# Approximations for the Queue Length Distributions of Time-Varying Many-Server Queues

Jamol Pender

School of Operations Research and Information Engineering

Cornell University

jjp274@cornell.edu

Young Myoung Ko

Department of Industrial and Management Engineering

Pohang University of Science and Technology

youngko@postech.ac.kr

January 6, 2016

### Abstract

This paper presents a novel methodology for approximating the queue length (the number of customers in the system) distributions of time-varying non-Markovian many-server queues (e.g., $G_t/G_t/n_t$ queues), where the number of servers ($n_t$) is large. Our methodology consists of two steps. The first step uses phase-type distributions to approximate the general inter-arrival and service times, thus generating an approximating $Ph_t/Ph_t/n_t$ queue. The second step develops strong approximation theory to approximate the $Ph_t/Ph_t/n_t$ queue with fluid and diffusion limits. However, by naively representing the $Ph_t/Ph_t/n_t$ queue as a Markov process by expanding the state space, we encounter the lingering phenomenon even when the queue is overloaded. Lingering typically occurs when the mean queue length is equal or near the number of servers, however, in this case it also happens when the queue is overloaded and this time is not of zero measure. As a result, we develop an alternative representation for the queue length process that avoids the lingering problem in the overloaded case, thus allowing for the derivation of a Gaussian diffusion limit. Finally, we compare the effectiveness of our proposed method with discrete event simulation in a variety parameter settings and show that our approximations are very accurate.

## 1 Introduction

Real-world applications of large-scale queueing systems such as data centers, call centers, and healthcare centers have time-varying and dynamic behavior. Furthermore, the arrival and service processes are not necessarily Markovian in general (Brown et al. [8], Arfeen

et al. [1], Nelson and Taaffe [32]). Many of the recent studies on large-scale non-Markovian queues rely on the asymptotic approaches that utilize fluid and diffusion limits as described in Billingsley [5] and Whitt [46]. Research on non-Markovian systems has progressed to the point of analyzing underloaded systems (a.k.a. the offered-load model, infinite-server queues) due to their analytical or numerical tractability (Whitt [45], Glynn [14], Eick et al. [11], Nelson and Taaffe [32, 31]). Studies on the delay model, e.g., $M_t/G_t/n_t$, $G_t/M_t/n_t$, $G_t/G_t/n_t$ queues, have been conducted from the context of fluid queues or heavy traffic diffusion models in the Halfin-Whitt regime (Halfin and Whitt [15], Puhalskii and Reiman [42], Pang and Whitt [37], Reed [43], Whitt [47], Liu and Whitt [24, 26, 25]).

This paper uses the *uniform acceleration* method coupled with strong approximations and accelerates parameters while keeping the traffic intensity constant, see for example (Kurtz [23], Mandelbaum et al. [27], Hampshire et al. [18]). Kurtz [23] establishes strong approximation theorems for state-dependent continuous time Markov chains (CTMCs) having differentiable rate functions. Extending Kurtz [23], Mandelbaum et al. [27] consider time-varying parameters and non-differentiable rate functions such as $\min(\cdot, \cdot)$ that commonly occur in the analysis of queues. Mandelbaum et al. [28] prove that the strong approximation results developed in Kurtz [23] can also be applied when the fluid limit stays at the non-differentiable points of rate functions for a measure-zero amount of time. However, in some queueing processes, it is hard to avoid the measure-zero assumption. See for instance Niyirora and Pender [34] and Hampshire and Massey [17, 16], Hampshire et al. [20, 19] where optimal staffing methods force staffing at the non-differentiable points.

To address the issue of when the fluid limit is near the non-differentiable points of the rate functions for more than a measure zero amount of time Ko and Gautam [22] propose a Gaussian-based approximation method that achieves better approximation quality. Massey and Pender [29, 30] improve the result of Ko and Gautam [22] by incorporating the skewness of the queueing process and by expanding the queue length process in terms of Hermite polynomials, which are orthogonal with respect to the Gaussian distribution. In the same spirit, the work of Pender [38, 39, 40, 41] extends the results of Massey and Pender [30] and add the impact of the kurtosis through a Gram-Charlier expansion and using other distributions as closure approximations. More work by Engblom and Pender [12] also proves that spectral expansions as closure approximations for the functional Kolmogorov forward equations of the queue length process are provably optimal in an $L^2$ sense for approximating the moments of nonstationary birth-death processes. Although the spectral approach offers great insight especially for higher moments of the queue length processes, fluid and diffusions also offer complementary insight for the sample path behavior of the queueing process.

In the spirit of fluid and diffusion limits, Liu and Whitt [24] prove a weak law of large numbers limit for the $G_t/GI/n_t + GI$ queue and extend the work of Mandelbaum et al. [27] in the sense that they consider non-Markovian inter-arrival, service and abandonment times. However, the service times are not time-varying and the limit does not converge almost surely as the limit in this work. In a follow-up paper, Liu and Whitt [25] provide a heavy-traffic diffusion limit for $G_t/M/s_t + GI$ queues. The methodology used by Liu and Whitt [25] is to paste together the overloaded and underloaded intervals of the nonstationary queueing process. Thus, they explicitly avoid the case where the number of servers is equal to the fluid limit. As shown in Mandelbaum et al. [28], Ko and Gautam [22], Liu and Whitt [25], it appears reasonable to approximate the queue length process with a Gaussian

process. However, estimating the parameters of a Gaussian process depends on both fluid and diffusion limits. Lastly, work by Reed [43] and Dai et al. [10] use the continuous mapping approach to prove diffusion limits for queues with general and phase type service respectively. Although this work was a significant advance in the many server literature, Reed [43], Dai et al. [10] do not explore the impact of nonstationary arrival and service times and this work generalizes their work in this regard. Lastly, since our approximations are for nonstationary processes, the approximations are universally useful and apply in any regime.

Using phase-type distributions for approximating general distributions in queueing analysis is not new, see for example Barbour [4]. The matrix-geometric method (MGM) described in Neuts [33] is a well-known approach for the analysis of non-Markovian queues. MGM, however, can only handle phase-type distributions with a small number of phases due to state space explosion. Nelson and Taaffe [32] develop a method based on the partial-moment differential equations (PMDEs) for the analysis of $Ph_t/Ph_t/\infty$ queues that accurately estimates the moments of the number of entities in the system. The number of differential equations to evaluate the first two moments is $m_A + m_S - 1 + m_A m_S(m_S + 1)$, where $m_A$ and $m_S$ are the number of phases in the inter-arrival and service time distributions, respectively. The result, however, is not applicable to the delay models, such as $Ph_t/Ph_t/n_t$ queues studied in our paper. Creemers et al. [9] devise a phase-type approximation algorithm for small-to-medium-sized queues (2-10 servers) using two-moment matching procedures, however, the downfall is that the method does not scale well with the number servers and it has a high computational cost when the number of servers is large. Our goal is to remove this dependence on the number of servers since it is very limiting in a computational sense, especially in large-scale service systems.

## 1.1 Main Contributions of Paper

The contributions of this work can be summarized as follows. First, we consider the dynamics of a $G_t/G_t/n_t$ queue. The $G_t/G_t/n_t$ queueing model is relatively intractable since we are unable to derive the exact distribution of the queue length as a function of time. Thus, we first approximate the general and non-Markovian arrival and service distributions with phase-type distributions with an appropriate number of phases. This reduces our problem to analyzing the $Ph_t/Ph_t/n_t$ queue, which is more tractable than its general counterpart. Second, we derive fluid and diffusion limits for a $Ph_t/Ph_t/n_t$ queue using *uniform acceleration* coupled with strong approximations of time changed Poisson processes. Unfortunately, when we naively keep track of the number of customers being served in each phase and the number of customers in the system separately, we encounter the *lingering* issue; the fluid limit stays at non-differentiable points during some intervals having positive measure. This prevents us from deriving a Gaussian or continuous diffusion limit. Thus, another important contribution of our work is our proposal of an alternative Markovian formulation of the queueing process that enables us to successfully obtain the diffusion limit. One attractive feature of our method is that the number of differential equations to obtain the fluid and diffusion limits is $O([m_A + m_S]^2)$ and it does *not* depend on the number of servers, $n_t$ like other numerical methods by Creemers et al. [9]. The number of phases used for approximating inter-arrival and service time distributions is 8-10 and the numerical solution is reached in less than a minute using a commercial solver (e.g., MATLAB).

## 1.2 Organization of Paper

The remainder of this paper is organized as follows. Section 2 describes the $G_t/G_t/n_t$ queueing model and the problem settings. Section 3 builds a mathematical model for describing the dynamics of the system for the $Ph_t/Ph_t/n_t$ queue. We explain the impact of the lingering problem and introduce an alternative sample path representation for analyzing it. Section 4 constructs the fluid and diffusion limit theorems as approximations for the sample path dynamics of the queueing process in the finite server setting. Section 5 discusses the infinite server setting and provides the fluid and diffusion limits for the infinite server queueing model. Section 6 discusses the numerical examples used to validate the effectiveness of our proposed approach. Section 7 concludes and offers suggestions for future research.

# 2  Problem description

We consider a $G_t/G_t/n_t$ queue, a time-varying version of a $G/G/n$ queue, with a general time-varying arrival process, a general time-varying service time distribution, and a time-varying number of servers. The system has an infinite capacity of waiting space and customers in the waiting space are served under the first-come, first-served discipline. Let $X(t)$ denote the number of customers in the system at time $t$ and $\bar{x}(t)$ denote the corresponding fluid limit. We assume that the fluid limit $(\bar{x}(t))$ alternates between the underloaded (i.e., $\bar{x}(t) < n_t$) and overloaded (i.e.. $\bar{x}(t) > n_t$) regimes and hits the critically loaded regime (i.e. $\bar{x}(t) = n_t$) at most a countable number of times. The performance measures of interest are $\mathrm{E}[X(t)]$, $\mathrm{Var}[X(t)]$ and, if possible, the distribution of $X(t)$ for all time $0 \leq t \leq T$ and $T < \infty$.

More specifically, we analyze a $Ph_t/Ph_t/n_t$ queue as an approximation of the $G_t/G_t/n_t$ queue since phase-type distributions are dense in all positive-support distributions and the use of phase-type distribution in queueing analysis does not lose generality significantly (Barbour [4], Whitt [45], and Asmussen et al. [3]). A phase-type distribution with $m$ phases represents the time taken from an initial state to an absorbing state of a continuous time Markov chain with the following infinitesimal generator matrix:

$$\mathbf{Q} = \begin{pmatrix} 0 & \mathbf{0} \\ \mathbf{s} & \mathbf{S} \end{pmatrix},$$

where $\mathbf{0}$ is a $1 \times m$ zero vector, $\mathbf{s} =$ is an $m \times 1$ vector, and $\mathbf{S}$ is an $m \times m$ matrix. Note $\mathbf{s} = -\mathbf{Se}$ where $\mathbf{e}$ is an $m \times 1$ vector of ones. The matrix $\mathbf{S}$ and the initial distribution $\boldsymbol{\alpha}$ which is a $1 \times m$ vector identify the phase-type distributions. Finding the best phase-type distribution for approximating a general distribution is beyond the scope of this paper, and we refer to the reader to a large number of references [6, 21, 48, 7, 13, 36, 3, 35]. To give the reader a better understanding of our methodology, we describe the fitting algorithm that we use in Section 6.

We assume that our phase-type distributions have initial distributions, $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$, and infinitesimal generator matrices, $\mathbf{Q_A}$ and $\mathbf{Q_S}$, for the arrival process and service times respectively. The number of phases in $\mathbf{S_A}$ and $\mathbf{S_S}$ is $m_A$ and $m_S$ respectively. The matrices,

4

$\mathbf{S_A}$ and $\mathbf{S_S}$, and the vectors, $\mathbf{s_A}$ and $\mathbf{s_S}$ can be expressed as:

$$\mathbf{S_A} = \begin{pmatrix} \lambda_{11} & \cdots & \lambda_{1m_A} \\ \vdots & \vdots & \vdots \\ \lambda_{m_A 1} & \cdots & \lambda_{m_A m_A} \end{pmatrix}, \quad \mathbf{s_A} = (\lambda_{10}, \ldots, \lambda_{m_A 0})' \tag{2.1}$$

$$\mathbf{S_S} = \begin{pmatrix} \mu_{11} & \cdots & \mu_{1m_S} \\ \vdots & \vdots & \vdots \\ \mu_{m_S 1} & \cdots & \mu_{m_S m_S} \end{pmatrix}, \quad \mathbf{s_S} = (\mu_{10}, \ldots, \mu_{m_S 0})', \tag{2.2}$$

where $\lambda_{jk}$'s and $\mu_{il}$'s agree with the definition of the infinitesimal generator matrices, $\mathbf{Q_A}$ and $\mathbf{Q_S}$. Note that the time-varying extension can be achieved by replacing $\lambda_{jk}$ and $\mu_{il}$ with $\lambda_{jk}(t)$ and $\mu_{il}(t)$ and making sure that their integrals are locally bounded away from infinity.

# 3 The Queueing model

With the phase-type distributions described in Section 2, we build a mathematical queueing model to describe the dynamics of the $Ph_t/Ph_t/n_t$ queue. We assume that the system starts with no customers.
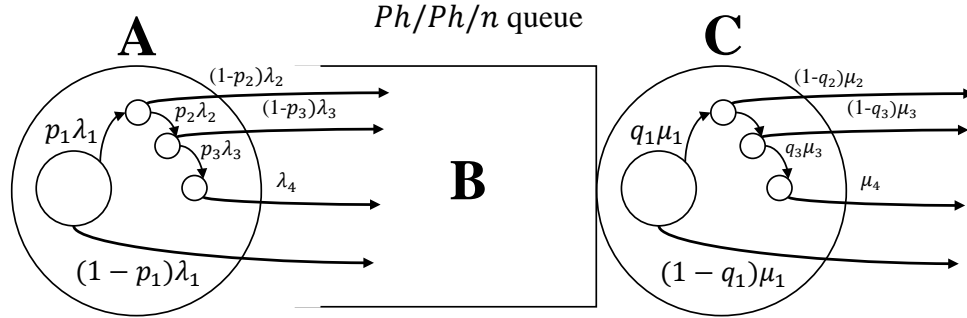


Figure 3.1: $Ph/Ph/n$ queue with Coxian distributions

Figure 3.1 illustrates an example of $Ph/Ph/n$ queue with Coxian inter-arrival and service times. In order to model the $Ph_t/Ph_t/n_t$ queue, we need to keep track of the phase in which the arriving customer is (area $\mathbf{A}$ in Figure 3.1), the number of customers being served in each phase (area $\mathbf{C}$), and the number of customers in the waiting space (area $\mathbf{B}$). We let $U_i(t)$ be the number of customers in phase $i$ of the arrival process at time $t$, $X_j(t)$ be the number of customers being served in phase $j$ of the service process, and $Z(t)$ be the total number of customers in the system. Note that the number of customers in the waiting space is $Z(t) - \sum_{i=1}^{m_S} X_i(t) \geq 0$ and $\sum_{i=1}^{m_A} U_i(t) = 1$ for all $t > 0$. Then, the state of the system $\mathbf{V}(t) = (U_1(t), \ldots, U_{m_A}, X_1(t), \ldots, X_{m_S}, Z(t))'$ is the solution to the following integral

equations:

$$U_j(t) = U_j(0) + \sum_{k \neq j}^{m_A} Y_{kj}^A \left( \int_0^t \lambda_{kj} U_k(s) ds \right) - \sum_{k \neq j}^{m_A} Y_{jk}^A \left( \int_0^t \lambda_{jk} U_j(s) ds \right) \tag{3.1}$$

$$- \sum_{k \neq j}^{m_A} \sum_{l=1}^{m_S} Y_{jkl}^I \left( \int_0^t \lambda_{j0} \alpha_k \beta_l U_j(s) \mathbf{1}_{\{Z(s) \leq n\}} ds \right) - \sum_{k \neq j}^{m_A} Y_{jk}^Q \left( \int_0^t \lambda_{j0} \alpha_k U_j(s) \mathbf{1}_{\{Z(s) > n\}} ds \right)$$

$$+ \sum_{k \neq j}^{m_A} \sum_{l=1}^{m_S} Y_{kjl}^I \left( \int_0^t \lambda_{k0} \alpha_j \beta_l U_k(s) \mathbf{1}_{\{Z(s) \leq n\}} ds \right)$$

$$+ \sum_{k \neq j}^{m_A} Y_{kj}^Q \left( \int_0^t \lambda_{k0} \alpha_j U_k(s) \mathbf{1}_{\{Z(s) > n\}} ds \right) \text{ for } 1 \leq j \leq m_A,$$

$$X_i(t) = \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} Y_{jki}^I \left( \int_0^t \lambda_{j0} \alpha_k \beta_i U_j(s) \mathbf{1}_{\{Z(s) \leq n\}} ds \right) + \sum_{l \neq i}^{m_S} Y_{li}^S \left( \int_0^t \mu_{li} X_l(s) ds \right) \tag{3.2}$$

$$- \sum_{l \neq i}^{m_S} Y_{il}^S \left( \int_0^t \mu_{il} X_i(s) ds \right) - Y_{i0}^D \left( \int_0^t \mu_{i0} X_i(s) \mathbf{1}_{\{Z(s) \leq n\}} ds \right)$$

$$- \sum_{l \neq i}^{m_S} Y_{il}^D \left( \int_0^t \mu_{i0} X_i(s) \mathbf{1}_{\{Z(s) > n\}} \beta_l ds \right)$$

$$+ \sum_{l \neq i}^{m_S} Y_{li}^D \left( \int_0^t \mu_{l0} X_l(s) \mathbf{1}_{\{Z(s) > n\}} \beta_i ds \right) \text{ for } 1 \leq i \leq m_S,$$

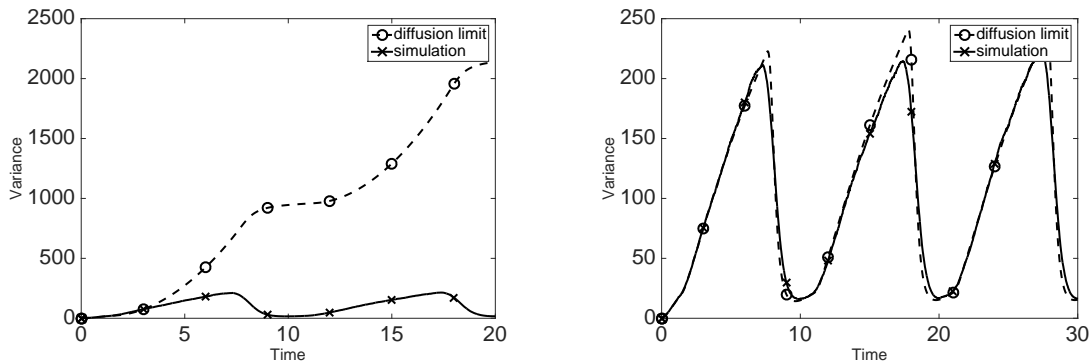$$Z(t) = \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \sum_{l=1}^{m_S} Y_{jkl}^I \left( \int_0^t \lambda_{j0} \alpha_k \beta_l U_j(s) \mathbf{1}_{\{Z(s) \leq n\}} ds \right) \tag{3.3}$$

$$+ \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} Y_{jk}^Q \left( \int_0^t \lambda_{j0} \alpha_k U_j(s) \mathbf{1}_{\{Z(s) > n\}} ds \right) - \sum_{i=1}^{m_S} Y_{i0}^D \left( \int_0^t \mu_{i0} X_i(s) \mathbf{1}_{\{Z(s) \leq n\}} ds \right)$$

$$- \sum_{i=1}^{m_S} \sum_{l=1}^{m_S} Y_{il}^D \left( \int_0^t \mu_{i0} X_i(s) \mathbf{1}_{\{Z(s) > n\}} \beta_l ds \right).$$

For notational convenience, the equations (3.1)-(3.3) represent the dynamics of a $Ph/Ph/n$ queue. As mentioned in Section 2, we can obtain the time-varying extension by replacing $\lambda_{jk}$, $\mu_{il}$ and $n$ with $\lambda_{jk}(t)$, $\mu_{il}(t)$, and $n(t)$ respectively under mild conditions given in Mandelbaum et al. [27]. Poisson processes, $Y_{kj}^A(\cdot)$'s count the number of transitions from phase $k$ to phase $j$ of the arrival process. When the waiting space is empty ($Z(t) \leq n$), Poisson processes, $Y_{jkl}^I(\cdot)$'s, count the number of departures from phase $j$ of the arrival process to phase $l$ of the service process according to the initial distribution $\boldsymbol{\beta}$ and the arrival process restarts from phase $k$ according to the initial distribution $\boldsymbol{\alpha}$. When the waiting space is not empty ($Z(t) > n$), Poisson processes, $Y_{jk}^Q(\cdot)$'s, count the number of departures from phase $j$

of the arrival process to the waiting space and a new arrival process begins in phase $k$. Poisson processes, $Y_{li}^S(\cdot)$'s, count the internal transitions from phase $l$ to phase $j$ of the service process. When the waiting space is empty, Poisson processes, $Y_{i0}^D(\cdot)$'s, count the number of departures from phase $i$ of the service process. When the waiting space is not empty, Poisson processes, $Y_{il}^D(\cdot)$'s, count the number of departures from phase $i$ and a new customer enters phase $l$ from the waiting space. Note that the Poisson processes explained above have rate 1 (with random time changes) and are mutually independent.

We can easily figure out that the rate functions in equations (3.1)-(3.3) (the integrands in Poisson processes) are not differentiable with respect to the elements of the state space vector, $\mathbf{V}(t)$. Thus, before applying the uniform acceleration, we conduct a quick check to find whether the time during which the fluid limit stays at the non-differentiable points has measure zero or not.

Let $\mathbf{v}(t) = (\bar{u}_1(t), \ldots, \bar{u}_{m_A}(t), \bar{x}_1(t), \ldots, \bar{x}_{m_S}(t), \bar{z}(t))'$ be the fluid limit of $\mathbf{V}(t)$. We check the Poisson process, $Y_{il}^D(\cdot)$ in equation (3.2). The fluid limit for $Y_{il}^D(\cdot)$ is $\mu_{i0}\bar{x}_i(t)\mathbf{1}_{\{\bar{z}(t)>n\}}$. When $\bar{z}(t)$ hits $n$, the non-differentiable point, $\sum_{i=1}^{m_S}\bar{x}(t) = n$. However, during the overloaded time $\{t : \bar{z}(t) > n\}$ which can have strictly positive measure in our setting, $\sum_{i=1}^{m_S}\bar{x}(t)$ remains unchanged (i.e., $\sum_{i=1}^{m_S}\bar{x}(t) = n$). This implies that the subvector $(\bar{x}_1(t), \ldots, \bar{x}_{m_S}(t))'$ moves on the hyperplane during the overloaded period and we cannot obtain the diffusion limit from the result of Kurtz [23] and Mandelbaum et al. [28]. When we try to apply fluid and diffusion limits with equations (3.1)-(3.3) just ignoring the issue, we observe a huge gap between simulation and the numerical solution. For example (Exp. 7 in Section 6), Figure 3.2 (a) shows the gap between the simulated variance and the variance from the diffusion limit. We devise an alternative formulation which can significantly improve the approximation accuracy (see Figure 3.2 (b)).



(a) Variance from the formulation in (3.1)-(3.3)    (b) Variance from the alternative formulation

Figure 3.2: Variance estimation of Exp. 7

The issue occurs because $\sum_{i=1}^{m_S}\bar{x}(t) = n$ during the overloaded period. The alternative formulation avoids this situation but requires an additional assumption that the phase-type distribution for service times has a unique initial state. Such distributions include the Erlang distribution and the Coxian distribution. According to Asmussen et al. [3],

the Coxian distribution provides almost the same quality of fit as the general phase-type distribution with the same number of phases. One reason is that the Coxian and generalized hyperexponential distribution, which are specific classes of phase-type distributions, are also dense in the class of positive-support distributions, see for example Sasaki et al. [44]. Thus, the additional assumption of restricting to the Coxian class, therefore, may not be quite restrictive. Without loss of generality, we assume the unique initial state is phase 1. The main idea is to maintain the waiting space inside phase 1 and control transition rates from phase 1 so that the system serves at most $n$ customers. We have the same state space except for $Z(t)$ because $X_1(t)$ accounts for customers in the waiting space. Using this representation, we can now write our new formulation of the queueing process as follows:

$$U_j(t) = U_j(0) + \sum_{k\neq j}^{m_A} Y_{kj}^A\left(\int_0^t \lambda_{kj} U_k(s)ds\right) - \sum_{k\neq j}^{m_A} Y_{jk}^A\left(\int_0^t \lambda_{jk} U_j(s)ds\right) \tag{3.4}$$
$$- \sum_{k\neq j}^{m_A} Y_{jk}^I\left(\int_0^t \lambda_{j0}\alpha_k U_j(s)ds\right) + \sum_{k\neq j}^{m_A} Y_{kj}^I\left(\int_0^t \lambda_{k0}\alpha_j U_k(s)ds\right) \text{ for } 1 \leq j \leq m_A,$$

$$X_1(t) = \sum_{j=1}^{m_A}\sum_{k=1}^{m_A} Y_{jk}^I\left(\int_0^t \lambda_{j0}\alpha_k U_j(s)ds\right) + \sum_{l\neq 1}^{m_S} Y_{l1}^S\left(\int_0^t \mu_{l1} X_l(s)ds\right) \tag{3.5}$$
$$- \sum_{l\neq 1}^{m_S} Y_{1l}^S\left(\int_0^t \mu_{1l}\left[\mathbf{1}_{\{\sum_{r=1}^{m_S} X_r(s)\leq n\}} X_1(s) + \mathbf{1}_{\{\sum_{r=1}^{m_S} X_r(s)>n\}}\left(n - \sum_{r=2}^{m_S} X_r(s)\right)^+\right]ds\right)$$
$$- Y_1^D\left(\int_0^t \mu_{10}\left[\mathbf{1}_{\{\sum_{r=1}^{m_S} X_r(s)\leq n\}} X_1(s) + \mathbf{1}_{\{\sum_{r=1}^{m_S} X_r(s)>n\}}\left(n - \sum_{r=2}^{m_S} X_r(s)\right)^+\right]ds\right).$$

$$X_i(t) = Y_{1i}^S\left(\int_0^t \mu_{1i}\left[\mathbf{1}_{\{\sum_{r=1}^{m_S} X_r(s)\leq n\}} X_1(s) + \mathbf{1}_{\{\sum_{r=1}^{m_S} X_r(s)>n\}}\left(n - \sum_{r=2}^{m_S} X_r(s)\right)^+\right]ds\right) \tag{3.6}$$
$$+ \sum_{l=2,l\neq i}^{m_S} Y_{li}^S\left(\int_0^t \mu_{li} X_l(s)ds\right) - \sum_{l\neq i}^{m_S} Y_{il}^S\left(\int_0^t \mu_{il} X_i(s)ds\right) - Y_i^D\left(\int_0^t \mu_{i0} X_i(s)ds\right)$$
$$\text{for } 2 \leq i \leq m_S.$$

Poisson processes, $Y_{kj}^A(\cdot)$'s and $Y_{li}^S(\cdot)$'s, are the same as those in equations (3.1) and (3.2). Poisson processes, $Y_{jkl}^I(\cdot)$'s in equation (3.1) are now replaced by $Y_{jk}^I(\cdot)$'s because the initial state of the service process is phase 1, that is, we do not need the index of the starting phase in the service process. Then, Poisson processes, $Y_{jk}^I(\cdot)$'s count the number of departures from phase $j$ that restart from phase $k$ of the arrival process according to the initial distribution $\boldsymbol{\alpha}$. Note that we do not have to count the number of departures that restart from the same phase, i.e. we do not count the case of $j = k$. Poisson processes, $Y_i^D(\cdot)$'s count departures from phase $i$ of the service process. Note that the Poisson processes explained above have rate 1 (with random time changes) and are mutually independent. We can verify that the issue is not incurred in equations (3.4)-(3.6). In the following section we describe the fluid and diffusion approximations.

## 3.1 Lipschitz Representation

It turns out that we can write our new formulation in terms of Lipschitz rate functions. This representation will aid us tremendously when proving the fluid and diffusion limit theorems for the queueing model.

$$
\begin{aligned}
U_j(t) &= U_j(0) + \sum_{k \neq j}^{m_A} Y_{kj}^A \left( \int_0^t \lambda_{kj} U_k(s) ds \right) - \sum_{k \neq j}^{m_A} Y_{jk}^A \left( \int_0^t \lambda_{jk} U_j(s) ds \right) \\
&\quad - \sum_{k \neq j}^{m_A} Y_{jk}^I \left( \int_0^t \lambda_{j0} \alpha_k U_j(s) ds \right) + \sum_{k \neq j}^{m_A} Y_{kj}^I \left( \int_0^t \lambda_{k0} \alpha_j U_k(s) ds \right) \text{ for } 1 \leq j \leq m_A, \\
X_1(t) &= \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} Y_{jk}^I \left( \int_0^t \lambda_{j0} \alpha_k U_j(s) ds \right) + \sum_{l \neq 1}^{m_S} Y_{l1}^S \left( \int_0^t \mu_{l1} X_l(s) ds \right) \\
&\quad - \sum_{l \neq 1}^{m_S} Y_{1l}^S \left( \int_0^t \mu_{1l} \left[ \left( X_1(s) \wedge \left( n - \sum_{r=2}^{m_S} X_r(s) \right)^+ \right) \right] ds \right) \\
&\quad - Y_1^D \left( \int_0^t \mu_{10} \left[ \left( X_1(s) \wedge \left( n - \sum_{r=2}^{m_S} X_r(s) \right)^+ \right) \right] ds \right). \\
X_i(t) &= Y_{1i}^S \left( \int_0^t \mu_{1i} \left[ \left( X_1(s) \wedge \left( n - \sum_{r=2}^{m_S} X_r(s) \right)^+ \right) \right] ds \right) + \sum_{l=2, l \neq i}^{m_S} Y_{li}^S \left( \int_0^t \mu_{li} X_l(s) ds \right) \\
&\quad - \sum_{l \neq i}^{m_S} Y_{il}^S \left( \int_0^t \mu_{il} X_i(s) ds \right) - Y_i^D \left( \int_0^t \mu_{i0} X_i(s) ds \right) \qquad \text{for } 2 \leq i \leq m_S.
\end{aligned}
$$

# 4  Fluid and diffusion approximations

In this section, we now provide our second main contribution of the paper, fluid and diffusion limit theorems for the queue length process. However, we first provide some definitions for notational convenience of the reader that will be used throughout the rest of the paper.

$\mathbf{V}(t) = (U_1(t), \ldots, U_{m_A}(t), X_1(t), \ldots, X_{m_S}(t))'.$

$\mathbf{v} = (u_1, \ldots, u_{m_A}, x_1, \ldots, x_{m_S})'.$

$\mathbf{d}_{jk}^A : (m_A + m_S) \times 1$ vector, $j^{\text{th}}$ element is -1, $k^{\text{th}}$ element is 1, and other elements are 0.

$\mathbf{d}_{jk}^I : (m_A + m_S) \times 1$ vector, $j^{\text{th}}$ element is -1, $k^{\text{th}}$ element is 1, and other elements are 0.

$\mathbf{d}_{il}^S : (m_A + m_S) \times 1$ vector, $i^{\text{th}}$ element is -1, $l^{\text{th}}$ element is 1, and other elements are 0.

$\mathbf{d}_i^D : (m_A + m_S) \times 1$ vector, $i^{\text{th}}$ element is -1, and other elements are 0.

$f_{jk}^A(t, \mathbf{v})$ : rate function (integrand) in $Y_{jk}^A(\cdot)$.

$f_{jk}^I(t, \mathbf{v})$ : rate function (integrand) in $Y_{jk}^I(\cdot)$.

$f_{il}^S(t, \mathbf{v})$ : rate function (integrand) in $Y_{il}^S(\cdot)$.

$f_i^D(t, \mathbf{v})$ : rate function (integrand) in $Y_i^D(\cdot)$.

$W_{jk}^A(t), W_{jk}^I(t), W_{il}^S(t), W_i^D(t)$ : mutually independent standard Brownian motions.

$$\mathbf{F}(t, \mathbf{v}) = \sum_{j=1}^{m_A} \sum_{k=1, k\neq j}^{m_A} \mathbf{d}_{jk}^A f_{jk}^A(t, \mathbf{v}) + \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I f_{jk}^I(t, \mathbf{v}) + \sum_{i=1}^{m_S} \sum_{l=1, l\neq i}^{m_S} \mathbf{d}_{il}^S f_{il}^S(t, \mathbf{v}) + \sum_{i=1}^{m_S} \mathbf{d}_i^D f_i^D(t, \mathbf{v}).$$

$$d\mathbf{H}(t, \mathbf{v}) = \sum_{j=1}^{m_A} \sum_{k=1, k\neq j}^{m_A} \mathbf{d}_{jk}^A \sqrt{f_{jk}^A(t, \mathbf{v})} dW_{jk}^A(t) + \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \sqrt{f_{jk}^I(t, \mathbf{v})} dW_{jk}^I(t)$$
$$+ \sum_{i=1}^{m_S} \sum_{l=1, l\neq i}^{m_S} \mathbf{d}_{il}^S \sqrt{f_{il}^S(t, \mathbf{v})} dW_{il}^S(t) + \sum_{i=1}^{m_S} \mathbf{d}_i^D \sqrt{f_i^D(t, \mathbf{v})} dW_i^D(t).$$

$$\mathbf{G}(t, \mathbf{v}) = \sum_{j=1}^{m_A} \sum_{k=1, k\neq j}^{m_A} \mathbf{d}_{jk}^A \mathbf{d}_{jk}^{A\,\prime} f_{jk}^A(t, \mathbf{v}) + \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \mathbf{d}_{jk}^{I\,\prime} f_{jk}^I(t, \mathbf{v}) + \sum_{i=1}^{m_S} \sum_{l=1, l\neq i}^{m_S} \mathbf{d}_{il}^S \mathbf{d}_{il}^{S\prime} f_{il}^S(t, \mathbf{v})$$
$$+ \sum_{i=1}^{m_S} \mathbf{d}_i^D \mathbf{d}_i^{D\prime} f_i^D(t, \mathbf{v}).$$

With the definitions above, we rewrite equations (3.4)-(3.6) in a vector form as follows:

$$\mathbf{V}(t) = \mathbf{V}(0) + \sum_{j=1}^{m_A} \sum_{k=1, k\neq j}^{m_A} \mathbf{d}_{jk}^A Y_{jk}^A \left( \int_0^t f_{jk}^A(s, \mathbf{V}(s))ds \right) + \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I Y_{jk}^I \left( \int_0^t f_{jk}^I(s, \mathbf{V}(s))ds \right)$$
$$+ \sum_{i=1}^{m_S} \sum_{l=1, l\neq i}^{m_S} \mathbf{d}_{il}^S Y_{il}^S \left( \int_0^t f_{il}^S(s, \mathbf{V}(s))ds \right) + \sum_{i=1}^{m_S} \mathbf{d}_i^D Y_i^D \left( \int_0^t f_i^D(s, \mathbf{V}(s))ds \right).$$

Following the procedure of the uniform acceleration in Mandelbaum et al. [27] and Kurtz [23], we define a sequence of processes $\{\mathbf{V}^\eta(t), \eta \geq 1, t \geq 0\}$, where

$$\mathbf{V}^\eta(t) = \mathbf{V}^\eta(0) + \sum_{j=1}^{m_A} \sum_{k=1, k\neq j}^{m_A} \mathbf{d}_{jk}^A Y_{jk}^A \left( \eta \int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s))ds \right)$$
$$+ \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I Y_{jk}^I \left( \eta \int_0^t f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s))ds \right) + \sum_{i=1}^{m_S} \sum_{l\neq i}^{m_S} \mathbf{d}_{il}^S Y_{il}^S \left( \eta \int_0^t f_{il}^S(s, \bar{\mathbf{V}}^\eta(s))ds \right)$$
$$+ \sum_{i=1}^{m_S} \mathbf{d}_i^D Y_i^D \left( \eta \int_0^t f_i^D(s, \bar{\mathbf{V}}^\eta(s))ds \right).$$

10

## 4.1 Fluid Limit Theorem

Then, we have the following proposition for the fluid limit:

**Theorem 4.1.** *Suppose* $\mathbf{v}^\eta(0) \to \mathbf{v}(0)$ *as* $\eta \to \infty$, *then*

$$\lim_{\eta \to \infty} \frac{\mathbf{V}^\eta(t)}{\eta} = \mathbf{v}(t) \text{ almost surely,}$$

*where* $\mathbf{v}(t)$ *is the solution to the following system of ordinary differential equations:*

$$\frac{d}{dt}\mathbf{v}(t) = \sum_{j=1}^{m_A}\sum_{k\neq j}^{m_A}\mathbf{d}_{jk}^A f_{jk}^A(t,\mathbf{v}(t)) + \sum_{j=1}^{m_A}\sum_{k\neq j}^{m_A}\mathbf{d}_{jk}^I f_{jk}^I(t,\mathbf{v}(t)) \tag{4.1}$$

$$+ \sum_{i=1}^{m_S}\sum_{l\neq i}^{m_S}\mathbf{d}_{il}^S f_{il}^S(t,\mathbf{v}(t)) + \sum_{i=1}^{m_S}\mathbf{d}_i^D f_i^D(t,\mathbf{v}(t)).$$

*Proof.* See Appendix. $\qquad\qquad\square$

## 4.2 Diffusion Limit Theorem

Now that we have the fluid limit, $\mathbf{v}(t)$, we can derive the diffusion limit as follows:

**Theorem 4.2.** *Let* $\mathbf{D}^\eta(t) = \sqrt{\eta}(\mathbf{V}^\eta(t)/\eta - \mathbf{v}(t))$, *then we have that*

$$\lim_{\eta \to \infty} \mathbf{D}^\eta(t) = \mathbf{D}(t) \text{ in distribution,}$$

*where* $\mathbf{D}(t)$ *is the solution to the following stochastic differential equation*

$$d\mathbf{D}(t) = d\mathbf{H}(t,\mathbf{v}(t)) + \partial\mathbf{F}(t,\mathbf{v}(t))\mathbf{D}(t)dt,$$

*and* $\partial\mathbf{F}(t,\mathbf{v})$ *is the gradient matrix of* $\mathbf{F}(t,\mathbf{v})$ *with respect to* $\mathbf{v}$. *If* $\mathbf{D}(0)$ *is a constant or normally distributed, then* $\{\mathbf{D}(t), t \geq 0\}$ *is a Gaussian process (Arnold [2]).*

*Proof.* See Appendix. $\qquad\qquad\square$

Now that we have fluid and diffusion limits for the queue length process, we can therefore, for a large $\eta$, give an approximation for the original model as

$$\mathbf{V}^\eta(t) \approx \eta\mathbf{v}(t) + \sqrt{\eta}\mathbf{D}(t).$$

One should note that by increasing $\eta$ also implies that we are effectively increasing the number of servers along with other parameters (Mandelbaum et al. [28]). Therefore, if the number of servers is sufficiently large in the original setting (i.e., $\eta = 1$), we can approximate $\mathbf{V}(t)$ as follows:

$$\mathbf{V}(t) \approx \mathbf{v}(t) + \mathbf{D}(t).$$

Since $\{\mathbf{D}(t), t \geq 0\}$ is a Gaussian process, $\{\mathbf{V}(t), t \geq 0\}$ is approximately a Gaussian process. If we have the mean vector and the covariance matrix of $\mathbf{D}(t)$, we can approximately identify the queue length distributions as follows:

**Proposition 4.3** (Mean and covariance matrix of $\mathbf{D}(t)$, Arnold [2]). *Let* $\mathbf{M}(t) = E[\mathbf{D}(t)]$ *and* $\boldsymbol{\Sigma}(t) = \mathrm{Cov}[\mathbf{D}(t), \mathbf{D}(t)]$. *Then,* $\mathbf{M}(t)$ *and* $\boldsymbol{\Sigma}(t)$ *are the unique solution to the following ordinary equations:*

$$\frac{d}{dt}\mathbf{M}(t) = \partial\mathbf{F}(t, \mathbf{v}(t))\mathbf{M}(t), \tag{4.2}$$

$$\frac{d}{dt}\boldsymbol{\Sigma}(t) = \partial\mathbf{F}(t, \mathbf{v}(t))\boldsymbol{\Sigma}(t) + \boldsymbol{\Sigma}(t)\partial\mathbf{F}(t, \mathbf{v}(t))' + \mathbf{G}(t, \mathbf{v}(t)). \tag{4.3}$$

*If* $\mathbf{M}(0) = \mathbf{0}$, $\mathbf{M}(t) = \mathbf{0}$ *for all* $t \geq 0$.

Recall that we start with an empty queue, which implies that we do not have to solve equation (4.2), i.e., $\mathbf{M}(t) = \mathbf{0}$ for all $t \geq 0$.

By solving differential equations (4.1) and (4.3), we can approximate $E[\mathbf{V}(t)]$ and $\mathrm{Cov}[\mathbf{V}(t), \mathbf{V}(t)]$ as follows:

$$E[\mathbf{V}(t)] \approx \mathbf{v}(t),$$
$$\mathrm{Cov}[\mathbf{V}(t), \mathbf{V}(t)] \approx \boldsymbol{\Sigma}(t).$$

Let $X(t)$ be the number of customers in the system at time $t$. Then,

$$X(t) = \sum_{i=1}^{m_S} X_i(t).$$

Note that $\{X(t), t \geq 0\}$ is a Gaussian process and we can obtain the mean and variance of $X(t)$ as follows:

$$E[X(t)] = \sum_{i=1}^{m_S} E[X_i(t)],$$

$$\mathrm{Var}[X(t)] = \sum_{i=1}^{m_S} \mathrm{Var}[X_i(t)] + 2\sum_{i=1}^{m_S-1}\sum_{l=i+1}^{m_S} \mathrm{Cov}[X_i(t), X_l(t)].$$

## 4.3 Probability of Delay

Now armed with our fluid and diffusion approximations, we can also approximate other performance measures other than the mean and variance of the queue length process. One of the most important performance measures is the probability of delay or the probability that a customer must wait for service when they arrive to the queue i.e.

$$\mathbb{P}(Delay) = \mathbb{P}(W(t) > 0) \tag{4.4}$$

where $W(t)$ is the waiting time of a customer that joins the queue at time t. Thus, given our fluid and diffusion approximations for the mean and variance of the queue length we can

derive a Gaussian approximation for the probability of delay as

$$\mathbb{P}(Delay) = \mathbb{P}(W(t) > 0) \tag{4.5}$$

$$\approx \mathbb{P}(\mathbf{V}(t) \geq n(t)) \tag{4.6}$$

$$\approx \mathbb{P}\left( \tilde{Z} \geq \frac{n(t) - \mathbf{v}(t)}{\sqrt{\Sigma(t)}} \right) \tag{4.7}$$

$$\approx \overline{\Phi}\left( \frac{n(t) - \mathbf{v}(t)}{\sqrt{\Sigma(t)}} \right) \tag{4.8}$$

where $\tilde{Z}$ is a standard Gaussian random variable, $\mathbf{v}(t)$ is the fluid limit mean, and $\sqrt{\Sigma(t)}$ is the standard deviation of the diffusion limit.

## 5   The Infinite Server Case

In this section, we demonstrate that we can also apply our fluid and diffusion limits in the infinite server setting as well. This provides first and second order approximations for the queue length process that was first studied by Nelson and Taaffe [32]. However, we rigorously justify our approximations by limit theorems.

### 5.1   Infinite Server Representation

In the infinite server setting we have the following representation for the queue length process,

$$
\begin{aligned}
U_j(t) &= U_j(0) + \sum_{k\neq j}^{m_A} Y_{kj}^A\left( \int_0^t \lambda_{kj} U_k(s)ds \right) - \sum_{k\neq j}^{m_A} Y_{jk}^A\left( \int_0^t \lambda_{jk} U_j(s)ds \right) \\
&\quad - \sum_{k\neq j}^{m_A} Y_{jk}^I\left( \int_0^t \lambda_{j0}\alpha_k U_j(s)ds \right) + \sum_{k\neq j}^{m_A} Y_{kj}^I\left( \int_0^t \lambda_{k0}\alpha_j U_k(s)ds \right) \text{ for } 1 \leq j \leq m_A, \\
X_1(t) &= \sum_{j=1}^{m_A}\sum_{k=1}^{m_A} Y_{jk}^I\left( \int_0^t \lambda_{j0}\alpha_k U_j(s)ds \right) + \sum_{l\neq 1}^{m_S} Y_{l1}^S\left( \int_0^t \mu_{l1} X_l(s)ds \right) \\
&\quad - \sum_{l\neq 1}^{m_S} Y_{1l}^S\left( \int_0^t \mu_{1l}\Big[ (X_1(s)) \Big]ds \right) - Y_1^D\left( \int_0^t \mu_{10}\Big[ (X_1(s)) \Big]ds \right), \\
X_i(t) &= Y_{1i}^S\left( \int_0^t \mu_{1i}\Big[ (X_1(s)) \Big]ds \right) + \sum_{l=2,l\neq i}^{m_S} Y_{li}^S\left( \int_0^t \mu_{li} X_l(s)ds \right) \\
&\quad - \sum_{l\neq i}^{m_S} Y_{il}^S\left( \int_0^t \mu_{il} X_i(s)ds \right) - Y_i^D\left( \int_0^t \mu_{i0} X_i(s)ds \right) \qquad \text{for } 2 \leq i \leq m_S.
\end{aligned}
$$

The major difference between the finite and infinite server settings is the rate functions for the $X_i$ Poisson processes. In the finite setting, at most $n_t$ customers can be processed at any

13

time $t$, however, in the infinite server setting, this is no longer a limitation. Thus, all of the rate functions with the terms $X_1 \wedge \left(n - \sum_{r=2}^{m_S} X_r(s)\right)^+$ since the term $\left(n - \sum_{r=2}^{m_S} X_r(s)\right)^+$ is equal to $\infty$.

## 5.2 Infinite Server Fluid Limit Theorem

We have the following proposition for the fluid limit for the $Ph_t/Ph_t/\infty$ queue

**Proposition 5.1.** *Suppose $\mathbf{v}^{\eta,\infty}(0) \to \mathbf{v}^\infty(0)$ as $\eta \to \infty$, then*

$$\lim_{\eta \to \infty} \frac{\mathbf{V}^{\eta,\infty}(t)}{\eta} = \mathbf{v}^\infty(t) \text{ almost surely,}$$

*where $\mathbf{v}^\infty(t)$ is the solution to the following system of ordinary differential equations:*

$$\frac{d}{dt}\mathbf{v}^\infty(t) = \sum_{j=1}^{m_A}\sum_{k \neq j}^{m_A}\mathbf{d}_{jk}^A f_{jk}^A(t, \mathbf{v}^\infty(t)) + \sum_{j=1}^{m_A}\sum_{k \neq j}^{m_A}\mathbf{d}_{jk}^I f_{jk}^I(t, \mathbf{v}^\infty(t))$$
$$+ \sum_{i=1}^{m_S}\sum_{l \neq i}^{m_S}\mathbf{d}_{il}^S f_{il}^S(t, \mathbf{v}^\infty(t)) + \sum_{i=1}^{m_S}\mathbf{d}_i^D f_i^D(t, \mathbf{v}^\infty(t))$$

*where the rate functions correspond to the infinite server representation given in Section 5.1.*

*Proof.* The proof of this result immediately follows from the proof of the finite case and setting $n = \infty$. $\square$

## 5.3 Infinite Server Diffusion Limit Theorem

Now that we have the fluid limit, $\mathbf{v}^\infty(t)$, we can derive the diffusion limit as follows:

**Proposition 5.2.** *Let $\mathbf{D}^{\eta,\infty}(t) = \sqrt{\eta}(\mathbf{V}^{\eta,\infty}(t)/\eta - \mathbf{v}^\infty(t))$, then we have that*

$$\lim_{\eta \to \infty} \mathbf{D}^{\eta,\infty}(t) = \mathbf{D}^\infty(t) \text{ in distribution,}$$

*where $\mathbf{D}^\infty(t)$ is the solution to the following stochastic differential equation*

$$d\mathbf{D}^\infty(t) = \mathbf{H}(t, \mathbf{v}^\infty(t)) + \partial\mathbf{F}(t, \mathbf{v}^\infty(t))\mathbf{D}^\infty(t)dt,$$

*and $\partial\mathbf{F}(t, \mathbf{v}^\infty(t))$ is the gradient matrix of $\mathbf{F}(t, \mathbf{v}^\infty(t))$ with respect to $\mathbf{v}^\infty(t)$ where the rate functions correspond to the infinite server representation given in Section 5.1.*

*Proof.* The proof of this result immediately follows from the proof of the finite case and setting $n = \infty$. $\square$

Figure 6.1: Overall flow of the numerical study

# 6 Numerical results

In this section, we provide some numerical results comparing the proposed method with the simulation results. Referring to the flow chart in Figure 6.1, we choose Coxian distributions to approximate Weibull and lognormal distributions for inter-arrival and service times. Coxian distributions have a unique initial state that the proposed method requires and the overall fitting quality is known to be good (Asmussen et al. [3]). We use the EM algorithm developed by Asmussen et al. [3], although other phase-type distributions and fitting algorithms can also be used. Since we want to approximate the distribution itself, we use 8-10 phases to fit the target distributions accurately. Figure 6.2 illustrates a density and distribution fitting with a Coxian distribution. In this example, we use 10 phases to approximate the Weibull distribution. We derive the ordinary differential equations (ODEs) from equations (4.1) and (4.3), and solve them using MATLAB. We write the simulation code in C++. In order to generate a general time-varying arrival process, we implement the algorithm based on the standard equilibrium renewal process (SERP) explained in the longer version of Liu and Whitt [24]. We use Weibull distributions with mean 1 as a base distribution in order to generate time-varying arrival times. We run 5,000 independent instances for each setting and estimate the mean and the variance of the number of customers in the system and the probability of delay over time.

We choose two Weibull distributions having the same mean 1 for the arrival processes: the squared coefficient of variation (SCoV) of Weibull(0.79,0.7) is 2.1387 which is greater

than one, and the SCoV of Weibull(1.1271,2.5) is 0.1831 which is less than one. Time-varying rates are applied to the base distributions for constructing the actual arrival processes. We do not consider the case when the SCoV is 1 since it is an exponential distribution and has been studied extensively in the literature. For the service times, we choose two lognormal distributions with the different SCoV values. Without loss of generality, the means of two service time distributions are 1. Increasing the number of servers makes us expect more accurate estimations since the fluid and diffusion limits are asymptotically exact. Therefore, we compare the cases when the number of servers is 50 and 200. The corresponding time-varying rates to the number of servers are $45 + 30\sin(2\pi t/10)$ and $180 + 120\sin(2\pi t/10)$ respectively. Then, we have 8 combinations of experiments: two distributions for arrivals, two distributions for services, two values of the number of servers:

Exp. 1: 50 servers, SCoV of inter-arrival times $> 1$ and SCoV of service times $> 1$

- Time-varying rate: $45 + 30\sin(2\pi t/10)$
- Base inter-arrival time distribution: Weibull$(0.79, 0.7)$, SCoV $= 2.1387$
- Service time distribution: Lognormal$(-0.5, 1)$, SCoV $= 1.7183$

Exp. 2: 200 servers, SCoV of inter-arrival times $> 1$ and SCoV of service times $> 1$

- Time-varying rate: $180 + 120\sin(2\pi t/10)$
- Base inter-arrival time distribution: Weibull$(0.79, 0.7)$, SCoV $= 2.1387$
- Service time distribution: Lognormal$(-0.5, 1)$, SCoV $= 1.7183$

Exp. 3: 50 servers, SCoV of inter-arrival times $> 1$ and SCoV of service times $< 1$

- Time-varying rate: $45 + 30\sin(2\pi t/10)$
- Base inter-arrival time distribution: Weibull$(0.79, 0.7)$, SCoV $= 2.1387$
- Service time distribution: Lognormal$(-0.2027, 0.6368)$, SCoV $= 0.5$

Exp.4: 200 servers, SCoV of inter-arrival times $> 1$ and SCoV of service times $< 1$

- Time-varying rate: $180 + 120\sin(2\pi t/10)$
- Base inter-arrival time distribution: Weibull$(0.79, 0.7)$, SCoV $= 2.1387$
- Service time distribution: Lognormal$(-0.2027, 0.6368)$, SCoV $= 0.5$

Exp. 5: 50 servers, SCoV of inter-arrival times $< 1$ and SCoV of service times $> 1$

- Time-varying rate: $45 + 30\sin(2\pi t/10)$
- Base inter-arrival time distribution: Weibull$(1.1271, 2.5)$, SCoV $= 0.1831$
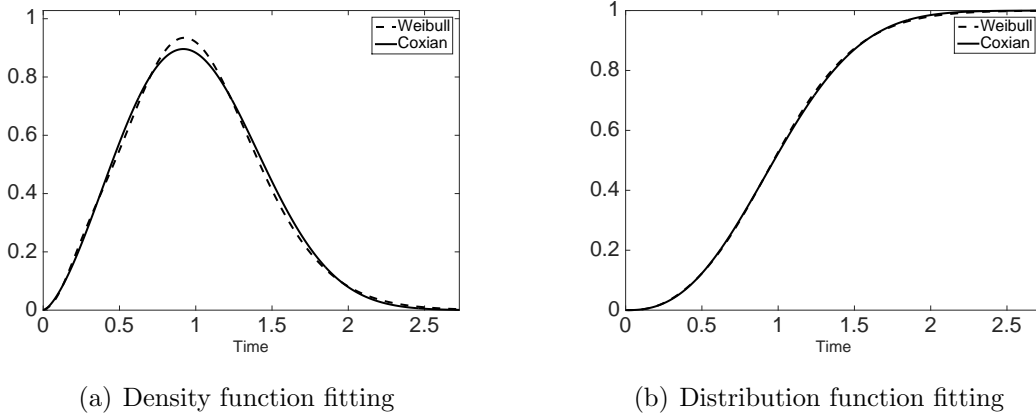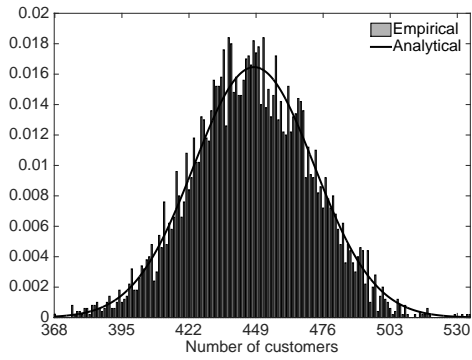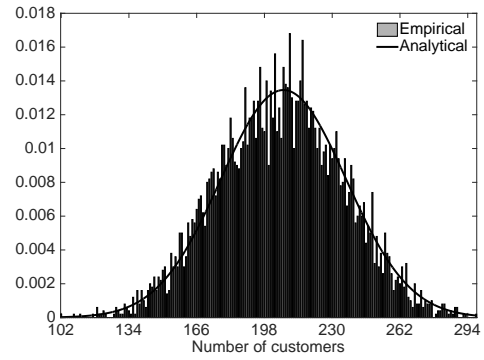- Service time distribution: Lognormal$(-0.5, 1)$, SCoV $= 1.7183$

Exp. 6: 200 servers, SCoV of inter-arrival times $> 1$ and SCoV of service times $> 1$

- Time-varying rate: $180 + 120\sin(2\pi t/10)$

    – Base inter-arrival time distribution: Weibull$(1.1271, 2.5)$, SCoV $= 0.1831$

    – Service time distribution: Lognormal$(-0.5, 1)$, SCoV $= 1.7183$

Exp. 7: 50 servers, SCoV of inter-arrival times $< 1$ and SCoV of service times $< 1$

    – Time-varying rate: $45 + 30\sin(2\pi t/10)$

    – Base inter-arrival time distribution: Weibull$(1.1271, 2.5)$, SCoV $= 0.1831$

    – Service time distribution: Lognormal$(-0.2027, 0.6368)$, SCoV $= 0.5$

Exp. 8: 200 servers, SCoV of inter-arrival times $> 1$ and SCoV of service times $> 1$

    – Time-varying rate: $180 + 120\sin(2\pi t/10)$

    – Base inter-arrival time distribution: Weibull$(1.1271, 2.5)$, SCoV $= 0.1831$

    – Service time distribution: Lognormal$(-0.2027, 0.6368)$, SCoV $= 0.5$



(a) Density function fitting        (b) Distribution function fitting

Figure 6.2: Weibull(1.1271, 2.5) and corresponding Coxian distributions.

We mention that the queue length distributions are approximately Gaussian in Section 4. Figure 6.3 compares the empirical density and the density from the diffusion limit at several time points (underloaded times 5 and 10, critically loaded times 7.5 and 17.5 and overloaded times 5 and 15). Although we observe some skewness in the empirical density, the Gaussian approximation seems to work well.

Figures 6.4-6.7 plot the mean and the variance of the number of customers and the probability of delay over time comparing the proposed method and the simulation results for the cases of 50 and 200 servers. Each figure represents a different combination of distributions for arrival processes and service times. Overall we observe that the proposed method provides accurate estimations of the mean and the variance of the number of customers and the probability of delay. Comparing Figures 6.4 (a) and (b), we observe that increasing the number of servers results in more accurate estimations of the mean as expected. We observe the same result for the variance (Figures 6.4 (c) and (d)) and the probability of

17

(a) Density at $t = 5$

(b) Density at $t = 7.5$

(c) Density at $t = 10$
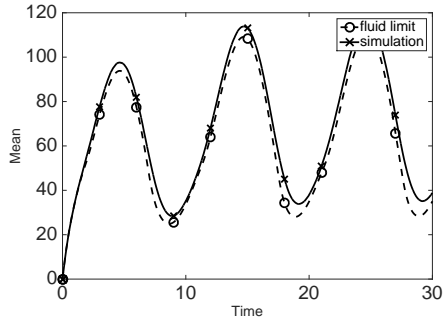
(d) Density at $t = 15$
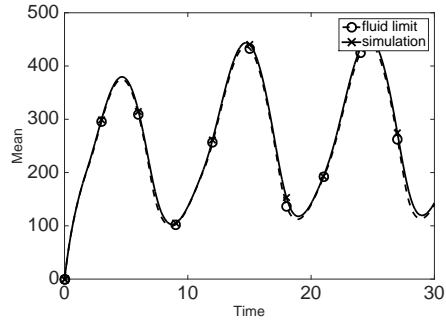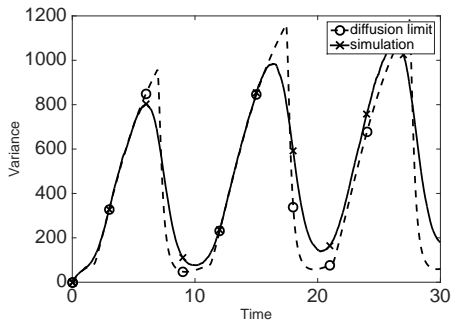
(e) Density at $t = 17.5$

(f) Density at $t = 20$

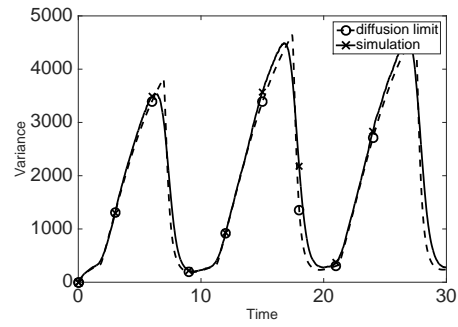Figure 6.3: Density of the number of customers at time 5, 7.5, 10, 15, 17.5 and 20.

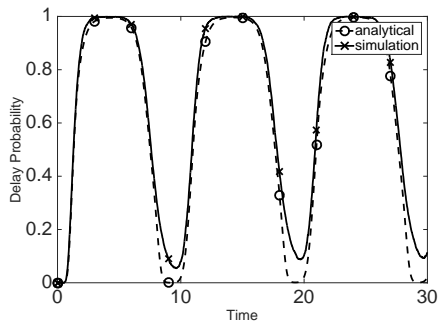(a) Mean number of customers, Exp. 1
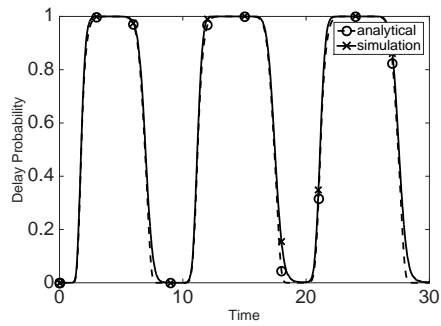


(b) Mean number of customers, Exp. 2



(c) Variance of the number of customers, Exp. 1



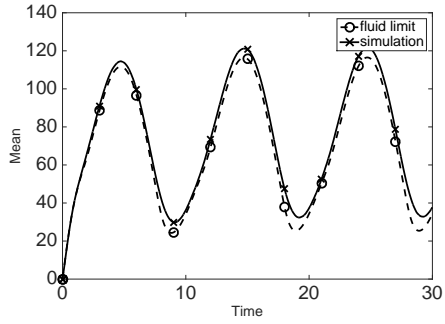(d) Variance of the number of customers, Exp. 2
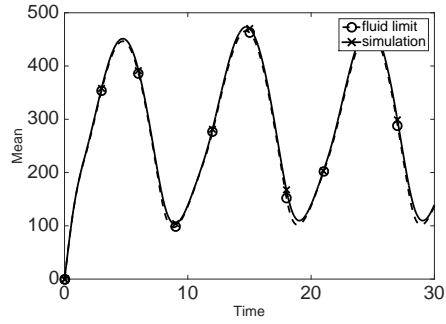


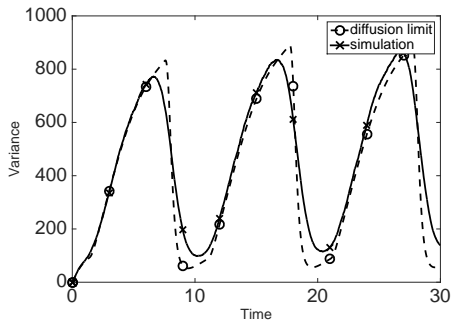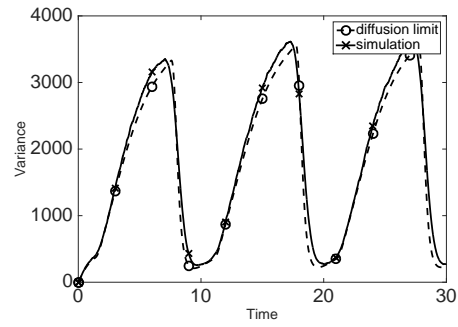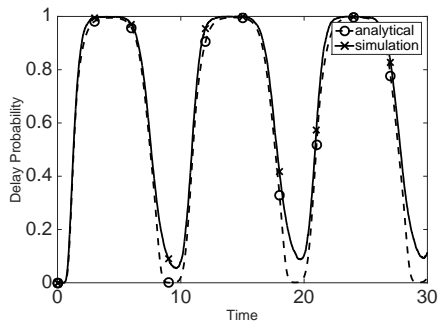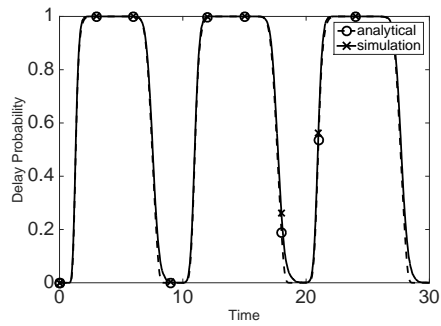(e) Delay probability, Exp. 1



(f) Delay probability, Exp. 2

Figure 6.4: Comparison between Exp. 1 and Exp. 2

(a) Mean number of customers, Exp. 3



(b) Mean number of customers, Exp. 4



(c) Variance of the number of customers, Exp. 3



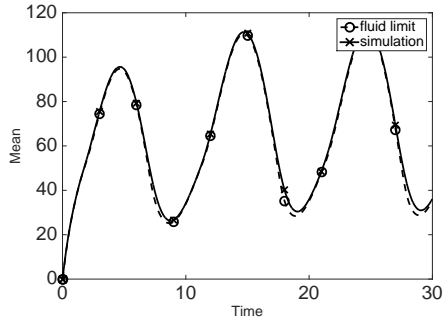(d) Variance of the number of customers, Exp. 4
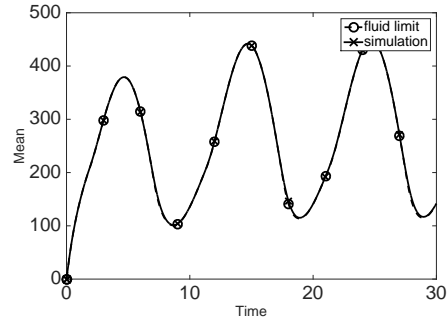


(e) Delay probability, Exp. 3
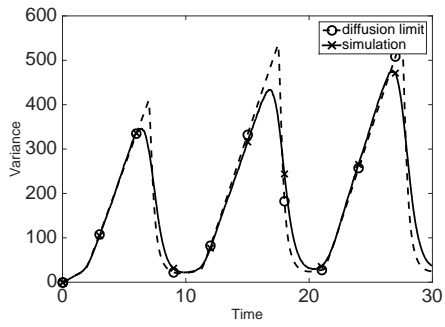


(f) Delay probability, Exp. 4

Figure 6.5: Comparison between Exp. 3 and Exp. 4

(a) Mean number of customers, Exp. 5



(b) Mean number of customers, Exp. 6



(c) Variance of the number of customers, Exp. 5



(d) Variance of the number of customers, Exp. 6



(e) Delay probability, Exp. 5



(f) Delay probability, Exp. 6

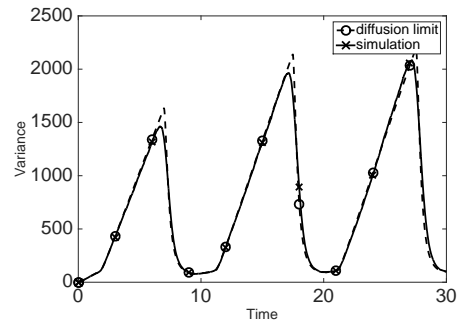Figure 6.6: Comparison between Exp. 5 and Exp. 6
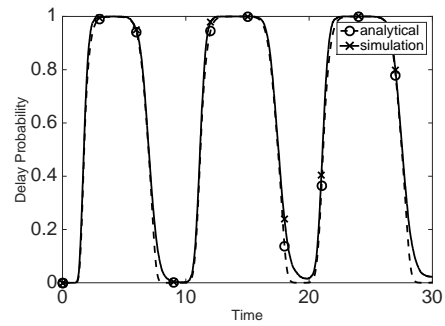
(a) Mean number of customers, Exp. 7
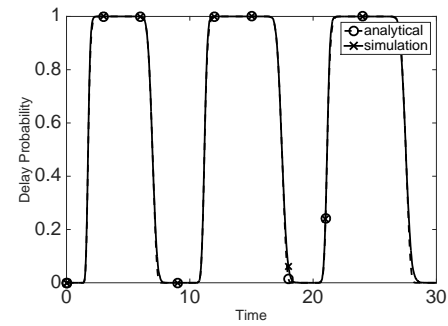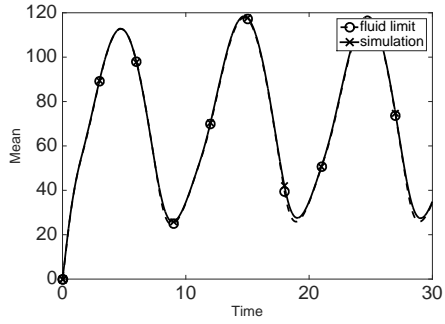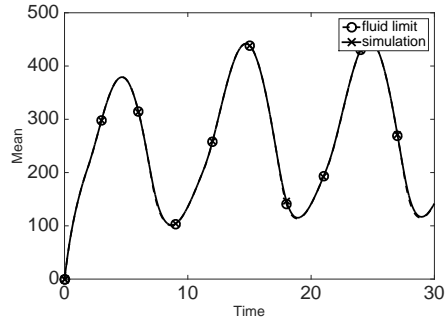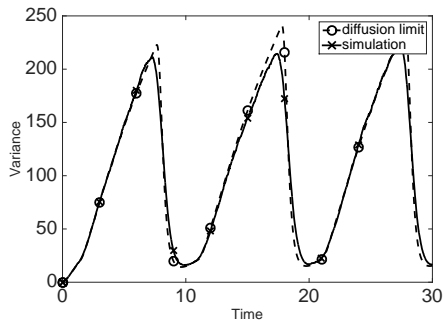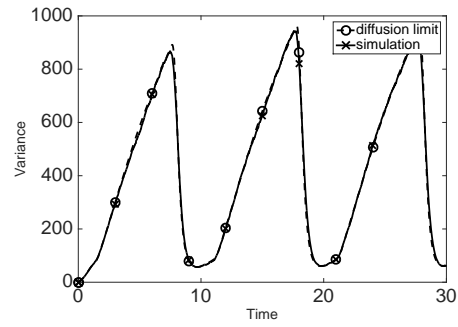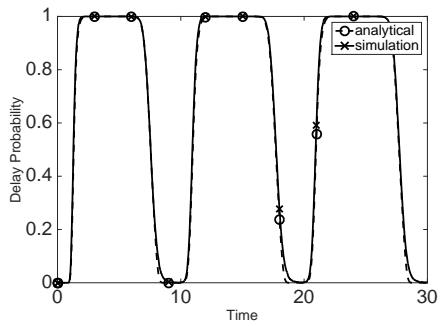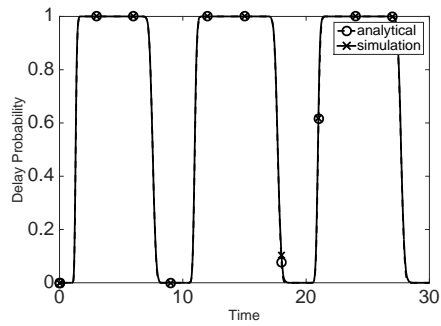
(b) Mean number of customers, Exp. 8

(c) Variance of the number of customers, Exp. 7

(d) Variance of the number of customers, Exp. 8

(e) Delay probability, Exp. 7

(f) Delay probability, Exp. 8

Figure 6.7: Comparison between Exp. 7 and Exp. 8

delay (Figures 6.4 (e) and (f)). The same results hold across different distribution settings (Figures 6.5-6.7). The distributions in Figure 6.4 have the largest SCoV values and those in Figure 6.7 have the smallest SCoV values. In Figures 6.4 and 6.7, we observe that the proposed method works better when the SCoV values are small.

# 7 Conclusion

This paper describes a new methodology to approximate the queue length distributions of large-scale $G_t/G_t/n_t$ queues. Instead of analyzing a $G_t/G_t/n_t$ directly, we study a $Ph_t/Ph_t/n_t$ queue since phase-type distributions can approximate positive-valued distributions in any level of accuracy. Applying the uniform acceleration and strong approximations to $Ph_t/Ph_t/n_t$ queues to obtain fluid and diffusion limits, we encounter the lingering problem in our formulation and cannot obtain the diffusion limit. To resolve the issue, we propose a new formulation with an additional condition that is not quite restrictive. The new formulation works well and we successfully derive the fluid and diffusion limits. We find that the queue length process is approximately a Gaussian process and we derive ordinary differential equations to obtain the mean and variance of the queue length over time.

From the numerical study, we observe that the proposed method works better when the distributions for arrival processes and service times have smaller SCoVs. Since the uniform acceleration method increases the number of servers to infinity, the estimations should become more accurate as the number of servers increases. We exactly observe this phenomenon as expected.

We suggest two directions for future research. For example, in order to obtain the diffusion limit, we put an additional condition (a unique initial state for phase-type distributions). Although it does not seem to be critical, the method will be improved if the restriction can be removed. Extending the proposed method to multi-dimensional queueing networks is another possible research direction that we plan to pursue in a follow-up paper.

# A  Appendix

Our convergence results will be based on the following strong approximation result which allows for a pathwise approximation of a Poisson process by a standard Brownian motion and linear drift that lives on the same probability space. Before we prove the main result, we need the following lemma.

**Lemma A.1** (Kurtz 1978). *A standard Poisson process $\{\Pi(t)\}_{t \geq 0}$ can be realized on the same probability space as a standard Brownian motion $\{W(t)\}_{t \geq 0}$ in such a way that the almost surely finite random variable*

$$Z \equiv \sup_{t \geq 0} \frac{|\Pi(t) - t - W(t)|}{\log(2 \vee t)}$$

*as finite moment generating function in the neighborhood of the origin and in particular finite mean.*

**Lemma A.2** (Mandelbaum et al. 1998). *Let $x, y$, and $z$ be measurable, non-negative functions on the reals. If $y$ is bounded and $z$ is integrable on $[0, T]$ and for all $0 \leq t \leq T$,*

$$x(t) \leq z(t) + \int_0^t x(s)y(s)ds, \tag{A.1}$$

*then*

$$x(t) \leq z(t) + \int_0^t z(s)y(s) \cdot \exp\left(\int_s^t y(r)dr\right) ds \tag{A.2}$$

*and*

$$\sup_{0 \leq t \leq T} x(t) \leq \sup_{0 \leq t \leq T} z(t) \cdot \exp\left(\int_0^T y(t)dt\right). \tag{A.3}$$

## A.1  Proof of Fluid Limit

Suppose $\mathbf{v}^\eta(0) \to \mathbf{v}(0)$ as $\eta \to \infty$, then

$$\lim_{\eta \to \infty} \frac{\mathbf{V}^\eta(t)}{\eta} = \mathbf{v}(t) \text{ almost surely,}$$

where $\mathbf{v}(t)$ is the solution to the following system of ordinary differential equations:

$$\frac{d}{dt}\mathbf{v}(t) = \sum_{j=1}^{m_A}\sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^A f_{jk}^A(t, \mathbf{v}(t)) + \sum_{j=1}^{m_A}\sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^I f_{jk}^I(t, \mathbf{v}(t))$$
$$+ \sum_{i=1}^{m_S}\sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S f_{il}^S(t, \mathbf{v}(t)) + \sum_{i=1}^{m_S} \mathbf{d}_i^D f_i^D(t, \mathbf{v}(t)).$$

*Proof.* In view of the strong approximation results given in Lemma A.1 and because the rate functions of the queueing process are Lipschitz continuous,

$$\sum_{j=1}^{m_A} \sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^A \left| Y_{jk}^A \left( \eta \int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) ds \right) - \eta \int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) ds - W_{jk}^A \left( \eta \int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|$$

$$+ \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \left| Y_{jk}^A \left( \eta \int_0^t f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) ds \right) - \eta \int_0^t f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) ds - W_{jk}^A \left( \eta \int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|$$

$$+ \sum_{i=1}^{m_S} \sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S \left| Y_{il}^S \left( \eta \int_0^t f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) ds \right) - \eta \int_0^t f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) ds - W_{il}^S \left( \eta \int_0^t f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|$$

$$+ \sum_{i=1}^{m_S} \mathbf{d}_i^D \left| Y_i^D \left( \eta \int_0^t f_i^D(s, \bar{\mathbf{V}}^\eta(s)) ds \right) - \eta \int_0^t f_i^D(s, \bar{\mathbf{V}}^\eta(s)) ds - W_i^D \left( \eta \int_0^t f_i^D(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|$$

is $\Theta(\log(\eta))$ almost surely. Moreover, the $W(\eta\cdot)$ terms are standard Brownian motions. Since the rate functions are Lipschitz continuous and scalable in the sense of Mandelbaum et al. [27], we know that know that the law of the iterated logarithm for Brownian motion yields,

$$\limsup_{\eta \to \infty} \sup_{t \leq T} \frac{1}{\eta} W_{jk}^A \left( \eta \int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) ds \right) = 0 \quad \text{almost surely,}$$

$$\limsup_{\eta \to \infty} \sup_{t \leq T} \frac{1}{\eta} W_{jk}^I \left( \eta \int_0^t f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) ds \right) = 0 \quad \text{almost surely,}$$

$$\limsup_{\eta \to \infty} \sup_{t \leq T} \frac{1}{\eta} W_{il}^S \left( \eta \int_0^t f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) ds \right) = 0 \quad \text{almost surely,}$$

$$\limsup_{\eta \to \infty} \sup_{t \leq T} \frac{1}{\eta} W_i^D \left( \eta \int_0^t f_i^D(s, \bar{\mathbf{V}}^\eta(s)) ds \right) = 0 \quad \text{almost surely.}$$

This implies that

$$\sum_{j=1}^{m_A} \sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^A \frac{1}{\eta} \left| W_{jk}^A \left( \eta \int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right| + \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \frac{1}{\eta} \left| W_{jk}^I \left( \eta \int_0^t f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|$$

$$+ \sum_{i=1}^{m_S} \sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S \frac{1}{\eta} \left| W_{il}^S \left( \eta \int_0^t f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right| + \sum_{i=1}^{m_S} \mathbf{d}_i^D \frac{1}{\eta} \left| W_i^D \left( \eta \int_0^t f_i^D(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|$$

converges to zero uniformly over compact sets of time as $\eta$ goes to $\infty$. Thus, for some

constant $\tilde{C}$, we have that

$$
\left| \frac{1}{\eta} \mathbf{V}^\eta(t) - \mathbf{v}(t) \right|
$$

$$
\leq \sum_{j=1}^{m_A} \sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^A \int_0^t \left| f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^A(s, \mathbf{v}(s)) \right| ds + \sum_{j=1}^{m_A} \sum_{k=1, k \neq j}^{m_A} \mathbf{d}_{jk}^A \frac{1}{\eta} \left| W_{jk} \left( \eta \int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|
$$

$$
+ \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \int_0^t \left| f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^I(s, \mathbf{v}(s)) \right| ds + \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \frac{1}{\eta} \left| W_{jk}^I \left( \eta \int_0^t f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|
$$

$$
+ \sum_{i=1}^{m_S} \sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S \int_0^t \left| f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) - f_{il}^S(s, \mathbf{v}(s)) \right| ds + \sum_{i=1}^{m_S} \sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S \frac{1}{\eta} \left| W_{il}^S \left( \eta \int_0^t f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|
$$

$$
+ \sum_{i=1}^{m_S} \mathbf{d}_i^D \int_0^t \left| f_i^D(s, \bar{\mathbf{V}}^\eta(s)) - f_i^D(s, \mathbf{v}(s)) \right| ds + \sum_{i=1}^{m_S} \mathbf{d}_i^D \frac{1}{\eta} \left| W_i^D \left( \eta \int_0^t f_i^D(s, \bar{\mathbf{V}}^\eta(s)) ds \right) \right|
$$

$$
+ \tilde{C} \cdot \left( \frac{\log(\eta)}{\eta} \right).
$$

Thus, if we fix $\varepsilon > 0$, we have that from Lemma A.1 and the law of the iterated logarithm for Brownian motion that there exists an $\eta^* \in \mathbb{N}$ such that for all $\eta \leq \eta^*$ and uniformly on compact sets of time, for $t \leq T$ we have that

$$
\begin{aligned}
\left| \frac{1}{\eta} \mathbf{V}^\eta(t) - \mathbf{v}(t) \right| &\leq \sum_{j=1}^{m_A} \sum_{k=1, k \neq j}^{m_A} \mathbf{d}_{jk}^A \int_0^t \left| f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^A(s, \mathbf{v}(s)) \right| ds \\
&+ \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \int_0^t \left| f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^I(s, \mathbf{v}(s)) \right| ds \\
&+ \sum_{i=1}^{m_S} \sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S \int_0^t \left| f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) - f_{il}^S(s, \mathbf{v}(s)) \right| ds \\
&+ \sum_{i=1}^{m_S} \mathbf{d}_i^D \int_0^t \left| f_i^D(s, \bar{\mathbf{V}}^\eta(s)) - f_i^D(s, \mathbf{v}(s)) \right| ds + \varepsilon \\
&\leq \left| \mathbf{F}(t, \bar{\mathbf{V}}^\eta(t)) - \mathbf{F}(t, \mathbf{v}(t)) \right| + \varepsilon.
\end{aligned}
$$

Now that we have the rate functions are Lipschitz continuous functions, we know that there exists a constant $M$ such that

$$
\sup_{0 \leq t \leq T} \left| \frac{1}{\eta} \mathbf{V}^\eta(t) - \mathbf{v}(t) \right| \leq M \int_0^t \sup_{0 \leq r \leq s} \left| \frac{1}{\eta} \mathbf{V}^\eta(r) - \mathbf{v}(r) \right| ds + \varepsilon.
$$

Now by applying Gronwall's lemma, we have that

$$
\sup_{0 \leq t \leq T} \left| \frac{1}{\eta} \mathbf{V}^\eta(t) - \mathbf{v}(t) \right| \leq \varepsilon e^{MT}
$$

which implies our fluid limit result since $\varepsilon$ was arbitrarily chosen. $\qquad \square$

## A.2 Proof of Diffusion Limit

Let $\mathbf{D}^\eta(t) = \sqrt{\eta}(\mathbf{V}^\eta(t)/\eta - \mathbf{v}(t))$, then we have that

$$\lim_{\eta \to \infty} \mathbf{D}^\eta(t) = \mathbf{D}(t) \text{ in distribution,}$$

where $\mathbf{D}(t)$ is the solution to the following stochastic differential equation

$$d\mathbf{D}(t) = \mathbf{H}(t, \mathbf{v}(t)) + \partial \mathbf{F}(t, \mathbf{v}(t))\mathbf{D}(t)dt, \tag{A.4}$$

Via the theory of strong approximations given in Lemma A.1, we can now represent $\mathbf{D}^\eta(t)$ by the following equality

$$
\begin{aligned}
\mathbf{D}^\eta(t) &= \sqrt{\eta}\left(\frac{1}{\eta}\mathbf{V}^\eta(t) - \mathbf{v}(t)\right) \\
&= \sqrt{\eta} \cdot \left(\mathbf{F}(t, \bar{\mathbf{V}}^\eta(t)) - \mathbf{F}(t, \mathbf{v})\right) + Z^\eta(t) + \Theta\left(\frac{\log(\eta)}{\sqrt{\eta}}\right) \\
&= \sum_{j=1}^{m_A}\sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^A \int_0^t \sqrt{\eta}\left(f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^A(s, \mathbf{v}(s))\right)ds \\
&+ \sum_{j=1}^{m_A}\sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \int_0^t \sqrt{\eta}\left(f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^I(s, \mathbf{v}(s))\right)ds \\
&+ \sum_{i=1}^{m_S}\sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S \int_0^t \sqrt{\eta}\left(f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) - f_{il}^S(s, \mathbf{v}(s))\right)ds \\
&+ \sum_{i=1}^{m_S} \mathbf{d}_i^D \int_0^t \sqrt{\eta}\left(f_i^D(s, \bar{\mathbf{V}}^\eta(s)) - f_i^D(s, \mathbf{v}(s))\right)ds + Z^\eta(t) + \Theta\left(\frac{\log(\eta)}{\sqrt{\eta}}\right).
\end{aligned}
$$

where

$$
\begin{aligned}
Z^\eta(t) &= \sum_{j=1}^{m_A}\sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^A B_{jk}^A\left(\int_0^t f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s))ds\right) + \sum_{j=1}^{m_A}\sum_{k=1}^{m_A} \mathbf{d}_{jk}^I B_{jk}^I\left(\int_0^t f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s))ds\right) \\
&+ \sum_{i=1}^{m_S}\sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S B_{il}^S\left(\int_0^t f_{il}^S(s, \bar{\mathbf{V}}^\eta(s))ds\right) + \sum_{i=1}^{m_S} \mathbf{d}_i^D B_i^D\left(\int_0^t f_i^D(s, \bar{\mathbf{V}}^\eta(s))ds\right).
\end{aligned}
$$

**Lemma A.3.** *The sequence of stochastic processes $Z^\eta(t)$ converges in distribution to the process $Z(t)$ where*

$$
\begin{aligned}
Z(t) &= \sum_{j=1}^{m_A}\sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^A B_{jk}^A\left(\int_0^t f_{jk}^A(s, \mathbf{v}(s))ds\right) + \sum_{j=1}^{m_A}\sum_{k=1}^{m_A} \mathbf{d}_{jk}^I B_{jk}^I\left(\int_0^t f_{jk}^I(s, \mathbf{v}(s))ds\right) \tag{A.5} \\
&+ \sum_{i=1}^{m_S}\sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S B_{il}^S\left(\int_0^t f_{il}^S(s, \mathbf{v}(s))ds\right) + \sum_{i=1}^{m_S} \mathbf{d}_i^D B_i^D\left(\int_0^t f_i^D(s, \mathbf{v}(s))ds\right).
\end{aligned}
$$

*Proof.* For this proof, we will use the Holder continuity of Brownian motion. We know that for any $\alpha \in (0, 1/2)$ and $T > 0$, there exists an integrable and hence almost surely finite random variable $M$ such that

$$|B(t_2) - B(t_1)| \leq M|t_2 - t_1|^\alpha$$

almost surely for all $t_1, t_2 \leq T$. Therefore, using the fluid limit and the Lipschitz continuity of the rate functions along with the Holder continuity of the Brownian motion, we have that $Z^\eta(t)$ converges in distribution to the process $Z(t)$. □

The following lemma shows that the sequence $D^\eta(t)$ is bounded in probability.

**Lemma A.4.** *For an $\epsilon > 0$, there exists $\eta^* \in \mathbb{N}$ and $K < \infty$ such that*

$$\mathbb{P}\left(\sup_{0 \leq t \leq T} |D^\eta(t)| > K\right) < \epsilon \quad \text{for all } \eta \geq \eta^*. \tag{A.6}$$

*Proof.* The strong approximation for the Brownian motion yields the following representation

$$
\begin{aligned}
\mathbf{D}^\eta(t) &= \sqrt{\eta}\left(\frac{1}{\eta}\mathbf{V}^\eta(t) - \mathbf{v}(t)\right) + Z^\eta(t) + \tilde{C} \cdot \frac{\log \eta}{\sqrt{\eta}} \\
&= \sum_{j=1}^{m_A} \sum_{k=1, k \neq j}^{m_A} \mathbf{d}_{jk}^A \int_0^t \sqrt{\eta}\left(f_{jk}^A(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^A(s, \mathbf{v}(s))\right) ds \\
&\quad + \sum_{j=1}^{m_A} \sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \int_0^t \sqrt{\eta}\left(f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^I(s, \mathbf{v}(s))\right) ds \\
&\quad + \sum_{i=1}^{m_S} \sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S \int_0^t \sqrt{\eta}\left(f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) - f_{il}^S(s, \mathbf{v}(s))\right) ds \\
&\quad + \sum_{i=1}^{m_S} \mathbf{d}_i^D \int_0^t \sqrt{\eta}\left(f_i^D(s, \bar{\mathbf{V}}^\eta(s)) - f_i^D(s, \mathbf{v}(s))\right) ds + Z^\eta(t) + \tilde{C} \cdot \frac{\log \eta}{\sqrt{\eta}}.
\end{aligned}
$$

We know that $Z^\eta(t)$ is tight by Lemma A.3 and hence is bounded in probability. Moreover, by using the Lipschitz continuity of the rate functions we have that

$$\sup_{0 \leq t \leq T} \mathbf{D}^\eta(t) \leq M \int_0^T \sup_{0 \leq t \leq s} \mathbf{D}^\eta(s) ds + \sup_{0 \leq t \leq T} Z^\eta(t) + \varepsilon$$

for some Lipschitz constant M. Thus, by Gronwall's inequality we have almost surely that

$$\sup_{0 \leq t \leq T} \mathbf{D}^\eta(t) \leq e^{MT} \sup_{0 \leq t \leq T} (Z^\eta(t) + \varepsilon)$$

and this concludes the proof. □

**Lemma A.5.** *If $f^\eta = \{f_t^\eta\}_{t \geq 0}$ be a sequence of non-negative random processes such that*

$$\lim_{\eta \to \infty} \int_0^T f^\eta(u) du = 0 \quad \text{in probability}, \tag{A.7}$$

*then, for all $\delta > 0$,*

$$\lim_{\eta \to \infty} \mathbb{P}\left(\sup_{0 \leq t \leq T} \left|\int_0^T f^\eta(u)\mathbf{D}^\eta(u) du\right| > \delta\right) = 0. \tag{A.8}$$

28

*Proof.* If we fix $\epsilon > 0$, then we know that there exists a constant $\eta^* \in \mathbb{N}$ such that for all all $\eta > \eta^*$, there exists sets $\Omega_{\eta,1}$ and $\Omega_{\eta,2}$ such that

$$\int_0^T f^\eta(u)du < \epsilon/2 \quad \text{on } \Omega_{\eta,1} \text{ and such that } \mathbb{P}(\Omega_{\eta,1}) \geq 1 - \epsilon/2, \tag{A.9}$$

and

$$\sup_{0 \leq t \leq T} |\mathbf{D}^\eta(t)| < K \quad \text{on } \Omega_{\eta,2} \text{ and such that } \mathbb{P}(\Omega_{\eta,2}) \geq 1 - \epsilon/2, \tag{A.10}$$

Therefore, we have that

$$\sup_{0 \leq t \leq T} \left| \int_0^T f^\eta(u)\mathbf{D}^\eta(u)du \right| \leq \sup_{0 \leq t \leq T} |\mathbf{D}^\eta(t)| \int_0^T f^\eta(u) < K\epsilon \quad \text{on } \Omega_{\eta,1} \cap \Omega_{\eta,2}. \tag{A.11}$$

This concludes the proof of the lemma. $\qquad\square$

**Theorem A.6** (Proof of Diffusion Limit). *Let us first defined a sequence of stochastic processes*

$$\tilde{\mathbf{D}}^\eta(t) \equiv \int_0^t \partial\mathbf{F}(s, \mathbf{v}(s))\tilde{\mathbf{D}}^\eta(s)ds + Z^\eta(t). \tag{A.12}$$

*By the continuous mapping theorem and Lemma A.3, which shows that $Z^\eta(t)$ converges to $Z(t)$ in Equation A.5, then we know that that $\tilde{\mathbf{D}}^\eta(t)$ converges to $\tilde{\mathbf{D}}(t)$ given in Equation A.4. It now suffices to show that*

$$\lim_{\eta \to \infty} \sup_{0 \leq t \leq T} |\mathbf{D}^\eta(t) - \tilde{\mathbf{D}}^\eta(t)| = 0 \quad \text{in probability.} \tag{A.13}$$

*To prove this, we will let*

$$E^\eta(t) \equiv \mathbf{D}^\eta(t) - \tilde{\mathbf{D}}^\eta(t).$$

*Therefore, from the definition of $\tilde{\mathbf{D}}^\eta(t)$ and the representation of $\mathbf{D}^\eta(t)$ we obtain the following equality for $E^\eta(t)$*

$$
\begin{aligned}
E^\eta(t) &= \sum_{j=1}^{m_A}\sum_{k \neq j}^{m_A} \mathbf{d}_{jk}^A \int_0^t \sqrt{\eta}\left(f_{jk}^A(u, \bar{\mathbf{V}}^\eta(u)) - f_{jk}^A(s, \mathbf{v}(u))\right) du \\
&+ \sum_{j=1}^{m_A}\sum_{k=1}^{m_A} \mathbf{d}_{jk}^I \int_0^t \sqrt{\eta}\left(f_{jk}^I(s, \bar{\mathbf{V}}^\eta(s)) - f_{jk}^I(s, \mathbf{v}(t))\right) ds \\
&+ \sum_{i=1}^{m_S}\sum_{l \neq i}^{m_S} \mathbf{d}_{il}^S \int_0^t \sqrt{\eta}\left(f_{il}^S(s, \bar{\mathbf{V}}^\eta(s)) - f_{il}^S(s, \mathbf{v}(t))\right) ds \\
&+ \sum_{i=1}^{m_S} \mathbf{d}_i^D \int_0^t \sqrt{\eta}\left(f_i^D(s, \bar{\mathbf{V}}^\eta(s)) - f_i^D(s, \mathbf{v}(t))\right) ds - \int_0^t \partial\mathbf{F}(u, \mathbf{v}(u))\mathbf{D}^\eta(u) \\
&= \int_0^t \partial\mathbf{F}(u, \mathbf{v}(u))E^\eta(u) + \sqrt{\eta}\int_0^t \left(\mathbf{F}(u, \bar{\mathbf{V}}^\eta(u)) - \mathbf{F}(u, \mathbf{v}(u))\right) du \\
&- \int_0^t \partial\mathbf{F}(u, \mathbf{v}(u))\mathbf{D}^\eta(u)
\end{aligned}
$$

*By the mean value theorem, there exist vectors $\zeta^\eta(u)$ that is in between $\mathbf{v}(u)$ and $\bar{\mathbf{V}}^\eta(u)$) such that*

$$\mathbf{F}(u, \zeta^\eta(u)) - \mathbf{F}(u, \mathbf{v}(u)) = \frac{1}{\sqrt{\eta}} \partial \mathbf{F}(u, \zeta^\eta(u)) \mathbf{D}^\eta(u)$$

*This implies that*

$$E^\eta(t) = \int_0^t \left( \partial \mathbf{F}(u, \zeta^\eta(u)) - \partial \mathbf{F}(u, \mathbf{v}(u)) \right) \mathbf{D}^\eta(u) + \int_0^t \partial \mathbf{F}(u, \mathbf{v}(u)) E^\eta(u) du. \quad \text{(A.14)}$$

*We also know that*

$$\lim_{\eta \to \infty} \sup_{0 \le t \le T} \| \partial \mathbf{F}(u, \zeta^\eta(u)) - \partial \mathbf{F}(u, \mathbf{v}(u)) \| = 0 \quad a.s \quad \text{(A.15)}$$

*in lieu of the fluid limit convergence. Moreover, since $\mathbf{D}^\eta(u)$ is bounded in probability and Lemma A.5 is true, we have that the process*

$$\lim_{\eta \to \infty} \sup_{0 \le t \le T} \int_0^t \left( \partial \mathbf{F}(u, \zeta^\eta(u)) - \partial \mathbf{F}(u, \mathbf{v}(u)) \right) \mathbf{D}^\eta(u) = 0 \quad \text{in probability.}$$

*Applying Gronwall's lemma and using Lemma A.5, we finally obtain our diffusion limit result.*
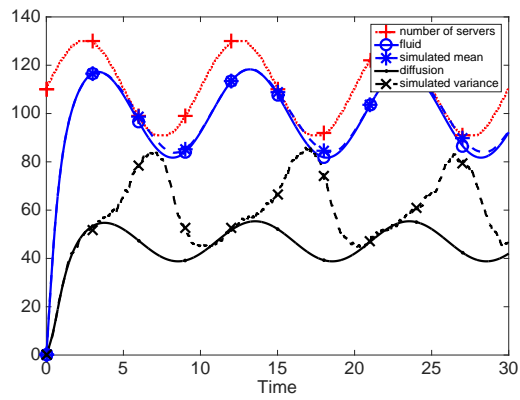
## A.3   Additional Numerical Examples

In this section, we provide some numerical examples where our proposed approach may not work well and discuss what causes this problem. We conduct 3 experiments sharing the same base inter-arrival time distribution, Weibull$(2.1271, 2.5)$, SCoV $= 0.1831$, and service time distribution, Lognormal$(-0.2027, 0.6368)$, SCoV $= 0.5$ with time-varying arrival rates and number of servers:

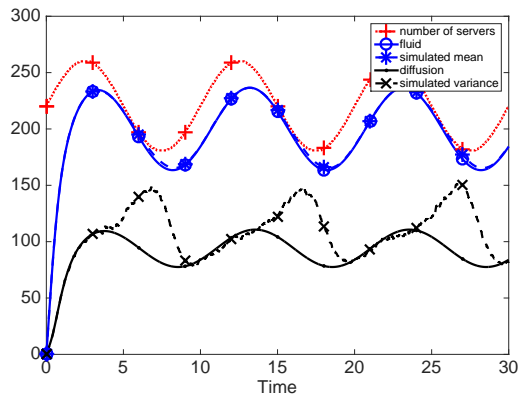Exp. A1: $\lceil 110 + 20 \sin(2\pi t/10) \rceil$ servers, $100 + 20 \sin(2\pi t/10)$ arrival rate

Exp. A2: $\lceil 220 + 40 \sin(2\pi t/10) \rceil$ servers, $200 + 40 \sin(2\pi t/10)$ arrival rate

Exp. A3: $\lceil 880 + 160 \sin(2\pi t/10) \rceil$ servers, $800 + 160 \sin(2\pi t/10)$ arrival rate
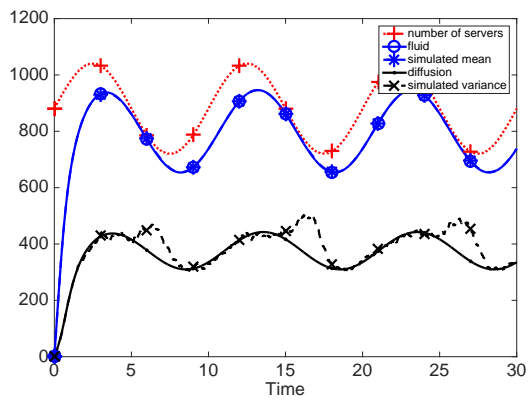
Figure A.1 (a) shows a large discrepancy between simulated variance and the variance from the diffusion limit on the time interval $(3, 7)$. It is because of the well-known lingering effect mentioned in Mandelbaum et al. [28] and addressed by Ko and Gautam [22]. Lingering occurs when the queue is critically loaded, i.e., the fluid limit stays close to the number of servers. We also observe that the discrepancy decreases as we increases the number of servers (the acceleration parameter). In order to reduce the discrepancy for the queues with a small number of servers, we might need to do some adjustments just as Ko and Gautam [22] did. However, the application of such an adjustment is not straightforward under multi-dimensional settings. We leave the issue for future research.

(a) Exp. A1



(b) Exp. A2



(c) Exp. A3

Figure A.1: Lingering effect when critically loaded

# References

[1] Muhammad Asad Arfeen, K. Pawlikowski, D. McNickle, and A. Willig. The role of the Weibull distribution in Internet traffic modeling. In *Proceedings of the 2013 25th International Teletraffic Congress (ITC)*, pages 1–8. IEEE, September 2013.

[2] Ludwig Arnold. *Stochastic Differential Equations: Theory and Applications*. Krieger Publishing Company, 1992.

[3] Sø ren Asmussen, Olle Nerman, and Marita Olsson. Fitting Phase-Type Distributions via the EM Algorithm. *Scandinavian Journal of Statistics*, 23(4):419–441, 1996.

[4] Andrew D Barbour. Networks of queues and the method of stages. *Advances in Applied Probability*, pages 584–591, 1976.

[5] Patrick Billingsley. *Convergence of Probability Measures*. A John Wiley & Sons, Inc., Publication, 1999.

[6] A. Bobbio, A. Horváth, and M. Telek. Matching Three Moments with Minimal Acyclic Phase Type Distributions. *Stochastic Models*, 21(2-3):303–326, 2005.

[7] R. F. Botta and C. M. Harris. Approximation with generalized hyperexponential distributions: Weak convergence results. *Queueing Systems*, 1(2):169–190, September 1986.

[8] Lawrence Brown, Noah Gans, Avishai Mandelbaum, Anat Sakov, Haipeng Shen, Sergey Zeltyn, and Linda Zhao. Statistical Analysis of a Telephone Call Center. *Journal of the American Statistical Association*, 100(469):36–50, March 2005.

[9] Stefan Creemers, Mieke Defraeye, and Inneke Van Nieuwenhuyse. G-RAND: A phase-type approximation for the nonstationary $G(t)/G(t)/s(t) + G(t)$ queue. *Performance Evaluation*, 80:102–123, August 2014.

[10] JG Dai, Shuangchi He, Tolga Tezcan, et al. Many-server diffusion limits for g/ph/n+ gi queues. *The Annals of Applied Probability*, 20(5):1854–1890, 2010.

[11] S. G. Eick, W. A. Massey, and W. Whitt. The Physics of the $M_t/G/\infty$ Queue. *Operations Research*, 41(4):731–742, July 1993.

[12] Stefan Engblom and Jamol Pender. Approximations for the moments of nonstationary and state dependent birth-death queues. 2014.

[13] Anja Feldmann and Ward Whitt. Fitting mixtures of exponentials to long-tail distributions to analyze network performance models. *Performance Evaluation*, 31(3-4): 245–279, 1998.

[14] Peter W. Glynn. On the Markov Property of the $GI/G/\infty$ Gaussian Limit. *Advances in Applied Probability*, 14(1):191–194, 1982.

[15] Shlomo Halfin and Ward Whitt. Heavy-Traffic Limits for Queues with Many Exponential Servers. *Operations Research*, 29(3):567–588, May 1981.

[16] R Hampshire and WA Massey. A tutorial on dynamic optimization and applications to queueing systems with time-varying rates. *Tutorials in Operations Research*, pages 208–247, 2010.

[17] Robert C Hampshire and William A Massey. Variational optimization for call center staffing. In *Proceedings of the 2005 conference on Diversity in computing*, pages 4–6. ACM, 2005.

[18] Robert C. Hampshire, Mor Harchol-Balter, and William A. Massey. Fluid and diffusion limits for transient sojourn times of processor sharing queues with time varying rates. *Queueing Systems*, 53(1-2):19–30, June 2006.

[19] Robert C Hampshire, Otis B Jennings, and William A Massey. A time-varying call center design via lagrangian mechanics. *Probability in the Engineering and Informational Sciences*, 23(02):231–259, 2009.

[20] Robert C Hampshire, William A Massey, and Qiong Wang. Dynamic pricing to control loss systems with quality of service targets. *Probability in the Engineering and Informational Sciences*, 23(02):357–383, 2009.

[21] Mary A. Johnson and Michael R. Taaffe. An investigation of phase-distribution moment-matching algorithms for use in queueing models. *Queueing Systems*, 8(1):129–147, December 1991.

[22] Young Myoung Ko and Natarajan Gautam. Critically Loaded Time-Varying Multi-server Queues: Computational Challenges and Approximations. *INFORMS Journal on Computing*, 25(2):285–301, May 2013.

[23] T Kurtz. Strong approximation theorems for density dependent Markov chains. *Stochastic Processes and their Applications*, 6(3):223–240, February 1978.

[24] Yunan Liu and Ward Whitt. The $G_t/GI/s_t + GI$ many-server fluid queue. *Queueing Systems*, 71(4):405–444, March 2012.

[25] Yunan Liu and Ward Whitt. Many-server heavy-traffic limit for queues with time-varying parameters. *The Annals of Applied Probability*, 24(1):378–421, February 2014.

[26] Yunan Liu and Ward Whitt. Algorithms for Time-Varying Networks of Many-Server Fluid Queues. *INFORMS Journal on Computing*, 26(1):59–73, February 2014.

[27] Avi Mandelbaum, William Massey, and Martin Reiman. Strong approximations for Markovian service networks. *Queueing Systems*, 30(1):149–201, 1998.

[28] Avi Mandelbaum, William A. Massey, Martin I. Reiman, and Alexander Stolyar. Queue Lengths and Waiting Times for Multiserver Queues with Abandonment and Retrials. *Telecommunication Systems*, 21(2-4):149–171, 2002.

[29] William A Massey and Jamol Pender. Poster: skewness variance approximation for dynamic rate multiserver queues with abandonment. *ACM SIGMETRICS Performance Evaluation Review*, 39(2):74–74, 2011.

[30] William A. Massey and Jamol Pender. Gaussian skewness approximation for dynamic rate multi-server queues with abandonment. *Queueing Systems*, 75(2-4):243–277, February 2013.

[31] Barry L. Nelson and Michael R. Taaffe. The $[Ph_t/Ph_t/\infty]^K$ Queueing System: Part II-The Multiclass Network. *INFORMS Journal on Computing*, 16(3):275–283, August 2004.

[32] Barry L. Nelson and Michael R. Taaffe. The $Ph_t/Ph_t/\infty$ Queueing System: Part I-The Single Node. *INFORMS Journal on Computing*, 16(3):266–274, August 2004.

[33] Marcel F. Neuts. *Matrix-Geometric Solutions In Stochastic Models: An Algorithmic Approach*. Dover Publication, Inc., 1981.

[34] Jerome Niyirora and Jamol Pender. Optimal staffing of clinical revenue centers in healthcare organizations. *Under Review*, 2015.

[35] Takayuki Osogami and Mor Harchol-Balter. Closed form solutions for mapping general distributions to quasi-minimal PH distributions. *Performance Evaluation*, 63(6):524–552, June 2006.

[36] Jihong Ou, Jingwen Li, and Süleyman Özekici. Approximating a Cumulative Distribution Function by Generalized Hyperexponential Distributions. *Probability in the Engineering and Informational Sciences*, 11(1):11–18, 1997.

[37] Guodong Pang and Ward Whitt. Heavy-traffic limits for many-server queues with service interruptions. *Queueing Systems*, 61(2-3):167–202, Mar 2009.

[38] Jamol Pender. Gram charlier expansion for time varying multiserver queues with abandonment. *SIAM Journal on Applied Mathematics*, 74(4):1238–1265, 2014.

[39] Jamol Pender. Nonstationary loss queues via cumulant moment approximations. *Probability in the Engineering and Informational Sciences*, 29(01):27–49, 2015.

[40] Jamol Pender. An analysis of nonstationary coupled queues. *Telecommunication Systems*, pages 1–16, 2015.

[41] Jamol Pender. The truncated normal distribution: Applications to queues with impatient customers. *Operations Research Letters*, 43(1):40–45, 2015.

[42] Anatolii A. Puhalskii and Martin I. Reiman. The multiclass GI/PH/N queue in the Halfin-Whitt regime. *Advances in Applied Probability*, 32(2):564–595, Jun 2000.

[43] Josh Reed. The g/gi/n queue in the halfin–whitt regime. *The Annals of Applied Probability*, 19(6):2211–2269, 2009.

[44] Yukie Sasaki, Hiroei Imai, Masahiro Tsunoyama, and Ikuo Ishii. Approximation of probability distribution functions by coxian distribution to evaluate multimedia systems. *Systems and Computers in Japan*, 35(2):16–24, 2004.

[45] Ward Whitt. On the Heavy-Traffic Limit Theorem for $GI/G/\infty$ Queues. *Advances in Applied Probability*, 14(1):171–190, 1982.

[46] Ward Whitt. *Stochastic Process Limits*. Springer, 1 edition, 2002.

[47] Ward Whitt. Fluid Models for Multiserver Queues with Abandonments. *Operations Research*, 54(1):37–54, January 2006.

[48] K. Yu, M.-L. Huang, and P. H. Brill. An Algorithm for Fitting Heavy-Tailed Distributions via Generalized Hyperexponentials. *INFORMS Journal on Computing*, 24(1): 42–52, March 2011.