ(4) scheme. This policy also belongs to the class of occupancy based policies that we study in this dissertation. Our simulation results show that the system exhibits insensitivity under this scheme when  $N \gg 1$ . In [24], insensitivity of the  $d$ -adaptive policy was not discussed.

For queueing systems with general service times, we need to use the occupancy state as well as the age or residual service times of the progressing jobs in service to model the system's evolution by a Markov process. For the FCFS models, this results in a bi-dimensional Markov process defined on  $Z_+ \times \mathbb{R}_+$ . On the other hand, for the PS model, the Markov process is defined on  $\bigcup_{n \in Z_+} (n; \mathbb{R}_+^n)$  as we need to track both occupancy and age of each progressing job. The analysis is much more difficult because the rate at which a job is served at a server depends on the server's occupancy. Similar to the loss model, we use measure-valued Markov processes to establish the mean-field limit under certain assumptions on initial conditions.

## Related Literature

In [24], the SQ( $d$ ) policy was investigated for the PS, FCFS, and LIFO systems with general service time distributions under the assumption of asymptotic independence of any finite set of servers and the existence of a unique limiting stationary distribution as stated in an *ansatz*. However, the proofs of the assumptions considered in the *ansatz* are still unknown. The asymptotic independence of finite set of servers has been shown for the PS and the Erlang loss models in [9, 10, 20] for the case of exponential distributions.

The assumption of asymptotic independence of servers for the PS model for the case of general JLDs renders the arrival process to each server as a state-dependent Poisson process. Then the resulting model satisfies the Whittle's condition [17], which implies insensitivity. A proof of the *ansatz* remains an open problem except for the FCFS models when JLDs have decreasing hazard rate functions [40] that give rise to monotonicity in the model. From [40], the proof techniques cannot be extended to the PS model as we do not have ordering between the server states, unlike the FCFS systems.

A mean-field analysis to characterize the transient and stationary distributions has been

extended to heterogeneous systems in which servers are classified into different classes based on their capacity to process jobs but with exponential service distributions. This was done in a series of works [9, 11, 37].

Recently, the SQ( $d$ ) policy for FCFS systems with generally distributed service times was studied in [26] by using a mean-field approach. However, the analysis was restricted to the transient case, and no results were given on the stationary regime. Another important parameter of a policy is its complexity. Hence, it is of interest to generalize the mean-field analysis for occupancy based policies where the number of samples is adaptive to the system congestion level, which is what we also consider in this chapter.

In [27, 28], measure-valued processes have been used to study PS models with general JLDs in a different context. The fluid limit was obtained in [27] for the GI/GI/1 system with a single PS queue by using measure-valued processes. They used residual service times of jobs to construct measure-valued processes. The analysis was extended to the case when the customers are impatient in [28].

## Summary of Contributions

Preliminary results were presented in [41, 42]. Our contributions are the following:

*For the Transient Regime:* The main contribution of our analysis is to show the existence of a mean-field limit for occupancy based randomized routing policies under the assumption of general JLDs. The MFEs are now PDEs, unlike the exponential case that results in ODEs. We also show asymptotic independence of servers for any given finite time. We observe that the MFEs represent dynamics of a classical single server system with Poisson job arrival process having intensity that depends on both the instantaneous occupancy of the server and the instantaneous distribution of the Markov process that models the system's dynamics. Any load balancing policy influences only the intensity of the job arrival process.

*For the Stationary Regime:* We show that every fixed-point of the PDEs that describe the mean-field limit corresponds to a probability measure on  $Z_+$ , is also a fixed point of the MFEs under exponential distributions having the same average job length. The fixed point is unique when the corresponding property can be shown in the exponential case. Therefore, insensitivity of the fixed-point holds if the exponential case has a unique fixed-



point. Finally, we present four occupancy based policies that have been discussed in the literature. We obtain the MFEs for these policies from our generalized mean-field analysis of occupancy based routing schemes. For the threshold based policies, Policy 2 and Policy 4 defined in Section 3.4, when job lengths are exponentially distributed, we show that the mean-field has a globally stable fixed-point. In these cases, it follows that the fixed point is the stationary distribution of the individual servers when  $N \rightarrow \infty$ . We also provide simulation results providing evidence for insensitivity of the stationary distribution of a server as  $N \rightarrow \infty$  for all the considered policies.

For the SQ( $d$ ) policy, under the assumption of mixed-Erlang JLDs, we numerically solve the MFEs. We observe that the mean-field is not quasi-monotonic while it is in the exponential case. Furthermore, if  $\mu^{(N)}$  is the stationary distribution of the empirical process and  $Z$  is a limit point of a subsequence  $\mu^{(N_k)}$ , then in the limiting model as  $N \rightarrow \infty$  we show that any finite set of servers are independent of each other if and only if  $Z = \delta_x$ , where  $x$  is the unique fixed-point of the mean-field. If  $Z$  is not a Dirac measure, then servers are coupled through the position of the mean-field where the initial point of the mean-field has distribution  $Z$ . This result is also true for other mean-field models.

### 1.1.3 A Functional Central Limit Theorem for Erlang Loss Systems Under the SQ( $d$ ) Policy

In the previous two problems stated in Sections 1.1.1 and 1.1.2, the objective was to study the mean-field and its fixed-point as an approximation of the transient and stationary distribution of a server's state when  $N$  is large. We need to characterize the resulting approximation error as a function of  $N$ . We focus on loss models with exponential JLDs because these results require a complete mean-field analysis that is lacking when general JLDs are assumed. The results that are obtained for the loss model case are suggestive of the more general type of results that can be obtained for other models.

#### Objective and Key Challenges

We consider large-scale multi-server loss systems with  $N$  servers where each server has capacity to process at most  $C$  jobs simultaneously at unit rate. We assume that the service times are exponentially distributed with unit mean. The central job dispatcher dispatches the incoming jobs according to the SQ( $d$ ) policy. For this model under the



assumption that the jobs arrive according to a Poisson process with rate  $N$ , from the previous works [8,10], the mean-field analysis shows that  $\lim_{N \rightarrow \infty} \mathbf{X}^{(N)}(t) = \mathbf{x}(t; \mathbf{u})$ ; where  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$  is the mean-field limit with initial point  $\mathbf{u}$ . Now the key question is, how do we characterize the stochastic fluctuation process  $(\mathbf{X}^{(N)}(t) - \mathbf{x}(t; \mathbf{u}); t \geq 0)$ ? Moreover, if  $P_{block}^{(N)}$  is the actual blocking probability in the system with  $N$  servers, can we quantify the gap between  $P_{block}^{(N)} - \frac{d}{C}$  as a function of  $N$ ?

Similarly, we study the Halfin-Whitt regime by assuming that the arrival process is a Poisson process with rate  $N \rho^{(N)}$  where  $\rho^{(N)}$  is a function of  $N$  with  $\lim_{N \rightarrow \infty} \rho^{(N)} \frac{N}{C} (1 - \frac{\rho^{(N)}}{C}) = b$  for some  $b \in 0$ . In this case, how do we obtain approximations to  $P_{block}^{(N)}$  when  $N$  is large? In this case, the system is heavily loaded when  $N$  is large. One way is first to obtain the asymptotic results when  $N \rightarrow \infty$ , and then estimate the error terms. The interest in the Halfin-Whitt regime is because, for  $d = 1$ , there is a phase change in the behavior of the blocking probabilities going from exponential (in  $C$ ) decrease to a  $\frac{1}{C}$  scaling as  $C$  becomes large (see [43], [44]). As explained later in this section, we also have the similar result for the complete resource pooling case, i.e., a loss system with a single server having arrival rate of  $N \rho^{(N)}$  with capacity  $NC$ , then the blocking probabilities satisfy  $\frac{1}{N}$  scaling as  $N$  becomes large. Although loss systems are stable for any finite average load, this phase change shows that the critically loaded regime is important for appropriate dimensioning of systems because of the change in the blocking sensitivity. Our objective is to investigate the existence of such a phenomenon for loss systems under the SQ( $d$ ) load balancing for  $d \geq 2$ . The main challenge is that we do not know  $P_{block}^{(N)}$  while it is computed using the Erlang-B formula for the case  $d = 1$ .

We address the two cases  $\rho^{(N)} = \frac{C}{N}$  and the Halfin-Whitt regime as special cases of our framework in which we assume that

$$\rho^{(N)} = \frac{C}{N} + \frac{\tilde{\rho}}{N} \tag{1.2}$$

for  $\tilde{\rho} \geq \mathbb{R}_+$  and  $\tilde{\rho} \geq \mathbb{R}$ . Here, if  $\tilde{\rho} = 0$ , then we get the first case in which the arrival process rate is  $N$ . On the other hand, if  $\tilde{\rho} = C$  and  $\tilde{\rho} \in 0$ , then we get the Halfin-Whitt regime.

## Related Literature

So far, in the literature [8,10,19], only the mean-field analysis was investigated for  $\rho^{(N)} = \frac{C}{N}$ . An important problem which was not studied in the literature is the quantification of the

error between the actual blocking probability for a system with large  $N$  and the asymptotic blocking probability, as a function of  $N$ .

We now discuss related existing works and provide some insights into the system performance in the Halfin-Whitt regime. The average blocking probability depends on how efficiently we use the system resources. For example, let us consider the random routing case where an arrival is routed to a randomly selected server in the Halfin-Whitt regime. Then the average blocking probability experienced by an arrival is the same as in the single server loss system with capacity  $C$  where the jobs arrive at a Poisson process with intensity  $N$ . In this case, since the arrival process to each server is a Poisson process, due to reversibility, the average blocking probability is equal to  $Er(N; C)$ , where  $Er(a; n)$  denotes the Erlang-B formula for Poisson arrivals with rate  $a$  and server's capacity  $n$ . Then since  $C$  is fixed and  $Er(N; C)$  is a continuous function of  $N$ , the average blocking probability converges to  $Er(C; C)$  when  $N \rightarrow \infty$ . On the other hand, if we consider the complete resource pooling case that corresponds to the single server system with Poisson job arrival process having rate  $N$  and server capacity  $NC$ , then the average blocking probability is given by  $Er(N; NC)$ . Then from [45],

$$\lim_{N \rightarrow \infty} \frac{1}{N} P_{block}^{(N)}(Er(N; NC)) = \frac{\phi_{(normal)}(\tilde{c})}{C \Phi_{(normal)}(\tilde{c})}; \quad (1.3)$$

where  $\phi_{(normal)}(\cdot)$  and  $\Phi_{(normal)}(\cdot)$  denote the density and distribution functions of the standard normal distribution, respectively.

It was shown in [46] that under the JSQ scheme, we obtain the same result (1.3) as in the case of complete resource pooling, i.e.,

$$\lim_{N \rightarrow \infty} \frac{1}{N} P_{block}^{(N)} = \frac{\phi_{(normal)}(\tilde{c})}{C \Phi_{(normal)}(\tilde{c})}; \quad (1.4)$$

This is expected since an arrival will not be blocked from service when there is an empty spot in the system similar to the case of complete resource pooling. As a result, the average blocking probability under the JSQ policy is equal to  $Er(N; NC)$ . This result was not explicitly stated in [46] for the system with  $N$  servers. Under the influence of the SQ( $d$ ) policy, an arrival could be blocked from service even if there is an empty spot in the system. Therefore, we expect a decrease in the system utilization when we use the SQ( $d$ ) policy, but it has less computational cost than the JSQ policy. However, it was shown in [46] that if  $d$  is a function of  $N$  denoted by  $d^{(N)}$ , and if  $\lim_{N \rightarrow \infty} \frac{d^{(N)}}{N \log(N)} = 1$ , then we still obtain (1.4) for the SQ( $d^{(N)}$ ) scheme.

Thus, we observe that even if we scale up the system capacity  $NC$  with  $N$  when the arrival rate of jobs scales as  $N^{-1}$ , if we do not assign available resources to job requests cleverly, the system performance may not be optimal. This is because for a fixed  $d$ , an arrival could be blocked from service even if there is an empty spot in the system under the SQ( $d$ ) scheme. We would like to understand the impact of the SQ( $d$ ) policy on blocking probabilities in the Halfin-Whitt regime.

The key challenge in our model with  $\rho = \frac{\lambda}{\mu} \approx \frac{\rho}{N}$  is that we cannot compute  $P_{block}^{(N)}$  using the Erlang formula. Moreover, we cannot obtain an explicit expression for  $P_{block}^{(N)}$  as the system is non-reversible and non-tractable.

In our analysis, we first establish a functional central limit theorem (FCLT) that characterizes the fluctuations of the empirical process around the mean-field of our model and then use this result to obtain efficient approximations for  $P_{block}^{(N)}$  when  $N$  is large. The limiting diffusion scaled fluctuation process is an Ornstein-Uhlenbeck (OU) process that is characterized by the mean-field limit. We then exploit the obtained FCLT to quantify the error between  $P_{block}^{(N)}$  and the asymptotic blocking probability  $\frac{d}{C}$ , where  $\rho = (\lambda, \mu; 0, \rho)$  is the fixed-point of the corresponding mean-field and  $\rho$  is the probability that a server is fully occupied when  $N \rightarrow \infty$ . Our results conclude that  $\lim_{N \rightarrow \infty} \frac{1}{N} \log \frac{P_{block}^{(N)}}{\frac{d}{C}}$  goes to a limit that can be explicitly characterized in terms of  $\rho$ ,  $\lambda$ ,  $\mu$ , and  $C$ . More importantly, we show that  $\lim_{N \rightarrow \infty} \frac{1}{N} \log \frac{P_{block}^{(N)}}{\frac{d}{C}}$  is non-zero if and only if  $\rho$  is non-zero. Graham [47] established a similar FCLT for the FCFS queueing model under the SQ( $d$ ) policy but did not exploit the FCLT result further to characterize the system performance. An FCLT approach was also used by Hunt [48] to analyze large symmetric star loss networks where asymptotic independence between nodes was established and also, the error between the actual blocking probabilities and approximate blocking probabilities was quantified.

A by-product of our result is to show that the empirical occupancy process converges to the fixed point of the MFE at a rate of  $O(\frac{1}{N})$ , a result that was also shown by Ying [49] using Stein's method and more recently by Gast [50] where a refined  $O(\frac{1}{N})$  approximation was also obtained by using Stein's method under the assumption that the mean-field is exponentially stable. More importantly, in [49, 50], the focus was on quantifying the mean-square error between  $\mathbf{X}^{(N)}(t)$  and  $\mathbf{x}(t; \mathbf{u})$  for both the transient and stationary regimes but they did not study the fluctuation process  $(\mathbf{X}^{(N)}(t) - \mathbf{x}(t; \mathbf{u}); t \geq 0)$ .

In [51], it was shown that even the the gap between  $\mathbb{E}[f(\mathbf{X}^{(N)}(t))]$  and  $f(\mathbf{x}(t; \mathbf{u}))$  is  $O(\frac{1}{N})$  for both the transient and stationary regimes where  $f$  is a twice differentiable performance function. As a result, it was concluded that for the FCFS model, the gap

between the average waiting time of the system with  $N$  servers and the asymptotic average waiting time computed as function of the fixed-point of the mean-field is  $O(\frac{1}{N})$ . The result of [51] can be applied to loss models to conclude that  $P_{block}^{(N)} = \frac{d}{c} = O(\frac{1}{N})$ . However, the analysis of [51] assumes  $\tilde{\rho} = 0$ .

Recently, Eschenfeldt and Gamarnik [52] studied the CLT scaling of the queue occupancy process for a system of  $M=M=1$  queues with the JSQ policy in the Halfin-Whitt where they showed that asymptotically the distribution concentrates on queues having up to two jobs. They did not study the stationary regime. In [53], it was shown that for the system with  $M=M=1$  queues, under the SQ( $d^{(N)}$ ) policy where  $\frac{d^{(N)}}{N \log(N)} \rightarrow 1$ , the limiting diffusion process is the same as in the case of the JSQ policy in the Halfin-Whitt regime. Their analysis was also restricted to the transient regime.

In [54], for the system with  $M=M=1$  queues and Poisson arrival process of jobs with rate  $N^{1-\tilde{\rho}}$ , the impact of the SQ( $d$ ) policy was studied under the assumption that there exists a sequence  $f_N g_{N-1}$  such that  $\lim_{N \rightarrow \infty} f_N = 1$  and  $\lim_{N \rightarrow \infty} (1 - f_N) = \tilde{\rho}$  where  $\tilde{\rho} > 0$ . They showed that in the stationary regime, majority of queues have queue lengths at least  $\log_d(1 - f_N) - 1 = O(1)$  and the fraction of such queues approaches one as  $N \rightarrow \infty$  with probability one. Therefore an incoming job experiences a delay of at least  $\log_d(1 - f_N) - 1 = O(1)$  with probability one as  $N \rightarrow \infty$ . For the case  $\tilde{\rho} = O(\frac{1}{N})$ , they showed that the behavior of queues that have queue lengths less than  $\log_d(1 - f_N) - 1 = O(1)$  can be studied by the time evolution of a deterministic process even under the diffusion scale. The complete diffusion limit analysis that characterizes queues of all sizes remains an open problem.

It is worth pointing out that the result we obtain is interesting: for the Halfin-Whitt regime, the effect of the randomized SQ( $d$ ) routing results in individual loss servers that are critically loaded but whose blocking cannot be obtained from the classical Halfin-Whitt blocking limit, instead, the blocking is obtained from the mean-field or a functional law of large numbers limit of the empirical occupancy distribution.

## Summary of Contributions

We characterize the asymptotics of the gap between the stochastic empirical process  $(\mathbf{X}^N(t); t \geq 0)$  and the corresponding mean-field limit  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$  when  $N \rightarrow \infty$  for both the transient and stationary regimes. To achieve this, we study the limit of the fluctuation process  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  where  $\mathbf{Z}^{(N)}(t) = \sqrt{N}(\mathbf{X}^{(N)}(t) - \mathbf{x}(t; \mathbf{u}))$  for both the

transient and stationary regimes when  $N \rightarrow \infty$ . We show that the limiting diffusion process is an OU process which depends on the mean-field limit in the transient regime and on the fixed-point of the mean-field in the stationary regime. We then use this result to show that  $\lim_{N \rightarrow \infty} \frac{1}{N} \overline{P}_{block}^{(N)}(\frac{d}{c})$  is a finite real number (Theorem 4.7). Since  $\frac{d}{c}$  is a finite constant real number, we conclude that  $\lim_{N \rightarrow \infty} P_{block}^{(N)} = 1$ . Preliminary results for the lightly loaded case were presented in [55].

## 1.2 Common Notation

In this section, we introduce the notation and terminology that is used in the rest of the dissertation. Let  $\mathbb{Z}, \mathbb{R}$  be the set of integers and real numbers, respectively. Further, let  $\mathbb{Z}_+, \mathbb{R}_+$  be the set of non-negative integers and non-negative real numbers, respectively.

For any given metric space  $E$ , let  $K_b(E); C_b(E); C_s(E)$  be the space of bounded measurable real valued functions, the space of bounded continuous real valued functions, and the space of continuous real valued functions with compact support, that are defined on  $E$ , respectively. Furthermore, for a metric space  $E$  of interest in this dissertation, let  $C^1(E)$  be the space of continuously differentiable real valued functions defined on  $E$  and let the subspace of functions in  $C^1(E)$  which have compact support be denoted by  $C_s^1(E)$ . The space of bounded functions in  $C^1(E)$  whose first derivatives are also bounded is denoted by  $C_b^1(E)$ . For any function  $f \in K_b(E), h \in C^1(E)$ , we define

$$\|f\| = \sup_{x \in E} |f(x)|; \|h\|_1 = \|h\| + \|h'\|;$$

where  $h'$  denotes the first derivative of  $h$ . The space  $C_b(E)$  is equipped with the uniform topology, i.e., we say that a sequence of functions  $\{f_n\}_{n=1}^\infty$  in  $C_b(E)$  converges to a function  $f \in C_b(E)$  if  $\|f_n - f\| \rightarrow 0$  as  $n \rightarrow \infty$ . The space  $C^1(E)$  is equipped with the topology induced by the norm  $\|\cdot\|_1$ . For two real valued functions  $f(\cdot)$  and  $h(\cdot)$ , we write  $f(x) = o(g(x))$  as  $x \rightarrow \infty$  if  $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = 0$ . Similarly, we write  $f(x) = O(g(x))$  as  $x \rightarrow \infty$  if there exists  $a > 0$  such that  $|f(x)| \leq a|g(x)|$  for sufficiently large values of  $x$ . For a given set  $A$ ,  $|A|$  denotes the number of elements in  $A$  and  $\bar{A}$  denotes the complement of  $A$ .

For a given metric space  $E$ , let the Borel  $\sigma$ -algebra be denoted by  $\mathcal{B}(E)$ . Let the space of finite non-negative measures on  $E$  be denoted by  $M_F(E)$ . We use the notation  $\mu(B)$  and  $\mu(y)$  to denote the measure of a Borel set  $B \in \mathcal{B}(E)$  and an element  $y \in E$  with respect to the measure  $\mu \in M_F(E)$ , respectively. The space of probability measures is denoted by

$M_1(E)$ . For a measure  $\mu \in M_1(E)$ , we write the product measure obtained from the  $k$ -fold product of  $\mu$  as  $\mu^{\otimes k}$ . Also, let  $M_1^{(N)}(E) \subset M_1(E)$  be the subspace of probability measures defined as

$$M_1^{(N)}(E) = \{ \mu \in M_1(E) : N(B) \in \mathbb{Z}_+; \forall B \in \mathcal{B}(E) \}$$

For any  $\mu \in M_b(E)$ ,  $\nu \in M_F(E)$ , we define

$$h \cdot \nu = \int_E h(y) \nu(dy)$$

The space of measures  $M_F(E)$  is equipped with the weak topology induced by the weak convergence of measures. For a signed finite measure  $\nu$ , the mapping  $\nu \mapsto h \cdot \nu$  is a continuous linear operator on the space of functions  $\mathcal{C}_b(E)$  which induces a norm for  $\nu$  as follows,

$$\|\nu\| = \sup_{h \in \mathcal{C}_b(E)} |h \cdot \nu| \tag{1.5}$$

For any Borel set  $B \in \mathcal{B}(E)$ , let us define the indicator function  $I_{fBg}$  of  $B$  as

$$I_{fBg}(u) = \begin{cases} 1 & \text{if } u \in B; \\ 0 & \text{otherwise;} \end{cases}$$

Let  $\mathbf{1}$  be the function defined such that for all  $u \in E$ , we have

$$\mathbf{1}(u) = 1$$

We use bold-faced Greek letters to write random measures in Chapters 2 and 3, and we use bold-faced capital letters to write other random variables. We also use small bold-faced letters to write vectors.

Consider a Polish space  $H$  and  $T \in \mathbb{R}_+$ . Let  $D_H([0; T])$  and  $D_H([0; \cdot])$  be the space of the càdlàg functions (right continuous functions with left limits) that are defined on  $[0; T]$  and  $[0; \cdot]$  with values in  $H$ , respectively. The càdlàg functions are also referred to as RCLL functions. Similarly, let  $C_H([0; T])$  and  $C_H([0; \cdot])$  be the space of continuous functions that take values in  $H$  defined on  $[0; T]$  and  $[0; \cdot]$ , respectively. The spaces  $D_H([0; T])$  and  $D_H([0; \cdot])$  are equipped with the Skorohod  $J_1$ -topology and hence, they are Polish spaces [56, Theorem 5.6, p.121]. Let the covariation of two local martingales  $(\mathbf{M}_t^1; t \geq 0)$  and  $(\mathbf{M}_t^2; t \geq 0)$  in  $D_R([0; T])$  be denoted by  $(\langle \mathbf{M}^1; \mathbf{M}^2 \rangle_t; t \geq 0)$  and the quadratic variation of  $(\mathbf{M}_t^1; t \geq 0)$  be denoted by  $(\langle \mathbf{M}^1 \rangle_t; t \geq 0)$ .

For our analysis, in Chapters 2 and 3, we define a set  $U$  that contains all the possible servers' states as its elements. In these two chapters, we study  $H$ -valued stochastic processes where  $H = M_F(U)$ . In Chapter 4, we define a set  $U$  that contains all the possible states of  $\mathbf{X}^{(N)}(t)$ . In this chapter, we study  $H$ -valued stochastic processes where  $H = U$ . The considered stochastic processes are random elements defined on  $(\Omega; F; P)$  with sample paths in  $D_H([0; T])$ , and are equipped with the Borel algebra generated by the open sets under the Skorohod  $J_1$ -topology [57]. We say that a sequence of random elements  $\{Y_n\}_{n \geq 1}$  defined on  $(\Omega; F; P)$  converges in distribution to  $\mathbf{Y}$  defined on  $(\Omega; F; P)$ , if for every bounded, continuous, and real valued functional  $f$ , we have  $\lim_{n \rightarrow \infty} E(f(Y_n)) = E(f(\mathbf{Y}))$ . We denote the convergence of a random element  $\{Y_n\}_{n \geq 1}$  in distribution to  $\mathbf{Y}$  by  $Y_n \xrightarrow{d} \mathbf{Y}$ . For a random element  $\mathbf{Y}$ , let us write its law as  $L(\mathbf{Y})$ .

## 1.3 An Introduction to Mean-field Techniques

In this section, we first provide a brief introduction to the mean-field techniques. We then summarize the mean-field analysis of the loss and PS models with the  $SQ(d)$  load balancing under the assumption of exponential JLDs.

### 1.3.1 An Overview of the Mean-field Techniques

Many models that arise in Computer, Information, and Societal systems are examples of complex stochastic systems can be modelled as interacting particle systems. A mathematical study of these systems to obtain some meaningful insights typically requires us to first model the time evolution of the system as the time evolution of a suitable Markov process. We can then understand the system's behavior by studying the Markov process. However, due to complex interactions between particles, the exact analysis of the Markov process is not tractable when the number of nodes or particles is finite. But as the number of particles becomes large, for some models, the influence of a particle on the rest of the system is negligible. In such a situation and under certain initial conditions, we can obtain a simpler limiting model referred to as the mean-field model by taking  $N \rightarrow \infty$  which can be used to approximate the behavior of the system with a large number of particles.

Consider a system with  $N$  nodes. For the sake of simplicity, we assume that the state of each particle lies in a finite set  $U = \{1; 2; \dots; ng\}$ . Let  $\mathbf{S}_i^{(N)}(t)$  be the state of the  $i^{\text{th}}$

particle at time  $t$ . Furthermore, assume that the process  $(\hat{\mathbf{Y}}^{(N)}(t); t \geq 0)$  is a Markov process where  $\hat{\mathbf{Y}}^{(N)}(t) = (\mathbf{S}_1^{(N)}(t); \dots; \mathbf{S}_N^{(N)}(t))$ . Then the state of  $\hat{\mathbf{Y}}^{(N)}(t)$  lies in  $U^N$  where  $U^N$  is the  $N$ -fold product of  $U$ . Let  $\hat{R}(\mathbf{x}; \mathbf{y})$  be the rate of transition of  $\hat{\mathbf{Y}}^{(N)}(t)$  from the state  $\mathbf{x} \geq U^N$  to  $\mathbf{y} \geq U^N$ . Then the time evolution of  $(\hat{\mathbf{Y}}^{(N)}(t); t \geq 0)$  is governed through a set of transition rates  $(\hat{R}(\mathbf{x}; \mathbf{y}); \mathbf{x}; \mathbf{y} \geq U^N)$ .

Recall that  $M_1(U)$  is the space of probability measures on  $U$ . For  $\mathbf{x} = (x_1; \dots; x_N) \geq U^N$ , let

$$\delta^{(N; \mathbf{x})} = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}. \quad (1.6)$$

Let  $M_1^{(N)}(U) \subset M_1(U)$  be defined such that

$$M_1^{(N)}(U) = \{ f^{(N; \mathbf{x})} : \mathbf{x} \geq U^N \mid g \}. \quad (1.7)$$

At time  $t$ , let  $\mu_t^{-(N)}$  be the empirical measure defined as

$$\mu_t^{-(N)} = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{s}_i^{(N)}(t)}. \quad (1.8)$$

Then the empirical process  $(\mu_t^{-(N)}; t \geq 0)$  is a measure-valued Markov process with values in  $M_1^{(N)}(U)$  [16, Proposition 2.3.3].

We assume that with probability one, only one particle's state jumps at any time  $t$ . For  $p \geq M_1(U)$ , let  $\tilde{R}(a; b; p)$  be the rate of transition of a particle from state  $a \geq U$  to  $b \geq U$  given that the empirical measure of the system is  $p$ . Then the evolution of the process is governed through a set of rates  $(\tilde{R}(p); p \geq M_1(U))$  where  $\tilde{R}(p) = (\tilde{R}(a; b; p); a; b \geq U)$  and  $\tilde{R}(a; b; p)$  is given by

$$\tilde{R}(a; b; p) = \lim_{h \downarrow 0} \frac{\mathbb{P}(\mathbf{S}_i^{(N)}(t+h) = b \mid \mathbf{S}_i^{(N)}(t) = a; \mu_t^{-(N)} = p)}{h}. \quad (1.9)$$

Because of the above assumption, a particle's behaviour depends on its own state and the state of  $\mu_t^{-(N)}$ . That is, for every measurable function  $f : U \rightarrow \mathbb{R}$ ,  $1 \leq i \leq N$ , and  $0 \leq s \leq t$ , we have

$$\mathbb{E} \left[ f(\mathbf{S}_i^{(N)}(t)) \mid \hat{\mathbf{Y}}^{(N)}(s) \right] = \mathbb{E} \left[ f(\mathbf{S}_i^{(N)}(t)) \mid \mathbf{S}_i^{(N)}(s); \mu_s^{-(N)} \right]. \quad (1.10)$$

This means all the particles influence the state of  $\mu_t^{-(N)}$  but the evolution of each particle's state is influenced by its own state and the state of  $\mu_t^{-(N)}$  at any time  $t$ . Here, the interactions



between particles take place through the empirical measure  $\mu_t^{(N)}$ . The system is a weakly interacting particle system because the contribution of each particle is of order  $\frac{1}{N}$  towards the empirical measure that captures the interactions between particles. As  $N \rightarrow \infty$ , the influence of any finite number of particles on the rest of the system is negligible.

For many systems, it is sufficient to study the empirical measure-valued process  $(\mu_t^{(N)}; t \geq 0)$ . Let  $A^{(N)}(\cdot)$  be the generator of the Markov process  $(\mu_t^{(N)}; t \geq 0)$ . Then for  $f: M_1(U) \rightarrow \mathbb{R}$  and  $\rho \in M_1(U)$ ,

$$A^{(N)}f(\rho) = \sum_{x,y \in U} N\rho_x \left( f\left(\rho + \frac{1}{N}\delta_y - \frac{1}{N}\delta_x\right) - f(\rho) \right) \tilde{R}(x,y;\rho) \quad (1.11)$$

Then although  $(\mu_t^{(N)}; t \geq 0)$  is a stochastic process for each  $N$ , this process converges to a deterministic process referred to as the mean-field limit when  $N \rightarrow \infty$  under certain initial conditions. Let  $(\mathbf{X}^{(N)}(t); t \geq 0)$  be the empirical process corresponding to the measure-valued process  $(\mu_t^{(N)}; t \geq 0)$ . From applying the Kurtz theorem [58, Theorem 2.11], we have

**Theorem 1.1.** *Let  $\tilde{R}(x,y;\cdot)$  be Lipschitz continuous for all  $x,y \in U$ . For a constant  $\mathbf{u}$  and  $T > 0$ , if  $\mathbf{X}^{(N)}(0) \rightarrow \mathbf{u}$  in probability as  $N \rightarrow \infty$ , then  $(\mathbf{X}^{(N)}(t); 0 \leq t \leq T) \rightarrow (\mathbf{x}(t;\mathbf{u}); 0 \leq t \leq T)$  in probability as  $N \rightarrow \infty$  uniformly in  $t$  where  $(\mathbf{x}(t;\mathbf{u}); t \geq 0)$  is the unique solution of*

$$\frac{d\mathbf{x}(t;\mathbf{u})}{dt} = \mathbf{h}(\mathbf{x}(t;\mathbf{u})) \quad \mathbf{x}(0;\mathbf{u}) = \mathbf{u}; \quad (1.12)$$

where  $\mathbf{h}(\cdot) = (h_i(\cdot); i = 0)$  is a function that depends on the system's dynamics.

The process  $(\mathbf{x}(t;\mathbf{u}); t \geq 0)$  corresponds to the distribution of a non-linear Markov process with distribution  $\mathbf{x}(t;\mathbf{u})$  and infinitesimal generator  $\mathbf{h}(\mathbf{x}(t;\mathbf{u}))$  at time  $t$ . Existence of the mean-field limit for an interacting particle system implies that if there is a law of large numbers effect in the system at time  $t = 0$ , it guarantees the functional law of large numbers (FLLN) result due to weak interactions between particles in the transient regime. The initial law of large numbers effect can be achieved if we assume all the particles' states are i.i.d. random variables.

A consequence of having the mean-field limit for an interacting particle system is the result of propagation of chaos. We now recall some results on the propagation of chaos from [47, 59].

**Definition 1.1.** For  $Q \in M_1(U)$ , a sequence of random variables  $(\tilde{\mathbf{W}}_i^{(N)})_{1 \leq i \leq N}$  on  $U^N$  is said to be  $Q$ -chaotic if for any  $k \geq 1$ ,

$$\lim_{N! \rightarrow \infty} L(\tilde{\mathbf{W}}_1^{(N)}; \dots; \tilde{\mathbf{W}}_k^{(N)}) = Q^{\otimes k}; \quad (1.13)$$

weakly in  $M_1(U^k)$ . Similarly, if a sequence of processes  $(\tilde{\mathbf{Z}}_i^{(N)}(t); t \geq 0)_{1 \leq i \leq N}$  is chaotic whenever  $(\tilde{\mathbf{Z}}_i^{(N)}(0))_{1 \leq i \leq N}$  is chaotic, then we say that there is a propagation of chaos.

The  $Q$ -chaoticity for random variables implies that asymptotically any finite number of random variables are independent of each other and furthermore, each random variable's law is  $Q$ . The propagation of chaos for processes implies that any finite collection of processes are asymptotically independent of each other whenever the initial states satisfy the chaoticity.

We define the notion of exchangeability of random variables.

**Definition 1.2.** *Exchangeability of Random Variables:* Let  $(\mathbf{Y}_k; 1 \leq k \leq N)$  denote a collection of  $N$  random variables. Then the collection is called exchangeable if the joint law of the collection is invariant under any permutation of indices,  $1 \leq k \leq N$ , of random variables, i.e., if  $\pi$  is a permutation on indices  $1, 2, \dots, N$ , then

$$L(\mathbf{Y}_1; \dots; \mathbf{Y}_N) = L(\mathbf{Y}_{\pi(1)}; \dots; \mathbf{Y}_{\pi(N)}); \quad (1.14)$$

From [59, Proposition 2.2, p.177], we have

**Proposition 1.1.** For  $Q \in M_1(U)$ , let a collection of random variables  $(\tilde{\mathbf{W}}_i^{(N)})_{1 \leq i \leq N}$  on  $U^N$  be  $Q$ -chaotic, then for  $\mathbf{w}^{(N)} = \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i^{(N)}$ , we have  $\lim_{N! \rightarrow \infty} L(\mathbf{w}^{(N)}) = Q$  weakly in  $M_1(M_1(U))$ . If  $(\tilde{\mathbf{W}}_i^{(N)})_{1 \leq i \leq N}$  is exchangeable for  $N \geq 1$  and  $\lim_{N! \rightarrow \infty} L(\mathbf{w}^{(N)}) = Q$  weakly, then  $(\tilde{\mathbf{W}}_i^{(N)})_{1 \leq i \leq N}$  is  $Q$ -chaotic.

Assuming that  $(\mathbf{S}_i^{(N)}(0))_{1 \leq i \leq N}$  is exchangeable, we have the exchangeability of  $(\mathbf{S}_i^{(N)}(t))_{1 \leq i \leq N}$  for any time  $t \geq 0$ . In this case, from Theorem 1.1 and Proposition 1.1, any finite set of particles are asymptotically independent of each other and furthermore, each particle's distribution is  $\mathbf{x}(t; \mathbf{u})$ . One can establish chaoticity not only for each  $t$  but also on the entire path space. For example, in [47, Theorem 3.4], chaoticity on path space was established for the FCFS model with exponential JLDs under the SQ( $d$ ) policy. Let  $\mathbf{S}_i^{(N)} = (\mathbf{S}_i^{(N)}(t); t \geq 0)$ . A random variable  $\mathbf{w}^{(N)}$  was defined as follows

$$\mathbf{w}^{(N)} = \frac{1}{N} \sum_{i=1}^N \mathbf{s}_i^{(N)}; \quad (1.15)$$

Here,  $\bar{\mu}^{(N)}$  has samples in  $M_1(D_{Z^+}([0; 1]))$ . Then Theorem 3.4 of [47] concludes that if  $(\mathbf{S}_i^{(N)}(0))_{1 \leq i \leq N}$  is  $\bar{\mu}_0$ -chaotic, then  $(\mathbf{S}_i^{(N)})_{1 \leq i \leq N}$  is  $\bar{\mu}$ -chaotic where  $\bar{\mu} = (\bar{\mu}_t; t \geq 0)$  is the mean-field with initial point  $\bar{\mu}_0$ .

So far, we have studied both the mean-field limit and the system with  $N$  particles in the transient regime. Often, we are interested in understanding the equilibrium or stationary behavior of complex stochastic systems. For many systems of interest, it is a challenging task to find the stationary distributions or their approximations. The key question is, can we use the mean-field limit to obtain the stationary distribution of a particle in the system with a large number of particles? Here, the main issue is that we have established the mean-field limit under a strong assumption on the initial measures, but we do not know whether such assumption holds in the stationary regime or not as it depends on the system's behavior. Suppose the mean-field has a unique fixed-point, then can we use the fixed-point of the mean-field to approximate the stationary distribution of a particle in the system with large  $N$ ? In the literature, there are some sufficient conditions under which the fixed-point of the mean-field can be used to approximate the stationary distribution of a particle. The relationship between the mean-field and the stationary system was studied in [49–51, 60–65] for different contexts.

In [62], by working with empirical processes, the mean-field limit was established for a generic context. Let  $\mathbf{x}(t; \mathbf{u})$  be the mean-field limit with initial point  $\mathbf{u}$ . In this case, the mean-field limit is a solution to a set of MFEs. If the system with parameter  $N$  is stationary, let  $\bar{\mu}^{(N)}$  be the stationary distribution of the empirical process, then the support of every limit point of  $\bar{\mu}^{(N)}$  is a compact set included in the Birkhoff center of the mean-field where the Birkhoff center is the closure of all the recurrent points of the mean-field [62, 65]. The Birkhoff center contains all the limit cycles and fixed-points of the mean-field. As a result, it was argued that if the mean-field has a unique fixed-point and furthermore if it is GAS, then  $\bar{\mu}^{(N)} \rightarrow \delta_{\bar{\mu}}$ . As a result, as  $N \rightarrow \infty$ , the distribution of a particle in the stationary regime is  $\delta_{\bar{\mu}}$ . In [61], it was shown that every limit of  $\bar{\mu}^{(N)}$  is an invariant distribution of the mean-field. That is, let  $Z \in M_1(V)$  be a limiting point of  $\bar{\mu}^{(N)}$  where  $V$  is the space in which the mean-field lies, if the initial point of the mean-field has distribution  $Z$ , then for every  $f \in C_b(M_1(V)) \rightarrow \mathbb{R}$ ,

$$\int_{M_1(V)} f(\mathbf{x}(t; \mathbf{u})) dZ(\mathbf{u}) = \int_{M_1(V)} f(\mathbf{u}) dZ(\mathbf{u}): \quad (1.16)$$

In [63], it was shown that if the system with  $N$  particles is reversible system, then the support of every limit point concentrates on the fixed-points of the mean-field. As

a result, if the fixed-point of the mean-field has a unique fixed-point  $\mu$ , then  $\mu^{(N)} \rightarrow \mu$ . In several contexts, the existence of multiple fixed-points and limit cycles was observed for the mean-field in the literature. In [60, 61], they give examples of models for which the mean-field limit has a unique fixed-point, but the mean-field limit starting from other than the fixed-point follows a cyclic path. It was also observed that the support of  $\mu^{(N)}$  lies closer to the limit cycle as  $N$  becomes large. In [65], a mean-field model was studied for which the Birkhoff center of the mean-field is a limit cycle, and it was shown that the support of  $\mu^{(N)}$  lies inside the limit cycle when  $N$  becomes large. In [66], a mean-field limit was found for which there exist two fixed-points of the mean-field, and a sample path of the empirical process oscillates between two regions with each region having a fixed-point supporting the meta-stability of the system. All these examples imply that although the mean-field limit can be used to approximate the distribution of a particle in the transient regime, it requires further investigation whether the obtained mean-field is useful or not to obtain approximations for the stationary distributions.

### 1.3.2 Mean-field Analysis of the SQ( $d$ ) Policy

We now provide a summary of the mean-field analysis of the SQ( $d$ ) policy for the loss and PS models under the assumption of exponential JLDs.

#### Analysis of the Loss System

In this section, we provide a summary of the analysis given in [8, 10]. Consider a system of  $N$  loss servers each with capacity  $C$ . Assume that jobs arrive according to a Poisson process with rate  $N\lambda$  and furthermore, jobs are dispatched according to the SQ( $d$ ) policy. The jobs have exponential JLDs with unit average job length. The first step is to obtain a mathematical formulation to the time evolution of the system. The empirical process  $(\mathbf{X}^{(N)}(t); t \geq 0)$  is a Markov process where  $\mathbf{X}^{(N)}(t) = (\mathbf{X}_i^{(N)}(t); i = 0, \dots, C)$  and  $\mathbf{X}_i^{(N)}(t)$  denotes the fraction of servers with at least  $i$  progressing jobs at time  $t$ . Assume that a job arrives to the system when the state of  $\mathbf{X}^{(N)}(t)$  is  $\mathbf{b} = (b_0; \dots; b_C)$  indicating that the fraction of servers with at least  $i$  progressing jobs is equal to  $b_i$  for  $0 \leq i \leq C$ . As servers are sampled with replacement uniformly at random, the probability that a sampled server has  $i$  jobs is equal to  $b_i - b_{i+1}$ . Since the dispatcher samples  $d$  servers with replacement to dispatch the arrival, the destination server of the job will have occupancy  $n$  with probability  $b_n^d - b_{n+1}^d$ . Then the probability that a tagged server with  $n$  jobs is selected as the destination server

is equal to  $\frac{1}{N} \frac{b_n^d}{b_n} \frac{b_{n+1}^d}{b_{n+1}}$ . As a result, the arriving job joins a particular server with a certain probability that depends on the occupancy of that server and the empirical measure  $\mathbf{b}$ . Hence, the time evolution of each server depends on its state and the state of the empirical measure. Therefore each server influences  $\mathbf{X}^{(N)}(t)$  and an arrival event to each server is influenced by the server's state and the state of  $\mathbf{X}^{(N)}(t)$ . Hence, the model is a weakly interacting particle system.

We now present results on the mean-field analysis of the model from [10]. The underlying space is

$$\mathbb{U} = \{ \mathbf{b} = (b_i)_{i=0}^C : b_0 = 1, b_i \geq 0, b_C = b_{C+1} = 0 \} \quad (1.17)$$

The space  $\mathbb{U}$  is equipped with the euclidean metric.

**Theorem 1.2.** *For  $\mathbf{u} \in \mathbb{U}$ , if  $\mathbf{X}^{(N)}(0) \rightarrow \mathbf{u}$  as  $N \rightarrow \infty$ , then  $(\mathbf{X}^{(N)}(t); t \geq 0) \rightarrow (\mathbf{x}(t; \mathbf{u}); t \geq 0)$  as  $N \rightarrow \infty$  where  $(\mathbf{x}(t; \mathbf{u}); t \geq 0) = (x_n(t; \mathbf{u}); t \geq 0; 0 \leq n \leq C)$  is the unique solution to the following equations:*

for  $\mathbf{h}(\mathbf{x}(t; \mathbf{u})) = (h_n(\mathbf{x}(t; \mathbf{u})))_{n=0}^C$ ,

$$\mathbf{x}(0; \mathbf{u}) = \mathbf{u}; \quad \frac{dx_n(t; \mathbf{u})}{dt} = h_n(\mathbf{x}(t; \mathbf{u})); \quad (1.18)$$

where

$$h_0(\mathbf{x}(t; \mathbf{u})) = 0; \quad (1.19)$$

and for  $n = 1, \dots, C$ ,

$$h_n(\mathbf{x}(t; \mathbf{u})) = (x_{n-1}^d(t; \mathbf{u}) - x_n^d(t; \mathbf{u})) - n(x_n(t; \mathbf{u}) - x_{n+1}(t; \mathbf{u})); \quad (1.20)$$

with  $x_0(t; \mathbf{u}) = 1$  and  $x_{C+1}(t; \mathbf{u}) = 0$ . The deterministic process  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$  is referred to as the mean-field limit and the equations (1.18)-(1.20) are referred to as the mean-field equations with the initial point  $\mathbf{u}$ .

The mean-field  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$  has a unique GAS fixed-point  $\mathbf{x}^* = (x_n^*)_{n=0}^C$  with  $x_0^* = 1$ . The proof that the fixed-point is GAS mainly uses the property that the mean-field is quasi-monotonic.

**Definition 1.3. Quasi-Monotonicity:** Let  $\mathbf{r}(t; \mathbf{u})$  be a deterministic process defined on  $\mathbb{R}^n$ ,  $n = 1, \dots, C$ , that is a solution to a set of ODEs, for  $\mathbf{f}(\cdot) = (f_i(\cdot))_{i=1}^n$ ,

$$\frac{d\mathbf{r}(t; \mathbf{u})}{dt} = \mathbf{f}(\mathbf{r}(t; \mathbf{u})); \quad \mathbf{r}(0; \mathbf{u}) = \mathbf{u};$$

Then the process is said to satisfy quasi-monotonicity if  $\mathbf{u} \leq \mathbf{v}$  element-wise, then  $\mathbf{r}(t; \mathbf{u}) \leq \mathbf{r}(t; \mathbf{v})$  element-wise.

It can be shown that the following exchange of limits holds

$$\lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \mathbf{X}^{(N)}(t) = \lim_{t \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbf{X}^{(N)}(t):$$

Using Theorem 1.1, under the assumption of the exchangeability of initial servers' states, we can show the independence of any finite set of servers as  $N \rightarrow \infty$ . Also, for a server,  $\mathbf{x}(t; \mathbf{u})$  denotes the distribution at time  $t$  and  $\pi$  denotes the stationary distribution, as  $N \rightarrow \infty$ . As a result,  $\pi_C$  denotes the stationary probability that a server is fully occupied when  $N \rightarrow \infty$ . Since the dispatcher samples  $d$  servers when a job arrives, the stationary average blocking probability of a job as  $N \rightarrow \infty$  is equal to  $\frac{d}{C}$ , where we use the fact that a sampled server will have  $C$  jobs with probability  $\pi_C$ .

We now focus on finding the fixed-point  $\pi$ . For loss models, a closed form for  $\pi$  is not known while it is known for the case of the FCFS models [7, 12]. We can find  $\pi$  numerically as explained below. The fixed-point  $\pi$  is the unique solution to the following equations

$$\binom{d}{n-1} \pi_{n-1} = n \binom{d}{n} \pi_n \quad (1.21)$$

for  $n \geq 1$  and  $\pi_{C+1} = 0$ . Then from (1.21), we can also write

$$\frac{\binom{d}{n-1} \pi_{n-1}}{\binom{d}{n} \pi_n} = n \quad (1.22)$$

for  $n \geq 1$  and  $\pi_{C+1} = 0$ . Let us define  $\hat{\pi}_n = \frac{\binom{d}{n} \pi_{n+1}}{\binom{d}{n-1} \pi_n}$ . Then from (1.22),  $\hat{\pi}_n$  is the stationary distribution of the single server loss model with a Poisson arrival process of jobs having rate  $\hat{\pi}_n$  when there are  $n$  progressing jobs, and  $\pi_n$  is the probability that the server has at least  $n$  progressing jobs. Let  $\mathcal{M}_1(\mathbb{R}_+; 1; \cdot; Cg)$  be the set of probability measures on  $\mathbb{R}_+; 1; \cdot; Cg$ . Then from [10], the fixed-point  $\pi$  can be computed using the formula for the stationary distribution of a single server loss system with state-dependent arrival rates. We first define two mappings,  $\Theta : \mathcal{M}_1(\mathbb{R}_+; 1; \cdot; Cg) \rightarrow \mathbb{R}_+^{C+1}$  and  $\hat{\Xi} : \mathbb{R}_+^{C+1} \rightarrow \mathcal{M}_1(\mathbb{R}_+; 1; \cdot; Cg)$  that are used in computing  $\pi$ . For every  $(p_0; \dots; p_C) \in \mathcal{M}_1(\mathbb{R}_+; 1; \cdot; Cg)$ , there exists  $(r_0; \dots; r_C) \in \mathbb{R}_+^{C+1}$  such that

$$\Theta((p_0; \dots; p_C)) = (r_0; \dots; r_C);$$

where

$$r_n = \frac{((\sum_{j=n}^C p_j)^d - (\sum_{i=n+1}^C p_i)^d)}{((\sum_{j=n}^C p_j) - (\sum_{i=n+1}^C p_i))}.$$

Similarly, for every  $(b_0; \dots; b_C) \in \mathbb{R}_+^{C+1}$ , there exists  $(a_0; \dots; a_C) \in \mathcal{M}_1(\mathbb{R}_+; 1; \cdot; Cg)$  such that

$$\hat{\Xi}((b_0; \dots; b_C)) = (a_0; \dots; a_C);$$

where

$$a_n = \left( \prod_{i=1}^n \left\{ \frac{b_{i-1}}{i} \right\} \right) a_0$$

for  $n \geq 1$  and  $\sum_{i=0}^C a_i = 1$ .

**Lemma 1.1.** *Since the mean-field is GAS, the mapping  $\widehat{\Xi}(\Theta)$  has a unique fixed-point denoted by  $\mathbf{u} = (u_i; i = 0, \dots, n = C)$ . Furthermore, the unique fixed-point  $\mathbf{u}$  of the mean-field is given by,  $n \geq 1$ ,*

$$u_n = \sum_{j=n}^C u_j$$

### Analysis of the PS System

In this section, we consider the mean-field analysis of the PS model with exponential JLDs from [7, 9, 11]. In [7], the FCFS model with exponential JLDs was studied. Heterogeneous PS model in which servers are categorized into different classes based on their capacities was studied in [9, 11]. We provide a summary of the analysis given in [9, 11] for our model. We first define the underlying spaces

$$U^0 = \{ (b_i; i = 0) : b_0 = 1; b_i \in [0, 1]; i \geq 1 \} \quad (1.23)$$

$$U = \{ (b_i; i = 0) : b_0 = 1; b_i \in [0, 1]; i \geq 1; \sum_{i=0}^C b_i < 1 \} \quad (1.24)$$

We equip the space  $U^0$  with the metric  $d(\cdot, \cdot)$  defined as

$$d(\mathbf{u}; \mathbf{v}) = \sup_{i \geq 1} \frac{|u_i - v_i|}{i+1}; \quad (1.25)$$

for all  $\mathbf{u} = (u_i; i = 0)$ ,  $\mathbf{v} = (v_i; i = 0)$  in  $U$ . The space  $U^0$  equipped with the metric  $d$  is compact, complete, and separable.

A mean-field limit exists in this case also similar to the loss model but now the mean-field lies in a countable infinite dimensional space.

**Theorem 1.3.** *For  $\mathbf{u} \in U^0$ , if  $\mathbf{X}^{(N)}(0) \rightarrow \mathbf{u}$  as  $N \rightarrow \infty$ , then  $(\mathbf{X}^{(N)}(t); t \geq 0) \rightarrow (\mathbf{x}(t; \mathbf{u}); t \geq 0)$  as  $N \rightarrow \infty$  where  $(\mathbf{x}(t; \mathbf{u}); t \geq 0) = (x_n(t; \mathbf{u}); t \geq 0; n = 0)$  is the unique solution to the following equations:*

for  $\mathbf{h}(\mathbf{x}(t; \mathbf{u})) = (h_n(\mathbf{x}(t; \mathbf{u})); n = 0)$ ,

$$\mathbf{x}(0; \mathbf{u}) = \mathbf{u}; \quad \frac{dx_n(t; \mathbf{u})}{dt} = h_n(\mathbf{x}(t; \mathbf{u})); \quad (1.26)$$

where

$$h_0(\mathbf{x}(t; \mathbf{u})) = 0; \quad (1.27)$$

and for  $n \geq 1$ ,

$$h_n(\mathbf{x}(t; \mathbf{u})) = (x_{n-1}^d(t; \mathbf{u}) - x_n^d(t; \mathbf{u}), (x_n(t; \mathbf{u}) - x_{n+1}(t; \mathbf{u})); \quad (1.28)$$

with  $x_0(t; \mathbf{u}) = 1$ . The deterministic process  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$  is referred to as the mean-eld limit and the equations (1.26)-(1.28) are referred to as the mean-eld equations with the initial point  $\mathbf{u}$ .

It was then shown that there exists a unique fixed-point of the mean-field  $\bar{\mathbf{x}} = (\bar{x}_i; i \geq 0)$  that belongs to the set  $\mathcal{U}$ . Furthermore, the quasi-monotonicity of the mean-field was used to show  $\lim_{t \rightarrow \infty} \mathbf{x}(t; \mathbf{u}) = \bar{\mathbf{x}}$  for all  $\mathbf{u} \in \mathcal{U}$ . After that the stochastic system was shown to be stable for  $\rho < 1$  and let  $\pi^{(N)}$  be the unique stationary distribution of the stochastic process  $\mathbf{X}^{(N)}$ . Since the space  $\mathcal{U}^\theta$  is compact, from Prohorov's theorem [57],  $\pi^{(N)} g_{N-1}$  is tight. Therefore, it is sufficient to show that every limit point of  $\pi^{(N)} g_{N-1}$  coincides with  $\bar{\mathbf{x}}$ . Let  $\pi$  be the limit point of a converging subsequence  $\pi^{(N_k)} g_{k-1}$ , then it was shown that  $\mathbb{E}[\sum_{i=1}^d \mathbf{V}_i] < 1$  where  $(\mathbf{V}_i; i \geq 0)$  is random quantity with distribution  $\pi$ . As a consequence, we have  $\pi(\mathcal{U}) = 1$ . Since  $\bar{\mathbf{x}}$  is an invariant distribution of the mean-field and  $\lim_{t \rightarrow \infty} \mathbf{x}(t; \mathbf{u}) = \bar{\mathbf{x}}$  for all  $\mathbf{u} \in \mathcal{U}$ , it implies  $\pi = \bar{\mathbf{x}}$ . Having  $\pi^{(N)} g_{N-1} \rightarrow \bar{\mathbf{x}}$  implies that any finite set of servers are asymptotically independent in the stationary regime and furthermore, each server's distribution coincides with the fixed-point  $\bar{\mathbf{x}}$  of the mean-field. More importantly, the fixed-point is given by  $\bar{x}_i = \frac{\rho^i}{\rho^i - 1}, i \geq 1$ .

## 1.4 Outline

The rest of the dissertation is organized as follows. In Chapter 2, we provide the mean-field analysis of the first problem discussed in Section 1.1.1. We then study the second problem stated in Section 1.1.2 in Chapter 3. The analysis of the FCLT for the loss model introduced in Section 1.1.3 is given in Chapter 4. We conclude in Chapter 5 with a discussion on future work. We provide additional background material in the Appendix.



## Chapter 2

# Insensitivity of the Mean-field Limit of Loss Systems Under $SQ(d)$ Load Balancing

In this chapter, we study a large multi-server loss model under a randomized dynamic load balancing scheme when the service time distributions are general with finite mean. In particular, we consider the  $SQ(d)$  load balancing scheme according to which an incoming job is dispatched to the server with the least number of progressing jobs among  $d$  randomly chosen servers. Previous works have addressed the exponential service time case when the number of servers goes to infinity, giving rise to a mean-field model. The fixed-point of the mean-field equations (MFEs) was seen to be insensitive to the service time distribution in simulations, but no proof was available. While insensitivity is well known for loss systems models even with state-dependent inputs such models belong to the class of linear Markov models. In the context of the  $SQ(d)$  load balancing, the resulting model belongs to the class of nonlinear Markov processes (processes whose generator itself depends on the distribution) for which traditional arguments do not directly apply. Showing insensitivity to the general service time distributions has thus remained an open problem. In this case, obtaining the MFEs poses a challenge due to the resulting Markov description of the system being in positive orthant as opposed to a finite chain in the exponential case. In this chapter, we first obtain the MFEs and then show that the MFEs have a unique fixed-point that corresponds to a distribution of occupancy coinciding with the fixed-point in the exponential case, thus establishing insensitivity. The approach is via a measure-valued Markov process representation and the martingale problem to establish the mean-field

limit.

**Organization of the Chapter:** The rest of the chapter is organized as follows: Section 2.1 describes the system model. In Section 2.2, we introduce the notation used in the rest of the chapter. In Section 2.3, we derive a measure-valued representation for the state of the system. The main results of the chapter are given in Section 2.4. We then establish the mean-field limit in Section 2.5. In Section 2.6, we prove the result on the uniqueness of the fixed-point of the MFEs and show that the fixed-point is insensitive to the distribution, i.e., it depends only on the mean service time. In Section 2.7, we provide numerical results that suggest the global asymptotic stability of the fixed-point of the MFEs and hence, the relation (1.1) indeed holds. In Section 2.8, we give proofs of the main results stated in Section 2.4. Finally, Section 2.9 concludes the chapter with some remarks and generalizations.

## 2.1 System Model

We consider a system consisting of a large number  $N$  of parallel servers. Jobs arrive according to a Poisson process with rate  $N\lambda$ , and the job lengths are assumed to be i.i.d. from a general distribution  $G(\cdot)$  defined on  $\mathbb{R}_+$ . Furthermore, we assume that the distribution  $G(\cdot)$  has a continuous density function  $g(\cdot)$ . A central job dispatcher routes an incoming job to a server according to the  $SQ(d)$  policy described below. We assume that each server has capacity to process up to a number  $C$  of jobs simultaneously, and each accepted job is processed at a unit rate. At any time  $t$ , if a server is currently serving  $i$  jobs, then we say that this server has occupancy  $i$  and vacancy  $C - i$  at time  $t$ . If an incoming job is routed to a server with occupancy  $C$ , then the routed job is blocked from service, and it is discarded from the system immediately. Otherwise, the processing of the job begins immediately, and it is processed at a unit rate.

**Definition 2.1.1.**  *$SQ(d)$  or Power-of- $d$  load balancing policy: An incoming job is routed to the server with the minimum occupancy among  $d$  servers that are selected randomly with replacement. Ties among servers are broken by choosing a server uniformly at random. The randomly chosen  $d$  servers are referred to as the potential destination servers and the server to which a job is routed called the destination server.*

**Remark 2.1.** *In Definition 2.1.1, we assume sampling with replacement because of notational convenience, and it can be shown that the asymptotic results that are of interest in*

this chapter are not affected whether we sample with or without replacement [32, 47].

We assume that the service times have finite mean  $\bar{1}$ . We make an assumption that  $\bar{G}(\cdot)$  is supported on  $[0; 1)$  where  $\bar{G}(\cdot)$  denotes the complementary distribution. Otherwise, if  $M$  denotes the finite support of  $\bar{G}(\cdot)$ , then our analysis easily extends to this case by using the fact that the ages of jobs are at most  $M$ . The hazard rate function of  $G(\cdot)$  is defined as

$$h(x) = \frac{g(x)}{\bar{G}(x)} = \frac{g(x)}{1 - G(x)} \quad (2.1)$$

for  $x \geq 0$ . The hazard rate function indicates the instantaneous rate at which the service of a job ends. More precisely, a job with age  $y$  (where  $y$  denotes the time since its arrival) at time  $t$  exits the server in the interval  $[t; t + dt)$  with probability  $h(y)dt + o(dt)$ .

**Assumption 2.1.1.** *The hazard rate function satisfies  $h \in C_b(\mathbb{R}_+)$ :*

**Remark 2.2.** *The Assumption 2.1.1 is true for several classes of distributions such as Phase-Type distributions, Gamma distributions, Log-Normal distributions, and any Pareto distribution with finite mean.*

## 2.2 Additional Notation and Terminology

In this section, we introduce the required additional notation and terminology that is specific to this chapter.

In this chapter, we define the term ‘age’ for a progressing job as the amount of time elapsed since its arrival. To model the dynamics of an Erlang loss system with capacity  $C$  for each server by a Markov process, we define the state of each server as  $(n; a_1; a_2; \dots; a_n)$  where  $n$  denotes the number of jobs that are in progress at the server and  $a_i$  denotes the age of the  $i^{\text{th}}$  progressing job. We now define a space  $\mathcal{U}$  that was used in earlier works to study queueing models with general service time distributions by using the classical supplemental variable method [67, 68] such that it contains all the possible servers’ states as its elements. The space  $\mathcal{U}$  is defined as

$$\mathcal{U} = \bigcup_{n=0}^C \mathcal{U}_n;$$

where  $\mathcal{U}_0 = \{0\}$  and an element in  $\mathcal{U}_n$  for  $n \geq 1$  is of the form  $(n; a_1; \dots; a_n)$  where  $1 \leq n \leq C$  and  $a_i \geq 0$ . The state of a server that has  $n$  jobs belongs to the space

$\cup_n$ . Here, one might omit the variable  $n$  and consider just  $(a_1; \dots; a_n)$  to denote a server's state, but such a representation does not account for idle servers while  $(0)$  is the state of idle servers in our representation. Furthermore, the variable  $n$ , directly gives us information about the number of progressing jobs at a server which changes upon every arrival and departure. Hence, it is convenient to work with the state representation that has a variable to denote the number of progressing jobs at a server.

It is also possible to define an element in  $\cup_n$  by  $(n; a_1; \dots; a_n; 0; \dots; 0)$  of size  $C + 1$ . This allows us to have constant size of  $C + 1$  for an element in  $\cup$ . Note that the zeros in the state  $(n; a_1; \dots; a_n; 0; \dots; 0)$  act as dummy variables as there are only  $n$  jobs. Hence, to make it simple, we consider that an element in  $\cup_n$  is of the form  $(n; a_1; \dots; a_n)$  with size  $n+1$ . Without loss of generality, we refer to an element in the set  $\cup$  by  $\mathbf{u}$  and an element in the set  $\cup_n$  by  $\mathbf{u}_n$ . Note that we have  $\mathbf{u}_0 = 0$ . For  $\mathbf{y}_n = (n; y_1; \dots; y_n); \mathbf{z}_m = (m; z_1; \dots; z_m)$ , we define the metric  $d_{\cup}(\mathbf{y}_n; \mathbf{z}_m)$  as

$$d_{\cup}(\mathbf{y}_n; \mathbf{z}_m) = \begin{cases} \sum_{i=1}^n j y_i & z_j & \text{if } n = m; \\ j m & n j & \text{otherwise.} \end{cases}$$

For  $n \geq 1$ ,  $\cup_n$  is a complete, separable, and Polish space. Furthermore,  $\cup$  is a Polish space as it is a union of a set of disjoint Polish spaces. Also,  $\cup$  is separable and complete.

For  $(n; u_1; \dots; u_n) \in \cup_n$  and  $y \geq 0$ , we use the following notation

$$\begin{aligned} \mathbf{u}_n &= (n; u_1; \dots; u_n); \\ \mathbf{u}_n^j &= (n-1; u_1; \dots; u_{j-1}; u_{j+1}; \dots; u_n); \\ (\mathbf{u}_n^j; y) &= (n+1; u_1; \dots; u_{j-1}; y; u_j; \dots; u_n); \\ (\mathbf{u}_n^j; y) &= (n; u_1; \dots; u_{j-1}; y; u_{j+1}; \dots; u_n); \end{aligned}$$

A function  $f : \cup \rightarrow \mathbb{R}$  is said to be differentiable if for every  $n \geq 1$ , the function  $\frac{\partial f(\mathbf{u}_n)}{\partial u_i}$  exists for all  $1 \leq i \leq n$  at every  $\mathbf{u}_n \in \cup_n$ . As a result, from (1.6), the function  $f|_{\cup_n}; n \geq 1$ , is differentiable. For a differentiable function  $f : \cup \rightarrow \mathbb{R}$ , we define

$$\|f\|_k = \max_{n \geq 1} \left( \sup_{\mathbf{u}_n \in \cup_n} \left( \max_{1 \leq i \leq n} \left| \frac{\partial f(\mathbf{u}_n)}{\partial u_i} \right| \right) \right);$$

Further, for a differentiable function  $f : \cup \rightarrow \mathbb{R}$ , let the function  $r_1 f$  be defined as

$$r_1 f(n; u_1; \dots; u_n) = r f(\mathbf{1}) = \sum_{i=1}^n \frac{\partial f(\mathbf{u}_n)}{\partial u_i};$$

A measure  $\mu \in \mathcal{M}_F(U)$  when it is restricted to  $U_0$  is a Dirac measure at  $f(0)g$  satisfying  $\mu(U_0) = \mu(f(0)g)$ . We say that a measure  $\mu$  is absolutely continuous with respect to the Lebesgue measure if  $\mu(\mathbf{x}_n) = 0$  at every  $\mathbf{x}_n \in U_n$  for all  $n \geq 1$ . For any Borel measurable function  $f$  that is defined on  $U$ , we define

$$h(\cdot; f) = f(0) \mu(f(0)g) + \sum_{n=1}^{\infty} \int_{U_n} f(\mathbf{z}_n) \mu(d\mathbf{z}_n):$$

We now define a function  $\mu : U \rightarrow \mathbb{R}$  as follows:

$$\mu(\mathbf{x}_n) = \begin{cases} \sum_{i=1}^n x_i & \text{if } n \geq 1; \\ 0 & \text{otherwise;} \end{cases}$$

For  $b \geq 0$ , let  $\mu_b^+ : U \rightarrow U$  be the transition operator defined as

$$\mu_b^+(\mathbf{x}_n) = \begin{cases} (n; x_1 + b; \dots; x_n + b) & n \geq 1; \\ 0 & \text{otherwise;} \end{cases}$$

Similarly, for any  $b \geq 0$  and  $f \in \mathcal{K}_b(U)$ , let the mapping  $\mu_b : \mathcal{K}_b(U) \rightarrow \mathcal{K}_b(U)$  be defined as

$$\mu_b f(\mathbf{u}) = f(\mu_b^+ \mathbf{u}):$$

Also, for  $b \geq 0$ , let  $\mu_b \in \mathcal{M}_F(U)$  be the measure defined such that for any Borel set  $B \in \mathcal{B}(U)$ , we have

$$\mu_b(B) = \mu_b^+(B):$$

For  $\mu \in \mathcal{M}_F(U)$ , the measure  $\mu_b \in \mathcal{M}_F(U)$  satisfies

$$h_{\mu_b}(\cdot; f) = h_{\mu}(\cdot; \mu_b f) \tag{2.2}$$

for all  $f \in \mathcal{K}_b(U)$  and the existence of the unique measure  $\mu_b$  follows from the Riesz-Markov-Kakutani theorem [69, Theorem 2.14] stated in Theorem A.4. By using (2.2), we map a change in the system state in a given small time interval to an equivalent change in the function  $f$ . By working with the class of functions of the type  $\mu_t^{-(N)} \in \mathcal{K}_t^{-(N)}$ ;  $i$  for  $\mu_t^{-(N)} \in \mathcal{C}_b^1(U)$ , we obtain the generator of the Markov process  $(\mu_t^{-(N)}; t \geq 0)$ .

## 2.3 System Dynamics

In this section, we first define a system state descriptor. Using this, we then describe the time evolution of the system.

We index a sequence of systems by  $N$  that denotes the total numbers of servers. Incoming jobs arrive according to a Poisson process with rate  $N\lambda$ , and the job lengths are i.i.d. from a common distribution  $G(\cdot)$  defined on  $\mathbb{R}_+$ . The state of a server is written as  $\mathbf{a}_n = (n; a_1; \dots; a_n) \in \mathcal{U}$  when there are  $n$  progressing jobs and the  $i^{\text{th}}$  job has age  $a_i$  for  $1 \leq i \leq n$ . A server with a state  $\mathbf{a}_n$  can be viewed as a particle with the given state. Therefore the system evolution can be considered as the evolution of a system with  $N$  particles where the interactions between particles take place while routing an arrival according to the SQ( $d$ ) load balancing policy.

The age of a progressing job indicating the time elapsed since its arrival increases linearly with time at a unit rate until its service is completed. We next describe the possible state of a server at time  $t+h$  ( $h > 0$ ) given that it has state  $\mathbf{a}_n$  at time  $t$ . When  $h$  is small enough, in the interval  $[t; t+h)$ , the probability of having multiple events of arrivals or departures is negligible. In the interval  $[t; t+h)$ , if there is no arrival or departure at the given server, then the server's state will be equal to  $\mathbf{a}_n^+$  at time  $t+h$ . On the other hand, if the  $i^{\text{th}}$  job expires in the interval  $[t; t+h)$ , then the server's state will be equal to  $\mathbf{a}_n^i$  at time  $t+h$ . Considering arrivals, suppose there is an arrival that is accepted at the given server at time  $t+r$  ( $0 \leq r < h$ ), then the arriving job chooses its position uniformly at random out of  $n+1$  possible positions. If the arrived job chooses the  $j^{\text{th}}$  position, then the server's state will be equal to  $(\mathbf{a}_n^+)^j; h-r$  at time  $t+h$ .

Let  $\mathbf{S}_i^{(N)}(t) \in \mathcal{U}$  be the random variable that indicates the state of the server  $i$  at time  $t$ . Although one can think of considering  $(\mathbf{S}_1^{(N)}; \dots; \mathbf{S}_N^{(N)}(t))$  to denote the system state at time  $t$  which is a Markovian representation of the system, the dimension of this state space increases with  $N$  as  $N! - 1$  which is inconvenient to work with since our focus of interest is to study the asymptotic behavior of the system as  $N \rightarrow \infty$ . Hence, we consider an alternative simple system state representation that can be used to describe the system evolution as the evolution of a Markov process. Note that the system is symmetric to the servers as they have the same parameters, and the servers' identities do not play any role in the time evolution of the system. Therefore, to model the system evolution by a Markov process, we will show that it is enough to keep track of the number of servers that lie in each state  $\mathbf{u} \in \mathcal{U}$  in order to establish the mean-field limit. Measure-valued Markov processes have also been used to study other interacting particle systems as in [70–72] where the state of each particle lies in the space  $\mathbb{R}^n$ ,  $n > 1$ . Following these works, we consider the following system state descriptor.

**Definition 2.3.1.** *System state descriptor: At time  $t$ , the state descriptor of the system*

with index  $N$  is a random measure given by

$$\mu_t^{(N)} = \sum_{i=1}^N \mathbf{s}_i^{(N)}(t). \quad (2.3)$$

The interpretation of  $\mu_t^{(N)}$  is that for any measurable function  $f$  defined on  $\mathbb{U}$ , we have

$$\int \mu_t^{(N)}(f) = \sum_{i=1}^N f(\mathbf{S}_i^{(N)}(t)).$$

At time  $t$ , conditioned on servers' states say  $\mathbf{S}_i^{(N)}(t) = \mathbf{s}_i(t)$ ,  $i = 1, \dots, N$ , the state of the system can be represented by a measure defined as

$$\mu_t = \sum_{i=1}^N \delta_{\mathbf{s}_i(t)}. \quad (2.4)$$

For  $\mu_t^{(N)} = \mu_t$ , an element  $\mathbf{y} \in \mathbb{U}$  is an atom of  $\mu_t$  if there exists at least one server with the state  $\mathbf{y}$  at time  $t$ . The mass of an atom of  $\mu_t$  is equal to the number of servers lying at that atom at time  $t$ . As a result, since the number of interacting particles in the system is equal to  $N$ , the measure  $\mu_t$  defined on  $\mathbb{U}$  contains a finite number of atoms which is bounded by  $N$ . If all the servers lie in different states, then the number of atoms is equal to  $N$ . Otherwise, the number of atoms is less than  $N$ . For  $\mu_t$ , let  $V(t)$  be the number of atoms at time  $t$  and let the  $i^{\text{th}}$  atom be denoted by  $\mathbf{v}^{(i)}(t)$ . Further, let the mass of the atom  $\mathbf{v}^{(i)}(t)$  be denoted by  $a^{(i)}(t)$ . Here,  $a^{(i)}(t)$  denotes the number of servers that lie in the state  $\mathbf{v}^{(i)}(t)$  at time  $t$  where  $a^{(i)}(t) \geq 1$ . For given time  $t$ , from (2.4), we can also write as

$$\mu_t = \sum_{i=1}^{V(t)} a^{(i)}(t) \delta_{\mathbf{v}^{(i)}(t)}. \quad (2.5)$$

For any Borel set  $B \in \mathcal{B}(\mathbb{U})$ , the number of servers with ages lying in the set  $B$  is equal to  $\mu_t^{(N)}(B) = \int_B \mu_t^{(N)}$ . We now define the measure of an element  $\mathbf{y}_n = (n; y_1; \dots; y_n)$  as below. Let  $B(\mathbf{y}_n) = \{ (n; r_1; \dots; r_n) : y_i \leq r_i < y_i + \Delta y_i, 1 \leq i \leq n \}$ . Then as in [27], we define

$$(\hat{f}_{\mathbf{y}_n} g) = \lim_{\Delta \mathbf{y}_n \rightarrow 0} \int_{B(\mathbf{y}_n)} g(\mathbf{y}_n) \mu_t^{(N)}(\mathbf{y}_n). \quad (2.6)$$

Essentially,  $(\hat{f}_{\mathbf{y}_n} g)$  indicates the number of servers with state  $\mathbf{y}_n$  at time  $t$  and can be viewed as an occupation count. The notation  $d(\mathbf{y}_n)$  denotes the number of servers with state lying in the interval  $[\mathbf{y}_n; \mathbf{y}_n + d\mathbf{y}_n)$ , where  $d\mathbf{y}_n = (dy_1; \dots; dy_n)$  and  $\mathbf{y}_n + d\mathbf{y}_n$  is the

vector addition of  $\mathbf{y}_n$  and  $d\mathbf{y}_n$ . If there is no server lying in the state  $\mathbf{y}_n$  at time  $t$ , then  $(f\mathbf{y}_ng) = 0$ , otherwise  $\mathbf{y}_n$  is an atom with mass  $(f\mathbf{y}_ng)$ . Note that the number of servers that have  $n$  progressing jobs at time  $t$  is given by  $(U_n) = h ; l_{fU_ng}!$

We now obtain the probability that the destination server of an arrival lies in a particular state.

**Lemma 2.1.** *At time  $t$ , given that the system state is  $\mathbf{z}_n$ , i.e.,  $\bar{N}_t = \mathbf{z}_n$ , under the SQ( $d$ ) load balancing policy, the probability that the destination server of an arrival at time  $t$  lies in the state  $\mathbf{z}_n = (n; z_1; \dots; z_n)$  where  $\mathbf{z}_n$  is an atom of  $\mathbf{z}$  is given by*

$$p_r(\mathbf{z}_n) = \left\{ \frac{(f\mathbf{z}_ng)}{N} \right\} \frac{f(\bar{R}_n(\bar{N}))^d}{f\bar{R}_n(\bar{N})} \frac{(\bar{R}_{n+1}(\bar{N}))^d g}{\bar{R}_{n+1}(\bar{N})g}; \quad (2.7)$$

where

$$\bar{R}_n\left(\frac{U_j}{N}\right) = \sum_{j:n \leq j} \frac{(U_j)}{N} \quad (2.8)$$

denotes the fraction of the servers with at least  $n$  jobs.

*Proof.* When a potential destination server is chosen uniformly at random from  $N$  servers, it will have state  $(n; z_1; \dots; z_n)$  with probability  $\frac{(f(n; z_1; \dots; z_n)g)}{N}$ . Suppose out of the  $d$  potential destination servers, say  $j$  servers have occupancy  $n$  and the remaining  $d - j$  servers have occupancy at least  $n + 1$ . Further, out of the  $j$  ( $j - 1$ ) potential destination servers with occupancy  $n$ , assume  $r$  ( $r - 1$ ) servers lie in the state  $\mathbf{z}_n$ . Then the probability that the destination server is a server with state  $\mathbf{z}_n$  is given by

$$\binom{d}{j} \binom{j}{r} \binom{r}{j} \left( \frac{(f\mathbf{z}_ng)}{N} \right)^r \left( \frac{(fU_ng)}{N} \frac{(f\mathbf{z}_ng)}{N} \right)^{j-r} \left( \sum_{i:n+1}^c \frac{(U_i)}{N} \right)^{d-j};$$

Finally, by summing over all the possible values of  $j$  ( $j - 1$ ) and  $r$  ( $r - 1$ ), we have

$$\begin{aligned} & \sum_{j=1}^d \sum_{r=1}^j \binom{d}{j} \binom{j}{r} \binom{r}{j} \left( \frac{(f\mathbf{z}_ng)}{N} \right)^r \left( \frac{(fU_ng)}{N} \frac{(f\mathbf{z}_ng)}{N} \right)^{j-r} \left( \sum_{i:n+1}^c \frac{(U_i)}{N} \right)^{d-j} \\ &= \sum_{j=1}^d \binom{d}{j} \frac{1}{j} \left( \sum_{i:n+1}^c \frac{(U_i)}{N} \right)^{d-j} \left( \frac{(fU_ng)}{N} \right)^j \\ & \quad \left[ \sum_{r=1}^j r \binom{j}{r} \left( \frac{(f\mathbf{z}_ng)}{N} \right)^r \left( \frac{(fU_ng)}{N} \frac{(f\mathbf{z}_ng)}{N} \right)^{j-r} \right]; \end{aligned}$$



The term inside the square bracket in the above equation is the average of a binomial random variable and hence, it is equal to  $j \left\{ \frac{\binom{fzng}{N}}{\binom{fUng}{N}} \right\}$ . As a result, the above expression simplifies to

$$p_r(\cdot : \mathbf{z}_n) = \left\{ \frac{\binom{fzng}{N}}{\binom{fUng}{N}} \right\} \sum_{j=1}^d \binom{d}{j} \left( \sum_{i:n+1}^c \frac{U_i}{N} \right)^{d-j} \left( \frac{fUng}{N} \right)^j :$$

We can further write

$$p_r(\cdot : \mathbf{z}_n) = \left\{ \frac{\binom{fzng}{N}}{\binom{fUng}{N}} \right\} \left[ \left( \sum_{j=0}^d \binom{d}{j} \left( \sum_{i:n+1}^c \frac{U_i}{N} \right)^{d-j} \left( \frac{fUng}{N} \right)^j \right) \left( \sum_{i:n+1}^c \frac{U_i}{N} \right)^d \right] :$$

After simplifications, we get (2.7). □

**Remark 2.3.** We can also interpret the expression of  $p_r(\cdot : \mathbf{z}_n)$  as follows: The probability that all the potential destination servers have occupancy at least  $n$  and there exists at least one potential destination server with occupancy  $n$ , is equal to  $(\bar{R}_n(\bar{N}))^d - (\bar{R}_{n+1}(\bar{N}))^d$ . From the SQ( $d$ ) policy, the probability that the destination server has occupancy  $n$  is equal to  $(\bar{R}_n(\bar{N}))^d - (\bar{R}_{n+1}(\bar{N}))^d$ . From the list of the servers with occupancy  $n$ , the fraction of the servers with the state  $\mathbf{z}_n$  is equal to  $\left\{ \frac{\binom{fzng}{N}}{\binom{fUng}{N}} \right\}$ . Therefore, the probability that the destination server lies in the state  $\mathbf{z}_n$  is equal to  $\left\{ \frac{\binom{fzng}{N}}{\binom{fUng}{N}} \right\} ((\bar{R}_n(\bar{N}))^d - (\bar{R}_{n+1}(\bar{N}))^d)$ .

For the case of exponential JLDs,  $U_n = fng$  and  $\mathbf{z}_n = n$ . Hence,  $p_r(\cdot : \mathbf{z}_n) = (\bar{R}_n(\bar{N}))^d - (\bar{R}_{n+1}(\bar{N}))^d$  coinciding with the analysis for the exponential case in [10, 12].

As it is clear from (2.7), the routing decision depends only on the number of servers lying in each possible server state. Hence, we get the evolution of the process  $(\binom{N}{t}; t \geq 0)$  by tracking arrival events, routing decisions, and departure events.

## 2.4 Summary of Main Results

In this section, we present the main results of this chapter.

Our aim is to study the limit as  $N \rightarrow \infty$  of the empirical probability measures of the states of the servers, i.e., the limit of the sequence of processes  $f(\frac{t}{N}; t \geq 0)g_{N-1}$  as  $N \rightarrow \infty$ . A system with index  $N$  has  $N$  servers that serve the incoming jobs arriving

according to a Poisson process with rate  $N^{-1}$ , and all other system parameters remain the same for all  $N$  as given in Section 2.1. For given  $N$ , the process  $(\tilde{X}_t^{(N)}; t \geq 0)$  defined in (2.3) describes the dynamics of the system with index  $N$ . The goal is to characterize the limit of the normalized process  $(\tilde{X}_t^{(N)}; t \geq 0)$  as  $N \rightarrow \infty$  where

$$\tilde{X}_t^{(N)} = \frac{X_t^{(N)}}{N}. \quad (2.9)$$

For a Borel set  $B \subseteq \mathcal{B}(U)$ ,  $\tilde{X}_t^{(N)}(B)$  is equal to fraction of the servers with state lying in the set  $B$  at time  $t$ .

### 2.4.1 An Overview of the Analysis

We now give a brief overview of the analysis of this chapter.

The mean-field limit corresponds to  $\lim_{N \rightarrow \infty} (\tilde{X}_t^{(N)}; t \geq 0) = (\bar{X}_t; t \geq 0)$  that lies in  $C_{M_1(U)}([0; \infty))$  and it is a deterministic measure-valued process satisfying a set of evolution equations referred to as the MFEs. We then obtain an alternative form of the evolution equations satisfied by the process  $(\tilde{X}_t^{(N)}; t \geq 0)$  for  $\varphi \in C_b(U)$ . This is stated in Lemma 2.2. Using these equations, we show in Theorem 2.1 that there exists a unique solution to the MFEs for a given initial point.

We then show that the sequence of processes  $\tilde{X}_t^{(N)}; t \geq 0)g_{N^{-1}}$  is tight. For this, we first study the Feller property of the Markov process  $(X_t^{(N)}; t \geq 0)$  and construct a martingale process in Theorem 2.2 by using the generator of the Markov process  $(X_t^{(N)}; t \geq 0)$  by employing the Dynkin's formula [56, Proposition 1.7, p.162].

We prove the tightness of the sequence of processes  $\tilde{X}_t^{(N)}; t \geq 0)g_{N^{-1}}$  using the result that the normalized version of the martingale process converges to the null process as  $N \rightarrow \infty$ . Furthermore, we show that any limit point of the sequence  $\tilde{X}_t^{(N)}; t \geq 0)g_{N^{-1}}$  coincides almost surely with the unique solution to the MFEs referred to as the mean-field limit. This is stated in Theorem 2.3.

Finally, we obtain a set of partial differential equations satisfied by the mean-field limit. We then prove the uniqueness of the fixed-point and its insensitivity in Theorem 2.5. The proofs of Theorems 2.3 and 2.5 are given in Sections 2.5 and 2.6, respectively. The remaining proofs are given in Section 2.8.

## 2.4.2 Transient Regime

In this section, we state the results that are related to the transient regime. We first study the proposed MFEs. For the given system parameters  $\lambda, C, d$ , and the probability density function  $g(\cdot)$  of the JLDs, we state the MFEs in Proposition 2.1. The dynamics of a mean-field solution  $(\bar{r}_t; t \geq 0)$  are described by using a set of evolution equations of the real valued processes  $(h_t^-; i; t \geq 0)$  for all  $\bar{r} \in C_b^1(\mathbb{U})$ , referred to as the MFEs.

**Proposition 2.1.** *Mean-field equations:*

For given system parameters  $(\lambda; C; d; g(\cdot))$ , the process  $(\bar{r}_t; t \geq 0)$  satisfies:

1. The mapping  $t \mapsto \bar{r}_t$  is a continuous.
2. For  $\bar{r} \in C_b^1(\mathbb{U})$ , the process  $(\bar{r}_t; t \geq 0)$  satisfies

$$\begin{aligned} h_t^-; i = h_0^-; i + \int_{s=0}^t h_s^-; r_1; i ds \\ + \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int_{\mathbb{U}_n} (x_j) (\bar{r}_s^j) (\bar{r}_s) d_s(\mathbf{x}_n) \right) ds \\ + \int_{s=0}^t \left( -\lambda \int_{\mathbb{U}} g(\bar{r}_s) \Phi_0(\bar{r}_s) (\bar{r}_s; 0) (\bar{r}_s) + \sum_{n=1}^C \sum_{i=1}^{n+1} \int_{\mathbb{U}_n} \left\{ \frac{1}{n+1} \right\} \right. \\ \left. \Phi_n(\bar{r}_s) (\bar{r}_s^i; 0) (\bar{r}_s) d_s(\mathbf{x}_n) \right) ds; \quad (2.10) \end{aligned}$$

where the index  $j$  is used to denote the position of the departing job when there are  $n$  progressing jobs and the index  $i$  denotes the position of the arriving job when there are already  $n$  progressing jobs at the server. Further,  $\Phi_n(\bar{r}_s) = \frac{f(\bar{R}_n(\bar{r}_s))^d (\bar{R}_{n+1}(\bar{r}_s))^d g}{f\bar{R}_n(\bar{r}_s) \bar{R}_{n+1}(\bar{r}_s) g}$  where  $\bar{R}_j(\bar{r}_s) = \sum_{n: n=j}^C \bar{r}_s(\mathbb{U}_n)$ .

In (2.10), the second term on the right-hand side is due to the increase of the ages of the progressing jobs linearly with time at unit rate. The third and fourth terms on the right-hand side of (2.10) are due to departure and arrival of a job, respectively.

**Remark 2.4.** The  $t$ -continuity of the mapping  $\bar{r}_t$  is equivalent to the continuity of the mapping  $t \mapsto h_t^-; i$  for all  $\bar{r} \in C_b^1(\mathbb{U})$  since  $C_b^1(\mathbb{U})$  is a separating class of  $M_1(\mathbb{U})$  [56, p. 111].

Although the MFEs (2.10) are defined for the class of functions  $\varphi \in C_b^1(U)$ , it is more useful to obtain an approximation of the process  $(h_t^{-(N)}; I_{fBg}^i; t \geq 0)$  for an open set  $B \subset B(U)$ . Therefore, we need to obtain the evolution equations of the real valued process  $(h_t^-; I_{fBg}^i; t \geq 0)$ . In this direction, we first obtain the evolution equations of the real valued process  $(h_t^-; i; t \geq 0)$  where  $\varphi \in C_b(U)$ . We then proceed to obtain the evolution equations of the process  $(h_t^-; I_{fBg}^i; t \geq 0)$  where  $B$  is an open set with the help of the monotone convergence theorem [69, Theorem 1.26] since there exists a sequence of functions in  $C_b(U)$  that increase point wise to  $I_{fBg}$ .

**Lemma 2.2.** *A process  $(\mu_t \in M_1(U); t \geq 0)$  with continuity of the mapping  $t \mapsto \mu_t$  satisfies the MFEs (2.10) if it satisfies the following equation for all  $\varphi \in C_b(U)$ ,*

$$\begin{aligned} h_t^-; i = h_0^-; i + \int_{r=0}^t & \left( \sum_{n=1}^C \sum_{j=1}^n \int \int_{U_n} (\chi_j) (\mu_{t-r}(\mathbf{x}_n^j) - \mu_{t-r}(\mathbf{x}_n)) d_r(\mathbf{x}_n) \right. \\ & + \left[ \int_r (f_0 g) \Phi_0(r) (\mu_{t-r}(1; 0) - \mu_{t-r}(0)) \right. \\ & \left. \left. + \sum_{n=1}^C \sum_{j=1}^{n+1} \int \int_{U_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(r) (\mu_{t-r}(\mathbf{x}_n^j; 0) - \mu_{t-r}(\mathbf{x}_n)) d_r(\mathbf{x}_n) \right] \right) dr. \end{aligned} \quad (2.11)$$

*Proof.* See Section 2.8.2. □

Using (2.11), we now show that starting with an initial measure  $\mu_0$ , for  $t \geq 0$ , there exists a unique measure  $\mu_t \in M_1(U)$  that satisfies (2.10).

**Theorem 2.1.** *There exists a unique solution in  $C_{M_1(U)}([0; 1])$  to the MFEs. In particular, if  $(\mu_t^1; t \geq 0)$  and  $(\mu_t^2; t \geq 0)$  are two mean-field solutions starting at initial measures  $\mu_0^1 \in M_1(U)$ ,  $\mu_0^2 \in M_1(U)$ , respectively, then*

$$\|\mu_t^1 - \mu_t^2\| \leq \|\mu_0^1 - \mu_0^2\| e^{(2Ck + 8d^2)t}. \quad (2.12)$$

*Proof.* See Section 2.8.3. □

We next focus on the proof of the convergence of the sequence of the processes  $(f_t^{-(N)}; t \geq 0)g_{N-1}$  as  $N \rightarrow \infty$ . We first give a result on the Feller property of  $(f_t^{(N)}; t \geq 0)$ . We consider the filtration

$$F_t^{(N)} = \sigma_{s \geq 0}(\mu_s^{(N)}; s \leq t): \quad (2.13)$$

Let  $A^{(N)}(\cdot)$  be the generator of the Markov process  $(f_t^{(N)}; t \geq 0)$ . Using the Dynkin's formula [56], we have the following result.

**Theorem 2.2.** *The process  $(\mathbf{h}_t^{(N)}; t \geq 0)$  is a Feller-Dynkin process in  $D_{M_F(U)}([0; \infty))$ . Let  $\mathbf{h}_0^{(N)} \in C_b^1(U)$ , then the process  $(\mathbf{M}_t^{(N)}(\mathbf{h}_0^{(N)}); t \geq 0)$  defined as*

$$\mathbf{M}_t^{(N)}(\mathbf{h}_0^{(N)}) = \mathbf{h}_t^{(N)} - \mathbf{h}_0^{(N)} - \int_{s=0}^t A^{(N)} \mathbf{h}_s^{(N)} ds \quad (2.14)$$

*is a square integrable  $F_t^{(N)}$ -martingale and it is an RCLL process.*

*Proof.* See Section 2.8.1. □

Our analysis requires the following assumptions.

**Assumption 2.4.1.** *The sequence of the initial random measures  $\mathbf{f}_0^{(N)} g_{N-1}$  satisfy*

$$(\mathbf{f}_0^{(N)}; \mathbf{h}_0^{(N)}; i) \in (\#; \mathbf{h}_\#; i); \quad (2.15)$$

*where  $\# \in M_1(U)$  is a probability measure that possesses a density (w.r.t. the Lebesgue measure) and  $\mathbf{h}_\#; i < 1$ .*

Here, we can interpret the condition  $\mathbf{h}_\#; i < 1$  as follows. If  $\#$  is the probability measure of the state of a server, then the average sum of ages of jobs is finite.

We use Theorem 2.2 and the Assumption 2.4.1 to establish the following main result.

**Theorem 2.3.** *Under the Assumptions 2.4.1 and 2.1.1, we show that  $(\mathbf{h}_t^{(N)}; t \geq 0)$  converges to  $(\mathbf{h}_t; t \geq 0)$  as  $N \rightarrow \infty$ , where  $(\mathbf{h}_t; t \geq 0)$  is the unique solution to (2.10) with the initial point  $\#$ . The process  $(\mathbf{h}_t; t \geq 0)$  is referred to as the mean-field limit.*

*Proof.* See Section 2.5. □

The existence of the mean-field limit allows us to show that any finite subset of servers become independent of each other in the limiting system.

**Theorem 2.4.** *If  $\mathbf{f}_{\mathbf{S}_k}^{(N)}(0); 1 \leq k \leq N$  are exchangeable and under the Assumptions 2.1.1 and 2.4.1, the following result holds:*

- *For each fixed  $k$  and  $t \in [0; \infty)$ ,  $L(\mathbf{S}_k^{(N)}(t)) \rightarrow \mathbf{h}_t$  as  $N \rightarrow \infty$ .*
- *For any fixed positive integer  $l$  and for each  $t \in [0; \infty)$ , we have  $\mathbf{f}_{\mathbf{S}_k}^{(N)}(t); 1 \leq k \leq l$  are independent random variables with  $L(\mathbf{S}_k(t))$  equal to  $\mathbf{h}_t$  for all  $1 \leq k \leq l$ .*

*Proof.* See Section 2.8.4. □

**Remark 2.5.** For any time  $t$ , a consequence of Theorem 2.4 is that as  $N \rightarrow \infty$ , due to Poisson thinning, the arrival process to each server is a Poisson process with rate  $\Phi_n(\bar{\nu}_t)$  when there are  $n \geq 0$  progressing jobs.

**Lemma 2.3.** For any time  $t$ , the measure  $\bar{\nu}_t$  has a density function w.r.t: the Lebesgue measure for almost all  $\mathbf{u}_n \geq \mathbb{U}_n$ ,  $n \geq 1$ .

*Proof.* See Section 2.8.6. □

For a given subset  $B \subseteq \mathbb{B}(\mathbb{U})$ , once we show that  $(\bar{\nu}_t^{(N)}; t \geq 0) \rightarrow (\bar{\nu}_t; t \geq 0)$  as  $N \rightarrow \infty$ , since  $\bar{\nu}_t$  is absolutely continuous w.r.t: Lebesgue measure for every  $t \geq 0$ , the continuous mapping theorem [57, Theorem 2.7] implies that  $(h_t^{(N)}; I_{fBg}; t \geq 0) \rightarrow (h_t; I_{fBg}; t \geq 0)$ . This shows that for large  $N$ , we can approximate  $h_t^{(N)}; I_{fBg}$  by  $h_t; I_{fBg}$ .

### 2.4.3 Stationary Regime

We now present the results related to the stationary behavior of the mean-field.

We first demonstrate an analogy between the MFEs of the considered multi-server Erlang loss system under the SQ( $d$ ) load balancing policy and the dynamics of a single server Erlang loss system with state-dependent arrival rates. We then exploit this analogy to prove the uniqueness of the fixed-point of the mean-field and its insensitivity. We first recall the dynamics of the probability measure of the server's state of a single server Erlang loss system with capacity  $C$ , where jobs arrive according to a Poisson process with pre-specified state-dependent arrival rates.

Consider a single server system with capacity  $C$  where jobs arrive according to a Poisson process at rate  $\lambda_n$  when there are  $n$  progressing jobs in the system. The service times are generally distributed as stated in the system model of Section 2.1. Let  $\bar{\nu}_t^{(single)}$  be the probability measure of the server's state at time  $t$  defined on  $\mathbb{U}$ . For  $\bar{\nu}_t \in \mathcal{C}_b^1(\mathbb{U})$ , it can be verified that the Kolmogorov equations are given by,

$$\begin{aligned} \dot{h}_t^{(single)}; i &= h_0^{(single)}; i + \int_{s=0}^t h_s^{(single)}; r_{-1}(i) ds \\ &\quad - \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int_{\mathbb{U}_n} (x_j) (x_n^j) - (x_n) \right) d_s^{(single)}(x_n) \end{aligned}$$

$$+ \left[ \left( \int_0^s \int_{\mathbb{U}_n} \left\{ \frac{1}{n+1} \right\} \right. \right. \\ \left. \left. \int_{\mathbb{U}_n} \left\{ \frac{1}{n+1} \right\} \right) d_s^{(single)}(\mathbf{x}_n) \right] ds: \quad (2.16)$$

On comparing the MFEs (2.10) with the Kolomogorov equation of a single-server system given by (2.16), it is clear that both the dynamics are similar except that  $\mathbb{U}_n$  in (2.16) is replaced by  $\Phi_n(\bar{s})$  when the probability measure of the server's state is  $\bar{s}$  at time  $s$ . Equation (2.10) differs from (2.16) in that the arrival rates at time  $t$  depend on  $\bar{t}$ . This is an example of a non-linear Markov process which means that the generator of the Markov process at time  $t$  depends on the current distribution  $\bar{t}$  of the Markov process [73] while the equation (2.16) for fixed  $(i; 0 \leq i \leq C)$  denotes a Markov process whose generator does not depend on the current distribution.

We now study the fixed-point of the mean-field. Let  $P_t(0)$  be equal to  $\int_0^t (f_0 g)$  and let  $\rho_t(\mathbf{x}_n)$  be the probability density of  $\int_0^t w(r; t)$  Lebesgue measure at  $\mathbf{x}_n$ . We obtain the differential equations satisfied by the process  $(P_t; t \geq 0)$  with  $P_t = (P_t(\mathbf{u}); \mathbf{u} \in \mathbb{U})$  where

$$P_t(\mathbf{y}_n) = \int_{x_1=0}^{y_1} \cdots \int_{x_n=0}^{y_n} \rho_t(\mathbf{x}_n) dx_1 \cdots dx_n: \quad (2.17)$$

Here, from the Remark 2.5, since  $\bar{t}$  is the distribution of a server's state as  $N \rightarrow \infty$ , it implies that  $P_t(\mathbf{y}_n)$  is the probability that a server has  $n$  jobs and the  $i^{\text{th}}$  job's age is at most  $y_i$  for  $1 \leq i \leq n$  as  $N \rightarrow \infty$ . Also, since  $\bar{t}^{(N)}(\cdot) \rightarrow \bar{t}(\cdot)$ , for a large value of  $N$ , the fraction of servers with  $n$  jobs and the  $i^{\text{th}}$  job's age is at most  $y_i$  for  $1 \leq i \leq n$  can be approximated by  $P_t(\mathbf{y}_n)$ .

**Lemma 2.4.** *The process  $(P_t; t \geq 0)$  satisfies*

$$\frac{dP_t(0)}{dt} = \int_{y=0}^1 (y) \left( \frac{\partial P_t(1; y)}{\partial y} \right) dy - \Phi_0(P_t)P_t(0); \quad (2.18)$$

for  $1 \leq n \leq C-1$ ,

$$\begin{aligned} \frac{dP_t(\mathbf{y}_n)}{dt} = & \sum_{i=1}^n \frac{\partial P_t(\mathbf{y}_n)}{\partial y_i} + \sum_{j=1}^{n+1} \int_{x_j=0}^1 (x_j) \left( \frac{\partial P_t(\mathbf{y}_n^j; x_j)}{\partial x_j} \right) dx_j \\ & \sum_{j=1}^n \int_{x_j=0}^{y_j} (x_j) \left( \frac{\partial P_t(\mathbf{y}_n^j; x_j)}{\partial x_j} \right) dx_j \\ & + \sum_{j=1}^n \frac{\Phi_{n-1}(P_t)}{n} P_t(\mathbf{y}_n^j) - \Phi_n(P_t)P_t(\mathbf{y}_n); \quad (2.19) \end{aligned}$$

and for  $n = C$ ,

$$\begin{aligned} \frac{dP_t(\mathbf{y}_n)}{dt} = & \sum_{i=1}^n \frac{\partial P_t(\mathbf{y}_n)}{\partial y_i} \sum_{j=1}^n \int_{x_j=0}^{y_j} (x_j) \left( \frac{\partial P_t(\mathbf{y}_n^j; x_j)}{\partial x_j} \right) dx_j \\ & + \sum_{j=1}^n \frac{\Phi_{n-1}(P_t)}{n} P_t(\mathbf{y}_n^j); \end{aligned} \quad (2.20)$$

where  $\Phi_n(P_t) = \frac{f(R_n(P_t))^d (R_{n+1}(P_t))^d g}{fR_n(P_t) R_{n+1}(P_t)g}$  and  $R_n(P_t) = \sum_{j:j=n}^C \lim_{b \uparrow} P_t(j; b; \dots; b)$ .

*Proof.* See Section 2.8.5. □

**Remark 2.6.** Specializing the results to the exponential case with mean  $\frac{1}{\lambda}$ ,  $\bar{G}(x) = e^{-\lambda x}$ , and denoting  $q_n(t) = \lim_{b \uparrow} P_t(n; b; \dots; b)$ , it can be verified that the process  $(\mathbf{q}(t); t \geq 0)$  where  $\mathbf{q}(t) = (q_n(t); 0 \leq n \leq C)$  is the unique solution of the MFEs given in [10] for the case of the exponential distributions with rate  $\lambda = 1$ .

We next state the the principal result on the insensitivity of the fixed-point of the MFEs.

**Theorem 2.5.** The process  $(P_t; t \geq 0) = (P_t(\mathbf{u}); \mathbf{u} \in \mathcal{U}; t \geq 0)$  has a unique fixed-point given by  $\bar{\mathbf{y}} = (\bar{y}_n; \mathbf{y} \in \mathcal{U})$  where

$$\bar{y}_n = \frac{(\text{exp})}{n} \prod_{i=1}^n \int_{x_i=0}^{y_i} \bar{G}(x_i) dx_i; \quad (2.21)$$

and  $\frac{(\text{exp})}{n} = (\frac{(\text{exp})}{n}; 0 \leq n \leq C)$  denotes the unique fixed-point of the mean-field when the service times are exponentially distributed with the mean  $\frac{1}{\lambda}$  and  $\frac{(\text{exp})}{n}$  is the stationary probability that there are  $n$  jobs in the limiting system. Further, since  $\int_{x=0}^{\infty} \bar{G}(x) dx = \frac{1}{\lambda}$ , the fixed-point of the mean-field is insensitive, i.e.,

$$\lim_{b \uparrow} (n; b; \dots; b) = \frac{(\text{exp})}{n}; \quad (2.22)$$

*Proof.* See Section 2.6. □

## 2.5 Proof of Theorem 2.3

From Theorem 2.2, we now show that the normalized process  $(\frac{N}{t}; t \geq 0)$  converges to the mean-field limit.



Let  $(\bar{F}_t^{(N)}; t \geq 0)$  be the right continuous filtration associated with the process  $(x_t^{(N)}; t \geq 0)$ . Note that we have  $(x_t^{(N)}; t \geq 0) \in \mathcal{D}_{M_1^1(\mathbb{U})}([0; 1])$ . We first show that the sequence of processes  $(x_t^{(N)}; t \geq 0)$  is relatively compact. We then prove that every limit point  $(x_t; t \geq 0)$  almost surely has sample paths that are continuous in  $t$  and furthermore, they coincide with the unique mean-field solution with the initial point  $\#$ . For every limit point  $(x_t; t \geq 0)$ ,  $\mu_0$  almost surely coincides with the measure  $\#$  from the Assumption 2.4.1. Further, we have that the mean-field solution is unique for the given initial measure. Hence, we conclude that for all the limit points, almost surely sample paths coincide with the unique mean-field solution  $(x_t; t \geq 0)$  with the initial point  $\#$ . The process  $(x_t; t \geq 0)$  is referred to as the mean-field limit. As a result,  $(x_t^{(N)}; t \geq 0)$  converges in distribution to the mean-field limit  $(x_t; t \geq 0)$ .

For  $\varphi \in C_b^1(\mathbb{U})$ , from Proposition 2.3, the process  $(\bar{M}_t^{(N)}(\varphi); t \geq 0)$  defined as follows is an RCLL square integrable  $\bar{F}_t^{(N)}$ -martingale

$$\begin{aligned} \bar{M}_t^{(N)}(\varphi) = & \varphi(x_t^{(N)}) - \varphi(x_0^{(N)}) - \int_{s=0}^t \varphi'(x_s^{(N)}) r_{1,i} ds - \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int_{\mathbb{U}_n} \varphi(x_j) \right. \\ & \left. - \varphi(x_n) \right) d^{-s(N)}(x_n) \\ & + \left[ \varphi(x_t^{(N)}) - \varphi(x_0^{(N)}) - \int_{s=0}^t \varphi'(x_s^{(N)}) r_{1,i} ds - \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{\mathbb{U}_n} \left\{ \frac{1}{n+1} \right\} \right. \right. \\ & \left. \left. \varphi_n(x_n^{(N)}; 0) - \varphi(x_n) \right) d^{-s(N)}(x_n) \right] ds: \quad (2.23) \end{aligned}$$

We further have

$$\begin{aligned} \langle \bar{M}^{(N)}(\varphi) \rangle_t = & \frac{1}{N} \left[ \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int_{\mathbb{U}_n} \varphi(x_j) - \varphi(x_n) \right)^2 d^{-s(N)}(x_n) \right. \\ & + \left[ \varphi(x_t^{(N)}) - \varphi(x_0^{(N)}) - \int_{s=0}^t \varphi'(x_s^{(N)}) r_{1,i} ds - \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{\mathbb{U}_n} \left\{ \frac{1}{n+1} \right\} \right. \right. \\ & \left. \left. \varphi_n(x_n^{(N)}; 0) - \varphi(x_n) \right)^2 d^{-s(N)}(x_n) \right] ds \Big]: \quad (2.24) \end{aligned}$$

Since the space  $\mathcal{D}_{M_1(\mathbb{U})}([0; 1])$  endowed with the Skorohod topology is complete and separable, by using the Prohorov's theorem (Theorem A.2) [57], establishing the relative compactness of the sequence of the processes  $\varphi(x_t^{(N)}; t \geq 0)g_{N-1}$  is equivalent to proving the tightness of the processes  $\varphi(x_t^{(N)}; t \geq 0)g_{N-1}$ . From Theorem 4.6 of [74], the Jakubowski's

criteria stated in Appendix A.3 can be used to establish the relative compactness of the sequence of the processes  $f(h_t^{-(N)}; t \geq 0)g_{N-1}$ . According to the Jakubowski's criteria, we need to show that the conditions J1 and J2 are true.

We next focus on the proof of the condition J2. For this, we prove the conditions C1 and C2 that are sufficient to prove the relative compactness of the sequence  $f(h_t^{-(N)}; i; t \geq 0)g_{N-1}$  for  $\varphi \in C_b^1(\mathbb{U})$  in  $D_{\mathbb{R}}([0; T])$ . For any  $T > 0$ ,  $t \geq [0; T]$ , we have  $h_t^{-(N)}; i = k_1 h_t^{-(N)}; \mathbf{1} i$  and since  $h_t^{-(N)}; \mathbf{1} i = 1$ , the condition C1 is trivially satisfied with  $b = k_1$ .

We next prove that the condition C2 holds. For  $\epsilon > 0$ , by using (2.24) and the Doob's inequality (Theorem A.8) [56, page 63], we have

$$\begin{aligned} \mathbb{P} \left( \sup_{t \geq T} \left| \overline{\mathbf{M}}_t^{(N)}(\varphi) \right| > \epsilon \right) &\leq \frac{4}{2} \mathbb{E} \left[ \langle \overline{\mathbf{M}}^{(N)}(\varphi) \rangle_T \right] \\ &\leq 4Tk^2 \frac{1}{N} (k_1 + d); \end{aligned}$$

and hence,  $\mathbb{P} \left( \sup_{t \geq T} \left| \overline{\mathbf{M}}_t^{(N)}(\varphi) \right| > \epsilon \right) \rightarrow 0$  as  $N \rightarrow \infty$ . Therefore, the sequence of processes  $f(\overline{\mathbf{M}}_t^{(N)}(\varphi); t \geq 0)g_{N-1}$  converges in distribution to the null process from the standard convergence criterion in  $D_{\mathbb{R}}([0; T])$  [56, Theorem 1.4, p.339]. Furthermore, the sequence of processes  $f(\overline{\mathbf{M}}_t^{(N)}(\varphi); t \geq 0)g_{N-1}$  is tight in  $D_{\mathbb{R}}([0; T])$  and hence, there exists  $\epsilon^0 > 0$  and  $N^0 > 0$  such that for all  $N \geq N^0$ , we have

$$\mathbb{P} \left( \sup_{u, v \in [T; T+\epsilon^0]} \left| \overline{\mathbf{M}}_v^{(N)}(\varphi) - \overline{\mathbf{M}}_u^{(N)}(\varphi) \right| > \frac{\epsilon}{2} \right) \leq \frac{\epsilon}{2} \quad (2.25)$$

For any  $u < v \in [T; T+\epsilon^0]$ , from (2.23), we have

$$\begin{aligned} \left| h_v^{-(N)}; i - h_u^{-(N)}; i \right| &\leq \int_{s=u}^v \left| h_s^{-(N)}; r_1 i \right| ds + 2k_1 k_1 C j u - v j + 2k_1 k_1 j u - v j \\ &\quad + \left| \overline{\mathbf{M}}_v^{(N)}(\varphi) - \overline{\mathbf{M}}_u^{(N)}(\varphi) \right|; \end{aligned} \quad (2.26)$$

Further, we can write

$$\left| h_v^{-(N)}; i - h_u^{-(N)}; i \right| \leq j v - u j C k_1 (1 + 2k_1 k_1 + 2d) + \left| \overline{\mathbf{M}}_v^{(N)}(\varphi) - \overline{\mathbf{M}}_u^{(N)}(\varphi) \right|; \quad (2.27)$$

Therefore, by using (2.25) and (2.27), there exists  $\epsilon > 0$  and  $N_1 > 0$  such that for  $N \geq N_1$ , we have  $\mathbb{P} \left( \sup_{u, v \in [T; T+\epsilon^0]} \left| h_v^{-(N)}; i - h_u^{-(N)}; i \right| > \epsilon \right) \rightarrow 0$ : This proves the condition C2. Since the conditions C1 and C2 hold, the condition J2 also holds.

We next prove the compact containment condition J1. Let  $(n_i(t); x_{i1}(t) \dots; x_{in_i(t)}(t))$  be the state of the  $i^{\text{th}}$  server at time  $t$  where  $x_{ij}(t)$  denotes the age of the  $j^{\text{th}}$  job at the  $i^{\text{th}}$  server. Clearly, we have  $h_t^{-(N)}; i = \frac{1}{N} \sum_{i=1; n_i(t) > 0}^N (x_{i1}(t) + \dots + x_{in_i(t)}(t))$ :

We can classify the progressing jobs into two classes. The jobs that are in service from the beginning ( $t = 0$ ) form the first class and the second class of jobs are the ones that entered the system in the interval  $(0; t]$ . At a server, the number of progressing jobs that belong to a class is upper bounded by  $C$ . Let  $\mathbf{Y}_t$  be a random variable representing the age of a job belonging to the second class that is in progress at time  $t$ , and  $\mathbf{Y}$  be a random variable with job length distribution  $G$ , then for any  $b \geq 0$ , we have

$$P(\mathbf{Y}_t \leq b) = P(\mathbf{Y} \leq b); \quad (2.28)$$

Therefore, using (2.28), since each server has capacity  $C$ , for any time  $t \geq 0$ , we can write

$$P(h_t^{-(N)}; i \leq b) = P\left(h_0^{-(N)}; i + Ct + \frac{1}{N} \sum_{i=1}^N (\mathbf{Y}_{i1} + \dots + \mathbf{Y}_{iC}) \leq b\right); \quad (2.29)$$

where  $(\mathbf{Y}_{ij}; 1 \leq i \leq N; 1 \leq j \leq C)$  are i.i.d. random variables with distribution  $G$ . Further, by the weak law of large numbers, we have  $\frac{1}{N} \sum_{i=1}^N (\mathbf{Y}_{i1} + \dots + \mathbf{Y}_{iC}) \xrightarrow{c} \mathbb{E}(\mathbf{Y})$  as  $N \rightarrow \infty$ . Therefore, by choosing  $Z_T = 2h_0; i + 2CT + \frac{2C}{N}$ , we have

$$P\left(\sup_{t \in [0; T]} h_t^{-(N)}; i > Z_T\right) \rightarrow 0 \quad (2.30)$$

as  $N \rightarrow \infty$ . Let us define  $L_T = \{f \in M_1(U) : h; i \leq Z_T\}$ . Since  $h; i \leq Z_T$  for  $f \in L_T$ , let  $B_r = \cup_{0 \leq t \leq r} \{f \in M_1(U) : \sup_{t \in [0; r]} h; i \leq Z_T\}$  and  $\bar{B}_r$  be the complement of  $B_r$ , then we have  $P(\bar{B}_r) \leq \frac{Z_T}{r}$ . Hence,  $\lim_{r \rightarrow \infty} P(\bar{B}_r) = 0$ . Therefore, from Lemma A7.5 of [75],  $L_T$  is relatively compact in  $M_1(U)$ . Furthermore, from (2.30), we have  $\liminf_{N \rightarrow \infty} P(L_T \cap \{h; i \leq Z_T\}) > 1 - \epsilon$ . Let  $K_T$  be the closure of  $L_T$ , then we have a compact set  $K_T$  such that  $\liminf_{N \rightarrow \infty} P(K_T \cap \{h; i \leq Z_T\}) > 1 - \epsilon$  for all  $0 < \epsilon < 1$ .

This establishes the condition J1 and hence, the proof of the tightness of the sequence of processes  $f_t^{-(N)}; t \geq 0)g_{N-1}$  is completed.

Let  $(t_k; t_k \geq 0)$  be a limit point of a converging subsequence  $f_t^{-(N_{i_k})}; t \geq 0)g_{k-1}$ . From the condition C2,  $t_k$  is continuous in  $t$ ,  $P = a.s.$ , where  $P$  is the probability law of  $(t; t \geq 0)$ . Furthermore, from [74, Theorem 1.7] for  $f \in C_b(U); t \in M_1(U)$ , it follows that for any  $T > 0$ , we have  $(t; 0 \leq t \leq T) \xrightarrow{d} (h; t; f; 0 \leq t \leq T)$  is continuous in the Skorokhod topology. Then since the sequence of martingales  $f(\bar{\mathbf{M}}_t^{(N_{i_k})}(\cdot)); t \geq 0)g_{k-1}$  converges to the null process, by the continuous mapping theorem [57, Theorem 2.7], we conclude

$$\begin{aligned}
h_t; i = h_0; i + \int_{s=0}^t h_s; r_1 i ds \\
\int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int \int_{\cup_n} (x_j) ( (\mathbf{x}_n^j) (\mathbf{x}_n) ) d_s(\mathbf{x}_n) \right. \\
\left. + \left[ ( s(f_0g) \Phi_0( s) ( (1;0) (0) ) + \sum_{n=1}^C \sum_{j=1}^{n+1} \int \int_{\cup_n} \left\{ \frac{1}{n+1} \right\} \right. \right. \\
\left. \left. \Phi_n( s) ( (\mathbf{x}_n^j; 0) (\mathbf{x}_n) ) d_s(\mathbf{x}_n) \right] \right) ds: \quad (2.31)
\end{aligned}$$

From the Assumption 2.4.1,  $\rho = \#$  almost surely and hence, the sample paths coincide almost surely with the unique mean-field solution with the initial point  $\#$ . This argument holds for every limit point, and hence, the sample paths of every limit point are almost surely the same as the deterministic mean-field solution with the initial point  $\#$ . This completes the proof.

## 2.6 Proof of Theorem 2.5

We now show that  $\rho = (\mathbf{u}; \mathbf{u} \geq \cup)$  is the unique fixed-point of the mean-field. From [10], we first recall that under the assumption of exponential service time distributions, there exists a unique probability measure of occupancy  $\rho^{(exp)} = (\rho_n^{(exp)}; 0 \leq n \leq C)$  on  $f_0; 1; \dots; Cg$  to the stationary MFEs given below,

$$\rho_n^{(exp)} (\rho^{(exp)}) \rho_n^{(exp)} = (n+1) \rho_{n+1}^{(exp)}; \quad (2.32)$$

where

$$\rho_n^{(exp)} (\rho^{(exp)}) = \frac{f(\sum_{j=n}^C \rho_j^{(exp)})^d (\sum_{j=n+1}^C \rho_j^{(exp)})^d g}{f(\sum_{j=n}^C \rho_j^{(exp)}) (\sum_{j=n+1}^C \rho_j^{(exp)}) g}; \quad (2.33)$$

Let  $\rho = (\mathbf{u}; \mathbf{u} \geq \cup)$  be a fixed-point of the MFEs of the process  $(P_t; t \geq 0)$  under general service time distributions. Using  $\rho$ , let the corresponding probability measure of occupancy be  $\mathbf{\Gamma} = (\Gamma_n; 0 \leq n \leq C)$  defined such that  $\Gamma_n = \lim_{b \uparrow \infty} \rho(n; b; \dots; b)$  and  $\Gamma_0 = \rho(0)$ . We now show that

$$(\mathbf{y}_n) = \frac{\left\{ \prod_{i=1}^n \frac{\rho_i^{(exp)}(\mathbf{\Gamma})}{i} \right\}}{1 + \sum_{m=1}^C \left\{ \prod_{i=1}^m \frac{\rho_i^{(exp)}(\mathbf{\Gamma})}{i} \right\}} n \prod_{i=1}^n \int_{x_i=0}^{y_i} \bar{G}(x_i) dx_i \quad (2.34)$$

and

$$(0) = \frac{1}{1 + \sum_{m=1}^C \left\{ \prod_{i=1}^m \frac{{}^{(exp)}(\mathbf{\Gamma})}{i} \right\}}. \quad (2.35)$$

Then it implies that  $\mathbf{\Gamma}$  also satisfies (2.32)-(2.33), and hence,  $\mathbf{\Gamma} = {}^{(exp)}$  concluding the insensitivity of the fixed-point. Furthermore, we have that  $(\mathbf{y}_n) = {}_n^{(exp)} \prod_{i=1}^n \int_{x_i=0}^{y_i} \bar{G}(x_i) dx_i$  concluding the uniqueness of the fixed-point of the mean-field under general service time distributions.

To complete the proof, it remains to show the validity of equations (2.34)-(2.35). We now recall the stationary distribution  $(^{single}) = ( (^{single})(\mathbf{u}); \mathbf{u} \geq \mathbf{U})$  of a single server loss system with state-dependent Poisson arrival process with rate  $\lambda_n(0 \leq n \leq C)$  when there are  $n$  progressing jobs and the service time distributions are as in the system model of Section 2.1. Then from [22], the stationary probability that the server has  $n$  progressing jobs and the  $i^{\text{th}}$  job has age at most  $y_i$  ( $1 \leq i \leq n$ ) is given by

$$(^{single})(\mathbf{y}_n) = \frac{\left\{ \prod_{i=1}^n \frac{\lambda_{i-1}}{i} \right\}}{1 + \sum_{m=1}^C \left\{ \prod_{i=1}^m \frac{\lambda_{i-1}}{i} \right\}} \prod_{i=1}^n \int_{x_i=0}^{y_i} \bar{G}(x_i) dx_i \quad (2.36)$$

and

$$(^{single})(0) = \frac{1}{1 + \sum_{m=1}^C \left\{ \prod_{i=1}^m \frac{\lambda_{i-1}}{i} \right\}}. \quad (2.37)$$

For the given fixed-point  $\mathbf{\Gamma}$  of the mean-field and its corresponding occupancy probability measure  $\mathbf{\Gamma}$ , consider a single server system under the assumption of a Poisson arrival process with state-dependent rate  $\lambda_n^{(exp)}(\mathbf{\Gamma})$  ( $0 \leq n \leq C$ ) when there are  $n$  progressing jobs. Then the unique stationary distribution is given by (2.36)-(2.37) with  $\lambda_n$  replaced by  $\lambda_n^{(exp)}(\mathbf{\Gamma})$  for all  $0 \leq n \leq C$ . But, from (2.10), (2.16), and Lemma 2.4, since  $\bar{R}_n(\cdot) = \sum_{j=n}^C \Gamma_j$ , we have that  $(^{single})$  is also another stationary distribution for the single server system with state dependent Poisson arrival process having rates  $\lambda_n^{(exp)}(\mathbf{\Gamma})$  for all  $0 \leq n \leq C$ . Since the stationary distribution must be unique, equations (2.34)-(2.35) must hold. This completes the proof.

## 2.7 Numerical Results

In this section, we provide some numerical results to support that the mean-field is globally asymptotically stable (GAS) under the assumption of mixed-Erlang JLDs.

Showing that the fixed-point of the mean-field approximates the stationary distribution of the system with large  $N$ , remains an open problem. If one can establish that the equilibrium or fixed-point of the MFEs is GAS, then the conclusion of the interchange of limits would follow from the Prohorov's theorem (Theorem A.2) [57]. Proving that the fixed-point of the MFEs is GAS is a challenging problem because the joint distribution of the occupancy and ages does not possess any monotonicity properties unlike the case of exponential service time distributions [10]. In this section, we present numerical results on the validity of the GAS of the mean-field for the case in which the service time distributions are mixed-Erlang. In this case, the state of a server is also multi-dimensional and the mean-field is also non-monotonic unlike the exponential case. It is numerically easier to solve the MFEs for the case of mixed-Erlang distributions as they are systems of ODEs unlike the case of general service time distributions for which the MFEs are PDEs as we have shown. One more reason for using mixed-Erlang distributions is that such distributions are dense in the set of all distributions that have support on  $\mathbb{R}_+$ , see [76]. Our numerical results show that the mean-field is GAS for the case of mixed-Erlang service time distributions.

We consider the system parameters as follows: The capacity of a server is assumed to be  $C = 5$ . The average job length is assumed to be equal to one, i.e.  $\bar{r} = 1$ . The service times have a Mixed-Erlang distribution given by sums of independent exponentially distributed random variables (known as an Erlang distribution) where the number of exponential phases is equal to  $i \in \{1, 2, \dots, M\}$  with probability  $p_i$  such that  $\sum_{i=1}^M p_i = 1$ . Each exponential phase is assumed to have rate  $\rho$ . Therefore, we have,

$$\frac{1}{\rho} = \frac{\sum_{i=1}^M i p_i}{\rho}.$$

We choose  $M = 3$ ,  $p_1 = 0.3$ ,  $p_2 = 0.3$ , and  $p_3 = 0.4$ .

Under the mixed-Erlang service time distribution assumptions, let  $S$  be the set of all possible server states defined as  $S = \bigcup_{n=0}^C S_n$  where  $S_0 = \{f(0)g\}$  and  $S_n = \{f(n; r_1, \dots, r_n) : 1 \leq r_i \leq M; 1 \leq i \leq n\}$ . We refer to an element in the set  $S$  by  $\mathbf{r}$  and an element in the set  $S_n$  by  $\mathbf{r}_n$ . The system dynamics can be modeled as a Markov process  $\mathbf{Y}^{(N)}(t) = (\mathbf{Y}_{\mathbf{r}}^{(N)}(t); \mathbf{r} \in S)$  where  $\mathbf{Y}_{\mathbf{r}}^{(N)}(t)$  denotes the fraction of servers with  $n$  jobs such that the  $i^{\text{th}}$  job has  $r_i$  remaining phases at time  $t$ . Since the Markov process  $(\mathbf{Y}^{(N)}(t); t \geq 0)$  is defined on a finite-dimensional space, we can establish the mean-field limit by using the same procedure as in the case of the exponential service times in [10]. Hence, we recall the following result without proof from [33].

**Proposition 2.2.** *If  $f\mathbf{Y}^{(N)}(0)g_{N-1}$  converges in distribution to a state  $\mathbf{u}$  as  $N \rightarrow \infty$ , then the sequence of processes  $f(\mathbf{Y}^{(N)}(t); t \geq 0)g_{N-1}$  converges in distribution to a deterministic process  $\mathbf{x}(t; \mathbf{u})$  as  $N \rightarrow \infty$  called the mean-field. The process  $\mathbf{x}(t; \mathbf{u})$  is the unique solution to the following system of differential equations*

$$\mathbf{x}(0; \mathbf{u}) = \mathbf{u}; \quad (2.38)$$

$$\dot{\mathbf{x}}_{\mathbf{r}_n}(t; \mathbf{u}) = h_{\mathbf{r}_n}(\mathbf{x}(t; \mathbf{u})); \quad (2.39)$$

and  $\mathbf{h}(\cdot) = (h_{\mathbf{r}}(\cdot); \mathbf{r} \in \mathcal{S})$  with the mapping  $h_{\mathbf{r}_n}(\cdot)$  given by

$$\begin{aligned} h_{\mathbf{r}_n}(\mathbf{y}) = & \sum_{b=1}^n \left\{ \frac{\rho_{r_b}}{n} \right\} y_{(\mathbf{r}_n^b)} \frac{(ME)}{n-1}(\mathbf{y}) - y_{\mathbf{r}_n} \frac{(ME)}{n}(\mathbf{y}) I_{r_n < Cg} \\ & + \sum_{b=1}^{n+1} \rho I_{r_n < Cg} y_{(\mathbf{r}_n^{b+1})} + \sum_{b=1}^n \rho y_{(n; r_1; \dots; r_{b-1}; r_{b+1}; \dots; r_n)} - n \rho y_{\mathbf{r}_n}; \end{aligned} \quad (2.40)$$

where

$$\frac{(ME)}{n}(\mathbf{v}) = \left\{ \frac{1}{\sum_{\mathbf{r}_n} v_{\mathbf{r}_n}} \right\} \left[ \left( \sum_{i=n}^C \sum_{\mathbf{b}_i} v_{\mathbf{b}_i} \right)^d - \left( \sum_{i=n+1}^C \sum_{\mathbf{b}_i} v_{\mathbf{b}_i} \right)^d \right]; \quad (2.41)$$

In Figure 2.1, we plot  $d_E^2(\mathbf{x}(t; \mathbf{u}); \cdot)$  as a function of  $t$  where  $d_E$  is the euclidean distance defined by

$$d_E(\mathbf{u}; \mathbf{v}) = \sqrt{\sum_{\mathbf{1} \leq \mathbf{2} \leq \mathbf{S}} j u_{\mathbf{1}} - v_{\mathbf{j}}^2};$$

It is observed that for  $d = 2$ ,  $\rho = 1$ , and for four different initial points  $\mathbf{u}_1; \mathbf{u}_2; \mathbf{u}_3$ , and  $\mathbf{u}_4$ , the mean-field  $\mathbf{x}(t; \mathbf{u})$  for the mixed-Erlang service time distribution converges to its unique fixed-point  $\cdot$ . Note that the computed  $\cdot$  depends on the chosen value of  $d$ . This supports that  $\cdot$  is globally stable.

We conclude with some numerical results for the blocking probability of the above system showing closeness to the theoretical lower bound. Similar to the case of exponential distributions as in [10], under the asymptotic independence, any finite set of servers are independent, and the fixed-point of the mean-field implies that the fixed point is the stationary distribution of the state of a server. The average blocking probability is then given by  $Q_C^d$  where  $Q_C = \sum_{\mathbf{v}_C \in \mathcal{S}_C} (\mathbf{v}_C)$ . Let us recall the lower bound on the average blocking probability denoted by  $P_{block}^{avg}$  for any load balancing scheme shown in [19]. From

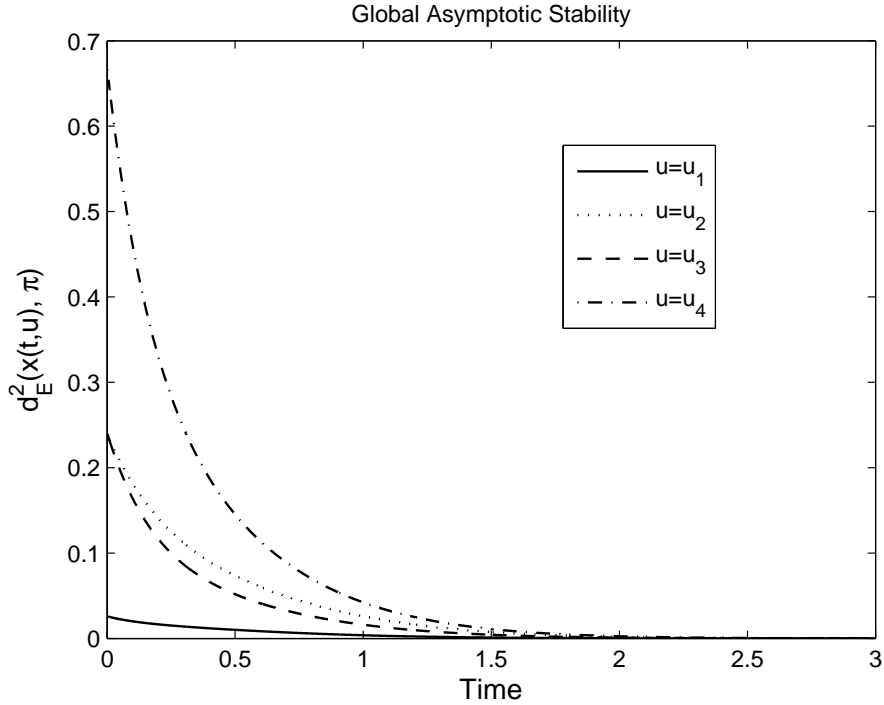


Figure 2.1: Convergence of the mean-field to the fixed-point

Little's law [76, Theorem 4.1], the average number of jobs in the system is equal to  $(1 - P_{block}^{avg})N$  which is upper bounded by  $NC$ . Hence,

$$P_{block}^{avg} \leq \left(1 - \frac{C}{N}\right)_+;$$

where  $(b)_+ = \max(b, 0)$ . In Figure 2.2, we plot the lower bound  $\left(1 - \frac{C}{N}\right)_+$  and the average blocking probability as a function of  $N$  when the dispatcher uses the  $SQ(d)$  load balancing and the state-independent random routing where a destination server is chosen uniformly at random. It is clear that the resulting average blocking probability under the  $SQ(d)$  policy is much lower than the resulting average blocking probability when pure random routing is employed. Furthermore, the average blocking probability under the  $SQ(d)$  load balancing approaches the lower bound as  $d$  increases.

## 2.8 Proofs of Main Results

In this section, we provide proofs of the results stated in Section 2.4.



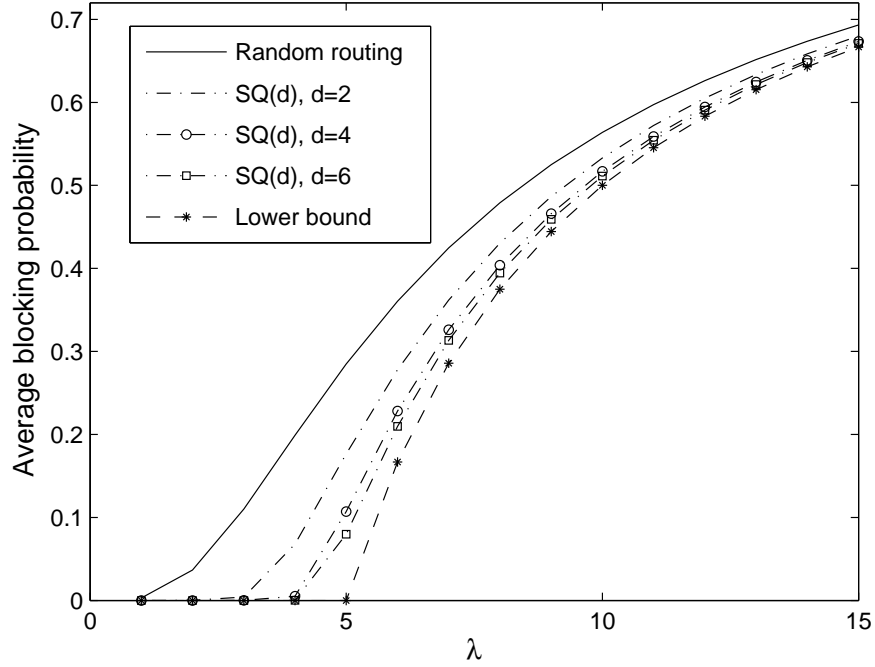


Figure 2.2: Comparison of the average blocking probability under  $SQ(d)$  with lower bound.

### 2.8.1 Proof Theorem 2.2

In this section, we first compute the semigroup of the Markov process  $(\mathbf{t}^{(N)}; t \geq 0)$  and then we show that the Markov process  $(\mathbf{t}^{(N)}; t \geq 0)$  is a Feller process. Finally, we construct a martingale using the generator of the Markov process  $(\mathbf{t}^{(N)}; t \geq 0)$  and the Dynkin's formula.

Conditioned on the initial state  $\mathbf{t}_0^{(N)}$ , let  $\mathbf{A}_h$  and  $\mathbf{D}_h$  be the number of arrivals and departures in the interval  $[0; h]$ , respectively. Note that a job with age  $x$  at time  $t$  departs from the system in the interval  $[t; t + h]$  with the probability  $\frac{G(x+h) - G(x)}{G(x)}$ . Further, from the definition of the hazard rate, we have that  $\lim_{h \downarrow 0} \frac{1}{h} \frac{G(x+h) - G(x)}{G(x)} = \mu(x)$  and hence,

$$\frac{G(x+h) - G(x)}{G(x)} = \mu(x)h + o(h); \quad (2.42)$$

Let  $T_h^{(N)}(\cdot)$  be the operator defined as

$$T_h^{(N)} f(\cdot) = \mathbb{E} \left[ f(\mathbf{t}_h^{(N)}) | \mathbf{t}_0^{(N)} = \cdot \right];$$

where  $f$  is a continuous bounded function  $f : M_F(\mathbb{U}) \rightarrow \mathbb{R}$  and the operator  $T_h^{(N)}(\cdot)$  is a semigroup operator when  $(\mathbf{t}^{(N)}; t \geq 0)$  is a Markov process. Before computing the

expression for  $T_h^{(N)} f(\cdot)$ , we first introduce the following notation. Suppose the measure  $\mu_0^{(N)}$  has mass at  $m$  atoms and let the  $i^{\text{th}}$  atom be  $\mathbf{v}^{(i)} = (n_i; v_1^{(i)}; \dots; v_{n_i}^{(i)})$  for  $1 \leq i \leq m$  and let the number of servers with the state  $\mathbf{v}^{(i)}$  be denoted by  $(f_{\mathbf{v}^{(i)}} g) = a^{(i)}$ . If a server lies in the state  $\mathbf{b}_n = (n; b_1; \dots; b_n)$  at time  $t$ , let the probability that there is no departure in the interval  $[t; t+h]$  be denoted by  $\rho_{ND}(\mathbf{b}_n; h)$ . We then have

$$\rho_{ND}(\mathbf{b}_n; h) = \prod_{i=1}^n \left\{ \frac{\bar{G}(b_i + h)}{G(b_i)} \right\}. \quad (2.43)$$

Note that using (2.42), we can write

$$\rho_{ND}(\mathbf{b}_n; h) = \prod_{j=1}^n (1 - (b_j)h) + o(h). \quad (2.44)$$

**Lemma 2.5.** *Let  $f$  be a real valued continuous bounded function defined on  $M_F(U)$ . Then the process  $(\mu_t^{(N)}; t \geq 0)$  is a Feller weak-homogeneous  $M_F(U)$ -valued Markov process with semigroup operator  $T_h^{(N)}(\cdot)$  given by*

$$\begin{aligned} T_h^{(N)} f(\cdot) &= (1 - N h) \left( \prod_{j=1; n_j > 0}^m (\rho_{ND}(\mathbf{v}^{(j)}; h))^{a^{(j)}} \right) f(\cdot) \\ &\quad + (1 - N h) \sum_{j=1; n_j > 0}^m \sum_{r=1}^{n_j} a^{(j)} \left\{ \frac{G(v_r^{(j)} + h) - G(v_r^{(j)})}{\bar{G}(v_r^{(j)})} \right\} \\ &\quad \left( \prod_{w=1; w \neq r}^{n_j} \left\{ \frac{\bar{G}(v_w^{(j)} + h)}{\bar{G}(v_w^{(j)})} \right\} \right) (\rho_{ND}(\mathbf{v}^{(j)}; h))^{(a^{(j)} - 1)} \left( \prod_{i=1; n_i > 0; i \neq j}^m (\rho_{ND}(\mathbf{v}^{(i)}; h))^{a^{(i)}} \right) \\ &\quad f\left(\cdot + \left( \sum_{i=1}^m \sum_{j=1}^{n_i} \mathbf{v}^{(j)} \right) \right) \\ &\quad + (N h) \int_{x=0}^h \frac{1}{h} \left( \sum_{i=1}^m \sum_{j=1}^{n_i+1} \frac{1}{n_i+1} \rho_r(x; \mathbf{v}^{(i)}) \right. \\ &\quad \left. \left[ I_{f_{n_i} < C_g} f\left(\cdot + \left( \sum_{i=1}^m \sum_{j=1}^{n_i} \mathbf{v}^{(j)} \right) \right) \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \bar{G}(h-x) \right. \right. \\ &\quad \left. \left. + I_{f_{n_i} = C_g} f\left(\cdot\right) \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \right] \right) dx + o(h); \quad (2.45) \end{aligned}$$

where  $o(h)$  is  $o(h)$  for all  $\cdot$ . Moreover, the process  $(\mu_t^{(N)}; t \geq 0)$  is a Feller-Dynkin process.

*Proof.* We can write

$$\begin{aligned}
T_h^{(N)} f(\cdot) &= \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=0; \mathbf{D}_h=0} j \binom{N}{0} = \right] + \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=0; \mathbf{D}_h=1} j \binom{N}{0} = \right] \\
&+ \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=1; \mathbf{D}_h=0} j \binom{N}{0} = \right] + \sum_{i=1}^m \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=i; \mathbf{D}_h=j} j \binom{N}{0} = \right]. \quad (2.46)
\end{aligned}$$

The proof mainly depends on the simplification of the four terms on the right side of (2.46).

We first simplify the first term on the right side of equation (2.46). In this case, since there are no arrivals or departures, we have  $\binom{N}{0} = 1$ . As a consequence, we have

$$\mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=0; \mathbf{D}_h=0} j \binom{N}{0} = \right] = f\left(\frac{\cdot}{h}\right) \mathbb{P}\left((f\mathbf{A}_h=0; \mathbf{D}_h=0) j \binom{N}{0} = \right): \quad (2.47)$$

Further, we can write

$$\mathbb{P}\left((f\mathbf{A}_h=0; \mathbf{D}_h=0) j \binom{N}{0} = \right) = \mathbb{P}\left(f\mathbf{A}_h=0 j \binom{N}{0} = \right) \mathbb{P}\left((f\mathbf{D}_h=0) j \mathbf{A}_h=0; \binom{N}{0} = \right):$$

Since the arrival process is a Poisson process with rate  $N$ , it is independent of the state  $\cdot$ . Therefore, we have  $\mathbb{P}\left(f\mathbf{A}_h=0 j \binom{N}{0} = \right) = \mathbb{P}(f\mathbf{A}_h=0) = e^{-N h}$ : On the other hand, the number of departures  $\mathbf{D}_h$  is influenced by the number of arrivals  $\mathbf{A}_h$ . Hence, we need to compute the expression for  $\mathbb{P}\left((f\mathbf{D}_h=j) j \mathbf{A}_h=i; \binom{N}{0} = \right)$  that gives the probability that there are  $j$  departures in the interval  $[0; h]$  conditioned on the event that there are  $i$  arrivals in the interval  $[0; h]$ . If the arrival process is a Poisson process, conditioned on the number of arrivals  $\mathbf{A}_h$ , the arrival instants are independent random variables with uniform distribution in the interval  $[0; h]$  [77, p. 325]. It can be seen that

$$\mathbb{P}\left((f\mathbf{D}_h=j) j \mathbf{A}_h=i; \binom{N}{0} = \right) = \prod_{j=1}^m (p_{ND}(\mathbf{v}^{(j)}; h))^{a^{(j)}}:$$

We can write

$$\begin{aligned}
\mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=0; \mathbf{D}_h=0} j \binom{N}{0} = \right] &= f\left(\frac{\cdot}{h}\right) e^{-N h} g \left( \prod_{j=1}^m (p_{ND}(\mathbf{v}^{(j)}; h))^{a^{(j)}} \right) f\left(\frac{\cdot}{h}\right) \\
&+ o(h);
\end{aligned}$$

where

$$o_1(\cdot; h) = f\mathbb{P}(f\mathbf{A}_h=0) \left(1 - N h\right) g \prod_{j=1}^m (p_{ND}(\mathbf{v}^{(j)}; h))^{a^{(j)}} f\left(\frac{\cdot}{h}\right)$$

is  $o(h)$  for all  $\cdot$ .

Similarly, we can write the second term of the right side of (2.46) as

$$\begin{aligned} \mathbb{E} \left[ f \binom{N}{h} I_{f\mathbf{A}_h=0; \mathbf{D}_h=1} g \binom{N}{0} = \right] &= (1 - N h) \sum_{j=1; n_j > 0}^m \sum_{r=1}^{n_j} a^{(j)} \left\{ \frac{G(v_r^{(j)} + h)}{\bar{G}(v_r^{(j)})} \right\} \\ &\left( \prod_{w=1; w \neq r}^{n_j} \left\{ \frac{\bar{G}(v_w^{(j)} + h)}{\bar{G}(v_w^{(j)})} \right\} \right) (\rho_{ND}(\mathbf{v}^{(j)}; h))^{(a^{(j)} - 1)} \left( \prod_{i=1; n_i > 0; i \neq j}^m (\rho_{ND}(\mathbf{v}^{(i)}; h))^{a^{(i)}} \right) \\ &f \left( h + \left( \binom{+}{h} \mathbf{v}^{(j)} \right)_r \quad \left( \binom{+}{h} \mathbf{v}^{(i)} \right) + 2(\cdot; h) \right); \end{aligned}$$

where we use  $r$  to denote the index of the departing job at a server with the state  $\mathbf{v}^{(j)}$  and  $2(\cdot; h)$  is  $o(h)$  for all  $\cdot$  given by

$$\begin{aligned} 2(\cdot; h) &= \mathbb{P}(f\mathbf{A}_h = 0) (1 - N h) g \sum_{j=1; n_j > 0}^m \sum_{r=1}^{n_j} a^{(j)} \\ &\left\{ \frac{G(v_r^{(j)} + h)}{\bar{G}(v_r^{(j)})} \right\} \left( \prod_{w=1; w \neq r}^{n_j} \left\{ \frac{\bar{G}(v_w^{(j)} + h)}{\bar{G}(v_w^{(j)})} \right\} \right) (\rho_{ND}(\mathbf{v}^{(j)}; h))^{(a^{(j)} - 1)} \\ &\left( \prod_{i=1; n_i > 0; i \neq j}^m (\rho_{ND}(\mathbf{v}^{(i)}; h))^{a^{(i)}} \right) f \left( h + \left( \binom{+}{h} \mathbf{v}^{(j)} \right)_r \quad \left( \binom{+}{h} \mathbf{v}^{(i)} \right) \right); \end{aligned}$$

We next compute the third term on the right side of (2.46). We can write

$$\begin{aligned} \mathbb{E} \left[ f \binom{N}{h} I_{f\mathbf{A}_h=1; \mathbf{D}_h=0} g \binom{N}{0} = \right] &= \mathbb{P}(f\mathbf{A}_h = 1) g \\ &\int_{x=0}^h \frac{1}{h} \left( \sum_{i=1}^m \sum_{j=1}^{n_i+1} \left\{ \frac{1}{n_i+1} \right\} \rho_r \left( x : \binom{+}{x} \mathbf{v}^{(i)} \right) \left[ I_{f n_i < C} g f \left( h + \left( \binom{+}{h} \mathbf{v}^{(i)} \right)_j; h - x \right) \quad \left( \binom{+}{h} \mathbf{v}^{(i)} \right) \right] \right. \\ &\left. \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \bar{G}(h - x) + I_{f n_i = C} g f \left( h \right) \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \right] \right) dx; \end{aligned}$$

where the arrival instant  $x$  is chosen uniformly in  $[0; h]$  given  $\mathbf{A}_h = 1$ ,  $i$  denotes the index of the atom corresponding to the state of the destination server, and  $j$  is the position of the routed job at the destination server chosen uniformly at random from  $n_i + 1$  positions.

Further, we write

$$\begin{aligned} \mathbb{E} \left[ f \binom{N}{h} I_{f\mathbf{A}_h=1; \mathbf{D}_h=0} g \binom{N}{0} = \right] &= (N h) \int_{x=0}^h \frac{1}{h} \left( \sum_{i=1}^m \sum_{j=1}^{n_i+1} \left\{ \frac{1}{n_i+1} \right\} \rho_r \left( x : \binom{+}{x} \mathbf{v}^{(i)} \right) \right. \\ &\left[ I_{f n_i < C} g f \left( h + \left( \binom{+}{h} \mathbf{v}^{(i)} \right)_j; h - x \right) \quad \left( \binom{+}{h} \mathbf{v}^{(i)} \right) \right] \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \bar{G}(h - x) \end{aligned}$$

$$+ I_{f n_i = C g} f(h) \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \Big] dx + \mathfrak{z}_3(\cdot; h);$$

where

$$\begin{aligned} \mathfrak{z}_3(\cdot; h) = & \mathbb{P}(f \mathbf{A}_h = 1g) \ N \ h g \int_{x=0}^h \frac{1}{h} \left( \sum_{i=1}^m \sum_{j=1}^{n_i+1} \left\{ \frac{1}{n_i+1} \right\} \rho_r(x : \mathbf{v}^{(i)}) \right. \\ & \left. \left[ I_{f n_i < C g} f \left( h + \left( \frac{x}{h} \mathbf{v}^{(i)} \right)^j ; h \right) \left( \frac{x}{h} \mathbf{v}^{(i)} \right) \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \overline{G}(h-x) \right. \right. \\ & \left. \left. + I_{f n_i = C g} f(h) \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \right] \right) dx: \end{aligned}$$

We now show that  $\mathfrak{z}_3(\cdot; h)$  is a  $o(h)$  term for all  $\cdot$ . For this, we apply the method of change of variables by replacing  $x$  with  $hy$ . As a consequence, we have

$$\begin{aligned} \mathfrak{z}_3(\cdot; h) = & \mathbb{P}(f \mathbf{A}_h = 1g) \ N \ h g h \int_{y=0}^1 \frac{1}{h} \left( \sum_{i=1}^m \sum_{j=1}^{n_i+1} \left\{ \frac{1}{n_i+1} \right\} \rho_r(hy : \mathbf{v}^{(i)}) \right. \\ & \left. \left[ I_{f n_i < C g} f \left( h + \left( \frac{hy}{h} \mathbf{v}^{(i)} \right)^j ; h \right) \left( \frac{hy}{h} \mathbf{v}^{(i)} \right) \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \overline{G}(h-hy) \right. \right. \\ & \left. \left. + I_{f n_i = C g} f(h) \prod_{k=1; n_k > 0}^m (\rho_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \right] \right) dy: \end{aligned}$$

By using the dominated convergence theorem [69, Theorem 1.34], we have  $\lim_{h \downarrow 0} \frac{\mathfrak{z}_3(\cdot; h)}{h} = 0$  for all  $\cdot$ .

Finally, by using the fact that  $f$  is a bounded function, we now prove that the fourth term on the right side of (2.46) is a  $o(h)$  term denoted by  $\mathfrak{z}_4(\cdot; h)$ . Since  $f \in C_b(M_F^N(U))$ , it is enough to prove that  $\sum_{i=1; j=1}^{n_i} \mathbb{P}(f \mathbf{A}_h = i; \mathbf{D}_h = jg \mid \mathcal{G}_0^{(N)} = \cdot)$  is a  $o(h)$  term. Due to the result that  $f \sum_{i=2; j=1}^{n_i} \mathbb{P}(f \mathbf{A}_h = i; \mathbf{D}_h = jg \mid \mathcal{G}_0^{(N)} = \cdot)g = \mathbb{P}(f \mathbf{A}_h = 2g)$  and  $\mathbb{P}(f \mathbf{A}_h = 2g)$  is a  $o(h)$  term, we conclude that  $\sum_{i=2; j=1}^{n_i} \mathbb{P}(f \mathbf{A}_h = i; \mathbf{D}_h = jg \mid \mathcal{G}_0^{(N)} = \cdot)$  is a  $o(h)$  term for all  $\cdot$ . We now show that  $f \sum_{j=1}^{n_1} \mathbb{P}(f \mathbf{A}_h = 1; \mathbf{D}_h = jg \mid \mathcal{G}_0^{(N)} = \cdot)g$  is a  $o(h)$  term. We can write

$$\begin{aligned} \left( \mathbb{P}(f \mathbf{A}_h = 1; \mathbf{D}_h = 1g \mid \mathcal{G}_0^{(N)} = \cdot) \right) &= \mathbb{P}(f \mathbf{A}_h = 1g \mid \mathcal{G}_0^{(N)} = \cdot) - \mathbb{P}(f \mathbf{A}_h = 1; \mathbf{D}_h = 0g \mid \mathcal{G}_0^{(N)} = \cdot) \\ &= \mathbb{P}(f \mathbf{A}_h = 1g) (1 - \mathbb{P}(f \mathbf{D}_h = 0g \mid \mathbf{A}_h = 1; \mathcal{G}_0^{(N)} = \cdot)): \end{aligned}$$

Again, by using the method of change of variables and the dominated convergence theorem as in the proof of the result that shows  $\mathfrak{z}_3(\cdot; h)$  is a  $o(h)$  term, we get that  $\lim_{h \downarrow 0} \mathbb{P}(f \mathbf{D}_h =$

$0g_j \mathbf{A}_h = 1; \binom{(N)}{0} = 1$  for all  $j$ . Since  $\lim_{h \downarrow 0} \frac{P(f \mathbf{A}_h = 1g)}{h} = N$ , for all  $f$ , we have that  $P(f \mathbf{A}_h = 1; \mathbf{D}_h = 1g_j \binom{(N)}{0} = 1)$  is  $o(h)$ . Therefore,  $Q_4(\cdot; h)$  is  $o(h)$  for all  $\cdot$ .

By combining the expressions for all the four terms on the right side of (2.46), and by defining  $Q(\cdot; h) = Q_1(\cdot; h) + Q_2(\cdot; h) + Q_3(\cdot; h) + Q_4(\cdot; h)$ ; we get the expression for  $T_h^{(N)} f(\cdot)$  as in (2.45). Finally, from [16, p.18],  $(\binom{(N)}{t}; t \geq 0)$  is a weak homogeneous Markov process.

We next give the proof of the Feller-Dynkin property of  $(\binom{(N)}{t}; t \geq 0)$ . From Lemma 3.5:1 and Corollary 3.5:2 of [16], we need to prove the following conditions to show the Feller-Dynkin property of  $(\binom{(N)}{t}; t \geq 0)$ . For  $f \in C_s^1(U)$ ,  $\varphi \in M_F(U)$ , let  $Q_f : M_F(U) \rightarrow \mathbb{R}$  be defined as  $Q_f(\varphi) = e^{-h} \int f(\varphi)$ , then we must have

1. For all  $f \in C_s^1(U)$  and  $h > 0$ , the mapping  $\varphi \mapsto \mathbb{E} \left[ Q_f(\binom{(N)}{h} j \binom{(N)}{0} = \varphi) \right]$  is continuous.

2. For all  $h > 0$ , we have

$$\mathbb{E} \left[ Q_1(\binom{(N)}{h} j \binom{(N)}{0} = \varphi) \right] \geq 0 \quad (2.48)$$

as  $\varphi \in M_F(U)$ .

3. For all  $\varphi \in M_F(U)$  and  $f \in C_s^1(U)$ , we have

$$\mathbb{E} \left[ Q_f(\binom{(N)}{h} j \binom{(N)}{0} = \varphi) \right] \geq Q_f(\varphi) \quad (2.49)$$

as  $h \downarrow 0$ .

From (2.45), we obtain

$$\begin{aligned} \mathbb{E} \left[ Q_f(\binom{(N)}{h} j \binom{(N)}{0} = \varphi) \right] &= e^{-h} \int f(\varphi) \left( (1 - N h) \left( \prod_{j=1; n_j > 0}^m (\rho_{ND}(\mathbf{v}^{(j)}; h))^{a^{(j)}} \right) \right. \\ &\quad \left. + (1 - N h) \sum_{j=1; n_j > 0}^m \sum_{r=1}^{n_j} a^{(j)} \left\{ \frac{G(v_r^{(j)} + h)}{\bar{G}(v_r^{(j)})} \right\} \right. \\ &\quad \left. \left( \prod_{w=1; w \neq r}^{n_j} \left\{ \frac{\bar{G}(v_w^{(j)} + h)}{\bar{G}(v_w^{(j)})} \right\} \right) (\rho_{ND}(\mathbf{v}^{(j)}; h))^{(a^{(j)} - 1)} \left( \prod_{i=1; n_i > 0; i \neq j}^m (\rho_{ND}(\mathbf{v}^{(i)}; h))^{a^{(i)}} \right) \right. \\ &\quad \left. + (N h) \int_{x=0}^h \frac{1}{h} \left( \sum_{i=1}^m \sum_{j=1}^{n_i+1} \left\{ \frac{1}{n_i + 1} \right\} p_r(x; \mathbf{v}^{(i)}) \right) \right. \\ &\quad \left. Q_f(\binom{(N)}{h} j \binom{(N)}{0} = \varphi) \right) \end{aligned}$$

$$\left[ I_{f_{n_i} < Cg} Q_f \left( ((\mathbf{v}^{(l)})^j; h, x) \mid (\mathbf{v}^{(l)}) \right) \prod_{k=1; n_k > 0}^m (p_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \bar{G}(h, x) + I_{f_{n_i} = Cg} \prod_{k=1; n_k > 0}^m (p_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \right] dx + f(\cdot; h) \quad (2.50)$$

where

$$f(\cdot; h) = {}_1f(\cdot; h) + {}_2f(\cdot; h) + {}_3f(\cdot; h) + {}_4f(\cdot; h)$$

such that

$${}_1f(\cdot; h) = (P(f\mathbf{A}_h = 0g) \mid (1 \ N \ h)) \left( \prod_{j=1; n_j > 0}^m (p_{ND}(\mathbf{v}^{(j)}; h))^{a^{(j)}} \right);$$

$${}_2f(\cdot; h) = (P(f\mathbf{A}_h = 0g) \mid (1 \ N \ h)) \sum_{j=1; n_j > 0}^m \sum_{r=1}^{n_j} a^{(j)} \left\{ \frac{G(v_r^{(j)} + h)}{\bar{G}(v_r^{(j)})} \right\} \left( \prod_{w=1; w \notin r}^{n_j} \left\{ \frac{\bar{G}(v_w^{(j)} + h)}{\bar{G}(v_w^{(j)})} \right\} \right) (p_{ND}(\mathbf{v}^{(j)}; h))^{(a^{(j)} - 1)} \left( \prod_{i=1; n_i > 0; i \notin j}^m (p_{ND}(\mathbf{v}^{(i)}; h))^{a^{(i)}} \right) Q_f \left( ((\mathbf{v}^{(j)})^r) \mid (\mathbf{v}^{(j)}) \right);$$

$${}_3f(\cdot; h) = (P(f\mathbf{A}_h = 1g) \mid N \ h) \int_{x=0}^h \frac{1}{h} \left( \sum_{i=1}^m \sum_{j=1}^{n_i+1} \left\{ \frac{1}{n_i + 1} \right\} p_r(x \mid \mathbf{v}^{(i)}) \right) \left[ I_{f_{n_i} < Cg} Q_f \left( ((\mathbf{v}^{(l)})^j; h, x) \mid (\mathbf{v}^{(l)}) \right) \prod_{k=1; n_k > 0}^m (p_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \bar{G}(h, x) + I_{f_{n_i} = Cg} \prod_{k=1; n_k > 0}^m (p_{ND}(\mathbf{v}^{(k)}; h))^{a^{(k)}} \right] dx;$$

and

$${}_4f(\cdot; h) = \sum_{i=1; j=1}^i \mathbb{E} \left[ Q_f \left( \sum_{r=1}^i (S((M_r; Z_{1r}; \dots; Z_{M_r r}); L_r; h; T_r) \mid (M_r; Z_{1r+h} \ T_r; \dots; Z_{M_r r+h} \ T_r)) + \sum_{l=1}^j (Y((n_l; X_{1l}; \dots; X_{n_l l}); K_l; h; W_l) \mid (n_l; X_{1l+h} \ W_l; \dots; X_{n_l l+h} \ W_l)) \right) I_{f\mathbf{A}_h=i; \mathbf{D}_h=jj} \binom{N}{0} = \right]; \quad (2.51)$$

where in equation (2.51),  $T_r$  is the arrival time of the  $r^{\text{th}}$  arrival and this job is routed to a server with state  $(M_r; Z_{1r}; \dots; Z_{M_r r})$  at time  $T_r$  and  $L_r$  is its position at the destination server. Further,  $S((M_r; Z_{1r}; \dots; Z_{M_r r}); L_r; h; T_r)$  denotes the possible state of the  $r^{\text{th}}$  arrival's destination at time  $t = h$ . Similarly, for the  $l^{\text{th}}$  departure,  $W_l$  is the departure time and this job departs from a server with state  $(n_l; X_{1l}; \dots; X_{n_l l})$  at time  $W_l$  and its position is  $K_l$ . Further,  $Y((n_l; X_{1l}; \dots; X_{n_l l}); K_l; h; W_l)$  is the possible state of the server where the  $l^{\text{th}}$  departure occurs at time  $t = h$ . It can be checked that  $Q_f(\cdot; h)$  is also  $o(h)$  and this follows from the same arguments as that of  $V(\cdot; h)$  in (2.45).

We now prove the first condition required to establish the Feller property. For this, we write (2.50) as

$$\mathbb{E} \left[ Q_f \left( \frac{(\cdot)^{(N)}}{h} \right) j \frac{(\cdot)^{(N)}}{0} = \cdot \right] = (e^{-h} \sum_{i=1}^N f_i) V(\cdot; h):$$

The mapping  $e^{-h} \sum_{i=1}^N f_i$  is a continuous mapping of  $\mathcal{U}$ . To prove the first condition, we need to show  $V(\cdot; h)$  is a continuous mapping of  $\mathcal{U}$ . Since  $\mathcal{U}$  is a point measure at finite  $N$ , the continuity of  $V(\cdot; h)$  w.r.t:  $t$  follows from the continuity of the routing probabilities and the departure probabilities. The second condition follows due to the fact that  $\sum_{i=1}^N (U) = \sum_{i=1}^N (U) = N$ . Finally, the third condition follows from the relation  $h \sum_{i=1}^N f_i = h \sum_{i=1}^N f_i$  and the dominated convergence theorem.

□

We next compute the infinitesimal generator of the Markov process  $(\frac{(\cdot)^{(N)}}{t}; t \geq 0)$ , and then we construct a martingale  $(\mathbf{M}_t^{(N)}(\cdot); t \geq 0) \in D_{\mathbb{R}}([0; \infty))$  where  $\mathcal{U} \in C_b^1(\mathcal{U})$  by applying the Dynkin's formula.

Since the set of linear combinations of  $Q_f : M_F(\mathcal{U}) \rightarrow \mathbb{R}$  for  $f \in C_b^1(\mathcal{U})$  defined by  $Q_f(\cdot) = e^{-h} \sum_{i=1}^N f_i$  is dense in the set  $C(M_F(\mathcal{U}))$  [78, proposition 7.10], by using  $A^{(N)} Q_f(\cdot)$ , for any continuous function  $F \in C(M_F(\mathcal{U}))$  such that the infinitesimal generator  $A^{(N)} F(\cdot) = \lim_{h \downarrow 0} \frac{\mathbb{E}[F(\frac{(\cdot)^{(N)}}{h}) j \frac{(\cdot)^{(N)}}{0} = \cdot] - F(\cdot)}{h}$  is well-defined for all  $\cdot$ , we have

$$\begin{aligned} A^N F(\cdot) &= \lim_{h \downarrow 0} \frac{F(\frac{(\cdot)}{h}) - F(\cdot)}{h} = N F(\cdot) - F(\cdot) \sum_{n=1}^C \sum_{j=1}^n \int \int_{\mathcal{U}_n} (x_j) d(\mathbf{x}_n) \\ &\quad + \sum_{n=1}^C \sum_{j=1}^n \int \int_{\mathcal{U}_n} (x_j) \left( F(\cdot +_{(\mathbf{x}_n^j)} \cdot_{(\mathbf{x}_n)}) \right) d(\mathbf{x}_n) \\ &+ N \left[ \left( \left\{ -\frac{(f_0 g)}{N} \right\} \Phi_0 \left( \frac{\cdot}{N} \right) \left( F(\cdot +_{(1,0)} \cdot_{(0)}) \right) \right) + \sum_{n=1}^C \sum_{j=1}^{n+1} \int \int_{\mathcal{U}_n} \left\{ \frac{1}{N(n+1)} \right\} \Phi_n \left( \frac{\cdot}{N} \right) \right] \end{aligned}$$



$$F(\cdot + (\mathbf{x}_n^j; 0) - (\mathbf{x}_n)) d(\mathbf{x}_n) + \int_{\cup_C} \int_{\cup_C} \frac{1}{N} \Phi_C \left( \frac{\cdot}{N} \right) F(\cdot) d(\mathbf{x}_C) \Big]: \quad (2.52)$$

For  $\varphi \in C_b^1(\mathbb{U})$ , it can be seen that the function  $\varphi \in M_F(\mathbb{U}) \cap \mathcal{H}^1; \varphi \in \mathbb{R}$  belongs to the domain of  $A^{(N)}(\cdot)$

**Proposition 2.3.** For all  $\varphi \in C_b^1(\mathbb{U})$ , the process  $(\mathbf{M}_t^{(N)}(\varphi); t \geq 0)$  given by

$$\mathbf{M}_t^{(N)}(\varphi) = h_t^{(N)}(\varphi); \varphi \in \mathcal{H}_0^{(N)}; \varphi \in \int_{s=0}^t A^{(N)} h_s^{(N)}(\varphi) ds \quad (2.53)$$

is an RCLL square integrable  $F_t^N$  martingale. For  $\varphi \in C_b^1(\mathbb{U})$ , the quadratic variation of  $(\mathbf{M}_t^{(N)}(\varphi); t \geq 0)$  is given by

$$\begin{aligned} \langle \mathbf{M}_t^{(N)}(\varphi) \rangle_{t=0} &= \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int_{\cup_n} \int_{\cup_n} (\varphi_j) (\varphi(\mathbf{x}_n^j) - \varphi(\mathbf{x}_n))^2 d_s^{(N)}(\mathbf{x}_n) \right. \\ &\quad \left. + N \left[ \left( \left\{ \frac{(\varphi)_0}{N} \right\} \Phi_0 \left( \frac{\cdot}{N} \right) (\varphi(1;0) - \varphi(0))^2 \right) \right. \right. \\ &\quad \left. \left. + \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{\cup_n} \left\{ \frac{1}{N(n+1)} \right\} \Phi_n \left( \frac{\cdot}{N} \right) (\varphi(\mathbf{x}_n^j; 0) - \varphi(\mathbf{x}_n))^2 d_s^{(N)}(\mathbf{x}_n) \right] \right) ds: \quad (2.54) \end{aligned}$$

*Proof.* From the Dynkin's formula [56], the process  $(\mathbf{M}_t^{(N)}(\varphi); t \geq 0)$  defined by

$$\mathbf{M}_t^{(N)}(\varphi) = h_t^{(N)}(\varphi); \varphi \in \mathcal{H}_0^{(N)}; \varphi \in \int_{s=0}^t A^{(N)} h_s^{(N)}(\varphi) ds \quad (2.55)$$

is an RCLL  $F_t^N$  local martingale. Therefore, by simplification we get

$$\begin{aligned} \mathbf{M}_t^{(N)}(\varphi) &= h_t^{(N)}(\varphi); \varphi \in \mathcal{H}_0^{(N)}; \varphi \in \int_{s=0}^t h_s^{(N)}(\varphi); r_1 \varphi ds \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int_{\cup_n} \int_{\cup_n} (\varphi_j) \right. \\ &\quad \left. (\varphi(\mathbf{x}_n^j) - \varphi(\mathbf{x}_n)) d_s^{(N)}(\mathbf{x}_n) \right. \\ &\quad \left. + N \left[ \left( \left\{ \frac{(\varphi)_0}{N} \right\} \Phi_0 \left( \frac{\cdot}{N} \right) (\varphi(1;0) - \varphi(0))^2 \right) \right. \right. \\ &\quad \left. \left. + \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{\cup_n} \left\{ \frac{1}{N(n+1)} \right\} \Phi_n \left( \frac{\cdot}{N} \right) (\varphi(\mathbf{x}_n^j; 0) - \varphi(\mathbf{x}_n))^2 d_s^{(N)}(\mathbf{x}_n) \right] \right) ds: \quad (2.56) \end{aligned}$$

By choosing  $F(\varphi) = h_t^{(N)}(\varphi); \varphi \in \mathcal{H}_0^{(N)}$ , from [79, Theorem 7.15], we have

$$\langle \mathbf{M}_t^{(N)}(\varphi) \rangle_{t=0} = \int_{s=0}^t (A^{(N)} F^2(\varphi) - 2F(\varphi) A^{(N)} F(\varphi)) ds: \quad (2.57)$$

After simplifications, we get (2.54). Finally, since  $\varphi \in C_b^1(\mathbb{U})$  and  $\varphi \in C_b(\mathbb{R}_+)$ , we have  $\mathbb{E}[\langle \mathbf{M}_t^{(N)}(\varphi) \rangle_t] < 1$  and hence,  $(\mathbf{M}_t^{(N)}(\varphi); t \geq 0)$  is a square integrable martingale.  $\square$

## 2.8.2 Proof of Lemma 2.2

We first show that any process  $(h_t; i; t \geq 0)$  that satisfies (2.10) also satisfies (2.11). By using the fundamental theorem of calculus [69, p.144], for  $\varphi \in C_b^1(U)$ , a real valued process  $(h_t; i; t \geq 0)$  satisfying equation (2.10) is a solution to the following differential equation (2.58) if the integrand in equation (2.10) is a continuous function of  $S$ ,

$$\begin{aligned} \frac{dh_t; i}{dt} = & h_t; r_1; i + \left( \sum_{n=1}^C \sum_{j=1}^n \int \int_{U_n} (\varphi_j) (\mathbf{x}_n^j) (\mathbf{x}_n) d_t(\mathbf{x}_n) \right. \\ & + \left[ \varphi_0(t) (\mathbf{1}; 0) (\mathbf{0}) + \sum_{n=1}^{C-1} \sum_{j=1}^{n+1} \int \int_{U_n} \left\{ \frac{1}{n+1} \right\} \right. \\ & \left. \left. \Phi_n(t) (\mathbf{x}_n^j; 0) (\mathbf{x}_n) d_t(\mathbf{x}_n) \right] \right); \end{aligned} \quad (2.58)$$

Therefore, we need to show that the two terms on the right side of (2.58) are continuous functions of  $t$ . Since  $\varphi \in C_b^1(U)$  and the mapping  $t \mapsto S_t$  is continuous, the first term  $h_t; r_1; i$  is a continuous function of  $t$ . In the second term, the expression related to the case of departures can be written as

$$\sum_{n=1}^C \sum_{j=1}^n \int \int_{U_n} (\varphi_j) (\mathbf{x}_n^j) (\mathbf{x}_n) d_t(\mathbf{x}_n) = h_t; \varphi_1; i;$$

where the function  $\varphi_1$  is defined as

$$\varphi_1(\mathbf{x}_n) = \begin{cases} 0 & \text{if } n = 0; \\ \sum_{j=1}^n (\varphi_j) (\mathbf{x}_n^j) (\mathbf{x}_n) & \text{otherwise;} \end{cases}$$

Since  $\varphi \in C_b^1(U)$  and  $\varphi \in C_b^1(\mathbb{R}_+)$ , we have that  $\varphi_1 \in C_b(U)$ . Therefore the mapping  $t \mapsto h_t; \varphi_1; i$  is continuous. The expression that corresponds to the case of arrivals can be written as

$$\begin{aligned} h_t; (\varphi); i = & (\varphi_0(t) (\mathbf{1}; 0) (\mathbf{0})) \\ & + \sum_{n=1}^{C-1} \sum_{j=1}^{n+1} \int \int_{U_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(t) (\mathbf{x}_n^j; 0) (\mathbf{x}_n) d_t(\mathbf{x}_n); \end{aligned}$$

where  $(\varphi)_t$  is defined as

$$(\varphi)_t(\mathbf{x}_n) = \begin{cases} 0 & \text{if } n = C; \\ \left\{ \frac{\varphi(\mathbf{x}_n; t)}{n+1} \right\} \sum_{j=1}^{n+1} (\varphi_j^j; 0) (\mathbf{x}_n) & \text{otherwise;} \end{cases} \quad (2.59)$$

For given  $t$ , since  $\mu \in C_b(U)$ , we have that  $\mu(\cdot) \in C_b(U)$ . Hence, for any constant  $a > 0$ , the mapping  $t \mapsto \int h_{t'; \mu(\cdot)} d\mu$  is continuous.

We next prove that the mapping  $t \mapsto \int h_{t'; \mu(\cdot)} d\mu$  is continuous, i.e., we need to prove that  $\int h_{t+b'; \mu(\cdot+b)} d\mu \rightarrow \int h_{t'; \mu(\cdot)} d\mu$  as  $b \rightarrow 0$ . We have

$$\left| \int h_{t+b'; \mu(\cdot+b)} d\mu - \int h_{t'; \mu(\cdot)} d\mu \right| = \left| \int h_{t+b'; \mu(\cdot+b)} d\mu - \int h_{t+b'; \mu(\cdot)} d\mu \right| + \left| \int h_{t+b'; \mu(\cdot)} d\mu - \int h_{t'; \mu(\cdot)} d\mu \right|: \quad (2.60)$$

Since  $\mu(\cdot) \in C_b(U)$ , we have that  $\lim_{b \rightarrow 0} \left| \int h_{t+b'; \mu(\cdot)} d\mu - \int h_{t'; \mu(\cdot)} d\mu \right| = 0$ : We next prove that  $\lim_{b \rightarrow 0} \left| \int h_{t+b'; \mu(\cdot+b)} d\mu - \int h_{t+b'; \mu(\cdot)} d\mu \right| = 0$ :

For  $L > 0$ , let

$$U^{(L)} = \{x_n \in U_n : n^{-1}x_i > L \text{ for all } 1 \leq i \leq n\}:$$

For given  $\epsilon > 0$ , since  $\mu$  is tight, we can find some  $L > 0$  such that  $\int_{U^{(L)}} h_{t'; \mu} d\mu < \epsilon$ : Furthermore, from the continuity of the mapping  $t \mapsto \int h_{t'; \mu} d\mu$ , we can find some  $h_1 > 0$  such that for all  $b \in [0, \min(t; h_1)]$ ,

$$\int_{U^{(L)}} h_{t+b'; \mu} d\mu < \epsilon: \quad (2.61)$$

By using the fact that the mapping  $t \mapsto \bar{R}_n(t) = \int h_{t'; \mu_{j=-n}^{j=n}} d\mu$  is continuous, we next show that the mapping  $t \mapsto \bar{R}_n(t)$  is continuous. For this, we need to show that  $\bar{R}_n(t+b) - \bar{R}_n(t) \rightarrow 0$  as  $b \rightarrow 0$ . From (2.59), we have

$$\begin{aligned} \bar{R}_n(t+b) - \bar{R}_n(t) &= 2 \int h_{t+b'; \mu} d\mu - \int h_{t+b'; \mu} d\mu - \int h_{t'; \mu} d\mu + \int h_{t'; \mu} d\mu \\ &= 4d \int h_{t+b'; \mu} d\mu - \int h_{t+b'; \mu} d\mu - \int h_{t'; \mu} d\mu: \end{aligned} \quad (2.62)$$

Since  $\left| \bar{R}_n(t+b) - \bar{R}_n(t) \right| \rightarrow 0$  as  $b \rightarrow 0$  for all  $n$ ,  $\bar{R}_n(t+b) - \bar{R}_n(t) \rightarrow 0$ . This proves that the mapping  $t \mapsto \bar{R}_n(t)$  is continuous. As a consequence, we have that  $\mu(\cdot+b)$  is uniformly continuous on the interval  $b \in [0, \min(t; h_1)]$  and  $\mathbf{u} \in \bar{U}^{(L)}$  (the complement of  $U^{(L)}$ ). As a result, there exists some  $h_2 \in (0; h_1)$  such that for  $b \in [0, \min(t; h_2)]$ ,  $\mathbf{u} \in \bar{U}^{(L)}$ , we have

$$\left| \int h_{t+b'; \mu} d\mu - \int h_{t'; \mu} d\mu \right| < \epsilon: \quad (2.63)$$

Using (2.61)-(2.63), for  $b \in [0, \min(t; h_2)]$ , we have

$$\left| \int h_{t+b'; \mu(\cdot+b)} d\mu - \int h_{t'; \mu(\cdot)} d\mu \right| \leq \int_{U^{(L)}} h_{t+b'; \mu} d\mu + 4d \int h_{t+b'; \mu} d\mu + 4d \int h_{t'; \mu} d\mu: \quad (2.64)$$

By letting  $b \rightarrow 0$  and then  $\epsilon \rightarrow 0$  in (2.60), we have the continuity of the mapping  $t \mapsto \int h_{t'; \mu(\cdot)} d\mu$ .

We next show that a solution to (2.58) is also a solution to an another differential equation obtained by applying a method of change of variables. For  $r = t$ , we have

$$\begin{aligned} \frac{dh_{r;t,r,i}}{dr} &= \lim_{h \downarrow 0} \frac{f_{h_{r+h;t,r,h,i}} - h_{r;t,r,i}g}{h} \\ &= \lim_{h \downarrow 0} \frac{f_{h_{r+h;t,r,h,i}} - h_{r+h;t,r,i}g}{h} + \lim_{h \downarrow 0} \frac{f_{h_{r+h;t,r,i}} - h_{r;t,r,i}g}{h} \end{aligned} \quad (2.65)$$

We now obtain the expression for the first term on the right side of (2.65). We can write

$$h_{r+h;t,r,h,i} - h_{r+h;t,r,i} = h_{r+h;i} \hat{W}_i;$$

where  $\hat{W}$  is defined such that  $\hat{W}(\mathbf{y}_n) = h_{r+h;t,r,h}(\mathbf{y}_n) - h_{r+h;t,r}(\mathbf{y}_n)$ : We further simplify the function  $\hat{W}$  by using the following definition, let

$$\frac{\partial}{\partial S_j}(\mathbf{y}_n) = \lim_{h \downarrow 0} \frac{(y_n^j; y_i + h) - (y_n)}{h};$$

We can write

$$\hat{W}(\mathbf{y}_n) = (h_{r+h;t,r,h}(\mathbf{y}_n) - (h_{r+h;t,r,h}(\mathbf{y}_n)^1; y_1 + t - r) + (h_{r+h;t,r,h}(\mathbf{y}_n)^1; y_1 + t - r) - h_{r+h;t,r}(\mathbf{y}_n));$$

Further, we have

$$(h_{r+h;t,r,h}(\mathbf{y}_n) - (h_{r+h;t,r,h}(\mathbf{y}_n)^1; y_1 + t - r)) = \int_{y_1+t-r-h}^{y_1+t-r} \frac{\partial}{\partial S_1}((h_{r+h;t,r,h}(\mathbf{y}_n)^1; S_1) ds_1;$$

By replacing  $S_1$  with  $y_1 + t - r - hv$ , we get

$$(h_{r+h;t,r,h}(\mathbf{y}_n) - (h_{r+h;t,r,h}(\mathbf{y}_n)^1; y_1 + t - r)) = h \int_{v=0}^1 \frac{\partial}{\partial S_1}((h_{r+h;t,r,h}(\mathbf{y}_n)^1; y_1 + t - r - hv) dv;$$

Similarly, we can write

$$\begin{aligned} &(n; y_1 + t - r; \dots; y_{i-1} + t - r; y_i + t - r - h; y_{i+1} + t - r - h; \dots; y_n + t - r - h) \\ &\quad - (n; y_1 + t - r; \dots; y_i + t - r; y_{i+1} + t - r - h; \dots; y_n + t - r - h) \\ &= h \int_{v=0}^1 \frac{\partial}{\partial S_i}(n; y_1 + t - r; \dots; y_{i-1} + t - r; y_i + t - r - hv; y_{i+1} + t - r - h; \dots; y_n + t - r - h) dv; \end{aligned}$$

For  $1 \leq i \leq n$ , let

$$W_{(i;t,r,h;v)}(\mathbf{y}_n) = \frac{\partial}{\partial S_i}(n; y_1 + t - r; \dots; y_{i-1} + t - r; y_i + t - r - hv; y_{i+1} + t - r - h; \dots; y_n + t - r - h);$$

As a consequence, after simplifications we have

$$\hat{w}(\mathbf{y}_n) = h \int_{v=0}^1 \sum_{i=1}^n (W_{(i;t;r;h;v)}(\mathbf{y}_n)) dv:$$

Let the function  $W_{(t;r;h;v)} \in C_b(U)$  be defined as

$$W_{(t;r;h;v)}(\mathbf{y}_n) = \begin{cases} 0 & \text{if } n = 0; \\ \sum_{i=1}^n (W_{(i;t;r;h;v)}(\mathbf{y}_n)) & \text{otherwise:} \end{cases}$$

Now we can see that

$$\lim_{h \downarrow 0} \frac{f h_{r+h; t r h} i h_{r+h; t r} i g}{h} = \lim_{h \downarrow 0} \int_{v=0}^1 h_{r+h; t r} W_{(t;r;h;v)} i dv:$$

Since  $h \nabla h_{r+h; t r} W_{(t;r;h;v)} i$  is continuous, by the dominated convergence theorem, we have

$$\lim_{h \downarrow 0} \frac{f h_{r+h; t r h} i h_{r+h; t r} i g}{h} = h_{r; r_1} \tilde{i}: \quad (2.66)$$

We now focus on the second term on the right side of (2.65). We have

$$h_{r+h; t r} i h_{r; t r} i = \int_{u=r}^{r+h} \frac{\partial}{\partial U} h_{u; t r} i du:$$

By using (2.58), we have

$$\begin{aligned} h_{r+h; t r} i h_{r; t r} i &= \int_{u=r}^{r+h} \left( h_{u; r_1} t r i \right. \\ &+ \sum_{n=1}^C \sum_{j=1}^n \int \int_{U_n} (X_j) (t r (\mathbf{x}_n^j) t r (\mathbf{x}_n)) d_u(\mathbf{x}_n) \\ &+ \left[ u(f_0 g) \Phi_0(u) (t r (1; 0) t r (0)) + \sum_{n=1}^C \sum_{j=1}^{n+1} \int \int_{U_n} \left\{ \frac{1}{n+1} \right\} \right. \\ &\quad \left. \Phi_n(u) (t r (\mathbf{x}_n^j; 0) t r (\mathbf{x}_n)) d_u(\mathbf{x}_n) \right] \Big) du: \end{aligned}$$

Again, by applying the method of change of variables, we have

$$\begin{aligned} h_{r+h; t r} i h_{r; t r} i &= h \int_{v=0}^1 h_{r+hv; r_1} t r i \\ &+ \left( \sum_{n=1}^C \sum_{j=1}^n \int \int_{U_n} (X_j) (t r (\mathbf{x}_n^j) t r (\mathbf{x}_n)) d_{r+hv}(\mathbf{x}_n) \right) \end{aligned}$$

$$+ \left[ \int_{r+h\nu}^r (f_0 g) \Phi_0(r+h\nu) (t, r; 1; 0) - \int_{r+h\nu}^r (f_0 g) \Phi_0(r+h\nu) (t, r; 0) + \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{U_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(r+h\nu) (t, r; \mathbf{x}_n^j; 0) - \int_{U_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(r+h\nu) (t, r; \mathbf{x}_n) d_{r+h\nu}(\mathbf{x}_n) \right] dv:$$

As a result, by using the dominated convergence theorem we have

$$\begin{aligned} \lim_{h \downarrow 0} \frac{f_h(r+h\nu; t, r; i) - f_h(r; t, r; i)}{h} &= f_{r; r; 1; t, r; i} \\ &+ \left( \sum_{n=1}^C \sum_{j=1}^n \int_{U_n} (X_j) (t, r; \mathbf{x}_n^j) - (t, r; \mathbf{x}_n) d_r(\mathbf{x}_n) \right. \\ &+ \left. \left[ \int_r (f_0 g) \Phi_0(r) (t, r; 1; 0) - \int_r (f_0 g) \Phi_0(r) (t, r; 0) + \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{U_n} \left\{ \frac{1}{n+1} \right\} \right. \right. \\ &\quad \left. \left. \Phi_n(r) (t, r; \mathbf{x}_n^j; 0) - \int_{U_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(r) (t, r; \mathbf{x}_n) d_r(\mathbf{x}_n) \right] \right): \quad (2.67) \end{aligned}$$

Finally, by using (2.66) and (2.67), we have

$$\begin{aligned} \frac{d f_{r; t, r; i}}{dr} &= \sum_{n=1}^C \sum_{j=1}^n \int_{U_n} (X_j) (t, r; \mathbf{x}_n^j) - (t, r; \mathbf{x}_n) d_r(\mathbf{x}_n) \\ &+ \left[ \int_r (f_0 g) \Phi_0(r) (t, r; 1; 0) - \int_r (f_0 g) \Phi_0(r) (t, r; 0) + \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{U_n} \left\{ \frac{1}{n+1} \right\} \right. \\ &\quad \left. \Phi_n(r) (t, r; \mathbf{x}_n^j; 0) - \int_{U_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(r) (t, r; \mathbf{x}_n) d_r(\mathbf{x}_n) \right]: \end{aligned}$$

By integrating  $\frac{d f_{r; t, r; i}}{dr}$  with respect to  $r$  from 0 to  $t$ , we get (2.11) for  $f \in C_b^1(U)$ . Then the result can be extended to the simple functions by using the monotone convergence theorem and then to the class of functions  $C_b(U)$  from the standard arguments by using the Dynkin - theorem [56, page 497] or from the fact that  $C_b^1(U)$  is dense in  $C_b(U)$ .

We next prove that for  $f \in C_b^1(U)$ , the solution  $(f_{t; i; t}^-)$  to (2.11) is also a solution to (2.10). For this, it is enough to prove the differentiability of  $f_{t; i}^-$  with respect to  $t$ . Since  $f \in C_b^1(U)$ , the existence of  $\frac{d f_{0; t; i}^-}{dt}$  follows from the dominated convergence theorem. By using the Leibniz integral rule, we now verify the existence of the derivative of the second term on the right side of (2.11) with respect to  $t$ . According to this rule, the first condition is that the integrand needs to be continuous with respect to both the variables  $r$  and  $t$ . This follows from the same arguments used in the proof of the continuity of the

integrand in equation (2.10). The second condition is that the derivative of the integrand with respect to  $t$  must exist and the differential should be continuous with respect to both  $r$  and  $t$ . The derivative of the integrand with respect to  $t$  exists from the dominated convergence theorem since  $\mathcal{L} \in C_b^1(U)$  and also, it is continuous with respect to  $r$  and  $t$  from the same arguments that we have used to prove the continuity of the integrand in (2.10). Therefore, any process  $(\mu_t; t \in [0, 1]) \in C_{M_1(U)}([0, 1])$  is a solution to (2.10) if and only if it is a solution to (2.11). Further, note that  $\mu_t$  need not be a differentiable function in (2.11).

### 2.8.3 Proof of Theorem 2.1

From equation (2.11), we first make it clear that for all  $\mathcal{L} \in C_b(U)$ , the operator  $\mathcal{V} h_t; i$  is a linear operator with  $\mathcal{V} 1(U) = 1$ . Hence from the Riesz-Markov-Kakutani theorem (Theorem A.4) [69, Theorem 2.14], for  $\mu_t \in M_1(U)$ , the existence of the unique operator  $\mathcal{V} h_t; i$  implies the existence of the unique probability measure  $\mu_t$ . The uniqueness of  $\mu_t$  also follows from the fact that  $C_b(U)$  is a separating class of  $M_1(U)$  [56, p.111]. If  $\mu_1; \mu_2 \in M_1(U)$  satisfies  $\mathcal{V} h_t; i = \mathcal{V} h_t; i$  and  $\mathcal{V} h_t; i = \mathcal{V} h_t; i$  for all  $\mathcal{L} \in C_b(U)$ , then we have  $\mu_1 = \mu_2$ .

Given an initial measure  $\mu_0$ , we next prove that there exists at most one mean-field solution by showing that there exists at most one real valued process  $h_t; i$  corresponding to the mean-field. Suppose  $(\mu_t^1; t \in [0, 1]); (\mu_t^2; t \in [0, 1])$  are two solutions to the MFEs with initial points  $\mu_0^1; \mu_0^2$ , respectively. For  $\mathcal{L} \in C_b(U)$ , we then have

$$\begin{aligned}
h_t^1; i - h_t^2; i &= h_0^1; i - h_0^2; i + \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int_{U_n} \int_{U_n} (X_j) (\mu_s^1(\mathbf{x}_n^j) - \mu_s^2(\mathbf{x}_n)) \right. \\
&\quad \left. d(\mu_s^1 - \mu_s^2)(\mathbf{x}_n) \right) ds \\
&\quad + \int_{s=0}^t \left( \left[ \mu_s^1(\mathcal{L}) - \mu_s^2(\mathcal{L}) \right] \Phi_0(\mu_s^1) (\mu_s^1(1;0) - \mu_s^2(1;0)) \right. \\
&\quad \left. + \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{U_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(\mu_s^1) (\mu_s^1(\mathbf{x}_n^j; 0) - \mu_s^2(\mathbf{x}_n)) d\mu_s^1(\mathbf{x}_n) \right] \\
&\quad \left[ \mu_s^2(\mathcal{L}) - \mu_s^1(\mathcal{L}) \right] \Phi_0(\mu_s^2) (\mu_s^2(1;0) - \mu_s^1(1;0)) \\
&\quad \left. + \sum_{n=1}^C \sum_{j=1}^{n+1} \int_{U_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(\mu_s^2) (\mu_s^2(\mathbf{x}_n^j; 0) - \mu_s^1(\mathbf{x}_n)) d\mu_s^2(\mathbf{x}_n) \right] \right) ds: \quad (2.68)
\end{aligned}$$

The first term on the right side of (2.68) can be bounded as  $\int_0^t \sum_{k=1}^n (x_k) (t, s; \mathbf{x}_n^j) (t, s; \mathbf{x}_n) d(\frac{1}{s}, \frac{2}{s})(\mathbf{x}_n) ds$ . To simplify the second term corresponding to departures, we define a function  $h_{t,s}$  as follows:

$$h_{t,s}(\mathbf{x}_n) = \begin{cases} 0 & \text{if } n = 0; \\ \sum_{k=1}^n (x_k) (t, s; \mathbf{x}_n^j) (t, s; \mathbf{x}_n) & \text{otherwise:} \end{cases}$$

Then since  $\sum_{k=1}^n (x_k) \in C_b(\mathbb{U})$  and  $\sum_{k=1}^n (x_k) \in C_b(\mathbb{R}_+)$ , we have  $h_{t,s} \in C_b(\mathbb{U})$ . Further, we have  $\|h_{t,s}\| \leq C \|k\|$ . Using the definition of  $h_{t,s}$ , we can write

$$\begin{aligned} \int_{s=0}^t \left( \sum_{n=1}^C \sum_{j=1}^n \int_{\mathbb{U}_n} (x_j) (t, s; \mathbf{x}_n^j) (t, s; \mathbf{x}_n) d(\frac{1}{s}, \frac{2}{s})(\mathbf{x}_n) ds \right) \\ = \int_{s=0}^t h_{t,s}(\frac{1}{s}, \frac{2}{s}; h_{t,s}) ds \end{aligned}$$

To simplify the third term that corresponds to arrivals, we define a function  $f_{t,s}$  as follows: for  $0 \leq n \leq C-1$ ,

$$f_{t,s}(\mathbf{x}_n) = \begin{cases} 0 & \text{if } n = C; \\ \sum_{j=1}^{n+1} \left\{ \frac{1}{n+1} \right\} \Phi_n(\cdot) (t, s; \mathbf{x}_n^j; 0) (t, s; \mathbf{x}_n) & \text{otherwise:} \end{cases}$$

Then the third term is equal to  $\int_{s=0}^t (h_{t,s}(\frac{1}{s}, \frac{2}{s}; f_{t,s}(\frac{1}{s}, \frac{2}{s}; i)) - h_{t,s}(\frac{2}{s}, \frac{1}{s}; f_{t,s}(\frac{2}{s}, \frac{1}{s}; i))) ds$ . Further, we can write

$$\begin{aligned} |h_{t,s}(\frac{1}{s}, \frac{2}{s}; f_{t,s}(\frac{1}{s}, \frac{2}{s}; i)) - h_{t,s}(\frac{2}{s}, \frac{1}{s}; f_{t,s}(\frac{2}{s}, \frac{1}{s}; i))| \\ \leq |h_{t,s}(\frac{1}{s}, \frac{2}{s}; f_{t,s}(\frac{1}{s}, \frac{2}{s}; i))| + |h_{t,s}(\frac{2}{s}, \frac{1}{s}; f_{t,s}(\frac{2}{s}, \frac{1}{s}; i))| \\ \leq \|k_{t,s}(\frac{1}{s}, \frac{2}{s})\| \|k_{t,s}(\frac{1}{s}, \frac{2}{s})\| \|f_{t,s}(\frac{1}{s}, \frac{2}{s}; i)\| + \|k_{t,s}(\frac{2}{s}, \frac{1}{s})\| \|k_{t,s}(\frac{2}{s}, \frac{1}{s})\| \|f_{t,s}(\frac{2}{s}, \frac{1}{s}; i)\| \end{aligned}$$

Since  $\frac{2}{s}$  is a probability measure,  $\|k_{t,s}(\frac{2}{s})\| = 1$ . Furthermore,  $\|k_{t,s}(\frac{1}{s})\| \leq 2\|k\|$  and

$$|f_{t,s}(\frac{1}{s})(\mathbf{x}_n) - f_{t,s}(\frac{2}{s})(\mathbf{x}_n)| \leq 2\|k\| (|\bar{R}_n(\frac{1}{s}) - \bar{R}_n(\frac{2}{s})| + |\bar{R}_{n+1}(\frac{1}{s}) - \bar{R}_{n+1}(\frac{2}{s})|) :$$

We can write  $\bar{R}_n(\frac{1}{s}) = h_{t,s}(\frac{1}{s}, \frac{2}{s}; f)$  where  $f$  is a function defined as

$$f(\mathbf{x}_m) = \begin{cases} 1 & \text{if } m = n; \\ 0 & \text{otherwise:} \end{cases}$$

We then have  $|\bar{R}_n(\frac{1}{s}) - \bar{R}_n(\frac{2}{s})| \leq \|k_{t,s}(\frac{1}{s}, \frac{2}{s})\| \|k_{t,s}(\frac{1}{s}, \frac{2}{s})\| \|f\| = \|k_{t,s}(\frac{1}{s}, \frac{2}{s})\|^2 \|k\|$ :

Finally, by using bounds for all the terms, we get

$$|h_{t,s}(\frac{1}{s}, \frac{2}{s}; i) - h_{t,s}(\frac{2}{s}, \frac{1}{s}; i)| \leq \left( \|k_{t,s}(\frac{1}{s}, \frac{2}{s})\|^2 \|k\| + \int_{s=0}^t 2\|k_{t,s}(\frac{1}{s}, \frac{2}{s})\| \|k_{t,s}(\frac{1}{s}, \frac{2}{s})\| \|k\| ds + \int_{s=0}^t 8\|k_{t,s}(\frac{1}{s}, \frac{2}{s})\|^2 \|k\| ds \right) \|k\|$$



Therefore we have

$$k_t^1 - \frac{2}{t}k - k_0 - \frac{2}{0}k + (2Ck - k + 8d^2) \int_{s=0}^t k_s^1 - \frac{2}{s}k ds: \quad (2.69)$$

From the Gronwall's inequality (Theorem A.7) [56, Theorem 5.1, p.498], for some  $b, c > 0$ ,  $t \geq [0; T]$ , if  $k_t^1 - \frac{2}{t}k \leq b + c \int_{s=0}^t k_s^1 - \frac{2}{s}k ds$ ; then it follows that  $k_t^1 - \frac{2}{t}k \leq b e^{ct}$ . Therefore, from (2.69), we have  $k_t^1 - \frac{2}{t}k - k_0 - \frac{2}{0}k \leq e^{(2Ck - k + 8d^2)t}$ : Hence, starting from an initial measure  $\mu_0$ , there exists at most one solution for the MFEs.

We now prove that there exists a process  $(k_t; t \in [0; 1]) \in C_{M_1(U)}([0; 1])$  satisfying the mean-field model equations. This follows from the relative compactness of the sequence  $\{k_t^{(N)}; t \in [0; 1]\}$  in  $D_{M_1(U)}([0; 1])$  from the proof of Theorem 2.3. In particular, we have that every limit point of the sequence  $\{k_t^{(N)}; t \in [0; 1]\}$  satisfies (2.11). Further, each limiting point is almost surely continuous. This concludes that there exists a solution to the MFEs.

#### 2.8.4 Proof of Theorem 2.4

The first part of Theorem 2.4 is a special case of the second part. Hence, it is sufficient to prove the second part.

For  $t \geq 0$ ,  $\mu_t^{(N)}$  is the empirical probability measure on  $U$  such that  $\mu_t^{(N)}(\mathbf{u})$  for  $\mathbf{u} \in U$  denotes the fraction of servers lying in state  $\mathbf{u} \in U$  at time  $t$ . From the dynamics of the system under the SQ( $d$ ) scheme and the exchangeability of  $\{S_k(N)(0); 1 \leq k \leq N\}$ , the collection  $\{S_k^{(N)}(t); 1 \leq k \leq N\}$  is also exchangeable for all  $t \in [0; 1]$ . Further, from Theorem 2.2, we have  $\mu_t^{(N)} \rightarrow \mu_t$  for  $t \in [0; 1]$  as  $N \rightarrow \infty$ .

To prove the result, it is sufficient to show that the following holds:

$$\mathbb{E} \left[ \prod_{k=1}^I k(S_k^{(N)}(t)) \right] \rightarrow \prod_{k=1}^I h_{t; k}^i; \quad (2.70)$$

for all continuous bounded mappings  $k: U \rightarrow \mathbb{R}$  as  $N \rightarrow \infty$ .

We can write

$$\left| \mathbb{E} \left[ \prod_{k=1}^I k(S_k^{(N)}(t)) \right] - \prod_{k=1}^I h_{t; k}^i \right| = \left| \mathbb{E} \left[ \prod_{k=1}^I k(S_k^{(N)}(t)) \right] - \mathbb{E} \left[ \prod_{k=1}^I h_{t; k}^{(N)} \right] \right| + \left| \mathbb{E} \left[ \prod_{k=1}^I h_{t; k}^{(N)} \right] - \prod_{k=1}^I h_{t; k}^i \right|; \quad (2.71)$$

From Theorem 2.2, the second term on the right hand side of the above inequality vanishes as  $N \rightarrow \infty$ . Let  $Q(r; n)$  be the set of all permutations of the numbers  $1; 2; \dots; n$  taken  $r$  at a time. Now, due to exchangeability, the permutation of states between servers does not affect the joint distribution. Then  $|Q(r; n)| = \binom{n}{r} r!$  where  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ . Hence, we have

$$\mathbb{E} \left[ \prod_{k=1}^l \mathbf{s}_k^{(N)}(t) \right] = \frac{1}{\binom{n}{l} l!} \mathbb{E} \left[ \sum_{2Q(l; n)} \prod_{k=1}^l \mathbf{s}_{(k)}^{(N)}(t) \right].$$

Also, by the definition of  $\bar{h}_t^{(N)}$ , we have

$$\mathbb{E} \left[ \prod_{k=1}^l \bar{h}_t^{(N)}(k) \right] = \mathbb{E} \left[ \left( \prod_{k=1}^l \frac{1}{N} \sum_{j=1}^N \mathbf{s}_j^{(N)}(t) \right)^l \right]. \quad (2.72)$$

Also, let  $D(r; n)$  be the set of all  $r$ -tuples that are formed by elements chosen from  $1; 2; \dots; n$ . Then  $|D(r; n)| = n^r$ . Let  $\overline{Q}(r; n)$  be the set of elements that are present in  $D(r; n)$  but not in  $Q(r; n)$ . Then from (2.72), we can write

$$\begin{aligned} \mathbb{E} \left[ \prod_{k=1}^l \bar{h}_t^{(N)}(k) \right] &= \frac{1}{N^l} \mathbb{E} \left[ \sum_{2D(l; n)} \prod_{k=1}^l \mathbf{s}_{(k)}^{(N)}(t) \right] \\ &= \frac{1}{N^l} \mathbb{E} \left[ \sum_{2Q(l; n)} \prod_{k=1}^l \mathbf{s}_{(k)}^{(N)}(t) \right] + \frac{1}{N^l} \mathbb{E} \left[ \sum_{2\overline{Q}(l; n)} \prod_{k=1}^l \mathbf{s}_{(k)}^{(N)}(t) \right]. \end{aligned}$$

Let  $\max_k \mathbf{s}_k = B$ , then we have the bound

$$\left| \mathbb{E} \left[ \prod_{k=1}^l \mathbf{s}_k^{(N)}(t) \right] - \mathbb{E} \left[ \prod_{k=1}^l \bar{h}_t^{(N)}(k) \right] \right| \leq \left( \frac{1}{\binom{n}{l} l!} - \frac{1}{N^l} \right) |Q(l; n)| B^l + \frac{1}{N^l} |\overline{Q}(l; n)| B^l \leq 2B^l \left( 1 - \frac{\binom{n}{l} l!}{N^l} \right).$$

The result follows since  $\left( 1 - \frac{\binom{n}{l} l!}{N^l} \right) \rightarrow 0$  as  $N \rightarrow \infty$ . This completes the proof.

### 2.8.5 Proof of Lemma 2.4

The mean-field  $(\bar{h}_t; t \geq 0)$  is a solution to (2.11). Let us consider the function  $\hat{h}_t = \int_{\mathbb{R}^n} \bar{h}_t(\mathbf{y}) \delta_{\mathbf{y}}$ . For  $s \geq 0$ , let  $\bar{h}_s$  be an absolutely continuous measure which has no atoms, we have  $\bar{h}_s^+ \hat{h}_t = \bar{h}_s^- \hat{h}_t$  where  $\bar{h}_s^\pm = \int_{\mathbb{R}^n} \bar{h}_s(\mathbf{y}) \delta_{\mathbf{y}}$ . Since there exists a sequence of functions  $f_n g_n \in C_b(\mathbb{U})$  that increase point wise to  $f_B g_B$  where  $B$  is an open set in

$\cup_n, n \geq 1$ , by using the monotone convergence theorem [69, Theorem 1.26] and (2.11), we have that (2.11) is true even for the function  $\hat{\cdot}$ . Furthermore, since the measure  $\bar{\nu}_s$  is absolutely continuous for all  $s \geq 0$ , we have that equation (2.11) is true even for the function  $\hat{\cdot}$ . Therefore, using (2.11), we can obtain the evolution equations of the process  $(P_t; t \geq 0)$  where  $P_t = (P_t(\mathbf{u}); \mathbf{u} \geq \cup)$  and  $P_t(\mathbf{y}_n) = h_{t; \hat{i}}$ . We can further simplify the expression of the process  $(P_t(\mathbf{u}); \mathbf{u} \geq \cup; t \geq 0)$  obtained from (2.11) using the fact that

$$\begin{aligned} h_{s; b}^{f_{\mathbf{x}_n \geq \cup_n; 0} x_i y_i; \delta i g^i} &= h_{s; b}^{f_{\mathbf{x}_n \geq \cup_n; 0} x_i + b y_i; \delta i g^i} \\ &= h_{s; b}^{f_{\mathbf{x}_n \geq \cup_n; 0} x_i y_i; b; \delta i g^i} \end{aligned}$$

By differentiating  $P_t(\mathbf{y}_n)$  with respect to  $t$  and after simplifications, it is verified that the process  $(P_t; t \geq 0)$  satisfies equations (2.18)-(2.20).

### 2.8.6 Proof of Lemma 2.3

From Remark 2.5, we recall that the MFEs are the dynamics of the probability distribution of a single server loss system with capacity  $C$  in which jobs arrive according to a Poisson process with rate  $\Phi_n(\bar{\nu}_t)$  ( $n \geq 0$ ) when there are  $n$  progressing jobs. We have that the initial distribution  $\#$  has a density function and our objective is to show that for given  $t = r$ ,  $\bar{\nu}_r$  has a density function w.r.t. the Lebesgue measure. For  $n \geq 1$ ,  $\mathbf{u}_n = (n; u_1; \dots; u_n)$ , we now prove that  $\bar{\nu}_r$  has density at  $\mathbf{u}_n$ . For  $\epsilon_i > 0$  and  $1 \leq i \leq n$ , let

$$B_i = \{f(n; y_1; \dots; y_n) : u_i - y_i < u_i + \epsilon_i; 1 \leq i \leq n\} \quad (2.73)$$

At time  $t = r$ , the probability that there are  $n$  progressing jobs and the  $i^{\text{th}}$  job has age  $y_i$  such that  $y_i \geq [u_i; u_i + \epsilon_i)$ ,  $1 \leq i \leq n$ , is equal to  $\bar{\nu}_r(B)$ . Out of the  $n$  progressing jobs that are present at time  $t = r$ , let  $\mathcal{J}_1$  be the set of indices of all the progressing jobs that entered the system at a time  $t > 0$  and  $\mathcal{J}_2$  be the set of indices of all the progressing jobs which are present in the system from time  $t = 0$ . Precisely,

$$\mathcal{J}_1 = \{i : r > u_i; 1 \leq i \leq n\}$$

and

$$\mathcal{J}_2 = \{i : r \leq u_i; 1 \leq i \leq n\}$$

Essentially, if  $i \in \mathcal{J}_1$ , it implies that the age of the  $i^{\text{th}}$  job is less than  $r$ . As the ages of progressing jobs increase linearly with time at a unit rate, the  $i^{\text{th}}$  job must have entered

the system at a time  $t > 0$ . Precisely, if the  $i^{\text{th}}$  job's age  $y_i$  satisfies  $y_i \in [u_i; u_i + \tau_i)$  and  $i \in \mathcal{J}_1$  at time  $t = r$ , then it implies that the  $i^{\text{th}}$  job must have entered the system in the time interval  $(r - \tau_i - \tau_i; r - \tau_i]$  and its service is not finished by the time  $t = r$ . On the other hand, if  $j \in \mathcal{J}_2$ , then it means that the  $j^{\text{th}}$  job is present in the system from time  $t = 0$ . At time  $t = r$ , if the  $j^{\text{th}}$  job's age  $y_j$  satisfies  $y_j \in [u_j; u_j + \tau_j)$  and  $j \in \mathcal{J}_2$ , then its age should lie in the interval  $[u_j - r; u_j + \tau_j - r)$  at time  $t = 0$ .

Using the sets  $\mathcal{J}_1$  and  $\mathcal{J}_2$ , we now obtain an upper bound on  $\mathbb{P}_r^-(B)$  from which we conclude that there exists a density function. Let  $D_1 = \#\mathcal{J}_1$  and  $D_2 = \#\mathcal{J}_2$ .

We next obtain bounds on the probability that there exists  $D_2$  jobs at time  $t = r$  such that their ages lie in the set  $\mathcal{J}_2$ . Let  $B_1$  be the event that there exists  $D_2$  jobs at time  $t = r$  such that their ages lie in the set  $\mathcal{J}_2$ . Note that the total number of jobs say  $q$  that are present at time  $t = 0$  can be more than  $D_2$  but only  $D_2$  of them should not expire by the time  $t = r$ . Let  $l_i$  be the  $i^{\text{th}}$  smallest element of the set  $\mathcal{J}_2$ . Also, let  $i_j$  be the index of the job out of  $q$  jobs which will not expire by the time  $t = r$  and its age will lie in the interval  $[u_{l_j}; u_{l_j} + \tau_{l_j})$  at time  $t = r$ . For the event  $B_1$  to occur, the state of the server at time  $t = 0$  should belong to the set  $V$  where

$$V = \{f(q; x_1, \dots, x_q) : x_m \in \mathbb{R}_+ \text{ for } m \in \{1, \dots, i_{D_2}\} \text{ and } x_m \in [u_{l_a} - r; u_{l_a} - r + \tau_{l_a}) \text{ for } m = i_a + 1, \dots, q\}$$

A job with age  $x$  at time  $t = 0$  will stay in the system at time  $t = r$  with probability  $\frac{\overline{G}(x+r)}{\overline{G}(x)}$ . Then by using all the above arguments, we get the following bound where  $f_{\#} = (f_{\#}(\mathbf{u}); \mathbf{u} \in \mathbb{U})$  is the pdf of  $\#$ ,

$$\mathbb{P}(B_1) = \sum_{q=D_2}^{\infty} \left( \sum_{(i_1, \dots, i_{D_2}) \in \mathcal{I}_{1,2}; q} \int_V \int_{\mathbb{U}} f_{\#}(q; x_1, \dots, x_q) \left( \left\{ \prod_{m=1}^{D_2} \frac{\overline{G}(x_{i_m} + r)}{\overline{G}(x_{i_m})} \right\} \right) dx_1 \dots dx_q \right) : \quad (2.74)$$

We now focus on the jobs that belong to the set  $\mathcal{J}_1$ . Let  $B_2$  be the event that for each  $j \in \mathcal{J}_1$ , there is an arrival in the time interval  $(r - \tau_j - \tau_j; r - \tau_j]$  and furthermore, this job is not expired by the time  $t = r$ . Since the arrival process is a Poisson process with rate  $\Phi_n(-t)$  when there are  $n$  jobs and  $\Phi_n(-t) = d$  for all  $n \geq 0$ , for any time interval  $[t_1; t_2]$ , we have

$$\mathbb{P}(\mathbf{Z}_1) = \mathbb{P}(\mathbf{Z}_2);$$

where  $\mathbf{Z}_1$  denotes the number of arrivals to the server in the interval  $[t_1; t_2]$  and  $\mathbf{Z}_2$  denotes the number of arrivals in the interval  $[t_1; t_2]$  when the arrival process is a Poisson process

with rate  $d$ . Let  $k_i$  be the  $i^{\text{th}}$  smallest element of the set  $\mathcal{J}_1$ . Then since the arrival instants have uniform distribution conditioned on the number of arrivals over a time interval [77, page 325], we get

$$P(B_2) = \frac{(d)^{D_1}}{D_1!} \left( \prod_{j=1}^{D_1} \bar{G}(u_{k_j}) \right) \quad (2.75)$$

Finally, from (2.74) and (2.75), we have

$$\begin{aligned} \bar{\pi}_t(B) &= \left( \sum_{q:q=D_2}^C \left( \sum_{(i_1, \dots, i_{D_2}) \in \mathcal{I}_{D_2}} \int \int_V f_{\#}(n; x_1, \dots, x_n) \right. \right. \\ &\quad \left. \left. \left( \prod_{m=1}^{D_2} \left\{ \frac{\bar{G}(x_{i_m} + r)}{\bar{G}(x_{i_m})} \right\} \right) dx_1 \dots dx_q \right) \right) \left( \frac{(d)^{D_1}}{D_1!} \left( \prod_{j=1}^{D_1} \bar{G}(u_{k_j}) \right) \right) \quad (2.76) \end{aligned}$$

Clearly,  $\bar{\pi}_t$  has density at  $\mathbf{u}$  since  $\bar{\pi}_t(B) \neq 0$  as  $j \neq 0$  for  $1 \leq j \leq n$ .

## 2.9 Conclusions

In this chapter, we have provided a measure-valued process approach to establish the mean-field behavior of loss systems with  $\text{SQ}(d)$  load balancing and general service time requirements. The extension of these results to multi-class systems where servers are classified into different classes based on their capacities and jobs are classified into different classes based on their service requirements follows in a similar manner *mutatis mutandis* from the approach used here. Establishing the global asymptotic stability of the unique fixed-point remains an open problem.

## Chapter 3

# On Occupancy Based Randomized Load Balancing for Large Processor Sharing Systems

In this chapter, we show that the mean-field analysis of the SQ( $d$ ) randomized load balancing schemes can be extended to randomized schemes that use the occupancy information of a finite number of randomly sampled servers to dispatch an incoming arrival in a large-scale system of parallel servers. Such policies besides the SQ( $d$ ) policy also contain threshold based policies, and  $d$ -adaptive policies and generalizations of them. In this chapter, we present a mean-field analysis of occupancy based routing schemes for a system with a large number of processor sharing (PS) servers as an archetype of shared resource systems. The PS case is interesting not only because of its use in server-farm applications but also because PS scheduling is also known to be ordinarily insensitive. The analysis of PS schemes is harder than for loss models. As in the loss model case, the MFEs are PDEs. We show that the probability measure of occupancy defined on the set of non-negative integers  $Z_+$  obtained from a fixed-point of the mean-field also satisfies the stationary mean-field equations when the job lengths are exponentially distributed with the same average length. If the mean-field in the exponential case has a unique fixed-point, then the fixed point is *insensitive* to the job length distribution. The approach is also via a measure-valued Markov process approach. We also provide simulation results to justify our analysis.

**Organization of the Chapter** The rest of the chapter is organized as follows. In Section 3.1, we first introduce the system model, and then we present a generic framework

for occupancy based randomized routing policies. We give additional notation used in this chapter in Section 3.2. After that, we provide a mathematical formulation to the time evolution of the system in Section 3.3. We then present the main results of this chapter in Section 3.4, where we also illustrate how to use our analysis to study four occupancy based randomized load balancing policies. We study the SQ( $d$ ) policy and three threshold based policies that are adaptive to variations in the system load. We give proofs of the main results in Section 3.7. In Section 3.5, we provide numerical results to study the impact of the routing policies on the average response time, implementation complexity, and the insensitivity of the stationary distribution as the number of servers converges to  $\infty$ . In Section 3.6, we provide numerical results which support that the mean-field is globally asymptotically stable (GAS) for the SQ( $d$ ) policy and also, we give some insights into the behavior of the system in the stationary regime. Finally, we conclude in Section 3.8.

## 3.1 System Model and Routing Policy

In this section, we introduce the system model along with a general framework for a class of occupancy based randomized routing policies which use the information about the occupancy of a finite number of randomly selected servers upon an arrival to dispatch the arrival to a server.

### 3.1.1 System Model

Consider a large-scale system with  $N$  servers where each server has an infinite buffer and jobs are served according to the PS service discipline. Jobs arrive according to a Poisson process with rate  $\lambda$  to a central job dispatcher. Upon an arrival, the job dispatcher collects information about the occupancy of a finite number of randomly chosen servers to decide the destination server based on a predefined routing policy. We assume that the dispatcher has no buffer and an arrival is assumed to join a server's buffer immediately although the routing decision incurs some delay overhead in practice.

We assume that every server has capacity to process jobs at a unit rate and hence, from the assumption that the service discipline is PS, if a server has  $n$  jobs in its buffer, then every job is served at a rate of  $\frac{1}{n}$ . Furthermore, a job that is routed to a server with occupancy  $n$  joins at a position chosen uniformly at random from  $n + 1$  positions. We

assume that the job lengths are distributed according to a general distribution function  $G(\cdot)$  with density function  $g(\cdot)$  having the mean  $\frac{1}{\mu}$ . To ensure the stability of the system, we assume that  $\rho < 1$ . The hazard rate function of the job length distribution is denoted by  $\lambda(\cdot)$  defined in (2.1). Furthermore, we assume that  $\lambda(\cdot)$  satisfies Assumption 2.1.1.

### 3.1.2 Routing Policy

Before defining the routing policy precisely, we first introduce a general framework for occupancy dependent randomized routing policies, according to which the number of randomly selected servers also referred to as potential destination servers is not a fixed value as in the SQ( $d$ ) policy, but it is a dynamic value that takes into account the system's congestion level. In this framework, we assume that there are at most  $M$  stages to decide the destination server for an arrival. In each stage, the dispatcher first samples a finite number of servers and then, it decides whether to route the arrival to one of the sampled servers or not based on the occupancies of the sampled servers. In stage  $i$  of the routing policy, the dispatcher selects  $d_i$  servers uniformly at random from all the  $N$  servers, and then, it decides whether to select or not to select a server from the list of the potential destination servers as the destination server, based on the occupancy information of the potential destination servers that are sampled up to stage  $i$ . In stage  $i$ , if the occupancies of the potential destination servers satisfy certain criteria as described in the routing policy, then the dispatcher selects a potential destination server as the destination server. Otherwise, the routing policy enters stage  $i + 1$  in which the dispatcher samples  $d_{i+1}$  servers and the same procedure as in stage  $i$  is repeated until stage  $M$ . The value of  $d_{i+1}$  may depend on the occupancies of the potential destination servers that are sampled in stages 1 to  $i$  to take into account the congestion level of the system. We assume that  $d_1$  is a fixed-value but  $d_i$  for  $i \geq 2$  may depend on the occupancies of the servers that are sampled up to stage  $i - 1$ . The routing policy is stopped at stage  $M$  in which a potential destination server is selected as the destination server.

Let  $\mathbf{n}^{[i]} = (n_1^{[i]}; \dots; n_{d_i}^{[i]})$  be the vector of occupancies of servers that are randomly selected in stage  $i$ , where  $n_j^{[i]}$  is the occupancy of the  $j^{\text{th}}$  potential destination server selected in stage  $i$ . Let  $S^{[i]}$  be the set of occupancies of all the potential destination servers selected up to the  $i^{\text{th}}$  stage of the routing policy given by

$$S^{[i]} = f(\mathbf{n}^{[1]}; \dots; \mathbf{n}^{[i]}) : \text{occupancies of sampled servers in stage } j \text{ is } \mathbf{n}^{[j]}; 1 \leq j \leq i; \quad (3.1)$$



We assume that the dispatcher can sample at most  $d^{(max)}$  servers to decide the destination server for an arrival. Therefore, we have  $\sum_{i=1}^M d_i \leq d^{(max)}$ . Without loss of generality, from now onwards, for a vector of real numbers, we refer to the number of elements in the given vector as its size. Let  $S^{[i;j]}$  be the set of all occupancy vectors in  $S^{[i]}$  with size  $j$ . Then for  $1 \leq i \leq M$ ,

$$S^{[i]} = \bigcup_{j=1}^{d^{(max)}} S^{[i;j]}. \quad (3.2)$$

In the  $i^{\text{th}}$  stage, based on the occupancies of the potential destination servers, the dispatcher may decide to route the arrival to one of the potential destination servers if the occupancies meet certain criteria as stated in the predefined routing policy. We indicate whether the dispatcher selects the destination server in stage  $i$  or not by a vector of non-negative real numbers. If the number of sampled potential destination servers up to stage  $i$  is  $k$ , then the outcome of stage  $i$  is represented by a vector of size  $k+1$ . If  $(n_1; \dots; n_k)$  is the occupancy vector of potential destination servers where  $k$  is the total number of potential destination servers selected up to stage  $i$ , then the outcome at the end of stage  $i$  according to the routing policy is given by  $(m; \rho_1; \dots; \rho_k)$  where  $m$  is a non-negative integer value and  $\rho_i$  is a non-negative real number for  $1 \leq i \leq k$ . If the dispatcher decides to route the arrival to one of the sampled potential destination servers in stage  $i$  as per the routing policy, then  $m = 0$  and  $(\rho_1; \dots; \rho_k)$  indicates the routing probability vector satisfying  $\sum_{j=1}^k \rho_j = 1$  where  $\rho_i$  is the probability with which the  $i^{\text{th}}$  potential destination server is chosen as the destination server for the arrival. In this case, we say that the outcome of stage  $i$  is a routing probability vector. On the other hand, if the dispatcher decides to enter stage  $i+1$  according to the routing policy, then  $m = d_{i+1} - 1$  which indicates the number of servers to be sampled in stage  $i+1$  and  $(\rho_1; \dots; \rho_k)$  is the null-vector satisfying  $\sum_{j=1}^k \rho_j = 0$ . In this case, we say that the outcome of stage  $i$  is the null-vector.

Let  $\Upsilon^{[i]}$  be the set of all occupancy vectors in the set  $S^{[i]}$  for which the outcome of  $i^{\text{th}}$  stage is a routing probability vector. Further, let  $\Upsilon^{[i;j]}$  be the set of elements of size  $j$  in  $\Upsilon^{[i]}$ . Then for  $1 \leq i \leq M$ ,

$$\Upsilon^{[i]} = \bigcup_{j=1}^{d^{(max)}} \Upsilon^{[i;j]}. \quad (3.3)$$

Let  $\overline{\Upsilon^{[i]}}$  be the complement of the set  $\Upsilon^{[i]}$  in  $S^{[i]}$ . In the  $i^{\text{th}}$  stage, the routing policy is a mapping  $\Lambda^{(i)}$  that maps  $\Upsilon^{[i;j]}$  as follows:

$$\Lambda^{(i)} : \Upsilon^{[i;j]} \rightarrow (0; M_1(f_1; 2; \dots; jg)); \quad (3.4)$$

where  $M_1(f_1; 2; \dots; jg)$  denotes the set of probability measures on the set  $f_1; 2; \dots; jg$ .

For  $1 \leq k \leq d^{(max)}$ , if  $(r_1; \dots; r_k) \in \overline{Y^{[k]}}$ , then

$$\Lambda^{(k)}((r_1; \dots; r_k)) = (d_{i+1}; 0; \dots; 0); \quad (3.5)$$

where  $d_{i+1}$  is the number of servers to be sampled in stage  $i+1$  and the null vector  $(0; \dots; 0)$  has size  $k$ .

Note that in the  $M^{\text{th}}$  stage, the outcome is always a routing probability vector since the arrival is always routed to one of the potential destination servers. We now state precisely the class of routing policies that we study in this chapter under a common framework.

**Definition 3.1.** *Occupancy Based Randomized Routing Policy:*

Upon an arrival, in stage  $i$ , the routing policy samples  $d_i$  servers uniformly at random with replacement from  $N$  servers. The randomly selected or sampled servers are referred to as the potential destination servers for the arrival.<sup>1</sup> Then the routing policy in stage  $i$  is a mapping function for  $1 \leq i \leq M$  and  $1 \leq j \leq d^{(max)}$ ,

$$\Lambda^{(i)} : Y^{[i;j]} \rightarrow (0; p_1; \dots; p_j) \in (0; M_1(f_1; 2; \dots; jg)); \quad (3.6)$$

where  $p_i$  denotes the probability with which the  $i^{\text{th}}$  potential destination server is chosen as the destination server. Also,

$$\Lambda^{(i)} : \overline{Y^{[i]}} \setminus S^{[i;j]} \rightarrow (d_{i+1}; 0; \dots; 0) \in \mathbb{R}_+^{j+1}; \quad (3.7)$$

where  $d_{i+1}$  denotes the number of servers to be sampled in stage  $i+1$ . In stage 1, the dispatcher always samples a fixed number of  $d_1$  servers. If the routing policy enters the last stage  $M$ , then one of the potential destination servers is chosen as the destination server based on the mapping  $\Lambda^{(M)}$ .

## 3.2 Additional Notation and Terminology

In this section, we provide the required additional notation and terminology that is specific to this chapter.

For any  $f \in C_b^1(\mathbb{R}_+^n)$ , the function  $f^\Phi$  which will appear in the mean-field equation is defined as

$$f^\Phi(x_1; \dots; x_n) = \frac{1}{n} \sum_{i=1}^n \frac{\partial f(x_1; \dots; x_n)}{\partial x_i}$$

---

<sup>1</sup> As  $N \rightarrow \infty$ , sampling potential servers with or without replacement would yield same mean-field limit and hence, we use sampling with replacement to simplify the notation.

Since the rate at which a job is served depends on the number of progressing jobs of the server that processes it, we keep track of the status of a progressing job by using its age defined as the amount of service received since its arrival. We then use the notation  $(n; a_1; \dots; a_n)$  to keep track of the status of a server where  $n$  is the number of progressing jobs and  $a_i$  is the age of the  $i^{\text{th}}$  progressing job. If a job arrives at time  $t$ , then its age  $a$  at time  $t + h$  is given by

$$a = \int_{s=t}^{t+h} \frac{1}{r(s)} ds;$$

where  $r(s)$  denotes the instantaneous number of progressing jobs at the server at time  $s$ . Let  $U_n$  be the set of all possible states of a server when there are  $n$  progressing jobs, i.e.,

$$U_n = \{f(n; a_1; \dots; a_n) : a_i \geq 0, 1 \leq i \leq n\}.$$

If a server is idle, then its state lies in the set

$$U_0 = \{f(0)\}.$$

Let  $U$  be the set defined as

$$U = \bigcup_{n=0}^{\infty} U_n;$$

and hence,  $U$  is the set of all possible server states. Without loss of generality, we write an element of the form  $(n; u_1; \dots; u_n)$  as  $\mathbf{u}_n$  to indicate that it belongs to the set  $U_n$  for  $n \geq 0$  and also, we write  $(n; u_1; \dots; u_n)$  as  $\mathbf{u}$  to say that it is an element in  $U$ . We associate the space  $U$  with the metric  $d_U$  where for  $\mathbf{u} = (n; u_1; \dots; u_n)$  and  $\mathbf{v} = (m; v_1; \dots; v_m)$ ,  $d_U(\mathbf{u}; \mathbf{v})$  is defined as

$$d_U(\mathbf{u}; \mathbf{v}) = \begin{cases} \sum_{i=1}^n |u_i - v_i| & \text{if } n = m \\ |n - m| & \text{otherwise.} \end{cases}$$

For  $n \geq 1$ ,  $U_n$  is a complete and separable, and Polish space. Furthermore,  $U$  is a Polish space as it is the union of a countable set of disjoint Polish spaces. Also,  $U$  is separable and complete.

The restriction of a measure  $\mu \in \mathcal{M}_F(U)$  to the space  $U_0$  is a Dirac measure with mass at  $(0)$ . Also, if  $\mu(\{f(\mathbf{x}_n)\}) = 0$  for  $\mathbf{x}_n \in U_n$ ,  $n \geq 1$ , then we say that  $\mu$  is absolutely continuous at  $\mathbf{x}_n$ . If  $\mu(\{f(\mathbf{y}_n)\}) = 0$  for all  $\mathbf{y}_n \in U_n$ ,  $n \geq 1$ , then we say that  $\mu$  is absolutely continuous with respect to Lebesgue measure. For every  $i \geq 1$  and  $(x_1; \dots; x_i) \in \mathbb{R}_+^i$ , if  $\frac{\partial f(i; x_1; \dots; x_i)}{\partial x_j}$  is defined for all  $1 \leq j \leq i$ , then  $f$  is said to be differentiable. From this definition, our analysis uses the fact that  $f = I_{f|_{U_n}}$  is differentiable for  $n \geq 1$ .

Let  $\Xi : \mathbb{U} \rightarrow \mathbb{R}$  and  $\mathcal{H} : \mathbb{U} \rightarrow \mathbb{R}$  be defined as for  $(n; x_1; \dots; x_n) \in \mathbb{U}_n$ ,

$$\Xi(n; x_1; \dots; x_n) = n;$$

and

$$(n; x_1; \dots; x_n) = \begin{cases} 0 & \text{for } n = 0; \\ \sum_{i=1}^n x_i & \text{otherwise;} \end{cases}$$

Note that if a server state has distribution  $\mu$ , then  $h \mu$  and  $h \mu$  indicate the average number of progressing jobs and the average sum of ages of the progressing jobs, respectively. The following transition operators on functions and measures play a crucial role in the mathematical modeling of the time evolution of the system. For  $b > 0$ , let  $\tilde{b}^+ : \mathbb{U} \rightarrow \mathbb{U}$  be the transition operator defined as

$$\tilde{b}^+(n; u_1; \dots; u_n) = \begin{cases} 0 & \text{if } n = 0; \\ (n; v_1; \dots; v_n) & \text{otherwise;} \end{cases}$$

where  $v_i = u_i + \frac{b}{n}$ , for  $1 \leq i \leq n$ . Also, for any  $b > 0; f \in \mathcal{K}_b(\mathbb{U})$ , let  $\tilde{b} : \mathcal{K}_b(\mathbb{U}) \rightarrow \mathcal{K}_b(\mathbb{U})$  be the transition operator defined as

$$\tilde{b}f(\mathbf{u}) = f(\tilde{b}^+\mathbf{u});$$

for  $\mathbf{u} \in \mathbb{U}$ . Similarly, let  $\tilde{b} \in \mathcal{M}_F(\mathbb{U})$  be the shifted measure of  $\mu$  satisfying that for  $B \in \mathcal{B}(\mathbb{U})$

$$\tilde{b}(B) = \mu(\tilde{b}^+(B)).$$

For  $\mu \in \mathcal{M}_F(\mathbb{U})$ , the measure  $\tilde{b} \in \mathcal{M}_F(\mathbb{U})$  satisfies

$$h \tilde{b} \mu = h \mu; \tilde{b} \mu \tag{3.8}$$

for all  $f \in \mathcal{K}_b(\mathbb{U})$ . The uniqueness of the measure  $\tilde{b}$  in (3.8) follows from the Riesz-Markov-Kakutani theorem given in Appendix A.4 [80]. We use (3.8) in computing the generator of the measure-valued Markov process  $(\tilde{b}^{(N)}; t \geq 0)$  to be defined in Section 3.4 which describes the time evolution of the system.

### 3.3 System Dynamics

In this section, we first define the state of a server using which we then define the system state. The system state descriptor is defined such that the time evolution of the system

can be described by the time evolution of a measure-valued Markov process. As stated earlier, we keep track of the status of a progressing a job by using its age. We say that a server lies in the state  $(n; a_1; \dots; a_n)$  if it has  $n$  progressing jobs and the age of the  $i^{\text{th}}$  progressing job is equal to  $a_i$ ,  $1 \leq i \leq n$ . If a server is idle, then its state is equal to  $(0)$ . Since the routing policy is symmetric to servers as server identities have no role to play, it is sufficient to keep track of the total number of servers lying in each possible server state instead of each individual server's state to model the time evolution of the system by a measure-valued Markov process. Hence, we use the following system state descriptor.

**Definition 3.2.** *System State Descriptor:*

Let  $\mathbf{S}_i^{(N)}(t)$  be the state of the  $i^{\text{th}}$  server at time  $t$ , then the state of the system at time  $t$  is given by

$$\mathbf{S}_t^{(N)} = \sum_{i=1}^N \mathbf{s}_i^{(N)}(t). \quad (3.9)$$

Hence, the number of servers lying in a given state  $(m; z_1; \dots; z_m)$  at time  $t$  is equal to  $\mathbf{S}_t^{(N)}(f(m; z_1; \dots; z_m)g)$ .

For given  $\mathbf{S}_t^{(N)}$  and  $h > 0$ , the value of  $\mathbf{S}_{t+h}^{(N)}$  depends on the events that occur in the interval  $(t; t+h]$ . In the time interval  $(t; t+h]$ , there can be no event (arrival or departure) or some events (arrivals and departures) can occur in the system. By considering that  $h$  is sufficiently small, the probability that multiple events occur in the interval  $(t; t+h]$  can be neglected. The departure events can be modeled by using the hazard rate function  $\lambda(x) = \frac{g(x)}{G(x)}$  which indicates the instantaneous rate of departure of a job conditioned on its age is  $x$ . Precisely, if a job has age  $x$  at time  $t$ , then the probability that this job departs the system when its age is less than  $x+y$  is equal to  $\frac{G(x+y) - G(x)}{G(x)} = \lambda(x)y + o(y)$ . If the number of progressing jobs at a server with state say  $(n; a_1; \dots; a_n)$  at time  $t$  remains constant in the interval  $(t; t+h]$ , then its state becomes  $\tilde{\mathbf{S}}_h^+(n; a_1; \dots; a_n)$  at time  $t+h$ .

### 3.4 Main Results

In this section, we give the main results of this chapter.

Let us consider an occupancy based randomized routing policy which maps occupancies of the potential destination servers up to stage  $k$  to an outcome using a mapping function  $\Lambda^{(k)}$ . In stage  $k$ , for  $\mathbf{n} = (n_1; \dots; n_i) \in \mathcal{Y}^{[k;l]}$ , let  $q^{[k;l]}(\mathbf{n}; j)$  be the probability with which

the  $j^{\text{th}}$  potential destination server is chosen as the destination server out of  $i$  potential destination servers. We begin with the result on the probability with which the destination server of an arrival lies in the given state  $(m; x_1; \dots; x_m)$ .

**Lemma 3.1.** For  $\mathbf{x}_t^{(N)} = \mathbf{x}$ , let  $Q_n(\frac{\cdot}{N}) = \frac{f \cup_n g}{N}$  be the fraction of servers with  $n$  jobs at time  $t$ . Then, an arrival at time  $t$  is routed to a server with state  $(m; x_1; \dots; x_m)$  with probability equal to

$$p_r(\mathbf{x}; m; x_1; \dots; x_m) = \frac{f(m; x_1; \dots; x_m)g}{N} \Phi_m\left(\frac{\cdot}{N}\right); \quad (3.10)$$

where

$$\Phi_m\left(\frac{\cdot}{N}\right) = \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k; i]}} \sum_{j=1}^i I_{f \cap j = mg} q^{[k; i]}(\mathbf{n}; j) \prod_{r=1; r \neq j}^i Q_{n_r}\left(\frac{\cdot}{N}\right); \quad (3.11)$$

*Proof.* See Section 3.7.1. □

We consider the filtration

$$F_t^{(N)} = \sigma_{>0}(\mathbf{x}_s^{(N)}; s \leq t+); \quad (3.12)$$

Let  $A^{(N)}(\cdot)$  be the generator of the Markov process  $(\mathbf{x}_t^{(N)}; t \geq 0)$ . Using the Dynkin's formula [56], we have the following result.

**Theorem 3.1.** The process  $(\mathbf{x}_t^{(N)}; t \geq 0)$  is a Feller-Dynkin process in  $D_{M_F(U)}([0; \infty))$ . Let  $h \in C_b^1(U)$ , then the process  $(\mathbf{M}_t^{(N)}(\cdot); t \geq 0)$  defined as

$$\mathbf{M}_t^{(N)}(\cdot) = h(\mathbf{x}_t^{(N)}); \quad \mathbf{M}_0^{(N)}(\cdot) = h(\mathbf{x}_0^{(N)}); \quad \mathbf{M}_t^{(N)}(\cdot) = \int_{s=0}^t A^{(N)}h(\mathbf{x}_s^{(N)}); \quad ds \quad (3.13)$$

is a square integrable  $F_t^{(N)}$ -martingale and it is an RCLL process.

*Proof.* See Section 3.7.3. □

From the process  $(\mathbf{x}_t^{(N)}; t \geq 0)$ , let us define the normalized process  $(\mathbf{x}_t^{-(N)}; t \geq 0)$  where

$$\mathbf{x}_t^{-(N)} = \frac{\mathbf{x}_t^{(N)}}{N}; \quad (3.14)$$

The main objective of this chapter is to obtain the limiting process of the normalized sequence of processes  $f(\mathbf{x}_t^{-(N)}; t \geq 0)g_{N-1}$  when  $N \rightarrow \infty$ . For our analysis, we make the following assumption.

**Assumption 3.4.1.** The sequence  $f_0^{-(N)} g_{N-1}$  satisfies

$$(f_0^{-(N)}; h_0^{-(N)}; \Xi; h_0^{-(N)}; i) (f_0; h_0; \Xi; h_0; i); \quad (3.15)$$

where  $f_0 \in \mathcal{M}_1(\mathcal{U})$  is an absolutely continuous measure with  $h_0; \Xi < 1$  and  $h_0; i < 1$ .

We can interpret the conditions  $h_0; \Xi < 1$  and  $h_0; i < 1$  as follows. If  $f_0$  is the probability measure of the state of a server, then both the average number of jobs  $h_0; \Xi$  and the average sum of ages of jobs  $h_0; i$ , are finite.

On comparing Assumptions 3.4.1 and 2.4.1, we make an extra assumption to analyze the PS model that the average occupancy  $h_0; \Xi$  w.r.t. the limiting measure  $f_0$  is finite. This condition is used to show that the sequence  $f_t^{-(N)}; t \geq 0) g_{N-1}$  is tight.

We now state the main result of the chapter on the convergence of sequence  $f_t^{-(N)}; t \geq 0) g_{N-1}$  as  $N \rightarrow \infty$  to the mean-field limit.

**Theorem 3.2.** Under Assumptions 2.1.1 and 3.4.1, we show that  $(f_t^{-(N)}; t \geq 0) (f_t; t \geq 0)$  where the process  $(f_t; t \geq 0)$  is a unique solution to the following equation with the initial point  $f_0 \in \mathcal{M}_1(\mathcal{U})$ , for all  $f \in \mathcal{C}_b(\mathcal{U})$ ,

$$\begin{aligned} h_t; i = h_0; i + \int_{s=0}^t \left( \sum_{n=1}^{\infty} \sum_{i=1}^n \int_{x_1} \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} (f_{t-s}(n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) \right. \\ \left. - f_{t-s}(n; x_1; \dots; x_n) \right) d\tau_s(n; x_1; \dots; x_n) + \left[ (f_0; g) \Phi_0(f_s) (f_{t-s}(1; 0) - f_{t-s}(0)) \right. \\ \left. + \sum_{n=1}^{\infty} \sum_{i=1}^{n+1} \int_{x_1} \int_{x_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(f_s) (f_{t-s}(n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) \right. \\ \left. - f_{t-s}(n; x_1; \dots; x_n) \right) d\tau_s(n; x_1; \dots; x_n) \Big] ds; \quad (3.16) \end{aligned}$$

We refer to equation (3.16) as the mean-field equation and its unique solution is referred to as the mean-field limit.

*Proof.* See Section 3.7.4. □

For any  $T > 0$ , under the assumption of exchangeability of random variables that represent servers' states at time  $t = 0$ , any finite set of servers at time  $t = T$  are asymptotically independent. This proof is similar to the proof of Theorem 2.4 and hence, we omit the proof. Furthermore,  $f_t$  represents the distribution of a server's state as  $N \rightarrow \infty$  at time  $t$ .

We now show that the mean-field limit is the unique solution of a set of partial differential equations. We next show that the measure  $\bar{\nu}_t$  at time  $t$  is absolute continuous w.r.t. the Lebesgue measure and hence, it has a density function.

**Lemma 3.2.** *The measure  $\bar{\nu}_t$  at time  $t$  has a density function w.r.t. the Lebesgue measure for almost all  $\mathbf{u}_n \in \mathcal{U}_n$ ,  $n \geq 1$ .*

*Proof.* See Section 3.7.6. □

Let  $p_t(n; x_1; \dots; x_n)$  be the density of  $\bar{\nu}_t$  at  $(n; x_1; \dots; x_n)$  for  $n \geq 1$  and let  $P_t(0)$  be equal to  $\bar{\nu}_t(\mathbb{R}^0)$ . Then we construct a probability distribution function  $P_t = (P_t(\mathbf{u}); \mathbf{u} \in \mathcal{U})$  defined as

$$P_t(n; y_1; \dots; y_n) = \int_{x_1=0}^{y_1} \int_{x_n=0}^{y_n} p_t(n; x_1; \dots; x_n) dx_1 \dots dx_n; \quad (3.17)$$

where  $P_t(n; y_1; \dots; y_n)$  denotes the probability that at time  $t$ , a server has  $n$  jobs and the  $i^{\text{th}}$  job ( $1 \leq i \leq n$ ) has age at most  $y_i$  when  $N \rightarrow \infty$ . From the mean-field equation, since there exists a sequence of bounded continuous functions that converge monotonically to an indicator function of an open set, by using the monotone convergence theorem [69, Theorem 1.26], we obtain the following differential equations. The proof follows mutatis mutandis from the arguments of proof of Corollary 2.4 except that we replace  $\nu$  with  $\bar{\nu}$ .

**Corollary 3.1.** *The probability distribution function  $P_t = (P_t(\mathbf{u}); \mathbf{u} \in \mathcal{U})$  satisfies the PDEs*

$$\frac{dP_t(0)}{dt} = \int_{y=0}^1 \Phi(y) \left( \frac{\partial P_t(1; y)}{\partial y} \right) dy - \Phi_0(P_t)P_t(0) \quad (3.18)$$

and for  $n \geq 1$ ,

$$\begin{aligned} \frac{dP_t(n; y_1; \dots; y_n)}{dt} &= \sum_{i=1}^n \left\{ \frac{1}{n} \right\} \frac{\partial P_t(n; y_1; \dots; y_n)}{\partial y_i} \\ &+ \sum_{i=1}^{n+1} \int_{x_i=0}^1 \left\{ \frac{x_i}{n+1} \right\} \left( \frac{\partial P_t(n+1; y_1; \dots; y_{i-1}; x_i; y_i; \dots; y_n)}{\partial x_i} \right) dx_i \\ &\quad - \sum_{i=1}^n \int_{x_i=0}^{y_i} \left\{ \frac{x_i}{n} \right\} \left( \frac{\partial P_t(n; y_1; \dots; y_{i-1}; x_i; y_{i+1}; \dots; y_n)}{\partial x_i} \right) dx_i \\ &+ \sum_{i=1}^n \Phi_{n-1}(P_t)P_t(n-1; y_1; \dots; y_{i-1}; y_{i+1}; \dots; y_n) - \Phi_n(P_t)P_t(n; y_1; \dots; y_n); \quad (3.19) \end{aligned}$$



where  $\Phi_n(P_t)$  is the same as in equation (3.11) except that we replace  $Q_n(\bar{n})$  with  $Q_n(P_t) = \lim_{b! \rightarrow 1} P_t(n; b; \cdot; b)$ .

The above mean-field PDEs (3.18)-(3.19) correspond to the dynamics of a single server PS system in which the job arrival process is a Poisson process with rate  $\Phi_n(P_t)$  when the server has  $n$  jobs and the current distribution is  $P_t$  at time  $t$ . Hence, the MFEs represent the dynamics of a non-linear Markov process whose generator is a function of the current distribution of the Markov process.

The stationary behavior of the mean-field is discussed below. Let  $\bar{\gamma}$  be a fixed-point of the mean-field ( $P_t; t \rightarrow 0$ ) and let  $\bar{R}_n(\bar{\gamma}) = \sum_{j=n}^{\infty} \lim_{b! \rightarrow 1} (j; b; \bar{\gamma}; b)$ . Then we consider the following class of fixed-points  $\mathcal{Y}$  of the mean-field defined as

$$\mathcal{Y} = \left\{ \bar{\gamma} : \sum_{n=1}^{\infty} \bar{R}_n(\bar{\gamma}) < 1 \text{ and } \sum_{n=1}^{\infty} \prod_{i=1}^n \frac{\Phi_{i-1}(\bar{\gamma})}{\lambda} < 1 \right\}; \quad (3.20)$$

If  $\bar{\gamma} \in \mathcal{Y}$ , since  $\sum_{n=1}^{\infty} \bar{R}_n(\bar{\gamma}) < 1$ , the average occupancy obtained from  $\bar{\gamma}$  is finite. Furthermore, as  $\sum_{n=1}^{\infty} \prod_{i=1}^n \frac{\Phi_{i-1}(\bar{\gamma})}{\lambda} < 1$ , we can obtain a distribution from  $\bar{\gamma}$  as in (3.124) and (3.125) with arrival rate  $\Phi_n(\bar{\gamma})$  when there are  $n$  jobs. We exploit this logic to obtain the following result.

**Theorem 3.3.** Any fixed-point  $\bar{\gamma} = (\mathbf{u}; \mathbf{u} \geq \mathbf{U}) \in \mathcal{Y}$  of the mean-field ( $P_t; t \rightarrow 0$ ) satisfies

$$(n; y_1; \dots; y_n) = \Gamma_n \prod_{i=1}^n \int_{x_i=0}^{y_i} \bar{G}(x_i) dx_i; \quad (3.21)$$

where  $\Gamma = (\Gamma_n; n \geq 0)$  is a fixed-point of the mean-field when job lengths are exponentially distributed with the same mean  $\frac{1}{\lambda}$ . Under the assumption of exponential JLDs, if the mean-field of the occupancy process has a unique fixed-point, then the fixed-point is unique when job lengths are generally distributed. Furthermore, the fixed-point is insensitive since  $\int_{x=0}^{\infty} \bar{G}(x) dx = \frac{1}{\lambda}$  which implies

$$\lim_{b! \rightarrow 1} (n; b; \bar{\gamma}; b) = \Gamma_n; \quad (3.22)$$

*Proof.* See Section 3.7.2. □

The benefit of our mean-field analysis of occupancy based randomized routing schemes is that, to obtain the mean-field limit for any routing policy, we just need to obtain the rate function  $\Phi_m(\cdot)$  for  $m \geq 0$  given in (3.11) and the final mean-field equation is then given by (3.19).

We now conclude our analysis by applying our results to four policies. For a given vector  $\mathbf{n} = (n_1; \dots; n_k)$ , let  $\min(\mathbf{n})$  be the minimum value of the set  $\{n_1; \dots; n_k\}$  defined as

$$\min(\mathbf{n}) = \min_{j=1, \dots, k} n_j$$

and let

$$B_{\min}(\mathbf{n}) = \{j : n_j = \min(\mathbf{n}); 1 \leq j \leq k\}.$$

Suppose  $|B_{\min}(\mathbf{n})|$  be the size of the set  $B_{\min}(\mathbf{n})$ .

**Definition 3.3.** *Policy 1: This policy is the SQ( $d$ ) routing policy defined in Definition 2.1.1. For this policy, the number of stages is equal to one. For  $\mathbf{n} = (n_1; \dots; n_d)$ , we have*

$$q^{[1:d]}(\mathbf{n}; j) = \begin{cases} \frac{1}{|B_{\min}(\mathbf{n})|} & \text{if } n_j = \min(\mathbf{n}); \\ 0 & \text{otherwise;} \end{cases} \quad (3.23)$$

As a consequence, upon simplification, we get that

$$\Phi_m(P_t) = \frac{f(\bar{R}_m(P_t))^d (\bar{R}_{m+1}(P_t))^d g}{f\bar{R}_m(P_t) \bar{R}_{m+1}(P_t) g}, \quad (3.24)$$

where  $\bar{R}_m(P_t) = \sum_{j=m}^1 P_t(j; 1; \dots; 1)$ .

We presented the analysis of the SQ( $d$ ) policy for the case of general JLDs in [41].

**Definition 3.4.** *Policy 2: In this policy, upon an arrival,  $d$  servers are chosen uniformly at random. Out of these potential destination servers, a server is chosen as the destination server uniformly at random from the list of the potential destination servers whose occupancy is less than or equal to  $\frac{1}{a}$  if the list is not empty. Otherwise, the potential destination server that has the least occupancy is chosen as the destination server with ties broken uniformly at random. In this policy, there is only one stage and the number of sampled servers is equal to  $d$ . Then for  $\mathbf{n} = (n_1; \dots; n_d)$  and let  $a = \sum_{i=1}^d I_{\{n_i \leq g\}}$ , we have*

$$q^{[1:d]}(\mathbf{n}; j) = \begin{cases} \frac{1}{a} & \text{if } n_j \leq g; \\ \frac{1}{|B_{\min}(\mathbf{n})|} & \text{else if } n_j > g \text{ and } n_j = \min(\mathbf{n}); \\ 0 & \text{otherwise;} \end{cases} \quad (3.25)$$

As a result, upon simplification,

$$\Phi_m(P_t) = \begin{cases} \frac{f_1(\bar{R}_{m+1}(P_t))^d g}{f_1 \bar{R}_{m+1}(P_t) g} & \text{if } m = d; \\ \frac{f(\bar{R}_m(P_t))^d (\bar{R}_{m+1}(P_t))^d g}{f \bar{R}_m(P_t) \bar{R}_{m+1}(P_t) g} & \text{otherwise:} \end{cases} \quad (3.26)$$

**Definition 3.5.** *Policy 3:* In this policy, upon an arrival, the dispatcher selects a server uniformly at random. If the chosen server has occupancy  $i$ , then dispatcher selects further  $\min(i; d-1)$  servers uniformly at random, and the potential destination server with the minimum occupancy among all the potential destination servers is chosen as the destination server with ties broken uniformly at random. In this policy, there are two stages. In the first stage, the dispatcher first selects a server uniformly at random, and this server is chosen as the destination if it has no jobs. Otherwise, the policy enters the second stage in which  $\min(i; d-1)$  servers are selected uniformly at random. The destination server in the second stage is the server with the least occupancy among all the  $\min(i+1; d)$  potential destination servers. Therefore, we have

$$q^{[1;1]}(n_1; 1) = \begin{cases} 1 & \text{if } n_1 = 0; \\ 0 & \text{otherwise;} \end{cases} \quad (3.27)$$

and for  $\mathbf{n} = (n_1; \dots; n_d)$ ,

$$q^{[2;1]}(\mathbf{n}; j) = \begin{cases} \frac{1}{j B_{\min(\mathbf{n})}^j} & \text{if } n_j = \min(\mathbf{n}); \\ 0 & \text{otherwise:} \end{cases} \quad (3.28)$$

Then  $\Phi_m(\cdot)$  is given by

$$\Phi_0(P_t) = 1 + \sum_{i=1}^{d-2} Q_i(P_t) \frac{f(\bar{R}_0(P_t))^i (\bar{R}_1(P_t))^i g}{f \bar{R}_0(P_t) \bar{R}_1(P_t) g} + \bar{R}_{d-1}(P_t) \frac{f(\bar{R}_0(P_t))^{d-1} (\bar{R}_1(P_t))^{d-1} g}{f \bar{R}_0(P_t) \bar{R}_1(P_t) g}; \quad (3.29)$$

for  $m = d-1$ ,

$$\Phi_m(P_t) = (\bar{R}_m(P_t))^m + \sum_{i=m+1}^{d-1} Q_i(P_t) \frac{f(\bar{R}_m(P_t))^i (\bar{R}_{m+1}(P_t))^i g}{f \bar{R}_m(P_t) \bar{R}_{m+1}(P_t) g} + \bar{R}_d(P_t) \frac{f(\bar{R}_m(P_t))^{d-1} (\bar{R}_{m+1}(P_t))^{d-1} g}{f \bar{R}_m(P_t) \bar{R}_{m+1}(P_t) g}; \quad (3.30)$$

and for  $m = d$ ,

$$\Phi_m(P_t) = \frac{f(\bar{R}_m(P_t))^d (\bar{R}_{m+1}(P_t))^d g}{f\bar{R}_m(P_t) \bar{R}_{m+1}(P_t)g}, \quad (3.31)$$

where  $Q_i(P_t) = \bar{R}_i(P_t) - \bar{R}_{i+1}(P_t)$ .

**Definition 3.6.** *Policy 4: In this policy, the dispatcher selects sequentially at most  $d$  servers uniformly at random, and the routing decision is made after every selection. For  $1 \leq i \leq d-1$ , after the  $i^{\text{th}}$  selection, the arrival is routed to the  $i^{\text{th}}$  potential destination server if its occupancy is less than or equal to  $\bar{R}_i(P_t) - \bar{R}_{i+1}(P_t)$ . Otherwise, the dispatcher selects a server uniformly at random for the  $(i+1)^{\text{th}}$  time, and the same procedure is repeated as in the case of the  $i^{\text{th}}$  potential destination server. If the  $(d-1)^{\text{th}}$  potential destination has occupancy greater than  $\bar{R}_{d-1}(P_t) - \bar{R}_d(P_t)$ , then another server is selected uniformly at random, and the potential destination server with the minimum occupancy among all the  $d$  potential destination servers is chosen as the destination server with ties broken uniformly at random. In this case, the policy has  $d$  stages. In the  $i^{\text{th}}$  ( $1 \leq i \leq d-1$ ) stage, the dispatcher selects a server uniformly at random, and this server is selected as the destination if its occupancy is less than or equal to  $\bar{R}_i(P_t) - \bar{R}_{i+1}(P_t)$ . In the  $d^{\text{th}}$  stage, the server with the minimum occupancy is chosen as the destination. For  $1 \leq i \leq d-1$  and  $\mathbf{n} = (n_1; \dots; n_i)$ , we have*

$$q^{[i;1]}(\mathbf{n}; j) = \begin{cases} 0 & \text{if } i \neq j; \\ 1 & \text{else if } i = j; n_i \leq \bar{R}_i(P_t) - \bar{R}_{i+1}(P_t); \\ 0 & \text{else if } i = j; n_i > \bar{R}_i(P_t) - \bar{R}_{i+1}(P_t); \end{cases} \quad (3.32)$$

and if  $i = d$ , then

$$q^{[d;d]}(\mathbf{n}; j) = \begin{cases} \frac{1}{jB_{\min(\mathbf{n})}^j} & \text{if } n_j = \min(\mathbf{n}) \\ 0 & \text{otherwise:} \end{cases} \quad (3.33)$$

As a result, upon simplification, if  $m = d$ , then

$$\Phi_m(P_t) = \frac{1}{1} \frac{\bar{R}_{+1}^d(P_t)}{\bar{R}_{+1}(P_t)}. \quad (3.34)$$

On the other hand, if  $m < d$ , then

$$\Phi_m(P_t) = \frac{f(\bar{R}_m(P_t))^d (\bar{R}_{m+1}(P_t))^d g}{f\bar{R}_m(P_t) \bar{R}_{m+1}(P_t)g}. \quad (3.35)$$

**Remark 3.1.** *Since the arrival rate function  $\Phi_m(\cdot)$  is the same for both Policy 2 and Policy 4, they result in the same mean-field equation. Hence, the fixed-point of the mean-field coincides for these two policies. However, the number of servers that are randomly selected to implement Policy 4 is adaptive to the system's traffic level whereas Policy 2 always selects  $d$  servers to decide the destination server for an arrival. Hence, Policy 4 is preferable over Policy 2 due to its lower complexity.*

The main challenge in the mean-field analysis of a routing policy is to show that the fixed-point of the mean-field is unique. Furthermore, we need to show that the stationary distribution of a server state as  $N \rightarrow \infty$  coincides with the unique fixed-point of the mean-field. Under the assumption of exponential JLDs, for Policy 1, it was shown in [9] (as in the FCFS case in [7]) that the stationary distribution of a server's state coincides with the unique fixed-point of the mean-field. In this chapter, we show that this result is true for the routing policies 2 and 4 based on the result that the mean-field is quasi-monotonic in these cases. For Policy 3, we could not show the uniqueness of a fixed-point of the mean-field. Furthermore, the mean-field equation is a complex equation to prove its monotonicity. For this policy, proving the uniqueness of a fixed-point of the mean-field and its global stability are still open problems.

**Theorem 3.4.** *Under the assumption of exponential JLDs, for Policy 2 and Policy 4, the stationary distribution of a server state as  $N \rightarrow \infty$  coincides with the unique global asymptotically stable fixed-point of the mean-field.*

*Proof.* See Section 3.7.8. □

**Remark 3.2.** *For Policies 1, 2, and 4, from Theorem 3.3 and Theorem 3.4, since the exponential case has a unique fixed-point, there exists a unique fixed-point for the case of general JLDs. Furthermore, the fixed-point is insensitive to the JLDs. We conjecture that this also holds for Policy 3.*

## 3.5 Numerical Results

In this section, we present numerical results to investigate the impact of the four routing policies on the system performance in terms of the average response time for an arrival, complexity in terms of the average number of sampled servers per arrival, and the insensitivity of the stationary distribution of the stochastic system when  $N \rightarrow \infty$ . Our numerical

results suggest that as  $N \rightarrow \infty$ , the stationary distribution of occupancy of a server is insensitive for all the four policies which we study in this chapter. We choose the following parameters for our results:  $d = 3$ ,  $\rho = 1$ ,  $N = 300$ , and  $\beta = 1$ .

## System Performance Versus Complexity

We now study the trade-off between the system performance and complexity of a routing policy. For this, we assume that job lengths are exponentially distributed with mean 1. The unique fixed-point of the mean-field can be computed numerically by computing the unique probability vector  $\Gamma = (\Gamma_n; n = 0)$  that satisfies for  $n = 0$ ,

$$\Phi_n(\Gamma) = \Gamma_{n+1} \quad ; \quad (3.36)$$

where  $\Phi_n(\Gamma)$  is the same as in equation (3.11) except that we replace  $Q_n(\bar{N})$  with  $\Gamma_n$ . Then the average response time  $E(T_s)$  is given by

$$E(T_s) = \frac{1}{\rho} \left\{ \sum_{i=1}^d \left( \sum_{j=i}^d \Gamma_j \right) \right\} \quad ; \quad (3.37)$$

In Figure 3.1, we plot the resulting average response time versus the system load  $\rho$  for all the four routing policies and the optimal JSQ policy.

We now present results on the average number of probed servers in deciding the destination server for an arrival in Table 3.1. From Figure 3.1 and Table 3.1, it is clear that a

Table 3.1: The average number of probed servers per arrival

|      | Policy 1 | Policy 2 | Policy 3 | Policy 4 | JSQ |
|------|----------|----------|----------|----------|-----|
| 0.6  | 3        | 3        | 1.7875   | 1.3511   | 300 |
| 0.7  | 3        | 3        | 2.0003   | 1.5238   | 300 |
| 0.8  | 3        | 3        | 2.2535   | 1.7823   | 300 |
| 0.9  | 3        | 3        | 2.5769   | 2.2020   | 300 |
| 0.95 | 3        | 3        | 2.7690   | 2.5238   | 300 |

policy achieves better system performance if the dispatcher probes higher number of servers and selects the server with the maximum number of available resources as the destination server. Hence, the JSQ policy which chooses the server with the least occupancy out of  $N$  servers as the destination for an arrival achieves the best system performance [5, 6]. We observe that Policy 1 that selects the server with the least occupancy from a set of  $d$

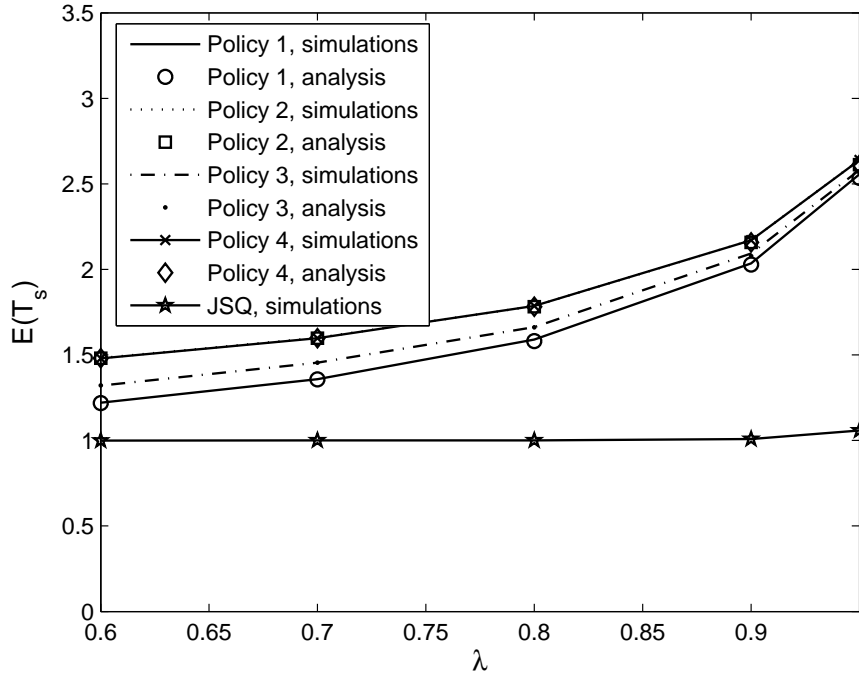


Figure 3.1: The average response time versus

randomly selected servers achieves good system performance but with a significant drop in complexity over the JSQ policy. From Figure 3.1, it can be seen that Policies 2 and 4 achieve the same system performance as they result in the same arrival rate function  $\Phi_n(\cdot)$  for  $n = 0$ , but Policy 4 probes a significantly smaller number of servers to dispatch an arrival when compared to Policy 2. On the other hand, Policy 3 probes higher number of servers when compared to Policy 4 and hence, it achieves better system performance. Both Policies 3 and 4 are robust to the variations in the system load as they probe lower number of servers when the system load is low, and they probe higher number of servers when the system load is high. Hence, all the Policies 1, 3, and 4 result in similar system performance when the system load is high. It is clear that Policy 3 is desirable in practice when there is a variation in the system load as it achieves the better trade-off between the performance and the complexity.

### Insensitivity of the System When $N \neq 1$

We now study numerically that for different classes of JLDs, when  $N \neq 1$ , the stationary distribution of occupancy of a server coincides with the unique fixed-point of the mean-field

of the case when job lengths are exponential with the same average job length.

Let  $\Gamma^{(N)} = (\Gamma_i^{(N)}; i = 0)$  be the stationary distribution of occupancy of a server in the system with parameter  $N$ . Here,  $\Gamma_i^{(N)}$  denotes the stationary probability that there are  $i$  jobs in progress at a server in the system with  $N$  servers in which arriving jobs are routed according to a predefined routing policy. In simulation results, we assume that servers are selected randomly without replacement, whereas in our analysis, we assume that servers are selected with replacement. We use PASTA property to compute  $\Gamma^{(N)}$  by simulating the system up to  $5 \cdot 10^6$  job arrivals. Let  $\Gamma^{(exp)} = (\Gamma_i^{(exp)}; i = 0)$  be the fixed point of the mean-field limit of the occupancy process and we compute  $\Gamma^{(exp)}$  numerically by the stationary mean-field equation. We now compute the total variation distance between  $\Gamma^{(N)}$  and  $\Gamma^{(exp)}$  defined as

$$\#_{tv}(\Gamma^{(N)}; \Gamma^{(exp)}) = \sum_i \left| \Gamma_i^{(N)} - \Gamma_i^{(exp)} \right| \quad (3.38)$$

We consider the following JLDs in our study: Exponential (Exp), Constant (Const), Power-law (PL), and mixed-Erlang (ME) distributions. The Power-law distribution has CDF

$$G(y) = \begin{cases} 1 - \frac{1}{3y^2} & \text{if } y \geq \frac{1}{3} \\ 0 & \text{otherwise:} \end{cases} \quad (3.39)$$

In the mixed-Erlang case, with probability  $p_i$  the distribution has  $i$  ( $1 \leq i \leq 2$ ) exponential phases and each exponential phase has intensity  $\rho$ . We choose  $p_1 = .4; p_2 = 0.6$ , and the value of  $\rho$  is computed by using the formula of the average job length given by

$$\sum_{i=1}^2 \frac{i p_i}{\rho} = \frac{1}{\rho}.$$

From Tables 3.2 and 3.3, it is evident that the fixed-point of the mean-field in the exponential case approximates  $\Gamma^{(N)}$  for different classes of JLDs having the same average job length.

### 3.6 On the Stationary Regime under the SQ( $d$ ) Policy

In this section, for the SQ( $d$ ) policy, we provide some numerical results to show that the fixed-point of the mean-field is GAS. The results of this section were presented in [41]. We consider the case of mixed-Erlang distributions that allow for efficient numerical evaluation



Table 3.2:  $\#_{tv}(\Gamma^{(N)}; \Gamma^{(exp)})$  for different JLDs for  $\rho = 0.7$  and  $N = 300$ .

| Policy   | <i>Exp</i> | <i>Const</i> | <i>PL</i> | <i>ME</i> |
|----------|------------|--------------|-----------|-----------|
| Policy 1 | 0.0020     | 0.0017       | 0.0018    | 0.0023    |
| Policy 2 | 0.0017     | 0.0017       | 0.0026    | 0.0015    |
| Policy 3 | 0.0029     | 0.0015       | 0.0016    | 0.0035    |
| Policy 4 | 0.0016     | 0.0021       | 0.0020    | 0.0016    |

Table 3.3:  $\#_{tv}(\Gamma^{(N)}; \Gamma^{(exp)})$  for different JLDs for  $\rho = 0.8$  and  $N = 300$ .

| Policy   | <i>Exp</i> | <i>Const</i> | <i>PL</i> | <i>ME</i> |
|----------|------------|--------------|-----------|-----------|
| Policy 1 | 0.0071     | 0.0044       | 0.0050    | 0.0061    |
| Policy 2 | 0.0038     | 0.0036       | 0.0038    | 0.0031    |
| Policy 3 | 0.0040     | 0.0041       | 0.0070    | 0.0066    |
| Policy 4 | 0.0049     | 0.0051       | 0.0040    | 0.0039    |

of the MFEs. We also discuss the propagation of chaos in the stationary regime and some existing related works. If one can prove the GAS of the mean-field, then we can exploit the Prohorov’s theorem [56] to conclude that the stationary distribution of a server’s occupancy as  $N \rightarrow \infty$  coincides with the fixed-point of the mean-field [9]. Proving the GAS of the mean-field is extremely difficult since the mean-field does not possess any monotonicity properties when job lengths are generally distributed unlike the exponential case studied in [9, 37]. We also provide simulation results to support the convergence of the stationary distribution of a server’s occupancy in a finite  $N$  system to the fixed-point of the mean-field when  $N \rightarrow \infty$ . Taken together these results provide evidence of the insensitivity of the stationary distribution of the limiting system as  $N \rightarrow \infty$ .

The numerical evaluation of the MFEs when JLDs are mixed-Erlang using the Euler’s method was performed with a step size of  $2 \times 10^{-3}$ . From the case of exponential distributions, we know that the stationary probability that there are at least  $k$  jobs at a server under the SQ( $d$ ) policy in the limiting system is given by  $(-\rho)^{\frac{d^k}{d-1}}$  [9] for given values of  $\rho < 1$ . We assume that the servers have a finite buffer size of  $C$  chosen such that  $(-\rho)^{\frac{d^C}{d-1}}$  is negligible. We consider the system parameters as follows. The job lengths have Mixed-Erlang distributions under which a job length is sampled with probability  $\rho_i$  ( $i = 1, 2, \dots, D$ ) from an Erlang distribution having  $i$  exponential phases with rate  $\mu_i$ .

Let us define the state of a server with  $n$  progressing jobs having  $l_j$  phases remaining

for the  $j^{\text{th}}$  progressing job by  $\mathbf{l} = (n; l_1; \dots; l_n)$  with  $1 \leq l_j \leq D$ ,  $1 \leq j \leq n$ . For  $n \geq 1$ , let  $S_n = f(n; l_1; \dots; l_n) : 1 \leq l_i \leq M; 1 \leq i \leq n$  be the set of all possible states of a server when there are  $n$  progressing jobs and  $S_0 = f(0)g$  denotes the state of a server when there are no progressing jobs. We then define  $S$  to be the set of all possible states of servers given by

$$S = \bigcup_{n=0}^C S_n. \quad (3.40)$$

We can model the system evolution by using a Markov process  $\mathbf{X}^N(t) = (\mathbf{X}_i^N(t); \mathbf{l} \in S)$  where  $\mathbf{X}_i^N(t)$  denotes the fraction of servers lying in state  $\mathbf{l}$  at time  $t$ . Since the underlying space  $S$  is countable and finite-dimensional, the mean-field limit can be established by the same procedure as that of the exponential case [9].

It can be shown that the mean-field equations are given by the following system of ordinary differential equations,

$$\dot{\mathbf{x}}(0; \mathbf{u}) = \mathbf{u}; \quad (3.41)$$

$$\dot{\mathbf{x}}_i(t; \mathbf{u}) = h_i(\mathbf{x}(t; \mathbf{u})); \quad (3.42)$$

and  $\mathbf{h} = (h_i; \mathbf{l} \in S)$  with the mapping  $h_i$  given by

$$\begin{aligned} h_i(\mathbf{x}(t; \mathbf{u})) &= \sum_{b=1}^n \left\{ \frac{\rho_{l_b}}{n} \right\} X_{(n-1; l_1; l_2; \dots; l_{b-1}; l_{b+1}; \dots; l_n)}(t; \mathbf{u}) \stackrel{(ME)}{n-1}(\mathbf{x}(t; \mathbf{u})) - X_i(t) \stackrel{(ME)}{n}(\mathbf{x}(t; \mathbf{u})) I_{fn < C} g \\ &\quad + \sum_{b=1}^{n+1} \left\{ \frac{\rho}{n+1} \right\} I_{fn < C} g X_{(n+1; l_1; \dots; l_{b-1}; l; l_{b+1}; \dots; l_n)}(t; \mathbf{u}) \\ &\quad + \sum_{b=1}^n \left\{ \frac{\rho}{n} \right\} X_{(n; l_1; \dots; l_{b-1}; l_{b+1}; l_{b+1}; \dots; l_n)}(t; \mathbf{u}) - \rho X_{(n; l_1; \dots; l_n)}(t; \mathbf{u}); \end{aligned} \quad (3.43)$$

where

$$\stackrel{(ME)}{n}(\mathbf{r}) = \left\{ \frac{\rho}{\sum_{\mathbf{l} \in S_n} r_{\mathbf{l}}} \right\} \left[ \left( \sum_{i=n}^C \sum_{\mathbf{l} \in S_i} r_{\mathbf{l}} \right)^d - \left( \sum_{j=n+1}^C \sum_{\mathbf{l} \in S_j} r_{\mathbf{l}} \right)^d \right]; \quad (3.44)$$

We numerically solved the MFEs by choosing the following parameters:  $d = 2$ ,  $\rho = 1$ ,  $C = 7$ ,  $D = 2$ ,  $\rho_1 = 0.4$  and  $\rho_2 = 0.6$ .

The unique fixed-point  $\mathbf{x}^* = (x_i^*; \mathbf{l} \in S)$  is given by

$$x_{(n; l_1; \dots; l_n)}^* = \stackrel{(exp)}{(n)} \prod_{i=1}^n \left\{ \frac{\sum_{j=l_i}^D \rho_j}{\sum_{r=1}^D r \rho_r} \right\}; \quad (3.45)$$

where  $\mathbf{y}^{(exp)} = (y_n^{(exp)}(t); 0 \leq n \leq C)$  is the unique fixed-point of the mean-field in the exponential case.

Furthermore, we define a process  $\mathbf{y}(t; \mathbf{v}) = (y_n(t; \mathbf{v}); 0 \leq n \leq C)$  referred to as the tail mean-field that satisfies  $\mathbf{y}(0; \mathbf{v}) = \mathbf{v}$  and  $y_j(t; \mathbf{v}) = \sum_{i=j}^C \sum_{\mathbf{1} \in \mathcal{S}_i} x_{\mathbf{1}}(t; \mathbf{u})$ , and its fixed-point is denoted by  $\mathbf{y}^{(exp)}$ . We now define following metrics,

$$d_{tv}(\mathbf{a}; \mathbf{b}) = \sum_{\mathbf{1} \in \mathcal{S}} j a_{\mathbf{1}} - b_{\mathbf{1}}; \quad (3.46)$$

$$\#_{tv}(\mathbf{w}; \mathbf{z}) = \sum_{0 \leq n \leq C} j W_n - Z_n; \quad (3.47)$$

$$\#_{ds}(\mathbf{w}; \mathbf{z}) = \sum_{0 \leq n \leq C} (W_n - Z_n); \quad (3.48)$$

We observed that the mean-field is GAS. In particular, in Figure 3.2, we plot  $d_{tv}(\mathbf{x}(t; \mathbf{u}); \mathbf{y}^{(exp)})$  as a function of  $t$  for three different initial points  $\mathbf{u}_1; \mathbf{u}_2$  and  $\mathbf{u}_3$ . Similarly, in Figure 3.3, we plot  $\#_{tv}(\mathbf{y}(t; \mathbf{v}); \mathbf{y}^{(exp)})$  as a function of  $t$  for three different initial points  $\mathbf{v}_1; \mathbf{v}_2$  and  $\mathbf{v}_3$ . The mean-field  $\mathbf{x}(t; \mathbf{u})$  and its tail mean-field  $\mathbf{y}(t; \mathbf{v})$  converge to their fixed-points for three different initial points. This provides evidence for the result that  $\mathbf{x}(t; \mathbf{u})$  and  $\mathbf{y}(t; \mathbf{v})$  are globally stable. However,  $\#_{tv}(\mathbf{y}(t; \mathbf{v}); \mathbf{y}^{(exp)})$  is not monotonically decreasing with  $t$ . Hence, the total variation distance cannot be used for constructing a Lyapunov function to show the global asymptotic stability of the mean-field  $\mathbf{x}(t; \mathbf{u})$  and its tail mean-field  $\mathbf{y}(t; \mathbf{v})$ .

We also observed that the tail mean-field is not quasi-monotonic, i.e., if  $\mathbf{v}_1 \leq \mathbf{v}_2$  (by component wise), then we do not always have  $\mathbf{y}(t; \mathbf{v}_1) \leq \mathbf{y}(t; \mathbf{v}_2)$ . We recall that when the job lengths are exponential, the tail mean-field is quasi-monotonic. Figure 3.4 plots  $\#_{tv}(\mathbf{y}(t; \mathbf{v}); \mathbf{y}^{(exp)})$  and Figure 3.5 plots  $\#_{ds}(\mathbf{y}(t; \mathbf{v}); \mathbf{y}^{(exp)})$  as a function of  $t$  for different initial points. In particular, we select  $v_4$  and  $v_7$  to be less than  $\rho$  whereas  $v_5; v_6; v_8$ , and  $v_9$  are chosen to be greater than  $\rho$ . The violation of the quasi-monotonicity property is evident for the case of  $\rho = 0.9$  with initial point  $v_7$ . This is because since  $v_7 < \rho$ , if we have quasi-monotonicity,  $\#_{ds}(\mathbf{y}(t; \mathbf{v}); \mathbf{y}^{(exp)})$  must be always a negative value. This property does not hold as we can see in Figure 3.5 and hence, the mean-field is not quasi-monotonic.

The SQ( $d$ ) policy even for small values of  $d$  improves the system performance significantly. We plot the average response time ( $E(T_s)$ ) of a job under the SQ( $d$ ) policy and the random routing scheme ( $d = 1$ ) in Figure 3.6. The expression for  $E(T_s)$  under the SQ( $d$ ) policy is given by

$$E(T_s) = \frac{1}{d} \sum_{i=1}^d \left( - \right)^{\frac{d}{d}-1} \quad (3.49)$$

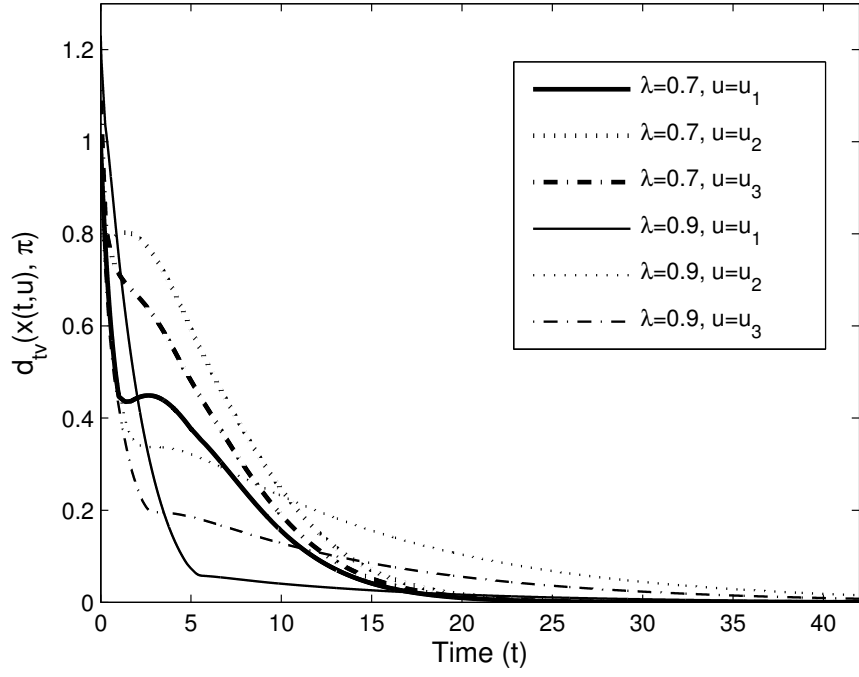


Figure 3.2:  $d_{tv}(x(t; u); \cdot)$  versus  $t$

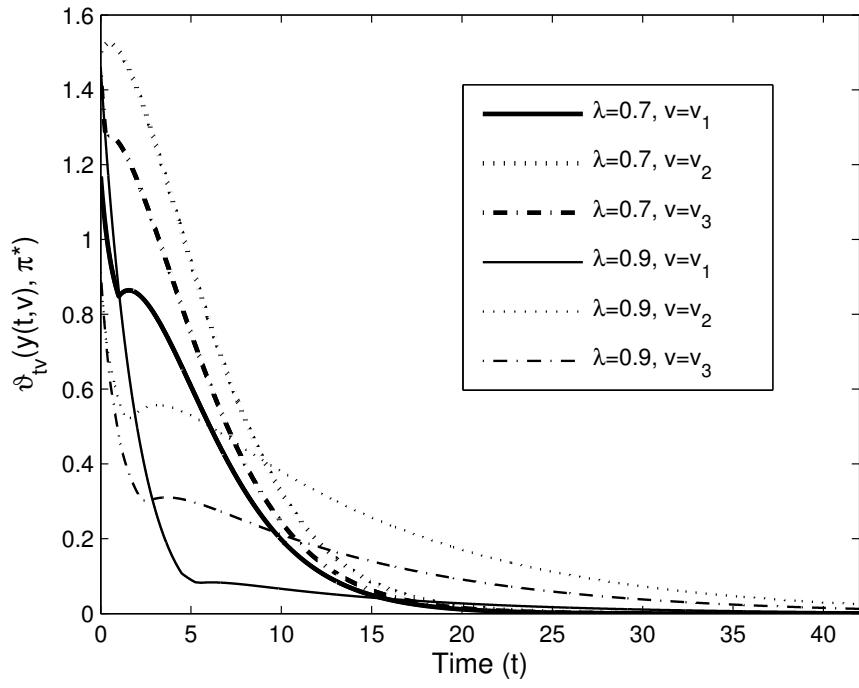


Figure 3.3:  $\phi_{tv}(y(t; v); \cdot)$  versus  $t$

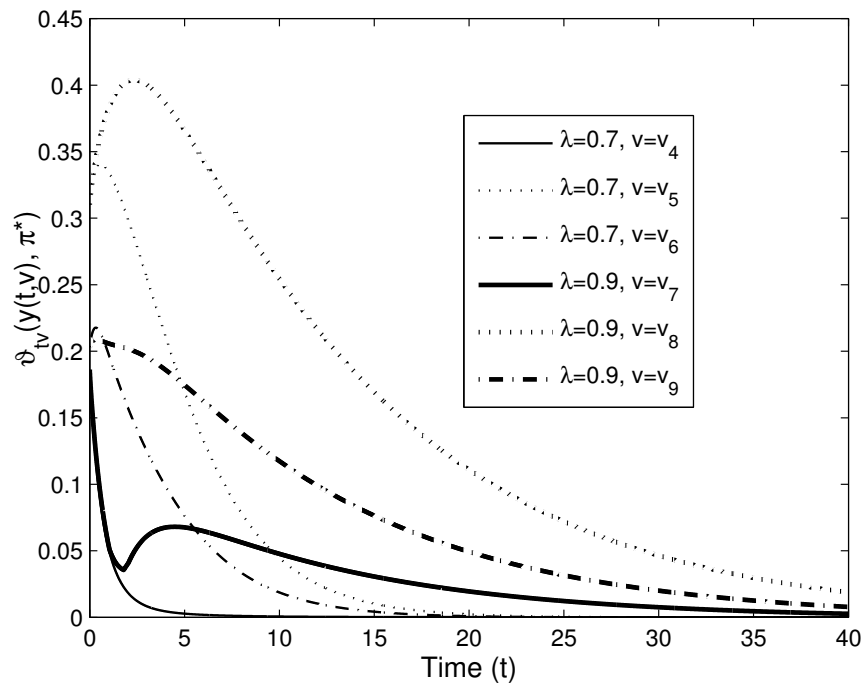


Figure 3.4:  $\vartheta_{tv}(y(t; \nu); \pi^*)$  versus  $t$

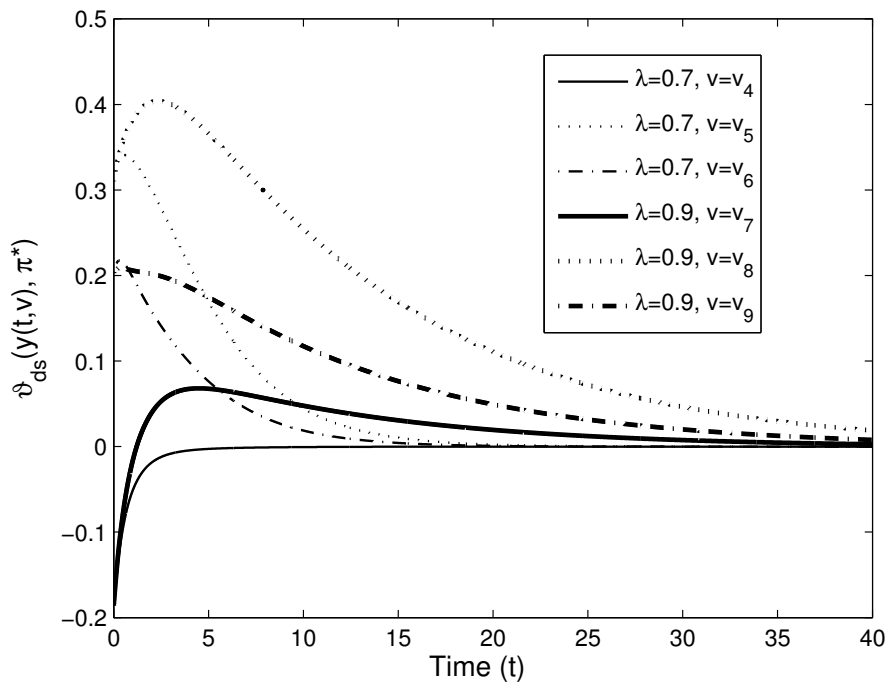


Figure 3.5:  $\vartheta_{ds}(y(t; \nu); \pi^*)$  versus  $t$

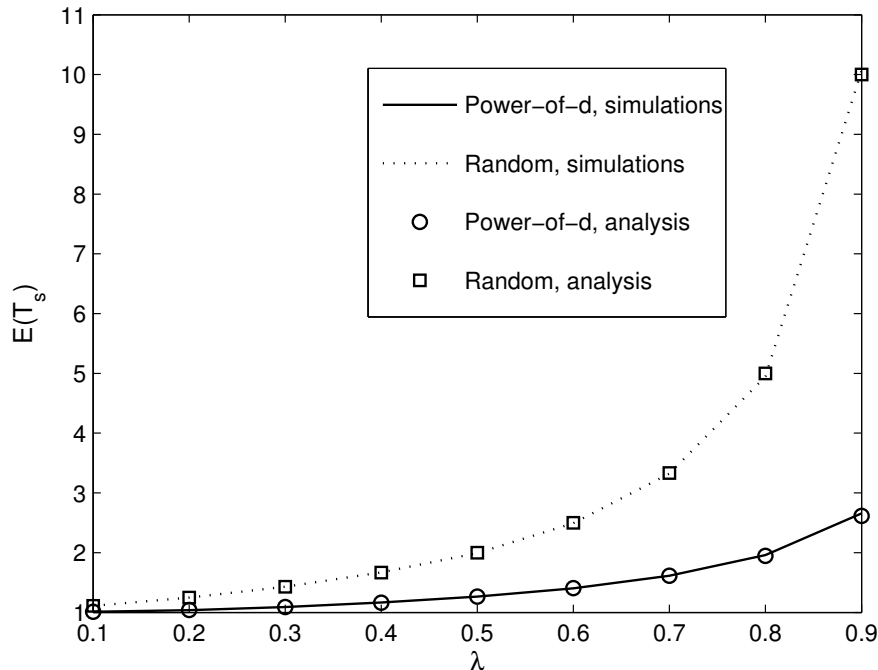


Figure 3.6: The average response time ( $E(T_s)$ ) versus

and for the random routing, we have

$$E(T_s) = \frac{1}{\lambda} \sum_{i=1}^{\infty} \binom{d-1}{i-1} \rho^i \quad (3.50)$$

We also plot the simulation results in Figure 3.6 by considering a system with  $N = 100$  and exponential JLDs. It is clear that the SQ( $d$ ) policy reduces the average response time significantly over the random routing policy ( $d = 1$ ).

### 3.6.1 On the Propagation of Chaos in the Stationary Regime

We now discuss the relationship between the propagation of chaos in the stationary regime, the tightness of the stationary distributions  $f^{(N)} g_{N-1}$  of the empirical process, and the global asymptotic stability of the fixed-point of the mean-field. For simplicity, we assume that the JLDs are mixed-Erlang and each server has finite buffer  $C$  and hence, there exists a unique invariant distribution  $\pi^{(N)}$  of the Markov process  $(\mathbf{X}^{(N)}(t); t \geq 0)$  where  $\mathbf{X}^{(N)}(t) = (\mathbf{X}_1^{(N)}(t); 1 \leq S)$ . In this case, the mean-field equations are given by equations (3.41)-(3.43). However, the system does not exhibit monotonicity properties unlike the simple exponential case, and thus establishing the propagation of chaos is a challenging problem [40]. When

the mean-field is GAS, by invoking the Prohorov's theorem (Theorem A.2), we can establish  $\lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \mathbf{X}^{(N)}(t) = \lim_{t \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbf{X}^{(N)}(t)$  where  $\lim_{t \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbf{X}^{(N)}(t)$  is the fixed-point of the mean-field [9]. This implies the validity of the interchange of limits

$$\lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \mathbf{X}^{(N)}(t) = \lim_{t \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbf{X}^{(N)}(t):$$

Furthermore, this would then imply that the propagation of chaos holds in the stationary regime [9, 61]. Thus it appears that the global stability of the mean-field, propagation of chaos, and the coincidence of the stationary distribution as  $N \rightarrow \infty$  with the fixed point of the mean-field are all inter-related. We now discuss what happens when we cannot show that the mean-field is GAS.

Since the space of probability measures on  $S$  denoted by  $M_1(S)$  with metric induced by total variation distance is compact, from the Prohorov's theorem [57], the sequence  $\{f^{(N)}g_{N-1}\}$  is tight in  $M_1(M_1(S))$  under the topology induced by the weak convergence. Let  $\{f^{(N_k)}g_{N_k-1}\}$  be a converging subsequence with limiting point  $Z \in M_1(M_1(S))$ . Then we say that for the sequence of systems with index  $f^{(N_k)}g_{N_k-1}$ , the limiting system is said to have the stationary distribution  $Z$  for the empirical random variable. Then from Theorem 1 of [61] and Section 1.3,  $Z$  is an invariant distribution of the mean-field  $\mathbf{x}(t; \cdot)$  that means

$$\int_{M_1(S)} f(\mathbf{x}(t; \mathbf{u})) dZ(\mathbf{u}) = \int_{M_1(S)} f(\mathbf{u}) dZ(\mathbf{u}): \quad (3.51)$$

Furthermore, from Theorem 3 of [62], the support of  $Z$  is a compact set included inside the Birkhoff center of the mean-field that includes the existing limit cycles, fixed-points of the mean-field.

Let  $\mathbf{S}_i^{(N_k)}(t)$  be the random variable that denotes the state of the  $i^{\text{th}}$  server in the stationary regime in a finite  $N_k$  system. Let  $\mathbf{V}^{(N_k)}(t), \mathbf{V}(t)$  be random variables with distributions  $\mu^{(N_k)}$  and  $Z$ , respectively. Note that since the system behavior is symmetric to servers as servers' labels do not play any role, the set  $(\mathbf{S}_i^{(N_k)}(t); 1 \leq i \leq N_k)$  is exchangeable irrespective of the initial conditions on  $(\mathbf{S}_i^{(N_k)}(t); 1 \leq i \leq N_k)$  in the transient regime. Let us consider continuous bounded mappings  $\phi_i: S \rightarrow \mathbb{R}_+, 1 \leq i \leq l$ .

**Theorem 3.5.** *If  $\mu^{(N_k)} \rightarrow Z$  as  $k \rightarrow \infty$ , then*

$$\mathbb{E} \left[ \prod_{i=1}^l \phi_i(\mathbf{S}_i^{(N_k)}(t)) \right] \rightarrow \mathbb{E} \left[ \prod_{i=1}^l \phi_i(\mathbf{V}(t)); \phi_i \right] \quad (3.52)$$

as  $k \rightarrow \infty$ . Any finite set of servers  $(n_i)_{1 \leq i \leq l}$  in the limiting system of the sequence  $\{f^{(N_k)}g_{N_k-1}\}$  are mutually independent if and only if  $Z$  is a Dirac measure. Furthermore, if

$Z = \delta_{\mathbf{u}^*}$  for some  $\mathbf{u}^* \in \mathcal{M}_1(S)$ , then each server's state is a random variable with distribution  $Z$ .

*Proof.* See Section 3.5. □

Since  $Z$  is an invariant distribution of the mean-field  $\mathbf{x}(t; \cdot)$ , from equation (3.51)

$$\mathbb{E} \left[ \prod_{i=1}^l \mathbf{S}_i^{(N_k)}(\mathbf{x}(t; \cdot)) \right] = \int_{\mathbf{u} \in \mathcal{M}_1(S)} \left( \prod_{i=1}^l h_{\mathbf{x}(t; \mathbf{u})}(\cdot; i) \right) dZ(\mathbf{u}) \quad (3.53)$$

as  $k \rightarrow \infty$ . From (3.53), at any time  $t$  in the stationary regime, servers are coupled through the position of the mean-field  $\mathbf{x}(t; \cdot)$  which is a random element since its initial point is random with distribution  $Z$ . Furthermore, in the limiting system, at any instant  $t$  in the stationary regime, conditioned on the position of the mean-field, each server's state is a random variable with distribution coinciding with the position of the mean-field, and any finite set of servers are independent. However, if the position of the initial point of the mean-field is random, then at any time  $t$ , the position of the mean-field is random thereby any finite set of servers are coupled through the position of the mean-field. For example, if the support of  $Z$  contains existing limit cycles or multiple fixed-points of the mean-field, then at any instant in the stationary regime, the position of the mean-field is random, as a consequence, any finite set of servers are coupled through the position of the mean-field. Since the mean-field has a unique fixed-point under the assumption of mixed-Erlang JLDs, we must have  $Z = \delta_{\mathbf{u}^*}$  to have the propagation of chaos in the stationary regime. If we show that the fixed-point  $\mathbf{u}^*$  is GAS, then we have  $Z = \delta_{\mathbf{u}^*}$ .

**Corollary 3.2.** *If  $(N_k) \rightarrow \infty$ ,  $Z$ , from (3.53), by taking  $l = 1$ , we get that  $\mathbb{E}[\mathbf{V}(1)]$  is the stationary distribution of a server in the limiting system of the sequence of systems  $fN_k g_k$ . That is, the average position of the mean-field with an invariant distribution  $Z$  indicates the stationary distribution of a single particle when  $Z$  is not a Dirac measure.*

**Remark 3.3.** *In [61], an example with the Birkhoff center containing a limit cycle and an unstable unique fixed-point was found. In this case, it was argued that the distributions of any two particles' states coincide and oscillate around the limit cycle. Hence, the two particles are correlated. But a proof was not provided. More importantly, they did not give any details about the distribution of a single particle when  $Z$  is not a Dirac measure.*



### 3.6.2 Discussion on a Related Work

In this section, we provide some insights by relating our work with an existing work for the sake of clarity. In the literature, insensitivity of the stationary distribution of the limiting system is claimed in [24] based on an *ansatz* that has not been shown to hold for PS models to the best of our knowledge. Using Theorem 3.5, we demonstrate that the *ansatz* in [24] implies the necessary and sufficient conditions required for the insensitivity of the limiting system. Therefore, we need to show that the *ansatz* holds to conclude insensitivity.

Let  $\mathbf{Y}^{(N)}(t) = (\mathbf{r}^{1:N}(t); \mathbf{r}^{2:N}(t); \dots; \mathbf{r}^{N:N}(t))$  be the joint queue-size process at time  $t$  where  $\mathbf{r}^{i:N}(t)$  (notation  $q^{i:n}(t)$  is used in [24]) denotes the number of jobs at server  $i$  at time  $t$ . Let  $\Upsilon^{(N)}$  ( $\Upsilon$  is replaced with  $\cdot$  in [24]) be the stationary distribution of  $\mathbf{Y}^{(N)}(t)$ . We now recall exactly the *ansatz* stated in [24].

**The Ansatz in [24]:** *Demonstrate  $\Upsilon^{(N)} \rightarrow \Upsilon$  as  $N \rightarrow \infty$ , where  $\Upsilon$  is a stationary and ergodic measure on  $Z_+^1$ . Show that the limit  $\Upsilon$  is unique, depending only on the service distribution, service discipline and load balancing rule. Let  $\Upsilon_{(k)}$  be the restriction of  $\Upsilon$  to its first  $k$  coordinates, with  $\Upsilon_{(1)} = \Upsilon_{(1)}$  being the one-dimensional marginal of  $\Upsilon$ . Show that, for every  $k$ ,*

$$\Upsilon_{(k)} = \mu_{(k)} \quad (3.54)$$

Let  $\overline{\mathbf{W}}^{(N)}(\gamma) = (\overline{\mathbf{W}}_i^{(N)}(\gamma); 0 \leq i \leq C)$  where  $\overline{\mathbf{W}}_i^{(N)}(\gamma)$  denotes the random variable in the stationary regime indicating the fraction of servers with  $i$  jobs. Then from Theorem 3.5 (also Proposition 2.1 of [32]),  $\overline{\mathbf{W}}^{(N)}(\gamma) \rightarrow \overline{\mathbf{W}}$  where  $\overline{\mathbf{W}}$  is a deterministic measure in  $M_1(S)$ . Since the SQ( $d$ ) policy uses the queue-size information of a finite set of  $d$  randomly sampled servers that are independent with identical distributions coinciding with  $\overline{\mathbf{W}}$ , the arrival process is a Poisson process to any particular server which is a necessary condition to have insensitivity in PS systems [76]. Furthermore, the arrival process to each server is a state-dependent Poisson process with rate  $\lambda_k = \frac{f(\binom{C}{j=k}^d) \binom{C}{j=k+1}^d g}{k}$  when there are  $k$  jobs at the server. Therefore, the set of arrival rates  $\Lambda = (\lambda_k; 0 \leq k \leq C)$  can be written as a function of  $\overline{\mathbf{W}}$  as

$$\Lambda = \tilde{F}_1(\overline{\mathbf{W}}) \quad (3.55)$$

Further, for a given set of arrival rates  $\Lambda$ , the stationary distribution for a server occupancy in the limiting system can be written through a mapping  $\tilde{F}_2$  as

$$\overline{\mathbf{W}} = \tilde{F}_2(\Lambda) \quad (3.56)$$

Therefore,  $\mu$  must be a unique fixed-point of the mapping  $\tilde{F}_2(\tilde{F}_1)$  for the case of general JLDs which was not shown in [24] except for the case of FCFS queues with service time distributions having decreasing hazard rate functions. To have insensitivity,  $\mu$  must be the same for all the general JLDs having the same average job length. In [24], from uniqueness of  $\mu$  in *ansatz*, insensitivity is concluded from reversibility since the arrival process to each server is a state-dependent Poisson process. Note that since the mappings  $\tilde{F}_2, \tilde{F}_1$  are the same for exponential and general distributions, the uniqueness of the fixed-point of  $\tilde{F}_2(\tilde{F}_1)$  follows from the result that the mean-field is GAS in the exponential case. Therefore, the fixed-point of the mapping  $\tilde{F}_2(\tilde{F}_1)$  is the same for both the exponential and general job length distributions when they have the same average job lengths. Since the *ansatz* implies the Poisson arrival process to servers and uniqueness of the stationary distribution in the limiting system, insensitivity follows immediately. However, the proof of the *ansatz* remains an open problem and has been shown only for the case of FCFS queueing models with service time distributions having decreasing hazard rate functions, and for more general systems with small arrival rates in [40]. It remains an open problem for loss and PS models. For these systems, we do not have ordering relationships on servers' states in  $\mathcal{U}$ . For FCFS models, an element in  $\mathcal{U}$  is two-dimensional of the form  $(n; \chi)$  where  $n$  indicates the number of jobs at the queue and  $\chi$  is the age or residual service time of the progressing job. Hence, we can define ordering relationships for elements in  $\mathcal{U}$  for FCFS models in that for two elements  $(n_1; \chi_1)$  and  $(n_2; \chi_2)$ , we say  $(n_1; \chi_1) \prec (n_2; \chi_2)$  if  $n_1 < n_2$  and  $\chi_1 < \chi_2$ . We can not define such relationships for the loss and PS models. Even the analysis for JLDs with decreasing hazard rate functions is also an open problem for loss and PS models.

Unlike the approach of [24], which makes several assumptions on the stationary regime, we first focus on establishing the mean-field limit by studying the transient regime. We then show that the probability distribution of occupancy obtained from every fixed-point of the mean-field is a fixed-point of the mapping  $\tilde{F}_2(\tilde{F}_1)$  which is the same mapping for the exponential case. Since we know that the exponential case has a unique fixed-point, the fixed-point is insensitive. If we can show that the mean-field is GAS, then the fixed-point is also the stationary distribution of a server when  $N \rightarrow \infty$ .

### 3.7 Proofs of Main Results

In this section, we provide proofs of the results stated in Section 3.4.

### 3.7.1 Proof of Lemma 3.1

*Proof.* An arrival can be routed in any one of the stages to a potential destination server. Let  $k$  be the index of the stage in which the arrival is routed to a server. Furthermore, let us assume that the dispatcher has sampled  $i$  potential destination servers by the end of stage  $k$  and let the vector of occupancies of the potential destination servers be equal to  $\mathbf{n} = (n_1; \dots; n_i) \in \mathcal{Y}^{[k;\ell]}$ . Also, let the  $j^{\text{th}}$  potential destination server be in the state  $(m; x_1; \dots; x_m)$ . Then since each server can be chosen as a potential destination server with probability  $\frac{1}{N}$ , the probability that the potential destination servers' occupancies form the vector  $\mathbf{n} = (n_1; \dots; n_i)$  and the  $j^{\text{th}}$  server lies in the state  $(m; x_1; \dots; x_m)$  is equal to

$$\left\{ \frac{(f(m; x_1; \dots; x_m)g)}{N} \right\} \prod_{r=1; r \neq j}^i \left\{ \frac{(U_{n_r})}{N} \right\};$$

Then the  $j^{\text{th}}$  potential destination server can be chosen as the destination server with probability  $q^{[k;\ell]}(\mathbf{n}; j)$  and hence, the probability that the occupancy vector of potential destination servers is  $\mathbf{n}$  and the destination server lies in the  $j^{\text{th}}$  position with state  $(m; x_1; \dots; x_m)$  is equal to

$$q^{[k;\ell]}(\mathbf{n}; j) \left\{ \frac{(f(m; x_1; \dots; x_m)g)}{N} \right\} \prod_{r=1; r \neq j}^i \left\{ \frac{(U_{n_r})}{N} \right\};$$

Finally, by summing over all possible values of  $k; \mathbf{n}; i$ , and  $j$ , we get (3.10).  $\square$

### 3.7.2 Proof of Theorem 3.3

Let  $\mathbf{v} = (\mathbf{v}; \mathbf{v} \in \mathcal{U})$  be a fixed-point of the process  $(P_t; t \geq 0)$  in the set  $\mathcal{Y}$ . Let  $v_n = \prod_{i=1}^n \frac{v_i}{1 + v_i}$ . Since  $\mathbf{v} \in \mathcal{Y}$ , we have  $\sum_{m=1}^n v_m < 1$ . We first show that

$$(n; y_1; \dots; y_n) = \frac{v_n}{1 + \sum_{m=1}^n v_m} \prod_{j=1}^n \int_{x_j=0}^{y_j} \bar{G}(x_j) dx_j; \quad (3.57)$$

and

$$(0) = \frac{1}{1 + \sum_{m=1}^n v_m}; \quad (3.58)$$

We prove (3.57)-(3.58) as follows. Consider a system that has one PS server where jobs arrive according to a Poisson process with rate  $\Phi_n(\cdot)$  when there are  $n$  jobs and the job length distribution is  $G$ . Then from [22], there exists a unique stationary distribution

(single). Then  $\pi^{(single)}(n; y_1; \dots; y_n)$  is equal to right side of (3.57). From the stationary mean-field dynamics, we have that  $\pi$  is also another stationary distribution of the system. Due to the stability of the single server system, we must have that  $\pi = \pi^{(single)}$ . This implies that (3.57)-(3.58) are true.

Let  $\Gamma = (\Gamma_n; n \geq 0)$  where  $\Gamma_n = \lim_{b \rightarrow \infty} \pi(n; b; \dots; b)$  and  $\Gamma_0 = \pi(0)$ . Then from (3.57)-(3.58), we have

$$\Gamma_n = \frac{v_n}{1 + \sum_{j=1}^n v_j}$$

and

$$\Gamma_0 = \frac{1}{1 + \sum_{j=1}^{\infty} v_j}.$$

Further, we can write  $v_n$  in terms of  $\Gamma$  as  $v_n = \prod_{i=1}^n \frac{\Gamma_{i-1}}{\Gamma_i}$  where

$$\Psi_l(\Gamma) = \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k;l]}} \sum_{r=1}^i l_{fn_r=lg} q^{[k;l]}(\mathbf{n}; r) \prod_{j=1; j \neq r}^i \Gamma_{n_j}.$$

It can be seen that for  $n \geq 0$ , we have

$$\Psi_n(\Gamma) \Gamma_n = \Gamma_{n+1}. \quad (3.59)$$

Then from (3.18)-(3.19), for exponential JLDs with the same mean  $\frac{1}{\mu}$ , a fixed-point  $\mathbf{w} = (w_n; n \geq 0)$  of the mean-field satisfies that for  $n \geq 0$ ,

$$\Psi_n(\mathbf{w}) w_n = w_{n+1}. \quad (3.60)$$

Therefore, from (3.59) and (3.60), we have that  $\Gamma$  is a fixed-point of the mean-field under the assumption of exponential JLDs with average job length  $\frac{1}{\mu}$ . Furthermore, from (3.57)-(3.58), we obtain (3.21).

### 3.7.3 Proof of Theorem 3.1

The proof is divided into three steps. In the first step, we obtain the semigroup operator of the Markov process  $(\binom{N}{t}; t \geq 0)$ . The Feller property of the Markov process  $(\binom{N}{t}; t \geq 0)$  is established in the second step. We then compute the generator  $A^{(N)}(\cdot)$  and prove the martingale property of the process defined in equation (3.69) by using the Dynkin's lemma [56, Proposition 1.7, p.162] in the third step.

We now compute the semigroup operator  $T_h^{(N)}(\cdot)$  which is defined as

$$T_h^{(N)} f(\cdot) = \mathbb{E} \left[ f(\cdot) \prod_{j=0}^{(N)} = \right];$$

where the mapping  $f : M_F(\mathbb{U}) \rightarrow \mathbb{R}$  is a bounded continuous mapping. For  $\mathbf{u}^{(N)} = \cdot$ , in the interval  $[0; h]$ , let  $\mathbf{A}_h$  and  $\mathbf{D}_h$  be the number of arrivals and departures, respectively. Let us assume that measure  $\cdot$  contains mass at  $m$  distinct points denoted by  $\mathbf{u}^{(l)} = (n_l; u_1^{(l)}; \dots; u_{n_l}^{(l)})$ ,  $1 \leq l \leq m$ . Then the number of servers with state  $\mathbf{u}^{(l)}$  is equal to  $(\bar{r}\mathbf{u}^{(l)}g)$ . If a server lies in the state  $\mathbf{b} = (n; b_1; \dots; b_n)$  at time  $t$ , conditioned on the event that there are no arrivals in the interval  $[0; h]$  at the server, let  $\rho_{ND}(\mathbf{b}; h)$  be the probability that there are no departures at the server in the interval  $[0; h]$ . Then

$$\rho_{ND}(\mathbf{b}; h) = \prod_{i=1}^n \left\{ \frac{\bar{G}(b_i + \frac{h}{n})}{\bar{G}(b_i)} \right\};$$

From the definition of the hazard rate function, we have

$$\rho_{ND}(\mathbf{b}; h) = \prod_{i=1}^n \left\{ 1 - \frac{(b_i)h}{n} \right\} + o(h); \quad (3.61)$$

In the following result,  $B(\mathbf{u}^{(j)}; i; y; h)$  denotes the state of a server at  $t = h$  given that its state is  $\mathbf{u}^{(j)}$  at time  $t = 0$  and the job at the  $i^{\text{th}}$  position departs when its age is  $y$ . Similarly,  $C(\mathbf{u}^{(j)}; i; x; h)$  denotes the state of a server at time  $t = h$  given that its state is  $\mathbf{u}^{(j)}$  at time  $t = 0$  and a job arrives to this server at time  $t = x$  which joins the server at the  $i^{\text{th}}$  position.

**Lemma 3.3.** For  $f \in C_b(M_F(\mathbb{U}))$ , the semigroup operator  $T_h^{(N)}(\cdot)$  is given by

$$\begin{aligned} T_h^{(N)} f(\cdot) &= (1 - N h) \prod_{j=1; n_j > 0}^m (\rho_{ND}(\mathbf{u}^{(j)}; h))^{(\bar{r}\mathbf{u}^{(j)}g)} f(\cdot) \\ &+ (1 - N h) \sum_{j=1; n_j > 0}^m (\bar{r}\mathbf{u}^{(j)}g) \sum_{i=1}^{n_j} \int_{y=u_i^{(j)}}^{u_i^{(j)} + \frac{h}{n_j}} \left\{ \frac{g(y)}{\bar{G}(u_i^{(j)})} \right\} \\ &\quad \prod_{k=1; k \neq i}^{n_j} \left\{ \frac{\bar{G}\left(u_k^{(j)} + (y - u_i^{(j)}) + \frac{h(y - u_i^{(j)})n_j}{n_j - 1}\right)}{\bar{G}(u_k^{(j)})} \right\} \\ &\quad \left( (\rho_{ND}(\mathbf{u}^{(j)}; h))^{(\bar{r}\mathbf{u}^{(j)}g - 1)} \right) \left( \prod_{r=1; r \neq j; n_r > 0}^m (\rho_{ND}(\mathbf{u}^{(r)}; h))^{(\bar{r}\mathbf{u}^{(r)}g)} \right) \\ &\quad f\left(\cdot + (B(\mathbf{u}^{(j)}; i; y; h)) \left( n_j; u_1^{(j)} + \frac{h}{n_j}; \dots; u_{n_j}^{(j)} + \frac{h}{n_j} \right)\right) dy \end{aligned}$$

$$\begin{aligned}
& + N \int_{x=0}^h \frac{1}{h} \left( \sum_{j=1}^m \sum_{i=1}^{n_j+1} \left\{ \frac{1}{n_j+1} \right\} p_r(\tilde{x} : \tilde{x}^+ \mathbf{u}^{(j)}) \left[ f\left( \tilde{x} + f_{\mathbf{C}(\mathbf{u}^{(j)}; i; x; h)g} \quad (\tilde{x}^+ \mathbf{u}^{(j)}) \right) \right. \right. \\
& \quad \left. \left. \left( \prod_{q=1}^{n_j} \left\{ \frac{\bar{G}(u_q^{(j)} + \frac{x}{n_j} + \frac{h-x}{n_j+1})}{\bar{G}(u_q^{(j)})} \right\} \right) \bar{G}\left( \frac{h-x}{n_j+1} \right) (\rho_{ND}(\mathbf{u}^{(j)}; h))^{(\bar{\mathbf{u}}^{(j)}g)} \right. \right. \\
& \quad \left. \left. \prod_{k=1; k \neq j; n_k > 0}^m (\rho_{ND}(\mathbf{u}^{(k)}; h))^{(\bar{\mathbf{u}}^{(k)}g)} \right] \right) dx + (\cdot; h); \quad (3.62)
\end{aligned}$$

where

$$\begin{aligned}
B(\mathbf{u}^{(j)}; i; y; h) &= \left( n_j - 1; u_1^{(j)} + (y - u_i^{(j)}) + \frac{fh}{n_j - 1} \frac{n_j(y - u_i^{(j)})g}{n_j - 1}; \quad ; \right. \\
& \quad \left. u_{i-1}^{(j)} + (y - u_i^{(j)}) + \frac{fh}{n_j - 1} \frac{n_j(y - u_i^{(j)})g}{n_j - 1}; u_{i+1}^{(j)} + (y - u_i^{(j)}) + \frac{fh}{n_j - 1} \frac{n_j(y - u_i^{(j)})g}{n_j - 1}; \quad ; \right. \\
& \quad \left. u_{n_j}^{(j)} + (y - u_i^{(j)}) + \frac{fh}{n_j - 1} \frac{n_j(y - u_i^{(j)})g}{n_j - 1} \right); \quad (3.63)
\end{aligned}$$

$$\begin{aligned}
C(\mathbf{u}^{(j)}; i; x; h) &= \left( n_j + 1; u_1^{(j)} + \frac{x}{n_j} + \frac{h-x}{n_j+1}; \quad ; u_{i-1}^{(j)} + \frac{x}{n_j} + \frac{h-x}{n_j+1}; \right. \\
& \quad \left. \frac{h-x}{n_j+1}; u_i^{(j)} + \frac{x}{n_j} + \frac{h-x}{n_j+1}; \quad ; u_{n_j}^{(j)} + \frac{x}{n_j} + \frac{h-x}{n_j+1} \right) \quad (3.64)
\end{aligned}$$

and  $(\cdot; h)$  is  $o(h)$ .

*Proof.* We write  $T_h^{(N)} f(\cdot)$  as

$$\begin{aligned}
T_h^{(N)} f(\cdot) &= \mathbb{E} \left[ f\left( \frac{\cdot}{h} \right) I_{f\mathbf{A}_h=0g} I_{f\mathbf{D}_h=0g} j \binom{N}{0} = \right] + \mathbb{E} \left[ f\left( \frac{\cdot}{h} \right) I_{f\mathbf{A}_h=0g} I_{f\mathbf{D}_h=1g} j \binom{N}{0} = \right] \\
&+ \mathbb{E} \left[ f\left( \frac{\cdot}{h} \right) I_{f\mathbf{A}_h=1g} I_{f\mathbf{D}_h=0g} j \binom{N}{0} = \right] + \sum_{i=1}^m \sum_{j=1}^{n_i} \mathbb{E} \left[ f\left( \frac{\cdot}{h} \right) I_{f\mathbf{A}_h=ig} I_{f\mathbf{D}_h=jg} j \binom{N}{0} = \right]; \quad (3.65)
\end{aligned}$$

We first obtain the expression for  $\mathbb{E} \left[ f\left( \frac{\cdot}{h} \right) I_{f\mathbf{A}_h=0g} I_{f\mathbf{D}_h=0g} j \binom{N}{0} = \right]$  which is associated with the case when there are no arrivals and no departures in the interval  $[0; h]$ . In this situation, since the number of progressing jobs is constant,  $\binom{N}{0} = \tilde{x}$ . Hence,

$$\mathbb{E} \left[ f\left( \frac{\cdot}{h} \right) I_{f\mathbf{A}_h=0g} I_{f\mathbf{D}_h=0g} j \binom{N}{0} = \right] = \mathbb{P}(f\mathbf{A}_h = 0; \mathbf{D}_h = 0g j \binom{N}{0} = \cdot) f(\tilde{x});$$

As the arrival process is a Poisson process, we can write

$$\begin{aligned} \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=0g} I_{f\mathbf{D}_h=0gj} \binom{N}{0} = \right] \\ = \mathbb{P}(f\mathbf{A}_h = 0g) \mathbb{P}(f\mathbf{D}_h = 0gj \mid \mathbf{A}_h = 0; \binom{N}{0} = ) f(\tilde{\cdot}_h): \end{aligned}$$

Let us define

$$\mathbf{1}(\cdot; h) = (\mathbb{P}(f\mathbf{A}_h = 0g) \quad (1 \quad N \quad h)) \mathbb{P}(f\mathbf{D}_h = 0gj \mid \mathbf{A}_h = 0; \binom{N}{0} = ) f(\tilde{\cdot}_h):$$

Then

$$\begin{aligned} \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=0g} I_{f\mathbf{D}_h=0gj} \binom{N}{0} = \right] \\ = (1 \quad N \quad h) \mathbb{P}(f\mathbf{D}_h = 0gj \mid \mathbf{A}_h = 0; \binom{N}{0} = ) f(\tilde{\cdot}_h) + \mathbf{1}(\cdot; h): \end{aligned}$$

In the time interval  $[0; h]$ , conditioned on the event that there are no arrivals, the probability that there are no departures is equal to

$$\mathbb{P}(f\mathbf{D}_h = 0gj \mid \mathbf{A}_h = 0; \binom{N}{0} = ) = \prod_{i=1; n_i > 0}^m (p_{ND}(\mathbf{u}^{(i)}; h))^{f\mathbf{u}^{(i)}g}.$$

Since  $f \in C_b(M_F(\mathbb{U}))$  and  $(\mathbb{P}(f\mathbf{A}_h = 0g) \quad (1 \quad N \quad h)) = o(h)$  as the arrival process is Poisson with rate  $N$ , we conclude  $\mathbf{1}(\cdot; h)$  is  $o(h)$ .

Next, consider the second case where there is no arrival and one job departs in the interval  $[0; h]$ . Let us define

$$\mathbf{2}(\cdot; h) = (\mathbb{P}(f\mathbf{A}_h = 0g) \quad (1 \quad N \quad h)) \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{D}_h=1gj} \mathbf{A}_h = 0; \binom{N}{0} = \right]:$$

Then  $\mathbf{2}(\cdot; h)$  is  $o(h)$ . We can write

$$\begin{aligned} \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{A}_h=0g} I_{f\mathbf{D}_h=1gj} \binom{N}{0} = \right] = (1 \quad N \quad h) \mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{D}_h=1gj} \mathbf{A}_h = 0; \binom{N}{0} = \right] \\ + \mathbf{2}(\cdot; h): \end{aligned}$$

We next simplify the expression of  $\mathbb{E} \left[ f\left(\frac{\cdot}{h}\right) I_{f\mathbf{D}_h=1gj} \mathbf{A}_h = 0; \binom{N}{0} = \right]$ . Suppose that a job in the  $i^{\text{th}}$  position with age  $x$  departs at a server whose state was  $(n; b_1; \dots; b_n)$  at time  $t = 0$ . This implies that the departure has occurred at time  $(x - b_i)n$ . In this case, at time  $h$ , the server where a job departs will lie in the state  $(n - 1; b_1 + (x - b_i) + \frac{fh}{n-1} \frac{n(x - b_i)g}{n-1}; \dots; b_i - 1 + (x - b_i) + \frac{fh}{n-1} \frac{n(x - b_i)g}{n-1}; b_{i+1} + (x - b_i) + \frac{fh}{n-1} \frac{n(x - b_i)g}{n-1}; \dots; b_n + (x - b_i) + \frac{fh}{n-1} \frac{n(x - b_i)g}{n-1})$ . Using this, we can write the following equation by assuming that a job at the  $i^{\text{th}}$  position with age  $y$  departs from a server that had state  $\mathbf{u}^{(j)} = (n_j; u_1^{(j)}; \dots; u_{n_j}^{(j)})$  at time  $t = 0$ ,

$$\begin{aligned}
\mathbb{E} \left[ f \left( \frac{(\cdot)}{h} \right) I_{f\mathbf{A}_h=0g} I_{f\mathbf{D}_h=1g} j_0^{(N)} = \right] &= (1 - N h) \sum_{j=1; n_j > 0}^m (\bar{f}_{\mathbf{u}^{(j)}} g) \sum_{i=1}^{n_j} \int_{y=u_i^{(j)}}^{u_i^{(j)} + \frac{h}{n_j}} \left\{ \frac{g(y)}{\bar{G}(u_i^{(j)})} \right\} \\
&\quad \prod_{k=1; k \neq i}^{n_j} \left\{ \frac{\bar{G} \left( u_k^{(j)} + (y - u_i^{(j)}) + \frac{h(y - u_i^{(j)})n_j}{n_j - 1} \right)}{\bar{G}(u_k^{(j)})} \right\} \left( (\rho_{ND}(\mathbf{u}^{(j)}; h))^{(\bar{f}_{\mathbf{u}^{(j)}} g) - 1} \right) \\
&\quad \left( \prod_{l=1; l \neq j; n_l > 0}^m (\rho_{ND}(\mathbf{u}^{(l)}; h))^{(\bar{f}_{\mathbf{u}^{(l)}} g)} \right) f \left( \tilde{h} + (B(\mathbf{u}^{(j)}; i; y; h))_{n_j; u_1^{(j)} + \frac{h}{n_j}; \dots; u_{n_j}^{(j)} + \frac{h}{n_j}} \right) dy \\
&\quad + 2(\cdot; h):
\end{aligned}$$

Let us now consider the third case that corresponds to the event that there is no departure but an arrival occurs in the interval  $[0; h]$ . Let  $\mathfrak{z}_3(\cdot; h)$  be defined as

$$\mathfrak{z}_3(\cdot; h) = (\mathbb{P}(f\mathbf{A}_h = 1g) - N h) \mathbb{E} \left[ f \left( \frac{(\cdot)}{h} \right) I_{f\mathbf{D}_h=0g} j\mathbf{A}_h = 1; j_0^{(N)} = \right]:$$

Since  $\mathbb{P}(f\mathbf{A}_h = 1g) - N h = o(h)$  and  $f \in C_b(M_F(U))$ ,  $\mathfrak{z}_3(\cdot; h)$  is  $o(h)$ . We can write

$$\begin{aligned}
\mathbb{E} \left[ f \left( \frac{(\cdot)}{h} \right) I_{f\mathbf{A}_h=1g} I_{f\mathbf{D}_h=0g} j_0^{(N)} = \right] &= N h \mathbb{E} \left[ f \left( \frac{(\cdot)}{h} \right) I_{f\mathbf{D}_h=0g} j\mathbf{A}_h = 1; j_0^{(N)} = \right] \\
&\quad + \mathfrak{z}_3(\cdot; h); \quad (3.66)
\end{aligned}$$

We next obtain the expression for  $\mathbb{E} \left[ f \left( \frac{(\cdot)}{h} \right) I_{f\mathbf{D}_h=0g} j\mathbf{A}_h = 1; j_0^{(N)} = \right]$ . We recall that since the arrival process is Poisson, conditioned on the event that there is one arrival in  $[0; h]$ , the time at which the job arrives is uniformly distributed in  $[0; h]$  [77]. Let  $x$  be the time at which the job arrives and let  $i$  be its position at the destination server that had state  $(n; b_1; \dots; b_n)$  at time  $t = 0$ . Then the state of the destination server of the arrival at time  $h$  will be equal to  $(n+1; b_1 + \frac{x}{n} + \frac{h-x}{n+1}; \dots; b_{i-1} + \frac{x}{n} + \frac{h-x}{n+1}; \frac{h-x}{n+1}; b_i + \frac{x}{n} + \frac{h-x}{n+1}; \dots; b_n + \frac{x}{n} + \frac{h-x}{n+1})$ . Note that the state of the system at time  $x$  is equal to  $\tilde{x}$  and hence, the dispatcher uses  $\tilde{x}$  as the system state while implementing the routing policy. As a consequence, we get the following equation where we use the notation that a job arrives at time  $x$  and it is routed to a server that was lying in state  $\mathbf{u}^{(j)}$  at time  $t = 0$ , and the job joins at the  $i^{\text{th}}$  position. Then

$$\begin{aligned}
\mathbb{E} \left[ f \left( \frac{(\cdot)}{h} \right) I_{f\mathbf{D}_h=0g} j\mathbf{A}_h = 1; j_0^{(N)} = \right] &= \int_{x=0}^h \frac{1}{h} \left( \sum_{j=1}^m \sum_{i=1}^{n_j+1} \left\{ \frac{1}{n_j+1} \right\} \rho_r(\tilde{x} : \tilde{x}^+ \mathbf{u}^{(j)}) \right) \\
&\quad \left[ f \left( \tilde{h} + (f_{C(\mathbf{u}^{(j)}; i; x; h)g})_{(-h^+ \mathbf{u}^{(j)})} \right) \left( \prod_{q=1}^{n_j} \left\{ \frac{\bar{G}(u_q^{(j)} + \frac{x}{n_j} + \frac{h-x}{n_j+1})}{\bar{G}(u_q^{(j)})} \right\} \right) \bar{G} \left( \frac{h-x}{n_j+1} \right) \right]
\end{aligned}$$



$$\left( \rho_{ND}(\mathbf{u}^{(j)}; h) \right)^{(\bar{r}_{\mathbf{u}^{(j)}} g)^{-1}} \prod_{k=1; k \neq j; n_k > 0}^m \left( \rho_{ND}(\mathbf{u}^{(k)}; h) \right)^{(\bar{r}_{\mathbf{u}^{(k)}} g)} \Big] dx: \quad (3.67)$$

We now consider the case when there are multiple events in the interval  $[0; h]$ . We show that  $\mathbb{E} \left[ f \left( \frac{(\cdot)^{(N)}}{h} \right) I_{\mathbf{A}_h = 1; \mathbf{D}_h = 1} j_0^{(N)} = \cdot \right]$  is  $o(h)$  denoted by  $\phi_4(\cdot; h)$ . Since  $f$  is bounded, it is sufficient to show that  $\mathbb{P}(\mathbf{A}_h = 1; \mathbf{D}_h = 1 | j_0^{(N)} = \cdot)$  is  $o(h)$ . From the fact that  $\mathbb{P}(\mathbf{A}_h = 2)$  is  $o(h)$  due to the Poisson arrival process assumption, we get that  $\mathbb{P}(\mathbf{A}_h = 2; \mathbf{D}_h = 1 | j_0^{(N)} = \cdot)$  is  $o(h)$ . Also, we can write

$$\mathbb{P}(\mathbf{A}_h = 1; \mathbf{D}_h = 1 | j_0^{(N)} = \cdot) = \mathbb{P}(\mathbf{A}_h = 1) (1 - \mathbb{P}(\mathbf{D}_h = 0 | \mathbf{A}_h = 1; j_0^{(N)} = \cdot))$$

We show  $\lim_{h \downarrow 0} \mathbb{P}(\mathbf{D}_h = 0 | \mathbf{A}_h = 1; j_0^{(N)} = \cdot) = 1$ . From (3.67), we replace  $x$  with  $hz$ , then

$$\begin{aligned} \mathbb{E} \left[ f \left( \frac{(\cdot)^{(N)}}{h} \right) I_{\mathbf{D}_h = 0 | \mathbf{A}_h = 1; j_0^{(N)} = \cdot} \right] &= \int_{z=0}^1 \left( \sum_{j=1}^m \sum_{i=1}^{n_j+1} \left\{ \frac{1}{n_j+1} \right\} \rho_r(\tilde{h}z : \tilde{h}z \mathbf{u}^{(j)}) \right. \\ &\left. \left[ f \left( \tilde{h} + f_{C(\mathbf{u}^{(j)}; i; hz; h)g}(\tilde{h} \mathbf{u}^{(j)}) \right) \left( \prod_{q=1}^{n_j} \left\{ \frac{\bar{G}(u_q^{(j)} + \frac{hz}{n_j} + \frac{h \cdot hz}{n_j+1})}{\bar{G}(u_q^{(j)})} \right\} \right) \bar{G} \left( \frac{h \cdot hz}{n_j+1} \right) \right. \right. \\ &\left. \left. \left( \rho_{ND}(\mathbf{u}^{(j)}; h) \right)^{(\bar{r}_{\mathbf{u}^{(j)}} g)^{-1}} \prod_{k=1; k \neq j; n_k > 0}^m \left( \rho_{ND}(\mathbf{u}^{(k)}; h) \right)^{(\bar{r}_{\mathbf{u}^{(k)}} g)} \right] \right) dz: \quad (3.68) \end{aligned}$$

From (3.68), it can be seen that  $\lim_{h \downarrow 0} \mathbb{P}(\mathbf{D}_h = 0 | \mathbf{A}_h = 1; j_0^{(N)} = \cdot) = 1$ , and hence, we get  $\mathbb{P}(\mathbf{A}_h = 1; \mathbf{D}_h = 1 | j_0^{(N)} = \cdot) = o(h)$ .

Finally, by defining

$$\phi(\cdot; h) = \phi_1(\cdot; h) + \phi_2(\cdot; h) + \phi_3(\cdot; h) + \phi_4(\cdot; h);$$

we get (3.62). □

By the same arguments as in the proof of Theorem 2.8.1, we can show that the process  $(\frac{(\cdot)^{(N)}}{t}; t \geq 0)$  is a Feller-Dynkin process.

We now focus on establishing the third result. We first recall the definition of  $A^{(N)}(\cdot)$ . For any  $F \in C(M_F(U))$ , the generator  $A^{(N)}(\cdot)$  is defined as

$$A^{(N)}F(\cdot) = \lim_{h \downarrow 0} \frac{\mathbb{E} \left[ F \left( \frac{(\cdot)^{(N)}}{h} \right) j_0^{(N)} = \cdot \right] - F(\cdot)}{h};$$

where  $F \in C(M_F(U))$  is chosen such that the limit is well defined.

**Lemma 3.4.** For  $f \in C_b^1(U)$ , the process  $(\mathbf{M}_t^{(N)}(f); t \geq 0)$  defined as

$$\mathbf{M}_t^{(N)}(f) = h_t^{(N)}(f); \quad h_0^{(N)}(f) = \int_{x_1, \dots, x_n} A^{(N)} h_s^{(N)}(f); \quad ds \quad (3.69)$$

is a square integrable  $F_t^{(N)}$ -martingale and it is an RCLL process. Furthermore, for  $f \in C_b^1(U)$ , the quadratic variation of  $(\mathbf{M}_t^{(N)}(f); t \geq 0)$  is given by

$$\begin{aligned} \langle \mathbf{M}^{(N)}(f) \rangle_t = & \int_{s=0}^t \left( \sum_{n=1}^{\infty} \sum_{i=1}^n \int_{x_1, \dots, x_n} \left\{ \frac{\partial f(x_i)}{\partial x_i} \right\} \right. \\ & \left( (n-1; x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) - (n; x_1, \dots, x_n) \right)^2 d_s^{(N)}(n; x_1, \dots, x_n) \\ & + N \left[ \left\{ \frac{\partial^2 f}{\partial x_i^2} \right\} \Phi_0 \left( \frac{\partial f}{\partial x_i} \right) (1; 0) - (0) \right]^2 \\ & + \sum_{n=1}^{\infty} \sum_{i=1}^{n+1} \int_{x_1, \dots, x_n} \left\{ \frac{1}{N(n+1)} \right\} \Phi_n \left( \frac{\partial f}{\partial x_i} \right) \\ & \left( (n+1; x_1, \dots, x_{i-1}, 0, x_i, \dots, x_n) - (n; x_1, \dots, x_n) \right)^2 \\ & \left. d_s^{(N)}(n; x_1, \dots, x_n) \right] ds: \quad (3.70) \end{aligned}$$

*Proof.* We first obtain the expression for the generator  $A^{(N)}(f)$ . From (3.62) and since the set of linear combinations of  $Q_f$  for  $f \in C_b^1(U)$  where  $Q_f(x) = e^{-h \cdot f(x)}$  is dense in the set  $C(M_F(U))$  [78, proposition 7.10], from the expression of  $A^{(N)}Q_f(x)$ , we get

$$\begin{aligned} A^{(N)}F(x) = & \lim_{h \downarrow 0} \left( \frac{F(x+h) - F(x)}{h} - N F(x) \right) \\ & \sum_{n=1}^{\infty} \sum_{i=1}^n \int_{x_1, \dots, x_n} \left\{ \frac{\partial f(x_i)}{\partial x_i} \right\} F(x) d(n; x_1, \dots, x_n) + \sum_{n=1}^{\infty} \sum_{i=1}^n \int_{x_1, \dots, x_n} \left\{ \frac{\partial^2 f(x_i)}{\partial x_i^2} \right\} \\ & \left( F(x + (n-1; x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)) - F(x) \right) d(n; x_1, \dots, x_n) \\ & + N \left[ \left\{ \frac{\partial^2 f}{\partial x_i^2} \right\} \Phi_0 \left( \frac{\partial f}{\partial x_i} \right) F(x + (1; 0)) - (0) \right. \\ & \left. + \sum_{n=1}^{\infty} \sum_{i=1}^{n+1} \int_{x_1, \dots, x_n} \left\{ \frac{1}{N(n+1)} \right\} \Phi_n \left( \frac{\partial f}{\partial x_i} \right) \right. \\ & \left. \left( F(x + (n+1; x_1, \dots, x_{i-1}, 0, x_i, \dots, x_n)) - F(x) \right) d(n; x_1, \dots, x_n) \right]: \end{aligned}$$

For  $f \in C_b^1(U)$ ,  $F \in M_F(U)$ , then  $A^{(N)}h_t^{(N)}(f)$  is well defined. From the Dynkin's formula [56, Proposition 1.7, p.162],  $(\mathbf{M}_t^{(N)}(f); t \geq 0)$  defined as

$$\mathbf{M}_t^{(N)}(f) = h_t^{(N)}(f); \quad h_0^{(N)}(f) = \int_{s=0}^t A^{(N)} h_s^{(N)}(f); \quad ds \quad (3.71)$$

is an RCLL  $F_t^{(N)}$ -local martingale. It can be checked that after simplifications, we get

$$\begin{aligned}
\mathbf{M}_t^{(N)}(\cdot) &= h_t^{(N)}; i - h_0^{(N)}; i - \int_{s=0}^t h_s^{(N)}; i ds - \int_{s=0}^t \left( \sum_{n=1}^1 \sum_{i=1}^n \int_{x_1} \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} \right. \\
&\quad \left. ( (n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) - (n; x_1; \dots; x_n) \right) d_s^{(N)}(n; x_1; \dots; x_n) \\
&\quad + N \left[ \left( \left\{ \frac{s^{(N)}(f)g}{N} \right\} \Phi_0 \left( \frac{s^{(N)}}{N} \right) ( (1; 0) - (0) ) \right) \right. \\
&\quad \left. + \sum_{n=1}^1 \sum_{i=1}^{n+1} \int_{x_1} \int_{x_n} \left\{ \frac{1}{N(n+1)} \right\} \Phi_n \left( \frac{s^{(N)}}{N} \right) \right. \\
&\quad \left. ( (n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) - (n; x_1; \dots; x_n) \right) d_s^{(N)}(n; x_1; \dots; x_n) \Big] ds;
\end{aligned} \tag{3.72}$$

where

$$h_s^{(N)}; i = \sum_{n=1}^1 \sum_{i=1}^n \int_{x_1} \int_{x_n} \frac{1}{n} \frac{\partial (n; x_1; \dots; x_n)}{\partial x_i} d_s^{(N)}(n; x_1; \dots; x_n); \tag{3.73}$$

By choosing  $F(\cdot) = h_t^{(N)}; i$ , from [79, Theorem 7.15], we have

$$\langle \mathbf{M}^{(N)}(\cdot) \rangle_{t=0} = \int_{s=0}^t (A^{(N)} F^2(\cdot) - 2F(\cdot) A^{(N)} F(\cdot)) ds; \tag{3.74}$$

We get (3.70) after simplifications. In the next section, we show in (3.78) that  $\sup_m \Phi_m(\cdot) < 1$ . Using this result, and from the fact that  $\mathcal{C}_b^1(\mathbb{U}) \subset \mathcal{C}_b(\mathbb{R}_+)$ , we have  $\mathbb{E} \left[ \langle \mathbf{M}^{(N)}(\cdot) \rangle_t \right] < 1$ . Hence,  $(\mathbf{M}_t^{(N)}(\cdot); t \geq 0)$  is a square integrable martingale.  $\square$

### 3.7.4 Proof of Theorem 3.2

We first show that there exists a unique mean-field solution. This step is required to prove the convergence of  $(\frac{t}{N}; t \geq 0)$  as  $N \rightarrow \infty$ . We then prove the convergence of  $(\frac{t}{N}; t \geq 0)$  as  $N \rightarrow \infty$  to the unique mean-field solution.

#### Existence and Uniqueness of a Mean-field Solution

From (3.16), for  $\varphi \in \mathcal{C}_b(\mathbb{U})$ , the operator  $\mathcal{H}_t^-; i$  is a linear operator and  $\mathcal{H}_t^-(\mathbb{U}) = 1$ . As a consequence, from the Riesz-Markov-Kakutani theorem [69, 80] given in Theorem A.4, for

$\mu_t \in M_1(U)$ , the uniqueness of the probability measure  $\mu_t$  is equivalent to the uniqueness of the operator  $\mathcal{V} h_t; i$ . Also, since  $C_b(U)$  is a separating class of  $M_1(U)$  [56, p. 111], if two measures  $\mu_1; \mu_2 \in M_1(U)$  satisfies  $h_{\mu_1}; i = h_{\mu_2}; i$  for all  $i \in C_b(U)$ , then  $\mu_1 = \mu_2$ . Therefore it is enough to show the uniqueness of the operator  $\mathcal{V} h_t; i$ .

For given  $\mu_0$ , we now show that there exists at most one real valued process  $h(\mu_t; i; t; 0)$  that satisfies the MFEs. Let  $(\mu_t^1; t = 0); (\mu_t^2; t = 0)$  be the two mean-field solutions with initial points  $\mu_0^1; \mu_0^2$ , respectively. Then for  $i \in C_b(U)$ , we get

$$\begin{aligned}
h_t^1 - \mu_t^2; i = h_0^1 - \mu_0^2; i + \int_{s=0}^t & \left( \sum_{n=1}^1 \sum_{i=1}^n \int_{x_1} \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} \right. \\
& \left. (\mu_{t-s}^1(n; x_1; \dots; x_i; \dots; x_{i-1}; x_{i+1}; \dots; x_n) - \mu_{t-s}^2(n; x_1; \dots; x_n)) \right. \\
& \left. d(\mu_s^{-1} - \mu_s^{-2})(n; x_1; \dots; x_n) \right) ds \\
& + \int_{s=0}^t \left( \left[ (\mu_s^{-1}(f_0 g) \Phi_0(\mu_s^{-1}) (\mu_{t-s}^1(1; 0) - \mu_{t-s}^2(0))) \right. \right. \\
& + \sum_{n=1}^1 \sum_{i=1}^{n+1} \int_{x_1} \int_{x_n} \frac{1}{(n+1)} \Phi_n(\mu_s^{-1}) (\mu_{t-s}^1(n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) \\
& \left. \left. - \mu_{t-s}^2(n; x_1; \dots; x_n)) d\mu_s^{-1}(n; x_1; \dots; x_n) \right] \right. \\
& \left. \left[ (\mu_s^{-2}(f_0 g) \Phi_0(\mu_s^{-2}) (\mu_{t-s}^2(1; 0) - \mu_{t-s}^2(0))) \right. \right. \\
& \left. \left. \sum_{n=1}^1 \sum_{i=1}^{n+1} \int_{x_1} \int_{x_n} \left\{ \frac{1}{(n+1)} \right\} \Phi_n(\mu_s^{-2}) (\mu_{t-s}^2(n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) \right. \right. \\
& \left. \left. - \mu_{t-s}^2(n; x_1; \dots; x_n)) d\mu_s^{-2}(n; x_1; \dots; x_n) \right] \right) ds: \quad (3.75)
\end{aligned}$$

Our objective is to obtain a result of the form

$$k_t^{-1} - \mu_t^2 k = b + c \int_{s=0}^t k_s^{-1} - \mu_s^2 k ds \quad (3.76)$$

for some  $b; c > 0$ ,  $t \in [0; T]$ . Then from Gronwall's inequality stated in Theorem A.7 [56], we would get

$$k_t^{-1} - \mu_t^2 k \leq b e^{ct} \quad (3.77)$$

for  $t \in [0; T]$ . In this direction, by using the first term on the right side of equation (3.75), we can write

$$|h_0^{-1} - \mu_0^2; i| \leq k_0^{-1} - \mu_0^2 k$$

We now simplify the second term on the right side of equation (3.75). For this, we define a function  $w_{t,s}$  as follows:

for  $n > 0$ ,

$$w_{t,s}(n; x_1, \dots, x_n) = \sum_{i=1}^n \left\{ \frac{(x_i)}{n} \right\} (\tilde{t}_s (n-1; x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) - \tilde{t}_s (n; x_1, \dots, x_n))$$

and for  $n = 0$ ,

$$w_{t,s}(0) = 0:$$

Then  $w_{t,s} \in C_b(U)$  since  $\tilde{t}_s \in C_b(U)$  and  $\tilde{t}_s \in C_b(\mathbb{R}_+)$ . Furthermore, we have

$$|kw_{t,s}k| \leq 2k \|\tilde{t}_s\|:$$

It can be seen that the second term can be written as  $\int_{s=0}^t h_s^{-1} \|\tilde{t}_s\| w_{t,s} ds$ . This implies that we can write

$$\left| \int_{s=0}^t h_s^{-1} \|\tilde{t}_s\| w_{t,s} ds \right| \leq \int_{s=0}^t 2k \|\tilde{t}_s\| k ds:$$

Let us now consider the third term on the right side of (3.75) that corresponds to the case of arrivals. Let  $h_{t,s}$  be the function defined as

$$h_{t,s}(n; x_1, \dots, x_n) = \sum_{i=1}^{n+1} \left\{ \frac{1}{n+1} \right\} \Phi_n(\tilde{t}_s (n+1; x_1, \dots, x_{i-1}, 0, x_i, \dots, x_n) - \tilde{t}_s (n; x_1, \dots, x_n))$$

for all  $(n; x_1, \dots, x_n)$ . Then the third term is equal to  $\int_{s=0}^t (h_s^{-1} h_{t,s}^{-1} i - h_s^{-2} h_{t,s}^{-2} i) ds$ . We can bound this term using the fact that

$$|h_s^{-1} h_{t,s}^{-1} i - h_s^{-2} h_{t,s}^{-2} i| = |h_s^{-1} \|\tilde{t}_s\| h_{t,s}^{-1} i| + |h_s^{-2} h_{t,s}^{-1} - h_{t,s}^{-2} i|:$$

Also, we can write

$$|h_s^{-1} h_{t,s}^{-1} i - h_s^{-2} h_{t,s}^{-2} i| \leq k_s^{-1} \|\tilde{t}_s\| k h_{t,s}^{-1} k + k_s^{-2} k h_{t,s}^{-1} - h_{t,s}^{-2} k:$$

We have that  $k_s^{-2} k = 1$  as  $\|\tilde{t}_s\|$  is a probability measure.

We now obtain an upper bound on  $k h_{t,s}^{-1} k$ . It can be seen that we have

$$k h_{t,s}^{-1} k \leq 2k \sup_{n \geq 0} j \Phi_n(\tilde{t}_s):$$

We now show that  $\sup_n \Phi_n(\cdot)$  is bounded by a finite value. For  $\mathbf{n} \in M_1(U)$ ,

$$\Phi_m(\mathbf{n}) = \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k;i]}} \sum_{j=1}^i I_{f\mathbf{n}_j = mg} q^{[k;i]}(\mathbf{n}; j) \prod_{r=1; r \neq j}^i Q_{n_r}(\cdot):$$

Let  $\mathcal{Y}_m^{[k;i]}$  be defined as

$$\mathcal{Y}_m^{[k;i]} = \{\mathbf{n} \in \mathcal{Y}^{[k;i]} : \exists j \text{ such that } n_j = mg\}$$

Then we can write

$$\Phi_m(\mathbf{n}) = \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}_m^{[k;i]}} \sum_{j=1}^i I_{f\mathbf{n}_j = mg} q^{[k;i]}(\mathbf{n}; j) \prod_{r=1; r \neq j}^i Q_{n_r}(\cdot):$$

We can bound the above term as

$$\Phi_m(\mathbf{n}) \leq \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}_m^{[k;i]}} \sum_{j=1}^i I_{f\mathbf{n}_j = mg} \prod_{r=1; r \neq j}^i Q_{n_r}(\cdot):$$

Out of  $i$  sampled servers upto stage  $k$ , let  $b$  be the number of servers with  $m$  jobs. Then we can write

$$\begin{aligned} \Phi_m(\mathbf{n}) &\leq \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{b=1}^i b \binom{i}{b} (Q_m(\cdot))^{b-1} (1 - Q_m(\cdot))^{i-b}; \\ &\leq \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{b=1}^i b \binom{i}{b}; \\ &\leq M (d^{(max)})^2 \sum_{b=1}^{d^{(max)}} \binom{d^{(max)}}{b}; \\ &\leq M (d^{(max)})^3 (d^{(max)}!); \end{aligned} \tag{3.78}$$

Therefore, we have

$$kh_{t;s}^{-1} \leq 2Mk (d^{(max)})^3 (d^{(max)}!):$$

Similarly, we can write

$$kh_{t;s}^{-1} \leq h_{t;s}^{-2} \leq 2k \sup_m |\Phi_m(-1_s) - \Phi_m(-2_s)|:$$

We have

$$|\Phi_m(-1_s) - \Phi_m(-2_s)|$$

$$= \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}_m^{[k;l]}} \sum_{j=1}^i I_{\bar{r}n_j = mg} q^{[k;l]}(\mathbf{n}; j) \left( \prod_{r=1; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) \prod_{r=1; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \right) \right| :$$

We can write

$$\begin{aligned} & \left| \Phi_m \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) \Phi_m \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \right| \\ & \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}_m^{[k;l]}} \sum_{j=1}^i I_{\bar{r}n_j = mg} q^{[k;l]}(\mathbf{n}; j) \left( \prod_{r=1; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) Q_{n_1} \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \prod_{r=2; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) \right) \right| \\ & + \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}_m^{[k;l]}} \sum_{j=1}^i I_{\bar{r}n_j = mg} q^{[k;l]}(\mathbf{n}; j) \left( Q_{n_1} \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \prod_{r=2; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) \prod_{r=1; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \right) \right| : \end{aligned} \quad (3.79)$$

We next bound the first term on the right side of (3.79). For each  $j$ , we have

$$\begin{aligned} & \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}_m^{[k;l]}} I_{\bar{r}n_j = mg} q^{[k;l]}(\mathbf{n}; j) \left( \prod_{r=1; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) Q_{n_1} \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \prod_{r=2; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) \right) \right| \\ & M d^{(max)} \sum_{m=0}^{\infty} \left| Q_m \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) Q_m \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \right| : \end{aligned} \quad (3.80)$$

Let us define a function  $r_{\left( \begin{smallmatrix} -1 \\ s, -2 \\ s \end{smallmatrix} \right)} \in \mathcal{C}_b(\mathbb{U})$  such that for  $\mathbf{u} \in \mathbb{U}_n$ ,

$$r_{\left( \begin{smallmatrix} -1 \\ s, -2 \\ s \end{smallmatrix} \right)}(\mathbf{u}) = \begin{cases} 1 & \text{if } Q_n \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) Q_n \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \\ 1 & \text{otherwise:} \end{cases}$$

Then  $\sum_{m=0}^{\infty} j Q_m \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) Q_m \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) j = h_s^{-1} \sum_{s; r_{\left( \begin{smallmatrix} -1 \\ s, -2 \\ s \end{smallmatrix} \right)}} i$ . Hence,

$$\sum_{m=0}^{\infty} \left| Q_m \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) Q_m \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \right| \leq k_s^{-1} \sum_{s; r_{\left( \begin{smallmatrix} -1 \\ s, -2 \\ s \end{smallmatrix} \right)}} k:$$

Therefore

$$\begin{aligned} & \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}_m^{[k;l]}} \sum_{j=1}^i I_{\bar{r}n_j = mg} q^{[k;l]}(\mathbf{n}; j) \left( \prod_{r=1; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) Q_{n_1} \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \prod_{r=2; r \neq j}^i Q_{n_r} \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) \right) \right| \\ & M (d^{(max)})^2 k_s^{-1} \sum_{s; r_{\left( \begin{smallmatrix} -1 \\ s, -2 \\ s \end{smallmatrix} \right)}} k: \end{aligned}$$

By similar arguments, the second term on the right side of (3.79) is bounded by  $M (d^{(max)})^2 (d^{(max)} - 1) k_s^{-1} \sum_{s; r_{\left( \begin{smallmatrix} -1 \\ s, -2 \\ s \end{smallmatrix} \right)}} k$ . Hence,

$$\left| \Phi_m \left( \begin{smallmatrix} -1 \\ s \end{smallmatrix} \right) \Phi_m \left( \begin{smallmatrix} -2 \\ s \end{smallmatrix} \right) \right| \leq M (d^{(max)})^3 k_s^{-1} \sum_{s; r_{\left( \begin{smallmatrix} -1 \\ s, -2 \\ s \end{smallmatrix} \right)}} k:$$

Therefore we obtain that

$$kh_{t,s}^{-1} \quad h_{t,s}^{-2}k \quad 2Mk \quad k(d^{(max)})^3k_s^{-1} \quad -^2_s k:$$

Finally, by using the above bounds, we get

$$|h_t^{-1} \quad -^2_t; \quad i| \quad \left( k_0^{-1} \quad -^2_0 k + \int_{s=0}^t 2k \quad k k_s^{-1} \quad -^2_s k \, ds \right. \\ \left. + 4M(d^{(max)})^3(d^{(max)}!) \int_{s=0}^t k_s^{-1} \quad -^2_s k \, ds \right) k \quad k:$$

Hence

$$k_t^{-1} \quad -^2_t k \quad k_0^{-1} \quad -^2_0 k + (2k \quad k + 4M(d^{(max)})^3(d^{(max)}!) \quad ) \int_{s=0}^t k_s^{-1} \quad -^2_s k \, ds:$$

From equation (3.76), we get

$$k_t^{-1} \quad -^2_t k \quad k_0^{-1} \quad -^2_0 k e^{(2k \quad k + 4M(d^{(max)})^3(d^{(max)}!) \quad )t}.$$

This implies that for given initial measure  $\bar{\nu}_0$ , there exists at most one mean-field solution.

The existence of a mean-field solution follows from the proof of the relative compactness of the sequence  $f_t^{(N)}; t \in [0; 1)$  in  $D_{M_1(U)}([0; 1))$ . We show that every limit point of the sequence  $f_t^{(N)}; t \in [0; 1)$  has sample paths that are  $a:s$ : solutions of equation (3.16). This implies that there exists a solution to the MFEs.

### Convergence of $(f_t^{(N)}; t \in [0; 1))$

We now study the convergence of  $(f_t^{(N)}; t \in [0; 1))$  in  $D_{M_1(U)}([0; 1))$  as  $N \rightarrow \infty$ . Let  $(\bar{F}_t^{(N)}; t \in [0; 1))$  be the right continuous filtration associated with the process  $(f_t^{(N)}; t \in [0; 1))$ . The main logic behind the proof is that we first show the relative compactness of the sequence  $f_t^{(N)}; t \in [0; 1))g_{N-1}$  and then we show every limit point  $(f_t; t \in [0; 1))$  has sample paths evolving almost surely according to the MFEs with the initial point  $\bar{\nu}_0$ . From the uniqueness of the mean-field solution for a given initial measure, all the limiting points have almost surely identical sample paths coinciding with the unique mean-field solution with the initial point  $\bar{\nu}_0$ . This implies that the sequence  $f_t^{(N)}; t \in [0; 1))g_{N-1}$  converges in distribution to the unique mean-field solution with the initial point  $\bar{\nu}_0$ , and this mean-field solution is referred to as the mean-field limit denoted by  $(f_t; t \in [0; 1))$ .



Let  $\bar{A}^{(N)}(\cdot)$  be the generator of the Markov process  $(\bar{X}_t^{(N)}; t \geq 0)$ . Then the following process  $(\bar{M}_t^{(N)}(\cdot); t \geq 0)$  for  $\bar{X} \in C_b^1(U)$  is an RCLL square integrable  $\bar{F}_t^{(N)}$ -martingale,

$$\bar{M}_t^{(N)}(\cdot) = h_t^{(N)}(\cdot) - h_0^{(N)}(\cdot) - \int_{s=0}^t \bar{A}^{(N)} h_s^{(N)}(\cdot) ds \quad (3.81)$$

By using the expression of  $\bar{A}^{(N)}(\cdot)$ , we get

$$\begin{aligned} \bar{M}_t^{(N)}(\cdot) &= h_t^{(N)}(\cdot) - h_0^{(N)}(\cdot) - \int_{s=0}^t h_s^{(N)}(\cdot) \mathbb{P} ds - \int_{s=0}^t \left( \sum_{n=1}^1 \sum_{i=1}^n \int_{x_1} \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} \right. \\ &\quad \left. ( (n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) - (n; x_1; \dots; x_n) \right) d_s^{-(N)}(n; x_1; \dots; x_n) \\ &+ \left[ \left( -\frac{(N)}{s} (f \circ g) \Phi_0 \left( -\frac{(N)}{s} \right) ( (1; 0) \dots (0) \right) + \sum_{n=1}^1 \sum_{i=1}^{n+1} \int_{x_1} \int_{x_n} \left\{ \frac{1}{n+1} \right\} \Phi_n \left( -\frac{(N)}{s} \right) \right. \\ &\quad \left. ( (n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) - (n; x_1; \dots; x_n) \right) d_s^{-(N)}(n; x_1; \dots; x_n) \left. \right] ds \quad (3.82) \end{aligned}$$

For  $\bar{X} \in C_b^1(U)$ , we get

$$\begin{aligned} \langle \bar{M}^{(N)}(\cdot) \rangle_t &= \frac{1}{N} \left[ \int_{s=0}^t \left( \sum_{n=1}^1 \sum_{i=1}^n \int_{x_1} \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} \right. \right. \\ &\quad \left. \left. ( (n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) - (n; x_1; \dots; x_n) \right)^2 d_s^{-(N)}(n; x_1; \dots; x_n) \right. \\ &+ \left[ \left( -\frac{(N)}{s} (f \circ g) \Phi_0 \left( -\frac{(N)}{s} \right) ( (1; 0) \dots (0) \right)^2 + \sum_{n=1}^1 \sum_{i=1}^{n+1} \int_{x_1} \int_{x_n} \left\{ \frac{1}{n+1} \right\} \Phi_n \left( -\frac{(N)}{s} \right) \right. \\ &\quad \left. \left. ( (n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) - (n; x_1; \dots; x_n) \right)^2 d_s^{-(N)}(n; x_1; \dots; x_n) \right] ds \right] \quad (3.83) \end{aligned}$$

From the Prohorov's theorem [57], since the space  $D_{M_1(U)}([0; T])$  is complete and separable, proving the relative compactness of  $(f(\bar{X}_t^{(N)}; t \geq 0))_{N \geq 1}$  is equivalent to showing the tightness of this sequence. To show the tightness, we now recall the Jakubowski's criteria (From Theorem 4.6 of [74]) in Appendix A.3, which give the necessary and sufficient conditions J1 and J2. We first show that the condition J2 is satisfied. For this, we prove that the two sufficient conditions C1 and C2 stated in Appendix A.3 are satisfied. For any  $T > 0$ ,  $t \in [0; T]$ , we have

$$h_t^{(N)}(\cdot) - k_1 h_t^{(N)}(\cdot) \leq k_2$$

and since  $h_t^{-(N)}; \mathbf{1}_i = 1$ , with  $b = k - k_1$ , the condition C1 is satisfied.

We now focus on the proof of the condition C2. From (3.83) and Doob's inequality given in Theorem A.8, for  $\epsilon > 0$ , we get

$$\begin{aligned} \mathbb{P} \left( \sup_{t \in [0, T]} \left| \overline{\mathbf{M}}_t^{(N)}(\cdot) \right| > \frac{\epsilon}{2} \right) &\leq \frac{4}{\epsilon} \mathbb{E} \left[ \langle \overline{\mathbf{M}}^{(N)}(\cdot) \rangle_T \right] \\ &\leq \frac{4T}{\epsilon} k - k^2 \frac{1}{N} (k - k_1 + d^{(max)2} (d^{(max)}!)) \leq \epsilon \end{aligned}$$

as  $N \rightarrow \infty$ . Therefore the sequence  $f(\overline{\mathbf{M}}_t^{(N)}(\cdot); t = 0) g_{N-1}$  converges in distribution to the null process from the convergence criterion in  $D_R([0; T])$  [56, Theorem 1.4, p.339]. Furthermore, since the sequence  $f(\overline{\mathbf{M}}_t^{(N)}(\cdot); t = 0) g_{N-1}$  is tight in  $D_R([0; T])$ , there exists  $\epsilon_1 > 0$  and  $N_1 > 0$  such that for  $N \geq N_1$ , we have

$$\mathbb{P} \left( \sup_{u, v \in [0, T]; |u - v| \leq \epsilon_1} \left| \overline{\mathbf{M}}_v^{(N)}(\cdot) - \overline{\mathbf{M}}_u^{(N)}(\cdot) \right| > \frac{\epsilon}{2} \right) \leq \epsilon \quad (3.84)$$

Also, for any  $u$  and  $v$  with  $u < v \leq T$ , from (3.82), we get

$$\begin{aligned} \left| h_v^{-(N)}; \mathbf{1}_i - h_u^{-(N)}; \mathbf{1}_i \right| &= \int_{s=u}^v \left| h_s^{-(N)}; \mathbf{1}_i \right| ds + 2k - k_1 k - k_j u - v_j \\ &\quad + 2Mk - k - (d^{(max)})^3 (d^{(max)}!) j u - v_j + \left| \overline{\mathbf{M}}_v^{(N)}(\cdot) - \overline{\mathbf{M}}_u^{(N)}(\cdot) \right|. \end{aligned}$$

The above equation can be simplified as

$$\begin{aligned} \left| h_v^{-(N)}; \mathbf{1}_i - h_u^{-(N)}; \mathbf{1}_i \right| &\leq j v - u j k - k_1 (1 + 2Mk - k + 2(d^{(max)})^3 (d^{(max)}!)) \\ &\quad + \left| \overline{\mathbf{M}}_v^{(N)}(\cdot) - \overline{\mathbf{M}}_u^{(N)}(\cdot) \right|. \quad (3.85) \end{aligned}$$

By using (3.84) and (3.85), we can find  $\epsilon_2 > 0$  and  $N_2 > 0$  such that for  $N \geq N_2$ , we have

$$\mathbb{P} \left( \sup_{u, v \in [0, T]; |u - v| \leq \epsilon_2} \left| h_v^{-(N)}; \mathbf{1}_i - h_u^{-(N)}; \mathbf{1}_i \right| > \epsilon \right) \leq \epsilon$$

This completes the proof of the condition J2.

We now focus on the proof of the condition J1. Let  $(\mathbf{n}_i(t); \mathbf{V}_{i1}(t); \dots; \mathbf{V}_{im_i(t)}(t))$  be the state of the  $i^{\text{th}}$  server where  $\mathbf{V}_{ij}(t)$  is the age of the  $j^{\text{th}}$  progressing job. Then

$$h_t^{-(N)}; \mathbf{1}_i = \frac{1}{N} \sum_{i=1; \mathbf{n}_i(t) > 0}^N (\mathbf{V}_{i1}(t) + \dots + \mathbf{V}_{im_i(t)}(t));$$

Let us consider a random variable  $\mathbf{H}_t$  to denote the age of a progressing job at time  $t$ . Further, let  $\mathbf{X}$  be a random variable with the same distribution as the job length distribution  $G$ . Then for any  $b > 0$ , we have

$$P(\mathbf{H}_t > b) = P(\mathbf{X} > b):$$

Also, we have

$$h_t^{-(N)}; \Xi_i = \sum_{n=0}^{\infty} \int_{x_1}^{\infty} \int_{x_n}^{\infty} n d_t^{-(N)}(n; x_1; \dots; x_n):$$

At any time  $t$ , we can classify the progressing jobs into two classes. The first class of jobs are the ones which are in service starting from the time  $t = 0$ , and the number of such jobs is upper bounded by its initial value say  $\mathbf{W}$ . If a job belonging to this class has age  $a$  at time  $t = 0$ , then its age is at most  $a + t^\theta$  at time  $t = t^\theta$ . The second class of jobs are the ones which entered the system after the time  $t = 0$ , and the number of such jobs is bounded by the total number of arrivals up to time  $t$  denoted by  $\mathbf{E}^{(N)}(t)$ . Hence, we can write

$$P((h_t^{-(N)}; \Xi_i + h_t^{-(N)}; i) > b) \\ = P\left(\left[h_0^{-(N)}; \Xi_i + \frac{\mathbf{E}^{(N)}(t)}{N} + h_0^{-(N)}; i + t h_0^{-(N)}; \Xi_i + \frac{\sum_{j=1}^{\mathbf{E}^{(N)}(t)} \mathbf{X}_j}{N}\right] > b\right)$$

where  $(\mathbf{X}_j; 1 \leq j \leq \mathbf{E}^{(N)}(t))$  are i.i.d. random variables with the distribution  $G$ . We now consider the convergence of  $\frac{\sum_{j=1}^{\mathbf{E}^{(N)}(t)} \mathbf{X}_j}{N}$  in distribution sense. Since the job arrival process is a Poisson process with rate  $N$ , it is the same as the sum of  $N$  independent Poisson processes with rate  $1$ . Hence,

$$\mathbf{E}^{(N)}(t) = \sum_{i=1}^N \mathbf{E}_{(i)}(t):$$

where  $\mathbf{E}_{(i)}(t)$  ( $1 \leq i \leq N$ ) is the number of arrivals from the  $i^{\text{th}}$  Poisson process whose rate is  $1$ . Therefore

$$\frac{\sum_{j=1}^{\mathbf{E}^{(N)}(t)} \mathbf{X}_j}{N} = \frac{1}{N} \sum_{i=1}^N \left( \sum_{k=1}^{\mathbf{E}_{(i)}(t)} \mathbf{X}_{ik} \right);$$

where  $f \mathbf{X}_{ik} g_{i-1; k-1}$  are i.i.d. random variables with distribution  $G$ . From the law of large numbers, we get

$$\left( \frac{\sum_{j=1}^{\mathbf{E}^{(N)}(t)} \mathbf{X}_j}{N} \right) \xrightarrow{P} \frac{t}{N};$$

and

$$\left( \frac{\mathbf{E}^{(N)}(t)}{N} \right) \xrightarrow{P} t:$$

From Assumption 3.4.1, we get

$$h_0^{-(N)}; i + (t+1)h_0^{-(N)}; \Xi i + \frac{\sum_{j=1}^{\mathbf{E}^{(N)}(t)} \mathbf{X}_j}{N} + \frac{\mathbf{E}^{(N)}(t)}{N} \Big) h_0^-; i + (t+1)h_0^-; \Xi i + t + \frac{t}{N}.$$

Also, using Assumption 3.4.1, we can find some constant  $M_0$  such that

$$\liminf_{N \uparrow} \mathbb{P}(\max(h_0^{-(N)}; \Xi i; h_0^{-(N)}; i) < M_0) > 1 \quad :$$

Let  $M_T = M_0(2 + T) + 2 \sqrt{T} + \frac{2T}{N}$ , then

$$\liminf_{N \uparrow} \mathbb{P}(\sup_{t \in [0; T]} (h_t^{-(N)}; \Xi i + h_t^{-(N)}; i) < M_T) > 1 \quad : \quad (3.86)$$

For given  $0 < \epsilon < 1$ , let

$$\mathcal{W}_T; \epsilon, f \in \mathcal{M}_1(\mathbb{U}) : (h_t; \Xi i + h_t; i) < M_T g:$$

For an integer  $n$ , let

$$\mathcal{C}_k^{(n)} = f(k; y_1; \dots; y_k) : y_i \in \mathbb{R}_+ \text{ and } y_i \leq n \text{ for all } 1 \leq i \leq k;$$

and let  $\mathcal{C}^{(n)} = f(0)g \cap (\bigcup_{k=1}^n \mathcal{C}_k^{(n)})$  and let  $\overline{\mathcal{C}^{(n)}}$  be the complement of  $\mathcal{C}^{(n)}$ , then

$$\overline{\mathcal{C}^{(n)}} \subset \frac{2M_T}{n};$$

and hence,

$$\lim_{n \uparrow} \sup_{\mathcal{W}_T} \overline{\mathcal{C}^{(n)}} = 0:$$

From Lemma A7.5 of [75], we have that  $\mathcal{W}_T; \epsilon$  is relatively compact in  $\mathcal{M}_1(\mathbb{U})$ . Further, from equation (3.86), we have

$$\liminf_{N \uparrow} \mathbb{P}(\overline{\mathcal{C}^{(N)}} \subset \mathcal{W}_T; \epsilon; \forall t \in [0; T]) > 1 \quad :$$

Let  $\mathcal{K}_T; \epsilon$  be the closure of  $\mathcal{W}_T; \epsilon$ , then  $\mathcal{K}_T; \epsilon$  is a compact set with

$$\liminf_{N \uparrow} \mathbb{P}(\overline{\mathcal{C}^{(N)}} \subset \mathcal{K}_T; \epsilon; \forall t \in [0; T]) = 1 \quad :$$

This proves the condition J1. Hence, we have proved that the sequence of processes  $f(\overline{\mathcal{C}^{(N)}}; t = 0)g_{N-1}$  is relatively compact.

Let  $(\overline{\mathcal{C}^{(N)}}; t = 0)$  be the limiting point of a converging subsequence of  $f(\overline{\mathcal{C}^{(N)}}; t = 0)g_{N-1}$ . Then  $\overline{\mathcal{C}^{(N)}}_0$  almost surely coincides with  $\overline{\mathcal{C}^{(N)}}_0$  from Assumption 3.4.1. The continuous mapping theorem implies that

$$\begin{aligned}
h_{t; i} = h_{0; i} + \int_{s=0}^t h_{s; i} \Phi_i ds &+ \int_{s=0}^t \left( \sum_{n=1}^1 \sum_{i=1}^n \int_{x_1} \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} \right. \\
& \left. ( (n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) (n; x_1; \dots; x_n) ) d_s(n; x_1; \dots; x_n) \right. \\
& + \left[ ( (f_0 g) \Phi_0(s) ( (1; 0) (0) ) ) + \sum_{n=1}^1 \sum_{i=1}^{n+1} \int_{x_1} \int_{x_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(s) \right. \\
& \left. \left. ( (n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) (n; x_1; \dots; x_n) ) d_s(n; x_1; \dots; x_n) \right] \right) ds: \quad (3.87)
\end{aligned}$$

Finally, we prove that the sample paths of  $(h_{t; i}; t \geq 0)$  coincide almost surely with the unique mean-field solution with the initial point  $\bar{h}_0$ . Since  $C_b^1(U)$  is a separating class of  $M_1(U)$ , we have that  $(h_{t; i}; t \geq 0) \in C_{M_1(U)}([0; 1])$ . We then prove that for  $\bar{h}_0$ , any process  $(h_{t; i}; t \geq 0) \in C_{M_1(U)}([0; 1])$  is a solution to equation (3.87) iff it is a solution to the mean-field equation (3.16). We give a proof of this in Section 3.7.5. Since there exists a unique solution to the mean-field equation for the given initial point  $\bar{h}_0$ , from Assumption 3.4.1, all the limiting points have almost surely identical sample paths coinciding with the unique mean-field solution. Therefore,  $f_t^{(-M)}; t \geq 0) g_{N-1}$  converges in distribution to the unique mean-field solution denoted by  $(\bar{h}_{t; i}; t \geq 0)$ .

### 3.7.5 Evolution of $(h_{t; i}; t \geq 0)$ for $(h_{t; i}; t \geq 0) \in C_b(U)$

In this section, we prove that  $(h_{t; i}; t \geq 0)$  is a solution to equation (3.87) if and only if it is a solution to equation (3.16).

We first show that any process  $(h_{t; i}; t \geq 0)$  that satisfies (3.87) also satisfies (3.16). For  $(h_{t; i}; t \geq 0) \in C_b^1(U)$ , if the integrand in equation (3.87) is a continuous function of  $S$ , any real valued process  $(h_{t; i}; t \geq 0)$  that satisfies (3.87) also satisfies the following differential equation

$$\begin{aligned}
\frac{dh_{t; i}}{dt} = h_{t; i} \Phi_i + \left( \sum_{n=1}^1 \sum_{i=1}^n \int_{x_1} \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} \right. \\
& \left. ( (n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) (n; x_1; \dots; x_n) ) d_t(n; x_1; \dots; x_n) \right. \\
& + \left[ ( (f_0 g) \Phi_0(t) ( (1; 0) (0) ) ) + \sum_{n=1}^1 \sum_{i=1}^{n+1} \int_{x_1} \int_{x_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(t) \right. \\
& \left. \left. ( (n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) (n; x_1; \dots; x_n) ) d_t(n; x_1; \dots; x_n) \right] \right): \quad (3.88)
\end{aligned}$$

Hence, we need to prove that the two terms on the right side of (3.88) are continuous in  $t$ . Since  $\mathcal{C}_b^1(\mathbb{U})$ , the first term  $h_{t; \Phi}^i$  is a continuous function of  $t$ . Let  $\tilde{\phantom{h}}$  be the function defined as

$$\tilde{\phantom{h}}(0) = 0$$

and for  $n \geq 1$

$$\tilde{\phantom{h}}(n; x_1; \dots; x_n) = \sum_{i=1}^n \left\{ \frac{(x_i)}{n} \right\} ((n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) - (n; x_1; \dots; x_n)):$$

Then the second term on the right side of (3.88) that corresponds to departures can be written as

$$h_{t; \tilde{\phantom{h}}}^i = \sum_{n=1}^{\infty} \sum_{i=1}^n \int_{x_1} \dots \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} ((n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) - (n; x_1; \dots; x_n)) d_t(n; x_1; \dots; x_n):$$

It can be seen that  $\tilde{\phantom{h}} \in \mathcal{C}_b(\mathbb{U})$  and hence,  $h_{t; \tilde{\phantom{h}}}^i$  is a continuous function of  $t$ .

The term related to arrivals can be written as below. Let  $\Phi_{(t)}$  be the function defined as

$$\begin{aligned} \Phi_{(t)}(n; x_1; \dots; x_n) \\ = \left\{ \frac{1}{n+1} \right\} \Phi_n(t) \sum_{i=1}^{n+1} ((n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) - (n; x_1; \dots; x_n)): \end{aligned}$$

Then

$$h_{t; \Phi_{(t)}}^i = \left[ (\Phi_{(t)}(0) - \Phi_{(t)}(1; 0; \dots; 0)) + \sum_{n=1}^{\infty} \sum_{i=1}^{n+1} \int_{x_1} \dots \int_{x_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(t) ((n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) - (n; x_1; \dots; x_n)) d_t(n; x_1; \dots; x_n): \right]$$

Since  $\Phi_{(t)} \in \mathcal{C}_b(\mathbb{U})$ ,  $t \mapsto h_{t; \Phi_{(t)}}^i$  is continuous when  $b$  is a fixed-value. The mapping  $t \mapsto h_{t; \Phi_{(t)}}^i$  is continuous if  $h_{t+b; \Phi_{(t+b)}}^i \rightarrow h_{t; \Phi_{(t)}}^i$  as  $b \rightarrow 0$ . We can write

$$|h_{t+b; \Phi_{(t+b)}}^i - h_{t; \Phi_{(t)}}^i| = |h_{t+b; \Phi_{(t+b)}}^i - h_{t+b; \Phi_{(t)}}^i| + |h_{t+b; \Phi_{(t)}}^i - h_{t; \Phi_{(t)}}^i|: \quad (3.89)$$

Since  $\Phi_{(t)} \in \mathcal{C}_b(\mathbb{U})$ , we get

$$\lim_{b \rightarrow 0} |h_{t+b; \Phi_{(t+b)}}^i - h_{t; \Phi_{(t)}}^i| = 0: \quad (3.90)$$

To prove the continuity of the mapping  $t \mapsto h_{t; (\cdot)}$ , it remains to be shown that

$$\lim_{b \downarrow 0} |h_{t+b; (\cdot)} - h_{t; (\cdot)}| = 0:$$

For an integer  $L > 0$ , let

$$V^{(L)} = \{f(n; x_1, \dots, x_n) \geq \frac{1}{L} \mid x_i \leq L \text{ for all } 1 \leq i \leq n \text{ where } 1 \leq n \leq L\}:$$

Let  $\bar{V}^{(L)}$  be the complement of  $V^{(L)}$ . Then for given  $\epsilon > 0$ , there exists  $L > 0$  such that

$$h_{t; \bar{V}^{(L)}} < \epsilon: \quad (3.91)$$

Since the mapping  $t \mapsto h_t$  is continuous, we can find some  $r_1 > 0$  such that for all  $b \geq [\min(t; r_1); r_1]$ ,

$$h_{t+b; \bar{V}^{(L)}} < \epsilon: \quad (3.92)$$

We now show that there exists  $r_2 \geq (0; r_1)$  such that for  $b \geq [\min(t; r_2); r_2]$  and  $\mathbf{u} \in V^{(L)}$ , the following result holds

$$|h_{t+b}(\mathbf{u}) - h_t(\mathbf{u})| < 6 M(d^{(max)})^3 k: \quad (3.93)$$

For  $\mathbf{u} \in V^{(L)}$ ,

$$|h_{t+b}(\mathbf{u}) - h_t(\mathbf{u})| < 2 \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k; i]}} j \Phi_m(t+b) - \Phi_m(t):$$

Let us now obtain bounds on  $\sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k; i]}} j \Phi_m(t+b) - \Phi_m(t)$ . We have

$$j \Phi_m(t+b) - \Phi_m(t) = \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k; i]}} \sum_{j=1}^i I_{f_{\mathbf{n}} = mg} q^{[k; i]}(\mathbf{n}; j) \prod_{r=1; r \neq j}^i (Q_{n_r}(t+b) - Q_{n_r}(t)) \right|:$$

For each  $j$ , we can write

$$\begin{aligned} & \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k; i]}} I_{f_{\mathbf{n}} = mg} q^{[k; i]}(\mathbf{n}; j) \prod_{r=1; r \neq j}^i (Q_{n_r}(t+b) - Q_{n_r}(t)) \right| \\ & \quad \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k; i]}} I_{f_{\mathbf{n}} = mg} q^{[k; i]}(\mathbf{n}; j) (Q_{n_1}(t+b) - Q_{n_1}(t)) \prod_{r=2; r \neq j}^i Q_{n_r}(t+b) \right| \\ & + \left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k; i]}} I_{f_{\mathbf{n}} = mg} q^{[k; i]}(\mathbf{n}; j) Q_{n_1}(t) \left( \prod_{r=2; r \neq j}^i (Q_{n_r}(t+b) - Q_{n_r}(t)) \right) \right|: \quad (3.94) \end{aligned}$$

We can bound the first term on the right side of (3.94) as

$$\left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k;l]}} I_{f_{n_j} = mg} q^{[k;l]}(\mathbf{n}; j) (Q_{n_1}(t+b) - Q_{n_1}(t)) \prod_{r=2; r \neq j}^i Q_{n_r}(t+b) \right|$$

$$M d^{(max)} \sum_{i=0}^j Q_i(t+b) - Q_i(t) j :$$

By similar arguments, we obtain

$$\left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k;l]}} I_{f_{n_j} = mg} q^{[k;l]}(\mathbf{n}; j) Q_{n_1}(t) \left( \prod_{r=2; r \neq j}^i (Q_{n_r}(t+b) - Q_{n_r}(t)) \right) \right|$$

$$M d^{(max)} (d^{(max)} - 1) \sum_{i=0}^j Q_i(t+b) - Q_i(t) j :$$

Hence, for each  $j$ , we have

$$\left| \sum_{k=1}^M \sum_{i=1}^{d^{(max)}} \sum_{\mathbf{n} \in \mathcal{Y}^{[k;l]}} I_{f_{n_j} = mg} q^{[k;l]}(\mathbf{n}; j) \prod_{r=1; r \neq j}^i (Q_{n_r}(t+b) - Q_{n_r}(t)) \right|$$

$$M (d^{(max)})^2 \sum_{i=0}^j Q_i(t+b) - Q_i(t) j :$$

As a result, we have

$$j \Phi_m(t) - \Phi_m(t+b) j \leq M (d^{(max)})^3 \sum_{i=0}^j Q_i(t+b) - Q_i(t) j :$$

We now show

$$\lim_{b \downarrow 0} \sum_{i=0}^j Q_i(t+b) - Q_i(t) j = 0 :$$

We can write

$$\sum_{i=0}^j Q_i(t+b) - Q_i(t) j = \sum_{i=L-1}^j Q_i(t+b) - Q_i(t) j + \sum_{i=L}^j Q_i(t+b) - Q_i(t) j$$

$$= \sum_{i=L-1}^j Q_i(t+b) - Q_i(t) j + \sum_{i=L}^j Q_i(t+b) j + \sum_{i=L}^j Q_i(t) j :$$

From (3.91)-(3.92), for  $b \geq [\min(t; r_1); r_1]$ , we have

$$\sum_{i=0}^j Q_i(t+b) - Q_i(t) j = \sum_{i=L-1}^j Q_i(t+b) - Q_i(t) j + \dots$$



Again, as  $\lim_{b \downarrow 0} jQ_i(t+b) - Q_i(t)j = 0$  for  $0 \leq i \leq L-1$ , there exists  $r_2 \geq (0; r_1)$  such that for  $b \geq [\min(t; r_2); r_2]$ ,

$$\sum_{i=0}^{L-1} jQ_i(t+b) - Q_i(t)j \leq 3$$

Hence, we have

$$\sum_{i=0}^{L-1} jQ_i(t+b) - Q_i(t)j \leq 3$$

for  $b \geq [\min(t; r_2); r_2]$ . Therefore,

$$\max_{m \in \mathbb{N}} j\Phi_m(t) - \Phi_m(t+b)j \leq 3M(d^{(max)})^3 :$$

From (3.91)-(3.93) and for  $b \geq [\min(t; r_2); r_2]$ , since  $(t+b) \in C_b(U)$ , we can write

$$\begin{aligned} |h_{t+b} - (t)j| &\leq 6M(d^{(max)})^3 k k h_{t+b} |f_{V(L)}g^j + 2(d^{(max)})^3 (d^{(max)})! k k \\ &\leq 8M(d^{(max)})^3 (d^{(max)})! k k : \end{aligned}$$

By taking  $b \downarrow 0$  and then  $t \downarrow 0$  in equation (3.89), the mapping  $t \mapsto h_{t; (t)j}$  is continuous.

We now show that any solution to (3.88) also satisfies an alternative equation. For this, we apply the method of change of variables. For  $r \in C_b^1(U)$  and  $r \geq t$ , let  $\tilde{\cdot}$  be defined as

$$\tilde{\cdot}(n; x_1; \dots; x_n) = \begin{cases} (0) & \text{if } n = 0; \\ \tilde{h}_{t;r}(n; x_1; \dots; x_n) & \text{otherwise;} \end{cases} \quad (3.95)$$

Let us consider the change in  $h_{r; \tilde{\cdot}i}$  due to change in 'r'. We can write

$$\frac{dh_{r; \tilde{\cdot}i}}{dr} = \lim_{h \downarrow 0} \frac{f h_{r+h; \tilde{h}_{t;r}i} - h_{r; \tilde{h}_{t;r}i} g}{h} + \lim_{h \downarrow 0} \frac{f h_{r+h; \tilde{h}_{t;r}h}i - h_{r+h; \tilde{h}_{t;r}i} g}{h} : \quad (3.96)$$

In the above equation, the first term on the right side can be computed by using (3.88) and the second term is equal to  $h_{t; \tilde{\cdot}i}$ . Hence, we obtain

$$\begin{aligned} \frac{dh_{r; \tilde{\cdot}i}}{dr} &= \left( \sum_{n=1}^L \sum_{i=1}^n \int_{x_1} \dots \int_{x_n} \left\{ \frac{(x_i)}{n} \right\} \right. \\ &\quad \left. \left( \tilde{h}_{r-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n} - \tilde{h}_{r; x_1; \dots; x_n} \right) d_r(n; x_1; \dots; x_n) \right. \\ &\quad \left. + \left[ \left( r(f \circ g) \Phi_0(r) \left( \tilde{h}_{(1;0)} - \tilde{h}_{(0)} \right) \right) + \sum_{n=1}^L \sum_{i=1}^{n+1} \int_{x_1} \dots \int_{x_n} \left\{ \frac{1}{n+1} \right\} \Phi_n(r) \right] \right) \end{aligned}$$

$$\left( \tilde{f}(n+1; x_1, \dots, x_{i-1}, 0, x_i, \dots, x_n) - \tilde{f}(n; x_1, \dots, x_n) \right) d_r(n; x_1, \dots, x_n) \Big] \Big) : \quad (3.97)$$

By integrating  $\frac{dh_{r;i}}{dr}$  with respect to  $r$  from 0 to  $t$ , we obtain (3.16).

For  $\tilde{f} \in C_b^1(U)$ , we now show that any solution  $(h_{t;i}; t \geq 0)$  to (3.16) also satisfies (3.87). For this, it is enough to show that the derivative of  $h_{t;i}$  with respect to  $t$  exists. Since  $\tilde{f} \in C_b^1(U)$ , the dominated convergence theorem [69, Theorem 1.34] implies that  $\frac{dh_{0;i}}{dt}$  is well defined. We can apply the Leibniz integral rule to prove that the derivative of the second term on the right side of equation (3.16) with respect to  $t$  exists. According to this rule, the integrand should be continuous with respect to both the variables  $r$  and  $t$ . This proof is similar to the proof of the integrand in equation (3.87). We then need to show that the derivative of the integrand with respect to  $t$  exists, and further, the differential should be continuous with respect to both the variables  $r$  and  $t$ . Since  $\tilde{f} \in C_b^1(U)$ , the dominated convergence theorem implies that the integrand is differentiable and further, it is continuous with respect to  $r$  and  $t$ . This proof is similar to that of the continuity of the integrand in (3.87). Hence, any solution  $(h_{t;i}; t \geq 0) \in C_{M_1(U)}([0; 1])$  is a solution to (3.87) if and only if it is a solution to (3.16). In equation (3.16),  $\tilde{f}$  need not be differentiable.

### 3.7.6 Proof of Lemma 3.2

We prove this result based on induction arguments. Let  $S(k; n)$  be the set of all  $k$  ordered numbers chosen from  $\{1; 2; \dots; n\}$ . For  $y_i \in \mathbb{R}_+$ ,  $i = 1, \dots, k$ ,  $\mathbf{b} = (b_1; \dots; b_k) \in S(k; n)$  and  $\mathbf{a} = (a_1; \dots; a_k)$ , let

$$\begin{aligned} B^{[k; \mathbf{a}; \mathbf{b}]}(y_1; \dots; y_k) \\ = \tilde{f}(n; x_1; \dots; x_n) : x_i \in \mathbb{R}_+ \text{ for } i \in \{b_1; \dots; b_k\} \text{ and } y_j = x_{b_j} - y_{j+1} \quad 1 \leq j < k; \end{aligned}$$

Further, we define

$$-_t(B^{[k; \mathbf{a}; \mathbf{b}]}) = \max_{\mathbf{b} \in S(k; n)} \sup_{(y_1; \dots; y_k) \in \mathbb{R}_+^k} -_t(B^{[k; \mathbf{a}; \mathbf{b}]}(y_1; \dots; y_k)) : \quad (3.98)$$

We now show that for each  $k, 1 \leq j < k$ ,

$$\lim_{j \rightarrow 0} \sum_{n=k} -_t(B^{[k; \mathbf{a}; \mathbf{b}]}) = 0 : \quad (3.99)$$

From (3.99), we obtain that the measure  $\bar{\nu}_t$  is absolute continuous w.r.t. the Lebesgue measure and hence, it has a density function. Here,  $\sum_n k \bar{\nu}_t(B^{[k; :n]})$  is well-defined since  $\bar{\nu}_t(B^{[k; :n]}) = \bar{\nu}_t(U_n)$  and  $\sum_n k \bar{\nu}_t(U_n) = 1$ .

Since there existence sequence of bounded continuous functions that decrease pointwise to indicators of closed sets, by using monotone convergence theorem, we can write from (3.16) as

$$\begin{aligned}
\bar{\nu}_t(B^{[k; :n; \mathbf{b}]}) (y_1; \dots; y_k) &= h_{0; \sim}^t I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g^i \\
&\int_{s=0}^t \left( \sum_{i=1}^n \int_{x_1} \int_{x_n} \frac{(x_i)}{n} \bar{\nu}_t^s I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g(n; x_1; \dots; x_n) d\bar{\nu}_s(n; x_1; \dots; x_n) \right. \\
&\quad \left. + \sum_{j=1}^{n+1} \int_{x_1} \int_{x_{n+1}} \frac{(x_j)}{n+1} \right. \\
&\quad \left. \bar{\nu}_t^s I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g(n; x_1; \dots; x_{j-1}; x_{j+1}; \dots; x_{n+1}) d\bar{\nu}_s(n+1; x_1; \dots; x_{n+1}) \right. \\
&\quad \left. \int_{x_1} \int_{x_n} \Phi_n(\bar{\nu}_s) \bar{\nu}_t^s I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g(n; x_1; \dots; x_n) d\bar{\nu}_s(n; x_1; \dots; x_n) \right. \\
&\quad \left. + \sum_{i=1}^n \int_{x_1} \int_{x_{i-1}} \int_{x_{i+1}} \int_{x_n} \frac{\Phi_{n-1}(\bar{\nu}_s)}{n} \right. \\
&\quad \left. \bar{\nu}_t^s I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g(n; x_1; \dots; x_{i-1}; 0; x_{i+1}; \dots; x_n) \right. \\
&\quad \left. d\bar{\nu}_s(n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) \right) ds: \quad (3.100)
\end{aligned}$$

We now obtain bounds on every term of the right side of (3.100). By using the fact that

$$\bar{\nu}_t^s I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g(n; x_1; \dots; x_n) = I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g \left( n; x_1 + \frac{t}{n}; \dots; x_n + \frac{t}{n} \right); \quad (3.101)$$

we obtain the following bounds

$$h_{0; \sim}^t I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g^i = h_{0; \sim}^t (B^{[k; :n]}); \quad (3.102)$$

$$\left\{ \sum_{i=1}^n \int_{x_1} \int_{x_n} \frac{(x_i)}{n} \bar{\nu}_t^s I_{f_{B^{[k; :n; \mathbf{b}]}}(y_1; \dots; y_k)} g(n; x_1; \dots; x_n) d\bar{\nu}_s(n; x_1; \dots; x_n) \right\} = k \bar{\nu}_s(B^{[k; :n]}); \quad (3.103)$$

$$\left\{ \sum_{j=1}^{n+1} \int_{x_1} \int_{x_{n+1}} \frac{(x_j)}{n+1} \right.$$

$$\left. \tilde{t} s \int_{f_{B^{[k]; [n; b]}(y_1; \dots; y_k)}} g(n; x_1; \dots; x_{j-1}; x_{j+1}; \dots; x_{n+1}) d_{-s}^-(n+1; x_1; \dots; x_{n+1}) \right\} \\ k \leq k_s(B^{[k]; [n+1]}); \quad (3.104)$$

and from (3.78) we obtain

$$\int_{x_1} \int_{x_n} \Phi_n(-s) \tilde{t} s \int_{f_{B^{[k]; [n; b]}(y_1; \dots; y_k)}} g(n; x_1; \dots; x_n) d_{-s}^-(n; x_1; \dots; x_n) \\ M(d^{(max)})^3 (d^{(max)}!)_{-s}^-(B^{[k]; [n]}); \quad (3.105)$$

We now obtain bounds on the last term. We have

$$\tilde{t} s \int_{f_{B^{[k]; [n; b]}(y_1; \dots; y_k)}} g(n; x_1; \dots; x_{i-1}; 0; x_{i+1}; \dots; x_n) \\ = \int_{f_{B^{[k]; [n; b]}(y_1; \dots; y_k)}} g\left(n; x_1 + \frac{t-s}{n}; \dots; x_{i-1} + \frac{t-s}{n}; \frac{t-s}{n}; x_{i+1} + \frac{t-s}{n}; \dots; x_n + \frac{t-s}{n}\right); \quad (3.106)$$

We now investigate the range of values of  $x_i$ ,  $i = 1$ , for which we get

$$\tilde{t} s \int_{f_{B^{[k]; [n; b]}(y_1; \dots; y_k)}} g(n; x_1; \dots; x_{i-1}; 0; x_{i+1}; \dots; x_n) = 1:$$

There are two possibilities: either  $i \geq fb_1; \dots; b_k g$  or  $i \leq fb_1; \dots; b_k g$ . We begin with the first case. Let us assume that  $i = b_r$  for some  $1 \leq r \leq k$ . Then we have

$$y_r = \frac{t-s}{n} = y_r + r$$

As a result, we conclude  $t - (ny_r + n - r) = s - t - ny_r$ . Also, for  $j = b_m$  with  $j \neq b_r$ , we obtain that  $x_j$  satisfies  $y_m = x_{b_m} + \frac{t-s}{n} = y_m + m$ . For  $j \leq fb_1; \dots; b_k g$ , we have  $x_j + \frac{t-s}{n} \geq R_+$ .

In the second case when  $i \leq fb_1; \dots; b_k g$ , we have  $\frac{t-s}{n} \geq R_+$  implying that  $s = t$ . Also, for  $j = b_m$  with  $1 \leq m \leq k$ , we obtain that  $x_j$  satisfies  $y_m = x_{b_m} + \frac{t-s}{n} = y_m + m$  and for  $j \leq fb_1; \dots; b_k g$ , we have  $x_j + \frac{t-s}{n} \geq R_+$ .

From the above discussion, we can show that

$$\left\{ \int_{s=0}^t \sum_{i=1}^n \left( \int_{x_1} \int_{x_{i-1}} \int_{x_{i+1}} \int_{x_n} \frac{\Phi_{n-1}(-s)}{n} \right. \right. \\ \left. \left. \tilde{t} s \int_{f_{B^{[k]; [n; b]}(y_1; \dots; y_k)}} g(n; x_1; \dots; x_{i-1}; 0; x_{i+1}; \dots; x_n) \right. \right. \\ \left. \left. d_{-s}^-(n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) \right) ds \right\}$$

$$M(d^{(max)})^3(d^{(max)})! \left( \sum_{l=1}^k \int_{s=0}^t \frac{1}{n} I_{\mathbb{F}[t-n(y_l+i);t-ny_l]g}(s)^-_{s} (B^{[k-1; l; n-1]}) ds + \frac{n-k}{n} \int_{s=0}^t -_s (B^{[k; ;n-1]}) ds \right); \quad (3.107)$$

where  $l = (1; ; l-1; l+1; ; k)$ . We next show the following result by using (3.102)-(3.107) for the case  $k = 1$  and  $l = (1)$ ,

$$\lim_{j \rightarrow 0} \sum_{n=1}^j -_t (B^{[1; ;n]}) = 0; \quad (3.108)$$

From (3.102)-(3.107), we obtain

$$\begin{aligned} -_t (B^{[1; ;n]}(y_1)) &= -_0 (B^{[1; ;n]}) + k \int_{s=0}^t -_s (B^{[1; ;n]}) ds + k \int_{s=0}^t -_s (B^{[1; ;n+1]}) ds \\ &\quad + M(d^{(max)})^3(d^{(max)})! \int_{s=0}^t -_s (B^{[1; ;n]}) ds \\ &\quad + M(d^{(max)})^3(d^{(max)})! \left( \sum_{l=1}^1 \int_{s=0}^t \frac{1}{n} I_{\mathbb{F}[t-n(y_l+i);t-ny_l]g}(s)^-_{s} (U_{n-1}) ds + \frac{n-1}{n} \int_{s=0}^t -_s (B^{[1; ;n-1]}) ds \right); \quad (3.109) \end{aligned}$$

We can also write

$$\int_{s=0}^t \frac{1}{n} I_{\mathbb{F}[t-n(y_l+i);t-ny_l]g}(s)^-_{s} (U_{n-1}) ds = \int_{s=0}^t \frac{1}{n} I_{\mathbb{F}[t-(y_l+i);t-y_l]g}(s)^-_{s} (U_{n-1}) ds; \quad (3.110)$$

From (3.109)-(3.110), we get

$$\begin{aligned} -_t (B^{[1; ;n]}) &= -_0 (B^{[1; ;n]}) + k \int_{s=0}^t (-_s (B^{[1; ;n]}) + -_s (B^{[1; ;n+1]})) ds \\ &\quad + M(d^{(max)})^3(d^{(max)})! \left( \int_{s=0}^t -_s (B^{[1; ;n]}) ds + \left( \sup_{y_1} \sum_{l=1}^1 \int_{s=0}^t I_{\mathbb{F}[t-(y_l+i);t-y_l]g}(s)^-_{s} (U_{n-1}) ds + \int_{s=0}^t -_s (B^{[1; ;n-1]}) ds \right) \right); \quad (3.111) \end{aligned}$$

From (3.111), by summing over  $n \geq 1$ , we obtain

$$\begin{aligned} \left( \sum_{n=1}^{\infty} -_t (B^{[1; ;n]}) \right) &= \left( \sum_{n=1}^{\infty} -_0 (B^{[1; ;n]}) \right) + 2k \int_{s=0}^t \left( \sum_{n=1}^{\infty} -_s (B^{[1; ;n]}) \right) ds \\ &\quad + 2 M(d^{(max)})^3(d^{(max)})! \left( \int_{s=0}^t \left( \sum_{n=1}^{\infty} -_s (B^{[1; ;n]}) \right) ds + \frac{1}{2} \right); \quad (3.112) \end{aligned}$$

By using Gronwall's inequality, we have

$$\left( \sum_{n-1}^{-}_t(B^{[1; :n]}) \right) \leq \left( \left( \sum_{n-1}^{-}_0(B^{[1; :n]}) \right) + 2 M(d^{(max)})^3(d^{(max)!}) \right) e^{(2k-k+2) M(d^{(max)})^3(d^{(max)!})t}. \quad (3.113)$$

Since the measure  $\bar{\nu}_0$  has a density function, we have

$$\lim_{r \downarrow 0} \left( \sum_{n-1}^{-}_t(B^{[1; :n]}) \right) = 0. \quad (3.114)$$

We now use induction arguments. Suppose, for  $\nu = [ \nu_1; \dots; \nu_m ]$ , we have

$$\lim_{r \downarrow 0} \left( \sum_{n-m}^{-}_t(B^{[m; :n]}) \right) = 0 \text{ and } \left( \sum_{n-m}^{-}_t(B^{[m; :n]}) \right) \leq \hat{z}_m e^{\hat{\nu}_m t}; \quad (3.115)$$

for all  $1 \leq r \leq m$  and  $\hat{z}_m; \hat{\nu}_m \geq 0$ . Then we will show that for all  $1 \leq r \leq m+1$ , there exists  $\hat{z}_{m+1}; \hat{\nu}_{m+1} \geq 0$  such that

$$\lim_{r \downarrow 0} \left( \sum_{n-m+1}^{-}_t(B^{[m+1; :n]}) \right) = 0 \text{ and } \left( \sum_{n-m+1}^{-}_t(B^{[m+1; :n]}) \right) \leq \hat{z}_{m+1} e^{\hat{\nu}_{m+1} t}. \quad (3.116)$$

From (3.102)-(3.107), we obtain

$$\begin{aligned} & \int_{n-m+1}^{-}_t(B^{[m+1; :n]})g \leq \int_{n-m+1}^{-}_0(B^{[m+1; :n]})g + 2k-k \int_{s=0}^t \int_{n-m+1}^{-}_s(B^{[m+1; :n]})g ds \\ & \quad + 2 M(d^{(max)})^3(d^{(max)!}) \int_{s=0}^t \int_{n-m+1}^{-}_s(B^{[m+1; :n]})g ds \\ & \quad + M(d^{(max)})^3(d^{(max)!}) \sup_{y_1; \dots; y_{m+1}} \left( \sum_{l=1}^{m+1} \int_{s=0}^t I_{\bar{\nu}_l}(y_l; \cdot; t-y_l)g(s) \left\{ \sum_{n-m+1}^{-}_s(B^{[m; :l; n-1]}) \right\} ds \right); \end{aligned} \quad (3.117)$$

We can also write

$$\begin{aligned} & \int_{n-m+1}^{-}_t(B^{[m+1; :n]})g \leq \int_{n-m+1}^{-}_0(B^{[m+1; :n]})g + 2k-k \int_{s=0}^t \int_{n-m+1}^{-}_s(B^{[m+1; :n]})g ds \\ & \quad + 2 M(d^{(max)})^3(d^{(max)!}) \int_{s=0}^t \int_{n-m+1}^{-}_s(B^{[m+1; :n]})g ds \\ & \quad + M(d^{(max)})^3(d^{(max)!}) \left( \sum_{l=1}^{m+1} \left\{ \sup_{0 \leq s \leq t} \sum_{n-m+1}^{-}_s(B^{[m; :l; n-1]}) \right\} \right); \end{aligned} \quad (3.118)$$

Now by using Gronwall's inequality, we have

$$\left\{ \sum_{n=m+1}^{\infty} -_t(B^{[m+1; \cdot; n]}) \right\} \left\{ \left( \sum_{n=m+1}^{\infty} -_0(B^{[m+1; \cdot; n]}) \right) + 2 M(d^{(max)})^3(d^{(max)!}) \left( \sum_{l=1}^{m+1} \left\{ \sup_{0 \leq s \leq t} \sum_{n=m+1}^{\infty} -_s(B^{[m; \cdot; n-1]}) \right\} \right) \right\} e^{(2k-k+2 M(d^{(max)})^3(d^{(max)!}))t}. \quad (3.119)$$

From (3.115), we conclude that for  $1 \leq r \leq m+1$ ,

$$\lim_{r \rightarrow 0} \left( \sum_{n=m+1}^{\infty} -_t(B^{[m+1; \cdot; n]}) \right) = 0. \quad (3.120)$$

This concludes the proof.

### 3.7.7 Single Server System with State-dependent Arrival Rates

Let us consider a system that has a single PS server in which jobs arrive according to a Poisson process with rate  $\lambda_n$  when there are  $n$  progressing jobs. The job length distribution is  $G(\cdot)$  which is the same as in the system model of our multi-server system. Let  $\mathbb{P}_t^{(single)}$  be the probability measure of the server state at time  $t$ , then we have the following Kolmogorov equations, for  $\mathbb{C}_b^1(U)$ ,

$$\begin{aligned} h_t^{(single)}; i = h_0^{(single)}; i + \int_{s=0}^t h_s^{(single)}; i \Phi_i ds \\ \int_{s=0}^t \left( \sum_{n=1}^{\infty} \sum_{i=1}^n \int_{x_1}^{\infty} \int_{x_n}^{\infty} \left\{ \frac{(x_i)}{n} \right\} ( (n-1; x_1; \dots; x_{i-1}; x_{i+1}; \dots; x_n) \right. \\ \left. (n; x_1; \dots; x_n) \right) d_s^{(single)}(n; x_1; \dots; x_n) \\ + \left[ ( \mathbb{P}_s^{(single)}(f)g ) ( (1; 0) \dots (0) ) + \sum_{n=1}^{\infty} \sum_{i=1}^{n+1} \int_{x_1}^{\infty} \int_{x_n}^{\infty} \left\{ \frac{1}{n+1} \right\} \right. \\ \left. n( (n+1; x_1; \dots; x_{i-1}; 0; x_i; \dots; x_n) \dots (n; x_1; \dots; x_n) \right) d_s^{(single)}(n; x_1; \dots; x_n) \Big] ds: \end{aligned} \quad (3.121)$$

Let  $\rho_t^{(single)}(\mathbf{u})$  be the density of the measure  $\mathbb{P}_t^{(single)}$  at  $\mathbf{u} \in U$ . Then we obtain the differential equations satisfied by the density function  $\rho_t^{(single)} = (\rho_t^{(single)}(\mathbf{u}); \mathbf{u} \in U)$

by using the same procedure as in [35, 67, 68, 81]. Let us define a process  $P_t^{(single)} = (P_t^{(single)}(\mathbf{u}); \mathbf{u} \geq \mathbf{U})$  where  $P_t^{(single)}(0)$  is equal to  $P_t^{(single)}(f0g)$  and

$$P_t^{(single)}(n; y_1; \dots; y_n) = \int_{x_1=0}^{y_1} \int_{x_n=0}^{y_n} p_t^{(single)}(n; x_1; \dots; x_n) dx_1 \dots dx_n:$$

Then we get the following differential equations

$$\frac{dP_t^{(single)}(0)}{dt} = \int_{y=0}^1 (y) \left( \frac{\partial P_t^{(single)}(1; y)}{\partial y} \right) dy - \lambda P_t^{(single)}(0); \quad (3.122)$$

for  $n \geq 1$ ,

$$\begin{aligned} \frac{dP_t^{(single)}(n; y_1; \dots; y_n)}{dt} = & \sum_{i=1}^n \frac{1}{n} \frac{\partial P_t^{(single)}(n; y_1; \dots; y_n)}{\partial y_i} \\ & + \sum_{i=1}^{n+1} \int_{x_i=0}^1 \frac{(x_i)}{n+1} \left( \frac{\partial P_t^{(single)}(n+1; y_1; \dots; y_{i-1}; x_i; y_i; \dots; y_n)}{\partial x_i} \right) dx_i \\ & - \sum_{i=1}^n \int_{x_i=0}^{y_i} \frac{(x_i)}{n} \left( \frac{\partial P_t^{(single)}(n; y_1; \dots; y_{i-1}; x_i; y_{i+1}; \dots; y_n)}{\partial x_i} \right) dx_i \\ & + \sum_{i=1}^n \left( \frac{n-1}{n} \right) P_t^{(single)}(n-1; y_1; \dots; y_{i-1}; y_{i+1}; \dots; y_n) \\ & - \lambda P_t^{(single)}(n; y_1; \dots; y_n); \quad (3.123) \end{aligned}$$

From [22], the single server PS system with state-dependent arrival rate  $\lambda_i$  when there are  $i$  jobs in progress and job lengths are generally distributed with finite mean  $\bar{1}$ , has a unique stationary distribution  $P^{(single)} = (P^{(single)}(\mathbf{u}); \mathbf{u} \geq \mathbf{U})$  given by,

$$P^{(single)}(n; y_1; \dots; y_n) = \frac{\left( \prod_{i=1}^n \frac{\lambda_{i-1}}{\lambda_i} \right)}{1 + \sum_{m=1}^{\infty} \left( \prod_{i=1}^m \frac{\lambda_{i-1}}{\lambda_i} \right)} \prod_{i=1}^n \int_{x_i=0}^{y_i} \bar{G}(x_i) dx_i; \quad (3.124)$$

and

$$P^{(single)}(0) = \frac{1}{1 + \sum_{m=1}^{\infty} \left( \prod_{i=1}^m \frac{\lambda_{i-1}}{\lambda_i} \right)}; \quad (3.125)$$

### 3.7.8 Proof of Theorem 3.4

To prove this result, we work with the mean-field of the tail occupancy process. The proof uses the result that the mean-field is quasi-monotonic for Policies 2 and 4. We show that the mean-field is GAS and we skip the rest of the details as they follow by similar



arguments of [11, Section 2.6]. We use the spaces  $\mathbb{U}^\theta$  and  $\mathbb{U}$  defined in (1.23) and (1.24), respectively. We equip the space  $\mathbb{U}^\theta$  with the metric  $!(\cdot; \cdot)$  defined in (1.25). The space  $\mathbb{U}^\theta$  equipped with the metric  $!$  is compact, complete, and separable. For the stochastic system with parameter  $N$ , let  $(\mathbf{X}_i^{(N)}(t); i = 0; t = 0)$  be the empirical process where  $\mathbf{X}_i^{(N)}(t)$  denotes the fraction of servers with at least  $i$  progressing jobs. Then if  $\mathbf{X}^{(N)}(0) = \mathbf{u}$ , then  $(\mathbf{X}^{(N)}(t); t = 0) = (\mathbf{x}(t; \mathbf{u}); t = 0)$  which is the unique solution to the following equations,

$$\mathbf{x}(0; \mathbf{u}) = \mathbf{u} \quad (3.126)$$

$$\frac{dx_i(t; \mathbf{u})}{dt} = F_i(\mathbf{x}(t; \mathbf{u})) \quad (3.127)$$

$$= (x_{i-1}(t; \mathbf{u}) - x_i(t; \mathbf{u})) \Phi_{i-1}(\mathbf{x}(t; \mathbf{u})) - x_i(t; \mathbf{u}) + x_{i+1}(t; \mathbf{u}); \quad (3.128)$$

where  $\Phi_i(\mathbf{x}(t; \mathbf{u}))$  is the arrival rate function given in (3.11) with  $Q_j(\frac{\cdot}{N})$  replaced by  $x_j(t; \mathbf{u}) - x_{j+1}(t; \mathbf{u})$ . Let  $\mathbf{F}(\mathbf{x}(t; \mathbf{u})) = (F_i(\mathbf{x}(t; \mathbf{u}); i = 1)$ . Then any fixed-point  $\mathbf{x} = (x_i; i = 0)$  of the mean-field satisfies  $\mathbf{F}(\mathbf{x}) = 0$ .

We show that  $\lim_{t \rightarrow \infty} \mathbf{x}(t; \mathbf{u}) = \mathbf{x}^*$  for all  $\mathbf{u} \in \mathbb{U}$ . The following result is the key step to conclude that the asymptotic stationary distribution of a server coincides with the unique fixed-point of the mean-field.

**Proposition 3.1.** *For Policies 2 and 4, the mean-field is quasi-monotonic. Furthermore,  $\lim_{t \rightarrow \infty} \mathbf{x}(t; \mathbf{u}) = \mathbf{x}^*$  for all  $\mathbf{u} \in \mathbb{U}$ .*

*Proof.* For Policy 2 and Policy 4, the mean-field equations are

$$\frac{dx_i(t; \mathbf{u})}{dt} = F_i(\mathbf{x}(t; \mathbf{u})); \quad (3.129)$$

where we have for  $1 \leq i \leq d-1$ ,

$$F_i(\mathbf{x}(t; \mathbf{u})) = \frac{(1 - x_{i+1}^d(t; \mathbf{u}))}{(1 - x_{i+1}(t; \mathbf{u}))} (x_{i-1}(t; \mathbf{u}) - x_i(t; \mathbf{u})) - x_i(t; \mathbf{u}) + x_{i+1}(t; \mathbf{u}); \quad (3.130)$$

and for  $i = d$ ,

$$F_d(\mathbf{x}(t; \mathbf{u})) = (x_{d-1}^d(t; \mathbf{u}) - x_d^d(t; \mathbf{u})) - x_d(t; \mathbf{u}) + x_{d+1}(t; \mathbf{u}); \quad (3.131)$$

Since  $F_i(\cdot)$  is non-decreasing in  $x_j(t; \mathbf{u})$  for  $i \neq j$ , the mean-field is quasi-monotone from [82, p. 70-74]. Hence if  $\mathbf{u} \leq \mathbf{v}$  by element wise, then  $\mathbf{x}(t; \mathbf{u}) \leq \mathbf{x}(t; \mathbf{v})$  by element wise. As a result, we have  $\mathbf{x}(t; \min(\mathbf{u}; \cdot)) \leq \mathbf{x}(t; \mathbf{u}) \leq \mathbf{x}(t; \max(\mathbf{u}; \cdot))$ . Therefore, to show that  $\lim_{t \rightarrow \infty} \mathbf{x}(t; \mathbf{u}) = \mathbf{x}^*$  for all  $\mathbf{u} \in \mathbb{U}$ , it is sufficient to show  $\lim_{t \rightarrow \infty} \mathbf{x}(t; \mathbf{u}) = \mathbf{x}^*$  for all  $\mathbf{u} \leq \mathbf{x}^*$  and for all  $\mathbf{u} \geq \mathbf{x}^*$ .

Let us now define  $z_n(t; \mathbf{u}) = \sum_{k=1}^n x_k(t; \mathbf{u})$  and  $z_n(\mathbf{u}) = \sum_{k=1}^n u_k$ . If  $\mathbf{u} \in U$ , then  $\mathbf{x}(t; \mathbf{u}) \in U$  and further, we have for  $1 \leq n \leq d+1$ ,

$$\frac{dz_n(t; \mathbf{u})}{dt} = \frac{(1 - x_{n+1}^d(t; \mathbf{u}))}{(1 - x_{n+1}(t; \mathbf{u}))} (x_{n-1}(t; \mathbf{u}) - x_{n+1}(t; \mathbf{u})) + x_{n+1}^d(t; \mathbf{u}) - x_n(t; \mathbf{u}); \quad (3.132)$$

and for  $n = d+2$ ,

$$\frac{dz_n(t; \mathbf{u})}{dt} = x_{n-1}^d(t; \mathbf{u}) - x_n(t; \mathbf{u}); \quad (3.133)$$

The fixed-point  $\mathbf{z}^*$  satisfies the following equations, for  $1 \leq n \leq d+1$ ,

$$\frac{(1 - x_{n+1}^d)}{(1 - x_{n+1})} (x_{n-1} - x_{n+1}) + x_{n+1}^d - z_n = 0; \quad (3.134)$$

and for  $n = d+2$ , we have

$$x_{n-1}^d - z_n = 0; \quad (3.135)$$

We first show that  $z_n(t; \mathbf{u})$  is bounded uniformly in  $t$  for  $n = 1$ . Note that

$$\frac{dz_1(t; \mathbf{u})}{dt} = x_1(t; \mathbf{u});$$

If  $\mathbf{u} \in U_1$ , then  $\mathbf{x}(t; \mathbf{u}) \in U_1$  and hence,  $z_1(t; \mathbf{u}) \leq z_1(\mathbf{u})$ . Since  $z_n(t; \mathbf{u}) \leq z_{n+1}(t; \mathbf{u})$ , we have that  $z_n(t; \mathbf{u})$  is uniformly bounded in  $t$ . On the other hand, if  $\mathbf{u} \in U_2$ , then  $\mathbf{x}(t; \mathbf{u}) \in U_2$  and hence,

$$\frac{dz_1(t; \mathbf{u})}{dt} - x_1 = 0;$$

This implies  $z_1(t; \mathbf{u}) \leq z_1(\mathbf{u})$  and hence,  $z_n(t; \mathbf{u})$  for  $n = 1$  is uniformly bounded in  $t$ . When  $\mathbf{u} \in U_2$ , since the derivative of  $x_n(t; \mathbf{u})$  is bounded for all  $n = 1, \dots, d+1$ , then  $\lim_{t \rightarrow \infty} x(t; \mathbf{u}) = \mathbf{z}^*$  follows if we show for all  $n = 1, \dots, d+1$ ,

$$\int_{t=0}^1 (x_n(t; \mathbf{u}) - z_n) dt < \epsilon;$$

Similarly, for the case of  $\mathbf{u} \in U_1$ , it is sufficient to show

$$\int_{t=0}^1 (z_n - x_n(t; \mathbf{u})) dt < \epsilon$$

for all  $n = 1, \dots, d+1$ . We only give a proof for the case when  $\mathbf{u} \in U_2$  and the proof for the other case is similar.

We first begin with the proof of

$$\int_{t=0}^1 (x_1(t; \mathbf{u}) - z_1) dt < \epsilon;$$

From (3.134) and (3.132), we have

$$\frac{dz_1(t; \mathbf{u})}{dt} = (x_1(t; \mathbf{u}) - 1):$$

Then we have

$$\int_{t=0}^T (x_1(t; \mathbf{u}) - 1) dt = \int_{t=0}^T \frac{dz_1(t; \mathbf{u})}{dt} dt:$$

This implies

$$\int_{t=0}^T (x_1(t; \mathbf{u}) - 1) dt = z_1(\mathbf{u}) - z_1(T; \mathbf{u}):$$

Since  $z_1(t; \mathbf{u})$  is uniformly bounded in  $t$ , we get

$$\int_{t=0}^T (x_1(t; \mathbf{u}) - 1) dt < 1$$

as  $T \neq 1$ .

We now focus on the proof of

$$\int_{t=0}^1 (x_n(t; \mathbf{u}) - n) dt < 1$$

for  $2 \leq n \leq K + 1$ . We prove this by using the induction method. For  $2 \leq K \leq K + 1$ , let

$$\int_{t=0}^1 (x_i(t; \mathbf{u}) - i) dt < 1$$

is true for  $1 \leq i \leq K - 1$ , then we prove that

$$\int_{t=0}^1 (x_K(t; \mathbf{u}) - K) dt < 1:$$

From (3.132), we have

$$\frac{dz_1(t; \mathbf{u})}{dt} - \frac{dz_K(t; \mathbf{u})}{dt} = \frac{(1 - x_{+1}^d(t; \mathbf{u}))}{(1 - x_{+1}(t; \mathbf{u}))} (1 - x_{K-1}(t; \mathbf{u})) - x_1(t; \mathbf{u}) + x_K(t; \mathbf{u}):$$

Then from (3.134), we have

$$\begin{aligned} \frac{dz_1(t; \mathbf{u})}{dt} - \frac{dz_K(t; \mathbf{u})}{dt} &= \frac{(1 - x_{+1}^d(t; \mathbf{u}))}{(1 - x_{+1}(t; \mathbf{u}))} (1 - x_{K-1}(t; \mathbf{u})) - \frac{(1 - x_{+1}^d)}{(1 - x_{+1})} (1 - x_{K-1}) \\ &\quad - x_1(t; \mathbf{u}) + x_K(t; \mathbf{u}) - K: \end{aligned}$$

Therefore, we obtain

$$\int_{t=0}^T (x_K(t; \mathbf{u}) - K) dt = \int_{t=0}^T \left( \frac{dz_1(t; \mathbf{u})}{dt} - \frac{dz_K(t; \mathbf{u})}{dt} - \frac{(1 - x_{+1}^d(t; \mathbf{u}))}{(1 - x_{+1}(t; \mathbf{u}))} (1 - x_{K-1}(t; \mathbf{u})) \right) dt$$

$$+ \frac{(1 - x_{+1}^d)}{(1 - x_{+1})} (1 - x_{K-1}) + x_1(t; \mathbf{u}) - 1) dt:$$

Since  $x_{+1}(t; \mathbf{u}) \leq 1$ , we have

$$\int_{t=0}^T (x_K(t; \mathbf{u}) - \kappa) dt = \int_{t=0}^T \left( \frac{dz_1(t; \mathbf{u})}{dt} - \frac{dz_K(t; \mathbf{u})}{dt} - \frac{(1 - x_{+1}^d(t; \mathbf{u}))}{(1 - x_{+1}(t; \mathbf{u}))} (1 - x_{K-1}(t; \mathbf{u})) + \frac{(1 - x_{+1}^d(t; \mathbf{u}))}{(1 - x_{+1}(t; \mathbf{u}))} (1 - x_{K-1}) + x_1(t; \mathbf{u}) - 1 \right) dt:$$

Further, as  $\frac{(1 - x_{+1}^d(t; \mathbf{u}))}{(1 - x_{+1}(t; \mathbf{u}))} \leq d$ , we get

$$\int_{t=0}^T (x_K(t; \mathbf{u}) - \kappa) dt \leq \int_{t=0}^T \left( \frac{dz_1(t; \mathbf{u})}{dt} - \frac{dz_K(t; \mathbf{u})}{dt} + d (x_{K-1}(t; \mathbf{u}) - x_{K-1}) + x_1(t; \mathbf{u}) - 1 \right) dt:$$

From the induction hypothesis and since  $z_n(t; \mathbf{u})$  is bounded uniformly in  $t$  for  $n \leq 1$ , we get

$$\int_{t=0}^1 (x_K(t; \mathbf{u}) - \kappa) dt < 1$$

as  $T \leq 1$ .

Finally, we show

$$\int_{t=0}^1 (x_n(t; \mathbf{u}) - \kappa) dt < 1$$

for  $n \leq K + 2$ . We prove this by using the induction method. For  $K \leq K + 2$ , let

$$\int_{t=0}^1 (x_i(t; \mathbf{u}) - \kappa) dt < 1$$

is true for  $1 \leq i \leq K - 1$ , then we prove that

$$\int_{t=0}^1 (x_K(t; \mathbf{u}) - \kappa) dt < 1:$$

From (3.133), we have

$$\int_{t=0}^T (x_K(t; \mathbf{u}) - \kappa) dt = \int_{t=0}^T \left( \frac{dz_K(t; \mathbf{u})}{dt} + (x_{K-1}^d(t; \mathbf{u}) - \kappa) \right) dt:$$

From the induction hypothesis and since  $z_K(t; \mathbf{u})$  is bounded uniformly in  $t$ , we get  $\int_{t=0}^1 (x_K(t; \mathbf{u}) - \kappa) dt < 1$  as  $T \leq 1$ . This completes the proof.

### 3.7.9 Proof of Theorem 3.5

The proof mainly uses the logic that in the stationary regime, the servers' states are exchangeable random variables. We can write

$$\begin{aligned} & \left| \mathbb{E} \left[ \prod_{i=1}^l \mathbf{s}_i^{(N_k)}(\tau) \right] - \mathbb{E} \left[ \prod_{i=1}^l h \mathbf{V}(\tau); i \right] \right| \\ & \quad \left| \mathbb{E} \left[ \prod_{i=1}^l \mathbf{s}_i^{(N_k)}(\tau) \right] - \mathbb{E} \left[ \prod_{i=1}^l h \mathbf{V}^{(N_k)}(\tau); i \right] \right| \\ & \quad + \left| \mathbb{E} \left[ \prod_{i=1}^l h \mathbf{V}^{(N_k)}(\tau); i \right] - \mathbb{E} \left[ \prod_{i=1}^l h \mathbf{V}(\tau); i \right] \right|; \end{aligned} \quad (3.136)$$

Note that since  $\mathbf{V}^{(N_k)}(\tau) \rightarrow \mathbf{V}(\tau)$ , the second term on the right hand side of the above inequality vanishes as  $N_k \rightarrow \infty$ . Now, due to exchangeability, the permutation of states between servers does not affect the joint distribution. Hence, we have

$$\mathbb{E} \left[ \prod_{i=1}^l \mathbf{s}_i^{(N_k)}(\tau) \right] = \frac{1}{(N_k)_l} \mathbb{E} \left[ \sum_{\sigma \in Q(l; N_k)} \prod_{i=1}^l \mathbf{s}_{\sigma(i)}^{(N_k)}(\tau) \right];$$

where  $(N)_j = N(N-1)\cdots(N-j+1)$ , and  $Q(r; n)$  denotes the set of all permutations of the numbers  $1, 2, \dots, n$  taken  $r$  at a time. Also, by the definition of  $\mathbf{V}^{(N_k)}(\tau)$ , we have

$$\mathbb{E} \left[ \prod_{i=1}^l h \mathbf{V}^{(N_k)}(\tau); i \right] = \mathbb{E} \left[ \left( \prod_{i=1}^l \frac{1}{N_k} \sum_{j=1}^{N_k} \mathbf{s}_j^{(N_k)}(\tau) \right) \right];$$

Hence, the first term on the right hand side of (3.136) can be bounded as follows

$$\left| \mathbb{E} \left[ \prod_{i=1}^l \mathbf{s}_i^{(N_k)}(\tau) \right] - \mathbb{E} \left[ \prod_{i=1}^l h \mathbf{V}^{(N_k)}(\tau); i \right] \right| \leq 2B^l \left( 1 - \frac{(N_k)_l}{(N_k)^l} \right);$$

where  $\max_j k_j = B$ . The result follows since  $\left( 1 - \frac{(N_k)_l}{(N_k)^l} \right) \rightarrow 0$  as  $N_k \rightarrow \infty$ .

Finally, from (3.52), any finite set of servers are independent of each other iff  $Z$  is a Dirac measure. Otherwise they are coupled through the sample value of the random variable  $\mathbf{V}(\tau)$ . If  $Z = \delta_{\mathbf{v}}$ , then it implies that in the limiting system, the stationary empirical random variable  $\mathbf{V}(\tau)$  is a deterministic value coinciding with  $\mathbf{v}$ . Then the following equation concludes that each server has distribution

$$\mathbb{E} \left[ \prod_{i=1}^l \mathbf{s}_i^{(N_k)}(\tau) \right] \rightarrow \prod_{i=1}^l h; i$$

as  $N_k \rightarrow \infty$ . This completes the proof.

## 3.8 Conclusions

In this chapter, we studied occupancy dependent policies that subsume threshold based policies and policies that are adaptive to variations in the system load using mean-field techniques. Our results conclude that the fixed-point of the mean-field is insensitive for any occupancy based policy as long as the mean-field when job lengths are exponentially distributed has a unique fixed-point. For general JLDs, when the fixed-point of the mean-field is unique, it is of great interest to show that the stationary distribution of a server as  $N \rightarrow \infty$  coincides with the unique fixed-point of the mean-field. This result seems to be true as it is observed in our simulation results for the policies studied in this chapter.

## Chapter 4

# A Functional Central Limit Theorem for Multi-Server Erlang Loss Systems Under SQ( $d$ ) Load Balancing Policy

In this chapter, we return to the multi-server Erlang loss system that has  $N$  servers, each server can serve only a finite number of  $C$  jobs simultaneously. We assume that the dispatcher employs the SQ( $d$ ) load balancing policy, according to which the destination server for an incoming job is the server with the least occupancy among  $d$  randomly sampled servers. The job arrival process is a Poisson process with rate  $N \lambda^{(N)}$ , where  $\lambda^{(N)} \geq \mathbb{R}_+$ .

The main contribution of this chapter is to derive a functional central limit theorem (FCLT) by considering the asymptotic regime,  $N \rightarrow \infty$ . The FCLT that we establish in this chapter is useful to obtain an approximation to the blocking probability for a job in the system with  $N$  servers when  $N$  is large. In earlier works [8, 10], they obtained the blocking probability for a job when  $N \rightarrow \infty$  with the help of mean-field techniques. Without loss of generality, we call the blocking probability obtained when  $N \rightarrow \infty$  as the asymptotic blocking probability. Quantifying the error between the actual blocking probability for the system with  $N$  servers and the asymptotic blocking probability is of great interest. We use the FCLT which we establish in this chapter to quantify the resulting error between the actual blocking probabilities and the asymptotic blocking probabilities, as a function of  $N$ . We assume that  $\lambda^{(N)} = \lambda \sqrt{N}$  for  $\lambda \in \mathbb{R}$  and  $\lambda \geq \mathbb{R}_+$ . Such an assumption includes the Halfin-Whitt regime [45] as a special case which corresponds to the case when  $\lambda = C$  and  $\lambda \neq 0$ . The Halfin-Whitt regime is so-called an efficiency driven heavy traffic regime that

allows us to study how the effects of increasing load affects performance.

## Organization of the Chapter

The rest of the chapter is organized as follows: We introduce the system model in Section 4.1. We then present the required notation and some preliminary results in Section 4.2. In Section 4.3, we state the main results of the chapter. In this section, we begin first with some preliminary results on the transient regime and then we state the functional central limit theorem for the transient regime in Theorem 4.3. After that we present some preliminary results for the stationary regime and then we give the functional central limit theorem for the stationary regime in Theorem 4.6. We conclude Section 4.3 with the result on approximations to the average blocking probability of the system with  $N$  servers. We give proofs of the main results in Section 4.4. Finally, we conclude in Section 4.5

## 4.1 System Model

In this section, we give a detailed description of the system model. We study a large-scale multi-server Erlang loss system with a central job dispatcher which routes an incoming job request to one of the servers according to the SQ( $d$ ) policy. The number of servers in the system is equal to  $N$ , a large value in the order of 10,000. Each server can serve at most  $C$  jobs simultaneously. As a result, the system's capacity indicating the maximum number of jobs that can be processed simultaneously in the system is equal to  $NC$ . When a job arrives, the dispatcher routes the arrived job to one of the servers, and the destination server of the arrival accepts the job if its occupancy is less than  $C$ . Otherwise, the destination server of the arrival blocks the job from service. The system immediately discards a blocked job. If a server accepts a job, it will serve or process the job at a constant unit rate until the service of the job is completed. We assume that the service time distributions are exponential with unit mean. The incoming jobs arrive according to a Poisson process with rate  $N^{(N)}$ , where  $N^{(N)} \geq \mathbb{R}_+$  and it is defined precisely later in this section.

**Remark 4.1.** *For our model, it does not matter whether the dispatcher samples  $d$  servers without replacement or with replacement to dispatch an arrival as they lead to the same asymptotic results. The proof follows by the same arguments as in [47, pages 11-12]. Therefore for simplicity, we assume that the dispatcher samples  $d$  servers with replacement upon an arrival.*



We next introduce the parameter  $\rho^{(N)}$ . For  $\lambda \in \mathbb{R}_+$  and  $\tilde{\rho} \in \mathbb{R}$ , the parameter  $\rho^{(N)}$  is defined as

$$\rho^{(N)} = \frac{\lambda}{c} + \frac{\tilde{\rho}}{N}. \quad (4.1)$$

We are interested in obtaining an FCLT for our model. The particular form of  $\rho^{(N)}$  in (4.1) allows us to study the following two regimes as special cases:

1. For  $\tilde{\rho} = 0$ : The considered case corresponds to the case when the job arrival process is a Poisson process with rate  $\lambda$ . The previous works [8, 10] focused on obtaining a functional law of large numbers result referred to as the mean-field limit. They derived the blocking probabilities for a job when  $N \rightarrow \infty$  as a function of the unique fixed-point of the mean-field. The FCLT that we establish in this chapter allows us to show that the error between the actual blocking probability for a job in the system with  $N$  servers and the asymptotic blocking probability obtained in terms of the fixed-point of the mean-field is  $o(N^{-\frac{1}{2}})$ .
2. For  $\tilde{\rho} \neq 0$  and  $c = C$ : In this case, the system is in a heavy-traffic regime when  $N$  is large. In fact, the resulting model corresponds to the Halfin-Whitt regime since the traffic intensity  $\rho^{(N)} = \frac{\lambda}{c} + \frac{\tilde{\rho}}{N}$  approaches one as  $N \rightarrow \infty$ . Then the arrival rate of jobs  $\lambda$  and the system capacity  $NC$  are related as  $NC = \lambda + \tilde{\rho}$ , and they converge to  $\lambda$  as  $N \rightarrow \infty$ . In the Halfin-Whitt regime, we obtain an approximation to the blocking probability for a job by exploiting the FCLT that we establish in this chapter.

## 4.2 Additional Notation and Preliminary Results

In this section, we introduce the required notation and provide some preliminary results.

### 4.2.1 Additional Notation

We now introduce the required additional notation.

Due to the fact that the SQ( $d$ ) policy uses only the occupancy information of servers but not their identities, we use the Markov process  $(\mathbf{X}^{(N)}(t); t \geq 0)$  to describe the system

evolution, where  $\mathbf{X}^{(N)}(t) = (\mathbf{X}_i^{(N)}(t); 0 \leq i \leq C)$  and  $\mathbf{X}_i^{(N)}(t)$  denotes the fraction of the servers with at least  $i$  progressing jobs at time  $t$ . Let us define a space  $\mathcal{U}$  as

$$\mathcal{U} = \{f(u_0; u_1; \dots; u_C) : u_0 = 1 - u_1 - \dots - u_C \geq 0\}. \quad (4.2)$$

It is evident that  $\mathbf{X}^{(N)}(t)$  lies in the space  $\mathcal{U}$ . Without loss of generality, we write an element of the form  $(u_0; \dots; u_C)$  as  $\mathbf{u}$ . The space  $\mathcal{U}$  is equipped with the metric generated by the euclidean norm  $\|\cdot\|_2$  defined as

$$\|\mathbf{u}\|_2 = \sqrt{\sum_{i=0}^C u_i^2}, \quad (4.3)$$

where  $\mathbf{u} = (u_0; \dots; u_C)$ . It can be verified that the space  $\mathcal{U}$  is a Polish space.

## 4.2.2 Preliminary Results

We next provide a mathematical framework for the main problem and present some preliminary results.

We first present a mathematical formulation to the time evolution of the process  $(\mathbf{X}^{(N)}(t); t \geq 0)$ . Upon an arrival, if the system state is  $\mathbf{b} = (b_0; \dots; b_C)$  implying that the fraction of servers with at least  $i$  progressing jobs is equal to  $b_i$  for  $0 \leq i \leq C$ , then the destination server of the job will have occupancy  $n$  with probability  $b_n^d = b_{n+1}^d$ . Since the rate of the arrival process is  $N^{(N)}$ , the total instantaneous rate of arrivals to servers having  $n$  jobs is  $N^{(N)}((\mathbf{X}_{n-1}^{(N)}(t))^d - (\mathbf{X}_n^{(N)}(t))^d)$ . Furthermore, as job lengths have exponential distributions with unit rate, the total instantaneous rate of departures from servers having  $n$  jobs is  $nN(\mathbf{X}_n^{(N)}(t) - \mathbf{X}_{n+1}^{(N)}(t))$ . Then we can write the system dynamics using random time change of a set of mutually independent unit rate Poisson processes as described in [83, Section 2.1].

Let  $f(N_i(t); t \geq 0)g_{i-1}$  be the set of mutually independent unit rate Poisson processes. We use  $(N_i(t); t \geq 0)$  to model the arrival process to servers that have  $i-1$  progressing jobs. Similarly, let  $f(D_i(t); t \geq 0)g_{i-1}$  be the collection of a set of mutually independent unit rate Poisson processes. Furthermore, the set of processes  $f(D_i(t); t \geq 0)g_{i-1}$  is independent of the set of processes  $f(N_i(t); t \geq 0)g_{i-1}$ . Also,  $f(N_i(t); t \geq 0)g_{i-1}$  and  $f(D_i(t); t \geq 0)g_{i-1}$  are independent of  $\mathbf{X}^{(N)}(0)$ . We use  $(D_i(t); t \geq 0)$  to model the departure process from servers that have  $i$  progressing jobs. Since the arrival process of jobs to the system is a

Poisson process with rate  $N^{(N)}$  and the service time distributions are exponential with unit mean, as in [83], we can write

$$\mathbf{X}_0^{(N)}(t) = 1; \quad (4.4)$$

and for  $n \geq 1$ ,

$$\begin{aligned} \mathbf{X}_n^{(N)}(t) = \mathbf{X}_n^{(N)}(0) + \frac{1}{N} N_n \left( N^{(N)} \int_{s=0}^t ((\mathbf{X}_n^{(N)}(s))^d - (\mathbf{X}_n^{(N)}(s))^d) ds \right) \\ \frac{1}{N} D_n \left( N n \int_{s=0}^t ((\mathbf{X}_n^{(N)}(s)) - (\mathbf{X}_{n+1}^{(N)}(s))) ds \right); \end{aligned} \quad (4.5)$$

We choose the filtration  $(F^{(N)}(t); t \geq 0)$  where

$$\begin{aligned} F^{(N)}(t) = \left( \mathbf{X}^{(N)}(0); N_n \left( N^{(N)} \int_{s=0}^r ((\mathbf{X}_n^{(N)}(s))^d - (\mathbf{X}_n^{(N)}(s))^d) ds \right); \right. \\ \left. D_n \left( N n \int_{s=0}^r ((\mathbf{X}_n^{(N)}(s)) - (\mathbf{X}_{n+1}^{(N)}(s))) ds \right); n \geq 1; 0 \leq r \leq t \right); \end{aligned} \quad (4.6)$$

augmented by all null sets.

We now present results on the mean-field analysis of the model without proofs as they directly follow from the case  $N^{(N)} = b$  for  $b \geq 2 \mathbb{R}_+$ , studied in [10]. The mean-field equations (MFEs) in our case are the same as the  $N^{(N)} =$  case except that  $\mathbf{X}$  replaces  $\mathbf{x}$  in the MFEs given in 1.3.2. We recall the MFEs in this chapter as we use them frequently.

**Theorem 4.1.** *For  $\mathbf{u} \geq \mathbf{0}$ , if  $\mathbf{X}^{(N)}(0) \geq \mathbf{u}$  as  $N \rightarrow \infty$ , then  $(\mathbf{X}^{(N)}(t); t \geq 0) \rightarrow (\mathbf{x}(t; \mathbf{u}); t \geq 0)$  as  $N \rightarrow \infty$  where  $(\mathbf{x}(t; \mathbf{u}); t \geq 0) = (x_n(t; \mathbf{u}); t \geq 0; 0 \leq n \leq C)$  is the unique solution to the following equations:*

for  $\mathbf{h}(\mathbf{x}(t; \mathbf{u})) = (h_n(\mathbf{x}(t; \mathbf{u})); 0 \leq n \leq C)$ ,

$$\mathbf{x}(0; \mathbf{u}) = \mathbf{u}; \quad \frac{dx_n(t; \mathbf{u})}{dt} = h_n(\mathbf{x}(t; \mathbf{u})); \quad (4.7)$$

where

$$h_0(\mathbf{x}(t; \mathbf{u})) = 0; \quad (4.8)$$

and for  $n \geq 1$ ,

$$h_n(\mathbf{x}(t; \mathbf{u})) = (x_{n-1}^d(t; \mathbf{u}) - x_n^d(t; \mathbf{u})) - n(x_n(t; \mathbf{u}) - x_{n+1}(t; \mathbf{u})); \quad (4.9)$$

with  $x_0(t; \mathbf{u}) = 1$  and  $x_{C+1}(t; \mathbf{u}) = 0$ . The deterministic process  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$  is referred to as the mean-field limit and equations (4.7)-(4.9) are referred to as the mean-field equations with initial point  $\mathbf{u}$ .

Without loss of generality, we say that a process  $(\mathbf{y}(t); t \geq 0)$  is a solution to the differential equations (4.7)-(4.9), it means that it is the unique generic solution with initial point  $\mathbf{y}(0)$ .

The mean-field  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$  has a unique global asymptotically stable fixed-point  $\mathbf{x} = (x_n; 0 \leq n \leq C)$  with  $x_0 = 1$ . Also, the following exchange of limits holds

$$\lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \mathbf{X}^{(N)}(t) = \lim_{t \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbf{X}^{(N)}(t): \quad (4.10)$$

Using (4.10), under the assumption of the exchangeability of initial servers' states, we can show the independence of any finite set of servers as  $N \rightarrow \infty$ . Also, for a server,  $\mathbf{x}(t; \mathbf{u})$  denotes the distribution at time  $t$  and  $\mathbf{x}$  denotes the stationary distribution, as  $N \rightarrow \infty$ . As a result,  $x_C$  denotes the stationary probability that a server is fully occupied when  $N \rightarrow \infty$ . Since the dispatcher samples  $d$  servers when a job arrives, the stationary average blocking probability of a job as  $N \rightarrow \infty$  is equal to  $\frac{d}{C}$ , where we use the fact that the chosen  $d$  servers are independent of each other. Our objective is to find the gap between the actual blocking probability of the system with  $N$  servers and  $\frac{d}{C}$  as a function of  $\rho, \lambda, \mu, N$ , and  $C$ .

We can find  $\mathbf{x}$  numerically as explained below. The fixed-point  $\mathbf{x}$  is the unique solution to the following equations

$$x_{n-1} \binom{d}{n-1} x_n = n(x_n - x_{n+1}) \quad (4.11)$$

for  $n \geq 1$  and  $x_{C+1} = 0$ . Then from (4.11), we can also write

$$\frac{\binom{d}{n-1} x_n}{\binom{d}{n-1} x_n} (x_{n-1} - x_n) = n(x_n - x_{n+1}) \quad (4.12)$$

for  $n \geq 1$  and  $x_{C+1} = 0$ . Let us define  $\hat{x}_n = \frac{\binom{d}{n} x_{n+1}}{\binom{d}{n} x_n}$ . Then from (4.12),  $\hat{x}_n$  is the stationary distribution of the single server loss model with a Poisson arrival process of jobs having rate  $\hat{x}_n$  when there are  $n$  progressing jobs, and  $x_n$  is the probability that the server has at least  $n$  progressing jobs. Let  $M_1(\mathcal{F}; 1; \rho; Cg)$  be the set of probability measures on  $\mathcal{F}; 1; \rho; Cg$ . Then from [10], the fixed-point  $\mathbf{x}$  can be computed using the formula for the stationary distribution of a single server loss system with state-dependent arrival rates. We first define two mappings,  $\Theta : M_1(\mathcal{F}; 1; \rho; Cg) \rightarrow \mathbb{R}_+^{C+1}$  and  $\hat{\Xi} : \mathbb{R}_+^{C+1} \rightarrow M_1(\mathcal{F}; 1; \rho; Cg)$  that are used in computing  $\mathbf{x}$ . For every  $(\rho_0; \dots; \rho_C) \in \mathbb{R}_+^{C+1}$ , there exists  $(r_0; \dots; r_C) \in \mathbb{R}_+^{C+1}$  such that

$$\Theta((\rho_0; \dots; \rho_C)) = (r_0; \dots; r_C); \quad (4.13)$$

where

$$r_n = \frac{((\sum_{j=n}^C \rho_j)^d - (\sum_{i=n+1}^C \rho_i)^d)}{((\sum_{j=n}^C \rho_j) - (\sum_{i=n+1}^C \rho_i))}. \quad (4.14)$$

Similarly, for every  $(b_0; \dots; b_C) \in \mathbb{R}_+^{C+1}$ , there exists  $(a_0; \dots; a_C) \in \mathcal{M}_1(\mathbb{R}_+; 1; \dots; Cg)$  such that

$$\widehat{\Xi}((b_0; \dots; b_C)) = (a_0; \dots; a_C); \quad (4.15)$$

where

$$a_n = \left( \prod_{i=1}^n \left( \frac{b_i - 1}{i} \right) \right) a_0 \quad (4.16)$$

for  $n \geq 1$  and  $\sum_{i=0}^C a_i = 1$ . Then  $\mathbf{z}$  can be computed from the fixed-point of the mapping  $\widehat{\Xi}(\Theta)$  as in Lemma 1.1 except that we replace  $\mathbf{b}$  by  $\mathbf{z}$ .

We first study the fluctuation process  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  as  $N \rightarrow \infty$  where

$$\mathbf{Z}^{(N)}(t) = \mathcal{P}_N^{-1} \mathbf{N}(\mathbf{X}^{(N)}(t) - \mathbf{x}(t; \mathbf{u})); \quad (4.17)$$

It can be checked that  $\mathbf{Z}^{(N)}(t)$  lies in the space  $\mathcal{V}$  defined as

$$\mathcal{V} = \{f(r_0; \dots; r_C) : r_0 = 0 \text{ and } r_i \in \mathbb{R}; 1 \leq i \leq Cg\}. \quad (4.18)$$

We equip the space  $\mathcal{V}$  with the topology induced by the euclidean norm (4.3). For a linear operator  $K : \mathcal{V} \rightarrow \mathcal{V}$ , let the operator norm  $\|K\|_2$  be defined as

$$\|K\|_2 = \sup_{\mathbf{v} \in \mathcal{V}} \frac{\|K\mathbf{v}\|_2}{\|\mathbf{v}\|_2}. \quad (4.19)$$

We next obtain the time evolution of the process  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  by using (4.5), (4.7), and (4.17). For this, we first define the following three useful operators  $W_1; W_2; W_3 : \mathcal{U} \rightarrow \mathcal{V}$  as follows: for  $\mathbf{b} \in \mathcal{U}$ ,

$$(W_1(\mathbf{b}))_0 = 0; (W_2(\mathbf{b}))_0 = 0; (W_3(\mathbf{b}))_0 = 0; \quad (4.20)$$

and for  $n \geq 1$ ,

$$(W_1(\mathbf{b}))_n = (b_n^d - b_n); (W_2(\mathbf{b}))_n = n(b_n - b_{n+1}); (W_3(\mathbf{b}))_n = \tilde{\sim}(b_n^d - b_n); \quad (4.21)$$

From (4.9) and (4.21), we have

$$h_n(\mathbf{b}) = (W_1(\mathbf{b}))_n - (W_2(\mathbf{b}))_n; \quad (4.22)$$

Furthermore, let  $W : \mathbb{U} \rightarrow \mathbb{V}$  be the operator defined as

$$W = W_1 \quad W_2: \quad (4.23)$$

The operator  $W$  is Lipschitz continuous satisfying the following inequality for all  $\mathbf{a}, \mathbf{b} \in \mathbb{U}$ ,

$$\|W(\mathbf{a}) - W(\mathbf{b})\|_2 \leq B_W \|\mathbf{a} - \mathbf{b}\|_2, \quad (4.24)$$

where  $B_W = 2d^{\frac{\rho}{2} + C^2}$ .

We now define a set of independent square-integrable martingales  $(\mathbf{M}^{(N)}(t); t \geq 0) = f(\mathbf{M}_i^{(N)}(t); t \geq 0)_{i \in \{1, \dots, d\}}$  adapted to the filtration  $(F^{(N)}(t); t \geq 0)$  such that  $(\mathbf{M}^{(N)}(t); t \geq 0)$  is independent of  $\mathbf{Z}^{(N)}(0)$  and for  $i = 1, \dots, d$ ,

$$\langle \mathbf{M}_i^{(N)} \rangle_t = \int_{s=0}^t \left( (W_1(\mathbf{X}^{(N)}(s)))_i + (W_2(\mathbf{X}^{(N)}(s)))_i - \frac{1}{N} (W_3(\mathbf{X}^{(N)}(s)))_i \right) ds \quad (4.25)$$

Then from (4.5)-(4.9), and (4.17), we get

$$\mathbf{Z}^{(N)}(t) = \mathbf{Z}^{(N)}(0) + \int_{s=0}^t \frac{\rho}{N} (W(\mathbf{X}^{(N)}(s)) - W(\mathbf{x}(s; \mathbf{u}))) ds + \int_{s=0}^t W_3(\mathbf{X}^{(N)}(s)) ds + \mathbf{M}^{(N)}(t); \quad (4.26)$$

The following result concludes the stochastic boundedness of the process  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  as  $N \rightarrow \infty$  and the stochastic boundedness is useful for proving the tightness of the sequence  $f(\mathbf{Z}^{(N)}(t); t \geq 0)_{N \geq 1}$ .

**Lemma 4.1.** *For any  $T > 0$ , if  $\limsup_{N \rightarrow \infty} \mathbb{E} [\|\mathbf{Z}^{(N)}(0)\|_2^2] < 1$ , then*

$$\limsup_{N \rightarrow \infty} \mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\mathbf{Z}^{(N)}(t)\|_2^2 \right] < 1; \quad (4.27)$$

*Proof.* See Section 4.4.2. □

### 4.3 Summary of Main Results

In this section, we give the main results and provide proofs in Section 4.4. We present the results related to the transient regime and the stationary regime in Sections 4.3.1 and 4.3.2, respectively.

### 4.3.1 Transient Regime

In this section, for the transient regime case, we show that the process  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  converges to an Ornstein-Uhlenbeck (OU) process as  $N \rightarrow \infty$ .

As we show later in this section, the limit of the sequence  $f(\mathbf{Z}^{(N)}(t); t \geq 0)g_{N-1}$  depends on the linearization of (4.7)-(4.9) around a solution  $(\mathbf{r}(t); t \geq 0)$  to (4.7)-(4.9) with an initial point  $\mathbf{r}(0)$  given by

$$\frac{d\mathbf{s}(t)}{dt} = H(\mathbf{r}(t))\mathbf{s}(t); \quad (4.28)$$

where for  $\mathbf{a} \geq \mathbb{U}$  and  $\mathbf{b} \geq \mathbb{V}$ , the linear operator  $H(\mathbf{a}) : \mathbb{V} \rightarrow \mathbb{V}$  is defined as

$$(H(\mathbf{a})\mathbf{b})_n = da_n^{d-1}b_{n-1} - (da_n^{d-1} + n)b_n + nb_{n+1}; \quad (4.29)$$

$n \geq 1$ . The limit of the sequence  $f(\mathbf{Z}^{(N)}(t); t \geq 0)g_{N-1}$  depends on the process  $(\mathbf{s}(t); t \geq 0)$  when  $(\mathbf{r}(t); t \geq 0)$  in (4.28) is replaced with the mean-field  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$ . Note that any solution  $(\mathbf{s}(t); t \geq 0)$  to (4.28) satisfies  $\mathbf{s}(t) = \mathbf{w}(t) - \mathbf{r}(t)$ , where  $(\mathbf{w}(t); t \geq 0)$  is a solution to the equations (4.7)-(4.9). The operator  $H(\mathbf{a})$  is a matrix in the canonical basis  $(0; 1; 0; \dots; 0), (0; 0; 1; 0; \dots; 0), \dots, (0; 0; \dots; 0; 1)$ , where the dimension of each vector is  $C + 1$ . We can write  $H(\mathbf{a})$  as the following matrix of size  $C \times C$ :

$$H(\mathbf{a}) = \begin{bmatrix} -a_1 & 1 & 0 & 0 & \dots & 0 & 0 \\ a_1 & -a_2 & 2 & 0 & \dots & 0 & 0 \\ 0 & a_2 & -a_3 & 3 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & a_{C-2} & a_{C-1} & C-1 \\ 0 & 0 & \dots & 0 & 0 & a_{C-1} & C \end{bmatrix};$$

where  $a_i = da_i^{d-1}$  and  $a_i = a_i + i, 1 \leq i \leq C$ .

Let  $(\mathbf{M}(t); t \geq 0) = f(\mathbf{M}_i(t); t \geq 0)g_{i2\mathbb{U}, 1; \dots; Cg}$  be a collection of mutually independent real valued continuous and centered Gaussian martingales, determined in law by their deterministic quadratic variation process

$$\langle \mathbf{M}_n \rangle_{t=0}^t = \int_{s=0}^t ((W_1(\mathbf{x}(s; \mathbf{u}))_n + (W_2(\mathbf{x}(s; \mathbf{u}))_n) ds; \quad (4.30)$$

$n \geq 0$ . Note that both  $\mathbf{M}(t)$  and  $(\langle \mathbf{M}_i \rangle_{t=0}^t; 0 \leq i \leq C)$  lie in  $\mathbb{V}$ . From (4.30), the martingale  $(\mathbf{M}(t); t \geq 0)$  is square integrable since  $((W_1(\mathbf{b}))_n + (W_2(\mathbf{b}))_n)$  for  $\mathbf{b} \geq \mathbb{U}$  is uniformly bounded in  $n$  and  $\mathbf{b}$  as  $0 \leq b_i \leq 1$ .

We now introduce the stochastic differential equation (SDE) satisfied by the limit of the process  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  as  $N \rightarrow \infty$  in the transient regime.

**Definition 4.1.** *Stochastic Differential Equation (SDE) for the Transient Regime:* Let  $(\mathbf{Z}(t); t \geq 0)$  be a solution to the following SDE

$$\mathbf{Z}(t) = \mathbf{Z}(0) + \int_{s=0}^t H(\mathbf{x}(s; \mathbf{u}))\mathbf{Z}(s) ds - \int_{s=0}^t W_3(\mathbf{x}(s; \mathbf{u})) ds + \mathbf{M}(t); \quad (4.31)$$

A solution to (4.31) is an OU process. We next study the SDE (4.31) below.

**Theorem 4.2.** *We show*

1. For  $\mathbf{a} \in \mathbb{U}$ , the linear operator  $H(\mathbf{a})$  satisfies  $\|H(\mathbf{a})\|_2 < B_H$ , where  $B_H = \sqrt{32(d^2 + C^2)}$ .
2. If  $\mathbb{E}[k\mathbf{Z}(0)k_2^2] < 1$ , then there exists a unique strong solution to (4.31) denoted by  $(\mathbf{Z}(t); t \geq 0)$  that satisfies  $\mathbb{E}[\sup_{t \leq T} k\mathbf{Z}(t)k_2^2] < 1$ .

*Proof.* See Section 4.4.1. □

We now present the main result of this chapter for the transient regime.

**Theorem 4.3.** *If  $\mathbf{Z}^{(N)}(0) \rightarrow \mathbf{Z}(0)$ , then  $(\mathbf{Z}^{(N)}(t); t \geq 0) \rightarrow (\mathbf{Z}(t); t \geq 0)$  where  $(\mathbf{Z}(t); t \geq 0)$  is the unique solution to (4.31) with initial point  $\mathbf{Z}(0)$ .*

*Proof.* See Section 4.4.3. □

**Remark 4.2.** *For a constant  $r \in \mathbb{R}_+$ , if  $\mathbf{u}^{(N)}$  is independent of  $N$  fixed at  $\mathbf{u}^{(N)} = r$ , then  $\mathbf{u} = r$  and  $\tilde{\mathbf{u}} = 0$ . As a result, for  $n \geq 1$ ,*

$$(W_1(\mathbf{b}))_n = r(b_{n-1}^d - b_n^d); (W_2(\mathbf{b}))_n = n(b_n - b_{n+1}); (W_3(\mathbf{b}))_n = 0;$$

We recover the SDE obtained in [55] for  $\mathbf{u}^{(N)} = r$  case from the SDE (4.31) given by

$$\mathbf{Z}(t) = \mathbf{Z}(0) + \int_{s=0}^t H(\mathbf{x}(s; \mathbf{u}))\mathbf{Z}(s) ds + \mathbf{M}(t);$$



### 4.3.2 Stationary Regime

In this section, we present results pertaining to the stationary regime. In the stationary regime, the mean-field is located at  $\bar{\mathbf{x}}$ , and hence we assume  $\mathbf{u} = \bar{\mathbf{x}}$ . We recall that satisfies

$$W(\bar{\mathbf{x}}) = W_1(\bar{\mathbf{x}}) \quad W_2(\bar{\mathbf{x}}) = 0: \quad (4.32)$$

Our objective is to obtain the limit of the process  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  defined in (4.17) as  $N \rightarrow \infty$  in the stationary regime. For this, by studying the process  $(\mathbf{Q}^{(N)}(t); t \geq 0)$  where

$$\mathbf{Q}^{(N)}(t) = \frac{1}{N} \overline{\mathbf{X}^{(N)}(t)}; \quad (4.33)$$

we conclude that the sequence  $f\mathbf{Z}^{(N)}(t)g_{N-1}$  is tight in the stationary regime. We introduce an SDE (4.38) and show that there exists a unique solution to this SDE with a unique invariant law. We then use this result to prove that the limit of the process  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  as  $N \rightarrow \infty$  in the stationary regime is a stationary OU process with the same invariant law as that of the solution to the proposed SDE (4.38).

We next state the exponential stability of the mean-field. We use the following result in the proof of the subsequent result stated in Lemma 4.3.

**Lemma 4.2.** *There exists  $\delta_1 > 0$  and  $D_3 < 1$  such that for all  $\mathbf{u} \in \mathcal{U}$ , the mean-field  $(\mathbf{x}(t; \mathbf{u}); t \geq 0)$  satisfies*

$$\|\mathbf{x}(t; \mathbf{u}) - \bar{\mathbf{x}}\| \leq k_2 e^{-\delta_1 t} D_3 k_1 \|\mathbf{u} - \bar{\mathbf{x}}\|; \quad (4.34)$$

*Proof.* See Section 4.4.7. □

The following result shows the tightness of  $f\mathbf{Z}^{(N)}(t)g_{N-1}$  in the stationary regime.

**Lemma 4.3.** *If  $\limsup_{N \rightarrow \infty} \mathbb{E} [k\mathbf{Q}^{(N)}(0)k_2^2] < 1$ , then*

$$\limsup_{N \rightarrow \infty} \sup_{t \geq 0} \mathbb{E} [k\mathbf{Q}^{(N)}(t)k_2^2] < 1; \quad (4.35)$$

*Consequently, in the stationary regime, we have*

$$\limsup_{N \rightarrow \infty} \mathbb{E} [k\mathbf{Z}^{(N)}(1)k_2^2] < 1; \quad (4.36)$$

*Proof.* See Section 4.4.4. □

We next state the SDE that is used to obtain the limit of the sequence  $f(\mathbf{Z}^{(N)}(t); t)g_{N-1}$  in the stationary regime. For this, we first claim that by linearizing (4.9) around  $\mathbf{x}$ , we get

$$\frac{d\mathbf{s}(t)}{dt} = H(\mathbf{x})\mathbf{s}(t): \quad (4.37)$$

Let  $\mathbf{B}(t) = (\mathbf{B}_i(t); 0 \leq i \leq C)$  with  $\mathbf{B}_0(t) = 0$  where  $f(\mathbf{B}_i(t); t)g_{0 \leq i \leq C}$  are independent centered Brownian motions and  $E[\mathbf{B}_i^2(1)] = V_i = \text{var}(\mathbf{B}_i(1)) = 2n(i-1)_{i+1}$ ,  $i = 1$ . The infinitesimal covariance matrix of  $(\mathbf{B}(t); t \geq 0)$  is diagonal  $\text{diag}(\mathbf{V})$ , where  $\mathbf{V} = (V_i; 0 \leq i \leq C)$ . From (4.30), the martingales  $(\mathbf{M}(t); t \geq 0)$  with  $\mathbf{x}(0; \mathbf{x}) = \mathbf{x}$  has the same law as  $(\mathbf{B}(t); t \geq 0)$ . We now define the following SDE which is used to study the process  $(\mathbf{Z}(t); t \geq 0)$  in the stationary regime.

**Definition 4.2.** *An SDE for the Stationary Regime: The process  $(\mathbf{Q}(t); t \geq 0)$  is a solution to the following SDE,*

$$\mathbf{Q}(t) = \mathbf{Q}(0) + \int_{s=0}^t H(\mathbf{x})\mathbf{Q}(s) ds - \int_{s=0}^t W_3(\mathbf{x}) ds + \mathbf{B}(t): \quad (4.38)$$

Any solution to (4.38) is an OU process.

For an arbitrary  $\mathbf{Q}(0)$  in (4.38), we have the following result and the proof follows by the same arguments as in the proof of Theorem 4.2. Hence, we omit the proof.

**Theorem 4.4.** *We show*

1.  $kH(\mathbf{x})k_2$  is bounded.
2. If  $E[k\mathbf{Q}(0)k_2^2] < 1$ , then  $(\mathbf{Q}(t); t \geq 0)$  given by

$$\mathbf{Q}(t) = e^{H(\mathbf{x})t}\mathbf{Q}(0) - \int_{s=0}^t e^{H(\mathbf{x})(t-s)}W_3(\mathbf{x}) ds + \int_{s=0}^t e^{H(\mathbf{x})(t-s)}d\mathbf{B}(s); \quad (4.39)$$

is the unique strong solution to (4.38). Furthermore,

$$E\left[\sup_{t \leq T} k\mathbf{Q}(t)k_2^2\right] < 1: \quad (4.40)$$

We point out that the transpose  $H(\mathbf{x})^T$  of  $H(\mathbf{x})$  is the generator of a finite state birth-death process and the birth, death, and killing rates in state  $i$  ( $1 \leq i \leq C$ ) are  $\lambda_i$ ,  $\mu_i$ , and  $\nu_i$ , respectively. Let  $I$  be the identity matrix of dimension  $C \times C$ . Then since  $H(\mathbf{x})^T + I$  is the generator of a birth-death process with zero killing rates, all the eigenvalues of  $H(\mathbf{x})^T + I$  are negative [84]. Hence, all the eigenvalues of  $H(\mathbf{x})$  are less than  $-1$ . As a consequence, we have the following result due to the fact that all the eigenvalues are negative.

**Lemma 4.4.** *The unique solution to (4.37) is given by  $(\mathbf{s}(t); t \geq 0)$  where  $\mathbf{s}(t) = e^{H(\cdot)t}\mathbf{s}(0)$ . Furthermore,  $(\mathbf{s}(t); t \geq 0)$  satisfies that for some  $\epsilon_2 > 0$  and  $D_4 < 1$ ,*

$$\|\mathbf{s}(t)\|_2 \leq e^{-\epsilon_2 t} D_4 \|\mathbf{s}(0)\|_2 \quad (4.41)$$

From Lemma 4.4 and the unique solution given in Theorem 4.4, the following result follows immediately. Hence, we omit the proof.

**Theorem 4.5.** *The unique solution to (4.38) as  $t \rightarrow \infty$  has the invariant law coinciding with the law of a stationary Gaussian process with mean  $\int_0^1 e^{H(\cdot)s} W_3(\cdot) ds$  and covariance matrix  $\int_0^1 e^{H(\cdot)s} \text{diag}(\mathbf{V}) e^{H(\cdot)s} ds$ .*

We are now ready to state the main result on the FCLT for the stationary regime.

**Theorem 4.6.** *Under the assumption that the system with index  $N$  is in the stationary regime, from Theorem 4.5, we prove that the sequence  $f(\mathbf{Z}^{(N)}(t); t \geq 0)g_{N-1}$  as  $N \rightarrow \infty$  converges in law to the unique stationary OU process which solves (4.38). Furthermore, the limit of the sequence  $f(\mathbf{Z}^{(N)}(0))g_{N-1}$  in the stationary regime has the same law as the invariant law of the solution to (4.38).*

*Proof.* See Section 4.4.5. □

We now use Theorem 4.6 to provide an approximation to the average blocking probability in the system with  $N$  servers.

**Theorem 4.7.** *Let  $P_{block}^{(N)}$  be the average blocking probability in the system with  $N$  servers, then*

$$P_{block}^{(N)} = \frac{d}{c} - \frac{1}{N} \left( \sum_{i=0}^c i(i-i+1) \right) - \frac{\tilde{P}}{N} (1 - \frac{d}{c}) + o(N^{-\frac{1}{2}}); \quad (4.42)$$

where the vector  $\tilde{P} = (i; 0 \leq i \leq c) = \int_0^1 e^{H(\cdot)s} W_3(\cdot) ds$  is the mean of the unique solution to (4.38) in the stationary regime.

*Proof.* See Section 4.4.6. □

**Remark 4.3.** *From Theorem 4.7, we have*

1. *If  $d = c$ , the result (4.42) corresponds to the Halfin-Whitt regime. In this case, we have*

$$\lim_{N \rightarrow \infty} \frac{P_{block}^{(N)}}{N} (P_{block}^{(N)} - \frac{d}{c}) = \frac{1}{N} \left( \sum_{i=0}^c i(i-i+1) \right) - \frac{\tilde{P}}{N} (1 - \frac{d}{c});$$

2. If  $\tilde{\rho} = 0$ , then  $W_3(\cdot) = 0$ . As a result,

$$\lim_{N \rightarrow \infty} \overline{P}_{block}^{(N)}(\tilde{\rho}, C) = 0:$$

We presented the results for the case  $\tilde{\rho} = 0$  in [55]. For this case, we can also apply the results of [51] to conclude that  $\overline{P}_{block}^{(N)}(\tilde{\rho}, C) = O(\frac{1}{N})$ . In [51], an FCLT result was not studied..

The significance of Theorem 4.7 is that although the exact blocking formula for  $\overline{P}_{block}^{(N)}$  is not known, and also it is difficult to characterize due to complex interactions between servers, but as  $N$  becomes large we can compute approximations to the blocking probability as a function of  $\tilde{\rho}$ ,  $\tilde{\rho}$ ,  $N$ ,  $\tilde{\rho}$ , and  $C$ .

## 4.4 Proofs of Main Results

### 4.4.1 Proof of Theorem 4.2

For  $\mathbf{y} \in \mathcal{V}$ , we have

$$\|kH(\mathbf{a})\mathbf{y}\|_2^2 = \sum_{n=0}^C \left| \frac{d^n}{n!} y_{n-1} - \left( \frac{d^n}{n!} + n \right) y_n + n y_{n+1} \right|^2:$$

Then it can be checked that  $\|kH(\mathbf{a})\mathbf{y}\|_2^2 \leq 32(\tilde{\rho}^2 + C^2)\|\mathbf{y}\|_2^2$  as  $0 \leq a_i \leq 1$ . This implies  $\|kH(\mathbf{a})\|_2^2 \leq B_H^2 = 32(\tilde{\rho}^2 + C^2)$  and hence,  $\|kH(\mathbf{a})\|_2 \leq B_H$ .

We next show the uniqueness of a solution to (4.31). For two solutions  $(\mathbf{Z}^1(t); t \geq 0)$  and  $(\mathbf{Z}^2(t); t \geq 0)$  with initial points  $\mathbf{Z}^1(0)$  and  $\mathbf{Z}^2(0)$ , we obtain

$$\|\mathbf{Z}^1(t) - \mathbf{Z}^2(t)\|_2 = \|\mathbf{Z}^1(0) - \mathbf{Z}^2(0)\|_2 + \int_{s=0}^t \|H(\mathbf{x}(s; \mathbf{u}))(\mathbf{Z}^1(s) - \mathbf{Z}^2(s))\|_2 ds:$$

Then

$$\|\mathbf{Z}^1(t) - \mathbf{Z}^2(t)\|_2 \leq \|\mathbf{Z}^1(0) - \mathbf{Z}^2(0)\|_2 + B_H \int_{s=0}^t \|\mathbf{Z}^1(s) - \mathbf{Z}^2(s)\|_2 ds:$$

By using the Gronwall's Lemma (Theorem A.7) [56, Page 498], we get

$$\|\mathbf{Z}^1(t) - \mathbf{Z}^2(t)\|_2 \leq e^{B_H t} \|\mathbf{Z}^1(0) - \mathbf{Z}^2(0)\|_2:$$

Hence,

$$\mathbb{E} [\|\mathbf{Z}^1(t) - \mathbf{Z}^2(t)\|_2^2] \leq e^{2B_H t} \mathbb{E} [\|\mathbf{Z}^1(0) - \mathbf{Z}^2(0)\|_2^2]:$$

If  $\mathbf{Z}^1(0) = \mathbf{Z}^2(0)$ , then we obtain  $\mathbf{Z}^1(t) = \mathbf{Z}^2(t)$  a.s. for all rational  $t$ . Finally, since  $(\mathbf{Z}^1(t); t \geq 0)$  and  $(\mathbf{Z}^2(t); t \geq 0)$  have continuous sample paths, we conclude that  $\mathbf{Z}^1(t) = \mathbf{Z}^2(t)$  a.s.; for all  $t \geq 0$ .

We next use the Gronwall's lemma and the Doob's inequality (Theorem A.8) [56, p. 63] to show that  $\mathbb{E} \left[ \sup_{0 \leq t \leq T} k\mathbf{Z}(t)k_2^2 \right] < 1$ . From (4.31), we have

$$k\mathbf{Z}(t)k_2 = k\mathbf{Z}(0)k_2 + B_H \int_{s=0}^t k\mathbf{Z}(s)k_2 ds + \int_{s=0}^t kW_3(\mathbf{x}(s; \mathbf{u}))k_2 + k\mathbf{M}(t)k_2:$$

However, for  $\mathbf{a} \in \mathcal{U}$ , we have  $kW_3(\mathbf{a})k_2 \leq 2j^* j^{\rho} \bar{C}$ . Therefore, we get

$$k\mathbf{Z}(t)k_2^2 \leq 4 \left( k\mathbf{Z}(0)k_2^2 + B_H^2 \left( \int_{s=0}^t k\mathbf{Z}(s)k_2 ds \right)^2 + 4t^{2-\alpha} C + k\mathbf{M}(t)k_2^2 \right):$$

We can apply the Holder's inequality to  $\left( \int_{s=0}^t k\mathbf{Z}(s)k_2 ds \right)^2$  to obtain the following inequality

$$k\mathbf{Z}(t)k_2^2 \leq 4 \left( k\mathbf{Z}(0)k_2^2 + B_H^2 t \int_{s=0}^t k\mathbf{Z}(s)k_2^2 ds + 4t^{2-\alpha} C + k\mathbf{M}(t)k_2^2 \right):$$

By using the Gronwall's Lemma, we obtain

$$k\mathbf{Z}(t)k_2^2 \leq 4 \left( k\mathbf{Z}(0)k_2^2 + 4t^{2-\alpha} C + k\mathbf{M}(t)k_2^2 \right) e^{4B_H^2 t^\alpha}:$$

Therefore for any  $T > 0$ , we get

$$\sup_{0 \leq t \leq T} k\mathbf{Z}(t)k_2^2 \leq 4 \left( k\mathbf{Z}(0)k_2^2 + 4T^{2-\alpha} C + \sup_{0 \leq t \leq T} k\mathbf{M}(t)k_2^2 \right) e^{4B_H^2 T^\alpha}:$$

Hence,

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} k\mathbf{Z}(t)k_2^2 \right] \leq 4e^{4B_H^2 T^\alpha} \left( \mathbb{E} [k\mathbf{Z}(0)k_2^2] + 4T^{2-\alpha} C + \mathbb{E} \left[ \sup_{0 \leq t \leq T} k\mathbf{M}(t)k_2^2 \right] \right):$$

The Doob's inequality implies

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} k\mathbf{Z}(t)k_2^2 \right] \leq 4e^{4B_H^2 T^\alpha} \left( \mathbb{E} [k\mathbf{Z}(0)k_2^2] + 4T^{2-\alpha} C + 4\mathbb{E} [k\mathbf{M}(T)k_2^2] \right):$$

Therefore

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} k\mathbf{z}(t)k_2^2 \right] \leq 4e^{4B_H^2 T^\alpha} \left( \mathbb{E} [k\mathbf{z}(0)k_2^2] + 4T^{2-\alpha} C + 4\mathbb{E} \left[ \sum_{i=1}^C \langle \mathbf{M}_i \rangle_T \right] \right):$$

For every  $T$  and  $i$ ,  $\langle \mathbf{M}_i \rangle_T$  is bounded. Hence,  $\mathbb{E} \left[ \sup_{0 \leq t \leq T} k\mathbf{Z}(t)k_2^2 \right] < 1$ .

The fact that there exists a solution to the SDE (4.31) follows from [85, page 354]. Let  $(\tilde{F}(t); t \geq 0)$  be the unique solution of the following equation

$$\frac{dY(t)}{dt} = H(\mathbf{x}(t; \mathbf{u}))Y(t); \quad Y(0) = I;$$

where  $I$  is the identity matrix. Then the unique solution to the SDE (4.31) is given by

$$\mathbf{Z}(t) = \tilde{F}(t) \left[ \mathbf{Z}(0) - \int_{s=0}^t (\tilde{F}(s))^{-1} W_3(\mathbf{x}(s; \mathbf{u})) ds + \int_{s=0}^t (\tilde{F}(s))^{-1} d\mathbf{M}(s) \right]; \quad (4.43)$$

#### 4.4.2 Proof of Lemma 4.1

We are given

$$\mathbf{Z}^{(N)}(t) = \mathbf{Z}^{(N)}(0) + \int_{s=0}^t \frac{1}{N} (W(\mathbf{X}^{(N)}(s)) - W(\mathbf{x}(s; \mathbf{u}))) ds + \int_{s=0}^t W_3(\mathbf{X}^{(N)}(s)) ds + \mathbf{M}^{(N)}(t); \quad (4.44)$$

and for  $n \geq 1$ ,

$$\langle \mathbf{M}_n^{(N)} \rangle_t = \int_{s=0}^t ((W_1(\mathbf{X}^{(N)}(s)))_n + (W_2(\mathbf{X}^{(N)}(s)))_n - \frac{1}{N} (W_3(\mathbf{X}^{(N)}(s)))_n) ds;$$

From (4.44), we obtain

$$\|\mathbf{Z}^{(N)}(t)\|_2 \leq \|\mathbf{Z}^{(N)}(0)\|_2 + B_W \int_{s=0}^t \|\mathbf{Z}^{(N)}(s)\|_2 ds + 2C\sqrt{t} + \|\mathbf{M}^{(N)}(t)\|_2;$$

By using the Gronwall's Lemma,

$$\|\mathbf{Z}^{(N)}(t)\|_2 \leq (\|\mathbf{Z}^{(N)}(0)\|_2 + 2C\sqrt{t} + \|\mathbf{M}^{(N)}(t)\|_2) e^{B_W t};$$

As a result, we get

$$\|\mathbf{Z}^{(N)}(t)\|_2^2 \leq 3(\|\mathbf{Z}^{(N)}(0)\|_2^2 + 4C^2 t + \|\mathbf{M}^{(N)}(t)\|_2^2) e^{2B_W t};$$

For  $T > 0$ , we have

$$\sup_{0 \leq t \leq T} \|\mathbf{Z}^{(N)}(t)\|_2^2 \leq 3(\|\mathbf{Z}^{(N)}(0)\|_2^2 + 4C^2 T^2 + \sup_{0 \leq t \leq T} \|\mathbf{M}^{(N)}(t)\|_2^2) e^{2B_W T};$$

Finally, the Doob's inequality implies the following inequality

$$\mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\mathbf{Z}^{(N)}(t)\|_2^2 \right] \leq 3e^{2B_W T} \left( \mathbb{E} [\|\mathbf{Z}^{(N)}(0)\|_2^2] + 4C^2 T^2 + 4\mathbb{E} \left[ \sum_{i=1}^C \langle \mathbf{M}_i^{(N)} \rangle_T \right] \right);$$

Since  $\sup_{N \geq 1} \mathbb{E} \left[ \sum_{i=1}^C \langle \mathbf{M}_i^{(N)} \rangle_T \right] < 1$ , we conclude that

$$\limsup_{N \rightarrow \infty} \mathbb{E} \left[ \sup_{0 \leq t \leq T} \|\mathbf{Z}^{(N)}(t)\|_2^2 \right] < 1;$$

### 4.4.3 Proof of Theorem 4.3

We recall that since the space  $\mathbb{V}$  is Polish, the space of càdlàg functions under the Skorohod topology is a Polish space [56, Theorem 5.6, p.121]. Hence from the Prohorov's theorem (Theorem A.2) [56], tightness is equivalent to relative compactness. Therefore we first need to show the tightness and then we need to show that every limiting point has the same law as the unique OU process with the initial point  $\mathbf{Z}(0)$ .

To show the tightness of  $(\mathbf{Z}^{(N)}(t); t \geq 0)$ , we need to show that Theorem 4.1 of [56, page 354] holds. For this, we first establish several useful results.

Since  $\mathbf{Z}^{(N)}(0) \rightarrow \mathbf{Z}(0)$ , it implies that the sequence  $\mathbf{Z}^{(N)}(0)$  is tight. Let  $\tilde{B}(r)$  be the closed ball with radius  $r$  centered at  $\mathbf{0}$ . For every  $\epsilon > 0$ , there exists  $r < 1$  such that  $\mathbb{P}(\mathbf{Z}^{(N)}(0) \in \tilde{B}(r)) > 1 - \epsilon$  for all  $N \geq 1$ . We now define a random variable  $\mathbf{X}^{(N;\cdot)}(0)$  such that it coincides with  $\mathbf{X}^{(N)}(0)$  on  $\{ \mathbf{Z}^{(N)}(0) \in \tilde{B}(r) \}$  and  $\mathbf{X}^{(N;\cdot)}(0)$  is uniformly bounded in  $N$  on  $\{ \mathbf{Z}^{(N)}(0) \notin \tilde{B}(r) \}$ . Then by using coupling arguments, the processes  $(\mathbf{Z}^{(N;\cdot)}(t); t \geq 0)$  and  $(\mathbf{Z}^{(N)}(t); t \geq 0)$  coincide on  $\{ \mathbf{Z}^{(N)}(0) \in \tilde{B}(r) \}$ . Hence, without loss of generality, we assume that  $\mathbf{Z}^{(N)}(0)$  is uniformly bounded in  $N$ . As a consequence, the result stated in Lemma 4.1 can be used in the rest of the proof.

We next recall the following useful result from [47, Lemma 3.3]. For  $a$  and  $h$  in  $\mathbb{R}$ , let

$$\hat{B}(a; h) = (a + h)^d - a^d - da^{d-1}h.$$

Then if both  $a$  and  $a + h$  lie in  $[0; 1]$ , we have

$$0 \leq \hat{B}(a; h) \leq h^d + (2^d - d - 2)ah^2. \quad (4.45)$$

We now define a mapping  $\hat{G} : \mathbb{U} \times \mathbb{V} \rightarrow \mathbb{V}$  as follows: for  $\mathbf{r} \in \mathbb{U}$  and  $\mathbf{y} \in \mathbb{V}$ ,

$$(\hat{G}(\mathbf{r}; \mathbf{y}))_n = \hat{B}(r_{n-1}; y_{n-1}) - \hat{B}(r_n; y_n). \quad (4.46)$$

Then if  $\mathbf{r} + \mathbf{y} \in \mathbb{U}$ , we have

$$W(\mathbf{r} + \mathbf{y}) - W(\mathbf{r}) = H(\mathbf{r})\mathbf{y} + \hat{G}(\mathbf{r}; \mathbf{y}). \quad (4.47)$$

Note that since  $\mathbf{Z}^{(N)}(t) = \frac{\rho}{N} \mathbf{X}^{(N)}(t) - \mathbf{x}(t; \mathbf{u})$ , we have

$$\mathbf{X}^{(N)}(t) = \mathbf{x}(t; \mathbf{u}) + \frac{\mathbf{Z}^{(N)}(t)}{\frac{\rho}{N}}. \quad (4.48)$$

Here,  $\mathbf{X}^{(N)}(t)$ ,  $\mathbf{x}(t; \mathbf{u}) \in \mathcal{U}$  and  $\frac{\mathbf{Z}^{(N)}(t)}{\rho \overline{N}} \in \mathcal{V}$ . Hence, from (4.47), we have

$$W(\mathbf{X}^{(N)}(t)) - W(\mathbf{x}(t; \mathbf{u})) = H(\mathbf{x}(t; \mathbf{u})) \frac{\mathbf{Z}^{(N)}(t)}{\rho \overline{N}} + \hat{G}\left(\mathbf{x}(t; \mathbf{u}); \frac{\mathbf{Z}^{(N)}(t)}{\rho \overline{N}}\right): \quad (4.49)$$

We now show that the conditions of [56, Theorem 4.1, page 354] are satisfied. Let us write  $\mathbf{B}_n(t)$  and  $\mathbf{A}_n^{ij}(t)$  of [56, Theorem 4.1, page 354] as  $\mathbf{D}_{(N)}(t)$  and  $\mathbf{A}_{(N)}^{ij}(t)$ , respectively. Then from (4.26)

$$\mathbf{D}_{(N)}(t) = \int_{s=0}^t \rho \overline{N} (W(\mathbf{X}^{(N)}(s)) - W(\mathbf{x}(s; \mathbf{u}))) ds - \int_{s=0}^t W_3(\mathbf{X}^{(N)}(s)) ds: \quad (4.50)$$

From (4.49), we can write

$$\begin{aligned} \mathbf{D}_{(N)}(t) &= \int_{s=0}^t H(\mathbf{x}(s; \mathbf{u})) \mathbf{Z}^{(N)}(s) ds + \int_{s=0}^t \rho \overline{N} \hat{G}\left(\mathbf{x}(s; \mathbf{u}); \frac{\mathbf{Z}^{(N)}(s)}{\rho \overline{N}}\right) ds - \int_{s=0}^t W_3(\mathbf{X}^{(N)}(s)) ds: \end{aligned} \quad (4.51)$$

Also, from (4.25)

$$\mathbf{A}_{(N)}^{ij}(t) = \int_{s=0}^t \left( (W_1(\mathbf{X}^{(N)}(s)))_i + (W_2(\mathbf{X}^{(N)}(s)))_i - \frac{1}{\overline{N}} (W_3(\mathbf{X}^{(N)}(s)))_i \right) ds: \quad (4.52)$$

and  $\mathbf{A}_{(N)}^{ij}(t) = 0$  for  $i \neq j$ . Let us define

$$\mathbf{D}(t) = \int_{s=0}^t H(\mathbf{x}(s; \mathbf{u})) \mathbf{Z}(s) ds - \int_{s=0}^t W_3(\mathbf{x}(s; \mathbf{u})) ds: \quad (4.53)$$

and

$$A^{ii}(t) = \int_{s=0}^t ((W_1(\mathbf{x}(s; \mathbf{u})))_i + (W_2(\mathbf{x}(s; \mathbf{u})))_i) ds: \quad (4.54)$$

with  $A^{ij}(t) = 0$  for  $i \neq j$ .

Since  $(\mathbf{Z}_i^{(N)}(t); t \geq 0)$  has jumps of size  $\frac{1}{\overline{N}}$  and from the continuity of  $\mathbf{D}_{(N)}(t)$  and  $\mathbf{A}_{(N)}^{ij}(t)$  in  $t$ , the conditions (4.3)-(4.5) of [56, Theorem 4.1, p. 354] are valid. From the condition (4.6) of [56, Theorem 4.1, p. 354], we need to show that for  $T > 0$ ,  $\mathbf{V}_{(N)}(t) \rightarrow 0$  in probability where  $\mathbf{V}_{(N)}(t) = \mathbf{D}_{(N)}(t) - \mathbf{D}(t)$  given by

$$\begin{aligned} \mathbf{V}_{(N)}(t) &= \left( \int_{s=0}^t H(\mathbf{x}(s; \mathbf{u})) \mathbf{Z}^{(N)}(s) ds + \int_{s=0}^t \rho \overline{N} \hat{G}\left(\mathbf{x}(s; \mathbf{u}); \frac{\mathbf{Z}^{(N)}(s)}{\rho \overline{N}}\right) ds - \int_{s=0}^t W_3(\mathbf{X}^{(N)}(s)) ds \right) \end{aligned}$$



$$\left( \int_{s=0}^t H(\mathbf{x}(s; \mathbf{u})) \mathbf{Z}^{(N)}(s) ds - \int_{s=0}^t W_3(\mathbf{x}(s; \mathbf{u})) ds \right): \quad (4.55)$$

From (4.27) and (4.45), we have  $\int_{s=0}^T \rho \overline{N} \hat{G} \left( \mathbf{x}(s; \mathbf{u}); \frac{\mathbf{z}^{(N)}(s)}{\rho \overline{N}} \right) ds \xrightarrow{P} 0$  in probability. From the existence of the mean-field limit,  $\int_{s=0}^T W_3(\mathbf{X}^{(N)}(s)) ds \xrightarrow{P} \int_{s=0}^T W_3(\mathbf{x}(s; \mathbf{u})) ds \xrightarrow{P} 0$  in probability. We conclude that the condition (4.6) of [56, Theorem 4.1, p. 354] is also true. The condition (4.7) of [56, Theorem 4.1, p. 354] follows as  $\mathbf{A}_{(N)}^{i,j}(T) \xrightarrow{P} A^{i,j}(T) \xrightarrow{P} 0$  in probability from the existence of the mean-field limit. From the fact that the SDE (4.31) has a unique solution, the sequence  $(\mathbf{Z}^{(N)}(t); t \in [0, T])_{N \geq 1}$  converges to the unique solution of the SDE (4.31) as  $N \rightarrow \infty$ .

#### 4.4.4 Proof of Lemma 4.3

The proof is based on Lemma 4.2. Although our main motivation is to establish the convergence of  $\rho \mathbf{Z}^{(N)}(t)_{N \geq 1}$  as  $N \rightarrow \infty$  in the stationary regime, we prove the tightness of the stationary sequence  $\rho \mathbf{Z}^{(N)}(t)_{N \geq 1}$  by studying an alternative process  $(\mathbf{Q}^{(N)}(t); t \in [0, T])$  with the help of Lemma 4.2, where

$$\mathbf{Q}^{(N)}(t) = \rho \overline{N}(\mathbf{X}^{(N)}(t) - \mathbf{y}(h; \mathbf{X}^{(N)}(t_0))):$$

Let us write the solution to the mean-field equation (4.9) at time  $h$  with the initial point  $\mathbf{v}$  as  $\mathbf{y}(h; \mathbf{v})$ . We have

$$\mathbf{y}(h; \mathbf{v}) = \mathbf{v} + \int_{s=0}^h W(\mathbf{y}(s; \mathbf{v})) ds: \quad (4.56)$$

Also, the process  $(\mathbf{X}^{(N)}(t); t \in [0, T])$  satisfies

$$\mathbf{X}^{(N)}(t) = \mathbf{X}^{(N)}(0) + \int_{s=0}^t W(\mathbf{X}^{(N)}(s)) ds - \int_{s=0}^t \frac{1}{\rho \overline{N}} W_3(\mathbf{X}^{(N)}(s)) ds + \frac{\mathbf{M}^{(N)}(t)}{\rho \overline{N}}: \quad (4.57)$$

For  $t_0 = 0$ , we obtain

$$\mathbf{Q}^{(N)}(t_0 + h) = \rho \overline{N}(\mathbf{X}^{(N)}(t_0 + h) - \mathbf{y}(h; \mathbf{X}^{(N)}(t_0))) + \rho \overline{N}(\mathbf{y}(h; \mathbf{X}^{(N)}(t_0)) - \mathbf{X}^{(N)}(t_0)):$$

By defining  $\mathbf{Z}^{(N)}(t_0; h) = \rho \overline{N}(\mathbf{X}^{(N)}(t_0 + h) - \mathbf{y}(h; \mathbf{X}^{(N)}(t_0)))$ , we obtain

$$\mathbf{Q}^{(N)}(t_0 + h) = \mathbf{Z}^{(N)}(t_0; h) + \rho \overline{N}(\mathbf{y}(h; \mathbf{X}^{(N)}(t_0)) - \mathbf{X}^{(N)}(t_0)): \quad (4.58)$$

Also, from (4.56),

$$\mathbf{Z}^{(N)}(t_0; h) = \rho \overline{N}(\mathbf{X}^{(N)}(t_0 + h) - \mathbf{X}^{(N)}(t_0)) - \rho \overline{N} \int_{s=0}^h W(\mathbf{y}(s; \mathbf{X}^{(N)}(t_0))) ds: \quad (4.59)$$

From (4.57), we have

$$\mathbf{X}^{(N)}(t_0 + h) - \mathbf{X}^{(N)}(t_0) = \int_{s=t_0}^{t_0+h} W(\mathbf{X}^{(N)}(s)) ds + \frac{\mathbf{M}^{(N)}(t_0 + h) - \mathbf{M}^{(N)}(t_0)}{\rho \bar{N}} + \frac{1}{\rho \bar{N}} \int_{s=t_0}^{t_0+h} W_3(\mathbf{X}^{(N)}(s)) ds: \quad (4.60)$$

Then by using (4.59) and (4.60), we get

$$\mathbf{Z}^{(N)}(t_0; h) = \rho \bar{N} \int_{s=0}^h W(\mathbf{X}^{(N)}(t_0 + s)) ds - \rho \bar{N} \int_{s=0}^h W(\mathbf{y}(s; \mathbf{x}^{(N)}(t_0))) ds + (\mathbf{M}^{(N)}(t_0 + h) - \mathbf{M}^{(N)}(t_0)) \int_{s=t_0}^{t_0+h} W_3(\mathbf{X}^{(N)}(s)) ds:$$

After simplifications, we get

$$k\mathbf{Z}^{(N)}(t_0; h)k_2 = \rho \bar{N} \int_{s=0}^h kW(\mathbf{X}^{(N)}(t_0 + s)) - W(\mathbf{y}(s; \mathbf{X}^{(N)}(t_0)))k_2 ds + k\mathbf{M}^{(N)}(t_0 + h) - \mathbf{M}^{(N)}(t_0)k_2 + j^- j^+ d^{\rho \bar{N}} Ch:$$

Hence, we obtain

$$k\mathbf{Z}^{(N)}(t_0; h)k_2 \leq B_W \int_{s=0}^h k\mathbf{Z}^{(N)}(t_0; s)k_2 ds + j^- j^+ d^{\rho \bar{N}} Ch + k\mathbf{M}^{(N)}(t_0 + h) - \mathbf{M}^{(N)}(t_0)k_2:$$

For any  $T > 0$ , the Gronwall's inequality implies that there exists a constant  $S_T$  such that

$$\sup_{0 \leq h \leq T} k\mathbf{Z}^{(N)}(t_0; h)k_2 \leq S_T(S_T + \sup_{0 \leq h \leq T} k\mathbf{M}^{(N)}(t_0 + h) - \mathbf{M}^{(N)}(t_0)k_2): \quad (4.61)$$

By using Lemma 4.2, from (4.58), we get

$$k\mathbf{Q}^{(N)}(t_0 + h)k_2 \leq k\mathbf{Z}^{(N)}(t_0; h)k_2 + e^{-\rho h} D_3 k\mathbf{Q}^{(N)}(t_0)k_2: \quad (4.62)$$

From (4.61) and (4.62), we obtain

$$k\mathbf{Q}^{(N)}(t_0 + h)k_2 \leq S_T(S_T + \sup_{0 \leq h \leq T} k\mathbf{M}^{(N)}(t_0 + h) - \mathbf{M}^{(N)}(t_0)k_2) + e^{-\rho h} D_3 k\mathbf{Q}^{(N)}(t_0)k_2:$$

As a result, we can find some constant  $L_T$  as a function of  $T$  such that for  $0 \leq h \leq T$ , we have

$$\mathbb{E} [k\mathbf{Q}^{(N)}(t_0 + h)k_2^2] \leq L_T + 3e^{-2\rho h} D_3^2 \mathbb{E} [k\mathbf{Q}^{(N)}(t_0)k_2^2]: \quad (4.63)$$

We now select a large value of  $T$  such that  $3e^{-2\rho T} D_3^2 < 1$ . Then for all  $N \geq 1$  and an integer  $m$ , we get

$$\mathbb{E} [k\mathbf{Q}^{(N)}((m+1)T)k_2^2] \leq L_T + \mathbb{E} [k\mathbf{Q}^{(N)}(mT)k_2^2]:$$

By using the induction method, we obtain

$$\begin{aligned} \mathbb{E} [k\mathbf{Q}^{(N)}(mT)k_2^2] &= L_T \left( \sum_{j=1}^m j^{-1} \right) + {}^m\mathbb{E} [k\mathbf{Q}^{(N)}(0)k_2^2] \\ &= \frac{L_T}{1} + \mathbb{E} [k\mathbf{Q}^{(N)}(0)k_2^2] : \end{aligned} \quad (4.64)$$

However, from (4.63),

$$\sup_{0 \leq h \leq T} \mathbb{E} [k\mathbf{Q}^{(N)}(mT+h)k_2^2] \leq L_T + 3D_3^2 \mathbb{E} [k\mathbf{Q}^{(N)}(mT)k_2^2] :$$

As a consequence, (4.64) implies

$$\sup_{0 \leq h \leq T} \mathbb{E} [k\mathbf{Q}^{(N)}(mT+h)k_2^2] \leq L_T + 3D_3^2 \left( \frac{L_T}{1} + \mathbb{E} [k\mathbf{Q}^{(N)}(0)k_2^2] \right) :$$

Since  $m$  is arbitrary, we conclude

$$\sup_{t \geq 0} \mathbb{E} [k\mathbf{Q}^{(N)}(t)k_2^2] \leq L_T + 3D_3^2 \left( \frac{L_T}{1} + \mathbb{E} [k\mathbf{Q}^{(N)}(0)k_2^2] \right) :$$

From ergodicity and the fatou's Lemma [56, p. 492], the stationary random variable  $\mathbf{Z}^{(N)}(1)$  satisfies

$$\begin{aligned} \mathbb{E} [k\mathbf{Z}^{(N)}(1)k_2^2] &= \liminf_{t \rightarrow 0} \mathbb{E} [k\mathbf{Q}^{(N)}(t)k_2^2] \\ &= \sup_{t \geq 0} \mathbb{E} [k\mathbf{Q}^{(N)}(t)k_2^2] : \end{aligned}$$

Finally, to show that  $\limsup_{N \rightarrow \infty} \mathbb{E} [k\mathbf{Z}^{(N)}(1)k_2^2] < 1$ , we need to find an  $\mathbf{X}^{(N)}(0)$  such that  $\limsup_{N \rightarrow \infty} \mathbb{E} [k\mathbf{Q}^{(N)}(0)k_2^2] < 1$ . For  $n \geq 1$ , if we select  $\mathbf{X}_n^{(N)}(0) = \frac{j}{N}$  so that  $\frac{1}{2N} \leq n \leq \frac{j}{N} \leq \frac{1}{2N}$ , then  $\mathbb{E} [k\mathbf{Q}^{(N)}(0)k_2^2] < \frac{C}{4N}$ . Hence,  $\limsup_{N \rightarrow \infty} \mathbb{E} [k\mathbf{Q}^{(N)}(0)k_2^2] = 0$ . This completes the proof.

#### 4.4.5 Proof of Theorem 4.6

From Lemma 4.3 and the Markov inequality, the sequence  $f\mathbf{Z}^{(N)}(0)g_{N-1}$  is tight. As a result, from the Prohorov theorem [56, Page 104], the sequence  $f\mathbf{Z}^{(N)}(0)g_{N-1}$  is relatively compact. Consider a converging subsequence and let  $\mathbf{Z}^{(1)}(0)$  be its limiting point, which is square integrable. Then from Theorem 4.3, the considered converging subsequence converges in law to the unique OU process  $(\mathbf{Z}^{(1)}(t); t \geq 0)$  satisfying the SDE (4.38) with initial point  $\mathbf{Z}^{(1)}(0)$ . But, we know from [56, Lemma 7.7 and Theorem 7.8, page

131] that the limit of a sequence of stationary processes is stationary. Hence the law of  $(\mathbf{Z}^{(1)}(t); t \geq 0)$  should be the unique law of the stationary OU process solving the SDE (4.38). This argument applies for every converging subsequence. Hence, the sequence  $f(\mathbf{Z}^{(N)}(t); t \geq 0)g_{N-1}$  in the stationary regime converges to the unique stationary OU process solving the SDE (4.38). This completes the proof.

#### 4.4.6 Proof of Theorem 4.7

The proof is based on Little's law [76, Theorem 4.1], Theorem 4.5, and Theorem 4.6. Let us consider a random variable  $\tilde{\mathbf{S}}^{(N)}(\tau)$  which denotes the number of progressing jobs in the system in stationary. From Little's law, we obtain

$$(N^{-1})P_{block}^{(N)} = \mathbb{E} \left[ \tilde{\mathbf{S}}^{(N)}(\tau) \right].$$

We now consider a random variable  $\mathbf{X}^{(N)}(\tau) = (\mathbf{X}_n^{(N)}(\tau); 0 \leq n \leq C)$ , where  $\mathbf{X}_n^{(N)}(\tau)$  is the fraction of servers with at least  $n$  progressing jobs in stationary. Then we can write

$$\tilde{\mathbf{S}}^{(N)}(\tau) = \sum_{n=0}^C N n (\mathbf{X}_n^{(N)}(\tau) - \mathbf{X}_{n+1}^{(N)}(\tau)).$$

Hence, we have

$$(N^{-1})P_{block}^{(N)} = \mathbb{E} \left[ \sum_{n=0}^C N n (\mathbf{X}_n^{(N)}(\tau) - \mathbf{X}_{n+1}^{(N)}(\tau)) \right].$$

Therefore, we obtain

$$(N^{-1})P_{block}^{(N)} = \sum_{n=0}^C n \mathbb{E} \left[ \mathbf{X}_n^{(N)}(\tau) - \mathbf{X}_{n+1}^{(N)}(\tau) \right]. \quad (4.65)$$

Also, from Theorem 4.5, Lemma 4.3, and Theorem 4.6, since the diffusion limit in stationary regime has the mean vector  $\bar{\mathbf{x}} = \int_0^1 e^{H(\cdot)s} W_3(\cdot) ds$  and  $\limsup_{N \rightarrow \infty} \mathbb{E} [k \mathbf{Z}^{(N)}(\tau) k^2] < \infty$ , we have

$$\lim_{N \rightarrow \infty} \rho_{\bar{N}} \left( \mathbb{E} \left[ \mathbf{X}^{(N)}(\tau) \right] \right) = \bar{\mathbf{x}}.$$

Therefore

$$\mathbb{E} \left[ \mathbf{X}^{(N)}(\tau) \right] = \rho_{\bar{N}} + o(N^{-\frac{1}{2}}). \quad (4.66)$$

Then from (4.65) and (4.66),

$$(N^{-1})P_{block}^{(N)} = \sum_{n=0}^C n (\bar{x}_n - \bar{x}_{n+1}) + \rho_{\bar{N}} \sum_{n=0}^C n (\bar{x}_n - \bar{x}_{n+1}) + o(N^{-\frac{1}{2}}). \quad (4.67)$$

But, from the stationary mean-field equations, the fixed-point satisfies

$$0 = 1 \text{ and } \sum_{n=0}^C n(n - n_{+1}) = (1 - \frac{d}{C}): \quad (4.68)$$

Then from (4.67) and (4.68),

$$1 - P_{block}^{(N)} = \frac{1}{(N)} \left[ (1 - \frac{d}{C}) + \frac{1}{N} \sum_{n=0}^C n(n - n_{+1}) \right] + o(N^{-\frac{1}{2}}):$$

However, by using  $\tilde{P}^{(N)} = \frac{\tilde{p}}{N}$  and the fact that  $\tilde{p} < 1$ , we obtain

$$1 - P_{block}^{(N)} = \frac{1}{(N)} \left( 1 + \frac{\tilde{p}}{N} + o(N^{-\frac{1}{2}}) \right) \left[ (1 - \frac{d}{C}) + \frac{1}{N} \sum_{n=0}^C n(n - n_{+1}) \right] + o(N^{-\frac{1}{2}}):$$

After simple calculations, we obtain

$$P_{block}^{(N)} = \frac{d}{C} - \frac{1}{N} \left( \sum_{n=0}^C n(n - n_{+1}) \right) - \frac{\tilde{p}}{N} (1 - \frac{d}{C}) + o(N^{-\frac{1}{2}}):$$

This completes the proof.

#### 4.4.7 Proof of Lemma 4.2

The proof relies on the quasi-monotonicity of the mean-field. Let us write the unique solution to the MFEs with the initial point  $\mathbf{v}$  as  $(\mathbf{y}(t; \mathbf{v}); t \geq 0)$ . From the quasi-monotonicity of the mean-field, we have

$$\mathbf{y}(t; \min(\mathbf{v}; \mathbf{a})) \leq \mathbf{y}(t; \mathbf{v}) \leq \mathbf{y}(t; \max(\mathbf{v}; \mathbf{b})): \quad (4.69)$$

For  $\mathbf{a}; \mathbf{b} \geq \mathbf{0}$ , let

$$k_{\mathbf{a}} = \mathbf{b} k_1 = \sum_{i=0}^C j a_i - b_j:$$

From Lemma 4 of [8], since  $\min(\mathbf{v}; \mathbf{a}) \geq \mathbf{0}$  and  $\max(\mathbf{v}; \mathbf{b}) \geq \mathbf{0}$ , we have

$$\begin{aligned} k_{\mathbf{y}}(t; \min(\mathbf{v}; \mathbf{a})) &\leq k_1 e^{-t k_{\min(\mathbf{v}; \mathbf{a})}} \leq k_1; \\ k_{\mathbf{y}}(t; \max(\mathbf{v}; \mathbf{b})) &\leq k_1 e^{-t k_{\max(\mathbf{v}; \mathbf{b})}} \leq k_1; \end{aligned} \quad (4.70)$$

For  $t \geq 0$ , let us define two sets  $V_+(t)$  and  $V_-(t)$  as

$$V_+(t) = \{i : y_i(t; \mathbf{v}) \geq g_i\}$$

$$V(t) = \{i : y_i(t; \mathbf{v}) < g\}$$

Then we can write

$$k_{\mathbf{y}}(t; \mathbf{v}) = k_1 = \sum_{i \in V_+(t)} (y_i(t; \mathbf{v}) - g) + \sum_{j \in V_-(t)} (g - y_j(t; \mathbf{v})).$$

From (4.69) and (4.70), we write

$$\begin{aligned} k_{\mathbf{y}}(t; \mathbf{v}) &= k_1 \left( \sum_{i \in V_+(t)} (y_i(t; \max(\mathbf{v}; \mathbf{g})) - g) + \sum_{j \in V_-(t)} (g - y_j(t; \min(\mathbf{v}; \mathbf{g}))) \right); \\ &= k_{\mathbf{y}}(t; \max(\mathbf{v}; \mathbf{g})) - k_1 + k_{\mathbf{y}}(t; \min(\mathbf{v}; \mathbf{g})) - k_1; \\ &= e^{-t} k_{\max(\mathbf{v}; \mathbf{g})} - k_1 + e^{-t} k_{\min(\mathbf{v}; \mathbf{g})} - k_1; \\ &= e^{-t} k_{\mathbf{v}} - k_1. \end{aligned}$$

Finally, the result follows from the fact that the norms  $k_{k_1}$  and  $k_{k_2}$  are equivalent.

## 4.5 Conclusions

In this chapter, we studied a FCLT that characterizes the fluctuations of the stochastic empirical process for both the transient and stationary regimes. Our analysis covers the Halfin-Whitt regime as a special case. We then used this result to quantify the error between actual blocking probabilities and asymptotic blocking probabilities.

# Chapter 5

## Summary and Future Research

In this dissertation, we studied the robustness and accuracy of a mean-field approach in analyzing the performance of two important classes of models arising in applications.

We mainly focused on two themes. In the first theme, we addressed the case of general JLDs for the loss and PS models. For these systems, we showed that there exists a mean-field limit that is a unique solution of a set of PDEs. Furthermore, we characterized the fixed-points of the mean-field. We have shown that the robustness measured through insensitivity extends to any occupancy based randomized load balancing policy for which the mean-field possesses a globally stable fixed-point. In the second theme, we studied the accuracy of the mean-field approximations of the loss model with exponential JLDs under the assumption that the dispatcher uses the  $SQ(d)$  load balancing policy. We studied the fluctuation of the stochastic empirical process around the mean-field limit. We then used this analysis to quantify the error between the actual blocking probability and the asymptotic blocking probability obtained as a function of the fixed-point of the mean-field.

In Chapter 2, we investigated the impact of the  $SQ(d)$  policy for the loss model that can be used to model Infrastructure-as-a-Service (IaaS) clouds. Under certain assumptions on the initial empirical random measures, we showed that there exists a mean-field. We then exploited the existence of the mean-field limit to show the asymptotic independence of any finite set of servers as  $N \rightarrow \infty$  when the initial states of servers are exchangeable. After that, we proved the uniqueness of the fixed-point of the mean-field. The proof uses two existing results, the form of the stationary distribution of a single server system and the uniqueness of the fixed-point in the exponential case. From the stationary distribution of a single-server loss system with a Poisson arrival process having constant state-dependent

arrival rates, the distribution of occupancy obtained from any fixed-point of the mean-field satisfies the corresponding fixed-point equations in the exponential job length distributions (JLDs) case having the same average job length. Since the exponential case has a unique fixed-point, the occupancy distribution obtained from any fixed-point must coincide with the fixed-point in the exponential case implying insensitivity. Again from the form of the stationary distribution of state-dependent loss models, the fixed-point of the mean-field corresponds to a joint distribution of occupancy and ages that satisfies a product form. We also studied the mean-field equations (MFEs) numerically for mixed-Erlang JLDs, and we observed the GAS of the fixed-point. The numerical evidence for the GAS of the fixed-point supports the hypothesis that the fixed-point approximates the stationary distribution of the state of a server.

In Chapter 3, we assumed that the servers are PS servers. For this model, we studied occupancy based randomized load balancing policies under a common framework as various parameters that include complexity and variations in the traffic level influence the choice of the routing policy in practice. We showed that all occupancy based randomized load balancing policies could be studied in the transient regime under a common framework by using mean-field techniques. We proved the existence of a mean-field limit and the corresponding MFEs correspond to the dynamics of the distribution of a single server system with a Poisson job arrival process having a rate that depends both on the instantaneous occupancy as well as the instantaneous distribution. Every occupancy based policy affects only the job arrival rate function. Furthermore, we showed that every fixed-point corresponds to a distribution on  $Z_+$ , which is a fixed-point in the exponential case having the same average job length. If the exponential case has a unique fixed-point, then the fixed-point in the general case is unique and insensitive. Although one can study occupancy based policies under a common framework, the stationary analysis depends strongly on the chosen routing policy. Hence, for a routing policy that falls into our framework, one can directly obtain the mean-field equations from our analysis, and then the existence, uniqueness, and stability of the fixed-point should be studied separately. We applied our results to four policies and observed that all of them exhibit insensitivity.

Finally, we investigated the stationary behavior of the mean-field for the  $SQ(d)$  policy by using numerical methods for the case of mixed-Erlang JLDs. We observed that the mean-field is not quasi-monotonic, but still, the fixed-point is GAS. We also showed a generic result that if a sequence of the stationary distributions of the empirical measures converges to a distribution, then any finite set of servers in the corresponding limit system



will be independent of each other if and only if the limiting measure is Dirac measure. Otherwise, servers are coupled through the position of the mean-field whose position is random. Conditioned on the position of the mean-field, any finite set of servers are independent of each other, and their distributions coincide with the position of the mean-field.

In Chapter 4, we established an FCLT to characterize the fluctuation of the empirical process around the mean-field. We showed that in the transient regime, the limiting diffusion process is an OU process whose drift and diffusion coefficients depend on the mean-field. In the stationary regime, the limiting diffusion is a stationary OU process whose drift and diffusion coefficients depend on the fixed-point of the mean-field. For the analysis of the stationary regime, we used the exponential stability of the fixed-point of the mean-field. We then used the FCLT and Little's law to show that  $\lim_{N \rightarrow \infty} \rho \overline{N}(P_{block}^{(N)} \quad \frac{d}{C}) \notin 0$  in the Halfin-Whitt regime and  $\lim_{N \rightarrow \infty} \rho \overline{N}(P_{block}^{(N)} \quad \frac{d}{C}) = 0$  if  $\tilde{\rho} = 0$ . This implies  $\lim_{N \rightarrow \infty} \rho \overline{N} P_{block}^{(N)} = 1$  whereas for the complete resource pooling case and for the JSQ policy, we have  $\lim_{N \rightarrow \infty} \rho \overline{N} P_{block}^{(N)} < 1$ .

## 5.1 Future Research

A significant open problem in Chapters 2 and 3 is, showing the GAS of the fixed-point of the mean-field. It is easier to investigate this problem for the simple case of mixed-Erlang distributions. Showing this result will help us to better understand the general case due to the fact that mixed-Erlang distributions are dense in the class of general JLDs.

Our numerical and simulation results provide evidence that the fixed-point of the mean-field is GAS and the limiting stationary distribution coincides with the fixed-point of the mean-field. Since the fixed-point corresponds to the stationary distribution of a reversible Markov process, our investigations support the hypothesis that the non-reversible Markov process which models the dynamics with a finite number of servers is converging to a reversible Markov process as  $N \rightarrow \infty$ . Showing this result will lend credence to the results in [63] where it was shown that for any reversible Markov process that corresponds to a mean-field model, any limit point of the stationary distribution is concentrated at fixed-points of the mean-field.

In Chapter 4, we studied the gap between  $P_{block}^{(N)}$  and  $\frac{d}{C}$  as a function of  $N$ . Future work should also address the Halfin-Whitt regime for the PS model. For this case, in the limiting system as  $N \rightarrow \infty$ , each server is critically loaded. Hence, the average response

time of a job approaches  $1/\lambda$  as  $N \rightarrow \infty$ . But for each  $N$ , the average response time is finite. Hence, it is of interest to find the rate of convergence of the average response times to  $1/\lambda$  as a function of  $N$  when  $N \rightarrow \infty$ . It is of interest to see whether the randomized schemes provide a slower increase over pure random routing or not.

# References

- [1] “Amazon EC2,” <http://aws.amazon.com/ec2/>.
- [2] “Microsoft Azure,” <http://www.microsoft.com/windowsazure/>.
- [3] Y. Lu, Q. Xie, G. Kliot, A. Geller, J. Larus, and A. Greenberg, “Join-idle-queue: A novel load balancing algorithm for dynamically scalable web services,” in *The 29th International Symposium on Computer Performance, Modeling, Measurements and Evaluation*, October 2011.
- [4] V. Gupta, M. Harchol Balter, K. Sigman, and W. Whitt, “Analysis of Join-the-shortest-queue routing for web server farms,” *Perform. Eval.*, vol. 64, no. 9-12, pp. 1062–1081, Oct. 2007.
- [5] W. Winston, “Optimality of the shortest line discipline,” *Journal of Applied Probability*, vol. 14(1), pp. 181–189, 1977.
- [6] R. R. Weber, “On the optimal assignment of customers to parallel servers,” *J. Appl. Probability*, vol. 15, no. 2, pp. 406–413, 1978.
- [7] N. D. Vvedenskaya, R. L. Dobrushin, and F. I. Karpelevich, “Queueing system with selection of the shortest of two queues: an asymptotic approach,” *Problems of Information Transmission*, vol. 32, no. 1, pp. 20–34, 1996.
- [8] Q. Xie, X. Dong, Y. Lu, and R. Srikant, “Power of d choices for large-scale bin packing: A loss model,” in *Proceedings of the 2015 ACM SIGMETRICS*, 2015, pp. 321–334.
- [9] A. Mukhopadhyay and R. R. Mazumdar, “Analysis of randomized join-the-shortest-queue (JSQ) schemes in large heterogeneous processor sharing systems,” *IEEE Transactions on Control of Network Systems*, vol. 3(2), pp. 116–126, 2016.

- [10] A. Mukhopadhyay, R. R. Mazumdar, and F. Guillemin, “The power of randomized routing in heterogeneous loss systems,” in *Teletra c Congress (ITC 27), 2015 27th International*, 2015, pp. 125–133.
- [11] Mukhopadhyay, Arpan, “Mean field interactions in heterogeneous networks,” Ph.D. dissertation, 2016. [Online]. Available: <http://hdl.handle.net/10012/10253>
- [12] M. Mitzenmacher, “The power of two choices in randomized load balancing,” *PhD Thesis, Berkeley*, 1996.
- [13] L. Brown, N. Gans, A. Mandelbaum, A. Sakov, H. Shen, S. Zeltyn, and L. Zhao, “Statistical analysis of a telephone call center,” *Journal of the American Statistical Association*, vol. 100, no. 469, pp. 36–50, 2005.
- [14] P. Kolesar, “Stalking the endangered cat: A queueing analysis of congestion at automatic teller machines,” *Interfaces*, vol. 14, no. 6, pp. 16–26, 1984.
- [15] S. R. E. Turner, “Resource pooling in stochastic networks,” *Ph.D. dissertation, University of Cambridge*, 1996.
- [16] D. A. Dawson, *Measure-valued Markov processes*, ser. École d’Été de Probabilités de Saint-Flour XXI—1991. Berlin: Springer, 1993, vol. 1541.
- [17] P. Whittle, “Partial balance and insensitivity,” *J. Appl. Probab.*, vol. 22, no. 1, pp. 168–176, 1985.
- [18] T. Bonald and A. Proutiere, “Insensitivity in processor sharing networks,” *Performance Evaluation*, vol. 49, no. 1, pp. 193–209, 2002.
- [19] A. Mukhopadhyay, A. Karthik, R. R. Mazumdar, and F. M. Guillemin, “Mean field and propagation of chaos in multi-class heterogeneous loss models,” *Performance Evaluation*, vol. 91, pp. 117–131, September 2015.
- [20] C. Graham and S. Méléard, “Propagation of chaos for a fully connected loss network with alternate routing,” *Stochastic Processes and their Applications*, vol. 44, no. 1, pp. 159–180, 1993.
- [21] —, “Stochastic particle approximations for generalized boltzmann models and convergence estimates,” *The Annals of Probability*, vol. 28, no. 1, pp. 115–132, 1997.

- [22] S. L. Brumelle, “A generalization of Erlang’s loss system to state dependent arrival and service rates,” *Mathematics of Operations Research*, vol. 3, no. 1, pp. 10–16, 1978.
- [23] T. Bonald, M. Jonckheere, and A. Proutière, “Insensitive load balancing,” in *Proceedings of the Joint International Conference on Measurement and Modeling of Computer Systems*, ser. SIGMETRICS ’04/Performance ’04. New York, NY, USA: ACM, 2004, pp. 367–377.
- [24] M. Bramson, Y. Lu, and B. Prabhakar, “Randomized load balancing with general service time distributions,” in *Proceedings of ACM SIGMETRICS*, 2010, pp. 275–286.
- [25] F. I. Karpelevich and A. N. Rybko, “Thermodynamic limit for the mean field model of simple symmetrical closed queueing network,” *Markov Processes and Related Fields*, vol. 6, pp. 89–105, 2000.
- [26] R. Aghajani and K. Ramanan, “The hydrodynamic limit of a randomized load balancing network,” *Ann. Appl. Probab.*, vol. 29, no. 4, pp. 2114–2174, 08 2019.
- [27] H. C. Gromoll, A. L. Puha, and R. J. Williams, “The fluid limit of a heavily loaded processor sharing queue,” *Ann. Appl. Probab.*, vol. 12, no. 3, pp. 797–859, 08 2002.
- [28] H. C. Gromoll, P. Robert, and B. Zwart, “Fluid limits for processor-sharing queues with impatience,” *Math. Oper. Res.*, vol. 33, no. 2, pp. 375–402, May 2008.
- [29] L. Decreusefond and P. Moyal, “A functional central limit theorem for the M/GI/1 queue,” *Ann. Appl. Probab.*, vol. 18, no. 6, pp. 2156–2178, 12 2008.
- [30] J. Zhang, “Fluid models of many-server queues with abandonment,” *Queueing Systems*, vol. 73, no. 2, pp. 147–193, 2013. [Online]. Available: <http://dx.doi.org/10.1007/s11134-012-9307-9>
- [31] H. Kaspi and K. Ramanan, “Law of large numbers limits for many-server queues,” *Ann. Appl. Probab.*, vol. 21, no. 1, pp. 33–114, 02 2011.
- [32] C. Graham, “Chaoticity on path space for a queueing network with selection of the shortest queue among several,” *Journal of Applied Probability*, vol. 37, no. 1, pp. 198–211, 2000.
- [33] T. Vasantam, A. Mukhopadhyay, and R. R. Mazumdar, “Mean-field analysis of loss models with mixed-Erlang distributions under Power-of-d routing,” in *2017 29th International Teletra c Congress (ITC 29)*, vol. 1, Sept 2017, pp. 250–258.

- [34] —, “Insensitivity of the mean-field limit of loss systems under SQ( $d$ ) routing,” *Advances in Applied Probability*, vol. 51, no. 4, pp. 1–40, December 2019.
- [35] S. F. Yashkov and A. S. Yashkova, “Processor sharing: A survey of the mathematical theory,” *Automation and Remote Control*, vol. 68, no. 9, pp. 1662–1731, 2007.
- [36] L. Kleinrock, “Time-shared systems: A theoretical treatment,” *J. ACM*, vol. 14, no. 2, pp. 242–261, Apr. 1967. [Online]. Available: <http://doi.acm.org/10.1145/321386.321388>
- [37] A. Mukhopadhyay and R. R. Mazumdar, “Rate-based randomized routing in large heterogeneous processor sharing systems,” in *Proceedings of 26th International Teletraffic Congress (ITC 26)*, 2014.
- [38] M. Bramson, “Stability of join the shortest queue networks,” *Ann. Appl. Probab.*, vol. 21, no. 4, pp. 1568–1625, 08 2011.
- [39] M. Bramson, Y. Lu, and B. Prabhakar, “Randomized load balancing with general service time distributions,” in *Proceedings of the ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, ser. SIGMETRICS ’10. New York, NY, USA: ACM, 2010, pp. 275–286.
- [40] —, “Asymptotic independence of queues under randomized load balancing,” *Queueing Systems*, vol. 71, no. 3, pp. 247–292, 2012.
- [41] T. Vasantam, A. Mukhopadhyay, and R. R. Mazumdar, “The mean-field behavior of processor sharing systems with general job lengths under the SQ( $d$ ) policy,” *Performance Evaluation*, vol. 127-128, pp. 120 – 153, 2018.
- [42] T. Vasantam and R. Mazumdar, “On occupancy based randomized routing schemes in large systems of shared servers,” *2018 30th International Teletraffic Congress (ITC 30)*, vol. 01, pp. 28–36, 2018.
- [43] P. Gazdzicki, I. Lambadaris, and R. Mazumdar, “Blocking probabilities for large multi-rate Erlang loss systems,” *Adv. Appl. Probab.*, vol. 25, pp. 997–1009, 1993.
- [44] J. van Leeuwen, B. Mathijssen, and B. Zwart, “Economies-of-scale in many-server queueing systems: Tutorial and partial review of the qed halfin–whitt heavy-traffic regime,” *SIAM Review*, vol. 61, no. 3, pp. 403–440, 2019. [Online]. Available: <https://doi.org/10.1137/17M1133944>

- [45] W. Whitt, “Heavy-traffic approximations for service systems with blocking,” *AT&T Bell Laboratories Technical Journal*, vol. 63, no. 5, pp. 689–708, May 1984.
- [46] D. Mukherjee, S. Borst, J. S. H. van Leeuwaarden, and P. Whiting, “Asymptotic optimality of Power-of- $d$  load balancing in large-scale systems,” *CoRR*, 12/2016. [Online]. Available: <http://arxiv.org/abs/1612.00722>
- [47] C. Graham, “Functional central limit theorems for a large network in which customers join the shortest of several queues,” *Probability Theory and Related Fields*, vol. 131, no. 1, pp. 97–120, Jan 2005.
- [48] P. J. Hunt, “Loss networks under diverse routing: The symmetric star network,” *Advances in Applied Probability*, vol. 27, no. 1, pp. 255–272, 1995.
- [49] L. Ying, “On the approximation error of mean-field models,” in *Proceedings of the 2016 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Science*, ser. SIGMETRICS ’16. New York, NY, USA: ACM, 2016, pp. 285–297.
- [50] N. Gast and B. Van Houdt, “A refined mean field approximation,” *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 1, no. 2, pp. 33:1–33:28, Dec. 2017.
- [51] N. Gast, “Expected values estimated via mean-field approximation are  $1/n$ -accurate,” *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 1, no. 1, pp. 17:1–17:26, Jun. 2017. [Online]. Available: <http://doi.acm.org/10.1145/3084454>
- [52] P. Eschenfeldt and D. Gamarnik, “Join the shortest queue with many servers. The heavy-traffic asymptotics,” *Math. Oper. Res.*
- [53] D. Mukherjee, S. Borst, J. van Leeuwaarden, and P. Whiting, “Universality of Power-of- $d$  load balancing in many-server systems,” *Stochastic Systems*, vol. 8, no. 4, pp. 265–292, 2018.
- [54] P. Eschenfeldt and D. Gamarnik, “Supermarket queueing system in the heavy traffic regime. short queue dynamics,” 2016.
- [55] T. Vasantam and R. R. Mazumdar, “Fluctuations around the mean-field for a large scale Erlang loss system under the SQ( $d$ ) load balancing, *accepted to Teletra c Congress (ITC 31), 2019 31st International*,” 2019.

- [56] S. N. Ethier and T. G. Kurtz, *Markov Processes: Characterization and Convergence*. John Wiley and Sons Ltd, 1985.
- [57] P. Billingsley, *Convergence of probability measures*, 2nd ed., ser. Wiley Series in Probability and Statistics: Probability and Statistics. John Wiley & Sons, Inc., New York, 1999, a Wiley-Interscience Publication.
- [58] T. G. Kurtz, “Solutions of ordinary differential equations as limits of pure jump markov processes,” *Journal of Applied Probability*, vol. 7, no. 1, p. 4958, 1970.
- [59] A.-S. Sznitman, “Topics in propagation of chaos,” in *Ecole d'Ete de Probabilites de Saint-Flour XIX / 1989*, P.-L. Hennequin, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 1991, pp. 165–251.
- [60] J.-w. Cho, J.-Y. Le Boudec, and Y. Jiang, “On the Asymptotic Validity of the Decoupling Assumption for Analyzing 802.11 MAC Protocol,” *arXiv e-prints*, p. arXiv:1106.6328, Jun 2011.
- [61] M. Benaim and J.-Y. Le Boudec, “On mean field convergence and stationary regime,” 11 2011.
- [62] M. Benam and J.-Y. L. Boudec, “A class of mean field interaction models for computer and communication systems,” *Performance Evaluation*, vol. 65, no. 11, pp. 823 – 838, 2008, performance Evaluation Methodologies and Tools: Selected Papers from ValueTools 2007.
- [63] J.-Y. L. Boudec, “The stationary behaviour of fluid limits of reversible processes is concentrated on stationary points,” p. 529, 2013. [Online]. Available: <http://aimsciences.org//article/id/ad4efa6f-9312-4213-a0c5-216623b07682>
- [64] N. Gast, L. Bortolussi, and M. Tribastone, “Size expansions of mean field approximation: Transient and steady-state analysis,” *Performance Evaluation*, vol. 129, pp. 60 – 80, 2019. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0166531618302633>
- [65] L. Bortolussi and N. Gast, “Mean field approximation of uncertain stochastic models,” in *2016 46th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, June 2016, pp. 287–298.



- [66] A. Mukhopadhyay, R. R. Mazumdar, and R. Roy, “Binary opinion dynamics with biased agents and agents with different degrees of stubbornness,” in *2016 28th International Teletraffic Congress (ITC 28)*, vol. 01, Sep. 2016, pp. 261–269.
- [67] Q.-L. Li and C. Lin, “The M/G/1 processor-sharing queue with disasters,” *Computers & Mathematics with Applications*, vol. 51, no. 6, pp. 987 – 998, 2006.
- [68] B. A. Sevastyanov, “An ergodic theorem for markov processes and its application to telephone systems with refusals,” *Theory of Probability & Its Applications*, vol. 2, no. 1, pp. 104–112, 1957.
- [69] W. Rudin, *Real and complex analysis*, 3rd ed. McGraw-Hill Book Co., New York, 1987.
- [70] J. Gärtner, “On the mckean-vlasov limit for interacting diffusions,” *Mathematische Nachrichten*, vol. 137, no. 1, pp. 197–248, 1988.
- [71] K. Oelschläger, “A martingale approach to the law of large numbers for weakly interacting stochastic processes,” *Ann. Probab.*, vol. 12, no. 2, pp. 458–479, 05 1984.
- [72] P. M. Kotelenez and T. G. Kurtz, “Macroscopic limits for stochastic partial differential equations of mckean–vlasov type,” *Probability Theory and Related Fields*, vol. 146, no. 1, p. 189, Dec 2008.
- [73] V. N. Kolokoltsov, *Nonlinear Markov Processes and Kinetic Equations*, ser. Cambridge Tracts in Mathematics. Cambridge University Press, 2010.
- [74] A. Jakubowski, “On the skorokhod topology,” *Annales de l’I.H.P. Probabilités et Statistiques*, vol. 22, no. 3, pp. 263–285, 1986.
- [75] O. Kallenberg, *Random measures*. Akademie-Verlag, 1983.
- [76] S. Asmussen, *Applied Probability and Queues*, ser. Stochastic Modelling and Applied Probability. Springer, New York, 2003, vol. 51.
- [77] S. M. Ross, *Introduction to Probability Models*. Academic Press; 10th edition, 2009.
- [78] P. Robert, *Stochastic Networks and Queues*, ser. Stochastic Modelling and Applied Probability Series. Springer-Verlag, 2003.

- [79] L. Decreusefond and P. M., *Stochastic Modeling and Analysis of Telecom Networks*, ser. ISTE, London, 2012.
- [80] V. Varadarajan, “On a theorem of F. Riesz concerning the form of linear functionals,” *Fund. Math.*, vol. 46, pp. 209–220, 1959).
- [81] S. F. Yashkov, “A note on application of the method of supplementary variables to the analysis of a processor sharing system,” *Automation and Remote Control*, vol. 69, no. 9, p. 1622, 2008.
- [82] K. Deimling, *Ordinary differential equations in Banach spaces*, ser. Lecture notes in mathematics. Springer-Verlag Berlin Heidelberg, 1977, vol. 596.
- [83] G. Pang, R. Talreja, and W. Whitt, “Martingale proofs of many-server heavy-traffic limits for markovian queues,” *Probab. Surveys*, vol. 4, pp. 193–267, 2007.
- [84] W. Ledermann and G. E. H. Reuter, “Spectral theory for the differential equations of simple birth and death processes,” *Philosophical Transactions of the Royal Society of London. Series A, Mathematical and Physical Sciences*, vol. 246, no. 914, pp. 321–369, 1954.
- [85] I. Karatzas and S. E. Shreve, *Brownian motion and stochastic calculus*, 2nd ed., ser. Graduate texts in mathematics; 113. Springer-Verlag, New York, 1991, vol. 113.

# Appendix A

## Some Key Background Results

In this Appendix, we collect some background material related to results that are frequently used in the dissertation. Consider a metric space  $(U; d_U)$ . Let  $\mathcal{B}(U)$  be the Borel  $\sigma$ -algebra of  $U$  and  $M_1(U)$  is the space of Borel probability measures on  $U$ .

## A.1 Weak Convergence and Prohorov's Theorem

We first define the concept of weak convergence of probability measures.

**Definition A.1.** ([56], p.107): For a metric space  $(U; d_U)$ , a sequence of probability measure in  $\mathcal{P}_n \subset M_1(U)$  is said to be weakly convergent to a probability measure  $P \in M_1(U)$  if

$$\lim_{n \rightarrow \infty} \int_{U} f(u) dP_n(u) = \int_{U} f(u) dP(u); \quad (\text{A.1})$$

for all  $f \in C_b(U)$ .

If  $(U; d_U)$  is separable, then the weak convergence of a sequence of probability measures is equivalent to the convergence w.r.t. a metric known as Prohorov metric. We next define the Prohorov metric.

**Definition A.2.** ([56], p.96): For  $P, Q \in M_1(U)$ , the Prohorov metric  $\rho(\cdot; \cdot)$  is defined as

$$\rho(P; Q) = \inf \{ \epsilon > 0 : P(A) \leq Q(A) + \epsilon; \forall A \in \mathcal{C}_\epsilon \}; \quad (\text{A.2})$$

where  $\mathcal{C}_\epsilon$  is the collection of all closed subsets of  $U$  and

$$A \in \mathcal{C}_\epsilon = \{ A \subset U : \inf_{y \in A} d_U(y; U \setminus A) > \epsilon \}; \quad (\text{A.3})$$

The next result shows that if  $(U; d_U)$  is separable and complete, then  $(M_1(U); \rho)$  is also separable and complete.

**Theorem A.1.** ([56], Theorem 1.7, p.101): If  $(U; d_U)$  is separable, then  $(M_1(U); \rho)$  is separable. In addition, if  $(U; d_U)$  is complete, then  $(M_1(U); \rho)$  is complete.

We now define the concept of tightness of a set of probability measures that relates to the compactness of the set of measures.

**Definition A.3.** ([56], p.103): A probability measure  $P \in M_1(U)$  is said to be tight if for every  $\epsilon > 0$ , there exists a compact set  $A \subset U$  such that  $P(A) \geq 1 - \epsilon$ . A family of probability measures  $\mathcal{M} \subset M_1(U)$  is tight, if for every  $\epsilon > 0$ , there exists a compact set  $A$  in  $U$  such that  $\inf_{Q \in \mathcal{M}} Q(A) \geq 1 - \epsilon$ .

We next see conditions on  $(U; d_U)$  that guarantee the tightness of a set of probability measures.

**Lemma A.1.** ([56], Lemma 2.1, p.104) *If  $(U; d_U)$  is complete and separable, then every  $P \in M_1(U)$  is tight.*

We now state an important theorem known as Prohorov's theorem from [56, Theorem 2.2, p.104]. For a given topological space  $X$ , a set  $B \subset X$  is said to be relatively compact if the closure of  $B$  is compact.

**Theorem A.2.** *Prohorov's Theorem: Let  $(U; d_U)$  be complete and separable. If  $M \subset M_1(U)$ , then  $M$  is tight if and only if  $M$  is relatively compact.*

Finally, we state the result on the equivalence of the weak convergence of a sequence of probability measures and their convergence w.r.t. Prohorov metric.

**Theorem A.3.** ([56], Theorem 3.1, p.108): *If  $(U; d_U)$  is separable, then for a sequence of probability measures  $\{P_n\}_{n=1}^\infty$  in  $M_1(U)$  and  $P \in M_1(U)$ , we have  $P_n \rightarrow P$  if and only if  $\lim_{n \rightarrow \infty} \rho(P_n; P) = 0$ .*

## A.2 Riesz Markov Kakutani Theorem

**Theorem A.4.** ([69], Theorem 2.14): *Let  $X$  be a locally compact Hausdorff space and let  $\Lambda$  be a positive linear functional on  $C_c(X)$ , then there exists a unique regular measure  $\mu$  on  $X$  such that*

$$\Lambda f = \int_{x \in X} f(x) d\mu(x); \quad (\text{A.4})$$

## A.3 Jakubowski's Criteria

**Theorem A.5.** ([74], Theorem 4.6): *Let  $(U; d_U)$  be a completely regular topological space. Then a sequence of stochastic processes  $\{\mathbf{X}^N\}_{N=1}^\infty$  of  $D_{M_1(U)}([0; 1])$ -valued random elements defined on  $(\Omega; \mathcal{F}; P)$  is tight if and only if the following two conditions are satisfied:*

J1: *For each  $T > 0$  and  $\epsilon > 0$ , there exists a compact set  $K_T \subset M_1(U)$  such that*

$$\liminf_{N \rightarrow \infty} P(\mathbf{X}_t^N \in K_T; \forall t \in [0; T]) > 1 - \epsilon; \quad (\text{A.5})$$

*This condition is called the compact-containment condition.*

J2: There exists a family  $\mathcal{Q}$  of real valued continuous functions  $F$  defined on  $M_1(U)$  that separates points in  $M_1(U)$  and is closed under addition such that for every  $F \in \mathcal{Q}$ , the sequence  $\int F(\mathbf{X}_t^N); t \in [0, 1] g_{N-1}$  is tight in  $D_{\mathbb{R}}([0; 1])$ .

To prove the condition J2, we define a class of functions  $\mathcal{Q}$  as follows:

$$\mathcal{Q} = \{fF : f \in C_b^1(U) \text{ such that } F(\cdot) = h(\cdot); h \in M_1(U)\} \quad (\text{A.6})$$

It can be seen that every function  $F \in \mathcal{Q}$  is continuous w.r.t. the weak topology on  $M_1(U)$  and furthermore, the class of functions  $\mathcal{Q}$  separates points in  $M_1(U)$  and is also closed under addition. We next recall the following result (From Theorem C.9, [78]) that is sufficient to prove the condition J2.

**Theorem A.6.** *Tightness in  $D_{\mathbb{R}}([0; T])$ : Let  $\{P_n\}_{n \in \mathbb{Z}_+}$  be a sequence of probability distributions on  $D_{\mathbb{R}}([0; T])$ , then  $\{P_n\}_{n \in \mathbb{Z}_+}$  is tight if for any  $\epsilon > 0$ ,*

C1: *There exists  $b > 0$  such that*

$$P_n(\|X(0)\| > b) \quad (\text{A.7})$$

*for all  $n \in \mathbb{Z}_+$ .*

C2: *For any  $\epsilon > 0$ , there exists  $\delta > 0$  such that*

$$P_n(w_X(\delta) > \epsilon) \quad (\text{A.8})$$

*for  $n$  sufficiently large, where*

$$w_X(\delta) = \sup_{\substack{t \in [0; T] \\ |s-t| \leq \delta}} \|X(t) - X(s)\| \quad (\text{A.9})$$

*and any limiting point  $P$  satisfies  $P(C_{\mathbb{R}}([0; T])) = 1$ .*

## A.4 Gronwall's Inequality

**Theorem A.7.** ([56], Theorem 5.1, p.498): *Let  $f$  be a function on  $[0; 1]$  that is bounded on bounded intervals. For  $\alpha \geq 0$  and  $K > 0$ , if*

$$f(t) \leq \alpha + K \int_0^t f(s) ds; \quad (\text{A.10})$$

*then*

$$f(t) \leq e^{Kt} \alpha; \quad (\text{A.11})$$

## A.5 Doob's Inequality

**Theorem A.8.** ([\[56\]](#), Corollary 2.17, p.64): Let  $(\mathbf{X}(t); t \geq 0)$  be a right continuous martingale. Then for  $\lambda > 1$  and  $T > 0$ ,

$$\mathbb{E} \left[ \sup_{t \leq T} |\mathbf{X}(t)|^p \right] \leq \left( \frac{\lambda}{\lambda - 1} \right)^p \mathbb{E} [|\mathbf{X}(T)|^p] \quad (\text{A.12})$$