

# STEIN'S METHOD FOR STEADY-STATE DIFFUSION APPROXIMATIONS OF $M/Ph/N + M$ SYSTEMS\*

BY ANTON BRAVERMAN AND J. G. DAI

*Cornell University*

We consider  $M/Ph/n + M$  queueing systems in steady state. We prove that the Wasserstein distance between the stationary distribution of the normalized system size process and that of a piecewise Ornstein-Uhlenbeck (OU) process is bounded by  $C/\sqrt{\lambda}$ , where the constant  $C$  is independent of the arrival rate  $\lambda$  and the number of servers  $n$  as long as they are in the Halfin-Whitt parameter regime. For each integer  $m > 0$ , we also establish a similar bound for the difference of the  $m$ th steady-state moments. For the proofs, we develop a modular framework that is based on Stein's method. The framework has three components: Poisson equation, generator coupling, and state space collapse. The framework, with further refinement, is likely applicable to steady-state diffusion approximations for other stochastic systems.

**1. Introduction.** This paper focuses on  $M/Ph/n + M$  systems, which serve as building blocks to model large-scale service systems such as customer contact centers [1, 23] and hospital operations [2, 48]. In such a system, there are  $n$  identical servers, the arrival process is Poisson (the symbol  $M$ ) with rate  $\lambda$ , the service times are i.i.d. having a phase-type distribution (the symbol  $Ph$ ) with mean  $1/\mu$ , the patience times of customers are i.i.d. having an exponential distribution (the symbol  $+M$ ) with mean  $1/\alpha < \infty$ . When the waiting time of a customer in queue exceeds her patience time, the customer abandons the system without service; once the service of a customer is started, the customer does not abandon.

Let  $X_i(t)$  be the number of customers in phase  $i$  at time  $t$  for  $i = 1, \dots, d$ , where  $d$  is the number of phases in the service time distribution. Let  $X(t)$  be the corresponding vector. Then the system size process  $X = \{X(t), t \geq 0\}$  has a unique stationary distribution for any arrival rate  $\lambda$  and any server number  $n$  due to customer abandonment; although  $X$  is not a Markov chain, it is a function of a Markov chain with a unique stationary distribution, see Section 4 for details. In this paper, we prove, in Theorem 1, that

$$(1.1) \quad \sup_{h \in \mathcal{H}} \left| \mathbb{E}[h(\tilde{X}^{(\lambda)}(\infty))] - \mathbb{E}[h(Y(\infty))] \right| \leq \frac{C}{\sqrt{\lambda}} \quad \text{for any } \lambda > 0 \text{ and } n \geq 1$$

---

\*This research is supported in part by NSF Grants CMMI-1030589, CNS-1248117 and CMMI-1335724.  
*MSC 2010 subject classifications:* Primary 60K25; secondary 90B20, 60F99, 60J60.

*Keywords and phrases:* Stein's method, diffusion approximation, steady-state, many servers, state space collapse, convergence rate.

satisfying

$$(1.2) \quad n\mu = \lambda + \beta\sqrt{\lambda},$$

where  $\beta \in \mathbb{R}$  is some constant and  $\mathcal{H}$  is some class of functions  $h : \mathbb{R}^d \rightarrow \mathbb{R}$ . In (1.1),  $\tilde{X}^{(\lambda)}(\infty)$  is a random vector having the stationary distribution of a properly scaled version of  $X = X^{(\lambda)}$  that depends on the arrival rate  $\lambda$ , number of servers  $n$ , the service time distribution, and the abandonment rate  $\alpha$ , and  $Y(\infty)$  is a random vector having the stationary distribution of a piecewise Ornstein-Uhlenbeck (OU) process  $Y = \{Y(t), t \geq 0\}$ . The stationary distribution of  $X^{(\lambda)}$  exists even when  $\beta$  is negative because  $\alpha$  is assumed to be positive. The constant  $C$  depends on the service time distribution, abandonment rate  $\alpha$ , the constant  $\beta$  in (1.2), and the choice of  $\mathcal{H}$ , but  $C$  is independent of the arrival rate  $\lambda$  and the number of servers  $n$ . Two different classes  $\mathcal{H}$  will be used in our Theorem 1. First, we take  $\mathcal{H}$  to be the class of polynomials up to a certain order. In this case, (1.1) provides rates of convergence for steady-state moments. Second,  $\mathcal{H}$  is taken to be  $\mathcal{W}^{(d)}$ , the class of all 1-Lipschitz functions

$$(1.3) \quad \mathcal{W}^{(d)} = \{h : \mathbb{R}^d \rightarrow \mathbb{R} : |h(x) - h(y)| \leq |x - y|\}.$$

In this case, (1.1) provides rates of convergence for stationary distributions under the Wasserstein metric [47]; convergence under Wasserstein metric implies the convergence in distribution [24].

In [14], an algorithm was developed to compute the stationary distribution of the diffusion process  $Y$ . The distribution of  $Y(\infty)$  is then used to approximate the stationary distribution of  $X^{(\lambda)}$ . The approximation is remarkably accurate; see, for example, Figure 1 there. It was demonstrated that computational efficiency, in terms of both time and memory, can be achieved by diffusion approximations. For example, in an  $M/H_2/500 + M$  system studied in [14], where the system has 500 servers and a hyper-exponential service time distribution, it took around 1 hour and peak memory usage of 5 GB to compute the stationary distribution of  $X^{(\lambda)}$  using an algorithm that fully explores the special structure of a three-dimensional Markov chain. On the same computer, to compute the stationary distribution of the corresponding two-dimensional diffusion process it took less than 1 minute and peak memory usage was less than 200 MB. The computational saving by the diffusion model is achieved partly through *state space collapse* (SSC), a phenomenon that causes dimension reduction in state space. Theorem 1 quantifies the steady-state diffusion approximations developed in [14].

In [29], the authors prove a version of (1.1) for the  $M/M/n + M$  system, a special case of the  $M/Ph/n + M$  system where the service time distribution is exponential. They do not impose assumption (1.2) on the relationship between the arrival rate  $\lambda$  and number of servers  $n$ , resulting in a *universal* approximation that is accurate in any parameter regime, from underloaded, to critically loaded, and to overloaded. To our knowledge, this is the first paper to study convergence rates of steady state diffusion approximations. Their method

relies on analyzing excursions of a one-dimensional Markov chain and the corresponding diffusion process. It is unclear how to generalize their method to the multi-dimensional setting.

To prove Theorem 1, we develop a framework that is based on Stein's method [49, 50]. The framework is modular and relies on three components: a Poisson equation, generator coupling, and SSC. The framework itself is an important part of our contribution, in addition to Theorem 1. We expect the framework will be refined and used to prove rates of convergence of steady-state diffusion approximations for many other stochastic systems. This framework is closely related to a recent paper [27] by Gurvich. We will discuss his work after giving an overview of the framework.

We consider two sequences of stochastic processes  $\{X^{(\ell)}\}_{\ell=1}^{\infty}$  and  $\{Y^{(\ell)}\}_{\ell=1}^{\infty}$  indexed by  $\ell$ , where  $X^{(\ell)} = \{X^{(\ell)}(t) \in \mathbb{R}^d, t \geq 0\}$  is a continuous-time Markov chain (CTMC) and  $Y^{(\ell)} = \{Y^{(\ell)}(t) \in \mathbb{R}^d, t \geq 0\}$  is a diffusion process. Suppose  $X^{(\ell)}(\infty)$  and  $Y^{(\ell)}(\infty)$  are two random vectors having the stationary distributions of  $X^{(\ell)}$  and  $Y^{(\ell)}$ , respectively. Let  $G_{X^{(\ell)}}$  and  $G_{Y^{(\ell)}}$  be the generators of  $X^{(\ell)}$  and  $Y^{(\ell)}$ , respectively; for a diffusion process,  $G_{Y^{(\ell)}}$  is the second order elliptic operator as in (5.3). For a function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  in a "nice" (but large enough) class, we wish to bound

$$\left| \mathbb{E}h(X^{(\ell)}(\infty)) - \mathbb{E}h(Y^{(\ell)}(\infty)) \right|.$$

**Component 1.** The first step is to set up the Poisson equation

$$(1.4) \quad G_{Y^{(\ell)}} f_h(x) = h(x) - \mathbb{E}h(Y^{(\ell)}(\infty))$$

and obtain various estimates of a solution  $f_h$  to the Poisson equation. Once we have  $f_h$ , one can take the expectation of both sides above to see that

$$(1.5) \quad \mathbb{E}h(X^{(\ell)}(\infty)) - \mathbb{E}h(Y^{(\ell)}(\infty)) = \mathbb{E}G_{Y^{(\ell)}} f_h(X^{(\ell)}(\infty)).$$

The Poisson equation (1.4) is a partial differential equation (PDE). Even when  $Y^{(\ell)}(\infty) = Y(\infty)$  (i.e. independent of  $\ell$ ), one of the biggest challenges is obtaining bounds on the partial derivatives of  $f_h(x)$  (usually up to third order). We refer to these as gradient bounds. In the one-dimensional case, (1.4) is an ordinary differential equation (ODE) that usually has a closed form expression that one can analyze directly, see for instance [12, Lemma 13.1]. However, when  $d > 1$  obtaining these gradient bounds becomes significantly harder. By exploiting probabilistic solutions to the Poisson equation, gradient bounds were established for cases when  $Y(\infty)$  is a multivariate normal [4], multivariate Poisson [5] and multivariate Gamma [41].

**Component 2.** The next step is to produce the generator coupling. For that, we use the basic adjoint relationship (BAR) for the stationary distribution of  $X^{(\ell)}(\infty)$ . One can check that a random vector  $X^{(\ell)}(\infty) \in \mathbb{R}^d$  has the stationary distribution of the CTMC  $X^{(\ell)}$  if and only if

$$(1.6) \quad \mathbb{E}G_{X^{(\ell)}} f(X^{(\ell)}(\infty)) = 0$$

for all functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  that have compact support. For a given  $h$ , the corresponding Poisson equation solution  $f_h$  does not have compact support. An important part of this step is to prove that (1.6) continues to hold for  $f_h$ . Thus, it follows from (1.5) and (1.6) that

$$(1.7) \quad \mathbb{E}h(X^{(\ell)}(\infty)) - \mathbb{E}h(Y^{(\ell)}(\infty)) = \mathbb{E}[G_{Y^{(\ell)}}f_h(X^{(\ell)}(\infty)) - G_{X^{(\ell)}}f_h(X^{(\ell)}(\infty))].$$

Note that two random variables in the left side of (1.7) are typically defined on two different probability spaces, whereas two random variables in the right side of (1.7) are all defined in terms of  $X^{(\ell)}(\infty)$ , thus producing a coupling on a common probability space.

To bound the right side of (1.7), we study

$$(1.8) \quad G_{X^{(\ell)}}f_h(x) - G_{Y^{(\ell)}}f_h(x)$$

for each  $x$  in the state space of  $X^{(\ell)}$ . By performing Taylor expansion on  $G_{X^{(\ell)}}f_h(x)$ , we find that the difference involves the product of partial derivatives of  $f_h$  and a term bounded by a polynomial of  $x$ . Therefore, in addition to gradient bounds on  $f_h$ , in a lot of cases we need bounds on various moments of  $|X^{(\ell)}(\infty)|$  which we refer to as moment bounds. The main challenge is that both gradient and moment bounds must be *uniform* in  $\ell$ .

**Component 3.** In the last step, SSC comes into play when  $X^{(\ell)}$  itself is not a CTMC, but a projection of some higher dimensional CTMC  $U^{(\ell)} = \{U^{(\ell)}(t) \in \mathcal{U}, t \geq 0\}$ , where the dimension of the state space  $\mathcal{U}$  is strictly greater than  $d$ . This is the case, for example, in the  $M/Ph/n + M$  system. It is this difference in dimensions that is responsible for most of the computational speedup in diffusion approximations; most complex stochastic processing systems exhibit some form of SSC [6, 9, 16, 18, 20, 32, 33, 45, 52, 54]. Let  $G_U$  be the generator of  $U^{(\ell)}$  and  $U^{(\ell)}(\infty)$  have its stationary distribution. Now, BAR (1.6) becomes  $G_{U^{(\ell)}}F(U^{(\ell)}(\infty)) = 0$  for each ‘nice’  $F : \mathcal{U} \rightarrow \mathbb{R}$ . Furthermore, (1.7) becomes

$$(1.9) \quad \mathbb{E}h(X^{(\ell)}(\infty)) - \mathbb{E}h(Y^{(\ell)}(\infty)) = \mathbb{E}[G_{Y^{(\ell)}}f_h(X^{(\ell)}(\infty)) - G_{U^{(\ell)}}F_h(U^{(\ell)}(\infty))],$$

where  $F_h : \mathcal{U} \rightarrow \mathbb{R}$  is the lifting of  $f_h : \mathbb{R}^d \rightarrow \mathbb{R}$  defined by letting  $x \in \mathbb{R}^d$  be the projection of  $u \in \mathcal{U}$  and then setting

$$(1.10) \quad F_h(u) = f_h(x).$$

As before, we can perform Taylor expansion on  $G_{U^{(\ell)}}F_h(u)$  to simplify the difference  $G_{U^{(\ell)}}F_h(u) - G_{Y^{(\ell)}}f_h(x)$ . To use this difference to bound the right side of (1.9), we need a steady-state SSC result for  $U^{(\ell)}(\infty)$ , which tells us how to approximate  $U^{(\ell)}(\infty)$  from  $X^{(\ell)}(\infty)$  and guarantees that this approximation error is small. To obtain our SSC result, we need to rely heavily on the structure of the  $M/Ph/n + M$  system.

In [27], Gurvich develops methodologies to prove statements similar to (1.1) for various queueing systems. In particular, Gurvich develops important elements of the first two components of our framework in the special case when  $\dim(\mathcal{U}) = d$ . Along the way,

he independently rediscovers many of the ideas central to Stein's method in the setting of steady-state diffusion approximations. He relies on the existence of uniform Lyapunov functions for the diffusion processes. Putting the Lyapunov functions together with the probabilistic solution for (1.4) and a-priori Schauder estimates for elliptic PDEs (see [25]), he is able to obtain uniform gradient bounds for a large class of Poisson equations. Furthermore, he also obtains the necessary uniform moment bounds using these Lyapunov functions by showing that uniform moment bounds for the diffusion process imply the same moments are uniformly bounded for the CTMC. However, his result on uniform moment bounds no longer holds when  $\dim(\mathcal{U}) > d$  due to the need for SSC, which poses an additional technical challenge. We overcome this challenge for the  $M/Ph/n + M$  system in Lemma 7, in which moment bounds are established *recursively*.

The work in [27] is conceptually close to this paper. In that paper, Gurvich packages all the components required to prove his results into several conditions, with the main condition being the existence of uniform Lyapunov functions for the diffusion processes. In contrast, a key contribution of our framework is its *modular* nature. The immediate benefit we gain is the ability to apply this framework to cases when SSC occurs ( $\dim(\mathcal{U}) > d$ ). Moreover, although we also rely on Lyapunov functions to establish both moment and gradient bounds in our particular setting, our framework clearly illustrates that Lyapunov functions are merely tools one can use to establish these moment and gradient bounds; the bounds themselves are the actual drivers of our main results.

We have already mentioned that Lemma 13.1 of [12] presents a systematic way to establish gradient bounds in the one-dimensional setting ( $d = 1$ ), and [4, 5, 41] establish gradient bounds in the multi-dimensional setting ( $d > 1$ ) for a few special cases of  $Y^{(\ell)}(\infty)$ . However, establishing multi-dimensional gradient bounds remains a very difficult problem that usually requires using structural properties of the distribution of  $Y^{(\ell)}(\infty)$ . Gurvich's use of a-priori Schauder estimates [25] together with Lyapunov functions represents the first systematic approach to establishing multi-dimensional gradient bounds.

With regards to using Lyapunov functions to establish moment bounds, certain systems may not require moment bounds at all. For example, approximating the stationary distribution of the simple birth-death process corresponding to a single-server queue does not require the use of moment bounds (although we do not consider the  $M/M/1$  queue in this paper, Stein's method is easily applicable to it). Thus, the modularity of our framework presents the components one needs to justify approximations for various systems, and promotes the view that Lyapunov functions are merely one of many tools to tackle the difficulties in these components.

It is useful to compare the challenge level of each component in our framework. The generator coupling is the least challenging component, because the class of functions for which (1.6) holds is usually rich enough. The remaining major difficulties are moment bounds, gradient bounds and SSC. Moment bounds and SSC are a property of the CTMC sequence  $\{X^{(\ell)}\}_{\ell=1}^{\infty}$ , and the difficulty in establishing them will depend heavily on the CTMCs. On the other hand, gradient bounds are tied to the diffusion processes  $\{Y^{(\ell)}\}_{\ell=1}^{\infty}$ , and are typ-

ically only difficult to establish when the diffusion processes are multi-dimensional. One important class of multi-dimensional diffusion processes for which we do not have gradient bounds are semi-martingale reflected Brownian motions (SRBMs) [34]. An SRBM can approximate networks of single-server queues, such as generalized Jackson networks. The Schauder gradient bounds of [27] are not immediately applicable to SRBMs, because the corresponding Poisson equation is defined on the non-negative orthant, and has oblique reflection boundary conditions.

Stein’s method is a powerful method that has been widely used in probability, statistics, and their wide range of applications such as bioinformatics; see, for example, the survey papers [11, 47], the recent book [12] and the references within. The connection between Stein’s method and diffusion processes was first made by Barbour in [4, 5]. In the context of Stein’s method, generator coupling is a realization of an abstract concept that first appeared in the famous commutative diagram in (28) of [50]; a more refined explanation of which is provided in (4) of [11]. In particular, using Chatterjee’s notation in [11], our  $\mathbb{E}G_{X^{(\ell)}}f_h(X^{(\ell)}(\infty))$  in (1.7) is his  $\mathbb{E}T\alpha f(W)$ .

Diffusion approximations are usually “justified” by heavy traffic limit theorems. It is proved in [15] that for our  $M/Ph/n + M$  systems,

$$(1.11) \quad \tilde{X}^{(\lambda)} = \{\tilde{X}^{(\lambda)}(t), t \geq 0\} \implies Y = \{Y(t), t \geq 0\}$$

as  $\lambda$  goes to infinity while satisfying (1.2) (we use the arrival rate  $\lambda$  to index these systems instead of the abstract  $\ell$  as before). Proving these limit theorems has been an active area of research in the last 50 years; see, for example, [7, 8, 31, 36, 37, 46] for single-class queueing networks, [9, 43, 54] for multiclass queueing networks, [38, 55] for bandwidth sharing networks, [15, 30, 44] for many-server queues. The convergence used in these limit theorems is the convergence in distribution on the path space  $\mathbb{D}([0, \infty), \mathbb{R}^d)$ , endowed with Skorohod  $J_1$ -topology [19, 53]. The  $J_1$ -topology on  $\mathbb{D}([0, \infty), \mathbb{R}^d)$  essentially means convergence in  $\mathbb{D}([0, T], \mathbb{R}^d)$  for each  $T > 0$ . In particular, it says nothing about the convergence at “ $\infty$ ”. Therefore, these limit theorems do not justify the steady-state convergence.

In [13], the authors prove the convergence of distribution  $\tilde{X}^{(\lambda)}(\infty)$  to that of  $Y(\infty)$  by proving an interchange of limits. The proof technique follows that of the seminal paper [22], where the authors prove an interchange of limits for generalized Jackson networks of single-server queues. The results in [22] were improved and extended by various authors for networks of single-servers [10, 39, 56], for bandwidth sharing networks [55], and for many-server systems [21, 28, 51]. These “interchange limits theorems” are qualitative and thus do not provide rates of convergence as in (1.1).

1.1. *Notation.* All random variables and stochastic processes are defined on a common probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  unless otherwise specified. For a stochastic process  $X = \{X(t), t \geq 0\}$  that has a unique stationary distribution we let  $X(\infty)$  be the random element having the stationary distribution of  $X$ . For a sequence of random variables  $\{X^n\}_{n=1}^\infty$ , we write  $X^n \Rightarrow X$  to denote convergence in distribution (also known as weak convergence) of

$X^n$  to some random variable  $X$ . If  $a > b$ , we adopt the convention that  $\sum_{i=a}^b (\cdot) = 0$ . For an integer  $d \geq 1$ ,  $\mathbb{R}^d$  denotes the  $d$ -dimensional Euclidean space and  $\mathbb{Z}_+^d$  denotes the space of  $d$ -dimensional vectors whose elements are non-negative integers. For  $a, b \in \mathbb{R}$ , we define  $a \vee b = \max\{a, b\}$  and  $a \wedge b = \min\{a, b\}$ . For  $x \in \mathbb{R}$ , we define  $x^+ = x \vee 0$  and  $x^- = (-x) \vee 0$ . For  $x \in \mathbb{R}^d$ , we use  $x_i$  to denote its  $i$ th entry and  $|x|$  to denote its Euclidean norm. For  $x, y \in \mathbb{R}^d$ , we write  $x \leq y$  when  $x_i \leq y_i$  for all  $i$  and when  $x \leq y$  we define the vector interval  $[x, y] = \{z : x \leq z \leq y\}$ . All vectors are assumed to be column vectors. We let  $x^T$  and  $A^T$  denote the transpose of a vector  $x$  and matrix  $A$ , respectively. For a matrix  $A$ , we use  $A_{ij}$  to denote the entry in the  $i$ th row and  $j$ th column. We reserve  $I$  for the identity matrix,  $e$  for the vector of all ones and  $e^{(i)}$  for the vector that has a one in the  $i$ th element and zeroes elsewhere; the dimensions of these vectors will be clear from the context.

**1.2. Outline for Rest of Paper.** The rest of the paper is structured as follows. Section 2 formally defines the  $M/Ph/n + M$  system as well as the diffusion process whose steady-state distribution will approximate the system. Section 3 states our main results. Section 4 describes the CTMC representation of the  $M/Ph/n + M$  system. Section 5 introduces the first two components of our framework; the Poisson equation and generator coupling. Section 6 describes the SSC result, illustrating the third component of our framework. It is here that the reader may see the reason behind our slower rate of convergence. This framework is then used in Section 7 to prove our main results. Appendix A contains the proofs for most of the lemmas.

**2. Models.** In this section, we give additional description of the  $M/Ph/n + M$  system and the corresponding diffusion model.

**2.1. The  $M/Ph/n + M$  System.** The basic description of the  $M/Ph/n + M$  queueing system was given in the first paragraph of the introduction. Here, we describe the dynamics of the system. Upon arrival to the system with idle servers, a customer begins service immediately. Otherwise, if all servers are busy, the customer enters an infinite capacity queue to wait for service. When a server completes serving a customer, the server becomes idle if the queue is empty, or takes a customer from the queue under the first-in-first-out (FIFO) service policy if it is nonempty. Recall that the  $Ph$  indicates that customer service times are i.i.d. following a phase-type distribution. We shall provide a definition of a phase-type distribution shortly below. The phase-type distribution can approximate any positive-valued distribution [3, Theorem III.4.2].

Recall that  $\lambda$  denotes the arrival rate of the system. We use  $1/\alpha$  to denote the mean patience time. In our study, we take the service time distribution and  $\alpha$  fixed, but allow the arrival rate  $\lambda$  and the number of servers  $n$  to grow without bound. Throughout this paper, we assume that  $n$  follows the square-root-safety staffing rule in (1.2). In the pioneering paper of [30], the authors studied these systems as  $\lambda \rightarrow \infty$  and  $n$  grows to infinity following

(1.2). This parameter regime is now known as the Halfin-Whitt regime. In this regime, the system has high server utilization and at the same time has small customer waiting time and abandonment fraction. Therefore, this regime is also known as the quality- and efficiency-driven (QED) regime, a term coined by [23].

*Phase-type Service Time Distribution.* A phase-type distribution is assumed to have  $d \geq 1$  phases. Each phase-type distribution is determined by the tuple  $(p, \nu, P)$ , where  $p \in \mathbb{R}^d$  is a vector of non-negative entries whose sum is equal to one,  $\nu \in \mathbb{R}^d$  is a vector of positive entries and  $P$  is a  $d \times d$  sub-stochastic matrix. We assume that  $P$  is transient, i.e.

$$(2.1) \quad (I - P)^{-1} \quad \text{exists,}$$

and without loss of generality, we also assume that the diagonal entries of  $P$  are zero ( $P_{ii} = 0$ ).

A random variable is said to have a phase-type distribution with parameters  $(p, \nu, P)$  if it is equal to the absorption time of the following CTMC. The state space of the CTMC is  $\{1, \dots, d + 1\}$ , with  $d + 1$  being the absorbing state. The CTMC starts off in one of the states in  $\{1, \dots, d\}$  according to distribution  $p$ . For  $i = 1, \dots, d$ , the time spent in state  $i$  is exponentially distributed with mean  $1/\nu_i$ . Upon leaving state  $i$ , the CTMC transitions to state  $j = 1, \dots, d$  with probability  $P_{ij}$ , or gets absorbed into state  $d + 1$  with probability  $1 - \sum_{j=1}^d P_{ij}$ .

The CTMC above is a useful way to describe the service times in the  $M/Ph/n + M$  system. Upon arrival to the system, a customer is assigned her first service phase according to distribution  $p$ . If the customer is forced to wait in queue because all servers are busy, she is still assigned a first service phase, but this phase of service will not start until a server takes on this customer for service. Once a customer with initial phase  $i$  enters service, her service time is the time until absorption to state  $d + 1$  by the CTMC. We assume without loss of generality that for each service phase  $i$ , either

$$(2.2) \quad p_i > 0 \text{ or } P_{ji} > 0 \text{ for some } j.$$

This simply means that there are no redundant phases.

We now define some useful quantities for future use. Define

$$(2.3) \quad R = (I - P^T)\text{diag}(\nu) \quad \text{and} \quad \gamma = \mu R^{-1}p,$$

where the matrix  $\text{diag}(\nu)$  is the  $d \times d$  diagonal matrix with diagonal entries given by the components of  $\nu$ . One may verify that  $\sum_{i=1}^d \gamma_i = 1$ . One can interpret  $\gamma_i$  to be the fraction of phase  $i$  service load on the  $n$  servers.

For concreteness, we provide two examples of phase-type distributions when  $d = 2$ . The first example is the two-phase hyper-exponential distribution, denoted by  $H_2$ . The corresponding tuple of parameters is  $(p, \nu, P)$ , where

$$p = (p_1, p_2)^T, \quad \nu = (\nu_1, \nu_2)^T, \quad \text{and} \quad P = 0.$$

Therefore, with probability  $p_i$ , the service time follows an exponential distribution with mean  $1/\nu_i$ .

The second example is the Erlang-2 distribution, denoted by  $E_2$ . The corresponding tuple of parameters is  $(p, \nu, P)$ , where

$$p = (1, 0)^T, \quad \nu = (\theta, \theta)^T, \quad \text{and} \quad P = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}.$$

An  $E_2$  random variable is a sum of two i.i.d. exponential random variables, each having mean  $1/\theta$ .

*2.2. System Size Process and Diffusion Model.* Before we state the main results, we introduce the process we wish to approximate, as well as the approximating diffusion process – the piecewise OU process. Recall that  $X = \{X(t) \in \mathbb{R}^d, t \geq 0\}$  is the system size process, where

$$X(t) = (X_1(t), \dots, X_d(t))^T,$$

and  $X_i(t)$  is the number of customers of phase  $i$  in the system (queue + service) at time  $t$ . We emphasize that  $X$  is not a CTMC, but it is a deterministic function of a higher-dimensional CTMC, which will be described in Section 4.

The process  $X$  depends on  $\lambda, n, \alpha, p, P$ , and  $\nu$ . However, in this paper we keep  $\alpha, p, P$ , and  $\nu$  fixed, and allow  $\lambda$  and  $n$  to vary according to (1.2). For the remainder of the paper we write  $X^{(\lambda)}$  to emphasize the dependence of  $X$  on  $\lambda$ ; the dependence of  $X^{(\lambda)}$  on  $n$  is implicit through (1.2).

Recall the definition of  $\gamma$  in (2.3) and define the scaled random variable

$$(2.4) \quad \tilde{X}^{(\lambda)}(\infty) = \delta(X^{(\lambda)}(\infty) - \gamma n),$$

where, for convenience, we let

$$(2.5) \quad \delta = 1/\sqrt{\lambda}.$$

To approximate  $\tilde{X}^{(\lambda)}(\infty)$ , we introduce the piecewise OU process  $Y = \{Y(t), t \geq 0\}$ . This is a  $d$ -dimensional diffusion process satisfying

$$(2.6) \quad Y(t) = Y(0) - p\beta t - R \int_0^t (Y(s) - p(e^T Y(s))^+) ds - \alpha p \int_0^t (e^T Y(s))^+ ds + \sqrt{\Sigma} B(t).$$

Above,  $B(t)$  is the  $d$ -dimensional standard Brownian motion and  $\sqrt{\Sigma}$  is any  $d \times d$  matrix satisfying

$$(2.7) \quad \sqrt{\Sigma} \sqrt{\Sigma}^T = \Sigma = \text{diag}(p) + \sum_{k=1}^d \gamma_k \nu_k H^k + (I - P^T) \text{diag}(\nu) \text{diag}(\gamma) (I - P),$$

where the matrix  $H^k$  is defined as

$$H_{ii}^k = P_{ki}(1 - P_{ki}), \quad H_{ij}^k = -P_{ki}P_{kj} \quad \text{for } j \neq i.$$

Comparing the form of  $\Sigma$  above to (2.24) of [14] confirms that it is positive definite. Thus  $\sqrt{\Sigma}$  exists. Observe that  $Y$  depends only on  $\beta, \alpha, p, P$ , and  $\nu$ , all of which are held constant throughout this paper.

The diffusion process in (2.6) has been studied by [17]. They prove that  $Y$  is positive recurrent by finding an appropriate Lyapunov function. In particular, this means that  $Y$  admits a stationary distribution.

### 3. Main Results.

We now state our main results.

**THEOREM 1.** *For every integer  $m > 0$ , there exists a constant  $C_m = C_m(\beta, \alpha, p, \nu, P) > 0$  such that for all locally Lipschitz functions  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfying*

$$|h(x)| \leq |x|^{2m} \quad \text{for } x \in \mathbb{R}^d,$$

we have

$$\left| \mathbb{E}h(\tilde{X}^{(\lambda)}(\infty)) - \mathbb{E}h(Y(\infty)) \right| \leq \frac{C_m}{\sqrt{\lambda}} \quad \text{for all } \lambda > 0$$

satisfying (1.2), which we recall below as

$$n\mu = \lambda + \beta\sqrt{\lambda}.$$

Theorem 1 will be proved in Section 7. As a consequence of the theorem, we immediately have the following corollary.

**COROLLARY 1.** *There exists a constant  $C_1 = C_1(\beta, \alpha, p, \nu, P) > 0$  such that*

$$\sup_{h \in \mathcal{W}^{(d)}} \left| \mathbb{E}h(\tilde{X}^{(\lambda)}(\infty)) - \mathbb{E}h(Y(\infty)) \right| \leq \frac{C_1}{\sqrt{\lambda}} \quad \text{for all } \lambda > 0$$

satisfying (1.2), where  $\mathcal{W}^{(d)}$  is defined in (1.3). In particular,

$$\tilde{X}^{(\lambda)}(\infty) \Rightarrow Y(\infty) \quad \text{as } \lambda \rightarrow \infty.$$

**PROOF.** Suppose  $h \in \mathcal{W}^{(d)}$ . Without loss of generality, we may assume that  $h(0) = 0$ , otherwise we may simply consider  $h(x) - h(0)$ . By definition of  $\mathcal{W}^{(d)}$ ,

$$|h(x)| \leq |x| \quad \text{for } x \in \mathbb{R}^d$$

and the result follows from Theorem 1 with  $m = 1$ . □

**REMARK 1.** For any fixed  $\beta \in \mathbb{R}$ , there are only finitely many combinations of  $\lambda \in (0, 4)$  and integer  $n \geq 1$  satisfying (1.2). Therefore, it suffices to prove Theorem 1 by restricting  $\lambda \geq 4$ , a convenience for technical purposes.

**4. Markov Representation.** The  $M/Ph/n + M$  system can be represented as a CTMC

$$U^{(\lambda)} = \{U^{(\lambda)}(t), t \geq 0\}$$

taking values in  $\mathcal{U}$ , the set of finite sequences  $\{u_1, \dots, u_k\}$ . The sequence  $u = \{u_1, \dots, u_k\}$  encodes the service phase of each customer and their order of arrival to the system. For example, the sequence  $\{5, 1, 4\}$  corresponds to 3 customers in the system, with the service phases of the first, second and third customers (in the order of their arrival to the system) being 5, 1 and 4, respectively. We use  $|u|$  to denote the length of the sequence  $u$ . The irreducibility of the CTMC  $U^{(\lambda)}$  is guaranteed by (2.1) and (2.2).

We remark here that  $U^{(\lambda)}$  is not the simplest Markovian representation of the  $M/Ph/n + M$  system. Another way to represent this system would be to consider a  $d + 1$  dimensional CTMC that keeps track of the total number of customers in the system, as well as the total number of customers in each phase that are currently in service; this  $d + 1$  dimensional CTMC is used in [15]. In this paper we use the infinite dimensional CTMC  $U^{(\lambda)}$  because the system size process  $X^{(\lambda)}$  cannot be recovered sample path wise from the  $d + 1$  dimensional CTMC, it can only be recovered from  $U^{(\lambda)}$ . Also, the CTMC  $U^{(\lambda)}$  will play an important role in our SSC argument in Section 6.

In addition to the system size process  $X^{(\lambda)}$ , we define the queue size process  $Q^{(\lambda)} = \{Q^{(\lambda)}(t) \in \mathbb{Z}_+^d, t \geq 0\}$ , where

$$Q^{(\lambda)}(t) = (Q_1^{(\lambda)}(t), \dots, Q_d^{(\lambda)}(t))^T,$$

and  $Q_i^{(\lambda)}(t)$  is the number of customers of phase  $i$  in the queue at time  $t$ . Then  $X_i^{(\lambda)}(t) - Q_i^{(\lambda)}(t) \geq 0$  is the number phase  $i$  customers in service at time  $t$ .

To recover  $X^{(\lambda)}(t)$  and  $Q^{(\lambda)}(t)$  from  $U^{(\lambda)}(t)$ , we define the projection functions  $\Pi_X : \mathcal{U} \rightarrow \mathbb{R}^d$  and  $\Pi_Q : \mathcal{U} \rightarrow \mathbb{R}^d$ . For each  $u \in \mathcal{U}$  and each phase  $i \in \{1, \dots, d\}$ ,

$$(\Pi_X(u))_i = \sum_{k=1}^{|u|} 1_{\{u_k=i\}} \quad \text{and} \quad (\Pi_Q(u))_i = \sum_{k=n+1}^{|u|} 1_{\{u_k=i\}}.$$

It is clear that on each sample path

$$(4.1) \quad X^{(\lambda)}(t) = \Pi_X(U^{(\lambda)}(t)) \quad \text{and} \quad Q^{(\lambda)}(t) = \Pi_Q(U^{(\lambda)}(t)) \quad \text{for } t \geq 0.$$

Because there is customer abandonment the Markov chain  $U^{(\lambda)}$  can be proved to be positive recurrent with a unique stationary distribution [13]. We use  $U^{(\lambda)}(\infty)$  to denote the random element that has the stationary distribution. It follows that  $X^{(\lambda)}(\infty) = \Pi_X(U^{(\lambda)}(\infty))$  has the stationary distribution of  $X^{(\lambda)}$ , and  $\tilde{X}^{(\lambda)}(\infty)$  in (2.4) is given by

$$(4.2) \quad \tilde{X}^{(\lambda)}(\infty) = \delta(\Pi_X(U^{(\lambda)}(\infty)) - \gamma n).$$

For  $u \in \mathcal{U}$ , we define

$$(4.3) \quad x = \delta(\Pi_X(u) - \gamma n), \quad q = \Pi_Q(u) \quad \text{and} \quad z = \Pi_X(u) - q.$$

When the CTMC is in state  $u$ , we interpret  $(\Pi_X(u))_i$ ,  $q_i$ , and  $z_i$  as the number of the phase  $i$  customers in system, in queue, and in service, respectively. It follows that  $z \geq 0$ .

Let  $G_{U^{(\lambda)}}$  be the generator of the CTMC  $U^{(\lambda)}$ . To describe it, we introduce the lifting operator  $A$ . For any function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , we define  $Af : \mathcal{U} \rightarrow \mathbb{R}$  by

$$(4.4) \quad Af(u) = f(\delta(\Pi_X(u) - \gamma n)) = f(x).$$

Hence, for any function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , the generator acts on the lifted version  $Af$  as follows:

$$(4.5) \quad \begin{aligned} G_{U^{(\lambda)}} Af(u) &= \sum_{i=1}^d \lambda p_i (f(x + \delta e^{(i)}) - f(x)) + \sum_{i=1}^d \alpha q_i (f(x - \delta e^{(i)}) - f(x)) \\ &\quad + \sum_{i=1}^d \nu_i z_i \left[ \sum_{j=1}^d P_{ij} f(x + \delta e^{(j)} - \delta e^{(i)}) \right. \\ &\quad \left. + (1 - \sum_{j=1}^d P_{ij}) f(x - \delta e^{(i)}) - f(x) \right]. \end{aligned}$$

Observe that  $G_{U^{(\lambda)}} Af(u)$  does not depend on the entire sequence  $u$ ; it depends on  $x$ ,  $q$ , and the function  $f$  only.

**5. The Generator Coupling of Stein's Method.** This section is devoted to developing a generator coupling of Stein's method. This framework will be used in Section 7 to prove Theorem 1.

5.1. *Poisson Equation.* The main idea behind Stein's method is that instead of bounding

$$(5.1) \quad \mathbb{E}h(\tilde{X}^{(\lambda)}(\infty)) - \mathbb{E}h(Y(\infty)),$$

one solves the Poisson equation

$$(5.2) \quad G_Y f_h(x) = h(x) - \mathbb{E}h(Y(\infty)),$$

where the generator  $G_Y$  of the diffusion process  $Y$ , applied to a function  $f \in C^2(\mathbb{R}^d)$ , is given by

$$(5.3) \quad \begin{aligned} G_Y f(x) &= \sum_{i=1}^d \partial_i f(x) \left[ p_i \beta - \nu_i (x_i - p_i(e^T x)^+) - \alpha p_i (e^T x)^+ + \sum_{j=1}^d P_{ji} \nu_j (x_j - p_j(e^T x)^+) \right] \\ &\quad + \frac{1}{2} \sum_{i,j=1}^d \Sigma_{ij} \partial_{ij} f(x) \quad \text{for } x \in \mathbb{R}^d. \end{aligned}$$

Then, to bound the difference in (5.1), it is sufficient to find a bound on

$$(5.4) \quad \mathbb{E}G_Y f_h(\tilde{X}^{(\lambda)}(\infty)).$$

The following lemma, based on the results of [27], guarantees the existence of a solution to (5.2) and provides gradient bounds for it. The proof of this lemma is given in Section A.1.

LEMMA 1. *For any locally Lipschitz function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfying  $|h(x)| \leq |x|^{2m}$ , equation (5.2) has a solution  $f_h$ . Moreover, there exists a constant  $C(m, 1) > 0$  (depending only on  $(\beta, \alpha, p, \nu, P)$ ) such that for  $x \in \mathbb{R}^d$*

$$(5.5) \quad |f_h(x)| \leq C(m, 1)(1 + |x|^2)^m,$$

$$(5.6) \quad |\partial_i f_h(x)| \leq C(m, 1)(1 + |x|^2)^m(1 + |x|),$$

$$(5.7) \quad |\partial_{ij} f_h(x)| \leq C(m, 1)(1 + |x|^2)^m(1 + |x|)^2,$$

$$(5.8) \quad \sup_{y \in \mathbb{R}^d: |y-x| < 1} \frac{|\partial_{ij} f_h(y) - \partial_{ij} f_h(x)|}{|y-x|} \leq C(m, 1)(1 + |x|^2)^m(1 + |x|)^3.$$

5.2. *Generator Coupling.* Let  $W^{(\lambda)}$  denote the random variable  $G_Y f_h(\tilde{X}^{(\lambda)}(\infty))$  in (5.4). To prove  $|\mathbb{E}W^{(\lambda)}|$  small, a common approach in using the Stein's method is to find a coupling  $\tilde{W}^{(\lambda)}$  for  $W^{(\lambda)}$  so that

$$(5.9) \quad \left| \mathbb{E}\tilde{W}^{(\lambda)} \right| \text{ is small, and}$$

$$(5.10) \quad \mathbb{E} \left| W^{(\lambda)} - \tilde{W}^{(\lambda)} \right| \text{ is small.}$$

Constructing an effective coupling is an art that is problem specific. See [47] for a recent survey that includes examples of various couplings.

We use  $\tilde{W}^{(\lambda)} = G_{U^{(\lambda)}} A f_h(U^{(\lambda)}(\infty))$  to construct the coupling, where  $A$  is the lifting operator defined in (4.4). The following lemma justifies the coupling property (5.9).

LEMMA 2. *Let  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfy  $|h(x)| \leq |x|^{2m}$ . The function  $f_h$  given by (5.2) satisfies*

$$(5.11) \quad \mathbb{E}G_{U^{(\lambda)}} A f_h(U^{(\lambda)}(\infty)) = 0.$$

To prove the lemma, we need finite moments of the steady-state system size.

LEMMA 3. (a) *Let  $L(u) = \exp(e^T \Pi_X(u))$  for  $u \in \mathcal{U}$ . Then*

$$(5.12) \quad \mathbb{E}L(U^{(\lambda)}(\infty)) < \infty.$$

(b) *all moments of  $e^T X^{(\lambda)}(\infty)$  are finite.*

PROOF. One may verify that

$$G_{U^{(\lambda)}}L(u) \leq \lambda(\exp(1) - 1)L(u) - \alpha(e^T \Pi_X(u) - n)^+(1 - \exp(-1))L(u).$$

It follows that there exist a positive constant  $C = C(\lambda, n, \alpha)$  such that, whenever  $e^T \Pi_X(u)$  is large enough,

$$(5.13) \quad G_{U^{(\lambda)}}L(u) \leq -CL(u) + 1.$$

Part (a) follows from [42, Theorem 4.2]. Part (b) follows from (5.12) and the equality  $e^T \Pi_X(U^{(\lambda)}(\infty)) = e^T X^{(\lambda)}(\infty)$ .  $\square$

The function  $L(u)$  is said to be a Lyapunov function. Inequality (5.13) is known as a Foster-Lyapunov condition and guarantees that the CTMC is positive recurrent; see, for example, [42].

PROOF OF LEMMA 2. A sufficient condition for (5.11) to hold is given by [35, Proposition 1.1] (alternatively, see [26, Proposition 3]), namely

$$(5.14) \quad \mathbb{E} \left[ \left| G_{U^{(\lambda)}}(U^{(\lambda)}(\infty), U^{(\lambda)}(\infty)) \right| \left| Af_h(U^{(\lambda)}(\infty)) \right| \right] < \infty.$$

Above,  $G_{U^{(\lambda)}}(u, u)$  is the  $u$ th diagonal entry of the generator matrix  $G_{U^{(\lambda)}}$ . In our case, the left side of (5.14) is equal to

$$\begin{aligned} &= \mathbb{E} \left[ \left| G_{U^{(\lambda)}}(U^{(\lambda)}(\infty), U^{(\lambda)}(\infty)) \right| \left| f_h(\tilde{X}^{(\lambda)}(\infty)) \right| \right] \\ &= \mathbb{E} \left[ \left| \lambda + \alpha(e^T X^{(\lambda)}(\infty) - n)^+ + \sum_{i=1}^d \nu_i(X_i^{(\lambda)}(\infty) - Q_i^{(\lambda)}(\infty)) \right| \left| f_h(\tilde{X}^{(\lambda)}(\infty)) \right| \right] \\ &\leq \mathbb{E} \left[ \left| \lambda + (\alpha \vee \max_i \{\nu_i\})e^T X^{(\lambda)}(\infty) \right| \left| f_h(\tilde{X}^{(\lambda)}(\infty)) \right| \right], \end{aligned}$$

where the first equality follows from (4.2) and (4.4). One may apply (5.5) and (5.12) to see that the quantity above is finite.  $\square$

5.3. *Taylor Expansion.* To prove that the coupling  $\tilde{W}^{(\lambda)} = G_{U^{(\lambda)}}Af_h(U^{(\lambda)}(\infty))$  satisfies the coupling property (5.10), we need to prove that

$$\mathbb{E} \left| W^{(\lambda)} - \tilde{W}^{(\lambda)} \right| = \mathbb{E} \left| G_{U^{(\lambda)}}Af_h(U^{(\lambda)}(\infty)) - G_Y f_h(\tilde{X}^{(\lambda)}(\infty)) \right|$$

is small. For that, we compare the generator  $G_{U^{(\lambda)}}$  of the CTMC with  $G_Y$ . By performing Taylor expansion on  $G_{U^{(\lambda)}}Af_h(u)$  in (4.5), one has

$$\begin{aligned}
G_{U^{(\lambda)}}Af_h(u) &= \sum_{i=1}^d \lambda p_i \left( \delta \partial_i f_h(x) + \frac{\delta^2}{2} \partial_{ii} f_h(\xi_i^+) \right) + \alpha q_i \left( -\delta \partial_i f_h(x) + \frac{\delta^2}{2} \partial_{ii} f_h(\xi_i^-) \right) \\
&+ \sum_{i=1}^d \nu_i z_i \left[ \left( 1 - \sum_{j=1}^d P_{ij} \right) \left( -\delta \partial_i f_h(x) + \frac{\delta^2}{2} \partial_{ii} f_h(\xi_i^-) \right) + \sum_{j=1}^d P_{ij} \left( -\delta \partial_i f_h(x) \right. \right. \\
(5.15) \quad &\left. \left. + \delta \partial_j f_h(x) + \frac{\delta^2}{2} \partial_{ii} f_h(\xi_{ij}) + \frac{\delta^2}{2} \partial_{jj} f_h(\xi_{ij}) - \delta^2 \partial_{ij} f_h(\xi_{ij}) \right) \right],
\end{aligned}$$

where  $\xi_i^+ \in [x, x + \delta e^{(i)}]$ ,  $\xi_i^- \in [x - \delta e^{(i)}, x]$  and  $\xi_{ij}$  lies somewhere between  $x$  and  $x - \delta e^{(i)} + \delta e^{(j)}$ . Using the gradient bounds in Lemma 1, we have the following lemma, which will be proved in Section A.2.

LEMMA 4. *There exists a constant  $C(m, 2) > 0$  (depending only on  $(\beta, \alpha, p, \nu, P)$ ) such that for any  $u \in \mathcal{U}$ ,*

$$\begin{aligned}
(5.16) \quad &G_{U^{(\lambda)}}Af_h(u) - G_Y f_h(x) \\
&= \sum_{i=1}^d \partial_i f_h(x) \left[ (\nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j) (\delta q_i - p_i (e^T x)^+) \right] + E(u),
\end{aligned}$$

where  $q$  and  $x$  are as in (4.3),  $\delta$  as in (2.5), and  $E(u)$  is an error term that satisfies

$$|E(u)| \leq \delta C(m, 2) (1 + |x|^2)^m (1 + |x|)^4.$$

**6. State Space Collapse.** One of the challenges we face comes from the fact that our CTMC  $U^{(\lambda)}$  is infinite-dimensional, while the approximating diffusion process is only  $d$ -dimensional. Recall the process  $(X^{(\lambda)}, Q^{(\lambda)})$  defined in (4.1) and the lifting operator  $A$  acting on functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , as defined in (4.4). When acting on the lifted functions  $Af(U^{(\lambda)}(\infty))$ , the CTMC generator  $G_{U^{(\lambda)}}$  depends on both  $\tilde{X}^{(\lambda)}(\infty)$  and  $Q^{(\lambda)}(\infty)$ , but its approximation  $G_Y f(\tilde{X}^{(\lambda)}(\infty))$  only depends on  $\tilde{X}^{(\lambda)}(\infty)$ . This is captured in (5.16) by the term

$$\sum_{i=1}^d \partial_i f_h(x) \left[ (\nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j) (\delta q_i - p_i (e^T x)^+) \right].$$

To bound this term, observe that for any  $1 \leq i \leq d$ ,

$$\begin{aligned}
& \left( \nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j \right) \partial_i f_h(x) (\delta q_i - p_i (e^T x)^+) \\
&= \left( \nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j \right) \left( \partial_i f_h(x) - \partial_i f_h(x - \delta q + p(e^T x)^+) \right) (\delta q_i - p_i (e^T x)^+) \\
&\quad + \left( \nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j \right) \partial_i f_h(x - \delta q + p(e^T x)^+) (\delta q_i - p_i (e^T x)^+) \\
&= \left( \nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j \right) \sum_{k=1}^d \partial_{ik} f_h(\xi) (\delta q_k - p_k (e^T x)^+) (\delta q_i - p_i (e^T x)^+) \\
(6.1) \quad & + \left( \nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j \right) \partial_i f_h(\delta(z - \gamma n) + p(e^T x)^+) (\delta q_i - p_i (e^T x)^+),
\end{aligned}$$

where  $z$ , defined in (4.3), is a vector that represents the number of customers of each type in service, and  $\xi$  is some point between  $x$  and  $x - \delta q + p(e^T x)^+$ . In particular, there exists some constant  $C$  that doesn't depend on  $\lambda$  and  $n$ , such that

$$(6.2) \quad |\xi| \leq |x| + \delta |q| + |p| (e^T x)^+ \leq C |x|,$$

because  $\delta q_i \leq (e^T x)^+$  for each  $1 \leq i \leq d$  (i.e. the number of phase  $i$  customers in queue can never exceed the queue size).

In order to bound the expected value of (6.1), we must prove a relationship between  $\tilde{X}^{(\lambda)}(\infty)$  and  $Q^{(\lambda)}(\infty)$ . Intuitively, the number of customers of phase  $i$  waiting in the queue should be approximately equal to a fraction  $p_i$  of the total queue size. The following two lemmas bound the error caused by the SSC approximation. They are proved at the end of this section.

**LEMMA 5.** *Let  $Z^{(\lambda)}(\infty) = X^{(\lambda)}(\infty) - Q^{(\lambda)}(\infty)$  be the vector representing the number of customers of each type in service in steady-state. Then conditioned on  $(e^T \tilde{X}^{(\lambda)}(\infty))^+$ , the random vectors  $Q^{(\lambda)}(\infty)$  and  $Z^{(\lambda)}(\infty)$  are independent. Furthermore,*

$$(6.3) \quad \mathbb{E} \left[ \delta Q^{(\lambda)}(\infty) - p(e^T \tilde{X}^{(\lambda)}(\infty))^+ \mid (e^T \tilde{X}^{(\lambda)}(\infty))^+ \right] = 0,$$

and for any integer  $m > 0$ , there exists  $C(m, 3) > 0$  (depending only on  $(\beta, \alpha, p, \nu, P)$ ) such that for all  $\lambda > 0$  and  $n \geq 1$  satisfying (1.2),

$$(6.4) \quad \mathbb{E} \left[ \left| \delta Q^{(\lambda)}(\infty) - p(e^T \tilde{X}^{(\lambda)}(\infty))^+ \right|^{2m} \right] \leq \delta^m C(m, 3) \mathbb{E}[(e^T \tilde{X}^{(\lambda)}(\infty))^+]^m,$$

where  $\delta = 1/\sqrt{\lambda}$  as in (2.5).

LEMMA 6. For any integer  $m > 0$ , there exists  $C(m, 4) > 0$  (depending only on  $(\beta, \alpha, p, \nu, P)$ ) such that for any locally Lipschitz function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfying  $|h(x)| \leq |x|^{2m}$ , and all  $\lambda > 0$  and  $n \geq 1$  satisfying (1.2)

$$(6.5) \leq \delta C(m, 4) \mathbb{E} \left[ \left( (e^T \tilde{X}^{(\lambda)}(\infty))^+ \right)^2 \right] \sqrt{\mathbb{E} \left[ 1 + \left| \tilde{X}^{(\lambda)}(\infty) \right|^8 \right]}$$

where  $f_h(x)$  is the solution to the Poisson equation (5.2).

PROOF OF LEMMA 5. We begin by proving (6.4), for which it suffices to show that for all  $\lambda > 0$  and  $n \geq 1$  satisfying (1.2)

$$\mathbb{E} \left[ \left| Q^{(\lambda)}(\infty) - p(e^T X^{(\lambda)}(\infty) - n)^+ \right|^{2m} \right] \leq C(m, 3) \mathbb{E}[(e^T X^{(\lambda)}(\infty) - n)^+]^m.$$

We first prove a version of (6.4) for any finite time  $t \geq 0$ . Then,  $(e^T X^{(\lambda)}(t) - n)^+$  is the total number of customers waiting in queue at time  $t$ . Assume that the system is empty at time  $t = 0$ , i.e.  $X^{(\lambda)}(0) = 0$ . Fix a phase  $i$ . Upon arrival to the system, a customer is assigned to service phase  $i$  with probability  $p_i$ . Consider the sequence  $\{\xi_j : j = 1, 2, \dots\}$ , where  $\xi_j$  is one if the  $j$ th customer to enter the system was assigned to phase  $i$ , and zero otherwise. Then  $\{\xi_j : j = 1, 2, \dots\}$  is a sequence of iid Bernoulli random variables with  $\mathbb{P}(\xi_j = 1) = p_i$ . For  $t > 0$ , define  $A(t)$  and  $B(t)$  to be the total number of customers to have entered the system, and entered service by time  $t$ , respectively. Also let  $\zeta_j(t)$  be the indicator of whether customer  $j$  is still waiting in queue at time  $t$ . Then

$$(6.6) \quad (e^T X^{(\lambda)}(t) - n)^+ = \sum_{j=B(t)+1}^{A(t)} \zeta_j(t),$$

$$(6.7) \quad Q_i^{(\lambda)}(t) = \sum_{j=B(t)+1}^{A(t)} \xi_j \zeta_j(t).$$

Let  $Z^{(\lambda)}(t) = X^{(\lambda)}(t) - Q^{(\lambda)}(t)$  be the vector keeping track of the customer types in service at time  $t$  and let  $B(\ell, p_i)$  be a binomial random variable with  $\ell \in \mathbb{Z}_+$  trials and success probability  $p_i$ . Assuming  $X^{(\lambda)}(0) = 0$ , by a sample path construction of the process  $U^{(\lambda)}$  one can verify that for any time  $t \geq 0$ , the following three properties hold. First, for any  $z \in \mathbb{Z}_+^d$ ,  $a, b \in \mathbb{Z}_+$  with  $a \geq 1$ , and  $x_1, \dots, x_a, y_1, \dots, y_a \in \{0, 1\}$ ,

$$(6.8) \quad \begin{aligned} & \mathbb{P}(\xi_{b+1} = x_1, \dots, \xi_{b+a} = x_a \mid A(t) = b+a, B(t) = b, Z^{(\lambda)}(t) = z, \\ & \quad \zeta_{b+1} = y_1, \dots, \zeta_{b+a} = y_a) \\ &= \mathbb{P}(\xi_1 = x_1) \mathbb{P}(\xi_2 = x_2) \dots \mathbb{P}(\xi_a = x_a) \\ &= p_i^{\sum_{i=1}^a x_i} (1 - p_i)^{a - \sum_{i=1}^a x_i}. \end{aligned}$$

The right side of (6.8) is independent of  $b, z, y_1, \dots, y_a$ . It then follows from (6.6), (6.7) and (6.8) that for any integer  $\ell \geq 1$ ,  $q_i \in \mathbb{Z}_+$ , and  $z \in \mathbb{Z}_+^d$ ,

$$\begin{aligned}
& \mathbb{P}(Q_i^{(\lambda)}(t) = q_i \mid (e^T X^{(\lambda)}(t) - n)^+ = \ell, Z^{(\lambda)}(t) = z) \\
&= \mathbb{P}(Q_i^{(\lambda)}(t) = q_i \mid (e^T X^{(\lambda)}(t) - n)^+ = \ell) \\
(6.9) \quad &= \mathbb{P}(B(\ell, p_i) = q_i).
\end{aligned}$$

Since (6.9) holds for all  $t \geq 0$ , it holds in stationarity as well.

We now say a few words about how to construct  $U^{(\lambda)}$  and argue (6.8)–(6.9). One would start with four primitive sequences: a sequence of inter-arrival times, potential service times, patience times, and routing decisions. The sequence of potential service times would hold all the service information about each customer provided they were patient enough to get into service. The routing sequence would represent the phase each customer is assigned upon entering the system.

To see why (6.8) is true, we first observe that at any time  $t > 0$ , the random variable  $A(t)$  depends only on the inter-arrival time primitives; in particular, it is independent of the routing sequence  $\{\xi_j, j \geq 1\}$ . Second, any customer to arrive after customer number  $B(t) = b$  has no impact on any of the servers at any point in time during  $[0, t]$ . In particular, the primitives including  $\{\xi_{b+j}, j \geq 1\}$  associated to those customers are independent of  $B(t) = b$  and  $Z^{(\lambda)}(t)$ . Lastly, the decisions of those customers whether to abandon or not by time  $t$  depends only on their arrival times, patience times, and the service history in the interval  $[0, t]$ . In particular, the sequence  $\{\zeta_{b+j}(t), j \geq 1\}$  is independent of  $\{\xi_{b+j}, j \geq 1\}$ . This proves the the first equality in (6.8).

We now move on to complete the proof of this lemma. We use (6.9) to see that for any positive integer  $N$ ,

$$\begin{aligned}
& \mathbb{E} \left( [Q_i^{(\lambda)}(t) - p_i(e^T X^{(\lambda)}(t) - n)^+]^{2m} \mathbf{1}_{\{(e^T X^{(\lambda)}(t) - n)^+ \leq N\}} \right) \\
&= \sum_{\ell=1}^N \mathbb{E} \left[ (B(\ell, p_i) - p_i \ell)^{2m} \right] \mathbb{P}((e^T X^{(\lambda)}(t) - n) = \ell) \\
&\leq \sum_{\ell=1}^N C(m, 6) \ell^m \mathbb{P}((e^T X^{(\lambda)}(t) - n) = \ell) \\
(6.10) \quad &= C(m, 6) \mathbb{E} \left( [(e^T X^{(\lambda)}(t) - n)^+]^m \mathbf{1}_{\{(e^T X^{(\lambda)}(t) - n)^+ \leq N\}} \right),
\end{aligned}$$

where we have used the fact that there is a constant  $C(m, 6) > 0$  such that

$$\mathbb{E} \left[ (B(\ell, p_i) - p_i \ell)^{2m} \right] \leq C(m, 6) \ell^m \quad \text{for all } \ell \geq 1;$$

see, for example, (4.10) of [40]. Letting  $t \rightarrow \infty$  in both sides of (6.10), by the dominated

convergence theorem, one has

$$\begin{aligned} & \mathbb{E}\left([Q_i^{(\lambda)}(\infty) - p_i(e^T X^{(\lambda)}(\infty) - n)^+]^{2m} \mathbf{1}_{\{(e^T X^{(\lambda)}(\infty) - n)^+ \leq N\}}\right) \\ & \leq C(m, 6) \mathbb{E}\left([e^T X^{(\lambda)}(\infty) - n]^m \mathbf{1}_{\{(e^T X^{(\lambda)}(\infty) - n)^+ \leq N\}}\right). \end{aligned}$$

Letting  $N \rightarrow \infty$ , by the monotone convergence theorem, one has

$$\mathbb{E}(Q_i^{(\lambda)}(\infty) - p_i(e^T X^{(\lambda)}(\infty) - n)^+)^{2m} \leq C(m, 6) \mathbb{E}[(e^T X^{(\lambda)}(\infty) - n)^+]^m.$$

Then (6.4) follows from this inequality for each  $i$  and the fact that there is a constant  $B_m > 0$  such that  $|x|^{2m} \leq B_m \sum_{i=1}^d (x_i)^{2m}$  for all  $x \in \mathbb{R}^d$ . One can check that (6.3) can be obtained by an argument very similar to the one used to prove (6.4).  $\square$

PROOF OF LEMMA 6. Recall that

$$Z^{(\lambda)}(\infty) = X^{(\lambda)}(\infty) - Q^{(\lambda)}(\infty)$$

is the vector representing the number of customers of each type in service in steady-state. Then from (6.1) we have

$$\begin{aligned} & \mathbb{E}\left[\partial_i f_h(\tilde{X}^{(\lambda)}(\infty))(\delta Q_i^{(\lambda)}(\infty) - p_i(e^T \tilde{X}^{(\lambda)}(\infty))^+)\right] \\ & = \sum_{k=1}^d \mathbb{E}\left[\partial_{ik} f_h(\xi)(\delta Q_k^{(\lambda)}(\infty) - p_k(e^T \tilde{X}^{(\lambda)}(\infty))^+)(\delta Q_i^{(\lambda)}(\infty) - p_i(e^T \tilde{X}^{(\lambda)}(\infty))^+)\right] \\ & \quad + \mathbb{E}\left[\partial_i f_h(\delta(Z^{(\lambda)}(\infty) - \gamma n) + p(e^T \tilde{X}^{(\lambda)}(\infty))^+)(\delta Q_i^{(\lambda)}(\infty) - p_i(e^T \tilde{X}^{(\lambda)}(\infty))^+)\right]. \end{aligned}$$

By Lemma 5, the second expected value equals zero. For the first term, one can use the Cauchy-Schwarz inequality, together with the gradient bound (5.7) and the SSC result (6.4) to see that for all  $1 \leq i, k \leq d$ ,

$$\begin{aligned} & \mathbb{E}\left[\partial_{ik} f_h(\xi)(\delta Q_k^{(\lambda)}(\infty) - p_k(e^T \tilde{X}^{(\lambda)}(\infty))^+)(\delta Q_i^{(\lambda)}(\infty) - p_i(e^T \tilde{X}^{(\lambda)}(\infty))^+)\right] \\ & \leq \sqrt{\mathbb{E}\left[(\partial_{ik} f_h(\xi))^2\right]} \sqrt{\mathbb{E}\left[(\delta Q_k^{(\lambda)}(\infty) - p_k(e^T \tilde{X}^{(\lambda)}(\infty))^+)^4\right]} \sqrt{\mathbb{E}\left[(\delta Q_i^{(\lambda)}(\infty) - p_i(e^T \tilde{X}^{(\lambda)}(\infty))^+)^4\right]} \\ & \leq \delta C(2, 3) \mathbb{E}\left[(e^T \tilde{X}^{(\lambda)}(\infty))^+\right]^2 \sqrt{\mathbb{E}\left[(\partial_{ik} f_h(\xi))^2\right]} \\ & \leq \delta C(2, 3) \mathbb{E}\left[(e^T \tilde{X}^{(\lambda)}(\infty))^+\right]^2 C(m, 1) \sqrt{\mathbb{E}\left[(1 + |\xi|^2)^2 (1 + |\xi|^4)\right]}. \end{aligned}$$

We now combine everything together with the fact that  $\xi$  satisfies (6.2) to conclude that there exists a constant  $C(m, 4)$  that does not depend on  $\lambda$  or  $n$ , such that

$$\begin{aligned} & \left| \sum_{i=1}^d \partial_i \mathbb{E} \left[ f_h(\tilde{X}^{(\lambda)}(\infty)) \left[ (\nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j) (\delta Q_i^{(\lambda)}(\infty) - p_i (e^T \tilde{X}^{(\lambda)}(\infty))^+ \right) \right] \right| \\ & \leq \delta C(m, 4) \mathbb{E} \left[ (e^T \tilde{X}^{(\lambda)}(\infty))^+ \right]^2 \sqrt{\mathbb{E} \left[ 1 + \left| \tilde{X}^{(\lambda)}(\infty) \right|^8 \right]}, \end{aligned}$$

which concludes the proof of the lemma.  $\square$

**7. Proof of Theorem 1.** To prove Theorem 1, we need an additional lemma on uniform bounds for moments of scaled system size. It will be proved in Section A.3.

LEMMA 7. *For any integer  $m \geq 0$ , there exists a constant  $C(m, 5) > 0$  (depending only on  $(\beta, \alpha, p, \nu, P)$ ) such that*

$$(7.1) \quad \mathbb{E} \left| \tilde{X}^{(\lambda)}(\infty) \right|^m \leq C(m, 5).$$

We remark that in the special case when the service time distribution is taken to be hyper-exponential, it is proved in [21] that

$$\limsup_{\lambda \rightarrow \infty} \mathbb{E} \exp \left( \theta \left| \tilde{X}^{(\lambda)}(\infty) \right| \right) < \infty$$

for  $\theta$  in some positive interval. The proof relies on a result that allows one to compare the system with an infinite-server system, whose stationary distribution is known to be Poisson.

PROOF OF THEOREM 1. It follows from Lemmas 4 and 6 that

$$\begin{aligned} & \left| \mathbb{E} h(\tilde{X}^{(\lambda)}(\infty)) - \mathbb{E} h(Y(\infty)) \right| = \left| \mathbb{E} G_{U^{(\lambda)}} A f_h(U^{(\lambda)}(\infty)) - \mathbb{E} G_Y f_h(\tilde{X}^{(\lambda)}(\infty)) \right| \\ & \leq \left| \sum_{i=1}^d \mathbb{E} \left[ \partial_i f_h(\tilde{X}^{(\lambda)}(\infty)) \left[ (\nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j) (\delta Q_i^{(\lambda)}(\infty) - p_i (e^T \tilde{X}^{(\lambda)}(\infty))^+ \right) \right] \right| \\ & \quad + \delta C(m, 2) \mathbb{E} \left[ (1 + \left| \tilde{X}^{(\lambda)}(\infty) \right|^2)^m (1 + \left| \tilde{X}^{(\lambda)}(\infty) \right|^4) \right] \\ & \leq \delta C(m, 4) \mathbb{E} \left[ ((e^T \tilde{X}^{(\lambda)}(\infty))^+)^2 \right] \sqrt{\mathbb{E} \left[ 1 + \left| \tilde{X}^{(\lambda)}(\infty) \right|^8 \right]} \\ (7.2) \quad & \quad + \delta C(m, 2) \mathbb{E} \left[ (1 + \left| \tilde{X}^{(\lambda)}(\infty) \right|^2)^m (1 + \left| \tilde{X}^{(\lambda)}(\infty) \right|^4) \right]. \end{aligned}$$

By Lemma 7, there are constants  $B_1(m), B_2(m) > 0$  (depending only on  $(\beta, \alpha, p, \nu, P)$ ) such that

$$\begin{aligned} \mathbb{E}\left[\left((e^T \tilde{X}^{(\lambda)}(\infty))^+\right)^2\right] \sqrt{\mathbb{E}\left[1 + \left|\tilde{X}^{(\lambda)}(\infty)\right|^8\right]} &\leq B_1(m), \\ \mathbb{E}\left[\left(1 + \left|\tilde{X}^{(\lambda)}(\infty)\right|^2\right)^m \left(1 + \left|\tilde{X}^{(\lambda)}(\infty)\right|^4\right)\right] &\leq B_2(m). \end{aligned}$$

Therefore, the right side of (7.2) is less than or equal to

$$\begin{aligned} &\delta C(m, 4)B_1(m) + \delta C(m, 2)B_2(m) \\ &\leq \left(C(m, 4)B_1(m) + C(m, 2)B_2(m)\right) \frac{1}{\sqrt{\lambda}} \quad \text{for } \lambda > 0. \end{aligned}$$

This concludes the proof of Theorem 1.  $\square$

**Acknowledgements.** The authors thank Shuangchi He, Josh Reed and John Pike for stimulating discussions. They also thank the participants of Applied Probability & Risk Seminar in Fall 2014 at Columbia University for their feedback on this research. This research is supported in part by NSF Grants CNS-1248117 and CMMI-1335724.

## APPENDIX A: PROOFS

**A.1. Proof of Lemma 1 (Gradient Bounds).** Before proving the lemma, we first state the common quadratic Lyapunov function introduced in [17]. This Lyapunov function plays a key role in our paper. As in (5.24) of [17], for  $x \in \mathbb{R}^d$ , define

$$(A.1) \quad V(x) = (e^T x)^2 + \kappa[x - p\phi(e^T x)]' M[x - p\phi(e^T x)],$$

where  $\kappa > 0$  is some constant,  $M$  is some  $d \times d$  positive definite matrix, and the function  $\phi$  is a smooth approximation to  $x \mapsto x^+$  and is defined by

$$\phi(x) = \begin{cases} x, & \text{if } x \geq 0, \\ -\frac{1}{2}\epsilon, & \text{if } x \leq -\epsilon, \\ \text{smooth,} & \text{if } -\epsilon < x < 0. \end{cases}$$

In (5.24) of [17], the authors use  $\tilde{Q}$  to represent the positive definite matrix that we called  $M$  in (A.1). We use  $M$  instead of  $\tilde{Q}$  on purpose, to avoid any potential confusion with the queue size  $Q(t)$ . For our purposes, “smooth” means that  $\phi$  can be anything as long as  $\phi \in C^3(\mathbb{R}^d)$ . We require that the “smooth” part of  $\phi$  also satisfies  $-\frac{1}{2}\epsilon < \phi(x) < x$  and  $0 \leq \phi'(x) \leq 1$ . For example,  $\phi$  can be taken to be a polynomial of sufficiently high degree on  $(-\epsilon, 0)$  and this will satisfy our requirements. The vector  $p$  is as in (2.6). The constant

$\kappa$  and matrix  $M$  are chosen just as in [17]; their exact values are not important to us. In their paper, they show that  $V$  satisfies

$$G_Y V(x) \leq -c_1 V(x) + c_2 \quad \text{for all } x \in \mathbb{R}^d$$

for some positive constants  $c_1, c_2$ ; this result requires  $\alpha > 0$ , i.e. a strictly positive abandonment rate. Before proceeding to the proof of Lemma 1, we state two bounds on  $V$  that shall be useful in the future. For some constant  $C > 0$ ,

$$(A.2) \quad V(x) \leq C(1 + |x|^2),$$

$$(A.3) \quad |x|^2 \leq C(1 + V(x)).$$

The first is immediate from the form of  $V$ , while the second is proved in [17].

PROOF OF LEMMA 1. Without loss of generality, we may assume that  $h(0) = 0$ , otherwise one may consider  $h(x) - h(0)$ . This lemma is essentially a restatement of equation (22) and equation (40) from the discussion that follows after [27, Theorem 4.1]. We verify that (22) and (40) are applicable in our case by first confirming that we have a function satisfying assumption 3.1 of [27]. Recalling the definition of  $V$  from (A.1), when  $\phi$  is taken to be a polynomial (of sufficiently high degree to guarantee  $V \in C^3(\mathbb{R}^d)$ ), the function

$$1 + V(x)$$

satisfies assumption 3.1. To verify condition (17) of Assumption 3.1, one observes that

$$X^{(\lambda)}(t) \leq X^{(\lambda)}(0) + n + A^{(\lambda)}(t),$$

where  $A^{(\lambda)}(t)$  is the total number of arrivals to the system by time  $t$  and it is a Poisson random variable with mean  $\lambda t$  for each  $t \geq 0$ . The properties of Poisson processes then yield (17). By [27, Remark 3.4],

$$C(1 + V(x))^m$$

also satisfies assumption 3.1 for any constant  $C > 0$ . Since we require that  $|h(x)| \leq |x|^m$ , by (A.3) we have

$$|h(x) - \mathbb{E}h(Y(\infty))| \leq |x|^m + \mathbb{E}|Y(\infty)|^m \leq C_m(1 + V(x))^m.$$

The finiteness of  $\mathbb{E}|Y(\infty)|^m$  is guaranteed because one of the conditions of assumption 3.1 is that

$$G_Y(1 + V(x))^m \leq -c_1(1 + V(x))^m + c_2$$

for some positive constants  $c_1$  and  $c_2$ . Therefore, equation (22) gives us (5.5) and equation (40) gives us (5.6) and (5.7). We get (5.8) by observing that in the discussion preceding (40), everything still holds if we replace  $B_x(\bar{l}/\sqrt{n})$  by an open ball of radius 1 centered at  $x$ . We wish to point out that the constants in (40) and (22) do not depend on the choice of function  $h$ .  $\square$

**A.2. Proof of Lemma 4 (Generator Difference).** The main idea here is that  $G_Y f_h(x)$  is hidden within  $G_{U^{(\lambda)}} A f_h(u)$ , where the lifting operator  $A$  is in (4.4). We algebraically manipulate the Taylor expansion of  $G_{U^{(\lambda)}} A f_h(u)$  to make this evident. First, we first rearrange the terms in the Taylor expansion (5.15) to group them by partial derivatives. Thus,  $G_{U^{(\lambda)}} A f_h(u)$  equals

$$\begin{aligned}
& \sum_{i=1}^d \delta \partial_i f_h(x) \left[ p_i \lambda - \alpha q_i - \nu_i z_i + \sum_{j=1}^d P_{ji} \nu_j z_j \right] \\
& + \sum_{i=1}^d \frac{\delta^2}{2} \partial_{ii} f_h(x) \left[ p_i \lambda + \alpha q_i + \nu_i z_i + \sum_{j=1}^d P_{ji} \nu_j z_j \right] - \sum_{i \neq j}^d \delta^2 \partial_{ij} f_h(x) [P_{ij} \nu_i z_i] \\
& + \sum_{i=1}^d \frac{\delta^2}{2} \left( \partial_{ii} f_h(\xi_i^-) - \partial_{ii} f_h(x) \right) \left[ \alpha q_i + \left( 1 - \sum_{j=1}^d P_{ij} \right) \nu_i z_i \right] \\
& + \sum_{i=1}^d \frac{\delta^2}{2} \left( \partial_{ii} f_h(\xi_i^+) - \partial_{ii} f_h(x) \right) [\lambda p_i] - \sum_{i \neq j}^d \delta^2 \left( \partial_{ij} f_h(\xi_{ij}) - \partial_{ij} f_h(x) \right) [P_{ij} \nu_i z_i] \\
& + \sum_{i=1}^d \sum_{j=1}^d \frac{\delta^2}{2} \left( \partial_{ii} f_h(\xi_{ij}) - \partial_{ii} f_h(x) \right) [P_{ij} \nu_i z_i + P_{ji} \nu_j z_j].
\end{aligned}$$

To proceed we observe that (2.3) gives us the identity

$$(A.4) \quad -\nu_i \gamma_i n + \sum_{j=1}^d P_{ji} \nu_j \gamma_j n = -n p_i.$$

Recall the form of  $G_Y f_h(x)$  from (5.3). From the form of  $\Sigma$  in (2.7), we see that

$$(A.5) \quad \Sigma_{ii} = 2 \left( p_i + \sum_{j=1}^d P_{ji} \gamma_j \nu_j \right), \quad \Sigma_{ij} = -(P_{ij} \nu_i \gamma_i + P_{ji} \nu_j \gamma_j) \text{ for } j \neq i$$

using (5.3), (A.4) and (A.5), the difference  $G_{U^{(\lambda)}} A f_h(u) - G_Y f_h(x)$  becomes

$$\begin{aligned}
\text{(A.6)} \quad & \sum_{i=1}^d \partial_i f_h(x) \left[ (\nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j) (\delta q_i - p_i (e^T x)^+) \right] \\
& + \sum_{i=1}^d \partial_{ii} f_h(x) \left[ \sum_{j=1}^d P_{ji} \nu_j \gamma_j \right] (n\delta^2 - 1) - \sum_{i \neq j}^d \partial_{ij} f_h(x) \left[ P_{ij} \nu_i \gamma_i + P_{ji} \nu_j \gamma_j \right] (n\delta^2 - 1) \\
& - \sum_{i=1}^d \frac{\delta^2}{2} \partial_{ii} f_h(x) \left[ p_i (\lambda - n) - \alpha q_i - \nu_i (z_i - \gamma_i n) - \sum_{j=1}^d P_{ji} \nu_j (z_j - \gamma_j n) \right] \\
& - \sum_{i \neq j}^d \frac{\delta^2}{2} \partial_{ij} f_h(x) \left[ P_{ij} \nu_i (z_i - \gamma_i n) + P_{ji} \nu_j (z_j - \gamma_j n) \right] \\
& + \sum_{i=1}^d \frac{\delta^2}{2} (\partial_{ii} f_h(\xi_i^-) - \partial_{ii} f_h(x)) \left[ \alpha q_i + (1 - \sum_{j=1}^d P_{ij}) \nu_i z_i \right] \\
& + \sum_{i=1}^d \frac{\delta^2}{2} (\partial_{ii} f_h(\xi_i^+) - \partial_{ii} f_h(x)) \left[ \lambda p_i \right] - \sum_{i \neq j}^d \delta^2 (\partial_{ij} f_h(\xi_{ij}) - \partial_{ij} f_h(x)) \left[ P_{ij} \nu_i z_i \right] \\
& + \sum_{i=1}^d \sum_{j=1}^d \frac{\delta^2}{2} (\partial_{ii} f_h(\xi_{ij}) - \partial_{ii} f_h(x)) \left[ P_{ij} \nu_i z_i + P_{ji} \nu_j z_j \right].
\end{aligned}$$

We remind the reader that our target is to prove that

$$\begin{aligned}
& G_{U(\lambda)} A f_h(u) - G_Y f_h(x) \\
& = \sum_{i=1}^d \partial_i f_h(x) \left[ (\nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j) (\delta q_i - p_i (e^T x)^+) \right] + E(u),
\end{aligned}$$

where  $E(u)$  is an error term that satisfies

$$|E(u)| \leq \delta C(m, 2)(1 + |x|^2)^m (1 + |x|)^4.$$

We choose  $E(u)$  to be all the terms in (A.6) except for the first line. We now describe how to bound  $|E(u)|$ . Most of the summands in (A.6) look as follows: a term in large square brackets multiplied by some partial derivative of  $f_h$ . The partial derivatives are very easy to bound; we simply use (5.6) - (5.8). We wish to point out that  $\xi_i^+$ ,  $\xi_i^-$  and  $\xi_{ij}$  lie within distance  $2\delta$  of  $x$ . When  $2\delta < 1$ , (5.8) implies

$$\text{(A.7)} \quad |\partial_{ij} f_h(\xi) - \partial_{ij} f_h(x)| \leq 2\delta C(1 + |x|^2)^m (1 + |x|)^3$$

for some constant  $C > 0$  (i.e. an extra  $\delta$  term is gained). When  $2\delta \geq 1$  (by Remark 1 this occurs in finitely many cases), we may use (5.7) to obtain (A.7) with a redefined  $C$ . From

here on out, we shall let  $C > 0$  be a generic positive constant that will change from line to line, but will always be independent of  $\lambda$  and  $n$ .

Now we shall list the facts needed to bound all the square bracket terms in (A.6) except for the very first one. Recall that we are operating in the Halfin-Whitt regime as defined by (1.2). Therefore,

$$(n\delta^2 - 1) = \delta\beta \text{ and } \delta(\lambda - n) = -\beta.$$

Furthermore, it must be true that

$$\delta q_i \leq (e^T x)^+ \leq C|x|,$$

as the number of phase  $i$  customers may never exceed the total queue size. Next,

$$|\delta(z_i - \gamma_i n)| = |x_i - \delta q_i| \leq C|x|$$

and lastly,

$$|\delta^2 z_i| \leq |\delta^2 \gamma_i n| + |\delta^2(z_i - \gamma_i n)| \leq C(1 + |x|).$$

It is now a simple matter to verify that the inequalities above, combined with the bounds on the partials of  $f_h$  are all that it takes to achieve our desired upper bound.

**A.3. Proof of Lemma 7 (Moment Bounds).** We first provide an intuitive roadmap for the proof. The goal is to show that a Lyapunov function for the diffusion process is also a Lyapunov function for the CTMC; this has two parts to it. In the first part of this proof, we compare how the two generators  $G_{U(\lambda)}$  and  $G_Y$  act on this Lyapunov function, obtaining an upper bound for the difference  $G_{U(\lambda)} - G_Y$  in (A.12). One notes that the right hand side of (A.12) is unbounded. This is due to the difference in dimensions of the CTMC and diffusion process. To overcome this difficulty, we move on to the second part of the proof, which exploits our SSC result in Lemma 5 to bound the expectation of the right hand side of (A.12). We end up with a recursive relationship that guarantees the  $2m$ th moment is bounded (uniformly in  $\lambda$  and  $n$  satisfying (1.2)) provided that the  $m$ th moment is. Finally, we rely on prior results obtained in [13] for a uniform bound on the first moment.

We remark that a version of this lemma was already proved [27, Theorem 3.3] for the case where the dimension of the CTMC equals the dimension of the diffusion process. However, the difference in dimensions poses an additional technical challenge, which is overcome in the second part of this proof.

Its enough to prove (7.1) for the cases when  $m = 2^j$  for some  $j \geq 0$ . Furthermore, we may assume that  $\lambda \geq 4$  because by Remark 1, there are only finitely many cases when  $\lambda < 4$ . In all those cases,  $\mathbb{E} \left| \tilde{X}^{(\lambda)}(\infty) \right|^m < \infty$  by (5.12). Throughout the proof, we shall use  $C, C_1, C_2, C_3, C_4$  to denote generic positive constants that may change from line to line. They may depend on  $(m, \beta, \alpha, p, \nu, P)$ , but will be independent of both  $\lambda$  and  $n$ . Define

$$V_m(x) = (1 + V(x))^m,$$

where  $V$  is as in (A.1). By [27, Remark 3.4],  $V_m$  also satisfies

$$G_Y V_m(x) \leq -C_1 V_m(x) + C_2$$

as long as  $V \in C^3(\mathbb{R}^d)$  and satisfies condition (30) of [27], which is easy to verify. To prove the lemma, we will show that for large enough  $\lambda$ ,  $V$  satisfies

$$\mathbb{E}G_{U^{(\lambda)}} AV_m(U^{(\lambda)}(\infty)) \leq -C_1 \mathbb{E}V_m(\tilde{X}^{(\lambda)}(\infty)) + C_2,$$

where  $A$  is the lifting operator defined in (4.4). We begin by observing

$$(A.8) \quad G_{U^{(\lambda)}} AV_m \leq G_{U^{(\lambda)}} AV_m - G_Y V_m + G_Y V_m \leq G_{U^{(\lambda)}} AV_m - G_Y V_m - C_1 V_m + C_2.$$

Using (A.6), we write  $G_{U^{(\lambda)}} AV_m - G_Y V_m$  as

$$\begin{aligned} & \sum_{i=1}^d \partial_i V_m(x) \left[ (\nu_i - \alpha - \sum_{j=1}^d P_{ji} \nu_j) (\delta q_i - p_i (e^T x)^+) \right] \\ & + \sum_{i=1}^d \partial_{ii} V_m(x) \left[ \sum_{j=1}^d P_{ji} \nu_j \gamma_j \right] (n\delta^2 - 1) - \sum_{i \neq j}^d \partial_{ij} V_m(x) \left[ P_{ij} \nu_i \gamma_i + P_{ji} \nu_j \gamma_j \right] (n\delta^2 - 1) \\ & - \sum_{i=1}^d \frac{\delta^2}{2} \partial_{ii} V_m(x) \left[ p_i (\lambda - n) - \alpha q_i - \nu_i (z_i - \gamma_i n) - \sum_{j=1}^d P_{ji} \nu_j (z_j - \gamma_j n) \right] \\ & - \sum_{i \neq j}^d \frac{\delta^2}{2} \partial_{ij} V_m(x) \left[ P_{ij} \nu_i (z_i - \gamma_i n) + P_{ji} \nu_j (z_j - \gamma_j n) \right] \\ & + \sum_{i=1}^d \frac{\delta^2}{2} (\partial_{ii} V_m(\xi_i^-) - \partial_{ii} V_m(x)) \left[ \alpha q_i + (1 - \sum_{j=1}^d P_{ij}) \nu_i z_i \right] \\ & + \sum_{i=1}^d \frac{\delta^2}{2} (\partial_{ii} V_m(\xi_i^+) - \partial_{ii} V_m(x)) \left[ \lambda p_i \right] - \sum_{i \neq j}^d \delta^2 (\partial_{ij} V_m(\xi_{ij}) - \partial_{ij} V_m(x)) \left[ P_{ij} \nu_i z_i \right] \\ & + \sum_{i=1}^d \sum_{j=1}^d \frac{\delta^2}{2} (\partial_{ii} V_m(\xi_{ij}) - \partial_{ii} V_m(x)) \left[ P_{ij} \nu_i z_i + P_{ji} \nu_j z_j \right]. \end{aligned}$$

Now we wish to bound the derivatives of  $V_m$ . By [27, Remark 3.4],  $V_m$  satisfies (16) and (30) of [27], namely

$$(A.9) \quad \sup_{|y| \leq 1} \frac{V_m(x+y)}{V_m(x)} \leq C$$

and

$$(A.10) \quad (|\partial_i V_m(x)| + |\partial_{ij} V_m(x)| + |\partial_{ijk} V_m(x)|)(1 + |x|) \leq C V_m(x).$$

For  $\xi$  being one of  $\xi_i^+$ ,  $\xi_i^-$  or  $\xi_{ij}$ ,

$$(A.11) \quad |\partial_{ij} V_m(\xi) - \partial_{ij} V_m(x)| (1 + |x|) \leq \delta |\partial_{ijj} V_m(\eta) + \partial_{ijj} V_m(\eta)| (1 + |x|) \leq C \delta V_m(x),$$

where the first inequality comes from a Taylor expansion and the second inequality follows by (A.10), the fact that  $|\eta - x| \leq 2\delta < 1$  and by (A.9). Following the exact same argument that we used to bound (A.6) in the proof of Lemma 4 (with (A.10) and (A.11) replacing the gradient bounds of  $f_h$  there), we get

$$G_{U^{(\lambda)}} A V_m - G_Y V_m \leq C \delta V_m(x) + C \sum_{i=1}^d |\partial_i V_m(x)| \left[ |q_i - p_i(e^T x)^+| \right].$$

Differentiating  $V$ , we see that

$$(\nabla V(x))^T = 2(e^T x)e^T + 2\kappa(x^T - p^T \phi(e^T x))\tilde{Q}(I - pe^T \phi'(e^T x)).$$

Combined with the fact that  $0 \leq \phi'(x) \leq 1$ , it is clear that

$$|\partial_i V(x)| \leq C(1 + |x|).$$

Therefore,

$$(A.12) \quad G_{U^{(\lambda)}} A V_m - G_Y V_m \leq C \delta V_m(x) + C \sum_{i=1}^d m V_{m-1}(x)(1 + |x|) \left[ |q_i - p_i(e^T x)^+| \right].$$

It remains to find an appropriate bound for

$$V_{m-1}(x)(1 + |x|) \left[ |q_i - p_i(e^T x)^+| \right] = \delta V_{m-1}(x)(1 + |x|) \left[ \frac{|q_i - p_i(e^T x)^+|}{\delta} \right].$$

We have

$$\begin{aligned}
& \delta V_{m-1}(x)(1+|x|) \left[ \frac{|q_i - p_i(e^T x)^+|}{\delta} \right] \\
& \leq \sqrt{\delta} V_{m-1}(x)(1+|x|)^2 + \sqrt{\delta} V_{m-1}(x) \left[ \frac{|q_i - p_i(e^T x)^+|^2}{\delta} \right] \\
& \leq C\sqrt{\delta} V_m(x) + \sqrt{\delta} V_{m-2}(x) V_2(x) + \sqrt{\delta} V_{m-2}(x) \left[ \frac{|q_i - p_i(e^T x)^+|^2}{\delta} \right]^2 \\
& \leq C\sqrt{\delta} V_m(x) + \sqrt{\delta} V_m(x) + \sqrt{\delta} V_{m-4}(x) V_4(x) + \sqrt{\delta} V_{m-4}(x) \left[ \frac{|q_i - p_i(e^T x)^+|^2}{\delta} \right]^4 \\
& \leq \dots \\
\text{(A.13)} \leq & C\sqrt{\delta} V_m(x) + \sqrt{\delta} \left[ \frac{|q_i - p_i(e^T x)^+|^2}{\delta} \right]^m,
\end{aligned}$$

where in the last inequality, we used the fact that  $m = 2^j$ . Using (A.8), (A.12) and (A.13),

$$G_{U^{(\lambda)}} AV_m(u) \leq -V_m(x)(C_1 - \sqrt{\delta} C_3) + C_2 + \sqrt{\delta} C_4 \sum_{i=1}^d \left[ \frac{|q_i - p_i(e^T x)^+|^2}{\delta} \right]^m,$$

where  $x$  and  $q$  are related to  $u$  by (4.3). The arguments in the proof of Lemma 2 can be used to show

$$\mathbb{E} G_{U^{(\lambda)}} AV_m(U^{(\lambda)}(\infty)) = 0.$$

Therefore, for  $\delta$  small enough,

$$\begin{aligned}
\mathbb{E} \left| \tilde{X}^{(\lambda)}(\infty) \right|^{2m} & \leq C \mathbb{E} V_m(\tilde{X}^{(\lambda)}(\infty)) \\
& \leq \frac{C}{(C_1 - \sqrt{\delta} C_3)} \left( C_2 + \sqrt{\delta} C_4 \sum_{i=1}^d \frac{\mathbb{E} \left| \delta Q_i^{(\lambda)}(\infty) - p_i(e^T \tilde{X}^{(\lambda)}(\infty))^+ \right|^{2m}}{\delta^m} \right).
\end{aligned}$$

By (6.4), it follows that

$$\mathbb{E} \left| \tilde{X}^{(\lambda)}(\infty) \right|^{2m} \leq \frac{C}{C_1 - \sqrt{\delta} C_3} \left( 1 + \sqrt{\delta} \mathbb{E} [(e^T \tilde{X}^{(\lambda)}(\infty))^+]^m \right).$$

Hence, we have a recursive relationship that guarantees

$$\sup_{\lambda > 0} \mathbb{E} \left| \tilde{X}^{(\lambda)}(\infty) \right|^{2m} < \infty$$

whenever

$$\sup_{\lambda>0} \mathbb{E}[(e^T \tilde{X}^{(\lambda)}(\infty))^+]^m < \infty.$$

To conclude, we need to verify that

$$\sup_{\lambda>0} \mathbb{E}[(e^T \tilde{X}^{(\lambda)}(\infty))^+] < \infty,$$

but this was proved in equation (5.2) of [13].

## REFERENCES

- [1] AKSIN, Z., ARMONY, M. and MEHROTRA, V. (2007). The modern call center: a multi-disciplinary perspective on operations management research. *Production and Operations Management* **16** 665–688.
- [2] ARMONY, M., ISRAELIT, S., MANDELBAUM, A., MARMOR, Y. N., TSEYTLIN, Y. and YOM-TOV, G. B. (2011). Patient flow in hospitals: A data-based queueing-science perspective. *working paper*.
- [3] ASMUSSEN, S. (2003). *Applied probability and queues*, second ed. *Applications of Mathematics (New York)* **51**. Springer-Verlag, New York. Stochastic Modelling and Applied Probability. [MR1978607 \(2004f:60001\)](#)
- [4] BARBOUR, A. D. (1990). Stein's method for diffusion approximations. *Probability Theory and Related Fields* **84** 297–322.
- [5] BARBOUR, A. D. (1988). Stein's Method and Poisson Process Convergence. *Journal of Applied Probability* **25** pp. 175–184.
- [6] BELL, S. L. and WILLIAMS, R. J. (2005). Dynamic scheduling of a parallel server system in heavy traffic with complete resource pooling: asymptotic optimality of a threshold policy. *Electronic Journal of Probability* **10** 1044–1115.
- [7] BOROVKOV, A. (1964). Some limit theorems in the theory of mass service, I. *Theory of Probability and its Applications* **9** 550–565.
- [8] BOROVKOV, A. (1965). Some limit theorems in the theory of mass service, II. *Theory of Probability and its Applications* **10** 375–400.
- [9] BRAMSON, M. (1998). State space collapse with application to heavy traffic limits for multiclass queueing networks. *Queueing Systems* **30** 89–140.
- [10] BUDHIRAJA, A. and LEE, C. (2009). Stationary distribution convergence for generalized Jackson networks in heavy traffic. *Mathematics of Operations Research* **34** 45–56.
- [11] CHATTERJEE, S. (2014). A short survey of Stein's method. To appear in Proceedings of ICM 2014.
- [12] CHEN, L. H. Y., GOLDSTEIN, L. and SHAO, Q.-M. (2011). *Normal approximation by Stein's method*. *Probability and its Applications (New York)*. Springer, Heidelberg. [MR2732624 \(2012b:60103\)](#)
- [13] DAI, J. G., DIEKER, A. B. and GAO, X. (2014). Validity of heavy-traffic steady-state approximations in many-server queues with abandonment. *Queueing Systems* **78** 1–29.
- [14] DAI, J. G. and HE, S. (2013). Many-server queues with customer abandonment: Numerical analysis of their diffusion model. *Stochastic Systems* **3** 96–146.
- [15] DAI, J. G., HE, S. and TEZCAN, T. (2010). Many-server diffusion limits for  $G/Ph/n + GI$  queues. *Annals of Applied Probability* **20** 1854–1890.
- [16] DAI, J. G. and TEZCAN, T. (2011). State space collapse in many-server diffusion limits of parallel server systems. *Mathematics of Operations Research* **36** 271–320.
- [17] DIEKER, A. B. and GAO, X. (2013). Positive recurrence of piecewise Ornstein–Uhlenbeck processes and common quadratic Lyapunov functions. *The Annals of Applied Probability* **23** 1291–1317.
- [18] ERYILMAZ, A. and SRIKANT, R. (2012). Asymptotically tight steady-state queue length bounds implied by drift conditions. *Queueing Systems* **72** 311–359.
- [19] ETHIER, S. N. and KURTZ, T. G. (1986). *Markov Processes: Characterization and Convergence*. Wiley, New York.

- [20] FOSCHINI, G. J. and SALZ, J. (1978). A basic dynamic routing problem and diffusion. *IEEE Transactions on Communications* **26** 320–327.
- [21] GAMARNIK, D. and STOLYAR, A. L. (2012). Multiclass multiserver queueing system in the Halfin-Whitt heavy traffic regime: asymptotics of the stationary distribution. *Queueing Systems* **71** 25–51.
- [22] GAMARNIK, D. and ZEEVI, A. (2006). Validity of heavy traffic steady-state approximation in generalized Jackson networks. *Ann. Appl. Probab.* **16** 56–90. [MR2209336](#)
- [23] GANS, N., KOOLE, G. and MANDELBAUM, A. (2003). Telephone call centers: Tutorial, review, and research prospects. *Manufacturing & Service Operations Management* **5** 79–141.
- [24] GIBBS, A. L. and SU, F. E. (2002). On Choosing and Bounding Probability Metrics. *International Statistical Review / Revue Internationale de Statistique* **70** pp. 419–435.
- [25] GILBARG, D. and TRUDINGER, N. S. (1983). *Elliptic Partial Differential Equations of Second Order*, 2nd ed. Springer, New York.
- [26] GLYNN, P. W. and ZEEVI, A. (2008). Bounding stationary expectations of Markov processes. In *Markov processes and related topics: a Festschrift for Thomas G. Kurtz. Inst. Math. Stat. Collect.* **4** 195–214. Inst. Math. Statist., Beachwood, OH. [MR2574232 \(2011b:60283\)](#)
- [27] GURVICH, I. (2014). Diffusion models and steady-state approximations for exponentially ergodic Markovian queues. *The Annals of Applied Probability* **24** 2527–2559.
- [28] GURVICH, I. (2014). Validity of Heavy-Traffic Steady-State Approximations in Multiclass Queueing Networks: The Case of Queue-Ratio Disciplines. *Mathematics of Operations Research* **39** 121–162.
- [29] GURVICH, I., HUANG, J. and MANDELBAUM, A. (2014). Excursion-Based Universal Approximations for the Erlang-A Queue in Steady-State. *Mathematics of Operations Research* **39** 325–373.
- [30] HALFIN, S. and WHITT, W. (1981). Heavy-traffic limits for queues with many exponential servers. *Oper. Res.* **29** 567–588. [MR629195 \(82i:90046\)](#)
- [31] HARRISON, J. M. (1978). The diffusion approximation for tandem queues in heavy traffic. *Advances in Applied Probability* **10** 886–905.
- [32] HARRISON, J. M. (1998). Heavy traffic analysis of a system with parallel servers: asymptotic analysis of discrete-review policies. *Annals of Applied Probability* **8** 822–848.
- [33] HARRISON, J. M. and LÓPEZ, M. J. (1999). Heavy traffic resource pooling in parallel-server systems. *Queueing Systems* **33** 339–368.
- [34] HARRISON, J. M. and REIMAN, M. I. (1981). Reflected Brownian motion on an orthant. *Annals of Probability* **9** 302–308.
- [35] HENDERSON, S. G. (1997). Variance reduction via an approximating Markov process PhD thesis, Department of Operations Research, Stanford University <http://people.orie.cornell.edu/shane/pubs/thesis.pdf>.
- [36] IGLEHART, D. L. and WHITT, W. (1970). Multiple channel queues in heavy traffic I. *Advances in Applied Probability* **2** 150–177.
- [37] IGLEHART, D. L. and WHITT, W. (1970). Multiple channel queues in heavy traffic II: sequences, networks, and batches. *Advances in Applied Probability* **2** 355–369.
- [38] KANG, W., KELLY, F., LEE, N. and WILLIAMS, R. (2009). State space collapse and diffusion approximation for a network operating under a fair bandwidth sharing policy. *The Annals of Applied Probability* **19** 1719–1780.
- [39] KATSUDA, T. (2010). State-space collapse in Stationarity and its application to a multiclass single-server queue in heavy traffic. *Queueing Systems: Theory and Applications* **65** 237–273.
- [40] KNOBLAUCH, A. (2008). Closed-form Expressions for the Moments of the Binomial Probability Distribution. *SIAM Journal on Applied Mathematics* **69** 197–8.
- [41] LUK, H. M. (1994). Stein’s method for the gamma distribution and related statistical applications PhD thesis, University of Southern California.
- [42] MEYN, S. P. and TWEEDIE, R. L. (1993). Stability of Markovian processes III: Foster-Lyapunov Criteria for Continuous Time Processes. *Adv. Appl. Probab.* **25** 518–548.
- [43] PETERSON, W. P. (1991). A heavy traffic limit theorem for networks of queues with multiple customer

- types. *Mathematics of Operations Research* **16** 90–118.
- [44] REED, J. (2009). The  $G/GI/N$  queue in the Halfin-Whitt regime. *Annals of Applied Probability* **19** 2211–2269.
- [45] REIMAN, M. I. (1984). Some diffusion approximations with state space collapse. In *Modeling and Performance Evaluation Methodology* (F. Baccelli and G. Fayolle, eds.) 209–240. Springer, Berlin.
- [46] REIMAN, M. I. (1984). Open queueing networks in heavy traffic. *Mathematics of Operations Research* **9** 441–458.
- [47] ROSS, N. (2011). Fundamentals of Stein's method. *Probab. Surv.* **8** 210–293. [MR2861132 \(2012k:60079\)](#)
- [48] SHI, P., CHOU, M., DAI, J. G., DING, D. and SIM, J. (2014). Models and Insights for Hospital Inpatient Operations: Time-Dependent ED Boarding Time. *Management Science*. forthcoming.
- [49] STEIN, C. (1972). A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 2: Probability Theory* 583–602. University of California Press, Berkeley, Calif.
- [50] STEIN, C. (1986). Approximate Computation of Expectations. *Lecture Notes-Monograph Series* **7** pp. i-iii+1-7+9-51+53-57+59-93+95-103+105-123+125-135+137-143+145-159+161-164.
- [51] TEZCAN, T. (2008). Optimal Control of Distributed Parallel Server Systems Under the Halfin and Whitt Regime. *Mathematics of Operations Research* **33** 51-90.
- [52] WHITT, W. (1971). Weak convergence theorems for priority queues: preemptive-resume discipline. *J. Appl. Probab.* **8** 74–94.
- [53] WHITT, W. (2002). *Stochastic-process limits*. Springer, New York. [MR1876437 \(2003f:60005\)](#)
- [54] WILLIAMS, R. J. (1998). Diffusion approximations for open multiclass queueing networks: sufficient conditions involving state space collapse. *Queueing Systems* **30** 27–88.
- [55] YE, H.-Q. and YAO, D. D. (2012). A Stochastic Network Under Proportional Fair Resource Control—Diffusion Limit with Multiple Bottlenecks. *Operations Research* **60** 716-738.
- [56] ZHANG, J. and ZWART, B. (2008). Steady state approximations of limited processor sharing queues in heavy traffic. *Queueing Systems: Theory and Applications* **60** 227–246. [MR2461617 \(2010a:60335\)](#)

SCHOOL OF OPERATIONS RESEARCH  
AND INFORMATION ENGINEERING  
CORNELL UNIVERSITY  
ITHACA, NY 14850  
USA  
E-MAIL: [ab2329@cornell.edu](mailto:ab2329@cornell.edu); [jd694@cornell.edu](mailto:jd694@cornell.edu)