

Optimal Scheduling of Critically Loaded Multiclass GI/M/n+M Queues in an Alternating Renewal Environment

Ari Arapostathis¹ · Guodong Pang² · Yi Zheng²

© Springer Science+Business Media, LLC, part of Springer Nature 2020

Abstract

In this paper, we study optimal control problems for multiclass GI/M/n + M queues in an alternating renewal (up–down) random environment in the Halfin–Whitt regime. Assuming that the downtimes are asymptotically negligible and only the service processes are affected, we show that the limits of the diffusion-scaled state processes under non-anticipative, preemptive, work-conserving scheduling policies, are controlled jump diffusions driven by a compound Poisson jump process. We establish the asymptotic optimality of the infinite-horizon discounted and long-run average (ergodic) problems for the queueing dynamics. Since the process counting the number of customers in each class is not Markov, the usual martingale arguments for convergence of mean empirical measures cannot be applied. We surmount this obstacle by demonstrating the convergence of the generators of an augmented Markovian model which incorporates the age processes of the renewal interarrival times and downtimes. We also establish long-run average moment bounds of the diffusion-scaled queueing processes under some (modified) priority scheduling policies. This is accomplished via Foster–Lyapunov equations for the augmented Markovian model.

Keywords Multiclass many-server queues · Halfin–Whitt (QED) regime · Service interruptions · Renewal arrivals · Alternating renewal process · Jump diffusions · Discounted cost · Ergodic control · Asymptotic optimality

☐ Guodong Pang gup3@psu.edu

Ari Arapostathis ari@utexas.edu

Yi Zheng yxz282@psu.edu

Published online: 08 July 2020

The Harold and Inge Marcus Department of Industrial and Manufacturing Engineering, College of Engineering, Pennsylvania State University, University Park, PA 16802, USA



Department of Electrical and Computer Engineering, The University of Texas at Austin, 2501 Speedway, EERC 7.824, Austin, TX 78712, USA

Mathematics Subject Classification Primary $90B22 \cdot 90B36 \cdot 60K37$; Secondary $60K25 \cdot 60J75 \cdot 60F17$

1 Introduction

There has been a lot of research activity on scheduling control problems for queueing networks in the Halfin–Whitt regime. The discounted problem for multiclass many-server queues was first studied in [1]. See also the work in [2,3]. For the ergodic control problem in the case of Markovian queueing networks see [4–6]. Scheduling control problems for queueing networks in random environments have also attracted much attention recently [7–10]. It is worth noting that in the study of asymptotic optimality in Markov-modulated environments, the scaling parameter depends on the rate of the underlying Markov process; see, for example, [7,10,11].

In this paper we consider queueing networks operating in alternating renewal (up-down) random environments, modeling service interruptions, and with renewal arrivals. It is well known that for large-scale service systems, service interruptions can have a dramatic impact on system performance [12]. For single class queues and networks in an alternating renewal environment, limit theorems have been studied in [12–16]. To the best of our knowledge, there are no studies on optimal scheduling control for multiclass many-server queues in alternating renewal environments, or even ergodic control in the Halfin–Whitt regime with arrivals that are renewal processes.

Specifically, we consider multiclass (d classes) GI/M/n + M queues with service interruptions in the Halfin–Whitt regime, where the arrival rate in each class and the number of servers in the pool are large, with a scaling parameter n, and the service interruptions are asymptotically negligible of order $n^{-1/2}$. The service interruption is modeled as an alternating renewal process constructed by regenerative 'up' and 'down' cycles. In the 'down' state, all servers stop functioning, and new customers arrive, which may abandon the queue. In the 'up' state, the queueing system functions normally. We assume that at least one class of customers has a strictly positive abandonment rate. The scheduling policy determines the allocation of servers to different classes of customers. We approximate the scheduling problem via the corresponding control problem of the limiting jump diffusion in the heavy-traffic regime, for which a sharp characterization of optimal Markov controls is available [17], and use this to exhibit matching upper and lower bounds on the optimal scheduling performance for the queueing dynamics.

In Proposition 3.1, we establish a functional central limit theorem (FCLT) for the *d*-dimensional diffusion-scaled state processes under work-conserving scheduling policies. The limiting controlled processes are jump diffusions with piecewise linear drift and compound Poisson jumps. The proof of weak convergence relies on the construction of a modified diffusion-scaled state process, where we add the cumulative downtime to a diffusion-scaled state process without interruptions. We show that the modified and original diffusion-scaled state processes have the same weak limits, which are governed by the jump diffusions described above.

The discounted and ergodic control problems for a large class of jump diffusions arising from queueing networks in the Halfin–Whitt regime have been studied in [17],



and these results are essential for establishing asymptotic optimality in the present paper. In Theorem 3.1, we show that the optimal value functions of the discounted problem for the diffusion-scaled processes converge to the corresponding function for the limiting jump diffusion. The proof of asymptotic optimality for the discounted problem follows the approach in [1], which deals with the discounted problem for multiclass GI/M/n + M queues. An essential part of this proof involves moment bounds for the diffusion-scaled state process, and the cumulative downtime process.

Asymptotic optimality for the ergodic control problem is more challenging. The result is stated in Theorem 3.2. Here, long-run average moment bounds for the diffusion-scaled state processes play a crucial role (see Proposition 4.2). Typically, such bounds are obtained in the literature via Foster-Lyapunov inequalities [4-6,10,18]. However, since the process counting the number of customers in each class, referred to as the queueing process, or state process, is not Markov, we first construct a sequence of auxiliary diffusion-scaled processes by adding the scaled residual time process of the alternating renewal process in the 'down' state to the original process, taking advantage of the fact that the long-run average moments of the scaled residual time process are negligible as the scaling parameter n tends to infinity (see equation (4.25)). We then consider the joint Markov process comprised of the auxiliary diffusion-scaled state process and the age processes of renewal arrival and alternating renewal processes, and construct Foster-Lyapunov functions, which bear a resemblance to the Lyapunov functions in [19]. In this part, we assume that the mean residual life functions are bounded, and use the criterion in [20, Theorem 4.2] to show that the joint Markov processes are positive Harris recurrent for all large enough n under some (modified) priority scheduling policy. We apply a two-step scheduling: first, the servers are allocated to the classes of customers with zero abandonment rate in such a manner that the servers used for each class do not exceed a certain proportion dictated by the traffic intensity; second, a static priority rule is applied to allocate the remaining servers. We show that the long-run average moments of the auxiliary diffusion-scaled state processes are bounded under this scheduling policy. We then establish a moment estimate for the difference between the auxiliary and original diffusion-scaled processes, and proceed to show that the analogous moment bounds hold for the original diffusion-scaled processes.

To prove asymptotic optimality for the ergodic control problem, we establish lower and upper bounds for the limits of the value functions (see Eqs. (5.10) and (5.28)). For the proof of the lower bound, we show that the sequence of mean empirical measures of the diffusion-scaled state processes is tight (see Lemma 5.2), and any limit of mean empirical measures is an ergodic occupation measure for the limiting jump diffusion. This is analogous to the technique used in [4–6,10]. However, characterizing the limits of mean empirical measures (see Theorem 5.2) is quite challenging here. Since we consider the diffusion-scaled processes with renewal arrivals in an alternating renewal environment, the martingale arguments in the above papers cannot be applied here. Instead, we develop a new approach. Following the technique of the proof of ergodicity under the specific scheduling policy described in the preceding paragraph, we consider the generator of the joint Markov process of the auxiliary diffusion-scaled state process, which incorporates the residual time process, and the associated age processes of the renewal arrivals and the alternating renewal environment. We construct suitable test



functions (see (5.12)) which involve the coefficients of variation of interarrival times, and proceed to show the convergence of generators.

For the proof of the upper bound, we adopt the spatial truncation technique developed in [4], which is also used in [5,6,10], and is extended to jump diffusions in [17]. This involves a concatenated scheduling policy. We first construct a continuous precise ϵ -optimal control for the ergodic control problem for the limiting jump diffusion (see Proposition 5.1). Then, inside a compact set, we map this control to a scheduling policy for the diffusion-scaled process. On the complement of this set, we apply the (modified) priority scheduling policy. We show that the long run average moments of the diffusion-scaled state process are bounded under this concatenated scheduling policy (see Proposition 4.3), and the limit of mean empirical measures is the ergodic occupation measure of the limiting jump diffusion governed by the ϵ -optimal control (see Lemma 5.3). Here, the techniques used in establishing the long-run average moment bounds under the (modified) priority scheduling policy, and the convergence of mean empirical measures, play an important role.

1.1 Organization of the Paper

The notation used in the paper is summarized in the next subsection. In Sect. 2, we describe the model of multiclass many-server queues with service interruptions. In Sect. 3, we define the diffusion-scaled processes and associated control problems, and state the main results on weak convergence and asymptotic optimality. In Sect. 4, we summarize the ergodic properties of the limiting controlled jump diffusion, and state the results concerning long-run average moment bounds for the diffusion-scaled processes. The proofs of Theorems 3.1 and 3.2 are given in Sect. 5. Appendix A is devoted to the proofs of Lemma 3.1 and Proposition 3.1. Appendix B contains the proofs of Lemmas 4.1 and 5.2.

1.2 Notation

We let $|\cdot|$ and $\langle\cdot,\cdot\rangle$ denote the standard Euclidean norm and the inner product in \mathbb{R}^d , respectively. For $x\in\mathbb{R}^d$, we let $\|x\|:=\sum_i|x_i|$, and x' denote the transpose of x. The symbols \mathbb{R}_+ , \mathbb{Z}_+ , \mathbb{N} , denote the set of nonnegative real numbers, nonnegative integers, and the set of natural numbers, respectively. The indicator function of a set $A\in\mathbb{R}^d$ is denoted by $\mathbb{1}_A$. Given $a,b\in\mathbb{R}$, the minimum (maximum) is denoted by $a\wedge b$ ($a\vee b$), respectively, $\lfloor a\rfloor$ denotes the integer part of a, and $a^\pm:=(\pm a)\vee 0$. The complement and closure of a set $A\subset\mathbb{R}^d$ are denoted by A^c and A, respectively. We use the notation e_i to denote the vector with ith entry equal to 1 and all other entries equal to 0. We also let $e:=(1,\ldots,1)^{\mathsf{T}}$. We let B_r denote the open ball of radius r in \mathbb{R}^d , centered at the origin. For a process $\{X_t\}_{t\geq 0}$, $\tau(A)$ denotes the first exit time from the set $A\subset\mathbb{R}^d$, defined by $\tau(A):=\inf\{t>0:X_t\notin A\}$, and we let $\tau_r:=\tau(B_r)$.

For a domain $D \subset \mathbb{R}^d$, the space $\mathcal{C}^k(D)$ ($\mathcal{C}^{\infty}(D)$), $k \geq 0$, stands for the class of all real-valued functions on D whose partial derivatives up to order k (of any order) exist and are continuous. $\mathcal{C}^{k,r}(D)$ stands for the set of functions that are k-times continuously differentiable and whose kth derivatives are locally Hölder continuous with exponent



r. We let $\mathcal{C}^k_c(D)$ denote the space of functions in $\mathcal{C}^k(D)$ with compact support, and \mathcal{C}^k_b the set of functions in $\mathcal{C}^k(D)$ whose partial derivatives up to order k are bounded. For a nonnegative function $g \in \mathcal{C}(\mathbb{R}^d)$, $\mathcal{O}(g)$ denotes the space of functions $f \in \mathcal{C}(\mathbb{R}^d)$ satisfying $\sup_{x \in \mathbb{R}^d} \frac{|f(x)|}{1+g(x)} < \infty$. By a slight abuse of notation, $\mathcal{O}(g)$ also denotes a generic member of these spaces.

For $k \in \mathbb{N}$, we let $\mathbb{D}^{\hat{k}} := \mathbb{D}(\mathbb{R}_+, \mathbb{R}^k)$ denote the space of \mathbb{R}^k -valued cádlág functions on \mathbb{R}_+ . When k = 1, we write \mathbb{D} for \mathbb{D}^k . Given a Polish space E, by $\mathcal{P}(E)$ we denote the space of probability measures on E, endowed with the Prokhorov metric.

2 Multiclass GI/M/N + M Queues with Service Interruptions

2.1 The Model and Assumptions

We consider a sequence of GI/M/n+M queueing models with d classes of customers. Let $\mathfrak{I}:=\{1,\ldots,d\}$. For the nth system, let $\{A_i^n(t)\}_{t\geq 0}$ denote the arrival process of class-i customers. We assume that the arrivals are mutually independent renewal processes defined as follows. Let $\{G_{i,j}\colon j\in\mathbb{N}\}, i\in \mathfrak{I}$, be an i.i.d. sequence of strictly positive random variables with mean $\mathbb{E}[G_i]=1$ and finite (squared) coefficient of variation $c_{a,i}^2:=\mathrm{Var}(G_i)/(\mathbb{E}[G_i])^2$, where $G_i\equiv G_{i,1}$. Then, we define

$$A_i^n(t) := \max \left\{ m \ge 0 : \sum_{i=1}^m G_{i,j} \le \lambda_i^n t \right\}, \quad t \ge 0, \ i \in \mathcal{I},$$
 (2.1)

where $\lambda_i^n > 0$ denotes the arrival rate. For each $n \in \mathbb{N}$, the service and patience times of the class-*i* customers are exponentially distributed with parameters μ_i^n and γ_i^n , respectively.

We adopt the following standard assumption on the parameters (see [1,4,13]).

Assumption 2.1 (*The Halfin–Whitt regime*) The parameters satisfy the following limits for each $i \in \mathcal{I}$ as $n \to \infty$:

$$n^{-1}\lambda_{i}^{n} \rightarrow \lambda_{i} > 0, \quad \mu_{i}^{n} \rightarrow \mu_{i} > 0, \quad \gamma_{i}^{n} \rightarrow \gamma_{i} \geq 0,$$

$$n^{-1/2}(\lambda_{i}^{n} - n\lambda_{i}) \rightarrow \hat{\lambda}_{i}, \quad n^{1/2}(\mu_{i}^{n} - \mu_{i}) \rightarrow \hat{\mu}_{i},$$

$$\frac{\lambda_{i}^{n}}{n\mu_{i}^{n}} \rightarrow \rho_{i} := \frac{\lambda_{i}}{\mu_{i}} < 1, \quad \sum_{i=1}^{d} \rho_{i} = 1.$$

We assume that $\inf_{n\in\mathbb{N}} \gamma_d^n > 0$. Assumption 2.1, which is also known as the Quality-and-Efficiency-Driven regime, implies that the system is critically loaded and

$$\rho^n \to \hat{\rho} := \sum_{i=1}^d \frac{\rho_i \hat{\mu}_i - \hat{\lambda}_i}{\mu_i} \in \mathbb{R}, \text{ where } \rho^n := \sqrt{n} \left(1 - \sum_{i=1}^d \frac{\lambda_i^n}{n \mu_i^n} \right).$$



All queues are in the same up-down alternating renewal random environment. Waiting customers may abandon at any time. In the 'up' state, the system functions normally, and in the 'down' state all servers stop, while customers keep joining the queues and any jobs that have started service will wait for the system to resume. For this reason, we also refer to this model as multiclass queues with service interruptions. Let $\{(u_k^n, d_k^n): k \in \mathbb{N}\}$ be a sequence of i.i.d. positive random vectors denoting the up-down cycles, and define the *counting process of downtimes* by

$$N^{n}(t) := \max\{k \ge 0 : T_{k}^{n} \le t\}, \text{ with } T_{k}^{n} := \sum_{i=1}^{k} (u_{i}^{n} + d_{i}^{n}), k \in \mathbb{N},$$
 (2.2)

and $T_0^n \equiv 0$. At time 0, the system is in the 'up' state.

Assumption 2.2 For each n and k in \mathbb{N} , u_k^n and d_k^n are independent, u_k^n is exponentially distributed with parameter β_u^n , which converges to $\beta > 0$ as $n \to \infty$. We assume that $d_1^n = \frac{1}{\vartheta^n} d_1$, with d_1 some nonnegative random variable satisfying $\mathbb{E}[d_1] = 1$, and $\frac{\vartheta^n}{\sqrt{n}} \to \vartheta > 0$ as $n \to \infty$.

For $k \in \mathbb{N}$, we let (\mathbb{D}^k, M_1) and (\mathbb{D}^k, J_1) denote the space \mathbb{D}^k endowed with the Skorokhod M_1 and J_1 topologies, respectively (see, for example, [21,22]). Assumption 2.2 implies that the service interruptions are asymptotically negligible, and

$$N^n \Rightarrow N$$
 in (\mathbb{D}, J_1) as $n \to \infty$,

where the limiting process N is a Poisson process with rate β . Define the *server* availability process $\Psi^n := \{\Psi^n(t) : t \ge 0\}$ by

$$\Psi^{n}(t) = \begin{cases} 1, & T_{k}^{n} \le t < T_{k}^{n} + u_{k+1}^{n}, \\ 0, & T_{k}^{n} + u_{k+1}^{n} \le t < T_{k+1}^{n}, \end{cases}$$
(2.3)

for $k \in \mathbb{N}$. We also define the cumulative up-time process $C_{\mathsf{u}}^n = \{C_{\mathsf{u}}^n(t)\}_{t\geq 0}$ by $C_{\mathsf{u}}^n(t) := \int_0^t \Psi^n(s) \, \mathrm{d}s$, and the cumulative down-time process by $C_{\mathsf{d}}^n(t) := t - C_{\mathsf{u}}^n(t)$. Let F^{d_1} denote the distribution function of d_1 . By [13, Lemma 2.2], we have

$$\sqrt{n}C_d^n \Rightarrow L \text{ in } (\mathbb{D}, M_1) \text{ as } n \to \infty,$$
 (2.4)

where $\{L_t\}_{t\geq 0}$ is a compound Poisson process with intensity $\Pi_L(dx)dt = \beta F^{d_1}(\vartheta dx)dt$, where β is given in Assumption 2.2.

For the *n*th system, we denote the processes counting the total number of customers, those in queue, and those in service, by $X^n = (X_1^n, \dots, X_d^n)'$, $Q^n = (Q_1^n, \dots, Q_d^n)'$, and $Z^n = (Z_1^n, \dots, Z_d^n)'$, respectively. These processes satisfy the following constraints:

$$X_i^n(t) = Q_i^n(t) + Z_i^n(t), \quad Q_i^n(t) \ge 0, \quad Z_i^n(t) \ge 0,$$



and
$$\langle e, Z^n(t) \rangle \le n$$
 (2.5)

for each $t \ge 0$ and $i \in \mathcal{I}$. We let

$$S_{i}^{n}(t,r) := S_{*,i}^{n} \left(\mu_{i}^{n} \int_{0}^{t} Z_{i}^{n}(s) \Psi^{n}(s) \, \mathrm{d}s + \mu_{i}^{n} r \right),$$

$$R_{i}^{n}(t,r) := R_{*,i}^{n} \left(\gamma_{i}^{n} \int_{0}^{t} Q_{i}^{n}(s) \, \mathrm{d}s + \gamma_{i}^{n} r \right),$$
(2.6)

for $i \in \mathcal{I}$, $t \geq 0$, and $r \geq 0$, where $\{S_{*,i}^n, R_{*,i}^n : i \in \mathcal{I}, n \in \mathbb{N}\}$ are Poisson processes with rate one. We assume that for each $n \in \mathbb{N}$, $\{X_i^n(0), A_i^n, S_{*,i}^n, R_{*,i}^n : i \in \mathcal{I}\}$ are mutually independent. These processes are governed by the equation

$$X_i^n(t) = X_i^n(0) + A_i^n(t) - S_i^n(t) - R_i^n(t)$$
 (2.7)

for each $t \ge 0$, $n \in \mathbb{N}$, and $i \in \mathcal{I}$, where $S_i^n(t) := S_i^n(t, 0)$ and $R_i^n(t) := R_i^n(t, 0)$.

2.2 Scheduling Policies

A scheduling policy is identified with a \mathbb{Z}_+^d -valued stochastic process \mathbb{Z}^n with cádlág sample paths, which satisfies (2.5). Let

$$\tilde{\tau}_i^n(t) := \inf\{r \ge t : A_i^n(r) - A_i^n(r-) > 0\},$$

and $\tilde{\tau}^n(t) := \inf\{r > t : \Psi^n(r) = 1\},$ (2.8)

for $i \in \mathcal{I}$. Recall the definitions of C_d^n in (2.4), and S^n and R^n in (2.6). Define the σ -fields

$$\mathcal{F}_{t}^{n} := \sigma \left\{ X^{n}(0), A_{i}^{n}(t), S_{i}^{n}(s), R_{i}^{n}(s), X_{i}^{n}(s), Z_{i}^{n}(s), Z_{i}^{n}(s), \Psi^{n}(s), N^{n}(s) : i \in \mathcal{I}, 0 \leq s \leq t \right\} \vee \mathcal{N}, \\
\mathcal{G}_{t}^{n} := \sigma \left\{ A_{i}^{n}(\tilde{\tau}_{i}^{n}(t) + r) - A_{i}^{n}(\tilde{\tau}_{i}^{n}(t)), S_{i}^{n}(\check{\tau}^{n}(t), r) - S_{i}^{n}(\check{\tau}^{n}(t)), \times R_{i}^{n}(\check{\tau}^{n}(t), r) - R_{i}^{n}(\check{\tau}^{n}(t)), C_{d}^{n}(\check{\tau}^{n}(t) + r) - C_{d}^{n}(\check{\tau}^{n}(t)) : i \in \mathcal{I}, r \geq 0 \right\} \vee \mathcal{N}, \tag{2.9}$$

for $t \geq 0$, where \mathcal{N} is the collection of all \mathbb{P} -null sets. We say that a scheduling policy Z^n is non-anticipative if

- (i) $Z^n(t)$ is adapted to \mathcal{F}_t^n ,
- (ii) \mathcal{F}_t^n and \mathcal{G}_t^n are independent at each time $t \geq 0$,
- (iii) for each $i \in \mathcal{I}$, and $t \geq 0$, the process $S_i^n(\check{\tau}^n(t), \cdot) S_i^n(\check{\tau}^n(t))$ agrees in law with $S_{*,i}^n(\mu_i^n)$, and the process $R_i^n(\check{\tau}^n(t), \cdot) R_i^n(\check{\tau}^n(t))$ agrees in law with $R_{*,i}^n(\gamma_i^n)$.

The information at time t is contained in \mathcal{F}^n_t , while \mathcal{G}^n_t represents the information about future increments. The renewal arrivals A^n_i , $i \in \mathcal{I}$, and the alternative renewal process Ψ^n are regenerative processes. So in \mathcal{G}^n_t , we use $\tilde{\tau}^n_i(t)$ and $\tilde{\tau}^n(t)$, respectively, instead of t. Note that parts (ii) and (iii) in the definition of non-anticipative scheduling policy are required so that the any limit of scheduling policies corresponds to a non-anticipative



control for the limiting controlled jump diffusion. See part (iii) of Proposition 3.1 for details.

Let $\tau^n_{i,k}$ denote the kth jump time of $A^n_i - S^n_i - R^n_i$, for each $n \in \mathbb{N}$ and $i \in \mathcal{I}$. Equation (2.7) implies that $X^n_i(t) = X^n_i(0)$ for $0 \le t \le \tau^n_{i,1}$, $X^n_i(t) = X^n_i(0) + \epsilon_1$ for $\tau^n_{i,1} \le t \le \tau^n_{i,2}$ and so forth, where ϵ_k denotes the jump size which takes values in a bounded set. Note that the integrals in (2.6) are finite by the definition of Ψ^n in (2.3) and (2.5). Thus, given any non-anticipative scheduling policy Z^n , and initial condition $X^n(0)$, there exists a unique solution to (2.7).

For $x \in \mathbb{Z}_+^d$, we define the action set $\mathbb{Z}^n(x)$ by

$$\mathcal{Z}^n(x) := \left\{ z \in \mathbb{Z}_+^d : z \le x, \ \langle e, z \rangle = \langle e, x \rangle \land n \right\}.$$

A scheduling policy Z^n is called *admissible* if $Z^n(t)$ takes values in $Z^n(X^n(t))$ at each t, and is non-anticipative. The set of admissible scheduling policies is denoted by \mathfrak{Z}^n . Note that an admissible policy allows preemption, that is, a server can interrupt service of a customer at any time to serve some other class of customers. In summary, given an admissible scheduling policy Z^n , the process X^n in (2.7) is well defined, and we say that X^n is governed by Z^n .

Next, we describe a well-known equivalent parameterization of the set of admissible policies. Let

$$\mathcal{S} := \{ u \in \mathbb{R}^{d}_{+} : \langle e, u \rangle = 1 \}.$$

We also define

$$\mathcal{S}^n(x) := \left\{ v \in \mathbb{Z}_+^d : v = \frac{y}{\langle e, x \rangle - n} \in \mathcal{S}, \ y \le x, \ y \in \mathbb{Z}_+^d \right\}, \quad \text{if } \langle e, x \rangle > n,$$

and $S^n(x) = \{e_d\}$, if $\langle e, x \rangle \leq n$. Let \mathfrak{U}^n denote the class of processes $\{U^n(t)\}_{t\geq 0}$ which are non-anticipative, in the sense of the definition given above, and $U^n(t)$ takes values in $S^n(X^n(t))$. Then, each $U^n \in \mathfrak{U}^n$ determines a policy $Z^n \in \mathfrak{Z}^n$ via

$$Z^{n}(t) = X^{n}(t) - Q^{n}(t)$$
, with $Q^{n}(t) = (\langle e, X^{n}(t) \rangle - n)^{+} U^{n}(t)$.

This map is invertible, and its inverse is given by

$$U^{n}(t) := \begin{cases} \frac{X^{n}(t) - Z^{n}(t)}{\langle e, X^{n}(t) \rangle - n} & \text{for } \langle e, X^{n}(t) \rangle > n, \\ e_{d} & \text{for } \langle e, X^{n}(t) \rangle \leq n. \end{cases}$$

Therefore, as far as control problems are concerned, we can use policies in \mathfrak{U}^n or \mathfrak{Z}^n interchangeably. Note that U_i^n can be considered as the proportion of class-i customers in the queue when there are waiting customers in the system.

Next, we augment the state space, and define the class of stationary Markov scheduling policies. Recall the definitions of A^n , N^n , and Ψ^n in (2.1)–(2.3), respectively.



Definition 2.1 Let $H_i^n(t)$ denote the age process for the class-i customers, that is,

$$H_i^n(t) := t - \frac{1}{\lambda_i^n} \sum_{j=1}^{A_i^n(t)} G_{i,j}, \quad t \ge 0, \quad i \in \mathcal{I},$$
 (2.10)

and define the age process K^n for the alternating renewal process in the 'down' state by

$$K^{n}(t) := \left(t - \sum_{k=1}^{N^{n}(t)} (u_{k}^{n} + d_{k}^{n}) - u_{N^{n}(t)+1}^{n}\right)^{+}, \quad t \ge 0.$$
 (2.11)

Then, (A_i^n, H_i^n) , $i \in \mathcal{I}$, and (Ψ^n, K^n) are strong Markov processes (see, e.g., [23]). We say that a scheduling policy $Z^n \in \mathfrak{Z}^n$ is (stationary) Markov if

$$Z^{n}(t) = z^{n}(X^{n}(t), H^{n}(t), \Psi^{n}(t), K^{n}(t))$$

for some $z^n : \mathbb{Z}_+^d \times \mathbb{R}_+^d \times \{0,1\} \times \mathbb{R}_+ \to \mathbb{Z}_+^d$, and we let $\mathfrak{Z}_{\mathrm{sm}}^n$ denote the class of these policies. Under a policy $Z^n \in \mathfrak{Z}_{\mathrm{sm}}^n$, the process (X^n, H^n, Ψ^n, K^n) is Markov with state space

$$\left\{ (x, h, \psi, k) \in \mathbb{Z}_+^d \times \mathbb{R}_+^d \times \{0, 1\} \times \mathbb{R}_+ \colon k \equiv 0 \text{ if } \psi = 1 \right\}.$$

Abusing the notation, when z^n depends only on its first argument, we simply write $Z^n(t) = z^n(X^n(t))$.

3 Diffusion-Scaled Processes and Control Problems

Let \hat{X}^n , \hat{Q}^n , and \hat{Z}^n denote the diffusion-scaled processes defined by

$$\hat{X}_i^n(t) := n^{-1/2} (X_i^n(t) - \rho_i n), \quad \hat{Q}_i^n(t) := n^{-1/2} Q_i^n(t), \quad \hat{Z}_i^n(t) := n^{-1/2} (Z_i^n(t) - \rho_i n),$$

respectively, for $t \ge 0$ and $i \in \mathcal{I}$. It follows by (2.7) that the process \hat{X}_i^n takes the form

$$\hat{X}_{i}^{n}(t) = \hat{X}_{i}^{n}(0) + \ell_{i}^{n}t + \hat{A}_{i}^{n}(t) - \hat{S}_{i}^{n}(t) - \hat{R}_{i}^{n}(t)
- \mu_{i}^{n} \int_{0}^{t} \hat{Z}_{i}^{n}(s) \Psi^{n}(s) \, \mathrm{d}s - \gamma_{i}^{n} \int_{0}^{t} \hat{Q}_{i}^{n}(s) \, \mathrm{d}s + \hat{L}_{i}^{n}(t), \quad t \ge 0,$$
(3.1)

where $\ell_i^n := n^{-1/2} (\lambda_i^n - n \mu_i^n \rho_i),$

$$\begin{split} \hat{A}_i^n(t) &:= n^{-1/2} \Big(A_i^n(t) - \lambda_i^n t \Big), \qquad \hat{S}_i^n(t) \, := \, n^{-1/2} \bigg(S_i^n(t) - \mu_i^n \int_0^t Z_i^n(s) \Psi^n(s) \, \mathrm{d}s \bigg), \\ \hat{R}_i^n(t) &:= n^{-1/2} \bigg(R_i^n(t) - \gamma_i^n \int_0^t Q_i^n(s) \, \mathrm{d}s \bigg), \quad \text{and } \hat{L}_i^n(t) \, := \, \sqrt{n} \mu_i^n \rho_i C_\mathsf{d}^n(t). \end{split}$$



Let \hat{W}^n and \hat{Y}^n , $n \in \mathbb{N}$, be d-dimensional processes defined by

$$\hat{W}_i^n := \hat{A}_i^n - \hat{S}_i^n - \hat{R}_i^n \quad \text{for } i \in \mathcal{I}, \tag{3.2}$$

and

$$\hat{Y}_{i}^{n}(t) := \ell_{i}^{n} t - \mu_{i}^{n} \int_{0}^{t} \hat{Z}_{i}^{n}(s) \Psi^{n}(s) \, \mathrm{d}s - \gamma_{i}^{n} \int_{0}^{t} \hat{Q}_{i}^{n}(s) \, \mathrm{d}s \quad \text{ for } i \in \mathcal{I}, \ t \geq 0,$$

respectively. Then, \hat{X}_{i}^{n} in (3.1) has the representation

$$\hat{X}_{i}^{n}(t) = \hat{X}_{i}^{n}(0) + \hat{Y}_{i}^{n}(t) + \hat{W}_{i}^{n}(t) + \hat{L}_{i}^{n}(t).$$

The initial condition $\hat{X}^n(0)$, $n \in \mathbb{N}$, is assumed to be deterministic throughout the paper.

3.1 The Limiting Controlled Diffusion with Compound Poisson Jumps

In Lemma 3.1 and Proposition 3.1 which follow, products or powers of the spaces (\mathbb{D}^d, J_1) and (\mathbb{D}^d, M_1) are viewed as metric spaces endowed with the maximum metric. The proofs of these results are given in Appendix A.

Lemma 3.1 Suppose that Assumptions 2.1 and 2.2 hold, and that $\{\hat{X}^n(0): n \in \mathbb{N}\}$ is bounded. Then, under any sequence of $U^n \in \mathcal{U}^n$, we have

$$(n^{-1}Q^n, n^{-1}Z^n) \Rightarrow (\mathfrak{e}_0, \mathfrak{e}_\rho) \text{ in } (\mathbb{D}^d, M_1)^2,$$

where $\mathfrak{e}_0(t) \equiv (0, \dots, 0)'$ for all $t \geq 0$, and $\mathfrak{e}_{\rho}(t) \equiv (\rho_1, \dots, \rho_d)'$.

Proposition 3.1 *Grant the assumptions in Lemma* 3.1. *Then, the following hold.*

(i) As $n \to \infty$,

$$(\hat{W}^n, \hat{L}^n) \Rightarrow (\Sigma W, \lambda L) \text{ in } (\mathbb{D}^d, J_1) \times (\mathbb{D}^d, M_1),$$

where the matrix Σ is given by $\Sigma := \operatorname{diag} \left(\sqrt{\lambda_1 (1 + c_{a,1}^2)}, \ldots, \sqrt{\lambda_d (1 + c_{a,d}^2)} \right)$, W is a d-dimensional standard Wiener process, $\lambda := (\lambda_1, \ldots, \lambda_d)'$, and $\{L_t\}_{t \geq 0}$ is the one-dimensional Lévy process in (2.4), and is independent of W.

- (ii) The sequence $(\hat{X}^n, \hat{Y}^n, \hat{W}^n, \hat{L}^n)$ is tight in $(\mathbb{D}^d, M_1) \times (\mathbb{D}^d, J_1)^2 \times (\mathbb{D}^d, M_1)$.
- (iii) Provided U^n is tight in (\mathbb{D}^d, J_1) , any limit X of \hat{X}^n is a strong solution to the stochastic differential equation

$$dX_t = b(X_t, U_t) dt + \sum dW_t + \lambda dL_t, \qquad (3.3)$$

with initial condition $X_0 = x \in \mathbb{R}^d$, where U is a limit of U^n , and $b(x, u) : \mathbb{R}^d \times \mathcal{S} \to \mathbb{R}^d$ takes the form

$$b(x, u) = \ell - M(x - \langle e, x \rangle^{+} u) - \langle e, x \rangle^{+} \Gamma u, \tag{3.4}$$

with $\ell := (\ell_1, \dots, \ell_d)'$, $M := \operatorname{diag}(\mu_1, \dots, \mu_d)$, and $\Gamma := \operatorname{diag}(\gamma_1, \dots, \gamma_d)$. Moreover, any such limit U is non-anticipative, that is, for s < t, $(W_t - W_s, L_t - L_s)$ is independent of

$$\mathcal{F}_s := \text{the completion of } \sigma\{X_0, U_r, W_r, L_r : r < s\}.$$

Throughout the paper, the time variable appears as a subscript in the processes governing the limiting controlled jump diffusion in order to distinguish them from the processes associated with the *n*th system.

3.2 The Control Problems

Define
$$\widetilde{\mathbb{R}} \colon \mathbb{R}^d_+ \to \mathbb{R}_+$$
 by
$$\widetilde{\widetilde{\mathbb{R}}}(x) := c|x|^m \tag{3.5}$$

for some c > 0 and $m \ge 1$. The running cost function $\Re: \mathbb{R}^d \times \mathcal{S} \to \mathbb{R}_+$ is defined by

$$\Re(x,u) := \widetilde{\Re}(\langle e, x \rangle^+ u).$$

Remark 3.1 We only choose a running cost function as in (3.5) to simplify the exposition. One may replace (3.5) with a function $\widetilde{\mathcal{R}}$, which is locally Lipschitz continuous, and satisfies

$$c_1|x|^m \le \widetilde{\Re}(x) \le c_2|x|^m \quad \forall x \in \mathbb{R}^d,$$
 (3.6)

for some positive constants c_1 , c_2 , and $m \ge 1$. All the results still hold with (3.6). Moreover, the lower bound in (3.6) is not needed for the discounted problem (see, e.g., [1]).

The α -discounted control problem for the *n*th system is given by

$$\hat{V}_{\alpha}^{n}(\hat{X}^{n}(0)) := \inf_{U^{n} \in \mathbb{N}^{n}} \hat{J}_{\alpha}(\hat{X}^{n}(0), U^{n}) \qquad \alpha > 0, \ n \in \mathbb{N},$$

where the cost criterion is defined by

$$\hat{J}_{\alpha}(\hat{X}^n(0), U^n) := \mathbb{E}\left[\int_0^{\infty} e^{-\alpha t} \,\Re(\hat{X}^n(s), U^n(s)) \,\mathrm{d}s\right] \quad \forall \alpha > 0.$$

For the controlled (jump) diffusion X in (3.3), we say that a control U is admissible if it takes values in S, and non-anticipative (see [17]). We denote the set of all admissible controls by $\mathfrak U$. The corresponding α -discounted cost criterion for the diffusion takes the form

$$J_{\alpha}(x, U) := \mathbb{E}_{x}^{U} \left[\int_{0}^{\infty} e^{-\alpha t} \, \Re(X_{s}, U_{s}) \, \mathrm{d}s \right] \quad \forall \, \alpha > 0,$$

and the optimal α -discounted value function is given by

$$V_{\alpha}(x) := \inf_{U \in \Omega} J_{\alpha}(x, U) \quad \forall \alpha > 0, \tag{3.7}$$

where \mathbb{E}_{x}^{U} denotes the expectation operator corresponding to the process under the control U, with initial condition $x \in \mathbb{R}^{d}$. We introduce the following assumption for the discounted problem.

Assumption 3.1 There exists a constant $m_A \ge m \lor 2$ with m as in (3.5) such that $\mathbb{E}[(G_i)^{m_A}] < \infty$, for all $i \in \mathcal{I}$, and $\mathbb{E}[(d_1)^{m_A \lor (m+1)}] < \infty$.

We state the main result for the discounted problem in the next theorem, whose proof is given in Sect. 5.2.

Theorem 3.1 *Grant the hypotheses in Assumptions* 2.1, 2.2, and 3.1, and suppose that $\hat{X}^n(0) \to x \in \mathbb{R}^d$ as $n \to \infty$. Then

$$\lim_{n \to \infty} \hat{V}_{\alpha}^{n} \left(\hat{X}^{n}(0) \right) = V_{\alpha}(x). \tag{3.8}$$

Remark 3.2 Note that in Theorem 3.1, we do not need to impose any restrictions on the limiting abandonment rates $\{\gamma_i : i \in \mathcal{I}\}$.

We define the ergodic control problem for the diffusion-scaled process by

$$\varrho^n(\hat{X}^n(0)) := \inf_{Z^n \in \mathfrak{Z}^n_{sm}} \hat{J}(\hat{X}^n(0), Z^n),$$

where the cost criterion \hat{J} is given by

$$\hat{J}(\hat{X}^n(0), Z^n) := \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}^{Z^n} \left[\int_0^T \widetilde{\mathcal{R}}(\hat{Q}^n(s)) \, \mathrm{d}s \right].$$

Here, the infimum is over all Markov scheduling policies, since for the ergodic control problem, we work with Markov processes. For the controlled jump diffusion in (3.3), the ergodic cost criterion, and the optimal ergodic value are defined by

$$J(x, U) := \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}_x^U \left[\int_0^T \Re(X_s, U_s) \, \mathrm{d}s \right],$$

and

$$\varrho_*(x) := \inf_{U \in \Omega} J(x, U),$$
(3.9)

respectively. By [17, Theorem 4.1], it follows that ϱ_* is independent of x, and optimality is attained by a stationary Markov control.

We introduce the following assumption on G_i and d_1 for the ergodic control problem.



Assumption 3.2 The following hold.

- (i) The right derivative of $F_i(t)$ is finite, and $F_i(t) < 1$, for all $t \ge 0$ and $i \in \mathcal{I}$. The distribution function F^{d_1} of d_1 satisfies the same property.
- (ii) The mean residual life functions of G_i and d_1 are bounded, that is, there exists some positive constant \widehat{C} such that

$$\frac{\int_{t}^{\infty} \left(1 - F^{d_1}(y)\right) dy}{1 - F^{d_1}(t)} \le \widehat{C}, \quad \text{and} \quad \frac{\int_{t}^{\infty} \left(1 - F_i(y)\right) dy}{1 - F_i(t)} \le \widehat{C} \quad \forall i \in \mathcal{I}, \quad (3.10)$$

and for all t > 0.

Assumption 3.2 implies that all absolute moments of G_i , $i \in \mathcal{I}$, and d_1 are finite. The main result of the ergodic control problem is stated in the next theorem, whose proof is given in Sect. 5.3.

Theorem 3.2 Grant Assumptions 2.1, 2.2, and 3.2. In addition, suppose that m in (3.5) is larger than 1, and that $\hat{X}^n(0) \to x \in \mathbb{R}^d$ as $n \to \infty$. Then, we have

$$\lim_{n \to \infty} \varrho^n (\hat{X}^n(0)) = \varrho_*.$$

4 Ergodic Properties

In this section, we present some ergodicity results for the limiting jump diffusion and the diffusion-scaled processes. These results are used to prove Theorem 3.2 in Sect. 5.3.

4.1 The Limiting Controlled Diffusion with Compound Poisson Jumps

The controlled generator of the controlled limiting jump diffusion in (3.3) is given by

$$\mathcal{A}\varphi(x,u) = \sum_{i \in \mathcal{I}} b_i(x,u)\partial_i \varphi(x) + \frac{1}{2} \sum_{i \in \mathcal{I}} \lambda_i (1 + c_{a,i}^2) \partial_{ii} \varphi(x)$$

$$+ \int_{\mathbb{R}^d} (\varphi(x+y) - \varphi(x)) \nu_L(\mathrm{d}y)$$
(4.1)

for $\varphi \in C^2(\mathbb{R}^d)$, where the drift b satisfies (3.4), and $v_L(A) := \Pi_L(\{z \in \mathbb{R}_* : \lambda z \in A\})$ for any Borel measurable set A, with Π_L as in (2.4). We refer the reader to [20, Sect. 6] for the definition of exponential ergodicity. The following proposition is a direct consequence of [24, Theorem 3.5].

Proposition 4.1 Under any constant control v such that $\Gamma v \neq 0$, the controlled limiting jump diffusion in (3.3) is exponentially ergodic.

Remark 4.1 It is shown in [25Theorem 5*] that the limiting controlled jump diffusion is exponentially ergodic uniformly over all stationary Markov controls resulting in a locally Lipschitz continuous drift, if $\Gamma > 0$.



Proposition 4.1 implies that the optimal control problems for the limiting jump diffusion are well-posed.

4.2 Preliminaries

We denote the scaled hazard rate function of G_i by r_i^n . This is defined by

$$r_i^n(h_i) := \frac{\lambda_i^n \dot{F}_i(\lambda_i^n h_i)}{1 - F_i(\lambda_i^n h_i)}, \quad \forall h_i \in \mathbb{R}_+, \quad \forall i \in \mathcal{I},$$

where \dot{F}_i denotes the right derivative of F_i . Recall H^n in (2.10). The extended generator of (A^n, H^n) associated with the renewal arrival processes, denoted by \mathcal{H}^n , is given by

$$\mathcal{H}^n f(x,h) = \sum_{i \in \mathcal{I}} \frac{\partial f(x,h)}{\partial h_i} + \sum_{i \in \mathcal{I}} r_i^n(h_i) \left(f(x+e_i,h-h_ie_i) - f(x,h) \right) \tag{4.2}$$

for $f \in \mathcal{C}_b(\mathbb{R}^d \times \mathbb{R}^d_+)$.

Remark 4.2 We sketch the derivation of (4.2); see also [26Theorem 5.5*]. It is enough to consider one component (A_i^n, H_i^n) , $i \in \mathcal{I}$. We obtain

$$\begin{split} &\mathbb{E}_{x,h} \Big[f \Big(A_i^n(t+s), H_i^n(t+s) \Big) \Big] - f(x,h) \\ &= \mathbb{E}_{x,h} \Big[f \Big(A_i^n(t+s), H_i^n(t+s) \Big) \Big] - \mathbb{E}_{x,h} \Big[f \Big(A_i^n(t+s), h \Big) \Big] \\ &+ \mathbb{E}_{x,h} \Big[f \Big(A_i^n(t+s), h \Big) \Big] - f(x,h) \\ &= r_{i,0,s}^n(h) \Big(f(x,h+s) - f(x,h) \Big) + r_{i,1,s}^n(h) \Big(f(x+1,h) - f(x,h) \Big) \\ &+ \sum_{j \in \mathbb{N}} r_{i,j,s}^n(h) \mathbb{E}_{x,h} \Big[f \Big(x+j, H_i^n(t+s) \Big) - f(x+j,h) \mid A_i^n(t+s) = x+j \Big] \\ &+ \sum_{j \in \mathbb{N}, j \ge 2} r_{i,j,s}^n(h) \Big(f(x+j,h) - f(x,h) \Big) \quad \forall \, f \in \mathcal{C}_b(\mathbb{R} \times \mathbb{R}), \, \forall \, (x,h) \in \mathbb{R} \times \mathbb{R}_+, \end{split}$$

where

$$r_{i,j,s}^{n}(h) := \mathbb{P}\left(A_{i}^{n}(t+s) = x+j \mid A_{i}^{n}(t) = x, H_{i}^{n}(t) = h\right) = \mathbb{P}\left(A_{i}^{n}(s+h) = j \mid G_{i} \geq \lambda_{i}^{n}h\right)$$

by the regenerative property of renewal process. Since $\dot{F}_i(t)$ is finite for all $t \ge 0$, it follows that

$$r_i^n(h) \equiv \lim_{s \searrow 0} \frac{1}{s} r_{i,1,s}^n(h) = \frac{\lambda_i^n \dot{F}_i(\lambda_i^n h_i)}{1 - F_i(\lambda_i^n h_i)}, \text{ and } \lim_{s \searrow 0} \frac{1}{s} r_{i,j,s}^n(h) = 0 \text{ for } j \ge 2.$$

It is evident that $\lim_{s\searrow 0} r_{i,0,s}^n = 1$ and $\lim_{s\searrow 0} r_{i,j,s}^n = 0$ for $j \in \mathbb{N}$. Thus, we obtain (4.2).



We define (compare this with [19])

$$\eta_i^n(h_i) := 1 - \frac{\int_{\lambda_i^n h_i}^{\infty} (1 - F_i(y)) \, \mathrm{d}y}{1 - F_i(\lambda_i^n h_i)}, \quad h_i \in \mathbb{R}_+, \ i \in \mathcal{I}.$$
(4.3)

Note that η_i^n is bounded by (3.10). The following identity is frequently used throughout the paper.

$$\dot{\eta}_i^n(h_i) - \eta_i^n(h_i)r_i^n(h_i) = \lambda_i^n - r_i^n(h_i), \quad \forall h_i \in \mathbb{R}_+, \quad \forall i \in \mathcal{I}.$$
 (4.4)

Recall that $c_{a,i}^2$ denotes the squared coefficient of variation of G_i . Let

$$\kappa_i^n(h_i) := \frac{\int_{\lambda_i^n h_i}^{\infty} \int_t^{\infty} (1 - F_i(x)) \, \mathrm{d}x \, \mathrm{d}t}{1 - F_i(\lambda_i^n h_i)} - \frac{c_{a,i}^2 + 1}{2} \frac{\int_{\lambda_i^n h_i}^{\infty} (1 - F_i(x)) \, \mathrm{d}x}{1 - F_i(\lambda_i^n h_i)}$$
(4.5)

for $h_i \in \mathbb{R}_+$ and $i \in \mathcal{I}$. Note that the first term on the right-hand side of (4.5) is the second order residual life function. It follows by (3.10) that κ_i^n is bounded. Using (4.5), we obtain $\kappa_i^n(0) = 0$, and

$$\dot{\kappa}_{i}^{n}(h_{i}) - r_{i}^{n}(h_{i})\kappa_{i}^{n}(h_{i}) = \left(\eta_{i}^{n}(h_{i}) + \frac{c_{a,i}^{2} - 1}{2}\right)\lambda_{i}^{n}, \quad h_{i} \in \mathbb{R}_{+}, \ i \in \mathcal{I}.$$
 (4.6)

The scaled hazard rate function of d_1 is defined by

$$\beta_{\mathsf{d}}^{n}(k) := \frac{\vartheta^{n} \dot{F}^{d_{1}}(\vartheta^{n} k)}{1 - F^{d_{1}}(\vartheta^{n} k)}, \quad k \in \mathbb{R}_{+}.$$

Recall K^n in (2.11). The extended generator of (Ψ^n, K^n) associated with the alternating renewal process, denoted by \mathcal{K}^n , is given by

$$\mathcal{K}^{n} f(\psi, k) = \psi \, \beta_{\mathsf{u}}^{n} \big(f(0, 0) - f(1, 0) \big)$$

$$+ (1 - \psi) \bigg(\beta_{\mathsf{d}}^{n}(k) \big(f(1, 0) - f(0, k) \big) + \frac{\partial f(0, k)}{\partial k} \bigg)$$
(4.7)

for $f \in \mathcal{C}_b(\{0, 1\} \times \mathbb{R}_+)$, with β^n_u as in Assumption 2.2. In analogy to (4.4), we define

$$\alpha^{n}(k) := 1 - \frac{\int_{\vartheta^{n}k}^{\infty} \left(1 - F^{d_{1}}(x)\right) dx}{1 - F^{d_{1}}(\vartheta^{n}k)} \quad \forall k \in \mathbb{R}_{+}.$$
 (4.8)

The following identities hold: $\alpha^n(0) = 0$, and

$$\dot{\alpha}^{n}(k) - \beta_{\mathsf{d}}^{n}(k)\alpha^{n}(k) = \vartheta^{n} - \beta_{\mathsf{d}}^{n}(k) \quad \forall k \in \mathbb{R}_{+}. \tag{4.9}$$



Let $\tilde{\alpha}^n(\psi, k) := (\psi + \alpha^n(k))(\vartheta^n)^{-1}$. It follows by (4.9) that

$$\mathcal{K}^n \tilde{\alpha}^n(\psi, k) = -\frac{\beta_u^n}{\vartheta^n} \psi + (1 - \psi). \tag{4.10}$$

Note that $\tilde{\alpha}^n$ is bounded by (3.10).

4.3 Diffusion-Scaled Processes

To prove Theorem 3.2, we need to establish long-run average moment bounds for the diffusion-scaled processes under a class of scheduling policies, which agree with a proposed policy outside a compact set. We make this formal in Proposition 4.3. The proposed policy is given in the next definition.

Let $J_0 := \{i \in \mathcal{I}: \gamma_i = 0\}$. If $J_0 \neq \emptyset$, then, without loss of generality, we assume that $J_0 = \{1, \ldots, |J_0|\}$, where $|J_0|$ denotes the cardinality of the set J_0 . In Definition 4.1 below, we introduce a modified priority scheduling policy which can be described as follows: First, $\lfloor n\rho_i/\sum_{i\in J_0}\rho_i\rfloor \wedge x_i$ servers are allocated to each class $i\in J_0$. Then, the remaining servers are allocated following the static priority rule.

Definition 4.1 The Markov policy \check{z}^n is defined by

$$\begin{split} \check{z}_i^n(x) &= \left\lfloor \frac{n\rho_i}{\sum_{i \in \mathcal{I}_0} \rho_i} + \left(n - \sum_{j \in \mathcal{I}_0} \left(x_j \wedge \left\lfloor \frac{n\rho_j}{\sum_{i \in \mathcal{I}_0} \rho_i} \right\rfloor \right) \right. \\ &\left. - \sum_{i=1}^{i-1} \left(x_j - \left\lfloor \frac{n\rho_j}{\sum_{i \in \mathcal{I}_0} \rho_i} \right\rfloor \right)^+ \right)^+ \right\rfloor \wedge x_i, \quad \forall i \in \mathcal{I}_0, \end{split}$$

and

$$\check{z}_i^n(x) := x_i \wedge \left(n - \sum_{i=1}^{i-1} x_i\right)^+, \quad \forall i \in \mathcal{I} \setminus \mathcal{I}_0.$$

We let $\check{q}_i^n(x) := x_i - \check{z}_i^n(x), i \in \mathcal{I}$.

In obtaining long-run average moment bounds, since the queueing system is in an alternative renewal environment, we do not work with the diffusion-scaled processes directly. To utilize the fact that (Ψ^n, K^n) is a Markov process, we introduce the following auxiliary process. We define the 'unscaled' process \check{X}^n by

$$\check{X}_{i}^{n}(t) := X_{i}^{n}(0) + A_{i}^{n}(t) - S_{i}^{n}(t)
- R_{*,i}^{n} \left(\gamma_{i}^{n} \int_{0}^{t} \left(\check{X}_{i}^{n}(s) - n \mu_{i}^{n} \rho_{i} \mathcal{R}^{n}(s) - Z_{i}^{n}(s) \right) ds \right) + n \mu_{i}^{n} \rho_{i} \mathcal{R}^{n}(t)
= X_{i}^{n}(t) + n \mu_{i}^{n} \rho_{i} \mathcal{R}^{n}(t) \quad \text{a.s.}$$
(4.11)



for $i \in \mathcal{I}$ and $t \geq 0$, where $\mathcal{R}^n(t)$ is the residual time process for the system in the 'down' state given by

$$\mathcal{R}^{n}(t) = \sum_{k=1}^{N_{\mathbf{u}}^{n}(t)} d_{k}^{n} - \int_{0}^{t} (1 - \Psi^{n}(s)) \, \mathrm{d}s,$$

and $N_{\rm u}^n(t)$ is the process counting the number of completed 'up' periods by time t. Here, the second equality in (4.11) follows by the fact that given $X^n(0)$, Ψ^n and Z^n , the evolution equation in (2.7) admits a unique solution. Also, if $\Psi^n(t)=1$, then $\mathcal{R}^n(t)=0$ and thus $\check{X}^n(t)=X^n(t)$ a.s. Note that under a Markov policy $z^n\in\mathfrak{Z}^n_{\rm sm}$, the process $(\check{X}^n,H^n,\Psi^n,K^n)$ is Markov with state space

$$\mathfrak{D} := \left\{ (\check{x}, h, \psi, k) \in \mathbb{R}^d_+ \times \mathbb{R}^d_+ \times \{0, 1\} \times \mathbb{R}_+ \colon k \equiv 0 \text{ if } \psi = 1 \right\},\,$$

and

$$Z^{n}(t) = z^{n} (\check{X}^{n}(t) - n\mu_{i}^{n} \rho_{i} \mathcal{R}^{n}(t), H^{n}(t), \Psi^{n}(t), K^{n}(t))$$

Under $z^n \in \mathfrak{Z}^n_{sm}$, the generator of $(\check{X}^n, H^n, \Psi^n, K^n)$ denoted by $\check{\mathcal{L}}_n^{z^n}$ is given by

$$\check{\mathcal{L}}_{n}^{z^{n}} f(\check{x}, h, \psi, k) = \overline{\mathcal{L}}_{n, \psi}^{z^{n}} f(\check{x}, h, \psi, k) + \mathcal{I}_{n, \psi} f(\check{x}, h, \psi, k) + \mathcal{Q}_{n, \psi} f(\check{x}, h, \psi, k) \tag{4.12}$$

for $(\check{x}, h, \psi, k) \in \mathfrak{D}$ and $f \in \mathcal{C}_b(\mathbb{R}^d \times \mathbb{R}^d_+ \times \{0, 1\} \times \mathbb{R}_+)$. The operators on the right-hand side of (4.12) are defined by

$$\overline{\mathcal{L}}_{n,\psi}^{z^{n}} f(\check{x}, h, \psi, k)
:= \sum_{i \in \mathcal{I}} \frac{\partial f(\check{x}, h, \psi, k)}{\partial h_{i}} + \sum_{i \in \mathcal{I}} r_{i}^{n}(h_{i}) (f(\check{x} + e_{i}, h - h_{i}e_{i}, \psi, k)
- f(\check{x}, h, \psi, k))
+ \psi \sum_{i \in \mathcal{I}} (\mu_{i}^{n} z_{i}^{n} (\check{x}, h, 1, 0) + \gamma_{i}^{n} q_{i}^{n} (\check{x}, z^{n})) (f(\check{x} - e_{i}, h, 1, 0) - f(\check{x}, h, 1, 0))
+ (1 - \psi) \sum_{i \in \mathcal{I}} \gamma_{i}^{n} (f(\check{x} - e_{i}, h, 0, k) - f(\check{x}, h, 0, k))
\times \int_{\mathbb{R}_{*}} q_{i}^{n} (\check{x} - n\mu^{n} (y - k), z^{n}) \tilde{F}_{\check{x}, k}^{d_{1}^{n}} (dy)
- (1 - \psi) \sum_{i \in \mathcal{I}} n\rho_{i} \mu_{i}^{n} \frac{\partial f(\check{x}, h, 0, k)}{\partial \check{x}_{i}} \tag{4.13}$$

with $q^n(\breve{x}, z^n) = \breve{x} - z^n$,

$$\mathcal{I}_{n,\psi} f(\breve{x},h,\psi,k)$$

$$:= \psi \, \beta_{\mathsf{u}}^{n} \int_{\mathbb{R}_{+}} \left(f\left(\check{\mathsf{x}} + \frac{n}{\vartheta^{n}} \mu^{n} \mathsf{y}, h, 0, 0 \right) - f(\check{\mathsf{x}}, h, 1, 0) \right) F^{d_{1}}(\mathsf{d} \mathsf{y}), \tag{4.14}$$

and

$$Q_{n,\psi} f(\check{x}, h, \psi, k) = (1 - \psi) \left(\beta_{\mathsf{d}}^{n}(k) \left(f(\check{x}, h, 1, 0) - f(\check{x}, h, 0, k) \right) + \frac{\partial f(\check{x}, h, 0, k)}{\partial k} \right). \tag{4.15}$$

In (4.13), $\mu^n := (\mu_1^n \rho_1, \dots, \mu_d^n \rho_d)'$, $\tilde{F}_{\tilde{x},k}^{d_1^n}$ denotes the conditional distribution of d_1^n given $\{d_1^n > k\}$, and $\{n\mu_i^n \rho_i(d_1^n - k) \le \check{x}_i : i \in \mathcal{I}\}$.

The first two terms on the right-hand side of (4.13) correspond to the extended generator associated with the renewal arrival processes. Compare this with (4.2). Conditioning on the alternative renewal process Ψ^n in the 'up' state, the third term on the right-hand side of (4.13) corresponds to the service and abandonment processes, and $\mathcal{I}_{n,\psi}$ corresponds to the residual time process \mathcal{R}^n together with Ψ^n . Similarly, conditioning on the alternative renewal process in the 'down' state, the last two terms on the right-hand side of (4.13) correspond to the abandonment process and \mathcal{R}^n , respectively, and $\mathcal{Q}_{n,\psi}$ corresponds to (Ψ^n, K^n) . The generators in (4.14) and (4.15) are analogous to the extended generator associated with the alternating renewal process in (4.7).

Remark 4.3 We sketch the derivation of $\mathcal{I}_{n,\psi}$. The rest of the terms in (4.12) follow by the calculation below and Remark 4.2. To simplify the calculation, we assume that the arrival processes are Poisson, and only consider the *i*th component $(\check{X}_i^n, \Psi^n, K^n)$, $i \in \mathcal{I}$. Note that $K^n(t) = 0$ when $\Psi^n(t) = 1$. Since there are no simultaneous jumps w.p.1., here we only consider the jumps caused by Ψ^n , that is, we consider

$$\begin{split} &\sum_{j\in\mathbb{N}} \Big(\mathbb{E}_{\breve{x},1,0} \big[f(\breve{X}_i^n(t+s), \Psi^n(t+s), K^n(t+s)) \mid \breve{N}^n(t+s) - \breve{N}^n(t) = j \big] \\ &- f(\breve{x},1,0) \Big) p_j^n(t,s), \end{split}$$

for $s,t\geq 0$, where $\check{N}^n(t)$ denotes the number of jumps of Ψ^n up to time t, and $p^n_j(t,s)=\mathbb{P}\big(\check{N}^n(t+s)-\check{N}^n(t)=j\big),\ j\in\mathbb{N}.$ By the memoryless property of 'up' times, and using the same calculation as in Remark 4.2 for 'down' times, it is straightforward to check that

$$\lim_{s \to 0} \frac{1}{s} p_1^n(t, s) = \beta_{\mathsf{u}}^n, \text{ and } \lim_{s \to 0} \frac{1}{s} p_j^n(t, s) = 0 \text{ for } j \ge 2,$$

and for any $t \ge 0$. By the continuity of K^n , we have

$$\lim_{s \searrow 0} \mathbb{P} \big(\check{N}^n(t+s) - \check{N}^n(t) = 1, \, K^n(t+s) = 0 \, \big| \, K^n(t) = 0 \big) \, = \, 1.$$



Thus,

$$\lim_{s \searrow 0} \mathbb{E}_{\check{x},1,0} \Big[f(\check{X}_i^n(t+s), \Psi^n(t+s), K^n(t+s)) \mid \check{N}^n(t+s) - \check{N}^n(t) = 1 \Big]$$

$$= \mathbb{E}_{\check{x},1,0} \Big[f\Big(\check{x} + n\mu_i^n \rho_i \frac{1}{\vartheta^n} d_1, 0, 0\Big) \Big].$$

This proves (4.14).

Definition 4.2 We define $\bar{x}_i^n(\check{x}) := \check{x}_i - \rho_i n, i \in \mathcal{I}$,

$$\bar{x} = \bar{x}^n(\check{x}) := \left(\bar{x}_1^n(\check{x}), \dots, \bar{x}_d^n(\check{x})\right)', \quad \tilde{x} = \tilde{x}^n(\check{x}) := n^{-1/2}\bar{x}^n(\check{x}), \quad \check{x} \in \mathbb{R}^d,$$

and

$$\mathfrak{A}_R^n := \left\{ x \in \mathbb{R}^d \colon |x - \rho n| \le R\sqrt{n} \right\}$$

for a positive constant R.

Let $\widetilde{\mathcal{L}}_n^{z_n}$ denote the generator of the scaled joint process $\widetilde{\Xi}^n:=(\widetilde{X}^n,H^n,\Psi^n,K^n)$ with $\widetilde{X}^n:=n^{-1/2}(\check{X}^n-n\rho)$. The state space of $\widetilde{\Xi}^n$ is given by

$$\widetilde{\mathfrak{D}}^n := \big\{ (\widetilde{x}^n(\widecheck{x}), h, \psi, k) \in \mathbb{R}^d \times \mathbb{R}^d_+ \times \{0, 1\} \times \mathbb{R}_+ \colon \widecheck{x} \in \mathbb{R}^d_+, \ k \equiv 0 \text{ if } \psi = 1 \big\}.$$

Then, under any $z^n \in \mathfrak{Z}^n_{sm}$, we have

$$\widetilde{\mathcal{L}}_{n}^{z_{n}} f(\tilde{x}, h, \psi, k) = \check{\mathcal{L}}_{n}^{z_{n}} f(\tilde{x}^{n}(\check{x}), h, \psi, k), \tag{4.16}$$

for $f \in \mathcal{C}_b(\mathbb{R}^d \times \mathbb{R}^d_+ \times \{0, 1\} \times \mathbb{R}_+)$.

The next lemma concerns the ergodicity of the process $\widetilde{\Xi}^n$ under the modified priority policy in Definition 4.1. Let $\mathcal{V}_{\kappa,\xi}(x) := \sum_{i \in \mathcal{I}} \xi_i |x_i|^{\kappa}$ for $x \in \mathbb{R}^d$, where $\kappa > 0$, and ξ is a positive vector. Define the function $\widetilde{\mathcal{V}}_{\kappa,\xi}^n : \mathbb{R}^d \times \mathbb{R}_+^d \times \{0,1\} \times \mathbb{R}_+ \to \mathbb{R}$ by

$$\widetilde{\mathcal{V}}_{\kappa,\xi}^{n}(x,h,\psi,k) := \mathcal{V}_{\kappa,\xi}(x) + \sum_{i \in \mathcal{I}} \eta_{i}^{n}(h_{i}) \left(\mathcal{V}_{\kappa,\xi}(x+n^{-1/2}e_{i}) - \mathcal{V}_{\kappa,\xi}(x) \right)
+ \frac{\psi + \alpha^{n}(k)}{\vartheta^{n}} \sum_{i \in \mathcal{I}} \mu_{i}^{n} \xi_{i}
\times \left(\widetilde{\mathcal{V}}_{\kappa,i}^{n}(x_{i}) + \eta_{i}^{n}(h_{i}) \left(\widetilde{\mathcal{V}}_{\kappa,i}^{n}(x_{i}+n^{-1/2}) - \widetilde{\mathcal{V}}_{\kappa,i}^{n}(x_{i}) \right) \right),$$
(4.17)

where η_i^n and α^n are as in (4.3) and (4.8), respectively, and $\tilde{\mathcal{V}}_{\kappa,i}^n(x_i) := -|x_i|^{\kappa}$ for $x_i \in \mathbb{R}_+$ and $i \in \mathcal{I} \setminus \mathcal{I}_0$, and

$$\tilde{\mathcal{V}}_{\kappa,i}^{n}(x_{i}) := \begin{cases} -|x_{i}|^{\kappa}, & \text{for } x_{i} < \frac{\sqrt{n}\rho_{i}\sum_{j\in\mathbb{J}\setminus\mathbb{J}_{0}}\rho_{j}}{\sum_{j\in\mathbb{J}\setminus\mathbb{J}_{0}}\rho_{j}}, \\ -\frac{\sqrt{n}\rho_{i}\sum_{j\in\mathbb{J}\setminus\mathbb{J}_{0}}\rho_{j}}{\sum_{j\in\mathbb{J}_{0}}\rho_{j}}|x_{i}|^{\kappa-1}, & \text{for } x_{i} \geq \frac{\sqrt{n}\rho_{i}\sum_{j\in\mathbb{J}\setminus\mathbb{J}_{0}}\rho_{j}}{\sum_{j\in\mathbb{J}_{0}}\rho_{j}}, \end{cases} \forall i \in \mathbb{J}_{0}.$$



The function $\widetilde{\mathcal{V}}_{\kappa,\xi}^n$ is constructed in such a manner as to allow us to take advantage of the identities in (4.4) and (4.10). We define the set

$$\mathcal{K}_n(x) \,:=\, \bigg\{ i \in \mathcal{I}_0 \colon x_i \,\geq\, \frac{\sqrt{n} \, \rho_i \, \sum_{j \in \mathcal{I} \setminus \mathcal{I}_0} \, \rho_j}{\sum_{j \in \mathcal{I}_0} \, \rho_j} \bigg\}.$$

Note that $\widetilde{\mathcal{L}}_n^{\check{z}^n}$ denotes the generator of $\widetilde{\Xi}^n$ under the modified priority scheduling policy in Definition 4.1. We have the following lemma.

Lemma 4.1 Grant Assumptions 2.1, 2.2, and 3.2. For any even integer $\kappa \geq 2$, there exist positive constants \widetilde{C}_0 and \widetilde{C}_1 , a positive vector $\xi \in \mathbb{R}^d_+$, and $\widetilde{n} \in \mathbb{N}$ such that:

$$\widetilde{\mathcal{L}}_{n}^{\check{\xi}^{n}}\widetilde{\mathcal{V}}_{\kappa,\xi}^{l}(\tilde{x},h,\psi,k) \leq \widetilde{C}_{0} - \widetilde{C}_{1} \sum_{i \in \mathcal{I} \setminus \mathcal{K}_{n}(\tilde{x})} \mathcal{V}_{\kappa,\xi}(\tilde{x}) - \widetilde{C}_{1} \sum_{i \in \mathcal{K}_{n}(\tilde{x})} \mathcal{V}_{\kappa-1,\xi}(\tilde{x}) \quad (4.18)$$

for all $n > \tilde{n}$, and $(\tilde{x}, h, y, k) \in \widetilde{\mathfrak{D}}^n$. As a consequence, for all large enough $n, \widetilde{\Xi}^n$ is positive Harris recurrent under the modified priority scheduling policy \check{z}^n .

The proof of Lemma 4.1 is given in Appendix B. We continue with the following prop, which plays a crucial role in proving Proposition 4.3. In its proof, especially, equation (4.26), we show the relationship between the processes \hat{X}^n and \tilde{X}^n .

Proposition 4.2 *Grant Assumptions* 2.1, 2.2, and 3.2. *Under the scheduling policy* \check{z}^n *in Definition* 4.1, and for any $\kappa > 0$, there exists $\check{n} \in \mathbb{N}$ such that

$$\sup_{n>\check{n}} \limsup_{T\to\infty} \frac{1}{T} \mathbb{E}^{\check{z}^n} \left[\int_0^T |\hat{X}^n(s)|^{\kappa} \, \mathrm{d}s \right] < \infty. \tag{4.19}$$

Proof Let $\kappa \geq 2$ be an arbitrary even integer. By (4.18), we have

$$\mathbb{E}^{\check{z}^{n}} \left[\widetilde{\mathcal{V}}_{\kappa,\xi}^{n} \left(\widetilde{\Xi}^{n}(T) \right) \right] - \mathbb{E}^{\check{z}^{n}} \left[\widetilde{\mathcal{V}}_{\kappa,\xi}^{n} \left(\widetilde{\Xi}^{n}(0) \right) \right] \\
= \mathbb{E}^{\check{z}^{n}} \left[\int_{0}^{T} \widetilde{\mathcal{L}}_{n}^{\check{z}^{n}} \widetilde{\mathcal{V}}_{\kappa,\xi}^{n} \left(\widetilde{\Xi}^{n}(s) \right) ds \right] \\
\leq \widetilde{C}_{0} T - \widetilde{C}_{1} \mathbb{E}^{\check{z}^{n}} \left[\int_{0}^{T} \mathcal{V}_{\kappa-1,\xi} \left(\widetilde{X}^{n}(s) \right) ds \right]. \tag{4.20}$$

Since $(\vartheta^n)^{-1}$ is of order $n^{-1/2}$ by Assumption 2.2, it follows by Young's inequality together with (3.10) that there exist some positive constants c_0 and c_1 such that $c_0(\mathcal{V}_{\kappa,\xi}-1)\leq\widetilde{\mathcal{V}}^n_{\kappa,\xi}\leq c_1(1+\mathcal{V}_{\kappa,\xi})$ for all large n. Note that $\hat{X}^n(0)=\widetilde{X}^n(0)$. Thus, by (4.20), we obtain

$$\widetilde{C}_1 \operatorname{\mathbb{E}}^{\check{z}^n} \left[\int_0^T \mathcal{V}_{\kappa-1,\xi} \left(\widetilde{X}^n(s) \right) \mathrm{d}s \right] \le (\widetilde{C}_0 + c_0) T + c_1 \left(1 + \mathcal{V}_{\kappa,\xi} \left(\hat{X}^n(0) \right) \right) \tag{4.21}$$



for some positive constants C_3 and C_4 . By dividing both sides of (4.21) by T, and taking $T \to \infty$, we have

$$\sup_{n>\check{n}} \limsup_{T\to\infty} \frac{1}{T} \mathbb{E}^{\check{z}^n} \left[\int_0^T |\widetilde{X}^n(s)|^{\kappa-1} \, \mathrm{d}s \right] < \infty. \tag{4.22}$$

Let $\mathbb{E} \equiv \mathbb{E}^{U^n}$ for some admissible scheduling policy U^n . We have

$$\frac{1}{T} \mathbb{E} \left[\int_0^T |\hat{X}_i^n(s) - \widetilde{X}_i^n(s)|^{\kappa - 1} \, \mathrm{d}s \right] \\
= (\mu_i^n \rho_i)^{\kappa - 1} \frac{1}{T} \mathbb{E} \left[\int_0^T \left(\sqrt{n} \mathcal{R}^n(s) \right)^{\kappa - 1} \, \mathrm{d}s \right] \quad \forall i \in \mathcal{I}. \tag{4.23}$$

We use the identity

$$\mathbb{E}\left[\left(\sqrt{n}\mathcal{R}^n(s)\right)^{\kappa-1}\right] = \mathbb{E}\left[\left(\sqrt{n}\mathcal{R}^n(s)\right)^{\kappa-1} | \mathcal{R}^n(s) > 0\right] \mathbb{P}(\mathcal{R}^n(s) > 0) \tag{4.24}$$

for any $s \ge 0$. Here $\mathcal{R}^n(s)$ is the residual time of the system in the 'down' state, and thus $\mathbb{E}[(\sqrt{n}\mathcal{R}^n(s))^{\kappa-1}|\mathcal{R}^n(s)>0] \le \mathbb{E}[(\sqrt{n}d_1^n)^{\kappa-1}] \le c_2$ for some positive constant c_2 , by Assumption 2.2 and (3.10). Also, $\mathbb{P}(\mathcal{R}^n(s)>0)=\mathbb{P}(\Psi^n(s)=0)$, and it follows by [27, Theorem 3.4.4] that

$$\lim_{s \to \infty} \mathbb{P}(\Psi^{n}(s) = 0) = \frac{(\vartheta^{n})^{-1}}{(\beta_{\mathbf{u}}^{n})^{-1} + (\vartheta^{n})^{-1}},$$

which is of order $n^{-1/2}$ by Assumption 2.2. Therefore, applying (4.24), we obtain

$$\lim_{(n,T)\to\infty} \frac{1}{T} \mathbb{E}\left[\int_0^T \left(\sqrt{n}\mathcal{R}^n(s)\right)^{\kappa-1} ds\right] = 0. \tag{4.25}$$

It follows by (4.23) and (4.25) that

$$\lim_{(n,T)\to\infty} \frac{1}{T} \mathbb{E} \left[\int_0^T \|\hat{X}^n(s) - \widetilde{X}^n(s)\|^{\kappa - 1} \, \mathrm{d}s \right] = 0. \tag{4.26}$$

Thus (4.19) follows by (4.22) and (4.26). This completes the proof.

The next prop is used to prove the upper bound for the ergodic control problem in Sect. 5.3.2, where we adopt the spatial truncation technique developed in [4]. We first introduce a class of concatenated scheduling policies.

Definition 4.3 We define the *quantization function* $\varpi: \mathbb{R}^d_+ \to \mathbb{Z}^d_+$ by

$$\varpi(x) := \left(\lfloor x_1 \rfloor, \dots, \lfloor x_{d-1} \rfloor, \lfloor x_d \rfloor + \sum_{i=1}^d (x_i - \lfloor x_i \rfloor) \right).$$



For a sequence $v^n : \mathbb{R}^d \to \mathcal{S}$, $n \in \mathbb{N}$, of continuous functions satisfying $v^n(\tilde{x}^n(x)) = e_d$ if $x \notin \mathfrak{A}_R^n$, R > 1, with \mathfrak{A}_R^n as in Definition 4.2, we define the map

$$q^{n}[v^{n}](x) := \begin{cases} \varpi\left(\left(\langle e, x \rangle - n\right)^{+} v^{n}\left(\tilde{x}^{n}(x)\right)\right) & \text{for } \sup_{i \in \mathbb{J}} |\tilde{x}^{n}(x)| \leq \frac{1}{2d} \sqrt{n}\left(\min_{i} \rho_{i}\right), \\ \check{q}^{n}(x) & \text{for } \sup_{i \in \mathbb{J}} |\tilde{x}^{n}(x)| > \frac{1}{2d} \sqrt{n}\left(\min_{i} \rho_{i}\right), \end{cases}$$

and the scheduling policy $z^n[v^n](x) := x - q^n[v^n](x)$

Proposition 4.3 *Under the scheduling policy* $z^n[v^n]$ *in Definition 4.3, the conclusions in Lemma 4.1 and Proposition 4.2 hold.*

Proof For all sufficiently large n, we have $q_i^n[v^n](\check{x}) \leq 2dR\sqrt{n}$ for $\check{x} \in \mathfrak{A}_R^n$ (see also the proof of [4, Lemma 5.1]). If $\sup_{i \in \mathbb{J}} |\tilde{x}_i^n(\check{x})| \leq \frac{1}{d} \sqrt{n} (\min_i \rho_i)$, it is evident that $\sum_{i=1}^{d-1} \check{x}_i \leq n$, and thus $z^n[e_d]$ is equivalent to the modified priority policy on this set. Therefore, the result follows by the argument in Lemma 4.1 and Proposition 4.2. \square

5 Asymptotic Optimality

5.1 Results Concerning the Limiting Jump Diffusion

In this subsection, we present some optimality results for the limiting jump diffusion. These results are used in proving asymptotic optimality.

Recall that a stationary Markov control v is called stable if the process under v is positive recurrent, and the set of such controls is denoted by \mathfrak{U}_{ssm} . Let \mathfrak{G} denote the set of ergodic occupation measures, that is,

$$\mathfrak{G} := \left\{ \pi \in \mathcal{P}(\mathbb{R}^d \times \mathbb{U}) \colon \int_{\mathbb{R}^d \times \mathbb{U}} \mathcal{A}f(x, u) \, \pi(\mathrm{d}x, \mathrm{d}u) \, = \, 0 \quad \forall \, f \in \mathcal{C}_c^{\infty}(\mathbb{R}^d) \right\}. \tag{5.1}$$

See [28, Sect. 2.1] for more details.

We summarize the characterization of optimal controls for the limiting jump diffusion in the following theorem. Recall the definition of d_1 in Assumption 2.2.

Theorem 5.1 Assume that $\mathbb{E}[(d_1)^{m+1}] < \infty$ with m as in (3.5). The following hold:

(i) For $\alpha > 0$, V_{α} in (3.7) is the minimal nonnegative solution in $C^{2,r}(\mathbb{R}^d)$, $r \in (0, 1)$, to the HJB equation

$$\min_{u \in \mathbb{U}} \left[AV_{\alpha}(x, u) + \Re(x, u) \right] = \alpha V_{\alpha}(x) \quad a.e. \text{ in } \mathbb{R}^d.$$
 (5.2)

In addition, V_{α} has at most polynomial growth with degree m. Moreover, a stationary Markov control v is optimal for the α -discounted problem if and only if it is an a.e. measurable selector from the minimizer in (5.2).

(ii) There exists a solution $V \in \mathcal{C}^{2,r}(\mathbb{R}^d)$, $r \in (0,1)$, to the HJB equation

$$\min_{u \in \mathbb{U}} \left[AV(x, u) + \Re(x, u) \right] = \varrho_* \quad a.e. \text{ in } \mathbb{R}^d.$$
 (5.3)



Moreover, a stationary Markov control v is optimal for the ergodic control problem if and only if it is an a.e. measurable selector from the minimizer (5.3).

Proof We first consider (i). It follows by [24, Remark 5.1] and Proposition 4.1 that [17, Assumptions 2.1 and 2.2] hold with V_{\circ} and V having at most polynomial growth of degree m. Since $\mathbb{E}[(d_1)^{m+1}] < \infty$, then (4.1) satisfies [17, Assumption 5.1]. Therefore, the results in part (i) follow by [17, Theorems 5.1 and 5.3]. Note that by [17, (5.4)], V_{α} has at most polynomial growth of degree m. Similarly, the claim in part (ii) follows by [17, Theorems 5.2 and 5.3].

Remark 5.1 If there is no jump part in (4.1), then it corresponds to the controlled limiting diffusion for GI/M/n + M queues. If we define the optimal control problems for the limiting diffusion in the same way as in (3.7) and (3.9), then the results in Theorem 5.1 still hold when \mathcal{A} in (4.1) does not contain the jump component. As a consequence, part (i) of Theorem 5.1 corresponds to [1, Theorem 3].

If we consider (3.9) over all stable Markov controls, then the ergodic control problem is equivalent to $\min_{\pi \in \mathcal{G}} \int_{\mathbb{R}^d \times \mathbb{U}} \mathcal{R}(x, u) \, \pi(\mathrm{d}x, \mathrm{d}u)$, see, for example, [17, Sect. 4]. We summarize a result on ϵ -optimal controls for the ergodic problem in the next prop, which follows directly by Corollary 7.1 in [17]. Note that the constant control $v \equiv e_d$ also satisfies Proposition 4.1. Recall that a stationary Markov control v is called precise if it is a measurable map from \mathbb{R}^d to \mathbb{U} .

Proposition 5.1 Assume that $\mathbb{E}[(d_1)^m] < \infty$, with m as in (3.5). For any $\epsilon > 0$, there exist a continuous precise control $v_{\epsilon} \in \mathfrak{U}_{ssm}$, and $R \equiv R(\epsilon) \in \mathbb{N}$ such that $v_{\epsilon} \equiv e_d$ on \bar{B}_R^c , and v_{ϵ} is ϵ -optimal, that is,

$$\int_{\mathbb{R}^d\times\mathbb{U}} \Re(x,u)\,\pi_{v_\epsilon}(\mathrm{d} x,\mathrm{d} u)\,\leq\,\varrho_*+\epsilon.$$

5.2 Proof of Theorem 3.1

To prove Theorem 3.1, we use the approach developed in [1]. We first establish a key moment estimate for the diffusion-scaled process \hat{X}^n , whose proof is similar to that of [1, Lemma 3].

Lemma 5.1 *Grant the hypotheses in Theorem* 3.1. *Then*

$$\mathbb{E}[\|\hat{X}^n(t)\|^{m_A}] \le c_1(1+t^{m_1})(1+\|x\|^{m_1}) \quad \forall t \ge 0, \tag{5.4}$$

where c_1 and m_1 are some positive constants independent of n, x and t.

Proof Recall \hat{L}^n and \hat{X}^n in (3.1), and \hat{W}^n in (3.2). Let $\hat{\Phi}^n$ be a d-dimensional process defined by $\hat{\Phi}^n_i(\cdot) := \mu^n_i \int_0^{\cdot} \hat{Z}^n_i(s) (1 - \Psi^n(s)) ds$, for $i \in \mathcal{I}$. Then,

$$\mu_i^n \int_0^t \hat{Z}_i^n(s) \Psi^n(s) \, \mathrm{d}s = -\hat{\Phi}_i^n(t) + \mu_i^n \int_0^t \hat{Z}_i^n(s) \, \mathrm{d}s \qquad \forall \, t \geq 0.$$



Thus, we obtain

$$\hat{X}_{i}^{n}(t) = \hat{X}_{i}^{n}(0) + \ell_{i}^{n}t + \hat{W}_{i}^{n}(t) + \hat{\Phi}_{i}^{n}(t) + \hat{L}_{i}^{n}(t) - \mu_{i}^{n} \int_{0}^{t} \hat{Z}_{i}^{n}(s) \, ds - \gamma_{i}^{n} \int_{0}^{t} \hat{Q}_{i}^{n}(s) \, ds$$

for all $t \ge 0$ and $i \in J$. Following the same method as in [1, Lemma 3], we have

$$\|\hat{X}^{n}(t)\| \leq C \left[1 + t^{2} + \|\hat{X}^{n}(0)\| + \|\hat{W}^{n}(t) + \hat{L}^{n}(t) + \hat{\Phi}^{n}(t)\| \right]$$

$$+ \int_{0}^{t} \|\hat{W}^{n}(s) + \hat{L}^{n}(s) + \hat{\Phi}^{n}(s)\| \, ds$$

$$+ \int_{0}^{t} \int_{0}^{s} \|\hat{W}^{n}(r) + \hat{L}^{n}(r) + \hat{\Phi}^{n}(r)\| \, dr \, ds \right]$$

$$(5.5)$$

for some positive constant C. Let

$$\widehat{N}^n(t) := \max \left\{ k \ge 0 \colon \sum_{i=1}^k u_i^n \le t \right\}$$

with u^n as in (2.2). By Assumption 2.2, $\widehat{N}^n(t)$ is a Poisson process with rate β_u^n . Then, we obtain

$$\mathbb{E}\left[\|\hat{L}^{n}(t)\|^{m_{A}}\right] \leq C_{1} \mathbb{E}\left[\left(\sqrt{n}C_{\mathsf{d}}^{n}(t)\right)^{m_{A}}\right] \leq C_{1}\left(\frac{\sqrt{n}}{\vartheta^{n}}\right)^{m_{A}} \mathbb{E}\left[\left(\sum_{i=1}^{N^{n}(t)+1} d_{i}\right)^{m_{A}}\right]$$

$$\leq C_{2}(1+t^{m_{2}}) \tag{5.6}$$

for some positive constants $C_1 = \sup\{\mu_i^n \rho_i : n \in \mathbb{N}, i \in \mathcal{I}\}$, C_2 , and m_2 . The third inequality in (5.6) follows by the independence of \widehat{N}^n and d_i , and Assumption 3.1. On the other hand, for some positive constant C_3 , we have

$$|n^{-1/2}\hat{Z}_{i}^{n}(s)| \le C_{3}(1+n^{-1}A_{i}^{n}(s))$$
 a.s. $\forall s \ge 0$. (5.7)

Thus,

$$\mathbb{E}\Big[\big|\hat{\Phi}_{i}^{n}(t)\big|^{m_{A}}\Big] \leq \mu_{i}^{n} \,\mathbb{E}\Big[\bigg(\int_{0}^{t} \big|n^{-1/2}\hat{Z}_{i}^{n}(s)\big|\big|\sqrt{n}\big(1-\Psi^{n}(s)\big)\big|\,\mathrm{d}s\bigg)^{m_{A}}\Big] \\
\leq \mu_{i}^{n}(C_{3})^{m_{A}}\Big(1+\sup_{s\leq t} \mathbb{E}\big[n^{-1}A_{i}^{n}(s)\big]\Big)^{m_{A}} \,\mathbb{E}\Big[\big(\sqrt{n}C_{\mathsf{d}}^{n}(t)\big)^{m_{A}}\Big]$$

$$\leq C_{4}(1+t^{m_{3}})$$
(5.8)

for some positive constant C_4 , where the second inequality follows by (5.7) and the independence of A^n and Ψ^n , and the third inequality follows by [29, Theorem 4] and (5.6). Therefore, following the argument in the proof of [1, Lemma 3], and using (5.5), (5.6), and (5.8), we establish (5.4). This completes the proof.



Proof of Theorem 3.1 We first prove the lower bound:

$$\liminf_{n\to\infty} \hat{V}_{\alpha}^{n}(\hat{X}^{n}(0)) \geq V_{\alpha}(x).$$

By Theorem 5.1, the partial derivatives of $V_{\alpha}(x)$ up to order two are locally Hölder continuous. Let $V_{\alpha}^{l} := \chi_{l} \circ V_{\alpha} = \chi_{l}(V_{\alpha})$, where $\chi_{l} \in \mathcal{C}^{2}(\mathbb{R})$ satisfies $\chi_{l}(x) = x$ for $x \leq l$ and $\chi_{l}(x) = l + 1$ for $x \geq l + 2$. Let $\mathcal{L} : \mathcal{C}^{2}(\mathbb{R}^{d}) \to \mathcal{C}^{2}(\mathbb{R}^{d} \times \mathcal{S})$ be the local operator defined by

$$\mathcal{L}\varphi(x,u) := \langle b(x,u), \nabla \varphi(x) \rangle + \frac{1}{2} \sum_{i \in \mathbb{I}} \lambda_i (1 + c_{a,i}^2) \, \partial_{ii} \varphi(x), \qquad \varphi \in \mathcal{C}^2(\mathbb{R}^d).$$

Compare this with (4.1). We define $\mathcal{H}(x, p) := \min_{u \in \mathbb{U}} [\langle b(x, u), p \rangle + \mathcal{R}(x, u)]$, for $(x, p) \in \mathbb{R}^d \times \mathbb{R}^d$. By Itô's formula, for any $l > \sup_{B_B} V_{\alpha}$, it follows that

$$\begin{split} & e^{-\alpha(t\wedge\tau_R)} V_{\alpha}^l(X_{t\wedge\tau_R}) \\ &= V_{\alpha}^l(x) - \int_0^{t\wedge\tau_R} \alpha e^{-\alpha s} \, V_{\alpha}(X_s) \, \mathrm{d}s \\ &+ \int_0^{t\wedge\tau_R} e^{-\alpha s} \, \mathcal{L} V_{\alpha}(X_s, U_s) \, \mathrm{d}s \\ &+ \int_0^{t\wedge\tau_R} \langle e^{-\alpha s} \, \nabla V_{\alpha}(X_s), \, \Sigma \, \mathrm{d}W_s \rangle + \int_0^{t\wedge\tau_R} \int_{\mathbb{R}_*} e^{-\alpha s} \, \left(V_{\alpha}^l(X_{s-} + \lambda y) - V_{\alpha}(X_{s-}) \right) \mathcal{N}_L(\mathrm{d}s, \mathrm{d}y), \end{split}$$

where \mathcal{N}_L is the Poisson random measure of $\{L_t : t \ge 0\}$ with the intensity Π_L . Thus, applying (5.2), we obtain

$$\begin{split} \mathrm{e}^{-\alpha(t\wedge\tau_R)} V_{\alpha}^l(X_{t\wedge\tau_R}) &= V_{\alpha}^l(x) + \int_0^{t\wedge\tau_R} \mathrm{e}^{-\alpha s} \left\langle b(X_s, U_s), \nabla V_{\alpha}(X_s) \right\rangle \mathrm{d}s \\ &+ \int_0^{t\wedge\tau_R} \left\langle \mathrm{e}^{-\alpha s} \, \nabla V_{\alpha}(X_s), \, \Sigma \, \mathrm{d}W_s \right\rangle - \int_0^{t\wedge\tau_R} \mathrm{e}^{-\alpha s} \, \mathcal{H} \big(X_s, \, \nabla V_{\alpha}(X_s) \big) \, \mathrm{d}s \\ &+ \int_0^{t\wedge\tau_R} \int_{\mathbb{R}_*} \mathrm{e}^{-\alpha s} \left(V_{\alpha}^l(X_{s-} + \lambda y) - V_{\alpha}(X_{s-}) \right) \widetilde{\mathcal{N}}_L(\mathrm{d}s, \mathrm{d}y) \\ &+ \int_0^{t\wedge\tau_R} \int_{\mathbb{R}_*} \mathrm{e}^{-\alpha s} \left(V_{\alpha}^l(X_{s-} + \lambda y) - V_{\alpha}(X_{s-} + \lambda y) \right) \Pi_L(\mathrm{d}s, \mathrm{d}y), \end{split}$$

where $\widetilde{\mathcal{N}}_L(t, A) = \mathcal{N}_L(t, A) - t \Pi_L(A)$ for any Borel set $A \subset \mathbb{R}$. Repeating the same calculation as for the claim (71) in [1], we obtain

$$e^{-\alpha(t\wedge\tau_R)}V_{\alpha}^l(X_t) \geq V_{\alpha}^l(x) + \int_0^{t\wedge\tau_R} \langle e^{-\alpha s} \nabla V_{\alpha}^l(X_s), \Sigma dW_s \rangle$$

$$-\int_0^{t\wedge\tau_R} e^{-\alpha s} \mathcal{R}(X_s, U_s) ds$$
(5.9)



$$\begin{split} &+ \int_0^{t \wedge \tau_R} \int_{\mathbb{R}_*} \mathrm{e}^{-\alpha s} \left(V_\alpha^l(X_{s-} + \lambda y) - V_\alpha(X_{s-}) \right) \widetilde{\mathcal{N}}_L(\mathrm{d} s, \mathrm{d} y) \\ &+ \int_0^{t \wedge \tau_R} \int_{\mathbb{R}_*} \mathrm{e}^{-\alpha s} \left(V_\alpha^l(X_{s-} + \lambda y) - V_\alpha(X_{s-} + \lambda y) \right) \Pi_L(\mathrm{d} s, \mathrm{d} y). \end{split}$$

Note that $\widetilde{\mathcal{N}}_L$ is a martingale measure and V_α is nonnegative. Taking expectations on both sides of (5.9), the second and fourth terms on the right-hand side of (5.9) vanish. Thus, first taking limits as $l \to \infty$, and then as $R \to \infty$, it follows by the monotone convergence theorem that

$$\mathbb{E}\bigg[\int_0^t e^{-\alpha s} \, \Re(X_s, U_s) \, \mathrm{d}s\bigg] \geq V_\alpha(x) - \mathbb{E}\big[e^{-\alpha t} \, V_\alpha(X_t)\big].$$

Applying Theorem 5.1 it follows that solutions of (5.2) have at most polynomial growth of degree m, which corresponds to [1, Proposition 5 (i)]. Note that Lemma 5.1 corresponds to Lemma 3 in [1]. The rest of the proof of the lower bound follows exactly the proof of [1, Theorem 4 (i)].

To prove (3.8), we construct a sequence of asymptotically optimal scheduling policies U^n . Let v_{α} be an optimal control to (5.2). Recall the quantization function in Definition 4.3. We define a sequence of scheduling policies

$$\bar{z}^n[v_\alpha](\hat{x}) := \begin{cases} \varpi\left(\langle e, \hat{x} \rangle^+ v_\alpha(\hat{x})\right), & \text{if } \hat{x} \in \hat{\mathfrak{X}}^n, \\ \check{z}^n(\sqrt{n}\hat{x} + n\rho) & \text{if } \hat{x} \notin \hat{\mathfrak{X}}^n, \end{cases}$$

where \check{z}^n is the modified priority policy in Definition 4.1, and

$$\hat{\mathfrak{X}}^n := \left\{ n^{-1/2} (x - n\rho) \colon x \in \mathbb{R}^d, \ \langle e, x \rangle \le x_i \ \forall i \in \mathcal{I} \right\}.$$

Here the policy on $(\hat{x}^n)^c$ may be chosen arbitrarily. Let $U^n[v_\alpha]$ be the equivalent parameterization of $\bar{z}^n[v_\alpha]$. Following the proof of [1, Theorem 2(i)], we obtain

$$\int_0^{\infty} e^{-\alpha s} \Upsilon^n(s) ds \Rightarrow 0,$$

where

$$\Upsilon^{n}(s) := \langle b(\hat{X}^{n}(s), U^{n}[v_{\alpha}](s)), \nabla V_{\alpha}(\hat{X}^{n}(s)) \rangle + \Re(\hat{X}^{n}(s), U^{n}[v_{\alpha}](s)) \rangle - \Re(\hat{X}^{n}(s), \nabla V_{\alpha}(\hat{X}^{n}(s))).$$

Thus, by using the method in [1, Theorem 4(ii)], and repeating the above calculation, we obtain

$$\limsup_{n\to\infty} \hat{V}_{\alpha}^{n}(\hat{X}^{n}(0)) \leq V_{\alpha}(x).$$

This completes the proof.



5.3 Proof of Theorem 3.2

In this section, we prove Theorem 3.2 by establishing lower and upper bounds.

5.3.1 The Lower Bound

We show that

$$\liminf_{n \to \infty} \varrho^n \left(\hat{X}^n(0) \right) \ge \varrho_*. \tag{5.10}$$

The proof is given at the end of this subsection.

We need the following lemma whose proof is similar to that of Proposition 4.2, and is given in Appendix B.

Lemma 5.2 Grant the hypotheses in Assumptions 2.1, 2.2, and 3.2. For any m > 1, and any sequence $\{z^n \in \mathfrak{Z}^n_{sm} : n \in \mathbb{N}\}$ with $\sup_n \hat{J}(\hat{X}^n(0), z^n) < \infty$, there exists $n_0 > 0$ such that

$$\sup_{n>n_o} \limsup_{T\to\infty} \frac{1}{T} \mathbb{E}^{z^n} \left[\int_0^T |\hat{X}^n(s)|^m \, \mathrm{d}s \right] < \infty.$$
 (5.11)

The main challenge in the proof lies in approximating the generator of the diffusionscaled process with the generator of the limiting jump diffusion. Recall the extended generator \mathcal{H}^n of (A^n, H^n) in (4.2). We define the function $\phi^n[f]$ by

$$\phi^{n}[f](x,h) := f(x) + \sum_{j \in \mathbb{J}} \hat{\phi}_{1,j}^{n}[f](x,h) + \sum_{j \in \mathbb{J}} \frac{c_{a,j}^{2} - 1}{2\sqrt{n}} \partial_{j} f(x)$$

$$+ \sum_{j \in \mathbb{J}} \hat{\phi}_{2,j}^{n}[f](x,h) + \sum_{j \in \mathbb{J}} \frac{\kappa_{j}^{n}(h_{j})}{n} \partial_{jj} f(x)$$

$$+ \sum_{j=1}^{d-1} \hat{\phi}_{3,j}^{n}[f](x,h)$$
(5.12)

for any $f \in \mathcal{C}_c^{\infty}(\mathbb{R}^d)$, and $n \in \mathbb{N}$, where

$$\hat{\phi}_{1,j}^n[f](x,h) := \frac{1}{j!} \sum_{i_j \in \mathbb{J}} \sum_{i_{j-1} \neq i_j} \cdots \sum_{i_1 \notin \{i_l : l > 1\}} \prod_{r=1}^j \eta_{i_r}^n(h_{i_r}) [f]_{i_1 \cdots i_j}^{1,n}(x),$$

with

$$[f]_{i_{1}\cdots i_{j}}^{1,n}(x) := [f]_{i_{1}\cdots i_{j-1}}^{1,n}(x+n^{-1/2}e_{i_{j}}) - [f]_{i_{1}\cdots i_{j-1}}^{1,n}(x),$$

$$[f]_{i_{1}}^{1,n}(x) := f(x+n^{-1/2}e_{i_{1}}) - f(x).$$
(5.13)



The function $\hat{\phi}_{2,j}^n[f]$ is defined analogously to (5.13) with $[f]_{i_1\cdots i_j}^{1,n}$ and $[f]_{i_1}^{1,n}$ replaced by $[f]_{i_1\cdots i_j}^{2,n}$ and

$$[f]_{i_1}^{2,n}(x) := \sum_{j \in \mathcal{I}} \frac{c_{a,j}^2 - 1}{2\sqrt{n}} (\partial_j f(x + n^{-1/2} e_{i_1}) - \partial_j f(x)),$$

respectively. Also,

$$\hat{\phi}^n_{3,j}[f](x,h) := \frac{1}{j!} \sum_{i_j \in \mathbb{I}} \sum_{i_{j-1} \neq i_j} \cdots \sum_{i_1 \notin \{i_l \colon l > 1\}} \prod_{r=2}^{j+1} \eta^n_{i_r}(h_{i_r}) \frac{\kappa^n_{i_1}(h_{i_1})}{n} \big[f \big]^{3,n}_{i_1 \cdots i_{j+1}}(x)$$

with $[f]_{i_1 \cdots i_{i+1}}^{3,n}(x)$ defined analogously to (5.13), and

$$[f]_{i_1i_2}^{3,n}(x) := \partial_{i_1i_1}f(x+n^{-1/2}e_{i_2}) - \partial_{i_1i_1}f(x)$$
 for $i_1, i_2, \dots, i_j, j \in \mathcal{I}$.

Note that $\phi^n[f]$ is bounded by Assumption 3.2(i).

The extended generator $\widetilde{\mathcal{H}}^n$ of the scaled process (\hat{A}^n, H^n) is given by $\widetilde{\mathcal{H}}^n f(\tilde{x}, h) = \mathcal{H}^n f(\tilde{x}^n(x), h)$, for $f \in \mathcal{C}_b(\mathbb{R}^d \times \mathbb{R}^d_+)$. We have the following lemma.

Lemma 5.3 *Grant Assumptions* 2.1 *and* 3.2(*i*). *Then,*

$$\widetilde{\mathcal{H}}^{n}\phi^{n}[f](\tilde{x},h) = \sum_{i\in\mathcal{I}} \frac{\lambda_{i}^{n}}{\sqrt{n}} \partial_{i} f(\tilde{x}) + \sum_{i\in\mathcal{I}} \frac{\lambda_{i}^{n} c_{a,i}^{2}}{2n} \partial_{ii} f(\tilde{x})$$

$$+ \sum_{i\in\mathcal{I}} \frac{\lambda_{i}^{n}}{n} \sum_{j\in\mathcal{I}} \left(\eta_{j}^{n}(h_{j}) + \frac{c_{a,j}^{2} - 1}{2} \right) \partial_{ij} f(\tilde{x}) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$$
(5.14)

for all $f \in \mathcal{C}_c^{\infty}(\mathbb{R}^d)$ and $(\tilde{x}, h) \in \mathbb{R}^d \times \mathbb{R}^d_+$.

Proof Note that

$$\begin{split} \hat{\phi}_{1,1}^{n}[f] &= \sum_{i \in \mathcal{I}} \eta_{i}^{n}(h_{i}) \big(f(\tilde{x} + n^{-1/2}e_{i}) - f(\tilde{x}) \big), \\ \hat{\phi}_{2,1}^{n}[f] &= \sum_{i \in \mathcal{I}} \eta_{i}^{n}(h_{i}) \sum_{j \in \mathcal{I}} \frac{c_{a,j}^{2} - 1}{2\sqrt{n}} \big(\partial_{j} f(\tilde{x} + n^{-1/2}e_{i_{1}}) - \partial_{j} f(\tilde{x}) \big). \end{split}$$

Using (4.4) and (4.6), and the Taylor expansion, we have

$$\begin{split} \widehat{\mathcal{H}}^n \Big(f + \widehat{\phi}_{1,1}^n[f] + \sum_{j \in \mathbb{J}} \frac{c_{a,j}^2 - 1}{2\sqrt{n}} \partial_j f + \widehat{\phi}_{2,1}^n[f] + \sum_{j \in \mathbb{J}} \frac{\kappa_j^n(h_j)}{n} \partial_{jj} f \Big) (\tilde{x}, h) \\ &= \sum_{i \in \mathbb{J}} \frac{\lambda_i^n}{\sqrt{n}} \partial_i f(\tilde{x}) + \sum_{i \in \mathbb{J}} \frac{\lambda_i^n c_{a,i}^2}{2n} \partial_{ii} f(\tilde{x}) + \sum_{i \in \mathbb{J}} \frac{\lambda_i^n}{n} \sum_{j \neq i} \frac{c_{a,j}^2 - 1}{2} \partial_{ij} f(\tilde{x}) + \mathcal{O}\bigg(\frac{1}{\sqrt{n}}\bigg) \end{split}$$



$$+ \sum_{i \in \mathcal{I}} r_{i}^{n}(h_{i}) \sum_{j \neq i} \eta_{j}^{n}(h_{j}) \Big([f]_{ij}^{1,n}(\tilde{x}) + [f]_{ij}^{2,n}(\tilde{x}) \Big)$$

$$+ \sum_{i \in \mathcal{I}} \frac{\lambda_{i}^{n}}{n} \Big(\eta_{i}^{n}(h_{i}) + \frac{c_{a,i}^{2} - 1}{2} \Big) \partial_{ii} f(\tilde{x}) + \sum_{i \in \mathcal{I}} r_{i}^{n}(h_{i}) \sum_{j \neq i} \frac{\kappa_{j}^{n}(h_{j})}{n} [f]_{ij}^{3,n}(\tilde{x}).$$
(5.15)

It is straightforward to verify that

$$\widehat{\mathcal{H}}^{n}(\widehat{\phi}_{1,2}^{n}[f] + \widehat{\phi}_{2,2}^{n}[f] + \widehat{\phi}_{3,1}^{n}[f])(\tilde{x}, h) \\
= \sum_{i \in \mathcal{I}} (\dot{\eta}_{i}^{n}(h_{i}) - \eta_{i}^{n}(h_{i})r_{i}^{n}(h_{i})) \sum_{j \neq i} \eta_{j}^{n}(h_{j}) \Big([f]_{ij}^{1,n}(\tilde{x}) + [f]_{ij}^{2,n}(\tilde{x}) \Big) \\
+ \frac{1}{2} \sum_{i \in \mathcal{I}} r_{i}^{n}(h_{i}) \sum_{j \neq i} \sum_{k \neq i, j} \eta_{j}^{n}(h_{j}) \eta_{k}^{n}(h_{k}) \Big([f]_{ijk}^{1,n}(\tilde{x}) + [f]_{ijk}^{2,n}(\tilde{x}) \Big) \\
+ \sum_{i \in \mathcal{I}} \Big(\Big(\dot{\eta}_{i}^{n}(h_{i}) - \eta_{i}^{n}(h_{i})r_{i}^{n}(h_{i}) \Big) \sum_{j \neq i} \frac{\kappa_{j}^{n}(h_{j})}{n} + \Big(\dot{\kappa}_{i}^{n} - r_{i}^{n}(h_{i})\kappa_{i}^{n}(h_{i}) \Big) \sum_{j \neq i} \frac{\eta_{j}^{n}(h_{j})}{n} \Big) [f]_{ij}^{3,n}(\tilde{x}) \\
+ \sum_{i \in \mathcal{I}} r_{i}^{n}(h_{i}) \sum_{j \neq i} \eta_{j}^{n}(h_{j}) \sum_{k \neq i, j} \frac{\kappa_{k}^{n}(h_{k})}{n} [f]_{ijk}^{3,n}(\tilde{x}) \tag{5.16}$$

for any $(\tilde{x}, h) \in \mathbb{R}^d \times \mathbb{R}^d_+$. Applying (4.4) and (4.6), and combining the first term on the right-hand side of (5.16) with the third, fifth and sixth terms on the right-hand side of (5.15), we obtain the third term on the right-hand side of (5.14). We repeat this procedure until all the terms r_i^n are canceled. This proves (5.14).

Definition 5.1 We define the operator $\hat{A}^n : \mathcal{C}^2(\mathbb{R}^d) \to \mathcal{C}^2(\mathbb{R}^d \times \mathcal{S})$ by

$$\hat{\mathcal{A}}^n f(x, u) := \sum_{i \in \mathcal{I}} \left(\mathcal{A}^n_{1,i}(x, u) \partial_i f(x) + \frac{1}{2} \mathcal{A}^n_{2,i}(x, u) \partial_{ii} f(x) \right),$$

where $\mathcal{A}_{1,i}^n$, $\mathcal{A}_{2,i}^n$: $\mathbb{R}^d \times \mathcal{S} \to \mathbb{R}$, $i \in \mathcal{I}$, are given by

$$\begin{split} \mathcal{A}^{n}_{1,i}(x,u) &:= \ell^{n}_{i} - \mu^{n}_{i}(x_{i} - \langle e, x \rangle^{+} u_{i}) - \gamma^{n}_{i} \langle e, x \rangle^{+} u_{i}, \\ \mathcal{A}^{n}_{2,i}(x,u) &:= \frac{\lambda^{n}_{i}}{n} c_{a,i}^{2} + \rho_{i} \mu^{n}_{i} + \frac{\mu^{n}_{i}(x_{i} - \langle e, x \rangle^{+} u_{i}) + \gamma^{n}_{i} \langle e, x \rangle^{+} u_{i}}{\sqrt{n}}, \end{split}$$

respectively. Define the operator \hat{I}^n by

$$\hat{\mathcal{I}}^n f(x) := \int_{\mathbb{R}^d} \left(f(x+y) - f(x) \right) \nu_{d_1}^n(\mathrm{d}y),$$

where

$$\nu_{d_1}^n(A) := \Pi_{d_1}^n \Big(\big\{ y \in \mathbb{R}_* \colon \left(\frac{\sqrt{n}}{\vartheta^n} \mu_1^n \rho_1 y, \dots, \frac{\sqrt{n}}{\vartheta^n} \mu_d^n \rho_d y \right) \in A \big\} \Big),$$

with $\Pi_{d_1}^n(\mathrm{d}y) := \beta_{\mathsf{u}}^n F^{d_1}(\mathrm{d}y)$, and β_{u}^n as in Assumption 2.2.



Recall the generator $\widetilde{\mathcal{L}}_n^{z^n}$ of $\widetilde{\Xi}^n$ given in (4.16). The next lemma establishes the relation between the generator of the diffusion-scaled process and the operator in Definition 5.1.

Lemma 5.4 Grant Assumptions 2.1, 2.2, and 3.2. Then,

$$\widetilde{\mathcal{L}}_{n}^{z^{n}} \phi^{n}[f](\widetilde{x}, h, \psi, k) = \widehat{\mathcal{A}}^{n} f(\widetilde{x}, v^{n}(\widetilde{x}, h, \psi, k)) + \widehat{\mathcal{I}}^{n} f(\widetilde{x})
+ \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) (\|\widetilde{x}\|
+ \|\widetilde{q}^{n}\|) + \mathcal{O}(1)(1 - \psi) (\|\widetilde{x}\| + \|\widetilde{q}^{n}\| + 1),$$
(5.17)

for any $f \in \mathcal{C}_c^{\infty}(\mathbb{R}^d)$ and $z^n \in \mathfrak{Z}_{sm}^n$, where $\tilde{q}^n = n^{-1/2}q^n$, and

$$v^{n}(\tilde{x}, h, \psi, k) = \begin{cases} \frac{\tilde{x} - \tilde{z}^{n}(\sqrt{n}\tilde{x} + n\rho, h, \psi, k)}{\langle e, \tilde{x} \rangle}, & \text{if } \langle e, \tilde{x} \rangle > 0, \\ e_{d}, & \text{if } \langle e, \tilde{x} \rangle \leq 0, \end{cases}$$
(5.18)

for $(\tilde{x}, h, \psi, k) \in \widetilde{\mathfrak{D}}^n$, with $\tilde{z}^n := n^{-1/2}(z^n - n\rho)$.

Proof Note that Lemma 5.3 concerns the renewal arrival process in the diffusion-scale. Recall that $z_i^n = \sqrt{n}(\tilde{x}_i - \tilde{q}_i^n) + n\rho_i$ for $i \in \mathbb{J}$, and $\check{x} = \sqrt{n}\tilde{x} + n\rho$. We let $q^n \equiv q^n(\sqrt{n}\tilde{x} + n\rho, z^n)$ and $z^n \equiv z^n(\sqrt{n}\tilde{x} + n\rho, h, \psi, k)$. Applying Lemma 5.3 and the Taylor expansion, it follows by the definition of $\widetilde{\mathcal{L}}_n^{z^n}$ that

$$\widetilde{\mathcal{L}}_{n}^{z^{n}} \phi^{n}[f](\tilde{x}, h, \psi, k) = \sum_{i \in \mathcal{I}} \left[\left(\frac{(\lambda_{i}^{n} - n\rho_{i}\mu_{i}^{n})}{\sqrt{n}} - \mu_{i}^{n}(\tilde{x}_{i} - \tilde{q}_{i}^{n}) - \gamma_{i}^{n}\tilde{q}_{i}^{n} \right) \partial_{i} f(\tilde{x}) \right. \\
+ \frac{1}{2} \left(\frac{\lambda_{i}^{n} c_{a,i}^{2}}{n} + \rho_{i}\mu_{i}^{n} + \frac{\tilde{x}_{i} + (\mu_{i}^{n} - \gamma_{i}^{n})\tilde{q}_{i}^{n}}{\sqrt{n}} \right) \partial_{ii} f(\tilde{x}) \\
+ \frac{\lambda_{i}^{n} - n\rho_{i}\mu_{i}^{n}}{n} \sum_{j \in \mathcal{I}} \left(\eta_{j}^{n}(h_{j}) + \frac{c_{a,j}^{2} - 1}{2} \right) \partial_{ij} f(\tilde{x}) \\
+ (1 - \psi)\gamma_{i}^{n} \left(\phi^{n}[f](\tilde{x} - n^{-1/2}e_{i}, h) - \phi^{n}[f](\tilde{x}, h) \right) \\
\int_{\mathbb{R}_{*}} q_{i}^{n} \left(\sqrt{n}\tilde{x} + n\rho - n\mu^{n}(y - k), z^{n} \right) \tilde{F}_{\tilde{x},k}^{d_{1}^{n}} (dy) \\
+ (\psi - 1)(\mu_{i}^{n} z_{i}^{n} + \gamma_{i}^{n} q_{i}^{n}) \left(\phi^{n}[f](\tilde{x} - n^{-1/2}e_{i}, h) - \phi^{n}[f](\tilde{x}, h) \right) \\
- (1 - \psi)\sqrt{n}\mu_{i}^{n}\rho_{i} \frac{\partial \phi^{n}[f](\tilde{x}, h)}{\partial \tilde{x}_{i}} \right] + \psi \, \hat{\mathcal{I}}^{n}\phi^{n}[f](\tilde{x}, h) \\
+ \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) (\|\tilde{x}\| + \|\tilde{q}^{n}\|) \tag{5.19}$$

for any $f \in \mathcal{C}_c^{\infty}(\mathbb{R}^d)$, where

$$\hat{\mathcal{I}}^n \phi^n[f](\tilde{x}, h) = \int_{\mathbb{R}^d} \left(\phi^n[f](\tilde{x} + y, h) - \phi^n[f](\tilde{x}, h) \right) \nu_{d_1}^n(\mathrm{d}y)$$



by a slight abuse of notation. It is clear that

$$\lambda_i^n - n\mu_i^n \rho_i = \mathcal{O}(\sqrt{n}) \tag{5.20}$$

by Assumption 2.1, and thus the third term in the sum on the right-hand side of (5.19) is of order $n^{-1/2}$. We next consider the fifth and sixth terms in the sum on the right-hand side of (5.19). Using the fact that

$$\phi^{n}[f](\tilde{x}-n^{-1/2}e_{i},h)-\phi^{n}[f](\tilde{x},h)=-\frac{1}{\sqrt{n}}\frac{\partial\phi^{n}[f](\tilde{x},h)}{\partial\tilde{x}_{i}}+\mathcal{O}\left(\frac{1}{n}\right),$$

and $z_i^n = \sqrt{n}\tilde{x}_i + n\rho_i - \sqrt{n}\tilde{q}_i^n$, we obtain

$$\begin{split} &(\psi-1)(\mu_i^n z_i^n + \gamma_i^n q_i^n) \big(\phi^n [f] (\tilde{x} - n^{-1/2} e_i, h) - \phi^n [f] (\tilde{x}, h) \big) \\ &- (1 - \psi) \sqrt{n} \mu_i^n \rho_i \frac{\partial \phi^n [f] (x, h)}{\partial \tilde{x}_i} \\ &= (\psi - 1) \big(\mu_i^n \tilde{x}_i + (\mu_i^n - \gamma_i^n) \tilde{q}_i^n \big) \bigg(- \frac{\partial \phi^n [f] (\tilde{x}, h)}{\partial x_i} + \mathcal{O} \bigg(\frac{1}{\sqrt{n}} \bigg) \bigg). \end{split}$$

Recall the definition of $\tilde{F}_{\tilde{x}k}^{d_1^n}$ in (4.13). Note that

$$\int_{\mathbb{R}_{*}} n\mu_{i}^{n} \rho_{i}(y-k) \, \tilde{F}_{\check{x},k}^{d_{1}^{n}}(\mathrm{d}y) \, \leq \, \frac{n}{\vartheta^{n}} \mu_{i}^{n} \rho_{i} \, \mathbb{E} \left[d_{1} - \vartheta^{n}k \, | \, d_{1} > \vartheta^{n}k \right] \in \, \mathfrak{O}(\sqrt{n}), \tag{5.21}$$

where the second equality follows by Assumption 2.2 and (3.10). Note that $\tilde{q}_i^n \leq \langle e, \tilde{x} \rangle^+$ for $i \in \mathcal{I}$ and $(\tilde{x}, h, \psi, k) \in \widetilde{\mathfrak{D}}^n$. Thus, the fourth term in the sum on the right-hand side of (5.19) is bounded by $C(1 - \psi)(1 + \langle e, \tilde{x} \rangle^+)$ for some positive constant C. It is evident that $\phi^n[f] - f \in \mathfrak{O}(n^{-1/2})$, and

$$\psi \,\hat{\mathcal{I}}^n \phi^n[f](\tilde{x},h) = \hat{\mathcal{I}}^n f(\tilde{x}) + (\psi - 1) \,\hat{\mathcal{I}}^n f(\tilde{x}) + \psi \,\hat{\mathcal{I}}^n(\phi^n[f] - f)(\tilde{x},h).$$

Therefore, (5.17) follows by the boundedness of $\phi^n[f]$ and (5.19). This completes the proof.

Definition 5.2 The mean empirical measure $\hat{\zeta}_T^{z^n} \in \mathcal{P}(\mathbb{R}^d \times \mathcal{S})$ associated with \hat{X}^n and a stationary Markov policy $z^n \in \mathfrak{Z}_{sm}^n$ is defined by

$$\hat{\xi}_{T}^{z^{n}}(A \times B) := \frac{1}{T} \mathbb{E} \left[\int_{0}^{T} \mathbb{1}_{A \times B} (\hat{X}^{n}(s), v^{n} (\hat{X}^{n}(s), H^{n}(s), \Psi^{n}(s), K^{n}(s))) ds \right]$$

for any Borel sets $A \subset \mathbb{R}^d$ and $B \subset \mathcal{S}$, and with v^n as in (5.18).

The following theorem characterizes the limit points of mean empirical measures.



Theorem 5.2 Grant the hypotheses in Theorem 3.2. Let $\{z^n \in \mathfrak{Z}^n_{sm} : n \in \mathbb{N}\}\$ be a sequence of policies satisfying (5.11). Then any limit point $\pi \in \mathcal{P}(\mathbb{R}^d \times \mathcal{S})$ of $\hat{\zeta}_T^{z^n}$ as $(n, T) \to \infty$ lies in \mathfrak{G} .

Proof It follows directly by Assumptions 2.1 and 2.2 that, for any $f \in \mathcal{C}_c^{\infty}(\mathbb{R}^d)$, we have

$$\hat{\mathcal{A}}^n f(\hat{x}, u) + \hat{\mathcal{I}}^n f(\hat{x}) \to \mathcal{A} f(\hat{x}, u) \text{ as } n \to \infty$$
 (5.22)

uniformly over compact sets of $\mathbb{R}^d \times \mathcal{S}$. Thus, in view of (5.22) and (5.1), in order to prove the theorem, it is enough to show that

$$\lim_{(n,T)\to\infty} \int_{\mathbb{R}^d\times\mathcal{S}} \left(\hat{\mathcal{A}}^n f(\hat{x},u) + \hat{\mathcal{I}}^n f(\hat{x})\right) \hat{\zeta}_T^{z^n} (d\hat{x}, du) = 0 \quad \forall f \in \mathcal{C}_c^{\infty}(\mathbb{R}^d).$$
 (5.23)

Applying (5.11 and (4.26)), we obtain

$$\sup_{n>n_{\circ}} \limsup_{T\to\infty} \frac{1}{T} \mathbb{E}^{z^{n}} \left[\int_{0}^{T} |\widetilde{X}^{n}(s)|^{m} ds \right] < \infty.$$
 (5.24)

It follows by the same calculation as in (5.6) that, for some positive constant C_1 , we have

$$\mathbb{E}^{z^n} \left[\int_0^T \sqrt{n} (1 - \Psi^n(s)) \, \mathrm{d}s \right] \le C_1 (1 + T) \quad \forall \, T \ge 0.$$
 (5.25)

Using the facts that $\tilde{q}_i^n \leq \langle e, x \rangle^+$ and $\Psi^n(s) \in \{0, 1\}$, and Young's inequality, we obtain

$$\frac{1}{T} \mathbb{E}^{z^{n}} \left[\int_{0}^{T} n^{\frac{m-1}{4m}} (1 - \Psi^{n}(s)) n^{\frac{1-m}{4m}} (\|\widetilde{X}^{n}(s)\| + \|\widetilde{q}^{n}(\sqrt{n}\widetilde{X}^{n}(s) + n\rho, z^{n})\|) ds \right] \\
\leq \frac{1}{T} \mathbb{E}^{z^{n}} \left[\int_{0}^{T} n^{\frac{1}{4}} (1 - \Psi^{n}(s)) ds \right] + \frac{C_{2}}{T} \mathbb{E}^{z^{n}} \left[\int_{0}^{T} n^{\frac{1-m}{4}} |\widetilde{X}^{n}(s)|^{m} ds \right] \\
\leq \frac{1}{Tn^{\frac{1}{4}}} C_{1} (1 + T) + n^{\frac{1-m}{4}} \frac{C_{2}}{T} \mathbb{E}^{z^{n}} \left[\int_{0}^{T} |\widetilde{X}^{n}(s)|^{m} ds \right] \longrightarrow 0 \\
\text{as } (n, T) \to \infty, \tag{5.26}$$

where C_2 is a positive constant. In (5.26), the second inequality follows by (5.25), and the convergence follows by (5.24) and the fact that m > 1. Applying Itô's formula to $\phi^n[f]$, and using Lemma 5.4 and (5.24) and (5.26), it follows by the boundedness of $\phi^n[f]$ that

$$\lim_{(n,T)\to\infty}\frac{1}{T}\mathbb{E}^{z^n}\left[\int_0^T\hat{\mathcal{A}}^nf\big(\widetilde{X}^n(s),v^n\big(\widetilde{\Xi}^n(s)\big)\big)+\hat{\mathcal{I}}^nf\big(\widetilde{X}^n(s)\big)\,\mathrm{d}s\right]=0.$$

Therefore, using (4.26) again, we obtain (5.23). This completes the proof.



Proof of (5.10) Without loss of generality, suppose $\{n_j\} \subset \mathbb{N}$ is an increasing sequence such that $z^{n_j} \in \mathfrak{Z}_{sm}$ and $\sup_j \hat{J}(\hat{X}^{n_j}(0), z^{n_j}) < \infty$. Recall $\hat{\zeta}_T^{z^n}$ in Definition 5.2. There exists a subsequence of $\{n_j\}$, denoted as $\{n_l\}$, such that $T_l \to \infty$ as $l \to \infty$, and

$$\liminf_{j \to \infty} \hat{J}(\hat{X}^{n_j}(0), z^{n_j}) + \frac{1}{l} \ge \int_{\mathbb{R}^d \times \mathbb{U}} \mathcal{R}(\hat{x}, u) \, \hat{\zeta}_{T_l}^{z^{n_l}}(d\hat{x}, du). \tag{5.27}$$

Applying Lemma 5.2 and Theorem 5.2, any limit of $\hat{\zeta}_{T_l}^{z^{n_l}}$ along some subsequence is in \mathcal{G} . Choose any further subsequence of (T_l, n_l) , also denoted by (T_l, n_l) , such that $(T_l, n_l) \to \infty$ as $l \to \infty$, and $\hat{\zeta}_{T_l}^{z^{n_l}} \to \pi \in \mathcal{G}$. By letting $l \to \infty$ and using (5.27), we obtain

$$\liminf_{j\to\infty} \hat{J}(\hat{X}^{n_j}(0), z^{n_j}) \geq \int_{\mathbb{R}^d \times \mathbb{U}} \Re(\hat{x}, u) \, \pi(d\hat{x}, du) \geq \varrho_*.$$

This completes the proof.

5.3.2 The Upper Bound

In this subsection, we show that

$$\limsup_{n \to \infty} \varrho^n (\hat{X}^n(0)) \le \varrho_*. \tag{5.28}$$

The following lemma concerns the convergence of mean empirical measures for the diffusion-scaled state processes under the scheduling policies in Definition 4.3. Recall \mathfrak{A}_R^n in Definition 4.2 and $\hat{\zeta}_T^{z^n}$ in Definition 5.2.

Lemma 5.5 Grant the hypotheses in Theorem 3.2. For $\epsilon > 0$, let v_{ϵ} be a continuous ϵ -optimal precise control, whose existence is asserted in Proposition 5.1, and $\{z^n[v^n]: n \in \mathbb{N}\}$ be as in Definition 4.3, and such that $R \equiv R(\epsilon)$ and v^n agrees with v_{ϵ} on \mathfrak{A}_R^n . Then, the ergodic occupation measure $\pi_{v_{\epsilon}}$ of the controlled jump diffusion in (3.3) under the control v_{ϵ} is the unique limit point in $\mathcal{P}(\mathbb{R}^d \times \mathcal{S})$ of $\hat{\zeta}_T^{z^n[v^n]}$ as $(n, T) \to \infty$.

Proof Using Proposition 4.3 and Theorem 5.2, the proof of this lemma is the same as that of Lemma 7.2 in [5].

Proof of (5.28) Let $\kappa = 2\lfloor m \rfloor$ with m as in (3.5), and $z^n[v^n]$ be the scheduling policy in Lemma 5.5. By Proposition 4.3, there exist $\tilde{n}_{\circ} \in \mathbb{N}$, and positive constants \tilde{C}_0 and \tilde{C}_1 such that

$$\widetilde{\mathcal{L}}_{n}^{z^{n}[v^{n}]} \widetilde{\mathcal{V}}_{\kappa,\xi}^{n}(\tilde{x},h,\psi,k) \leq \widetilde{C}_{0} - \widetilde{C}_{1} \mathcal{V}_{\kappa-1,\xi}(\tilde{x}) \qquad \forall (\tilde{x},h,\psi,k) \in \widetilde{\mathfrak{D}}^{n},
\text{and for all } n > \widetilde{n}_{\circ}.$$
(5.29)

Recall the definition of $\widetilde{\mathbb{R}}$ in (3.5), and let $\hat{z}^n[v^n] = n^{-1/2}(z^n[v^n] - n\rho)$. Applying (4.26) and (5.29), we may select an increasing sequence T_n such that

$$\sup_{n\geq \tilde{n}_o} \sup_{T\geq T_n} \int_{\mathbb{R}^d \times \mathbb{U}} \mathcal{V}_{\kappa-1,\xi}(\hat{x}) \, \hat{\zeta}_T^{z^n[v^n]}(\mathrm{d}\hat{x},\mathrm{d}u) < \infty.$$



This implies that $\widetilde{\mathbb{R}}(\hat{x} - \hat{z}^n[v](\sqrt{n}\hat{x} + n\rho))$ is uniformly integrable. By Lemma 5.5, $\hat{\zeta}_T^{z^n[v^n]}$ converges in $\mathcal{P}(\mathbb{R}^d \times \mathcal{S})$ to $\pi_{v_{\epsilon}}$ as $(n,T) \to \infty$. Applying Proposition 5.1, we deduce that v_{ϵ} is an ϵ -optimal control for the running cost function. Since ϵ is arbitrary, (5.28) follows.

Acknowledgements This research was supported in part by the Army Research Office through Grant W911NF-17-1-001, in part by the National Science Foundation through Grants DMS-1715210, CMMI-1635410 and DMS-1715875, and in part by the Office of Naval Research through Grant N00014-16-1-2956 and was approved for public release under DCN #43-5442-19.

Appendix A. Proofs of Lemma 3.1 and Proposition 3.1

Proof of Lemma 3.1 By [13, Lemma 5.1], $\hat{S}_i^n(t)$ and $\hat{R}_i^n(t)$ in (3.1) are martingales with respect to the filtration \mathcal{F}_t^n in (2.9), having predictable quadratic variation processes given by

$$\langle \hat{S}_i^n \rangle(t) = \mu_i^n \int_0^t n^{-1} Z_i^n(s) \Psi^n(s) \, \mathrm{d}s \quad \text{and} \quad \langle \hat{R}_i^n \rangle(t) = \gamma_i^n \int_0^t n^{-1} Q_i^n(s) \, \mathrm{d}s, \quad t \ge 0,$$

respectively. By (2.7), we have the crude inequality

$$0 \le n^{-1}X_i^n(t) \le n^{-1}X_i^n(0) + n^{-1}A_i^n(t), \quad t \ge 0.$$

Using the balance equation in (2.5), we see that the same inequalities hold for $n^{-1}Z_i^n$ and $n^{-1}Q_i^n$. Since $\Psi^n(s) \in \{0, 1\}$, it follows by Lemma 5.8 in [30] that $\{\hat{W}_i^n : n \in \mathbb{N}\}$ is stochastically bounded in (\mathbb{D}^d, J_1) . Also, $\{\hat{L}_i^n : n \in \mathbb{N}\}$ is stochastically bounded in (\mathbb{D}^d, M_1) by (2.4). On the other hand, it is evident that

$$\hat{Y}_i^n(t) \le C \int_0^t (1 + ||n^{-1}X^n(s)||) \, \mathrm{d}s, \quad t \ge 0,$$

where C is some positive constant. Thus, we obtain

$$\|\hat{X}^{n}(t)\| \leq \|\hat{X}^{n}(0)\| + \|\hat{W}^{n}(t)\| + \|\hat{L}^{n}(t)\| + C \int_{0}^{t} (1 + \|\hat{X}^{n}(s)\|) \, \mathrm{d}s \quad \forall \, t \geq 0. \tag{A.1}$$

Since $\hat{X}^n(0)$ is uniformly bounded, applying Lemma 5.3 in [30] and Gronwall's inequality, we deduce that $\{\hat{X}^n : n \in \mathbb{N}\}$ is stochastically bounded in (\mathbb{D}^d, M_1) . Using Lemma 5.9 in [30], we see that

$$n^{-1/2}\hat{X}^n = n^{-1}X^n - \rho \implies \mathfrak{e}_0 \text{ in } (\mathbb{D}^d, M_1) \text{ as } n \to \infty,$$

which implies that $n^{-1}X^n \Rightarrow \mathfrak{e}_{\rho}$ in (\mathbb{D}^d, M_1) . By (2.5), and the fact $\langle e, n^{-1}Q^n \rangle = (\langle e, n^{-1}X^n \rangle - 1)^+ \Rightarrow \mathfrak{e}_0$, we have $n^{-1}Q^n \Rightarrow \mathfrak{e}_0$, and thus $n^{-1}Z^n \Rightarrow \mathfrak{e}_{\rho}$. This completes the proof.



To prove Proposition 3.1, we first consider a modified process. Let $\check{X}^n = (\check{X}^n_1, \ldots, \check{X}^n_d)'$ be the *d*-dimensional process defined by

$$\check{X}_{i}^{n}(t) := \hat{X}^{n}(0) + \ell_{i}^{n}t + \hat{W}_{i}^{n}(t) + \hat{L}_{i}^{n}(t)
- \int_{0}^{t} \mu_{i}^{n} (\check{X}_{i}^{n}(s) - \langle e, \check{X}^{n}(s) \rangle^{+} U_{i}^{n}(s)) ds
- \int_{0}^{t} \gamma_{i}^{n} \langle e, \check{X}^{n}(s) \rangle^{+} U_{i}^{n}(s) ds, \text{ for } i \in \mathcal{I}.$$
(A.2)

Lemma A.1 As $n \to \infty$, \check{X}^n and \hat{X}^n are asymptotically equivalent, that is, if either of them converges in distribution as $n \to \infty$, then so does the other, and both of them have the same limit.

Proof Let $K = K(\epsilon_1) > 0$ be the constant satisfying $\mathbb{P}(\|\hat{X}^n\|_T > K) < \epsilon_1$ for T > 0 and any $\epsilon_1 > 0$, where $\|\hat{X}^n\|_T := \sup_{0 \le t \le T} \|\hat{X}^n(t)\|$. Since $\hat{U}^n(s) \in \mathcal{S}$ for $s \ge 0$, on the event $\{\|\hat{X}^n\|_T \le K\}$, we obtain

$$\|\check{X}^{n}(t) - \hat{X}^{n}(t)\| \leq C_{1} \int_{0}^{t} \|\hat{X}^{n}(s)\| (1 - \Psi^{n}(s)) \, \mathrm{d}s + C_{2} \int_{0}^{t} \|\check{X}^{n}(s) - \hat{X}^{n}(s)\| \, \mathrm{d}s$$

$$\leq C_{1} K C_{\mathsf{d}}^{n}(t) + C_{2} \int_{0}^{t} \|\check{X}^{n}(s) - \hat{X}^{n}(s)\| \, \mathrm{d}s \quad \forall t \in [0, T],$$

where C_1 and C_2 are some positive constants. Then, by Gronwall's inequality, on the event $\{\|\hat{X}^n\|_T \leq K\}$, we have

$$\|\dot{X}^n(t) - \hat{X}^n(t)\| \le C_1 K C_d^n(t) e^{C_2 T} \quad \forall t \in [0, T].$$

Thus, applying [13, Lemma 2.2], we deduce that for any $\epsilon_2 > 0$, there exist $\epsilon_3 > 0$ and $n_0 = n_0(\epsilon_1, \epsilon_2, \epsilon_3, T)$ such that

$$\|\check{X}^n - \hat{X}^n\|_T \le \epsilon_2$$

on the event $\{\|\hat{X}^n\|_T \leq K\} \cap \{\|C^n_d\|_T \leq \epsilon_3\}$, for all $n \geq n_\circ$, which implies that

$$\mathbb{P}(\|\check{X}^n - \hat{X}^n\|_T > \epsilon_2) < \epsilon_1, \quad \forall n \ge n_\circ.$$

As a consequence, $\|\check{X}^n - \hat{X}^n\|_T \Rightarrow 0$, as $n \to \infty$, and this completes the proof.

Proof of Proposition 3.1 We first prove (i). Define the processes

$$\tau_{1,i}^n(t) := \frac{\mu_i^n}{n} \int_0^t Z^n(s) \Psi^n(s) \, \mathrm{d}s, \quad \tau_{2,i}^n(t) := \frac{\gamma_i^n}{n} \int_0^t Q^n(s) \, \mathrm{d}s,$$

 $\tilde{S}_i^n(t) := n^{-1/2}(S^n(nt) - nt)$, and $\tilde{R}_i^n(t) := n^{-1/2}(R^n(nt) - nt)$, for $i \in \mathcal{I}$. Then, since $\Psi^n(s) \in \{0,1\}$ for $s \geq 0$, applying Lemma 3.1 and Lemma 2.2 in [13], we have



$$\tau_{1,i}^n(\cdot) = \mu_i^n \int_0^{\cdot} (n^{-1} Z_i^n(s) - \rho_i) \Psi^n(s) \, \mathrm{d}s + \mu_i^n \int_0^{\cdot} \rho_i \Psi^n(s) \, \mathrm{d}s \implies \lambda_i \mathfrak{e}(\cdot).$$

in (\mathbb{D}, M_1) , as $n \to \infty$, and that $\tau^n_{2,i}$ weakly converges to the zero process. Since $\{A^n_i, S^n_i, R^n_i, \Psi^n \colon i \in \mathcal{I}, n \in \mathbb{N}\}$ are independent processes, and $\tau^n_{1,i}$ and $\tau^n_{2,i}$ converge to deterministic functions, we have joint weak convergence of $(\hat{A}^n, \hat{S}^n, \hat{R}^n, \hat{L}^n, \tau^n_1, \tau^n_2)$, where $\tau^n_1 := (\tau^n_{1,1}, \dots, \tau^n_{1,d})'$, and τ^n_2 is defined analogously. On the other hand, since the second moment of A^n is finite, it follows that \hat{A}^n converges weakly to a d-dimensional Wiener process with mean 0 and covariance matrix $\mathrm{diag}(\sqrt{\lambda_1 c^2_{a,1}}, \dots, \sqrt{\lambda_d c^2_{a,d}})$ (see, e.g., [31]). Therefore, by the FCLT for the Poisson processes \tilde{S}^n and \tilde{R}^n , and using the random time change lemma in [21, Page 151], we obtain (i).

Using (A.1) and Proposition 3.1(i), the proof of (ii) is same as the proof of [1, Lemma 4 (iii)].

To prove (iii), we first show any limit of \check{X}^n in (A.2) satisfies (3.3). Following an argument similar to the proof of Lemma 5.2 in [13], one can easily show that the d-dimensional integral mapping $x = \Lambda(y, u) \colon \mathbb{D}^d \times \mathbb{D}^d \to \mathbb{D}^d$ defined by

$$x(t) = y(t) + \int_0^t h(x(s), u(s)) ds$$

is continuous in (\mathbb{D}^d, M_1) , provided that the function $h: \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}^d$ is Lipschitz continuous in each coordinate. Since

$$\check{X}^n = \Lambda(\hat{X}^n(0) + \hat{W}^n + \hat{L}^n, U^n),$$

then, by the tightness of U^n and the continuous mapping theorem, any limit of \check{X}^n satisfies (A.2), and the same result holds for \hat{X}^n by Lemma A.1.

Recall the definition of $\check{\tau}^n$ in (2.8). It is evident that

$$\hat{L}_{i}^{n}(t+r) - \hat{L}_{i}^{n}(t) = \hat{L}_{i}^{n}(\check{\tau}^{n}(t)+r) - \hat{L}_{i}^{n}(\check{\tau}^{n}(t))
+ \hat{L}_{i}^{n}(t+r) - \hat{L}_{i}^{n}(\check{\tau}^{n}(t)+r) + \hat{L}_{i}^{n}(\check{\tau}^{n}(t)) - \hat{L}_{i}^{n}(t).$$
(A.3)

for all $t, r \ge 0$ and $i \in \mathcal{I}$. By Assumption 2.2, we have $\check{\tau}^n(t) \Rightarrow t$ as $n \to \infty$, for $t \ge 0$. Then, by the random time change lemma in [21, Page 151], we deduce that the last four terms on the right-hand side of (A.3) converge to 0 in distribution. It follows by Proposition 3.1 (i) and (A.3) that

$$\hat{L}^n(\check{\tau}^n(t)+r) - \hat{L}^n(\check{\tau}^n(t)) \Rightarrow \lambda L_{t+r} - \lambda L_t \text{ in } \mathbb{R}^d.$$

Repeating the same argument we establish convergence of \hat{S}^n and \hat{R}^n . Proving that U is non-anticipative follows exactly as in [1]*Lemma 6. This completes the proof of (iii).



B Proofs of Lemmas 4.1 and 5.2

In this section, we construct two functions, which are used to show the ergodicity of $\tilde{\Xi}^n$. We provide two lemmas concerning the properties of these functions, respectively. The proofs of Lemmas 4.1 and 5.2 are given at the end of this section.

Definition B.1 For $z^n \in \mathfrak{Z}^n_{sm}$, define the operator $\mathcal{L}_n^{z^n} : \mathcal{C}_b(\mathbb{R}^d \times \mathbb{R}^d) \to \mathcal{C}_b(\mathbb{R}^d \times \mathbb{R}^d)$ by

$$\mathcal{L}_{n}^{z^{n}} f(\check{x}, h) := \sum_{i \in \mathbb{J}} \frac{\partial f(\check{x}, h)}{\partial h_{i}} + \sum_{i \in \mathbb{J}} r_{i}^{n} (h_{i}) (f(\check{x} + e_{i}, h - h_{i} e_{i}) - f(\check{x}, h))$$

$$+ \sum_{i \in \mathbb{J}} \mu_{i}^{n} z_{i}^{n} (f(\check{x} - e_{i}, h) - f(\check{x}, h))$$

$$+ \sum_{i \in \mathbb{J}} \gamma_{i}^{n} q_{i}^{n} (f(\check{x} - e_{i}, h) - f(\check{x}, h))$$
(B.1)

for $f \in \mathcal{C}_b(\mathbb{R}^d \times \mathbb{R}^d)$ and any $(\check{x}, h) \in \mathbb{R}^d_+ \times \mathbb{R}^d_+$, with $q^n := \check{x} - z^n$.

Note that if $d_1^n \equiv 0$ for all n, the queueing system has no interruptions. In this situation, under a Markov scheduling policy, the (infinitesimal) generator of (X^n, H^n) takes the form of (B.1). Recall the scheduling policies \check{z}^n in Definition 4.1, and $\bar{x} = \check{x} - n\rho$ in Definition 4.2. We define the sets

$$\tilde{\mathcal{K}}_n(\check{x}) := \left\{ i \in \mathcal{I}_0 \colon \check{x}_i \geq \frac{n\rho_i}{\sum_{i \in \mathcal{I}_0} \rho_i} \right\} = \left\{ i \in \mathcal{I}_0 \colon \bar{x}_i \geq \frac{n\rho_i \sum_{j \in \mathcal{I} \setminus \mathcal{I}_0} \rho_j}{\sum_{i \in \mathcal{I}_0} \rho_i} \right\}.$$

We have the following lemma.

Lemma B.1 Grant Assumptions 2.1, 2.2, and 3.2. For any even integer $\kappa \geq 2$, there exist a positive vector $\xi \in \mathbb{R}^d_+$, $\check{n} \in \mathbb{N}$, and positive constants \check{C}_0 and \check{C}_1 , such that the functions f_n , $n \in \mathbb{N}$, defined by

$$f_n(\check{x},h) := \sum_{i \in \mathcal{I}} \xi_i |\bar{x}_i|^{\kappa}$$

$$+ \sum_{i \in \mathcal{I}} \eta_i^n(h_i) \xi_i \left(|\bar{x}_i + 1|^{\kappa} - |\bar{x}_i|^{\kappa} \right) \quad \forall (\check{x},h) \in \mathbb{R}_+^d \times \mathbb{R}_+^d, \quad (B.2)$$

with η_i^n as defined in (4.3), satisfy

$$\mathcal{L}_{n}^{\check{z}^{n}} f_{n}(\check{x}, h) \leq \check{C}_{0} n^{\kappa/2} \\
- \check{C}_{1} \sum_{i \in \mathcal{I} \setminus \check{\mathcal{K}}_{n}(\check{x})} \xi_{i} |\bar{x}_{i}|^{\kappa} - \check{C}_{1} \sum_{i \in \check{\mathcal{K}}_{n}(\check{x})} \left(\mu_{i}^{n} (\check{z}_{i}^{n} - n\rho_{i}) + \gamma_{i}^{n} \check{q}_{i}^{n} \right) |\bar{x}_{i}|^{\kappa-1} \\
+ \sum_{i \in \mathcal{I}} \left(\mathcal{O}(\sqrt{n}) \mathcal{O}(|\bar{x}_{i}|^{\kappa-1}) + \mathcal{O}(n) \mathcal{O}(|\bar{x}_{i}|^{\kappa-2}) \right) \tag{B.3}$$

for all $n \geq \check{n}$ and $(\check{x}, h) \in \mathbb{R}^d_+ \times \mathbb{R}^d_+$.



Proof Using the estimate

$$(a \pm 1)^m - a^{\kappa} = \pm \kappa a^{\kappa - 1} + \mathcal{O}(a^{\kappa - 2}) \quad \forall a \in \mathbb{R},$$
 (B.4)

an easy calculation shows that

$$\mathcal{L}_{n}^{\check{z}^{n}} f_{n}(\check{x}, h) = \sum_{i \in \mathcal{I}} \dot{\eta}_{i}^{n}(h_{i}) \xi_{i} \left(|\bar{x}_{i} + 1|^{\kappa} - |\bar{x}_{i}|^{\kappa} \right) + \sum_{i \in \mathcal{I}} r_{i}^{n}(h_{i}) \eta_{i}^{n}(0) \xi_{i} \left((\bar{x}_{i} + 2)^{\kappa} - (\bar{x}_{i} + 1)^{\kappa} \right)$$

$$- \sum_{i \in \mathcal{I}} r_{i}^{n}(h_{i}) \eta_{i}^{n}(h_{i}) \xi_{i} \left(|\bar{x}_{i} + 1|^{\kappa} - |\bar{x}_{i}|^{\kappa} \right)$$

$$+ \sum_{i \in \mathcal{I}} \eta_{i}^{n}(h_{i}) (\mu_{i}^{n} \check{z}_{i}^{n} + \gamma_{i}^{n} \check{q}_{i}^{n}) \mathcal{O}(|\bar{x}_{i}|^{\kappa - 2}) + \sum_{i \in \mathcal{I}} r_{i}^{n}(h_{i}) \xi_{i} (|\bar{x}_{i} + 1|^{\kappa} - |\bar{x}_{i}|^{\kappa})$$

$$+ \sum_{i \in \mathcal{I}} (\mu_{i}^{n} \check{z}_{i}^{n} + \gamma_{i}^{n} \check{q}_{i}^{n}) \xi_{i} (|\bar{x}_{i} - 1|^{\kappa} - |\bar{x}_{i}|^{\kappa}), \tag{B.5}$$

where for the fourth term on the right-hand side we also used the fact that

$$(|\bar{x}_i|^{\kappa} - |\bar{x}_i - 1|^{\kappa}) - (|\bar{x}_i + 1|^{\kappa} - |\bar{x}_i|^{\kappa}) = \mathcal{O}(|\bar{x}_i|^{\kappa - 2}).$$

It is clear that $\eta_i^n(0) = 0$, since $F_i(0) = 0$ and $\mathbb{E}[G_i] = 1$. On the other hand, $\eta_i^n(t)$ is bounded for all $n \in \mathbb{N}$ and $t \ge 0$ by Assumption 3.2. Thus, applying (B.4), (B.5) and (4.4), it follows that

$$\mathcal{L}_{n}^{z^{n}} f_{n}(\check{x}, h) = \sum_{i \in \mathbb{J}} \left[\xi_{i} (\lambda_{i}^{n} - \mu_{i}^{n} \check{z}_{i}^{n} - \gamma_{i}^{n} \check{q}_{i}^{n}) \left(\kappa(\bar{x}_{i})^{\kappa - 1} + \mathcal{O}(|\bar{x}_{i}|^{\kappa - 2}) \right) + \eta_{i}^{n}(h_{i}) (\mu_{i}^{n} \check{z}_{i}^{n} + \gamma_{i}^{n} \check{q}_{i}^{n}) \mathcal{O}(|\bar{x}_{i}|^{\kappa - 2}) \right].$$
(B.6)

Since $\eta_i^n(h_i)$ is uniformly bounded, and $\check{z}_i^n, \check{q}_i^n \leq \bar{x}_i + n\rho_i$, it follows that the last term in (B.6) is equal to $\mathcal{O}(n)\mathcal{O}(|\bar{x}_i|^{\kappa-2}) + \mathcal{O}(|\bar{x}_i|^{\kappa-1})$. Note that for $i \in \mathcal{I} \setminus \mathcal{I}_0, \check{z}_i^n$ is equivalent to the static priority scheduling policy. Note also, that

$$\bar{x}_i \geq \check{z}_i^n - n\rho_i \geq \frac{n\rho_i \sum_{j \in \Im \setminus \Im_0} \rho_j}{\sum_{j \in \Im_0} \rho_j} > 0 \quad \forall i \in \tilde{\mathcal{K}}_n(\check{x}),$$
 (B.7)

and for $i \in \mathcal{I}_0 \setminus \tilde{\mathcal{K}}_n(\check{x})$, we have $\check{z}_i^n - n\rho_i = \bar{x}_i$ and $\check{q}_i^n = 0$. By using (B.6), and the identity in (5.20), we obtain

$$\mathcal{L}_{n}^{\check{z}^{n}} f_{n}(\check{x}, h) \leq \sum_{i \in \Im \backslash \Im_{0}} \xi_{i} \left(-\mu_{i}^{n} \bar{x}_{i} + (\mu_{i}^{n} - \gamma_{i}^{n}) \check{q}_{i}^{n} \right) m(\bar{x}_{i})^{\kappa - 1} \\
- \sum_{i \in \tilde{\mathcal{K}}_{n}(\check{x})} \xi_{i} \left(\mu_{i}^{n} (\check{z}_{i}^{n} - n\rho_{i}) + \gamma_{i}^{n} \check{q}_{i}^{n} \right) |\bar{x}_{i}|^{\kappa - 1} \\
- \sum_{i \in \Im_{0} \backslash \tilde{\mathcal{K}}_{n}(\check{x})} \xi_{i} \mu_{i}^{n} |\bar{x}_{i}|^{\kappa} + \sum_{i \in \Im} \left(\Im(\sqrt{n}) \Im(|\bar{x}_{i}|^{\kappa - 1}) \right) \\
+ \Im(n) \Im(|\bar{x}_{i}|^{\kappa - 2}) \right).$$
(B.8)



Let $\check{c}_1 := \sup_{i,n} \{\gamma_i^n, \mu_i^n\}$, and \check{c}_2 be some constant such that $\inf \{\mu_i^n, \gamma_j^n : i \in \mathbb{J}, j \in \mathbb{J} \setminus \mathbb{J}_0, n \in \mathbb{N}\} \geq \check{c}_2 > 0$. We select a positive vector $\xi \in \mathbb{R}_+^d$ such that $\xi_1 := 1$, $\xi_i := \frac{\kappa_1^m}{d^K} \min_{i' \leq i-1} \xi_{i'}, i \geq 2$, with $\kappa_1 := \frac{\check{c}_1}{8\check{c}_2}$. Compared with [4, Lemma 5.1], the important difference here is that, for $i \in \mathbb{J} \setminus \mathbb{J}_0$, we have

$$\check{q}_{i}^{n} = \left(\check{x}_{i} - \left(n - \sum_{j \in \tilde{\mathcal{K}}_{n}(\check{x})} \check{z}_{j}^{n} - \sum_{j \in \mathcal{I}_{0} \setminus \tilde{\mathcal{K}}_{n}(\check{x})} x_{j} - \sum_{j = |\mathcal{I}_{0}| + 1}^{i - 1} x_{j}\right)^{+}\right)^{+}.$$

Repeating the argument in the proof of [4, Lemma 5.1], it follows by (B.8) that

$$\mathcal{L}_{n}^{\check{z}^{n}} f_{n}(\check{x}, h) \leq c_{3} n^{\kappa/2} - c_{4} \sum_{i \in \Im\backslash\tilde{\mathcal{K}}_{n}(\check{x})} \xi_{i} |\bar{x}_{i}|^{\kappa} - c_{5} \sum_{i \in \tilde{\mathcal{K}}_{n}(\check{x})} \xi_{i} \left(\mu_{i}^{n} (\check{z}_{i}^{n} - n\rho_{i}) + \gamma_{i}^{n} \check{q}_{i}^{n}\right) |\bar{x}_{i}|^{\kappa-1} + \frac{c_{5}}{2} \sum_{i \in \tilde{\mathcal{K}}_{n}(\check{x})} \xi_{i} \mu_{i}^{n} (\check{z}_{i}^{n} - n\rho_{i})^{\kappa} + \sum_{i \in \Im} \left(\mathcal{O}(\sqrt{n}) \mathcal{O}(|\bar{x}_{i}|^{\kappa-1}) + \mathcal{O}(n) \mathcal{O}(|\bar{x}_{i}|^{\kappa-2}) \right)$$

$$(B.9)$$

for some positive constants c_3 , c_4 and c_5 . Therefore, (B.3) follows by (B.7) and (B.9), and this completes the proof.

Let

$$\tilde{g}_{n}(\check{\mathbf{x}}, h, \psi, k) \\
:= \frac{\psi + \alpha^{n}(k)}{\vartheta^{n}} \sum_{i \in \mathcal{I}} \mu_{i}^{n} \xi_{i} \Big(\tilde{g}_{n,i}(\check{\mathbf{x}}_{i}) + \eta_{i}^{n}(h_{i}) \Big(\tilde{g}_{n,i}(\check{\mathbf{x}}_{i} + 1) - \tilde{g}_{n,i}(\check{\mathbf{x}}_{i}) \Big) \Big) \tag{B.10}$$

for $(\check{x}, h, \psi, k) \in \mathfrak{D}$, where $\tilde{g}_{n,i}(\check{x}_i) := -|\bar{x}_i|^{\kappa}$ for $i \in \mathfrak{I} \setminus \mathfrak{I}_0$, and

$$\tilde{g}_{n,i}(\check{x}_i) := \begin{cases} -|\bar{x}_i|^{\kappa}, & \text{if } \bar{x}_i < \frac{n\rho_i \sum_{j \in \mathbb{J} \setminus \mathbb{J}_0} \rho_j}{\sum_{j \in \mathbb{J}_0} \rho_j}, \\ -\frac{n\rho_i \sum_{j \in \mathbb{J}_0} \rho_j}{\sum_{j \in \mathbb{J}_0} \rho_j} |\bar{x}_i|^{\kappa-1}, & \text{if } \bar{x}_i \geq \frac{n\rho_i \sum_{j \in \mathbb{J} \setminus \mathbb{J}_0} \rho_j}{\sum_{j \in \mathbb{J}_0} \rho_j}. \end{cases} \forall i \in \mathbb{J}_0.$$

Recall $\overline{\mathcal{L}}_{n,\psi}^{z^n}$ in (4.13). We also define

$$\overline{q}_i^{n,k}(\check{x},z^n) = \int_{\mathbb{R}_*} q_i^n (\check{x} - n\mu^n(y-k), z^n) \, \tilde{F}_{\check{x},k}^{d_1^n}(\mathrm{d}y).$$

Lemma B.2 Grant Assumptions 2.1, 2.2 and 3.2, and let $\xi \in \mathbb{R}^d_+$ be as in (B.2). Then, for any even integer $\kappa \geq 2$ and any $\varepsilon > 0$, there exist a positive constant \overline{C} , and $\overline{n} \in \mathbb{N}$, such that

$$\overline{\mathcal{L}}_{n,\psi}^{\tilde{z}^{n}} \tilde{g}_{n}(\check{x},h,\psi,k) \leq \overline{C} n^{\kappa/2} + \varepsilon \sum_{i \in \mathcal{I} \setminus \tilde{\mathcal{K}}_{n}(\check{x})} |\tilde{x}_{i}|^{\kappa} + \sum_{i \in \tilde{\mathcal{K}}_{n}(\check{x})} \mathcal{O}(|\tilde{x}_{i}|^{\kappa-1}) \\
+ \frac{1}{\sqrt{n}} \sum_{i \in \tilde{\mathcal{K}}_{n}(\check{x})} (\psi \mu_{i}^{n}(|z_{i}^{n} - n\rho_{i}|) + \psi \gamma_{i}^{n} q_{i}^{n} + (1 - \psi) \gamma_{i}^{n} \overline{q}_{i}^{n,k}) \mathcal{O}(|\tilde{x}_{i}|^{\kappa-1})$$
(B.11)

for any $z^n \in \mathfrak{Z}^n_{sm}$, and all $(\check{x}, h, \psi, k) \in \mathfrak{D}$ and $n > \bar{n}$.



Proof It is straightforward to verify that

$$|g_{n,i}(\check{x}_i \pm 1) - g_{n,i}(\check{x}_i)| = \mathcal{O}(|\bar{x}_i|^{\kappa-1}),$$

$$|(g_{n,i}(\check{x}_i) - g_{n,i}(\check{x}_i - 1)) - (g_{n,i}(\check{x}_i + 1) - g_{n,i}(\check{x}_i))| = \mathcal{O}(|\bar{x}_i|^{\kappa-2}),$$
(B.12)

for $i \in \mathcal{I}$. Repeating the calculation in (B.5) and (B.6), and applying (B.4) and (B.12), we have

$$\begin{split} \overline{\mathcal{L}}_{n,\psi}^{z^{n}} \tilde{g}_{n}(\check{x},h,\psi,k) &\leq \frac{\psi + \alpha^{n}(k)}{\vartheta^{n}} \\ \left[\sum_{i \in \tilde{\mathcal{K}}_{n}(\check{x})} \mu_{i}^{n} \xi_{i} \Big[\Big(|\lambda_{i}^{n} - n\mu_{i}^{n}\rho_{i}| + \psi\mu_{i}^{n}|z_{i}^{n} - n\rho_{i}| + \psi\gamma_{i}^{n}q_{i}^{n} + (1 - \psi)\gamma_{i}^{n}\overline{q}_{i}^{n,k} \Big) \mathcal{O}(|\bar{x}_{i}|^{\kappa - 1}) \right. \\ &+ \eta_{i}^{n}(h_{i}) \Big(\psi\mu_{i}^{n}z_{i}^{n} + \psi\gamma_{i}^{n}q_{i}^{n} + (1 - \psi)\gamma_{i}^{n}\overline{q}_{i}^{n,k} \Big) \mathcal{O}(|\bar{x}_{i}|^{\kappa - 2}) \Big] \\ &+ \sum_{i \in \mathcal{I} \setminus \tilde{\mathcal{K}}_{n}(\check{x})} \mu_{i}^{n} \xi_{i} \Big[\Big(\lambda_{i}^{n} + (1 - \psi)n\mu_{i}^{n}\rho_{i} \\ &+ \Big(1 + \eta_{i}^{n}(h_{i}) \Big) (\psi\mu_{i}^{n}z_{i}^{n} + \psi\gamma_{i}^{n}q_{i}^{n} + (1 - \psi)\gamma_{i}^{n}\overline{q}_{i}^{n,k} \Big) \mathcal{O}(|\bar{x}_{i}|^{\kappa - 1}) \Big] \Big]. \end{split} \tag{B.13}$$

Note that $\overline{q}_i^{n,k} \le c(1+\langle e, \overline{x}\rangle^+)$ for some positive constant c, by (5.21). Since z_i^n , $q_i^n \le \overline{x}_i + n\rho_i$, $(\vartheta^n)^{-1}$ is of order $n^{-1/2}$ by Assumption 2.2, and η_i^n and α^n are bounded, it follows by (5.20) and (B.13) that

$$\begin{split} & \overline{\mathcal{L}}_{n,\psi}^{z^n} \tilde{g}_n(\check{x},h,\psi,k) \\ & \leq \sum_{i \in \mathbb{J} \setminus \tilde{\mathcal{K}}_n(\check{x})} \frac{1}{\sqrt{n}} \Big(\mathbb{O}(n) \mathbb{O}(|\bar{x}_i|^{\kappa-1}) + \mathbb{O}(|\bar{x}_i|^{\kappa}) \Big) + \sum_{i \in \tilde{\mathcal{K}}_n(\check{x})} \mathbb{O}(\sqrt{n}) \mathbb{O}(|\bar{x}_i|^{\kappa-2}) \\ & + \sum_{i \in \tilde{\mathcal{K}}_n(\check{x})} \frac{1}{\sqrt{n}} \Big(\mathbb{O}(\sqrt{n}) + \psi \mu_i^n |z_i^n - n \rho_i| + \psi \gamma_i^n q_i^n + (1 - \psi) \gamma_i^n \overline{q}_i^{n,k} \Big) \mathbb{O}(|\bar{x}_i|^{\kappa-1}). \end{split}$$

Thus, applying Young's inequality, we obtain (B.11), and this completes the proof. \Box

Proof of Lemma 4.1 We define the function $\tilde{f}_n \in \mathcal{C}(\mathbb{R}^d \times \mathbb{R}^d_+ \times \{0, 1\} \times \mathbb{R}_+)$ by

$$\tilde{f}_n(\check{x}, h, \psi, k) := f_n(\check{x}, h) + \tilde{g}_n(\check{x}, h, \psi, k),$$

with f_n and \tilde{g}_n in (B.2) and (B.10), respectively. Recall $\widetilde{\mathcal{V}}_{\kappa,\xi}^n$ in (4.17). With $\xi \in \mathbb{R}_+^d$ as in (B.2), we have

$$n^{\kappa/2} \widetilde{\mathcal{V}}_{\kappa,\xi}^{n}(\tilde{x}^{n}(\check{x}),h,\psi,k) = \tilde{f}_{n}(\check{x},h,\psi,k) \quad \forall (\check{x},h,\psi,k) \in \mathfrak{D}.$$



Hence, to prove (4.18), it suffices to show that

$$\check{\mathcal{L}}_{n}^{\check{z}^{n}} \tilde{f}_{n}(\check{x}, h, \psi, k) \leq \widetilde{C}_{0} n^{\kappa/2}
-\widetilde{C}_{1} \sum_{i \in \Im \setminus \check{\mathcal{K}}_{n}(x)} \xi_{i} |\bar{x}_{i}|^{\kappa} - \widetilde{C}_{1} \sqrt{n} \sum_{i \in \check{\mathcal{K}}_{n}(\check{x})} \xi_{i} |\bar{x}_{i}|^{\kappa-1} \quad \forall n > \check{n},$$
(B.14)

and all $(\check{x}, h, \psi, k) \in \mathfrak{D}$, where the generator $\check{\mathcal{L}}_n^{\check{z}^n}$ is given in (4.12). It is clear that $\mathcal{Q}_{n,\psi} f_n(\check{x}, h) = 0$. Since $(\vartheta^n)^{-1}$ is of order $n^{-1/2}$, it follows by (4.10) and (4.15) that

$$\begin{aligned} \mathcal{Q}_{n,0}\tilde{g}_{n}(\check{x},h,0,k) &\leq \sum_{i\in\mathcal{I}\setminus\tilde{\mathcal{K}}_{n}(\check{x})} -\mu_{i}^{n}\xi_{i}|\bar{x}_{i}|^{\kappa} \\ &+ \sum_{i\in\tilde{\mathcal{K}}_{n}(\check{x})} -\mu_{i}^{n}\xi_{i}\frac{n\rho_{i}\sum_{j\in\mathcal{I}\setminus\mathcal{I}_{0}}\rho_{j}}{\sum_{j\in\mathcal{I}_{0}}\rho_{j}}|\bar{x}_{i}|^{\kappa-1} \\ &+ \epsilon_{n}\sum_{i\in\mathcal{I}\setminus\tilde{\mathcal{K}}_{n}(\check{x})} \mathcal{O}(|\bar{x}_{i}|^{\kappa}) + \sum_{i\in\tilde{\mathcal{K}}_{n}(\check{x})} \mathcal{O}(\sqrt{n})\mathcal{O}(|\bar{x}_{i}|^{\kappa-1}), \quad (B.15) \end{aligned}$$

where C is some positive constant and $\epsilon_n \to 0$ as $n \to \infty$. Since all the moments of d_1 are finite by (3.10) and $(a+z)^{\kappa} - a^{\kappa} = \mathcal{O}(z)\mathcal{O}(a^{\kappa-1}) + \mathcal{O}(z^2)\mathcal{O}(a^{\kappa-2}) + \cdots + \mathcal{O}(z^{\kappa})$ for any $a, z \in \mathbb{R}$, it is easy to verify that

$$\mathcal{I}_{n,1}\hat{f}_n(\check{x},h,1,0) = \sum_{i \in \mathcal{I}} \sum_{j=1}^{\kappa} \mathcal{O}(n^{j/2}) \mathcal{O}(|\bar{x}_i|^{\kappa-j}),$$
 (B.16)

using also the fact that

$$\beta_{\mathsf{u}}^{n} \int_{R_{*}} \left(\frac{n}{\vartheta^{n}} \mu_{i}^{n} \rho_{i} z \right)^{j} F^{d_{1}}(\mathrm{d}z) = \beta_{\mathsf{u}}^{n} \left(\frac{n}{\vartheta^{n}} \right)^{j} (\mu_{i}^{n} \rho_{i})^{j} \mathbb{E} \left[(d_{1})^{j} \right] = \mathfrak{O}(n^{j/2}) \quad \forall j > 0,$$

which follows by Assumptions 2.1, 2.2, and (3.10). Then, for $\psi = 1$, it follows by (B.16) and Young's inequality that

$$\check{\mathcal{L}}_{n}^{\check{z}^{n}} \tilde{f}_{n}(\check{x}, h, 1, 0) \leq \mathcal{L}_{n}^{\check{z}^{n}} f_{n}(\check{x}, h) + \overline{\mathcal{L}}_{n, 1}^{\check{z}^{n}} \tilde{g}_{n}(\check{x}, h, 1, 0)
+ C n^{\kappa/2} + \epsilon_{n} \sum_{i \in \mathcal{I} \setminus \tilde{\mathcal{K}}_{n}(\check{x})} \mathcal{O}(|\bar{x}_{i}|^{\kappa})
+ \sum_{i \in \tilde{\mathcal{K}}_{n}(\check{x})} \mathcal{O}(\sqrt{n}) \mathcal{O}(|\bar{x}_{i}|^{\kappa-1}).$$
(B.17)

Note that the last two terms in (B.3) and the last term in (B.11) are of smaller order than the second and third terms on the right-hand side of (B.3), respectively. Thus,



applying Lemmas B.1 and B.2, and using (B.17), we obtain

$$n^{-\kappa/2} \check{\mathcal{L}}_{n}^{\check{z}^{n}} \tilde{f}_{n}(\check{x}, h, 1, 0) \leq \widetilde{C}_{0} - \widetilde{C}_{1} \sum_{i \in \mathcal{I} \setminus \check{\mathcal{K}}_{n}(\check{x})} |\bar{x}_{i}|^{\kappa}$$

$$-\widetilde{C}_{1} \sum_{i \in \check{\mathcal{K}}_{n}(\check{x})} n^{-1/2} \left(\mu_{i}^{n} (\check{z}_{i}^{n} - n\rho_{i}) + \gamma_{i}^{n} \check{q}_{i}^{n} \right) |\tilde{x}_{i}|^{\kappa - 1}$$
(B.18)

for all large enough n, where \tilde{x} is defined in Definition 4.2. Since $\check{q}_i^n \geq 0$ and $\check{z}_i^n - n\rho_i > 0$ for $i \in \tilde{\mathcal{K}}_n(\check{x})$, then by using (B.7) and (B.18), we see that (B.14) holds when y = 1. For $\psi = 0$, using (B.15), Young's inequality, and the fact that for $i \in \tilde{\mathcal{K}}_n(\check{x})$, $\bar{x}_i > 0$, we obtain

$$\begin{split} & \check{\mathcal{L}}_{n}^{\check{z}^{n}} \tilde{f}_{n}(\check{x},h,0,k) \\ & \leq \sum_{i \in \mathcal{I}} \mathcal{O}(\sqrt{n}) \mathcal{O}(|\bar{x}_{i}|^{\kappa-1}) + \sum_{i \in \mathcal{I}} \mathcal{O}(n) \mathcal{O}(|\bar{x}_{i}|^{\kappa-2}) + Cn^{\kappa/2} \\ & + (\epsilon + \epsilon_{n}) \sum_{i \in \mathcal{I} \setminus \tilde{\mathcal{K}}_{n}(\check{x})} \xi_{i} |\bar{x}_{i}|^{\kappa} \\ & + \sum_{i \in \mathcal{I} \setminus \tilde{\mathcal{K}}_{n}(\check{x})} \left(-\mu_{i}^{n} \xi_{i} |\bar{x}_{i}|^{\kappa} + \gamma_{i}^{n} \xi_{i} \overline{q}_{i}^{n,k} \left(-\kappa(\bar{x}_{i})^{\kappa-1} + \mathcal{O}(|\bar{x}_{i}|^{\kappa-2}) \right) \right) \\ & + \sum_{i \in \tilde{\mathcal{K}}_{n}(\check{x})} - \frac{n\rho_{i} \sum_{j \in \mathcal{I} \setminus \mathcal{I}_{0}} \rho_{j}}{\sum_{j \in \mathcal{I}_{0}} \rho_{j}} \mu_{i}^{n} \xi_{i} |\bar{x}_{i}|^{\kappa-1} + \overline{\mathcal{L}}_{n,0}^{\check{z}^{n}} \tilde{g}_{n}(\check{x},h,0,k) \end{split}$$

for some positive constant C and sufficiently small $\epsilon > 0$. We proceed by invoking the argument in the proof of [4, Lemma 5.1]. The important difference here is that

$$\check{q}_i^n (\check{x} - n\mu^n (z - k)) = \tilde{\epsilon}_i (\check{x} - n\mu^n (z - k)) (\bar{x}_i - n\mu_i \rho_i (z - k))
+ \bar{\epsilon}_i (\check{x} - n\mu^n (z - k)) \sum_{i=1}^{i-1} (\bar{x}_j - n\mu_j \rho_j (z - k)),$$

where the functions $\tilde{\epsilon}_i$, $\bar{\epsilon}_i$: $\mathbb{R}^d \to [0,1]$, for $i \in \mathcal{I}$. Since $\tilde{\epsilon}_i$ and $\bar{\epsilon}_i$ are bounded, we have some additional terms which are bounded by $C \int_{\mathbb{R}_*} n \mu_i \rho_i(y-k) \tilde{F}_{\check{\chi},k}^{d_1^n}(\mathrm{d}y)$ for some positive constant C. Therefore, these are of order \sqrt{n} by (5.21). Thus, repeating the argument in the proof of Lemma B.1, and applying Lemma B.2, we deduce that (B.14) holds with $\psi = 0$. This completes the proof.

Proof of Lemma 5.2 The proof mimics that of Proposition 4.2. We sketch the proof when \mathfrak{I}_0 is empty. Using the estimate

$$\mathcal{O}(q_i^n)\mathcal{O}(|\bar{x}_i|^{m-1}) \le \epsilon^{1-m} \left(\mathcal{O}(q_i^n)\right)^m + \epsilon \left(\mathcal{O}(|\bar{x}_i|^{m-1})\right)^{m/m-1} \tag{B.19}$$



for any $\epsilon > 0$, which follows by Young's inequality, we deduce that, for some positive constants $\{c_k : k = 1, 2, 3\}$, we have

$$\mathcal{L}_{n}^{z^{n}} f_{n}(\check{x}, h) \leq c_{1} n^{m/2} + c_{2} (\langle e, q^{n} \rangle)^{m} - c_{3} \sum_{i \in \mathcal{I}} \xi_{i} |\bar{x}_{i}|^{m} \quad \forall (\check{x}, h) \in \mathbb{R}_{+}^{d} \times \mathbb{R}_{+}^{d}, \text{ (B.20)}$$

and all large enough n. Note that Lemma B.2 holds for all $z^n \in \mathfrak{Z}_{sm}^n$. Then, we may repeat the steps in the proof of Lemma 4.1, except that here we use

$$(\tilde{x}_{i})^{m-1} \int_{\mathbb{R}_{*}} \hat{q}_{i}^{n} (\check{x} - n\mu^{n}(y - k), z^{n}) \, \tilde{F}_{\check{x}, k}^{d_{1}^{n}} (\mathrm{d}y)$$

$$\leq \epsilon |\bar{x}_{i}|^{m} + \epsilon^{1-m} \Big(\mathbb{E} \Big[\hat{q}_{i}^{n} (\check{x} - n\mu^{n}(d_{1}^{n} - k), z^{n}) \, | \, d_{1}^{n} > k \Big] \Big)^{m},$$
(B.21)

where $\hat{q}^n = n^{-1/2}q^n$, with $\epsilon > 0$ chosen sufficiently small. Since $\hat{q}_i^n(\check{x}, z^n) \leq \langle e, \tilde{x} \rangle^+$, it follows by (5.21) that

$$\mathbb{E}\left[\hat{q}_i^n\left(\check{x} - n\mu^n(d_1^n - k), z^n\right) \mid d_1^n > k\right] \le c_4(1 + \langle e, \tilde{x} \rangle^+). \tag{B.22}$$

Thus, by the same calculation in Proposition 4.2, and using (B.19)–(B.22), we obtain

$$\mathbb{E}^{z^n} \left[\int_0^T |\widetilde{X}^n(s)|^m \right] \le C_1 (T + |\widehat{X}^n(0)|^m)$$

$$+ C_2 \mathbb{E}^{z^n} \left[\int_0^T \left(1 + \langle e, \widetilde{X}^n(s) \rangle^+ \right)^m ds \right]$$
(B.23)

for all large enough n, and $\{z^n \in \mathfrak{Z}^n_{sm} : n \in \mathbb{N}\}$. Since $\sup_n \hat{J}(\hat{X}^n(0), z^n) < \infty$, it follows by (4.26) that

$$\sup_{n} \limsup_{T \to \infty} \frac{1}{T} \mathbb{E} \left[\int_{0}^{T} \left(\langle e, \widetilde{X}^{n}(s) \rangle^{+} \right)^{m} \mathrm{d}s \right] < \infty.$$

Therefore, dividing both sides of (B.23) by T, taking $T \to \infty$ and using (4.26) again, we obtain (5.11). We may show that the result also holds when \mathfrak{I}_0 is nonempty by repeating the above argument and applying Lemma B.2. This completes the proof. \square

References

- Atar, R., Mandelbaum, A., Reiman, M.I.: Scheduling a multi class queue with many exponential servers: asymptotic optimality in heavy traffic. Ann. Appl. Probab. 14(3), 1084–1134 (2004)
- Atar, R.: Scheduling control for queueing systems with many servers: asymptotic optimality in heavy traffic. Ann. Appl. Probab. 15(4), 2606–2650 (2005). https://doi.org/10.1214/105051605000000601
- Atar, R., Mandelbaum, A., Shaikhet, G.: Simplified control problems for multiclass many-server queueing systems. Math. Oper. Res. 34(4), 795–812 (2009). https://doi.org/10.1287/moor.1090.0404
- 4. Arapostathis, A., Biswas, A., Pang, G.: Ergodic control of multi-class M/M/N + M queues in the Halfin-Whitt regime. Ann. Appl. Probab. **25**(6), 3511–3570 (2015)



- Arapostathis, A., Pang, G.: Infinite-horizon average optimality of the N-network in the Halfin-Whitt regime. Math. Oper. Res. 43(3), 838–866 (2018)
- Arapostathis, A., Pang, G.: Infinite horizon asymptotic average optimality for large-scale parallel server networks. Stoch. Process. Appl. 129(1), 283–322 (2019)
- 7. Budhiraja, A., Ghosh, A., Liu, X.: Scheduling control for Markov-modulated single-server multiclass queueing systems in heavy traffic. Queueing Syst. **78**(1), 57–97 (2014)
- 8. Kumar, R., Lewis, M.E., Topaloglu, H.: Dynamic service rate control for a single-server queue with Markov-modulated arrivals. Naval Res. Logist. **60**(8), 661–677 (2013)
- 9. Xia, L., He, Q., Alfa, A.S.: Optimal control of state-dependent service rates in a MAP/M/1 queue. IEEE Trans. Autom. Control **62**(10), 4965–4979 (2017)
- Arapostathis, A., Das, A., Pang, G., Zheng, Y.: Optimal control of Markov-modulated multiclass many-server queues. Stoch. Syst. 9(2), 155–181 (2019)
- Jansen, H. M., Mandjes, M., De Turck, K., Wittevrongel, S.: Diffusion limits for networks of Markov-modulated infinite-server queues. Perform. Eval. 135 (2019). https://doi.org/10.1016/j.peva.2019. 102039.
- Pang, G., Whitt, W.: Service interruptions in large-scale service systems. Manag. Sci. 55(9), 1499–1512 (2009)
- Pang, G., Whitt, W.: Heavy-traffic limits for many-server queues with service interruptions. Queueing Syst. 61(2–3), 167–202 (2009). https://doi.org/10.1007/s11134-009-9104-2
- Lu, H., Pang, G., Zhou, Y.: G/GI/N(+GI) queues with service interruptions in the Halfin-Whitt regime. Math. Methods Oper. Res. 83(1), 127–160 (2016). https://doi.org/10.1007/s00186-015-0523-z
- Lu, H., Pang, G.: Heavy-traffic limits for an infinite-server fork-join queueing system with dependent and disruptive services. Queueing Syst. 85(1–2), 67–115 (2017). https://doi.org/10.1007/s11134-016-9505-y
- Pang, G., Zhou, Y.: G/G/∞ queues with renewal alternating interruptions. Adv. Appl. Probab. 48(3), 812–831 (2016). https://doi.org/10.1017/apr.2016.29
- 17. Arapostathis, A., Pang, G., Zheng, Y.: Ergodic control of diffusions with compound Poisson jumps under a general structural hypothesis. ArXiv e-prints. arXiv:1908.01068 (2019)
- 18. Atar, R., Giat, C., Shimkin, N.: On the asymptotic optimality of the $c\mu/\theta$ rule under ergodic cost. Queueing Syst. **67**(2), 127–144 (2011). https://doi.org/10.1007/s11134-010-9206-x
- Konstantopoulos, T., Last, G.: On the use of Lyapunov function methods in renewal theory. Stoch. Process. Appl. 79(1), 165–178 (1999)
- Meyn Sean, P., Tweedie, R.L.: Stability of Markovian processes. III. Foster–Lyapunov criteria for continuous-time processes. Adv. Appl. Probab. 25(3), 518–548 (1993)
- Billingsley, P.: Convergence of Probability Measures. Wiley Series in Probability and Statistics, Second edn. Wiley, New York (1999). https://doi.org/10.1002/9780470316962. A Wiley-Interscience Publication
- Whitt, W.: Stochastic-Process Limits. An Introduction to Stochastic-Process Limits and Their Application to Queues. Springer Series in Operations Research. Springer, New York (2002)
- 23. Dai, J.G.: On positive Harris recurrence of multiclass queueing networks: a unified approach via fluid limit models. Ann. Appl. Probab. 5(1), 49–77 (1995)
- 25. Arapostathis, A., Hmedi, H., Pang, G., Sandrić, N.: Uniform polynomial rates of convergence for a class of Lévy-driven controlled SDEs arising in multiclass many-server queues. In: Yin, G., Zhang, Q. (eds.) Modeling, Stochastic Control, Optimization, and Applications. The IMA Volumes in Mathematics and its Applications, vol. 164. Springer, Cham (2019)
- Davis, M.H.A.: Piecewise-deterministic Markov processes: a general class of nondiffusion stochastic models. J. R. Stat. Soc. Ser. B 46(3), 353–388 (1984). With discussion
- 27. Ross, Sheldon M.: Stochastic Processes, 2nd edn. Wiley, New York (1996)
- 28. Arapostathis, A., Caffarelli, L., Pang, G., Zheng, Y.: Ergodic control of a class of jump diffusions with finite Lévy measures and rough kernels. SIAM J. Control Optim. **57**(2), 1516–1540 (2019)
- Krichagina, E.V., Taksar, M.I.: Diffusion approximation for GI/G/1 controlled queues. Queueing Syst. Theory Appl. 12(3–4), 333–367 (1992). https://doi.org/10.1007/BF01158808
- Pang, G., Talreja, R., Whitt, W.: Martingale proofs of many-server heavy-traffic limits for Markovian queues. Probab. Surv. 4, 193–267 (2007). https://doi.org/10.1214/06-PS091



 Iglehart, D.L., Whitt, W.: The equivalence of functional central limit theorems for counting processes and associated partial sums. Ann. Math. Stat. 42, 1372–1378 (1971). https://doi.org/10.1214/aoms/ 1177693249

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

