

## Dynamic Scheduling of a Parallel Server System in Heavy Traffic with Complete Resource Pooling: Asymptotic Optimality of a Threshold Policy

S. L. Bell and R. J. Williams <sup>1</sup>

Department of Mathematics  
University of California, San Diego  
La Jolla CA 92093-0112

Email: [williams@math.ucsd.edu](mailto:williams@math.ucsd.edu)  
<http://www.math.ucsd.edu/~williams>

### Abstract

We consider a parallel server queueing system consisting of a bank of buffers for holding incoming jobs and a bank of flexible servers for processing these jobs. Incoming jobs are classified into one of several different classes (or buffers). Jobs within a class are processed on a first-in-first-out basis, where the processing of a given job may be performed by any server from a given (class-dependent) subset of the bank of servers. The random service time of a job may depend on both its class and the server providing the service. Each job departs the system after receiving service from one server. The system manager seeks to minimize holding costs by dynamically scheduling waiting jobs to available servers. We consider a parameter regime in which the system satisfies both a heavy traffic and a complete resource pooling condition. Our cost function is an expected cumulative discounted cost of holding jobs in the system, where the (undiscounted) cost per unit time is a linear function of normalized (with heavy traffic scaling) queue length. In a prior work [42], the second author proposed a continuous review threshold control policy for use in such a parallel server system. This policy was advanced as an “interpretation” of the analytic solution to an associated Brownian control problem (formal heavy traffic diffusion approximation). In this paper we show that the policy proposed in [42] is asymptotically optimal in the heavy traffic limit and that the limiting cost is the same as the optimal cost in the Brownian control problem.

**Keywords and phrases:** Stochastic networks, dynamic control, resource pooling, heavy traffic, diffusion approximations, Brownian control problems, state space collapse, threshold policies, large deviations.

**AMS Subject Classification (2000):** Primary 60K25, 68M20, 90B36; secondary 60J70.

Submitted to EJP on July 16, 2004. Final version accepted on August 17, 2005.

---

<sup>1</sup>Research supported in part by NSF Grants DMS-0071408 and DMS-0305272, a John Simon Guggenheim Fellowship, and a gift from the David and Holly Mendel Fund.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1046</b>
1.1	Notation and Terminology . . . . .	1048
<b>2</b>	<b>Parallel Server System</b>	<b>1049</b>
2.1	System Structure . . . . .	1049
2.2	Stochastic Primitives . . . . .	1050
2.3	Scheduling Control and Performance Measures . . . . .	1050
<b>3</b>	<b>Sequence of Systems, Heavy Traffic, and the Cost Function</b>	<b>1052</b>
3.1	Sequence of Systems and Large Deviation Assumptions . . . . .	1052
3.2	Heavy Traffic and Fluid Model . . . . .	1053
3.3	Diffusion Scaling and the Cost Function . . . . .	1055
<b>4</b>	<b>Brownian Control Problem</b>	<b>1056</b>
4.1	Formulation . . . . .	1056
4.2	Solution via Workload Assuming Complete Resource Pooling . . . . .	1056
4.3	Approaches to Interpreting the Solution . . . . .	1059
<b>5</b>	<b>Threshold Policy and Main Results</b>	<b>1060</b>
5.1	Tree Conventions . . . . .	1060
5.2	Threshold Policy . . . . .	1060
5.3	Threshold Sizes . . . . .	1061
5.4	Examples . . . . .	1062
5.5	Main Results . . . . .	1064
<b>6</b>	<b>Preliminaries and Outline of the Proof</b>	<b>1065</b>
6.1	The Server-Buffer Tree $\mathcal{G}$ : Layers and Buffer Renumbering . . . . .	1065
6.2	Threshold Sizes and Transient Nominal Activity Rates . . . . .	1066
6.3	State Space Collapse Result and Outline of Proof . . . . .	1068
6.4	Residual Processes and Shifted Allocation Processes . . . . .	1069
6.5	Preliminaries on Stopped Arrival and Service Processes . . . . .	1070
6.6	Large Deviation Bounds for Renewal Processes . . . . .	1071
<b>7</b>	<b>Proof of State Space Collapse</b>	<b>1072</b>
7.1	Auxiliary Constants for the Induction Proof . . . . .	1073
7.2	Induction Setup . . . . .	1076
7.3	Estimates on Allocation Processes Imply Residual Processes Stay Near Zero – Proof of Lemma 7.3 . . . . .	1079
7.4	Estimates on Allocations for Activities Immediately Below Buffers – Proof of Lemma 7.4 . . . . .	1085
7.5	Estimates on Allocations for Activities Immediately Above Buffers – Proof of Lemma 7.5 . . . . .	1087
7.6	Transition Between Layers in the Server-Buffer Tree — Proof of Lemma 7.6 . . . . .	1087
7.7	Proofs of Theorems 6.1, 7.1, and 7.7 . . . . .	1098
<b>8</b>	<b>Weak Convergence under the Threshold Policy</b>	<b>1101</b>
8.1	Fluid Limits for Allocation Processes . . . . .	1101
8.2	Convergence of Diffusion Scaled Performance Measures under the Threshold Policy—Proof of Theorem 5.3 . . . . .	1105
<b>9</b>	<b>Asymptotic Optimality of the Threshold Policy</b>	<b>1106</b>

# 1 Introduction

We consider a dynamic scheduling problem for a parallel server queueing system. This system might be viewed as a model for a manufacturing or computer system, consisting of a bank of buffers for holding incoming jobs and a bank of flexible servers for processing these jobs (see e.g., [23]). Incoming jobs are classified into one of several different classes (or buffers). Jobs within a class are served on a first-in-first-out basis and each job may be served by any server from a given subset of the bank of servers (this subset may depend on the class). In addition, a given server may be able to service more than one class. Jobs depart the system after receiving one service. Jobs of each class incur linear holding costs while present within the system. The system manager seeks to minimize holding costs by dynamically allocating waiting jobs to available servers.

The parallel server system is described in more detail in Section 2 below. With the exception of a few special cases, the dynamic scheduling problem for this system cannot be analyzed exactly and it is natural to consider more tractable approximations. One class of such approximations are the so-called Brownian control problems which were introduced and developed as formal heavy traffic approximations to queueing control problems by Harrison in a series of papers [12, 19, 16, 17]. Various authors (see for example [7, 8, 20, 21, 24, 29, 30, 39]) have used analysis of these Brownian control problems, together with clever interpretation of their optimal (analytic) solutions, to suggest “good” policies for certain queueing control problems. (We note that some authors have used other approximations for queueing control problems such as fluid models, see e.g. [32] and references therein, but here we restrict attention to Brownian control problems as approximations.)

For the parallel server system considered here, Harrison and López [18] studied the associated Brownian control problem and identified a condition under which the solution of that problem exhibits *complete resource pooling*, i.e., in the Brownian model, the efforts of the individual servers can be efficiently combined to act as a single pooled resource or “superserver”. Under this condition, Harrison and López [18] conjectured that a “discrete review” scheduling policy (for the original parallel server system), obtained by using the BIGSTEP discretization procedure of Harrison [14], is asymptotically optimal in the heavy traffic limit. They did not attempt to prove the conjecture, although, in a slightly earlier work, Harrison [15] did prove asymptotic optimality of a discrete review policy for a two server case with special distributional assumptions.

Here, focusing on the parameter regime associated with heavy traffic and complete resource pooling, we first review the formulation and solution of the Brownian control problem. Then we prove that a continuous review “tree-based” threshold policy proposed by the second author in [42] is asymptotically optimal in the heavy traffic limit. (Independently of [42], Squillante et al. [35] proposed a tree-based threshold policy for the parallel server system. However, their policy is in general different from the one described in [42]. An example illustrating this is given in Section 5.4.) Our treatment of the Brownian control problem is similar to that in [18, 42] and our description of the candidate threshold policy is similar to that in [42], although some more details are included here. On the other hand, the proof of asymptotic optimality of this policy is new. In a related work [3], we have already proved that this policy is asymptotically optimal for a particular two-server, two-buffer system. Indeed, techniques developed in [3] have been useful for analysis of the more complex multiserver case treated here. These techniques have also recently been used in part by Budhiraja and Ghosh [7] in treating a classic controlled queueing network example known as the criss-cross network.

Since we began this work, three related works have appeared [1, 36, 31]. In [1], Ata and Kumar consider a dynamic scheduling problem for an open stochastic processing network that allows feedback routing. A parallel server system is a special case of such networks in which no routing occurs. Under heavy traffic and complete resource pooling conditions, Ata and Kumar prove asymptotic optimality of a discrete review policy for an open stochastic processing network with

linear holding costs. Although this provides an asymptotically optimal policy for the parallel server problem, we think it is still of interest to establish asymptotic optimality of a simple continuous review threshold policy, as we do here. Indeed, the policy proposed in [1] and the method of proof are substantially different from those used here. The discrete review policy of [1] observes the state of the system at discrete review times and determines open loop controls to be used over the periods between review times. Their control typically changes with the state observed at the beginning of each of these “review periods”. Our policy is a simple threshold priority policy. It continuously monitors the state, but the control only changes when a threshold is crossed or a queue becomes empty. Furthermore, our policy exploits a tree structure of the parallel server system, whereas the policy of [1], which is designed for more general systems with feedback, does not exploit the tree structure when restricted to parallel server systems. Finally, our proof involves an induction procedure to prove estimates related to the tree structure; there is not an analogue of this inductive proof in the paper [1].

The other related works are by Stolyar [36], who considers a generalized switch, which operates in discrete time, and Mandelbaum and Stolyar [31], who consider a parallel server system. Although it does not allow routing, Stolyar’s generalized switch is somewhat more general than a parallel server system operating in discrete time. In particular, it allows service rates that depend on the state of a random environment. Assuming heavy traffic and a resource pooling condition, which is slightly more general than a complete resource pooling condition, Stolyar [36] proves asymptotic optimality of a MaxWeight policy, for holding costs that are positive linear combinations of the individual queue lengths raised to the power  $\beta + 1$  where  $\beta > 0$ . In particular, the holding costs are not linear in the queue length. An advantage of the MaxWeight policy (which exploits the non-linear nature of the holding cost function) is that it does not require knowledge of the arrival rates for its execution, although checking the heavy traffic and resource pooling conditions does involve these rates.

Following on from [36], in [31], Mandelbaum and Stolyar focus on a parallel server system (operating in continuous time). Assuming heavy traffic and complete resource pooling conditions they prove asymptotic optimality of a MaxWeight policy (called a generalized  $c\mu$ -rule there), for holding costs that are sums of strictly increasing, strictly convex functions of the individual queue lengths. (They also prove a related result where queue lengths are replaced by sojourn times.) Again, the nonlinear nature of the holding cost function allows the authors to specify a policy that does not require knowledge of the arrival rates (nor of a solution of a certain dual linear program). Although Mandelbaum and Stolyar [31] conjecture a policy for linear holding costs (which like ours makes use of the solution of a dual linear program), they stop short of proving asymptotic optimality of that policy.

Thus, assuming heavy traffic scaling, the current paper provides the only proof of asymptotic optimality of a continuous review policy for the parallel server system with linear holding costs under a complete resource pooling condition. An additional difference between our work and that in [1, 31, 36] is that we impose finite exponential moment assumptions on our primitive stochastic processes, whereas only finite moment assumptions (of order  $2 + \epsilon$ ) are needed for the results in [1, 31, 36]. We conjecture that our exponential moment assumptions could be relaxed to moment conditions of sufficiently high order, at the expense of an increase in the size of our (logarithmic) thresholds. We have not pursued this conjecture here, having chosen the tradeoff of smaller thresholds at the expense of exponential moment assumptions.

This paper is organized as follows. In Section 2, we describe the model of a parallel server system considered here. In Section 3 we introduce a sequence of such systems, indexed by  $r$  (where  $r$  tends to infinity through a sequence of values in  $[1, \infty)$ ), which is used in formulating the notion of heavy traffic asymptotic optimality. The cost function used in the  $r^{\text{th}}$  system is an expected cumulative discounted linear holding cost, where the linear holding cost is per unit of normalized

queue length (in diffusion scale). In Section 3, we also review the notion of heavy traffic defined in [16, 18] using a linear program, and recall its interpretation in terms of the behavior of an associated fluid model, as previously described in [42]. In Section 4, we describe the Brownian control problem associated with the sequence of parallel server systems, and, under the complete resource pooling condition of Harrison and López [18], we review the solution of the Brownian control problem obtained in [18] using a reduced form of the problem called the equivalent workload formulation [16, 19]. The complete resource pooling condition ensures that the Brownian workload process is one-dimensional. Moreover, from [18] we know that this condition is equivalent to uniqueness of a solution to the dual to the linear program described in Section 3. In Section 5, we describe the dynamic threshold policy proposed in [42] for use in the parallel server system. This policy exploits a tree structure of a graph containing the servers and buffers as nodes. We then state the main result (Theorem 5.4) which implies that this policy is asymptotically optimal in the heavy traffic limit and that the limiting cost is the same as the optimal cost in the Brownian control problem. An outline of our method of proof is given in Section 6. The details of the proof are contained in Sections 7–9. Here a critical role is played by our analysis in Section 7 of what we call the residual processes, which measure the deviations of the queue lengths from the threshold levels, or from zero if a queue does not have a threshold on it, when the threshold policy is used. This allows us to establish a form of “state space collapse” (see Theorems 6.1 and 5.3) under this policy. The techniques used in proving state space collapse build on and extend those introduced in [3]. In particular, a major new feature is the need to show that allocations of time to various activities stay close to their nominal allocations over sufficiently long time intervals (with high probability), which in turn is used to show that the residual processes stay close to zero (with high probability). Using a suitable numbering of the buffers, the proof of state space collapse proceeds by induction on the buffer number, highlighting the fact that the queue length for a particular buffer depends (via the threshold policy) on the queue lengths associated with lower numbered buffers. Another new aspect of our proof lies in Section 9, where a certain uniform integrability is used to prove convergence of normalized cost functions associated with the sequence of parallel server systems, operating under the threshold policy, to the optimal cost function for the Brownian control problem. To establish the uniform integrability, estimates of the probabilities that the residual processes deviate far from zero need to be sufficiently precise (cf. Theorem 7.7). This accounts for the appearance of the polynomial terms in (I)–(III) of Section 7.2. Also, in the proof of the uniform integrability of the normalized idletime processes, a technical point that was overlooked in the proof of the analogous result in [3] is corrected. Specifically, the proof is divided into two separate cases depending on the size of the time index (cf. (9.34)–(9.36)). (In the proof of Theorem 5.3 in [3], the estimates in (173) and (176) should have been divided into two cases corresponding to  $r^2t > 2/\tilde{\epsilon}$  and  $r^2t \leq 2/\tilde{\epsilon}$ .)

## 1.1 Notation and Terminology

The set of non-negative integers will be denoted by  $\mathbb{N}$  and the value  $+\infty$  will simply be denoted by  $\infty$ . For any real number  $x$ ,  $\lfloor x \rfloor$  will denote the integer part of  $x$ , i.e., the greatest integer that is less than or equal to  $x$ , and  $\lceil x \rceil$  will denote the smallest integer that is greater than or equal to  $x$ . We let  $\mathbb{R}_+$  denote  $[0, \infty)$ . The  $m$ -dimensional ( $m \geq 1$ ) Euclidean space will be denoted by  $\mathbb{R}^m$  and  $\mathbb{R}_+^m$  will denote the  $m$ -dimensional positive orthant,  $[0, \infty)^m$ . Let  $|\cdot|$  denote the norm on  $\mathbb{R}^m$  given by  $|x| = \sum_{i=1}^m |x_i|$  for  $x \in \mathbb{R}^m$ . We define a sum over an empty index set to be zero. Vectors in  $\mathbb{R}^m$  should be treated as column vectors unless indicated otherwise, inequalities between vectors should be interpreted componentwise, the transpose of a vector  $a$  will be denoted by  $a'$ , the diagonal matrix with the entries of a vector  $a$  on its diagonal will be denoted by  $\text{diag}(a)$ , and the dot product of two vectors  $a$  and  $b$  in  $\mathbb{R}^m$  will be denoted by  $a \cdot b$ . For any set  $\mathcal{S}$ , let  $|\mathcal{S}|$  denote the cardinality of  $\mathcal{S}$ .

For each positive integer  $m$ , let  $D^m$  be the space of “Skorokhod paths” in  $\mathbb{R}^m$  having time

domain  $\mathbb{R}_+$ . That is,  $D^m$  is the set of all functions  $\omega : \mathbb{R}_+ \rightarrow \mathbb{R}^m$  that are right continuous on  $\mathbb{R}_+$  and have finite left limits on  $(0, \infty)$ . The member of  $D^m$  that stays at the origin in  $\mathbb{R}^m$  for all time will be denoted by  $\mathbf{0}$ . For  $\omega \in D^m$  and  $t \geq 0$ , let

$$\|\omega\|_t = \sup_{s \in [0, t]} |\omega(s)|. \quad (1.1)$$

Consider  $D^m$  to be endowed with the usual Skorokhod  $J_1$ -topology (see [10]). Let  $\mathcal{M}^m$  denote the Borel  $\sigma$ -algebra on  $D^m$  associated with the  $J_1$ -topology. All of the continuous-time processes in this paper will be assumed to have sample paths in  $D^m$  for some  $m \geq 1$ . (We shall frequently use the term process in place of stochastic process.)

Suppose  $\{W^n\}_{n=1}^\infty$  is a sequence of processes with sample paths in  $D^m$  for some  $m \geq 1$ . Then we say that  $\{W^n\}_{n=1}^\infty$  is tight if and only if the probability measures induced by the  $W^n$ 's on  $(D^m, \mathcal{M}^m)$  form a tight sequence, i.e., they form a weakly relatively compact sequence in the space of probability measures on  $(D^m, \mathcal{M}^m)$ . The notation " $W^n \Rightarrow W$ ", where  $W$  is a process with sample paths in  $D^m$ , will mean that the probability measures induced by the  $W^n$ 's on  $(D^m, \mathcal{M}^m)$  converge weakly to the probability measure on  $(D^m, \mathcal{M}^m)$  induced by  $W$ . If for each  $n$ ,  $W^n$  and  $W$  are defined on the same probability space, we say that  $W^n$  converges to  $W$  uniformly on compact time intervals in probability (u.o.c. in prob.), if  $\mathbf{P}(\|W^n - W\|_t \geq \varepsilon) \rightarrow 0$  as  $n \rightarrow \infty$  for each  $\varepsilon > 0$  and  $t \geq 0$ . We note that if  $\{W^n\}$  is a sequence of processes and  $W$  is a continuous deterministic process (all defined on the same probability space) then  $W^n \Rightarrow W$  is equivalent to  $W^n \rightarrow W$  u.o.c. in prob. This is implicitly used several times in the proofs below to combine statements involving convergence in distribution to deterministic processes.

## 2 Parallel Server System

### 2.1 System Structure

Our parallel server system consists of  $\mathbf{I}$  infinite capacity buffers for holding jobs awaiting service, indexed by  $i \in \mathcal{I} \equiv \{1, \dots, \mathbf{I}\}$ , and  $\mathbf{K}$  (non-identical) servers working in parallel, indexed by  $k \in \mathcal{K} \equiv \{1, \dots, \mathbf{K}\}$ . Each buffer has its own stream of jobs arriving from outside the system. Arrivals to buffer  $i$  are called class  $i$  jobs and jobs are ordered within each buffer according to their arrival times, with the earliest arrival being at the head of the line. Each entering job requires a single service before it exits the system. Several different servers may be capable of processing (or serving) a particular job class. Service of a given job class  $i$  by a given server  $k$  is called a processing activity. We assume that there are  $\mathbf{J} \leq \mathbf{I} \cdot \mathbf{K}$  possible processing activities labeled by  $j \in \mathcal{J} \equiv \{1, \dots, \mathbf{J}\}$ . The correspondences between activities and classes, and activities and servers, are described by two deterministic matrices  $\mathbf{C}$ ,  $\mathbf{A}$ , where  $\mathbf{C}$  is an  $\mathbf{I} \times \mathbf{J}$  matrix with

$$\mathbf{C}_{ij} = \begin{cases} 1 & \text{if activity } j \text{ processes class } i, \\ 0 & \text{otherwise,} \end{cases} \quad (2.1)$$

and  $\mathbf{A}$  is a  $\mathbf{K} \times \mathbf{J}$  matrix with

$$\mathbf{A}_{kj} = \begin{cases} 1 & \text{if server } k \text{ performs activity } j, \\ 0 & \text{otherwise.} \end{cases} \quad (2.2)$$

Note that each column of  $\mathbf{C}$  contains the number one exactly once and similarly for  $\mathbf{A}$ , since each activity  $j$  has exactly one class  $i(j)$  and one server  $k(j)$  associated with it. We also assume that each row of  $\mathbf{C}$  and each row of  $\mathbf{A}$  contains the number one at least once (i.e., each job class is

capable of being processed by at least one activity and each server is capable of performing at least one activity).

Once a job has commenced service at a server, it remains there until its service is complete, even if its service is interrupted for some time (e.g., by preemption by a job of another class). A server may not start on a new job of class  $i$  until it has finished serving any class  $i$  job that it is working on or that is in suspension at the server. In addition, a server cannot work unless it has a job to work on. When taking a new job from a buffer, a server always takes the job at the head of the line. (For concreteness, we suppose that a deterministic tie-breaking rule is used when two (or more) servers want to simultaneously take jobs from the same buffer, e.g., there is an ordering of the servers and lower numbered servers take jobs before higher numbered ones.) This setup allows a job to be allocated to a server just before it begins service, rather than upon arrival to the system. We assume that the system is initially empty.

## 2.2 Stochastic Primitives

All random variables and stochastic processes used in our model description are assumed to be defined on a given complete probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ . The expectation operator under  $\mathbf{P}$  will be denoted by  $\mathbf{E}$ . For  $i \in \mathcal{I}$ , we take as given a sequence of strictly positive i.i.d. random variables  $\{u_i(\ell), \ell = 1, 2, \dots\}$  with mean  $\lambda_i^{-1} \in (0, \infty)$  and squared coefficient of variation (variance divided by the square of the mean)  $a_i^2 \in [0, \infty)$ . We interpret  $u_i(\ell)$  as the interarrival time between the  $(\ell - 1)$ <sup>st</sup> and the  $\ell$ <sup>th</sup> arrival to class  $i$  where, by convention, the “0<sup>th</sup> arrival” is assumed to occur at time zero. Setting  $\xi_i(0) = 0$  and

$$\xi_i(n) = \sum_{\ell=1}^n u_i(\ell), \quad n = 1, 2, \dots, \quad (2.3)$$

we define

$$A_i(t) = \sup\{n \geq 0 : \xi_i(n) \leq t\} \quad \text{for all } t \geq 0. \quad (2.4)$$

Then  $A_i(t)$  is the number of arrivals to class  $i$  that have occurred in  $[0, t]$ , and  $\lambda_i$  is the long run arrival rate to class  $i$ . For  $j \in \mathcal{J}$ , we take as given a sequence of strictly positive i.i.d. random variables  $\{v_j(\ell), \ell = 1, 2, \dots\}$  with mean  $\mu_j^{-1} \in (0, \infty)$  and squared coefficient of variation  $b_j^2 \in [0, \infty)$ . We interpret  $v_j(\ell)$  as the amount of service time required by the  $\ell$ <sup>th</sup> job processed by activity  $j$ , and  $\mu_j$  as the long run rate at which activity  $j$  could process its associated class of jobs  $i(j)$  if the associated server  $k(j)$  worked continuously and exclusively on this class. For  $j \in \mathcal{J}$ , let  $\eta_j(0) = 0$ ,

$$\eta_j(n) = \sum_{\ell=1}^n v_j(\ell), \quad n = 1, 2, \dots, \quad (2.5)$$

and

$$S_j(t) = \sup\{n \geq 0 : \eta_j(n) \leq t\} \quad \text{for all } t \geq 0. \quad (2.6)$$

Then  $S_j(t)$  is the number of jobs that activity  $j$  could process in  $[0, t]$  if the associated server worked continuously and exclusively on the associated class of jobs during this time interval. The interarrival time sequences  $\{u_i(\ell), \ell = 1, 2, \dots\}$ ,  $i \in \mathcal{I}$ , and service time sequences  $\{v_j(\ell), \ell = 1, 2, \dots\}$ ,  $j \in \mathcal{J}$ , are all assumed to be mutually independent.

## 2.3 Scheduling Control and Performance Measures

Scheduling control is exerted by specifying a  $\mathbf{J}$ -dimensional allocation process  $T = \{T(t), t \geq 0\}$  where

$$T(t) = (T_1(t), \dots, T_{\mathbf{J}}(t))' \quad \text{for } t \geq 0, \quad (2.7)$$

and  $T_j(t)$  is the cumulative amount of service time devoted to activity  $j \in \mathcal{J}$  by the associated server  $k(j)$  in the time interval  $[0, t]$ . Now  $T$  must satisfy certain properties that go along with its interpretation. Indeed, one could give a discrete-event type description of the properties that  $T$  must have, including any system specific constraints such as no preemption of service. However, for our analysis, we shall only need the properties of  $T$  described in (2.11)–(2.16) below.

Let

$$I(t) = \mathbf{1}t - \mathbf{A}T(t), \quad t \geq 0, \quad (2.8)$$

where  $\mathbf{1}$  is the  $\mathbf{K}$ -dimensional vector of all ones. Then for each  $k \in \mathcal{K}$ ,  $I_k(t)$  is interpreted as the cumulative amount of time that server  $k$  has been idle up to time  $t$ . A natural constraint on  $T$  is that each component of the cumulative idletime process  $I$  must be continuous and non-decreasing. Since each column of the matrix  $\mathbf{A}$  contains the number one exactly once, this immediately implies that each component of  $T$  is Lipschitz continuous with a Lipschitz constant of one. For each  $j \in \mathcal{J}$ ,  $S_j(T_j(t))$  is interpreted as the number of complete jobs processed by activity  $j$  in  $[0, t]$ . For  $i \in \mathcal{I}$ , let

$$Q_i(t) = A_i(t) - \sum_{j=1}^{\mathbf{J}} \mathbf{C}_{ij} S_j(T_j(t)), \quad t \geq 0, \quad (2.9)$$

which we write in vector form (with a slight abuse of notation for  $S(T(t))$ ) as

$$Q(t) = A(t) - \mathbf{C}S(T(t)), \quad t \geq 0. \quad (2.10)$$

Then  $Q_i(t)$  is interpreted as the number of class  $i$  jobs that are either in queue or “in progress” (i.e., being served or in suspension) at time  $t$ . We regard  $Q$  and  $I$  as performance measures for our system.

We shall use the following minimal set of properties of any scheduling control  $T$  with associated queue length process  $Q$  and idletime process  $I$ . For all  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$ ,  $k \in \mathcal{K}$ ,

$$T_j(t) \in \mathcal{F} \text{ for each } t \geq 0, \quad (2.11)$$

$$T_j \text{ is Lipschitz continuous with a Lipschitz constant of one,} \quad (2.12)$$

$$T_j \text{ is non-decreasing, and } T_j(0) = 0,$$

$$I_k \text{ is continuous, non-decreasing, and } I_k(0) = 0, \quad (2.13)$$

$$Q_i(t) \geq 0 \text{ for all } t \geq 0. \quad (2.14)$$

Properties (2.12) and (2.13) are for each sample path. For later reference, we collect here the queueing system equations satisfied by  $Q$  and  $I$ :

$$Q(t) = A(t) - \mathbf{C}S(T(t)), \quad t \geq 0, \quad (2.15)$$

$$I(t) = \mathbf{1}t - \mathbf{A}T(t), \quad t \geq 0, \quad (2.16)$$

where  $Q$ ,  $T$  and  $I$  satisfy properties (2.11)–(2.14). We emphasize that these are descriptive equations satisfied by the queueing system, given  $\mathbf{C}$ ,  $\mathbf{A}$ ,  $A$ ,  $S$  and a control  $T$ , which suffice for the purposes of our analysis. In particular, we do not intend them to be a complete, discrete-event type description of the dynamics.

**Remark.** The reader might expect that  $T$  should satisfy some additional non-anticipating property. Although this is a reasonable assumption to make, and indeed the policy we propose in Section 5 satisfies such a condition, we have not restricted  $T$  a priori in this way. Indeed, we shall see that, for the parallel server system under the complete resource pooling condition, our policy is asymptotically optimal even when anticipating policies are allowed. This is related to the fact that the Brownian control problem has a so-called “pathwise solution”, cf. [18].



The cost function we shall use involves linear holding costs associated with the expense of holding jobs of each class in the system until they have completed service. We defer the precise description of this cost function to the next section, since it is formulated in terms of normalized queue lengths, where the normalization is in diffusion scale. Indeed, in the next section, we describe the sequence of parallel server systems to be used in formulating the notion of heavy traffic asymptotic optimality.

### 3 Sequence of Systems, Heavy Traffic, and the Cost Function

For the parallel server system described in the last section, the problem of finding a control policy that minimizes a cost associated with holding jobs in the system is notoriously difficult. One possible means for discriminating between policies is to look for policies that outperform others in some asymptotic regime. Here we regard the parallel server system as a member of a sequence of systems indexed by  $r$  that is approaching heavy traffic (this notion is defined below). In this asymptotic regime, the queue length process is normalized with diffusive scaling – this corresponds to viewing the system over long intervals of time of order  $r^2$  (where  $r$  will tend to infinity in the asymptotic limit) and regarding a single job as only having a small contribution to the overall cost of storage, where this is quantified to be of order  $1/r$ . The setup in this section is a generalization of that used in [3].

#### 3.1 Sequence of Systems and Large Deviation Assumptions

Consider a sequence of parallel server systems indexed by  $r$ , where  $r$  tends to infinity through a sequence of values in  $[1, \infty)$ . These systems all have the same basic structure as that described in Section 2, except that the arrival and service rates, scheduling control, and form of the cost function (which is defined below in Section 3.3) are allowed to depend on  $r$ . Accordingly, we shall indicate the dependence of relevant parameters and processes on  $r$  by appending a superscript to them. We assume that the interarrival and service times are given for each  $r \geq 1$ ,  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$ , by

$$u_i^r(\ell) = \frac{1}{\lambda_i^r} \check{u}_i(\ell), \quad v_j^r(\ell) = \frac{1}{\mu_j^r} \check{v}_j(\ell), \quad \text{for } \ell = 1, 2, \dots, \quad (3.1)$$

where the  $\check{u}_i(\ell)$ ,  $\check{v}_j(\ell)$ , do not depend on  $r$ , have mean one and squared coefficients of variation  $a_i^2$ ,  $b_j^2$ , respectively. The sequences  $\{\check{u}_i(\ell), \ell = 1, 2, \dots\}$ ,  $\{\check{v}_j(\ell), \ell = 1, 2, \dots\}$ ,  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$  are all mutually independent sequences of i.i.d. random variables. (The above structure is a convenient means of allowing the sequence of systems to approach heavy traffic by simply changing arrival and service rates while keeping the underlying sources of variability  $\check{u}_i(\ell)$ ,  $\check{v}_j(\ell)$  fixed. This type of setup has been used previously by others in treating heavy traffic limits, see e.g., Peterson [34]. For a first reading, the reader may like to simply choose  $\lambda^r = \lambda$  and  $\mu^r = \mu$  for all  $r$ . Indeed, that simplification is made in the paper [42].)

We make the following assumption on the first order parameters associated with our sequence of systems.

**Assumption 3.1** *There are vectors  $\lambda \in \mathbb{R}_+^{\mathbf{I}}$ ,  $\mu \in \mathbb{R}_+^{\mathbf{J}}$ , such that*

- (i)  $\lambda_i > 0$  for all  $i \in \mathcal{I}$ ,  $\mu_j > 0$  for all  $j \in \mathcal{J}$ ,
- (ii)  $\lambda^r \rightarrow \lambda$ ,  $\mu^r \rightarrow \mu$ , as  $r \rightarrow \infty$ .

In addition, we make the following exponential moment assumptions to ensure that certain *large deviation estimates* hold for the renewal processes  $A_i^r$ ,  $i \in \mathcal{I}$ , and  $S_j^r$ ,  $j \in \mathcal{J}$  (cf. Lemma 6.7 below and Appendix A in [3]).

**Assumption 3.2** For  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$ , and all  $\ell \geq 1$ , let

$$u_i(\ell) = \frac{1}{\lambda_i} \tilde{u}_i(\ell), \quad v_j(\ell) = \frac{1}{\mu_j} \tilde{v}_j(\ell). \quad (3.2)$$

Assume that there is a non-empty open neighborhood  $\mathcal{O}_0$  of  $0 \in \mathbb{R}$  such that for all  $l \in \mathcal{O}_0$ ,

$$\Lambda_i^a(l) \equiv \log \mathbf{E}[e^{lu_i(1)}] < \infty, \text{ for all } i \in \mathcal{I}, \text{ and} \quad (3.3)$$

$$\Lambda_j^s(l) \equiv \log \mathbf{E}[e^{lv_j(1)}] < \infty, \text{ for all } j \in \mathcal{J}. \quad (3.4)$$

Note that (3.3) and (3.4) hold with  $\ell$  in place of 1 for all  $\ell = 1, 2, \dots$ , since  $\{u_i(\ell), \ell = 1, 2, \dots\}$ ,  $i \in \mathcal{I}$ , and  $\{v_j(\ell), \ell = 1, 2, \dots\}$ ,  $j \in \mathcal{J}$ , are each sequences of i.i.d. random variables.

**Remark.** This finiteness of exponential moments assumption allows us to prove asymptotic optimality of a threshold policy with thresholds of order  $\log r$ . We conjecture that this condition could be relaxed to a sufficiently high finite moment assumption, and our method of proof would still work, provided larger thresholds are used to allow for larger deviations of the primitive renewal processes from their rate processes. Here we have chosen the tradeoff of smaller thresholds and exponential moment assumptions, rather than larger thresholds and certain finite moment assumptions. We note that to adapt our proof to use finite moment assumptions, Lemma 6.7 would need to be modified to use estimates of the primitive renewal processes based on sufficiently high finite moments, rather than exponential moments (cf. [1, 5, 31, 36]). In addition, this lemma is used repeatedly in proving the main theoretical estimates embodied in Theorem 7.7. The latter is proved by induction and the proof involves a recursive increase in the size of the thresholds as well as an increase in the size of the error probability which involves powers of  $r^{2t}$ . Thus, in order to determine the necessary order of a finite moment condition and the size of accompanying thresholds, one would need to carefully examine this recursive proof.

### 3.2 Heavy Traffic and Fluid Model

There is no conventional notion of heavy traffic for our model, since the nominal (or average) load on a server depends on the scheduling policy. Harrison [16] (see also Laws [29] and Harrison and Van Mieghem [19]) has proposed a notion of heavy traffic for stochastic networks with scheduling control. For our sequence of parallel server systems, Harrison's condition is the same as Assumption 3.3 below. Here, and henceforth, we define

$$\mathbf{R} = \mathbf{C} \operatorname{diag}(\boldsymbol{\mu}).$$

**Assumption 3.3** There is a unique optimal solution  $(\rho^*, x^*, \cdot)$  of the linear program:

$$\text{minimize } \rho \quad \text{subject to } \mathbf{R}x = \boldsymbol{\lambda}, \quad \mathbf{A}x \leq \rho \mathbf{1} \quad \text{and } x \geq 0. \quad (3.5)$$

Moreover, that solution is such that  $\rho^* = 1$  and  $\mathbf{A}x^* = \mathbf{1}$ .

**Remark.** It will turn out that under Assumption 3.3,  $x_j^*$  is the average fraction of time that server  $k(j)$  should devote to activity  $j$ . For this reason,  $x^*$  is called the *nominal allocation vector*.

It was shown in [42] that Assumption 3.3 is equivalent to a heavy traffic condition for a fluid model (a formal law of large numbers approximation) associated with our sequence of parallel server systems. We summarize that result here since the fluid model plays a role in establishing asymptotic optimality of a control policy for our sequence of systems.

A *fluid model solution* (with zero initial condition) is a triple of continuous (deterministic) functions  $(\bar{Q}, \bar{T}, \bar{I})$  defined on  $[0, \infty)$ , where  $\bar{Q}$  takes values in  $\mathbb{R}^{\mathbf{I}}$ ,  $\bar{T}$  takes values in  $\mathbb{R}^{\mathbf{J}}$  and  $\bar{I}$  takes values in  $\mathbb{R}^{\mathbf{K}}$ , such that

$$\bar{Q}(t) = \lambda t - \mathbf{R}\bar{T}(t), \quad t \geq 0, \quad (3.6)$$

$$\bar{I}(t) = \mathbf{1}t - \mathbf{A}\bar{T}(t), \quad t \geq 0, \quad (3.7)$$

and for all  $i, j, k$ ,

$$\begin{aligned} \bar{T}_j &\text{ is Lipschitz continuous with a Lipschitz constant of one,} \\ &\text{it is non-decreasing, and } \bar{T}_j(0) = 0, \end{aligned} \quad (3.8)$$

$$\bar{I}_k \text{ is continuous, non-decreasing, and } \bar{I}_k(0) = 0, \quad (3.9)$$

$$\bar{Q}_i(t) \geq 0 \text{ for all } t \geq 0. \quad (3.10)$$

A continuous function  $\bar{T} : [0, \infty) \rightarrow \mathbb{R}^{\mathbf{J}}$  such that (3.8)–(3.10) hold for  $\bar{Q}, \bar{I}$  defined by (3.6)–(3.7) will be called a *fluid control*. For a given fluid control  $\bar{T}$ , we say the fluid system is *balanced* if the associated fluid “queue length”  $\bar{Q}$  does not change with time (cf. Harrison [13]). Here, since the system starts empty, that means  $\bar{Q} \equiv 0$ . In addition, we say the fluid system *incurs no idleness* (or all fluid servers are fully occupied) if  $\bar{I} \equiv 0$ , i.e.,  $\mathbf{A}\bar{T}(t) = \mathbf{1}t$  for all  $t \geq 0$ .

**Definition 3.4** *The fluid model is in heavy traffic if the following two conditions hold:*

- (i) *there is a unique fluid control  $\bar{T}^*$  under which the fluid system is balanced, and*
- (ii) *under  $\bar{T}^*$ , the fluid system incurs no idleness.*

Since any fluid control is differentiable at almost every time (by (3.8)), we can convert the above notion of heavy traffic into one involving the rates  $x^*(t) = \dot{\bar{T}}^*(t)$ , where “ $\dot{\cdot}$ ” denotes time derivative. This leads to the following lemma which is stated and proved in [42] (cf. Lemma 3.3 there).

**Lemma 3.5** *The fluid model is in heavy traffic if and only if Assumption 3.3 holds.*

We impose the following heavy traffic assumption on our sequence of parallel server systems, henceforth.

**Assumption 3.6 (Heavy Traffic)** *For the sequence of parallel server systems defined in Section 3.1 and satisfying Assumptions 3.1 and 3.2, assume that Assumption 3.3 holds and that there is a vector  $\theta \in \mathbb{R}^{\mathbf{I}}$  such that*

$$r(\lambda^r - \mathbf{R}^r x^*) \rightarrow \theta, \quad \text{as } r \rightarrow \infty, \quad (3.11)$$

where  $\mathbf{R}^r = \mathbf{C} \text{diag}(\mu^r)$ .

For the formulation of the Brownian control problem, it will be helpful to distinguish *basic activities*  $j$  which have a strictly positive nominal fluid allocation level  $x_j^*$  from *non-basic activities*  $j$  for which  $x_j^* = 0$ . By relabeling the activities if necessary, we may and do assume henceforth that  $x_j^* > 0$  for  $j = 1, \dots, \mathbf{B}$  and  $x_j^* = 0$  for  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ . Thus there are  $\mathbf{B}$  basic activities and  $\mathbf{J} - \mathbf{B}$  non-basic activities.

### 3.3 Diffusion Scaling and the Cost Function

For a fixed  $r$  and scheduling control  $T^r$ , the associated queue length process  $Q^r = (Q_1^r, \dots, Q_{\mathbf{I}}^r)'$  and idletime process  $I^r = (I_1^r, \dots, I_{\mathbf{K}}^r)'$  are given by (2.15)–(2.16) where the superscript  $r$  needs to be appended to  $A$ ,  $S$ ,  $Q$ ,  $I$  and  $T$  there. The diffusion scaled queue length process  $\hat{Q}^r$  and idletime process  $\hat{I}^r$  are defined by

$$\hat{Q}^r(t) = r^{-1}Q^r(r^2t), \quad \hat{I}^r(t) = r^{-1}I^r(r^2t), \quad t \geq 0. \quad (3.12)$$

We consider an expected cumulative discounted holding cost for the diffusion scaled queue length process and control  $T^r$ :

$$\hat{J}^r(T^r) = \mathbf{E} \left( \int_0^\infty e^{-\gamma t} h \cdot \hat{Q}^r(t) dt \right), \quad (3.13)$$

where  $\gamma > 0$  is a fixed constant (discount factor) and  $h = (h_1, \dots, h_{\mathbf{I}})'$ ,  $h_i > 0$  for all  $i \in \mathcal{I}$ , is a constant vector of holding costs per unit time per unit of diffusion scaled queue length. Recall that “ $\cdot$ ” denotes the dot product between two vectors.

To write equations for  $\hat{Q}^r, \hat{I}^r$ , we introduce centered and diffusion scaled versions  $\hat{A}^r, \hat{S}^r$  of the primitive processes  $A^r, S^r$ :

$$\hat{A}^r(t) = r^{-1} (A^r(r^2t) - \lambda^r r^2t), \quad \hat{S}^r(t) = r^{-1} (S^r(r^2t) - \mu^r r^2t), \quad (3.14)$$

a *deviation process*  $\hat{Y}^r$  (which measures normalized deviations of server time allocations from the nominal allocations given by  $x^*$ ):

$$\hat{Y}^r(t) = r^{-1} (x^* r^2t - T^r(r^2t)), \quad (3.15)$$

and a fluid scaled allocation process  $\bar{T}^r$ :

$$\bar{T}^r(t) = r^{-2}T^r(r^2t). \quad (3.16)$$

On substituting the above into (2.15)–(2.16), we obtain

$$\hat{Q}^r(t) = \hat{A}^r(t) - \mathbf{C}\hat{S}^r(\bar{T}^r(t)) + r(\lambda^r - \mathbf{R}^r x^*)t + \mathbf{R}^r \hat{Y}^r(t), \quad (3.17)$$

$$\hat{I}^r(t) = \mathbf{A}\hat{Y}^r(t), \quad (3.18)$$

where by (2.12)–(2.14) we have

$$\hat{I}_k^r \text{ is continuous, non-decreasing and } \hat{I}_k^r(0) = 0, \quad \text{for all } k \in \mathcal{K}, \quad (3.19)$$

$$\hat{Q}_i^r(t) \geq 0 \text{ for all } t \geq 0 \text{ and } i \in \mathcal{I}. \quad (3.20)$$

On combining Assumption 3.1 with the finite variance and mutual independence of the stochastic primitive sequences of i.i.d. random variables  $\{\tilde{u}_i(\ell)\}_{\ell=1}^\infty$ ,  $i \in \mathcal{I}$ ,  $\{\tilde{v}_j(\ell)\}_{\ell=1}^\infty$ ,  $j \in \mathcal{J}$ , we may deduce from renewal process functional central limit theorems (cf. [22]) that

$$(\hat{A}^r, \hat{S}^r) \Rightarrow (\tilde{A}, \tilde{S}), \quad \text{as } r \rightarrow \infty, \quad (3.21)$$

where  $\tilde{A}, \tilde{S}$  are independent,  $\tilde{A}$  is an  $\mathbf{I}$ -dimensional driftless Brownian motion that starts from the origin and has a diagonal covariance matrix whose  $i^{\text{th}}$  diagonal entry is  $\lambda_i a_i^2$ , and  $\tilde{S}$  is a  $\mathbf{J}$ -dimensional driftless Brownian motion that starts from the origin and has diagonal covariance matrix whose  $j^{\text{th}}$  diagonal entry is  $\mu_j b_j^2$ .

## 4 Brownian Control Problem

### 4.1 Formulation

Under the heavy traffic assumption of the previous section, to keep queue lengths from growing on average, it seems desirable to choose a control policy for the sequence of parallel server systems that asymptotically on average allocates service to the processing activities in accordance with the proportions given by  $x^*$ . To see how to achieve this and to do so in an optimal manner, following a method proposed by Harrison [12, 15, 18], we consider the following Brownian control problem which is a formal diffusion approximation to control problems for the sequence of parallel server systems. The relationship between the Brownian model and the fluid model is analogous to the relationship between the central limit theorem and the law of large numbers.

**Definition 4.1 (Brownian control problem)**

$$\text{minimize } \mathbf{E} \left( \int_0^\infty e^{-\gamma t} h \cdot \tilde{Q}(t) dt \right) \quad (4.1)$$

using a  $\mathbf{J}$ -dimensional control process  $\tilde{Y} = (\tilde{Y}_1, \dots, \tilde{Y}_{\mathbf{J}})'$  such that

$$\tilde{Q}(t) = \tilde{X}(t) + \mathbf{R}\tilde{Y}(t) \quad \text{for all } t \geq 0, \quad (4.2)$$

$$\tilde{I}(t) = \mathbf{A}\tilde{Y}(t) \quad \text{for all } t \geq 0, \quad (4.3)$$

$$\tilde{I}_k \text{ is non-decreasing and } \tilde{I}_k(0) \geq 0, \quad \text{for all } k \in \mathcal{K}, \quad (4.4)$$

$$\tilde{Y}_j \text{ is non-increasing and } \tilde{Y}_j(0) \leq 0, \quad \text{for } j = \mathbf{B} + 1, \dots, \mathbf{J}, \quad (4.5)$$

$$\tilde{Q}_i(t) \geq 0 \quad \text{for all } t \geq 0, \quad i \in \mathcal{I}, \quad (4.6)$$

where  $\tilde{X}$  is an  $\mathbf{I}$ -dimensional Brownian motion that starts from the origin, has drift  $\theta$  (cf. (3.11)) and a diagonal covariance matrix whose  $i^{\text{th}}$  diagonal entry is equal to  $\lambda_i a_i^2 + \sum_{j=1}^{\mathbf{J}} \mathbf{C}_{ij} \mu_j b_j^2 x_j^*$  for  $i \in \mathcal{I}$ .

The above Brownian control problem is a slight variant of that used by Harrison and López [18]. In particular, we allow  $\tilde{Y}$  to anticipate the future of  $\tilde{X}$ . The process  $\tilde{X}$  is the formal limit in distribution of  $\hat{X}^r$ , where for  $t \geq 0$ ,

$$\hat{X}^r(t) \equiv \hat{A}^r(t) - \mathbf{C}\hat{S}^r(\bar{T}^r(t)) + r(\lambda^r - \mathbf{R}^r x^*)t. \quad (4.7)$$

The control process  $\tilde{Y}$  in the Brownian control problem is a formal limit of the deviation processes  $\hat{Y}^r$ . (cf. (3.15)). The non-increasing assumption in property (4.5) corresponds to the fact that  $\hat{Y}_j^r(t) = -r^{-1}T_j^r(r^2t)$  is non-increasing whenever  $j$  is a non-basic activity. The initial conditions on  $\tilde{I}$  and  $\tilde{Y}_j$ ,  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ , in (4.4)–(4.5) are relaxed from those in the prelimit to allow for the possibility of an initial jump in the queue length in the Brownian control problem. (In fact, for the optimal solution of the Brownian control problem, under the complete resource pooling condition to be assumed later, such a jump will not occur and then the Brownian control problem is equivalent to one in which  $\tilde{I}(0) = 0$ ,  $\tilde{Y}(0) = 0$ .)

### 4.2 Solution via Workload Assuming Complete Resource Pooling

For open processing networks, of which parallel server systems are a special case, Harrison and Van Mieghem [19] showed that the Brownian control problem can be reduced to an equivalent formulation involving a typically lower dimensional state process. In this reduction, the Brownian analogue  $\tilde{Q}$  of the queue length process is replaced by a Brownian analogue  $\tilde{W} = M\tilde{Q}$  of the workload process, where  $M$  is a certain matrix called the workload matrix. The row dimension  $\mathbf{L} \leq \mathbf{I}$  of  $M$  is

called the workload dimension and the reduced form of the Brownian control problem is called an equivalent workload formulation. Intuitively, the reduction in [19] is achieved by ignoring certain “reversible directions” in which the process  $\tilde{Q}$  can be controlled to move instantaneously without incurring any idleness and without using any non-basic activities; in addition, such movements are instantaneously reversible. In a work [16] following on from [19], Harrison elaborated on an interpretation of the reversible displacements and proposed a “canonical” choice for the workload matrix  $M$ , which reduces the possibilities for  $M$  to a finite set. This choice uses extremal optimal solutions to the dual to the linear program used to define heavy traffic, cf. (3.5). For the parallel server situation considered here, that dual program has the following form.

### Dual Program

$$\text{maximize } y \cdot \lambda \quad \text{subject to } y' \mathbf{R} \leq z' \mathbf{A}, \quad z \cdot \mathbf{1} \leq 1 \quad \text{and } z \geq 0. \quad (4.8)$$

We shall focus on the case in which the workload dimension is one. This is equivalent to there being a unique optimal solution of the dual program (4.8). Indeed, for parallel server systems, Theorem 4.3 below gives several conditions that are equivalent to this assumption of a one-dimensional workload process. For the statement of this result, we need the following notion of communicating servers.

**Definition 4.2** Consider the graph  $\mathcal{G}$  in which servers and buffers form the nodes and (undirected) edges between nodes are given by basic activities. We say that all servers communicate via basic activities if, for each pair of servers, there is a path in  $\mathcal{G}$  joining all of the servers.

**Theorem 4.3** The following conditions are equivalent (for parallel server systems).

- (i) the workload dimension  $\mathbf{L}$  is one,
- (ii) the dual program (4.8) has a unique optimal solution  $(y^*, z^*)$ ,
- (iii) the number of basic activities  $\mathbf{B}$  is equal to  $\mathbf{I} + \mathbf{K} - 1$ ,
- (iv) all servers communicate via basic activities,
- (v) the graph  $\mathcal{G}$  is a tree.

**Proof.** The equivalence of the first four statements of the theorem was derived in Harrison and López [18], and the equivalence of these statements with (v) was noted in [42]. Also, as noted by Ata and Kumar [1], the equivalence of (i) with (iii) can also be seen by observing from Corollary 6.2 in Bramson and Williams [6] that since the matrix  $A$  has exactly one positive entry in each column, the workload dimension  $\mathbf{L} = \mathbf{I} + \mathbf{K} - \mathbf{B}$ , where  $\mathbf{B}$  is the number of basic activities.  $\square$

Henceforth we make the following assumption.

**Assumption 4.4** (Complete Resource Pooling) The equivalent conditions (i)–(v) of Theorem 4.3 hold.

The intuition behind the term “complete resource pooling” is that, under this condition, all servers can communicate via basic activities, and it is reasonable to conjecture that the efforts of the servers (or resources) can be combined in a cooperative manner so that they act effectively as a single pooled resource. A main aim of this work is to show that this intuition is correct.

Let  $(y^*, z^*)$  be the unique optimal solution of (4.8). By complementary slackness,  $((y^*)' \mathbf{R})_j = ((z^*)' \mathbf{A})_j$  for  $j = 1, \dots, \mathbf{B}$ . Let  $u^*$  be the  $(\mathbf{J} - \mathbf{B})$ -dimensional vector of dual “slack variables” defined by  $((y^*)' \mathbf{R} - (z^*)' \mathbf{A})_j + u_{j-\mathbf{B}}^* = 0$  for  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ .

**Lemma 4.5** We have  $y^* > 0$ ,  $z^* > 0$ ,  $u^* > 0$ ,

$$(y^*)' \mathbf{R} = (z^*)' \mathbf{A} - [0' (u^*)'], \quad \text{and} \quad z^* \cdot \mathbf{1} = 1, \quad (4.9)$$

where  $\mathbf{1}$  is the  $\mathbf{K}$ -dimensional vector of ones,  $0'$  is a  $\mathbf{B}$ -dimensional row vector of zeros.

**Proof.** This is proved in [18] (Corollary to Proposition 3) using the relation between the primal and dual linear programs.  $\square$

Now, we review a solution of the Brownian control problem obtained by Harrison and López [18]. This exploits the fact that the workload is one-dimensional.

For  $\tilde{Q}$  satisfying (4.2)–(4.6), define  $\tilde{W} = y^* \cdot \tilde{Q}$ , which Harrison [16] calls the (Brownian) workload. Let  $\tilde{Y}_N$  be the  $(\mathbf{J} - \mathbf{B})$ -dimensional process whose components are  $\tilde{Y}_j$ ,  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ . By Lemma 4.5 and (4.2)–(4.6),

$$\tilde{W}(t) = y^* \cdot \tilde{X}(t) + \tilde{V}(t) \quad \text{for all } t \geq 0, \quad (4.10)$$

where

$$\tilde{V} \equiv z^* \cdot \tilde{I} - u^* \cdot \tilde{Y}_N \text{ is non-decreasing and } \tilde{V}(0) \geq 0, \quad (4.11)$$

$$\tilde{W}(t) \geq 0 \text{ for all } t \geq 0. \quad (4.12)$$

Now, for each  $t \geq 0$ , since the holding cost vector  $h > 0$  and  $y^* > 0$ , we have

$$h \cdot \tilde{Q}(t) = \sum_{i=1}^{\mathbf{I}} \left( \frac{h_i}{y_i^*} \right) y_i^* \tilde{Q}_i(t) \geq h^* \tilde{W}(t) \quad (4.13)$$

where

$$h^* \equiv \min_{i=1}^{\mathbf{I}} \left( \frac{h_i}{y_i^*} \right). \quad (4.14)$$

It is well-known and straightforward to see that any solution pair  $(\tilde{W}, \tilde{V})$  of (4.10)–(4.12) must satisfy for all  $t \geq 0$ ,

$$\tilde{V}(t) \geq \tilde{V}^*(t) \equiv \sup_{0 \leq s \leq t} \left( -y^* \cdot \tilde{X}(s) \right), \quad (4.15)$$

and hence  $\tilde{W}(t) \geq \tilde{W}^*(t)$  where

$$\tilde{W}^*(t) = y^* \cdot \tilde{X}(t) + \tilde{V}^*(t). \quad (4.16)$$

The process  $\tilde{W}^*$  is a one-dimensional *reflected Brownian motion* driven by the one-dimensional Brownian motion  $y^* \cdot \tilde{X}$ , and  $\tilde{V}^*$  is its local time at zero (see e.g., [9], Chapter 8). In particular,  $\tilde{V}^*$  can have a point of increase at time  $t$  only if  $\tilde{W}^*(t) = 0$ .

Now, let  $i^*$  be a class index such that  $h^* = h_{i^*}/y_{i^*}^*$ , i.e., the minimum in (4.14) is achieved at  $i = i^*$ , and let  $k^*$  be a server that can serve class  $i^*$  via a basic activity. Note that neither  $i^*$  nor  $k^*$  need be unique in general. Then the following choices  $\tilde{Q}^*$  and  $\tilde{I}^*$  for  $\tilde{Q}$  and  $\tilde{I}$  ensure that for each  $t \geq 0$ , properties (4.4)–(4.6) hold and the inequality in (4.13) is an equality with  $\tilde{W}(t) = \tilde{W}^*(t)$  there:

$$\tilde{Q}_{i^*}^*(t) = \tilde{W}^*(t)/y_{i^*}^*, \quad \tilde{Q}_i^*(t) = 0 \text{ for all } i \neq i^*, \quad (4.17)$$

$$\tilde{I}_{k^*}^*(t) = \tilde{V}^*(t)/z_{k^*}^*, \quad \tilde{I}_k^*(t) = 0 \text{ for } k \neq k^*, \quad \tilde{Y}_N^* = \mathbf{0}. \quad (4.18)$$

A control process  $\tilde{Y}^*$  such that (4.2)–(4.6) hold with  $\tilde{Q}^*, \tilde{Y}^*, \tilde{I}^*$  in place of  $\tilde{Q}, \tilde{Y}, \tilde{I}$  there is given in [42]. It can be readily verified that this is an optimal solution for the Brownian control problem (cf. [18]) and the associated minimum cost is

$$J^* \equiv \mathbf{E} \left( \int_0^\infty e^{-\gamma t} h \cdot \tilde{Q}^*(t) dt \right) = h^* \mathbf{E} \left( \int_0^\infty e^{-\gamma t} \tilde{W}^*(t) dt \right). \quad (4.19)$$

The quantity  $J^*$  is finite and can be computed as in Section 5.3 of [11].

Now, even though the Brownian control problem can be analyzed exactly (as above), the solution obtained does not automatically translate to a policy for the sequence of parallel server systems. However, some desirable features are suggested by the form of (4.17)–(4.18), namely, (a) try to keep the bulk of the work in the class  $i^*$  with the lowest (or equal lowest) ratio of holding cost to workload contribution, i.e., the class  $i$  with the lowest value of  $h_i/y_i^*$ , (b) try to ensure that the bulk of the idleness is incurred only when there is almost no work in the entire system, and (c) try to ensure that the bulk of the idletime is incurred by server  $k^*$  alone.

### 4.3 Approaches to Interpreting the Solution

Harrison [14] has proposed a general scheme (called BIGSTEP) for obtaining candidate policies for a queueing control problem from a solution of the associated Brownian control problem. The policies obtained in this manner are so-called discrete-review policies which allow review of the system status and changes in the control rule only at a fixed discrete set of times. For a two-server example (with Poisson arrivals, deterministic service times and particular values for  $\lambda, \mu$ ), a discrete-review policy was constructed and shown to be asymptotically optimal in Harrison [15]. Based on their solution of the Brownian control problem and the general scheme laid out by Harrison [14], Harrison and López [18] proposed the use of a discrete review policy for the multiserver problem considered here, but they did not prove asymptotic optimality of this policy. Recently, Ata and Kumar [1] have proved asymptotic optimality of a discrete review policy for open stochastic processing networks that include parallel server systems, under heavy traffic and complete resource pooling conditions.

Another approach to translation of solutions of Brownian control problems into viable policies has been proposed by Kushner et al. (see e.g., [25, 26, 27]). However this also involves discretization by way of numerical approximation. We note that, of the works by Kushner et al. mentioned above, the paper by Kushner and Chen [25] is the closest to the current one in that it considers a parallel server model. However, it is in a very different parameter regime, namely one that corresponds to heavy traffic but with *no resource pooling*.

Assuming the complete resource pooling condition, in the next section we describe a simple “continuous review” policy for the sequence of parallel server systems, which allows changes in the control to be made at random times and in particular at times when the system status changes. This policy is a dynamic priority policy in which priorities for certain “transition” activities depend on the number of jobs in the associated class relative to certain threshold or “safety-stock” levels. Changes in the priorities only occur as a threshold is crossed. This threshold policy was proposed in [42]. A few more details of its description are given here to facilitate our analysis. We prove in Section 9 that this policy is asymptotically optimal. In [3], we have already proved that this is so for the special case of a two-server two-buffer system. An important feature of that proof was that the limiting (under diffusion scale) queue length and idleness processes were effectively one-dimensional, i.e., a form of state-space collapse occurred in the diffusion limit. A similar phenomenon occurs here for our multiserver system (cf. Theorem 5.3).



## 5 Threshold Policy and Main Results

In this section, we describe a dynamic threshold policy for our sequence of parallel server systems and state our main results, which imply that this policy is asymptotically optimal. Recall that we are assuming throughout that the heavy traffic (Assumption 3.6) and complete resource pooling (Assumption 4.4) conditions hold. The threshold policy takes advantage of the tree structure of the server-buffer tree  $\mathcal{G}$ . In reviewing our description of the policy, the reader may find it helpful to refer to the examples in [42], where a version of this policy was first described. We also include two additional examples here in Section 5.4 to help explain the policy. We note that there can be many asymptotically optimal policies. The one described here is simply proposed as one that is intuitively appealing and that is asymptotically optimal.

Independently of [42], Squillante et al. [35] proposed a tree-based threshold priority policy for a parallel server system. However, their policy is different from the one described here. The second example we give in Section 5.4 is used to illustrate this.

### 5.1 Tree Conventions

The threshold policy only involves the use of basic activities (and so in the following description of the policy, the word “activity” will be synonymous with “basic activity”). Also, the terms class and buffer will be used interchangeably. A key to the description of the policy is a hierarchical structure of the server-buffer tree  $\mathcal{G}$  and an associated protocol for the dynamic allocation of (preemptive-resume) class priorities at each server. This protocol is described in an iterative manner, working from the bottom of the tree up towards the root. (One should imagine a tree as growing downwards from its root and the root as being at the highest level.) A server tree  $\mathcal{S}$ , which results from suppressing the buffers in  $\mathcal{G}$ , will be helpful for describing the iterative procedure. Recall the solution of the Brownian control problem described in Section 4.2. The root of the server-buffer tree  $\mathcal{G}$  (and of the server tree  $\mathcal{S}$ ) is taken to be a server  $k^*$  which serves the “cheapest” class  $i^*$  via a (basic) activity. Classes (or buffers) that link one level of servers to those at the next highest level in the tree  $\mathcal{G}$  are called *transition classes* (or transition buffers).

### 5.2 Threshold Policy

To describe the threshold policy, we first focus on the server tree  $\mathcal{S}$  and imagine it arranged in levels with the root  $k^*$  at the highest level of  $\mathcal{S}$ . We proceed inductively up through the levels in the tree  $\mathcal{S}$ .

First, consider a server at the lowest level. As a server within the server-buffer tree  $\mathcal{G}$ , this server is to service its classes according to a priority scheme that gives lowest priority to the class that is immediately above the server in  $\mathcal{G}$ . The latter class is also served by a server in the next level up in  $\mathcal{S}$  and so is a transition class. There will always be such a class unless the server is at the root of the tree. On the other hand there may be zero, one or more classes lying below a server at the lowest level in the tree. The relative priority ranking of these classes (if any) is not so important. These are all terminal classes in that there are no servers below them. Here, for concreteness, we rank the classes so that for a given server, the lower numbered classes receive priority over the higher numbered ones. For future use, we place a threshold on the transition class immediately above each server in the lowest level of  $\mathcal{S}$ .

Now go to the next level up in the server tree  $\mathcal{S}$ . This level may have “terminal” server nodes and server nodes that lead to server nodes lower down the server tree. As servers in the server-buffer tree  $\mathcal{G}$ , each server at this level performs its activities in the following prioritized manner. Activities leading (via transition buffers) to server nodes lower down the tree are given highest priority (if there is more than one such activity, rank the activities so that activities serving lower numbered

classes are served first). However, if the number of jobs in a transition class associated with such an activity is at or below the threshold for that class, service of that activity is suspended. The next priority is given to activities that service (terminal) classes that are only served by that server (again ranked according to a scheme that gives lower numbered classes higher priority), and lowest priority is given to the activity serving the transition class that is immediately above the server in the server-buffer tree. This transition class should again have a threshold placed on it such that service of that class by the server at the next highest level will be suspended when the number of jobs in the class is equal to or below the threshold level. If two or more servers simultaneously begin to serve a particular transition buffer, a tie breaking rule is used to decide which server takes a job first. For concreteness we suppose that the lowest numbered server shall select a job before the next higher numbered server, and so on. Note that two servers in the same level cannot both serve the same buffer below them since  $\mathcal{G}$  is a tree.

This process is repeated until the root of the server tree is reached. At the root, the same procedure is applied as for lower server levels, except that there are no activities above the root server in the server-buffer tree and an overriding rule is that lowest priority is given to the “cheapest” class  $i^*$ .

The idea behind the threshold policy is to keep servers below the root level busy the bulk of the time (indeed, they should only be rarely idle and their idletimes should vanish on diffusion scale as the heavy traffic limit is approached), while simultaneously preventing the queue lengths of all classes except  $i^*$  from growing appreciably on diffusion scale. Transition buffers are used to achieve these two competing goals. Intuitively, when the queue length for a transition class gets to or below its threshold, any service of that class by the server immediately above it in the tree is suspended and this causes temporary overloading (on average) of the servers below, which prevents these servers from incurring much idleness. When the queue length for the transition class builds up to a level above its threshold, then assistance from the higher server is again permitted and the servers below are temporarily underloaded (on average) and so queue lengths for classes serviced by these servers are prevented from growing too large. The intended effect of this policy is to allow the movement (via the transition buffers) of excess work from lower level to higher level buffers (and eventually, by an upwards cascade, to the buffer for class  $i^*$ ), while simultaneously keeping all servers busy the bulk of the time, unless the entire system is nearly empty and even then to ensure that the vast majority of the idleness is incurred by the server  $k^*$  at the root of the tree.

### 5.3 Threshold Sizes

For each  $r$ , the size of the thresholds in the  $r^{\text{th}}$  parallel server system is to be of order  $\log r$ . However, while each threshold is of order  $\log r$ , for our proofs, the threshold sizes need to increase moderately as one moves up and across the tree to compensate for an associated accumulation of stochastic variability. This is related to the hierarchical structure of the threshold policy under which allocations to activities associated with transition buffers higher up the tree can depend on allocations to activities much farther down the tree. These transition buffers need larger thresholds than their counterparts below to allow enough time for their allocation processes to approach their long term averages before the associated queue lengths approach zero or twice the threshold size. For similar reasons, the threshold size for a buffer belonging to a group of transition buffers served by one server from above, should be larger the lower the priority of the buffer. To facilitate the description of this increase in threshold sizes, in Section 6.1 we show how the buffers can be renumbered, so that higher priority buffers for a given server have lower numbers and buffers higher up the tree are assigned higher numbers. This renumbering does not change the threshold policy, it simply allows us to streamline its detailed description. In particular, under this scheme, the size of the threshold for each transition buffer increases with its numbering. A detailed specification of the size of the thresholds is given in Section 6.2.

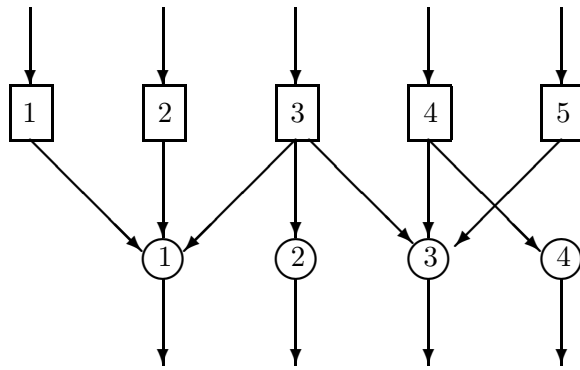


Figure 1: A parallel server system with four servers and five classes. Only basic activities are shown.

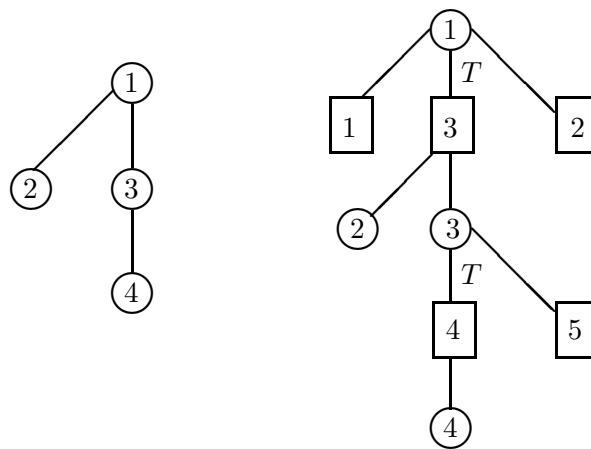


Figure 2: The server tree (on the left) and server-buffer tree (on the right) for the system pictured in Figure 1 with  $i^* = 2$  and  $k^* = 1$ . Activities subject to thresholding have a  $T$  beside them.

## 5.4 Examples

In this subsection we give two examples to illustrate our threshold policy. The second example is also used to illustrate that our policy differs from that of [35].

**Example 5.1** Consider the parallel server system pictured in Figure 1, where only basic activities are shown. Suppose for this that  $i^* = 2$  and hence  $k^* = 1$ . The server tree and server-buffer tree, with server 1 as root, are depicted in Figure 2. Activities subject to thresholding have a  $T$  beside them.

In the implementation of our policy, server 4 serves jobs from class 4 whenever there are class 4 jobs available for it to serve. Server 3 gives first priority to jobs in the transition class 4, except that when the number of class 4 jobs becomes equal to or goes below the threshold for that class, server 3 suspends service of the current job from that class and switches attention to class 5 jobs until those are exhausted and then finally server 3 turns its attention to class 3 jobs. Whenever the number of class 4 jobs again exceeds the threshold, server 3 preemptively switches its attention back to serving that class. Server 2 serves class 3 jobs whenever there are such jobs to be served and idles otherwise. Server 1 gives first priority to the transition class 3, except that server 1

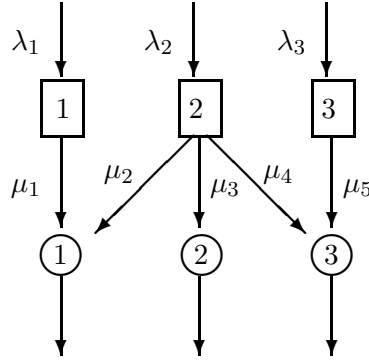


Figure 3: A parallel server system with three servers and three classes (or buffers).

suspends service of a job of class 3 whenever the number of jobs in that class reaches or goes below the threshold for the class. In that case, server 1 turns its attention to class 1 jobs until there are no jobs of that class to serve and then finally server 1 turns to serving jobs of the lowest priority class  $i^* = 2$ . Whenever the number of jobs in class 3 again goes above the threshold for that class, server 1 switches attention back to serving that class.

The policy proposed by Squillante et al. [35] involves the allocation of priorities for each server based on a local optimization of  $c\mu$ -type with suspension of certain thresholded activities when the associated class levels are at or below threshold. Their policy locally takes account of the tree structure, whereas our policy exploits the global tree structure. The following example illustrates how the two policies can yield different priority rules.

**Example 5.2** Consider the three-server, three-class parallel server system pictured in Figure 3. This example was considered by Squillante et al. [35] with two different sets of parameters in their Examples 4 and 5. For their Example 4 (using our notation),  $\lambda = (0.4, 1.4, 0.6)$  and  $\mu = (1, 0.5, 1, 0.25, 1)$ . Here

$$C = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \quad A = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}.$$

The heavy traffic Assumption 3.3 is readily verified to hold, where  $x^* = (0.4, 0.6, 1, 0.4, 0.6)$ , and hence all activities are basic. Since  $\mathbf{B} = 5 = \mathbf{I} + \mathbf{K} - 1$ , it follows from Theorem 4.3 that the complete resource pooling Assumption 4.4 holds. Noting Lemma 4.5, the unique solution of the dual program (4.8) is readily verified to be  $z^* = y^* = (2/7, 4/7, 1/7)$ . Let  $h = (1, 100, 10)$ . Then  $(h_i/y_i^*, i = 1, 2, 3) = (7/2, 175, 70)$  and so  $h^* = 7/2$ ,  $i^* = 1$  and  $k^* = 1$ , since class 1 is only served by server 1.

Under our threshold policy, server 3 gives priority to class 3 jobs over those of class 2, server 2 always serves class 2 jobs when there are such available, and service of class 2 by server 1 is subject to a threshold in the sense that server 1 gives priority to class 2 except when the number of jobs in class 2 is at or below its threshold and then server 1 suspends service of any class 2 job it is processing and switches attention to serving class 1 jobs until the number of jobs in class 2 once again exceeds the threshold for class 2 (cf. Figure 4).

The algorithm described by Squillante et al. [35] yields a similar policy for this example, with the exception that under their policy, activity 4, which links server 3 and class 2, is also subject to

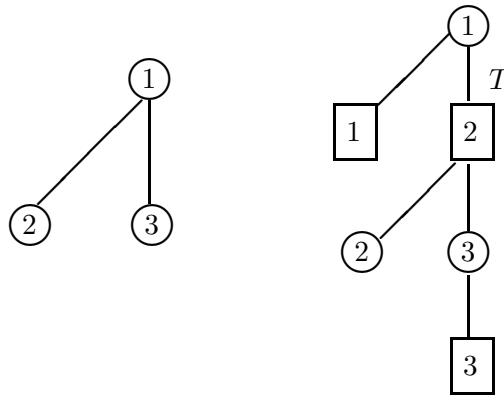


Figure 4: The server tree (on the left) and server-buffer tree (on the right) for the system pictured in Figure 3 with the data described for Example 5.2. The activity subject to thresholding under our policy is indicated with a  $T$ .

thresholding, such that if the number of jobs in class 2 is above threshold, server 3 serves jobs from class 2, whereas when the number of jobs in class 2 is at or below threshold, server 3 serves jobs from class 3.

## 5.5 Main Results

In Sections 7 and 8, we prove the following limit theorem for a certain sequence of allocation processes  $\{T^{r,*}\}$ . For each  $r$ ,  $T^{r,*}$  is obtained by applying the aforementioned threshold policy in the  $r^{\text{th}}$  parallel server system. The sizes of the thresholds on the transition buffers are of order  $\log r$ . The specific threshold sizes, which, as mentioned above, increase as one moves across and up the server-buffer tree, are specified precisely in the next section (cf. (6.4)).

**Theorem 5.3** Consider the sequence of parallel server systems indexed by  $r$ , where the  $r^{\text{th}}$  system operates under the allocation process  $T^{r,*}$  described above. Then the associated normalized queue length and idletime processes satisfy

$$(\hat{Q}^r, \hat{I}^r) \Rightarrow (\tilde{Q}^*, \tilde{I}^*), \quad \text{as } r \rightarrow \infty, \quad (5.1)$$

where  $\tilde{Q}_i^* = \mathbf{0}$  for all  $i \in \mathcal{I} \setminus \{i^*\}$ ,  $\tilde{I}_k^* = \mathbf{0}$  for all  $k \in \mathcal{K} \setminus \{k^*\}$ ,  $\tilde{Q}_{i^*}^*$  is a one-dimensional reflected Brownian motion that starts from the origin and has drift  $(y^* \cdot \theta)/y_{i^*}^*$  and variance parameter  $\sum_{i=1}^{\mathbf{I}} (y_i^*)^2 (\lambda_i a_i^2 + \sum_{j=1}^{\mathbf{J}} \mathbf{C}_{ij} \mu_j b_j^2 x_j^*) / (y_{i^*}^*)^2$ , and  $\tilde{I}_{k^*}^*$  is a specific multiple of the local time at the origin of  $\tilde{Q}_{i^*}^*$ . (In fact,  $\tilde{Q}_{i^*}^*$ ,  $\tilde{I}_{k^*}^*$  are equivalent in law to the processes given by (4.15)–(4.18)).

Recall the definitions of  $J^*$  and  $\hat{J}^r$  from (4.19) and (3.13), respectively. The following theorem is the main result of this paper. It is proved in Section 9 using Theorem 5.3. It shows that  $J^*$  is the best that one can achieve asymptotically and that this asymptotically minimal cost is achieved by the sequence of dynamic controls  $\{T^{r,*}\}$ . Thus we conclude that our threshold policy is asymptotically optimal.

**Theorem 5.4** (Asymptotic Optimality) Suppose that  $\{T^r\}$  is any sequence of scheduling controls (one for each member of the sequence of parallel server systems). Then

$$\liminf_{r \rightarrow \infty} \hat{J}^r(T^r) \geq J^* = \lim_{r \rightarrow \infty} \hat{J}^r(T^{r,*}), \quad (5.2)$$

and  $J^* < \infty$ .

**Remark.** Note that our threshold policy prescribes that a threshold be placed on class  $i^*$  if that class is a transition class. We conjecture that the policy which removes the threshold in this case is also asymptotically optimal, but we keep the threshold here as a means to simplify our proof. In particular, servers below  $i^*$  may experience significant idletime if the threshold is removed. Our threshold policy also involves preemption of service. There is a corresponding policy without preemption that we conjecture has the same behavior in the heavy traffic limit, since in that regime a maximum of  $\mathbf{J}$  jobs (in suspension or not) should not impact the asymptotic behavior of the system.

## 6 Preliminaries and Outline of the Proof

In this section, we will precisely specify the threshold sizes used for  $\{T^{r,*}\}$ , give some preliminary definitions and results, and outline the proofs of Theorems 5.3 and 5.4. Details of the proofs are contained in Sections 7–9.

Recall from Section 2 that  $\mathcal{I}$ ,  $\mathcal{K}$ , and  $\mathcal{J}$  index job classes, servers, and activities, respectively.

**Basic Activities Convention.** *Since our threshold policy only uses basic activities, to simplify notation, in this section and the next (Sections 6 and 7) only, the index set  $\mathcal{J}$  will just include the basic activities  $1, 2, \dots, \mathbf{B}$ . With this convention,  $\mathbf{J}$  will be synonymous with  $\mathbf{B}$ .*

### 6.1 The Server-Buffer Tree $\mathcal{G}$ : Layers and Buffer Renumbering

To facilitate our proof, we refine our description of the server-buffer tree,  $\mathcal{G}$ . We say that the server-buffer tree consists of one or more *layers*, where a layer consists of a server level along with the buffer level immediately below it. (At the lowest layer, the buffer level may be empty.) Activities that serve the buffers (from above and below) in a particular layer are also considered part of that same layer. We denote the number of layers by  $l^*$  and enumerate the layers in increasing order as one moves up the tree, so that layer 1 is the lowest layer and layer  $l^*$  is the top layer, consisting of server  $k^*$ , the buffers served by this server (including buffer  $i^*$ ), and the corresponding activities. Assuming  $l > 1$ , Figure 5 depicts layers  $l$  and  $l - 1$  in a server-buffer tree.

For each layer  $l = 1, \dots, l^*$ , we denote the collection of servers in layer  $l$  by  $\mathcal{K}^l$  and the collection of buffers in layer  $l$  by  $\mathcal{I}^l$ . For  $k \in \mathcal{K}^l$ , the collection of buffers in layer  $l$  served by server  $k$  is denoted by  $\underline{\mathcal{I}}_k$  (in Figure 5,  $i \in \underline{\mathcal{I}}_k$ ). The activity that serves buffer  $i \in \underline{\mathcal{I}}_k$  using server  $k$  is labeled  $a(i)$  (“a” is mnemonic for “above”), the collection of activities that serve buffer  $i$  from below is denoted by  $\underline{\mathcal{J}}_i$ , and the collection of servers in layer  $l - 1$  that serve buffer  $i$  using the activities in  $\underline{\mathcal{J}}_i$  is denoted by  $\underline{\mathcal{K}}_i$  (the underscores indicate that the quantities are “below”). We let  $\mathcal{J}_i$  denote the collection of all activities that service buffer  $i$ , i.e., activity  $a(i)$  together with  $\underline{\mathcal{J}}_i$ . Note that for each buffer  $i \in \mathcal{I}^1$ ,  $\underline{\mathcal{J}}_i = \emptyset$ ,  $\underline{\mathcal{K}}_i = \emptyset$ , i.e., layer 1 does not contain any activities that serve buffers in this layer from below. In fact, we may even have  $\mathcal{I}^1 = \emptyset$ .

Without loss of generality, we assume the following left-to-right arrangement and renumbering convention for buffers in the server-buffer tree  $\mathcal{G}$ . This simplifies the description of the priorities associated with our threshold policy. (The simplest way to think of doing the arrangement is to work from the top of the tree downwards.)

**Arrangement Convention.** *For  $k \in \mathcal{K} \setminus \{k^*\}$ , the buffers in  $\underline{\mathcal{I}}_k$  are arranged so that the transition buffers are positioned to the left of the non-transition buffers, and within the groups of transition and non-transition buffers, the lower numbered buffers are to the left of the higher numbered ones. For  $k = k^*$ , the arrangement in the previous sentence holds with the exception that buffer  $i^*$  is placed at the far right in level  $\mathcal{I}^{l^*} = \underline{\mathcal{I}}_{k^*}$ .*

**Renumbering Convention.** *Starting from layer 1, if  $\mathcal{I}^1 \neq \emptyset$ , we enumerate the buffers in  $\mathcal{I}^1$  in increasing order from left to right, i.e., the buffer farthest to the left is buffer 1 and the one farthest*

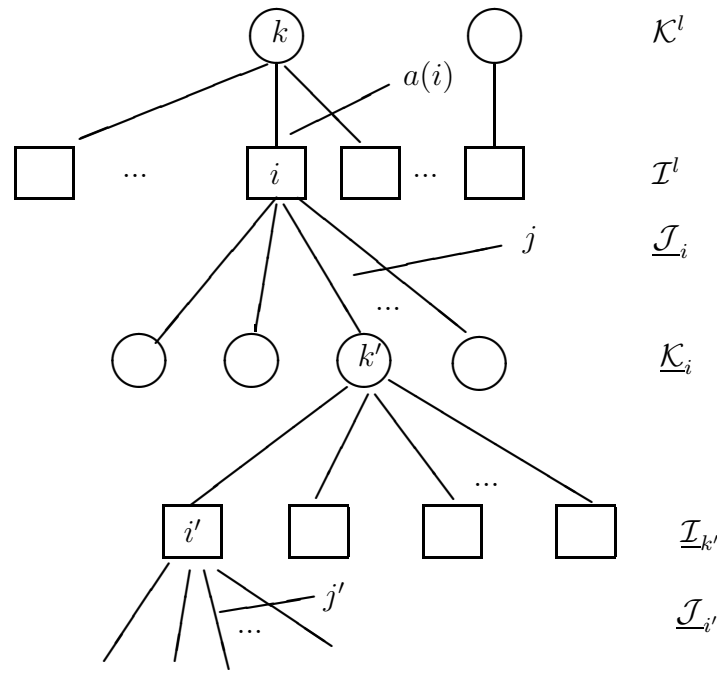


Figure 5: Layers  $l$  and  $l - 1$  of a server-buffer tree.

to the right is buffer  $|\mathcal{I}^1|$ . Continuing with layer 2, the buffer farthest to the left is labeled  $|\mathcal{I}^1| + 1$ , and the one farthest to the right is labeled  $|\mathcal{I}^1| + |\mathcal{I}^2|$ , and so on ending with layer  $l^*$ . If  $\mathcal{I}^1 = \emptyset$ , we use the same scheme except that buffer 1 will then be the buffer farthest to the left in layer 2.

Under the arrangement and renumbering conventions, for any server  $k$ , all of the buffers served by  $k$  are numbered such that lower priority buffers have higher numbers. (This is true whether the buffers are above or below the server and whether they are transition or non-transition buffers.) In particular, for  $k \neq k^*$ , the higher priority transition buffers in  $\underline{\mathcal{I}}_k$  are numbered lower than the non-transition buffers in  $\underline{\mathcal{I}}_k$  and the transition buffer above  $k$  has a higher number than all of the buffers in  $\underline{\mathcal{I}}_k$ . If  $k = k^*$ , the same statement holds with the exception that  $i^*$  is the highest numbered buffer and has the lowest priority among all buffers in  $\underline{\mathcal{I}}_{k^*}$ . It also follows from our numbering scheme that  $i^* = \mathbf{I}$ .

Figure 6 illustrates the server-buffer graph of Figure 2, first with the arrangement convention imposed and then with the renumbering convention applied, where the correspondence between the old and new numbering of classes (or buffers) is given by  $1 \rightarrow 4'$ ,  $2 \rightarrow 5'$ ,  $3 \rightarrow 3'$ ,  $4 \rightarrow 1'$ ,  $5 \rightarrow 2'$ .

## 6.2 Threshold Sizes and Transient Nominal Activity Rates

For  $k \in \mathcal{K}$  such that  $\underline{\mathcal{I}}_k \neq \emptyset$  and  $i \in \underline{\mathcal{I}}_k$ , let

$$x_i^+ = \sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' \leq i}} x_{a(i')}^*, \quad (6.1)$$

and for each  $i' \in \underline{\mathcal{I}}_k$ ,  $i' \leq i$ , let

$$\hat{x}_{i,i'} = \frac{x_{a(i')}^*}{x_i^+}. \quad (6.2)$$

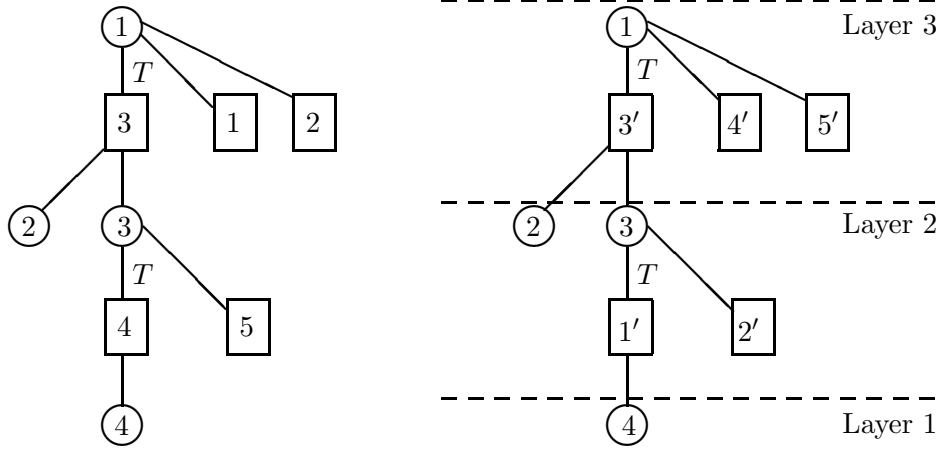


Figure 6: An illustration of the application of the arrangement (on the left) and renumbering (on the right) conventions to the server-buffer tree pictured in Figure 2.

(Note that this defines  $x_i^+$  and  $\hat{x}_{i,i}$  for all  $i \in \mathcal{I}$  since each buffer is below exactly one server.)

We refer to the  $\hat{x}_{i,i'}$  as *transient nominal activity rates* for the following reason. Since  $x_{a(i')}$  determines the overall average fraction of time that server  $k$  should devote to activity  $a(i')$ , and since activities  $a(i')$ ,  $i' > i$ ,  $i' \in \underline{\mathcal{I}}_k$ , will be turned off during a period in which buffer  $i$  either exceeds its threshold if it is a transition buffer or is non-empty if it is a non-transition buffer,  $\hat{x}_{i,i'}$ ,  $i' \leq i$ ,  $i' \in \underline{\mathcal{I}}_k$ , might be interpreted as the average fraction of time that server  $k$  should devote to activity  $a(i')$  during such a period of time.

We note that for  $i \in \underline{\mathcal{I}}_k$ ,

$$\sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' \leq i}} \hat{x}_{i,i'} = 1, \quad (6.3)$$

and if  $i \neq i^*$ ,  $i' \in \underline{\mathcal{I}}_k$ ,  $i' \leq i$ , then  $\hat{x}_{i,i'} > x_{a(i')}^*$  since then  $x_i^+ < 1$  (if  $k \in \mathcal{K}^l$  and  $l \neq l^*$ , server  $k$  serves a buffer from layer  $l+1$ , by an activity  $b(k)$  that is above  $k$ , and so by the heavy traffic condition,  $x_{b(k)}^* + \sum_{i' \in \underline{\mathcal{I}}_k} x_{a(i')}^* = 1$ , and if  $k \in \mathcal{K}^{l^*}$ , then  $k = k^*$  and  $\sum_{i' \in \underline{\mathcal{I}}_{k^*}} x_{a(i')}^* = 1$ , where  $i' \leq i^*$  for all  $i'$  in  $\underline{\mathcal{I}}_{k^*}$ ).

We now define the size of the thresholds to be used with our threshold policy. For each  $r \geq 1$ , let  $L_0^r = \lceil c \log r \rceil$  for a sufficiently large constant  $c$ . The minimum size of  $c$  is determined by the proofs of Lemmas 7.3–7.6 and Theorem 5.4 (see the remark below). For  $1 \leq i \leq \mathbf{I}$ , let

$$L_i^r = \left\lceil \frac{L_{i-1}^r}{\epsilon_{i-1}^3} \right\rceil, \quad (6.4)$$

where  $\{\epsilon_i\}_{i=0}^{\mathbf{I}-1}$  is defined as follows. We also define a constant  $\epsilon_{\mathbf{I}}$  in this process. First we choose  $\hat{\epsilon} > 0$  such that

$$\hat{\epsilon} < \min \left\{ \frac{d_{\min} \delta_{\min} \mu_{\min} \lambda_{\min} x_{\min}^*}{2048(\mathbf{J}+2) \mu_{\max} \lambda_{\max} \delta_{\max} \mu_{\text{sum}}}, \frac{\prod_{i=1}^{\mathbf{I}} \gamma_i}{\mathbf{I}} \right\}, \quad (6.5)$$

where  $\lambda_{\min} = \min\{1, \lambda_i : i \in \mathcal{I}\}$ ,  $\mu_{\min} = \min\{1, \mu_j : j \in \mathcal{J}\}$ ,  $\lambda_{\max} = \max\{1, \lambda_i : i \in \mathcal{I}\}$ ,  $\mu_{\max} = \max\{1, \mu_j : j \in \mathcal{J}\}$ ,  $d_{\min} = \min\{1, \lambda_i - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j : i \in \mathcal{I}\}$ ,  $\delta_{\min} = \min\{1, \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j +$



$\hat{x}_{i,i}\mu_{a(i)} - \lambda_i : i \in \mathcal{I} \setminus \{i^*\}$ ,  $\delta_{\max} = \max\{1, \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j + \hat{x}_{i,i}\mu_{a(i)} - \lambda_i : i \in \mathcal{I}\}$ ,  $x_{\min}^* = \min\{x_j^* : j \in \mathcal{J}\}$ ,  $\mu_{\text{sum}} = \max\{1, \sum_{j \in \mathcal{J}} \mu_j\}$ , and

$$\gamma_i = \frac{\mu_{a(i)}}{64 \sum_{j \in \underline{\mathcal{J}}_i} \mu_j}, \quad 1 \leq i \leq \mathbf{I}. \quad (6.6)$$

Then we define by backward induction on  $i$ ,

$$\epsilon_{\mathbf{I}} = \frac{\hat{c}}{\mathbf{I}}, \quad \epsilon_i = \frac{\gamma_{i+1}\epsilon_{i+1}}{\mathbf{I}}, \quad 0 \leq i < \mathbf{I}. \quad (6.7)$$

(Note that  $\gamma_i < 1$  for all  $i \in \mathcal{I}$ , and  $\epsilon_i < \epsilon_{i+1} \leq \hat{c} < 1$ ,  $i = 0, 1, 2, \dots, \mathbf{I} - 1$ .) Finally, in the  $r^{\text{th}}$  system, if  $i$  is a transition buffer, we let  $L_i^r$  be the threshold size for buffer  $i$ . If  $i$  is a non-transition buffer,  $L_i^r$  is not used to define a threshold, it is simply defined to facilitate the iterative definition of  $L_i^r$ ,  $i = 1, \dots, \mathbf{I}$ .

**Remark.** For our method of proof to work, the constant  $c$  must be sufficiently large. In the proofs of Lemmas 7.3–7.6 and of uniform integrability in the proof of Theorem 5.4, a means for determining a value  $c^*$  is described such that our method works provided  $c > c^*$ . This value is determined from several applications of large deviation estimates for the renewal processes associated with the interarrival and service time sequences (cf. Assumption 3.2). As in [3], we have not attempted to give a concise formula for  $c^*$  nor to optimize its value, since the relevant fact is that sufficiently large thresholds of order  $\log r$  work and this order is the smallest for which our proof works. (For analysis and approximate analysis of the effects of different threshold sizes for some dynamic scheduling problems, see for example, [37, 38] and [33, 35].)

### 6.3 State Space Collapse Result and Outline of Proof

A key element in the proof of Theorem 5.3 is to first show the following “state space collapse” result.

**Theorem 6.1** *Consider the sequence of parallel server systems indexed by  $r$ , where the  $r^{\text{th}}$  system operates under the scheduling control,  $T^{r,*}$ , described in Sections 5 and 6.2. Then*

$$\left( \hat{Q}_i^r, \hat{I}_k^r : i \in \mathcal{I} \setminus \{i^*\}, k \in \mathcal{K} \setminus \{k^*\} \right) \Rightarrow \mathbf{0} \quad \text{as } r \rightarrow \infty, \quad (6.8)$$

where  $\mathbf{0}$  is the function in  $D^{\mathbf{I}+\mathbf{K}-2}$  that remains at the origin of  $\mathbb{R}^{\mathbf{I}+\mathbf{K}-2}$  for all time.

The idea behind the proof of this theorem is that, for sufficiently large  $r$ , under the threshold control  $T^{r,*}$ , for a transition class  $i \in \mathcal{I} \setminus \{i^*\}$ , once the queue length process  $Q_i^r$  has first reached its threshold level  $L_i^r$  (cf. (6.4)), over time intervals of order  $r^2$  in length,  $Q_i^r$  rarely deviates as much as  $L_i^r - |\mathcal{J}_i|$  from the threshold level, since when it is above this level, it is driven down towards the level at an “average” rate of  $\hat{x}_{i,i}\mu_{a(i)}^r + \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j^r - \lambda_i^r > 0$ , and when it is below the level, it is driven up towards the level at an average rate of  $\lambda_i^r - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j^r > 0$ . Similarly, if  $i^*$  is a transition class, once  $Q_{i^*}^r$  has reached the level  $L_{i^*}^r$ ,  $\hat{Q}_{i^*}^r$  rarely reaches as low as the low level  $|\mathcal{J}_{i^*}|$ . However, since  $i^*$  is the cheapest buffer,  $Q_{i^*}^r$  may reach considerably above  $L_{i^*}^r$ . For a non-transition class  $i \neq i^*$ ,  $\underline{\mathcal{J}}_i = \emptyset$  and  $\hat{Q}_i^r$  is driven down towards zero at an average rate of  $\hat{x}_{i,i}\mu_{a(i)}^r - \lambda_i^r > 0$ . (The claimed positivity of the quantities above holds for large  $r$  since  $\mathbf{R}x^* = \lambda$  (cf. Assumptions 3.1 and 3.6) and  $1 \geq \hat{x}_{i,i} > x_{a(i)}^*$  for  $i \neq i^*$ .) The behavior of the queue length process  $Q_i^r$  for a transition buffer  $i \in \mathcal{I}$  has a consequential effect on the idleness processes  $I_k^r$ , for  $k \in \underline{\mathcal{K}}_i$ , since idletime for server  $k$  cannot increase when the queue length process  $Q_i^r$  is above the level  $|\mathcal{J}_i|$ .

Estimates associated with the above ideas, along with estimates on the cumulative idletime for server  $k \in \underline{\mathcal{K}}_i$  until  $Q_i^r$  reaches  $L_i^r$  (for a transition class  $i$ ), are stated formally in Theorem 7.1.

The proof of this theorem uses large deviation estimates for the primitive renewal processes and estimates for the cumulative allocation processes (cf. Lemma 7.3). The proof employs an induction on the buffers in the server-buffer tree, starting from  $i = 1$  and iterating to buffer  $i = \mathbf{I} - 1$ . The buffer  $i^*$  (when it is a transition buffer) is treated separately, although this treatment uses the consequences of the induction proof. The induction setup is described in Section 7.2 and the proofs are in Sections 7.3–7.7. Theorem 6.1 follows from Theorem 7.1 using the fact that the threshold  $L_i^r$  being of order  $\log r$  implies that  $(L_i^r - |\underline{\mathcal{J}}_i|)/r$  goes to zero as  $r$  goes to infinity.

Once Theorem 6.1 is established, one can show, using the model equations for queue length and idletime (cf. (2.15)–(2.16)), that the fluid scaled allocations  $\bar{T}^{r,*}(\cdot) \equiv r^{-2}T^{r,*}(r^2 \cdot)$  associated with  $T^{r,*}$  satisfy

$$\bar{T}^{r,*} \Rightarrow \bar{T}^*, \quad \text{as } r \rightarrow \infty, \quad (6.9)$$

where

$$\bar{T}^*(t) \equiv x^*t, \quad t \geq 0. \quad (6.10)$$

One can then combine the above to prove Theorem 5.3. These results are proved in Section 8.

For the proof of Theorem 5.4, we first show (cf. Lemma 9.1) that for any subsequence that achieves the “liminf” on the left side of (5.2) as a limit and for which the “liminf” is finite, the fluid level asymptotic behavior described in (6.9) must hold along the subsequence, with  $\{T^r\}$  in place of  $\{T^{r,*}\}$  there. This, together with a pathwise lower bound for  $h^r \cdot \hat{Q}^r$ , where  $h^r$  is a perturbation of  $h$  given in (9.12), allows us to establish the inequality in (5.2). The equality in (5.2) follows from Theorem 5.3 after showing that a certain uniform integrability condition holds.

## 6.4 Residual Processes and Shifted Allocation Processes

Key to our proof of Theorem 6.1 is the behavior of what we call the *residual processes* defined for  $i \in \mathcal{I}$ ,  $r \geq 1$ ,  $s \geq 0$ , by

$$R_i^r(s) = \begin{cases} Q_i^r(s) - L_i^r, & \text{if } i \text{ is a transition class,} \\ Q_i^r(s), & \text{otherwise.} \end{cases} \quad (6.11)$$

For the case when  $i \neq i^*$  is a transition buffer, the idea of our proof is to move the center of one’s attention to the threshold and to show that  $Q_i^r$  reaches the threshold level  $L_i^r$  relatively quickly and then “chatters” back and forth across this threshold, not frequently deviating “far” from it, so that  $Q_i^r$  rarely again goes as low as the level  $|\underline{\mathcal{J}}_i|$ , or as high as the level  $2L_i^r - |\underline{\mathcal{J}}_i|$ . When translated into the behavior of  $R_i^r$ , this means that we show that  $R_i^r$  reaches the level zero relatively quickly and then it chatters back and forth across the zero level, rarely deviating by as much as  $\pm(L_i^r - |\underline{\mathcal{J}}_i|)$  from this level (cf. Theorem 7.1). If  $i \neq i^*$  is a non-transition buffer, we show that  $Q_i^r$  (or equivalently  $R_i^r$ ) rarely goes above the level  $L_i^r$ .

For describing the excursions of  $R_i^r$  above zero, we introduce the following notation.

**Definition 6.2** For  $i \in \mathcal{I}$ , let  $\tau_{i,0}^r = \inf\{s \geq 0 : R_i^r(s) \geq 0\}$ ,  $\tau_{i,1}^r = \inf\{s \geq \tau_{i,0}^r : R_i^r(s) \geq 1\}$ ,  $\tau_{i,2}^r = \inf\{s \geq \tau_{i,1}^r : R_i^r(s) \leq 0\}$  and define recursively  $\tau_{i,2n-1}^r = \inf\{s \geq \tau_{i,2n-2}^r : R_i^r(s) \geq 1\}$ ,  $\tau_{i,2n}^r = \inf\{s \geq \tau_{i,2n-1}^r : R_i^r(s) \leq 0\}$ , for  $n = 2, 3, \dots$ . For each  $n \geq 1$ , we say that  $[\tau_{i,2n-1}^r, \tau_{i,2n}^r]$  is the  $n^{\text{th}}$  “up” excursion interval for  $R_i^r(\cdot)$ , and on  $\{\tau_{i,2n-1}^r < \infty\}$  we let  $\beta_{i,n}^r = \tau_{i,2n}^r - \tau_{i,2n-1}^r$ , and on  $\{\tau_{i,2n-1}^r = \infty\}$  we let  $\beta_{i,n}^r = 0$ .

For describing the “down” excursion intervals of  $R_i^r$  when  $i \in \mathcal{I}$  is a transition class, we define

$${}^dR_i^r \equiv -R_i^r = L_i^r - Q_i^r. \quad (6.12)$$

If  $i$  is not a transition class, we set  ${}^dR_i \equiv \mathbf{0}$ , and hence the following definition is only non-trivial for a transition class  $i$ .

**Definition 6.3** For  $i \in \mathcal{I}$ , let  $d_{\tau_{i,1}^r}^r = \inf\{s \geq \tau_{i,0}^r : dR_i^r(s) \geq 1\}$ ,  $d_{\tau_{i,2}^r}^r = \inf\{s \geq d_{\tau_{i,1}^r}^r : dR_i^r(s) \leq 0\}$  and define recursively  $d_{\tau_{i,2n-1}^r}^r = \inf\{s \geq d_{\tau_{i,2n-2}^r}^r : dR_i^r(s) \geq 1\}$ ,  $d_{\tau_{i,2n}^r}^r = \inf\{s \geq d_{\tau_{i,2n-1}^r}^r : dR_i^r(s) \leq 0\}$ , for  $n = 2, 3, \dots$ . For each  $n \geq 1$ , we say that  $[d_{\tau_{i,2n-1}^r}^r, d_{\tau_{i,2n}^r}^r]$  is the  $n^{\text{th}}$  down excursion interval for  $R_i^r(\cdot)$ , and on  $\{d_{\tau_{i,2n-1}^r}^r < \infty\}$  we let  $d\beta_{i,n}^r = d_{\tau_{i,2n}^r}^r - d_{\tau_{i,2n-1}^r}^r$ , and on  $\{d_{\tau_{i,2n-1}^r}^r = \infty\}$  we define  $d\beta_{i,n}^r = 0$ .

Estimates of the amount of time that activities in  $\mathcal{J}_i$  are on during the  $n^{\text{th}}$  (up/down) excursion interval for  $R_i^r$  are needed to obtain estimates of the value of  $R_i^r$  during such an interval (cf. (7.40)). Such estimates for the activities in  $\underline{\mathcal{J}}_i$  depend in turn on estimates for the on-time of activities farther down the tree, whereas estimates for the on-time of activity  $a(i)$  depend on estimates for the on-time of the (higher priority) activities that are served from above by the same server that serves buffer  $i$  from above. To keep track of all relevant on-times, for each  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$ , we define shifted allocation processes for  $n \geq 1$  and  $s \geq 0$ , by

$$T_{i,j}^{r,n}(s) = T_j^r(\tau_{i,2n-1}^r + s) - T_j^r(\tau_{i,2n-1}^r), \quad \text{on } \{\tau_{i,2n-1}^r < \infty\}, \quad (6.13)$$

$$d_{T_{i,j}^{r,n}}(s) = T_j^r(d_{\tau_{i,2n-1}^r}^r + s) - T_j^r(d_{\tau_{i,2n-1}^r}^r), \quad \text{on } \{d_{\tau_{i,2n-1}^r}^r < \infty\}, \quad (6.14)$$

and on  $\{\tau_{i,2n-1}^r = \infty\}$  let  $T_{i,j}^{r,n} \equiv \mathbf{0}$ , and on  $\{d_{\tau_{i,2n-1}^r}^r = \infty\}$  let  $d_{T_{i,j}^{r,n}} \equiv \mathbf{0}$ . We have that on  $\{\tau_{i,2n-1}^r < \infty\}$ ,  $T_{i,j}^{r,n}$  measures the on-time of activity  $j \in \mathcal{J}$  following an up-crossing to or above the level one by  $R_i^r$ , and, for a transition class  $i$ , on  $\{d_{\tau_{i,2n-1}^r}^r < \infty\}$ ,  $d_{T_{i,j}^{r,n}}$  measures the on-time of activity  $j$  following a down-crossing to or below level minus one by  $R_i^r$ . Note that  $d_{T_{i,j}^{r,n}} \equiv \mathbf{0}$  if  $i$  is not a transition class.

## 6.5 Preliminaries on Stopped Arrival and Service Processes

For the proof of (7.1) in Theorem 7.1, we need to establish some preliminary results concerning the properties of the arrival and service processes stopped at certain hitting times, so that we can apply the results of Appendix A in [3] to shifted versions of these processes.

Let  $\mathbb{N}_\infty = \mathbb{N} \cup \{\infty\}$ . Consider  $\mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}$  to be partially ordered by componentwise inequality, i.e.,  $(n, m) \leq (p, q)$  if and only if  $n_i \leq p_i$ , and  $m_j \leq q_j$  for all  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$ . Recall from Section 2.2 the definition, for the  $r^{\text{th}}$  system, of the cumulative interarrival time process for class  $i \in \mathcal{I}$ ,  $\xi_i^r$ , and the cumulative service time process for activity  $j \in \mathcal{J}$ ,  $\eta_j^r$ . For each  $(p, q) \in \mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}$  let

$$\mathcal{F}_{pq}^r = \sigma\{\xi_i^r(\cdot \wedge (p_i + 1)), \eta_j^r(\cdot \wedge (q_j + 1)) : i \in \mathcal{I}, j \in \mathcal{J}\} \vee \mathcal{P}_0,$$

where  $\mathcal{P}_0$  denotes the collection of  $\mathbf{P}$ -null sets in the complete probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ . Then  $\{\mathcal{F}_{pq}^r : (p, q) \in \mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}\}$  is a multiparameter filtration (cf. Ethier and Kurtz [10], p. 85).

**Definition 6.4** A (multiparameter) stopping time relative to  $\{\mathcal{F}_{pq}^r : (p, q) \in \mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}\}$  is a random variable  $\mathcal{T}$  taking values in  $\mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}$  such that

$$\{\mathcal{T} = (p, q)\} \in \mathcal{F}_{pq}^r \text{ for all } (p, q) \in \mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}. \quad (6.15)$$

The  $\sigma$ -algebra associated with such a stopping time  $\mathcal{T}$  is

$$\mathcal{F}_{\mathcal{T}}^r = \{B \in \mathcal{F} : B \cap \{\mathcal{T} = (p, q)\} \in \mathcal{F}_{pq}^r \text{ for all } (p, q) \in \mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}\}. \quad (6.16)$$

**Lemma 6.5** Suppose  $r \geq 1$  is such that  $L_0^r \geq \mathbf{J} + 1$ . Then, for each  $\iota \in \mathcal{I}$ ,  $n \geq 1$ ,

$$\begin{aligned} T_{n,\iota}^r &\equiv (A_i^r(\tau_{i,2n-1}^r), S_j^r(T_j^r(\tau_{i,2n-1}^r)) : i \in \mathcal{I}, j \in \mathcal{J}), \\ d_{T_{n,\iota}^r} &\equiv (A_i^r(d_{\tau_{i,2n-1}^r}^r), S_j^r(T_j^r(d_{\tau_{i,2n-1}^r}^r)) : i \in \mathcal{I}, j \in \mathcal{J}) \end{aligned}$$

are (multiparameter) stopping times relative to the filtration  $\{\mathcal{F}_{pq}^r : (p, q) \in \mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}\}$ , where we adopt the convention that each of  $A_i^r(\cdot)$ ,  $S_j^r(T_j^r(\cdot))$ ,  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$ , takes the value  $\infty$  when evaluated at the time  $\infty$ .

**Remark.** Here we need  $r$  large enough to simplify the proof of Lemma 6.5. Since, in the sequel, we let  $r$  approach infinity, this result will suffice for our purposes.

**Proof.** This lemma can be proved in a similar manner to Lemma 8.3 in [41]. A proof is given in the appendix to [2].  $\square$

**Lemma 6.6** *Let  $\mathcal{T}$  be a (multiparameter) stopping time relative to the filtration  $\{\mathcal{F}_{pq}^r : (p, q) \in \mathbb{N}_\infty^{\mathbf{I}} \times \mathbb{N}_\infty^{\mathbf{J}}\}$ . Let  $\mathcal{T}^{\mathcal{I}}$  denote the first  $\mathbf{I}$  components of  $\mathcal{T}$  and  $\mathcal{T}^{\mathcal{J}}$  denote the other  $\mathbf{J}$  ( $= \mathbf{B}$ ) components of  $\mathcal{T}$  so that  $\mathcal{T} = (\mathcal{T}^{\mathcal{I}}, \mathcal{T}^{\mathcal{J}})$ . In the following, for notational convenience, we make the convention that each of  $u_i^r(\cdot)$ ,  $v_j^r(\cdot)$ ,  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$ , takes the value  $\infty$  when its argument takes the value  $\infty$ . Then,*

$$(u_i^r(\mathcal{T}_i^{\mathcal{I}} + 1), v_j^r(\mathcal{T}_j^{\mathcal{J}} + 1) : i \in \mathcal{I}, j \in \mathcal{J}) \in \mathcal{F}_{\mathcal{T}}^r, \quad (6.17)$$

and, on  $\{\mathcal{T} \in \mathbb{N}^{\mathbf{I}} \times \mathbb{N}^{\mathbf{J}}\}$ , the conditional distribution of  $\{(u_i^r(\mathcal{T}_i^{\mathcal{I}} + n), v_j^r(\mathcal{T}_j^{\mathcal{J}} + n) : i \in \mathcal{I}, j \in \mathcal{J}), n = 2, 3, \dots\}$  given  $\mathcal{F}_{\mathcal{T}}^r$  is the same as the (unconditioned) distribution of the original family of i.i.d. random variables  $\{(u_i^r(n), v_j^r(n) : i \in \mathcal{I}, j \in \mathcal{J}), n = 1, 2, \dots\}$ .

**Proof.** For a proof with  $|\mathcal{I}| = 1$ ,  $|\mathcal{J}| = 2$ , see the proof of Lemma 7.6 in [3]. The general proof is similar.  $\square$

## 6.6 Large Deviation Bounds for Renewal Processes

The following lemma, which will be used extensively in proving state space collapse in Section 7, summarizes the results of the discussion in Appendix A in [3].

**Lemma 6.7** *Let  $\{\zeta(i)\}_{i=1}^\infty$  be a sequence of strictly positive independent random variables, where  $\{\zeta(i)\}_{i=2}^\infty$  are identically distributed with finite mean  $1/\nu$ , for some  $\nu \in (0, \infty)$ , and  $\zeta(1)$  may have a different distribution from  $\zeta(i)$  for  $i > 1$ . Assume that there is a nonempty open neighborhood  $\mathcal{O}$  of  $0 \in \mathbb{R}$  such that for  $i = 2, 3, \dots$ ,*

$$\Lambda(l) \equiv \log \mathbf{E}(e^{l\zeta(i)}) < \infty \quad \text{for all } l \in \mathcal{O}. \quad (6.18)$$

Let the values of  $r \geq 1$  range through a sequence that increases to infinity. For each  $r$ , let  $\nu^r > 0$  and suppose that  $\lim_{r \rightarrow \infty} \nu^r = \nu$ . For each  $r$  and  $i = 2, 3, \dots$ , let

$$\zeta^r(i) = \frac{\nu}{\nu^r} \zeta(i). \quad (6.19)$$

Given  $0 < \epsilon < \nu/2$ , let  $r_\epsilon \geq 1$  be such that for  $r \geq r_\epsilon$ ,

$$|\nu^r - \nu| < \epsilon, \quad (6.20)$$

$$\frac{\nu^r}{\nu} \left( \frac{1}{\nu^r + \frac{\epsilon}{2}} \right) \leq \frac{1}{\nu} \left( \frac{1}{1 + \frac{\epsilon}{3\nu}} \right) < \frac{1}{\nu}, \quad (6.21)$$

$$\frac{1}{\nu} \left( 1 + \frac{\epsilon}{2(\nu^r - \epsilon)} \right) \geq \frac{1}{\nu} \left( 1 + \frac{\epsilon}{2\nu} \right) > \frac{1}{\nu}. \quad (6.22)$$

For each  $r \geq 1$ ,  $s \geq 0$ , let

$$N^r(s) = \sup \left\{ n \geq 0 : \sum_{i=1}^n \zeta^r(i) \leq s \right\}. \quad (6.23)$$

Then for  $r \geq r_\epsilon$ ,  $s > 2/\epsilon$ ,

$$\begin{aligned} \mathbf{P}(N^r(s) > (\nu^r + \epsilon)s) &\leq \exp\left(-((\nu^r + \epsilon)s - 1)\Lambda^*\left(\frac{1}{\nu}\left(\frac{1}{1 + \frac{\epsilon}{3\nu}}\right)\right)\right) \\ &\leq \exp\left(-(\nu s - 1)\Lambda^*\left(\frac{1}{\nu}\left(\frac{1}{1 + \frac{\epsilon}{3\nu}}\right)\right)\right), \end{aligned} \quad (6.24)$$

and for  $r \geq r_\epsilon$ ,  $s \geq 0$ ,

$$\begin{aligned} \mathbf{P}(N^r(s) < (\nu^r - \epsilon)s) &\leq \exp\left(-(\nu^r - \epsilon)s\Lambda^*\left(\frac{1}{\nu}\left(1 + \frac{\epsilon}{2\nu}\right)\right)\right) \\ &\quad + \mathbf{P}\left(\zeta^r(1) > \frac{\epsilon}{2\nu^r}s\right) \\ &\leq \exp\left(-(\nu - 2\epsilon)s\Lambda^*\left(\frac{1}{\nu}\left(1 + \frac{\epsilon}{2\nu}\right)\right)\right) \\ &\quad + \mathbf{P}\left(\zeta^r(1) > \frac{\epsilon}{2\nu^r}s\right), \end{aligned} \quad (6.25)$$

where

$$\Lambda^*(x) \equiv \sup_{l \in \mathbb{R}}(lx - \Lambda(l)), \quad (6.26)$$

and where the values of the quantities involving  $\Lambda^*$  in the above are strictly positive. (The function  $\Lambda^*$  is called the Legendre-Fenchel transform of  $\Lambda$ .) Furthermore, if  $\zeta(1)$  has the same distribution as  $\{\zeta(i)\}_{i=2}^\infty$ , then for each  $r \geq 1$ ,  $s \geq 0$  and  $0 < l_0 \in \mathcal{O}$ , for any  $n \geq 1$ ,

$$\mathbf{P}\left(\max_{i=1}^n \zeta^r(i) > \frac{\epsilon}{2\nu^r}s\right) \leq n \exp\left(-\frac{l_0 \epsilon s}{2\nu}\right) \exp(\Lambda(l_0)). \quad (6.27)$$

## 7 Proof of State Space Collapse

Throughout this section, it is assumed that in the  $r^{\text{th}}$  parallel server system we use the allocation process  $T^{r,*}$  associated with the threshold policy described in Sections 5 and 6.2. To simplify notation, here we shall simply write  $T^r$  in place of  $T^{r,*}$ , since no other policy is considered in this section. The associated queue length and idletime processes will be denoted by  $Q^r$ ,  $I^r$ , respectively.

Recall the definition of the residual processes (cf. (6.11)), and the role of  $c$  in the definition of  $L_0^r$  (cf. Section 6.2). For the following theorem, which is the main technical result of this section and from which Theorem 6.1 will follow, there is  $c^* > 0$  such that the results hold provided the fixed constant  $c$  is greater than  $c^*$  (cf. the proof of Theorem 7.1).

**Theorem 7.1** For each  $i \in \mathcal{I} \setminus \{i^*\}$ ,  $k \in \underline{\mathcal{K}}_i$ ,  $t \geq 0$ , and  $\epsilon > 0$ ,

$$\mathbf{P}\left(\sup_{\tau_{i,0}^r \leq s \leq r^2 t} |R_i^r(s)| \geq L_i^r - |\underline{\mathcal{J}}_i|\right) \rightarrow 0 \quad \text{as } r \rightarrow \infty, \quad (7.1)$$

$$\mathbf{P}(I_k^r(\tau_{i,0}^r) \geq r\epsilon) \rightarrow 0 \quad \text{as } r \rightarrow \infty, \quad (7.2)$$

$$\mathbf{P}(I_k^r(r^2 t) - I_k^r(\tau_{i,0}^r) > 0, \tau_{i,0}^r < r^2 t) \rightarrow 0 \quad \text{as } r \rightarrow \infty, \quad (7.3)$$

and, in addition, if  $i^*$  is a transition class, then for each  $k \in \underline{\mathcal{K}}_{i^*}$ ,  $t \geq 0$ , and  $\epsilon > 0$ ,

$$\mathbf{P}\left(\inf_{\tau_{i^*,0}^r \leq s \leq r^2 t} Q_{i^*}^r(s) \leq |\underline{\mathcal{J}}_{i^*}|\right) \rightarrow 0 \quad \text{as } r \rightarrow \infty, \quad (7.4)$$

$$\mathbf{P}(I_k^r(\tau_{i^*,0}^r) \geq r\epsilon) \rightarrow 0 \quad \text{as } r \rightarrow \infty, \quad (7.5)$$

$$\mathbf{P}(I_k^r(r^2 t) - I_k^r(\tau_{i^*,0}^r) > 0, \tau_{i^*,0}^r < r^2 t) \rightarrow 0 \quad \text{as } r \rightarrow \infty. \quad (7.6)$$

Here,  $\tau_{i,0}^r = \inf\{s \geq 0 : Q_i^r(s) \geq L_i^r\}$  if class  $i$  is a transition class,  $\tau_{i,0}^r \equiv 0$  if class  $i$  is a non-transition class, and  $|\underline{\mathcal{J}}_i|$  is the number of basic activities that serve class  $i$  from below in the server-buffer tree (cf. Figure 5).

Here we have used the convention in (7.2)–(7.3) and (7.5)–(7.6) that  $I_k^r(\tau_{i,0}^r) = \lim_{t \rightarrow \infty} I_k^r(t)$ , on  $\{\tau_{i,0}^r = \infty\}$ , and in (7.1) (respectively, (7.4)) that the supremum (respectively, infimum) over an empty set is defined to equal  $-\infty$  (respectively,  $\infty$ ).

**Remark.** Although the results in Theorem 7.1 suffice for the proof of Theorem 6.1 and subsequently for Theorem 5.3, a refinement of Theorem 7.1, with estimates of the left members in (7.1)–(7.3), (cf. Theorem 7.7) is needed to establish certain uniform integrability used in the proof of Theorem 5.4.

## 7.1 Auxiliary Constants for the Induction Proof

In this subsection, we introduce various constants and establish various inequalities that hold for  $r$  sufficiently large. These are used in the proofs of Theorems 7.1 and 7.7 (see below). The reader may wish to simply skim this subsection at first, and only refer back to it as needed for reading proofs.

The following constants grow logarithmically with  $r$ . For  $i \in \mathcal{I}$ ,  $r \geq 1$ ,  $L_i^r$ , and  $\epsilon_i$  as defined in (6.4) and (6.7), respectively, let

$$s_i^r = \frac{L_i^r - (|\underline{\mathcal{J}}_i| + 2)}{(\lambda_i^r + \epsilon_i)}, \quad (7.7)$$

$$t_i^r = \frac{8L_i^r}{\lambda_i - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j}. \quad (7.8)$$

Note that if  $i$  is a non-transition class then the denominator in (7.8) is equal to  $\lambda_i > 0$  (by our convention that a sum over an empty set is zero), and if  $i$  is a transition class then the denominator is also positive since  $\lambda_i = x_{a(i)}^* \mu_{a(i)} + \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j$  (cf. (3.5)). Let

$$s_0^r = t_0^r = {}^d s_0^r = L_0^r, \quad (7.9)$$

and, for each  $i \in \mathcal{I}$ ,  $r \geq 1$ , let

$${}^d s_i^r = \frac{L_i^r - (|\underline{\mathcal{J}}_i| + 2)}{\sum_{j \in \underline{\mathcal{J}}_i} (\mu_j^r + \epsilon_i)(x_j^* + \epsilon_i)}, \quad \text{if } i \text{ is a transition class,} \quad (7.10)$$

otherwise let  ${}^d s_i^r = L_i^r$ . Finally, for each  $r \geq 1$ , let

$$M^r = \max \left\{ s_i^r, t_i^r, {}^d s_i^r : i \in \mathcal{I} \right\}. \quad (7.11)$$

Several lemmas are used in establishing Theorems 7.1 and 7.7 (including Lemmas 7.3–7.6 of Section 7.2). For the proofs of these lemmas we require that  $r \geq 1$  is large enough so that various relations involving the auxiliary constants defined above and the parameters for the  $r^{\text{th}}$  parallel server system hold. It is important to note that, for this, the size of  $r$  should not depend on the variable  $t$  referred to in those theorems and lemmas. We first require that  $r$  is large enough that

$$s_i^r \geq s_{i-1}^r, \quad t_i^r \geq t_{i-1}^r, \quad {}^d s_i^r \geq {}^d s_{i-1}^r, \quad \text{for } i = 1, \dots, \mathbf{I}. \quad (7.12)$$

For each  $i \in \mathcal{I}$ , let

$$\delta_i = \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j + \hat{x}_{i,i} \mu_{a(i)} - \lambda_i. \quad (7.13)$$

Then,  $\delta_i > 0$  for  $i \in \mathcal{I} \setminus \{i^*\}$  since  $\lambda_i = \sum_{j \in \mathcal{J}_i} x_j^* \mu_j + x_{a(i)}^* \mu_{a(i)} < \sum_{j \in \mathcal{J}_i} x_j^* \mu_j + \hat{x}_{i,i} \mu_{a(i)}$ , using the fact that  $x_{a(i)}^* < \hat{x}_{i,i}$  for  $i \neq i^*$  (cf. (6.2)). For part of the proof of Lemma 7.3 (associated with (I.1) and up excursions) we require that for all  $i \in \mathcal{I} \setminus \{i^*\}$ ,

$$|\lambda_i^r - \lambda_i| < \min \left\{ \frac{\epsilon_i}{4}, \frac{(\mu_{a(i)} - \epsilon_i) \epsilon_i}{32|\mathcal{J}_i|} \right\}, \quad (7.14)$$

$$|\mu_j^r - \mu_j| < \min \left\{ \frac{\epsilon_i}{4}, \frac{(\mu_{a(i)} - \epsilon_i) \epsilon_i}{32|\mathcal{J}_i|} \right\}, \quad \text{for all } j \in \mathcal{J}_i, \quad (7.15)$$

$$\mu_j^r - \epsilon_i > \frac{4\lambda_i + \delta_i}{4\lambda_i + 2\delta_i} \mu_j \quad \text{for all } j \in \mathcal{J}_i, \quad (7.16)$$

$$\mu_{a(i)}^r - \epsilon_i > \frac{4\lambda_i + \frac{3}{2}\delta_i}{4\lambda_i + 2\delta_i} \mu_{a(i)}, \quad (7.17)$$

$$\lambda_i^r + \epsilon_i < \left( \frac{2\lambda_i}{2\lambda_i + \delta_i} \right) \left( \sum_{j \in \mathcal{J}_i} x_j^* \mu_j + \hat{x}_{i,i} \mu_{a(i)} \right) = \lambda_i + \frac{\lambda_i \delta_i}{2\lambda_i + \delta_i}, \quad (7.18)$$

$$1 < \frac{\delta_i}{4\lambda_i + 2\delta_i} \hat{x}_{i,i} \mu_{a(i)} s_i^r. \quad (7.19)$$

(Note that (7.16)–(7.18) do not hold for  $i = i^*$ , for large  $r \geq 1$ , since, in this case,  $\delta_{i^*} = 0$ ,  $\epsilon_{i^*} > 0$ , and  $\lambda_{i^*}^r \rightarrow \lambda_{i^*}$ ,  $\mu_j^r \rightarrow \mu_j$ ,  $j \in \mathcal{J}_{i^*}$ , as  $r \rightarrow \infty$ .)

For each  $i \in \mathcal{I}$ , define

$$d\delta_i = \lambda_i - \sum_{j \in \mathcal{J}_i} x_j^* \mu_j, \quad (7.20)$$

where  $d\delta_i > 0$  since  $\lambda_i = \sum_{j \in \mathcal{J}_i} x_j^* \mu_j + x_{a(i)}^* \mu_{a(i)} > \sum_{j \in \mathcal{J}_i} x_j^* \mu_j \geq 0$ . For part of the proof of Lemma 7.3 (associated with (I.1) and down excursions) we require that if  $i$  is a transition class,

$$\mu_j^r + \epsilon_i < \frac{4\lambda_i - \frac{3}{2}d\delta_i}{4\lambda_i - 2d\delta_i} \mu_j, \quad \text{for all } j \in \mathcal{J}_i, \quad (7.21)$$

$$\lambda_i^r - \epsilon_i > \frac{2\lambda_i}{2\lambda_i - d\delta_i} \sum_{j \in \mathcal{J}_i} x_j^* \mu_j, \quad (7.22)$$

$$1 < \left( \frac{\lambda_i d\delta_i}{(4\lambda_i - d\delta_i)(2\lambda_i - d\delta_i)} \sum_{j \in \mathcal{J}_i} x_j^* \mu_j \right)^{d s_i^r}. \quad (7.23)$$

(Note that (7.23) does not hold for a non-transition class  $i$  since the right hand side there is zero for all  $r \geq 1$  by our convention that a sum over an empty index set is zero.)

For each transition class  $i \in \mathcal{I}$ , let

$$\tilde{\epsilon}_i = \frac{\epsilon_i \sum_{j \in \mathcal{J}_i} \mu_j}{1 - 2|\mathcal{J}_i| \epsilon_i}. \quad (7.24)$$

To use Lemma 6.7, we need  $\tilde{\epsilon}_i < \min\{\lambda_i/2, \mu_j/2 : j \in \mathcal{J}_i\}$  for each transition class  $i$ . To see that this condition is satisfied, we note that by (6.5) and (6.7),  $\epsilon_i < 1/4\mathbf{J} < 1/4|\mathcal{J}_i|$  so that  $\tilde{\epsilon}_i < 2\epsilon_i \sum_{j \in \mathcal{J}_i} \mu_j < \min\{\lambda_i/2, \mu_j/2 : j \in \mathcal{J}_i\}$ . For part of the proof of Lemma 7.3 (associated with (I.2)), we require that for each transition class  $i$ ,

$$\lambda_i^r - \sum_{j \in \mathcal{J}_i} x_j^* \mu_j^r \geq \frac{1}{2} \left( \lambda_i - \sum_{j \in \mathcal{J}_i} x_j^* \mu_j \right), \quad (7.25)$$

$$|\lambda_i - \lambda_i^r| < \tilde{\epsilon}_i, \quad |\mu_j - \mu_j^r| < \tilde{\epsilon}_i, \quad \text{for all } j \in \underline{\mathcal{J}}_i. \quad (7.26)$$

For the proof of Lemma 7.6, for all  $\iota \geq i$ ,  $i \in \mathcal{I}$ ,  $\iota \in \mathcal{I}$ , we require that

$$s_\iota^r > \left[ \frac{\epsilon_i}{4} \left( \mu_{a(i)} - \frac{\mu_{a(i)} \epsilon_i}{16} \right) \right]^{-1}, \quad \frac{\mu_{a(i)}^r}{\mu_{a(i)}} > \frac{1}{2}, \quad \frac{\lambda_\iota^r}{\lambda_\iota} < \frac{3}{2}, \quad (7.27)$$

$$\left| \lambda_i^r - x_{a(i)}^* \mu_{a(i)}^r - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j^r \right| < \frac{\mu_{a(i)} \epsilon_i}{32}, \quad (7.28)$$

$$L_i^r > \frac{\epsilon_i (|\underline{\mathcal{J}}_\iota| + 2) \max\{\mu_{a(i)}, \lambda_i - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j\}}{\min\{\lambda_\iota, \sum_{j \in \underline{\mathcal{J}}_\iota} (\mu_j + \epsilon_\iota)\}}, \quad (7.29)$$

and, in addition, if  $\iota$  is a transition class,

$$\frac{\sum_{j \in \underline{\mathcal{J}}_\iota} (\mu_j + \epsilon_\iota) (x_j^* + \epsilon_\iota)}{\sum_{j \in \underline{\mathcal{J}}_\iota} (\mu_j^r + \epsilon_\iota) (x_j^* + \epsilon_\iota)} > \frac{1}{2}, \quad d_{S_\iota}^r \geq \frac{2\mathbf{J}}{\mu_{a(i)}^r \epsilon_i}. \quad (7.30)$$

For each  $i \in \mathcal{I}$ , let

$$\tilde{\epsilon}_i = \frac{(\mu_{a(i)} - \epsilon_i) \epsilon_i}{16\mathbf{J}} < \min \left\{ \frac{\lambda_i}{2}, \frac{\mu_j}{2} : j \in \mathcal{J}_i \right\}. \quad (7.31)$$

The inequality here holds by (6.5). In addition, if  $i$  is a transition class, let

$$\epsilon_{1,i} = \left( \sum_{j \in \underline{\mathcal{J}}_i} (\mu_j + \epsilon_i) \right)^{-1} \frac{\epsilon_i}{16} \left( \mu_{a(i)} - \frac{\mu_{a(i)} \epsilon_i}{16} \right), \quad (7.32)$$

and set  $\epsilon_{1,i} = 1$  if  $i$  is not a transition class. In either case, we also define

$$\epsilon_{2,i} = \frac{\mu_{a(i)} \epsilon_i}{16} (x_{a(i)}^* + \epsilon_i)^{-1} < \frac{\mu_{a(i)}}{2}. \quad (7.33)$$

If  $i \in \mathcal{I}$  is a transition class, for each  $j \in \underline{\mathcal{J}}_i$  let

$$\epsilon_{3,j} = \frac{\mu_{a(i)} \epsilon_i}{16|\underline{\mathcal{J}}_i|} (x_j^* + \epsilon_{1,i})^{-1} < \frac{\mu_j}{2}. \quad (7.34)$$

The inequality (7.34) holds since  $\epsilon_i < \mu_{\min} x_{\min}^* / \mu_{\max}$ , by (6.5). To apply the large deviation bounds of Lemma 6.7 in the proofs of Lemmas 7.3 and 7.6, we require that

$$L_0^r > \max \left\{ \frac{2}{\epsilon_i}, \frac{2}{\tilde{\epsilon}_i}, \frac{2}{(x_{a(i)}^* - \epsilon_i) \epsilon_{2,i}}, \frac{4}{(x_j^* + \epsilon_i) \epsilon_i}, \frac{2}{(x_j^* + \epsilon_i) \tilde{\epsilon}_i}, \frac{2}{(x_j^* + \epsilon_{1,i}) \epsilon_{3,j}} : i \in \mathcal{I}, j \in \underline{\mathcal{J}}_i \right\}. \quad (7.35)$$

**Assumption 7.2** We henceforth assume that  $r^* \geq 1$  is fixed such that for all  $r \geq r^*$ , the following hold:



(i) conditions (7.12), (7.14)–(7.19), (7.21)–(7.23), (7.25)–(7.30), and (7.35) hold for all  $i \in \mathcal{I}$ ,  $\iota \in \mathcal{I}$ , with the exceptions that (7.14)–(7.19) do not need to hold if  $i = i^*$ , (7.21)–(7.23) and (7.25)–(7.26) do not need to hold if  $i$  is a non-transition class, and (7.30) does not need to hold if  $\iota$  is a non-transition class,

(ii) for each  $i \in \mathcal{I}$ ,  $j \in \underline{\mathcal{J}}_i$ , (6.20)–(6.22) hold with

- (a)  $\lambda_i^r$  in place of  $\nu^r$ ,  $\lambda_i$  in place of  $\nu$ , and any of  $\epsilon_i$ ,  $\tilde{\epsilon}_i$ ,  $\check{\epsilon}_i$ , or  $\check{\epsilon}_i/2$  in place of  $\epsilon$  there,
- (b)  $\mu_j^r$  in place of  $\nu^r$ ,  $\mu_j$  in place of  $\nu$ , and any of  $\epsilon_i$ ,  $\epsilon_i/2$ ,  $\tilde{\epsilon}_i$ ,  $\check{\epsilon}_i$ , or  $\epsilon_{3,j}$ , in place of  $\epsilon$  there, and
- (c)  $\mu_{a(i)}^r$  in place of  $\nu^r$ ,  $\mu_{a(i)}$  in place of  $\nu$ , and  $\epsilon_{2,i}$  in place of  $\epsilon$  there.

**Remark.** The condition before (6.20)–(6.22) that  $0 < \epsilon < \nu/2$  is automatically satisfied for the choices of  $\epsilon$ ,  $\nu$  in Assumption 7.2.

## 7.2 Induction Setup

In the sequel we will use induction on  $i$  to show that the following (I)–(II) hold for each  $i \in \mathcal{I} \setminus \{i^*\}$ , and we will show that (III) below holds when  $i^*$  is a transition buffer. Recall the definitions of  $r^*$ ,  $s_i^r$ ,  $t_i^r$ ,  $d_s^r$ ,  $M^r$  from Section 7.1. Note in particular that  $r^*$  is independent of  $t$ . We consider the following, (I)–(II), for  $i \in \mathcal{I} \setminus \{i^*\}$ .

(I) For all  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , and each  $k \in \underline{\mathcal{K}}_i$ ,

$$(I.1) \quad \mathbf{P} \left( \sup_{\tau_{i,0}^r \leq s \leq r^2 t} |R_i^r(s)| \geq L_i^r - |\underline{\mathcal{J}}_i| \right) \\ \leq p_{1,i}(r^2 t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-C_{1,i}^{(4)} r^2 t) \right),$$

$$(I.2) \quad \mathbf{P} \left( I_k^r(\tau_{i,0}^r) \geq t_i^r \right) \leq p_{2,i}(r^2 t) \left( C_{2,i}^{(1)} \exp(-C_{2,i}^{(2)} L_0^r) + C_{2,i}^{(3)} \exp(-C_{2,i}^{(4)} r^2 t) \right),$$

$$(I.3) \quad \mathbf{P} \left( I_k^r(r^2 t) - I_k^r(\tau_{i,0}^r) > 0, \tau_{i,0}^r < r^2 t \right) \\ \leq p_{3,i}(r^2 t) \left( C_{3,i}^{(1)} \exp(-C_{3,i}^{(2)} L_0^r) + C_{3,i}^{(3)} \exp(-C_{3,i}^{(4)} r^2 t) \right),$$

where  $p_{1,i}$ ,  $p_{2,i}$ ,  $p_{3,i}$  are polynomials with non-negative coefficients, and  $C_{l,i}^{(m)}$ ,  $l = 1, 2, 3$ ,  $m = 1, 2, 3, 4$  are positive constants; the polynomials and constants do not depend on  $t$  or  $r$ . The polynomials  $p_{1,i}$  and  $p_{3,i}$  have degree at most  $i + 1$  and  $p_{2,i}$  has degree at most  $i$ .

(II) For all  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , and each  $\iota \in \mathcal{I}$  with  $\iota > i$ ,

$$(II.1) \quad \sup_{n \geq 1} \mathbf{P} \left( T_{\iota, a(i)}^{r,n}(s_\iota^r) \geq (x_{a(i)}^* + \epsilon_i) s_\iota^r, s_\iota^r \leq \beta_{\iota,n}^r, \tau_{\iota, 2n-1}^r \leq r^2 t \right) \\ \leq p_{4,i}(r^2 t) \left( C_{4,i}^{(1)} \exp(-C_{4,i}^{(2)} L_0^r) + C_{4,i}^{(3)} \exp(-C_{4,i}^{(4)} r^2 t) \right),$$

$$(II.2) \quad \text{if } \iota \text{ is a transition class,} \\ \sup_{n \geq 1} \mathbf{P} \left( {}^d T_{\iota, a(i)}^{r,n}({}^d s_\iota^r) \leq (x_{a(i)}^* - \epsilon_i) {}^d s_\iota^r, {}^d \tau_{\iota, 2n-1}^r \leq r^2 t \right) \\ \leq p_{5,i}(r^2 t) \left( C_{5,i}^{(1)} \exp(-C_{5,i}^{(2)} L_0^r) + C_{5,i}^{(3)} \exp(-C_{5,i}^{(4)} r^2 t) \right),$$

$$(II.3) \quad \mathbf{P} \left( T_{a(i)}^r(t_\iota^r) \leq (x_{a(i)}^* - \epsilon_i) t_\iota^r \right) \\ \leq p_{6,i}(r^2 t) \left( C_{6,i}^{(1)} \exp(-C_{6,i}^{(2)} L_0^r) + C_{6,i}^{(3)} \exp(-C_{6,i}^{(4)} r^2 t) \right),$$

where  $p_{4,i}, p_{5,i}, p_{6,i}$  are polynomials (of degree at most  $i + 1$ ) with non-negative coefficients, and  $C_{l,i}^{(m)}, l = 4, 5, 6, m = 1, 2, 3, 4$  are positive constants; the polynomials and constants do not depend on  $t$  or  $r$ .

**Remark.** The appearance of the variable  $t$  in the right side of (II.3) stems from an estimate obtained in (7.115) involving the number of class  $i$  jobs in the system at time  $t_l^r \leq M^r \leq r^2 t$ .

We consider the following, (III), if  $i^*$  is a transition class.

(III) For all  $r \geq r^*, t > 0$  satisfying  $r^2 t \geq M^r$ , and each  $k \in \underline{K}_{i^*}$ ,

$$\begin{aligned} \text{(III.1)} \quad & \mathbf{P} \left( \inf_{\tau_{i^*,0}^r \leq s \leq r^2 t} R_{i^*}^r(s) \leq -L_{i^*}^r + |\underline{\mathcal{J}}_{i^*}| \right) \\ & \leq p_{1,i^*}(r^2 t) \left( C_{1,i^*}^{(1)} \exp(-C_{1,i^*}^{(2)} L_0^r) + C_{1,i^*}^{(3)} \exp(-C_{1,i^*}^{(4)} r^2 t) \right), \\ \text{(III.2)} \quad & \mathbf{P} \left( I_k^r(\tau_{i^*,0}^r) \geq t_{i^*}^r \right) \leq p_{2,i^*}(r^2 t) \left( C_{2,i^*}^{(1)} \exp(-C_{2,i^*}^{(2)} L_0^r) + C_{2,i^*}^{(3)} \exp(-C_{2,i^*}^{(4)} r^2 t) \right), \\ \text{(III.3)} \quad & \mathbf{P} \left( I_k^r(r^2 t) - I_k^r(\tau_{i^*,0}^r) > 0, \tau_{i^*,0}^r < r^2 t \right) \\ & \leq p_{3,i^*}(r^2 t) \left( C_{3,i^*}^{(1)} \exp(-C_{3,i^*}^{(2)} L_0^r) + C_{3,i^*}^{(3)} \exp(-C_{3,i^*}^{(4)} r^2 t) \right), \end{aligned}$$

where  $p_{1,i^*}, p_{2,i^*}, p_{3,i^*}$  are polynomials with non-negative coefficients, and  $C_{l,i^*}^{(m)}, l = 1, 2, 3, m = 1, 2, 3, 4$  are positive constants; the polynomials and constants do not depend on  $t$  or  $r$ . The polynomials  $p_{1,i^*}$  and  $p_{3,i^*}$  have degree at most  $\mathbf{I} + 1$  and  $p_{2,i^*}$  has degree at most  $\mathbf{I}$ .

Property (I) is used to obtain the conclusions (7.1)–(7.3) in Theorem 7.1. Property (II) describes the properties associated with buffer  $i$  that are carried along in the induction proof to prove (I) for buffers  $\iota > i$  and to prove Property (III), which in turn is used to prove (7.4)–(7.6). The induction proof depends on the following lemmas that are proved in Sections 7.3–7.6 below.

**Lemma 7.3** Fix  $i \in \mathcal{I} \setminus \{i^*\}$ . Suppose that for all  $r \geq r^*, t > 0$  satisfying  $r^2 t \geq M^r$ , the following properties (i)–(iii) hold for all  $j \in \underline{\mathcal{J}}_i$  whenever  $i$  is a transition buffer, and (iv) holds for  $i$  whether it is a transition buffer or not:

$$\begin{aligned} \text{(i)} \quad & \sup_{n \geq 1} \mathbf{P} \left( T_{i,j}^{r,n}(s_i^r) \leq (x_j^* - \epsilon_i) s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2 t \right) \\ & \leq p_{7,i}(r^2 t) \left( C_{7,i}^{(1)} \exp(-C_{7,i}^{(2)} L_0^r) + C_{7,i}^{(3)} \exp(-C_{7,i}^{(4)} r^2 t) \right), \\ \text{(ii)} \quad & \sup_{n \geq 1} \mathbf{P} \left( d_{i,j}^{r,n}(d_s^r) \geq (x_j^* + \epsilon_i) d_s^r, d_{i,2n-1}^r \leq r^2 t \right) \\ & \leq p_{8,i}(r^2 t) \left( C_{8,i}^{(1)} \exp(-C_{8,i}^{(2)} L_0^r) + C_{8,i}^{(3)} \exp(-C_{8,i}^{(4)} r^2 t) \right), \\ \text{(iii)} \quad & \mathbf{P} \left( T_j^r(t_i^r) \geq (x_j^* + \epsilon_i) t_i^r \right) \leq p_{9,i}(r^2 t) \left( C_{9,i}^{(1)} \exp(-C_{9,i}^{(2)} L_0^r) + C_{9,i}^{(3)} \exp(-C_{9,i}^{(4)} r^2 t) \right), \\ \text{(iv)} \quad & \sup_{n \geq 1} \mathbf{P} \left( T_{i,a(i)}^{r,n}(s_i^r) \leq (\hat{x}_{i,i} - \epsilon_i) s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2 t \right) \\ & \leq p_{10,i}(r^2 t) \left( C_{10,i}^{(1)} \exp(-C_{10,i}^{(2)} L_0^r) + C_{10,i}^{(3)} \exp(-C_{10,i}^{(4)} r^2 t) \right), \end{aligned}$$

where  $p_{7,i}, p_{8,i}, p_{9,i}, p_{10,i}$  are polynomials (of degree at most  $i$ ) with non-negative coefficients, and  $C_{l,i}^{(m)}, l = 7, 8, 9, 10, m = 1, 2, 3, 4$  are positive constants; the polynomials and constants are independent of  $t$  and  $r$ . Then, (I) holds for  $i$ .

**Lemma 7.4** Let  $i \in \mathcal{I} \setminus \{i^*\}$ . Suppose that (I) and (II) hold with  $i'$  in place of  $i$ , for all  $i' < i$ . Then, if  $i$  is a transition buffer, (i)–(iii) of Lemma 7.3 hold for all  $r \geq r^*, t > 0$  satisfying  $r^2 t \geq M^r$ , and each  $j \in \underline{\mathcal{J}}_i$ .

**Lemma 7.5** *Let  $i \in \mathcal{I} \setminus \{i^*\}$ . Suppose that (I) and (II) hold with  $i'$  in place of  $i$ , for all  $i' < i$ . Then (iv) of Lemma 7.3 holds for all  $r \geq r^*$  and  $t > 0$  satisfying  $r^2 t \geq M^r$ .*

**Lemma 7.6** *Let  $i \in \mathcal{I} \setminus \{i^*\}$ . Suppose that (I) and (II) hold with  $i'$  in place of  $i$ , for all  $i' < i$ . Then (II) holds for  $i$ .*

Lemmas 7.3–7.6 combined with a formal induction yield the following.

**Theorem 7.7** (I) and (II) hold for each  $i \in \mathcal{I} \setminus \{i^*\}$ . In addition, (III) holds if  $i^*$  is a transition class.

With Lemmas 7.3–7.6 in place, the steps in the induction argument used in proving Theorem 7.7 are as follows (assuming the ordering of the buffers is as described in Section 6.1).

1. Fix  $i \in \mathcal{I} \setminus \{i^*\}$ , and assume that (I) and (II) hold with  $i'$  in place of  $i$ , for all  $i' < i$  (for  $i = 1$ , this is a vacuous assumption).
2. Use Lemma 7.4 to show that (i)–(iii) in Lemma 7.3 hold for all  $j \in \underline{\mathcal{J}}_i$  if  $i$  is a transition buffer, and use Lemma 7.5 to show that (iv) holds for  $i$ .
3. Apply Lemma 7.3 to conclude that (I) holds for  $i$ .
4. Apply Lemma 7.6 to conclude that (II) holds for  $i$ .

It then follows that (I) and (II) hold for all  $i \in \mathcal{I} \setminus \{i^*\}$ . Then, if  $i^*$  is a transition buffer, we combine the proof of Lemma 7.4 and parts of the proof of Lemma 7.3 (adapted for  $i = i^*$ ) with the fact that (II) holds for  $i = i^*$  and  $i < i^*$ , to prove that (III) holds.

The formal proof of Theorem 7.7 is given in Section 7.7. As a guide to the reader, before beginning the proofs of the lemmas and Theorems 7.7 and 7.1, we briefly describe some of the ideas involved in these proofs.

Lemma 7.3 is the main lemma that drives the induction. The result of this lemma yields that all queue length processes and idletime processes (except  $Q_{i^*}^r$  and  $I_{k^*}^r$ ) vanish on diffusion scale as  $r$  goes to infinity (for  $c$  sufficiently large). To obtain (I) for  $i \in \mathcal{I} \setminus \{i^*\}$ , our proof requires that, in addition to large deviation estimates on the primitive renewal processes, there are estimates on the allocation processes (in each excursion interval for the residual process associated with buffer  $i$ ) corresponding to the activities which process class  $i$  jobs, i.e., (i)–(iv) in Lemma 7.3 hold. When  $i$  is a transition buffer, the estimates (i)–(iii) are derived from the induction assumptions for (II) (with  $i$  replaced by  $i' < i$  in (II)), using the fact that, for each  $j \in \underline{\mathcal{J}}_i$ , the utilization of activity  $j$  is constrained by the (higher priority) activities for server  $k(j)$  which serve buffers in the layer below that server. These estimates are obtained in the proof of Lemma 7.4, which is contained in Section 7.4. Similarly, for (iv), the on-time of the activity,  $a(i)$ , that serves buffer  $i$  from above can be estimated (in an up excursion for  $R_i^r$ ) by having estimates, derived from the induction assumption for (II), on the (higher priority) activities associated with server  $k(a(i))$  (recall that buffers which have higher priority for server  $k(a(i))$  are all numbered lower than  $i$ ). These estimates are obtained in the proof of Lemma 7.5 in Section 7.5. Finally, the proof of Lemma 7.6 in Section 7.6 (which uses the fact that (I) holds for buffer  $i$  together with the induction assumption for (II), with  $i$  replaced by  $i' < i$  in (II)), completes the induction step by showing how to transition between layers.

For (III.1)–(III.3), assuming that  $i^*$  is a transition buffer, we first use the proof of Lemma 7.4 (cf. Section 7.4) to show that (ii) and (iii) in Lemma 7.3 hold with  $i^*$  in place of  $i$ ,  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , and all  $j \in \underline{\mathcal{J}}_{i^*}$ , given that (I) and (II) hold for all  $i < i^*$  (cf. (7.69) and (7.70), respectively). Then, (III.1)–(III.3) can be proved in a similar manner to that in the proof of Lemma 7.3 in Section 7.3 (with  $i^*$  in place of  $i$  there).

Properties (I) and (III) are used to prove (7.1)–(7.3) and (7.4)–(7.6), respectively, of Theorem 7.1, for a sufficiently large constant  $c$  appearing in the definition of  $L_0^r$ .

### 7.3 Estimates on Allocation Processes Imply Residual Processes Stay Near Zero – Proof of Lemma 7.3

**Proof of Lemma 7.3.** Fix  $i \in \mathcal{I} \setminus \{i^*\}$ . Suppose that for all  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , and each  $j \in \underline{\mathcal{J}}_i$ , (i)–(iv) hold. The proof that follows is an extension of the proof of Theorem 7.2 in [3]. For the special case when class  $i$  is not a transition class, we have  $\underline{\mathcal{K}}_i = \emptyset$ ,  $\underline{\mathcal{J}}_i = \emptyset$ ,  $\tau_{i,0}^r \equiv 0$ , and  ${}^d R_i^r \equiv \mathbf{0}$ . This implies that (I.2) holds trivially and the part of the proof of (I.1) given below for down excursions of  $R_i^r$  is not needed.

*Proof of (I.1).* The idea of this proof is: (a) to show that the number of excursions of  $R_i^r$  from the zero level is at most of order  $r^2 t$ , with probability at least  $1 - K_1 \exp(-K'_1 r^2 t)$  where  $K_1, K'_1$  are constants not depending on  $r$  or  $t$ , and then (b) to estimate the probability that  $|R_i^r|$  reaches the level  $L_i^r - |\underline{\mathcal{J}}_i|$  or higher during any of the first  $O(r^2 t)$  excursions. Using large deviation estimates for the renewal processes  $A_i^r$  and  $S_j^r$ ,  $j \in \mathcal{J}_i$ , and the assumptions of Lemma 7.3, this probability will be shown to be dominated by an expression of the form  $p_i(r^2 t)(C \exp(-C' L_i^r) + C'' \exp(-C''' r^2 t))$ , for all  $r \geq r^*$  and  $t > 0$  satisfying  $r^2 t > M^r$ , where  $p_i$  is a polynomial of degree at most  $i$  with non-negative coefficients, and  $C, C', C'', C'''$  are non-negative constants not depending on  $r$  or  $t$  (cf. (7.57)). Then (a) and (b) are combined to yield (I.1).

We first consider the up excursions of  $R_i^r$ . For  $n \geq 1$ ,  $r \geq 1$ ,  $j \in \mathcal{J}_i$ , on  $\{\tau_{i,2n-1}^r < \infty\}$ , we define shifted renewal processes, for  $s \geq 0$ , as follows.

$$A_{i,i}^{r,n}(s) = A_i^r(\tau_{i,2n-1}^r + s) - A_i^r(\tau_{i,2n-1}^r), \quad (7.36)$$

$$S_{i,j}^{r,n}(s) = S_j^r(T_j^r(\tau_{i,2n-1}^r) + s) - S_j^r(T_j^r(\tau_{i,2n-1}^r)), \quad (7.37)$$

$$\check{A}_{i,i}^{r,n}(s) = \sup\{m \geq 0 : \xi_i^r(A_i^r(\tau_{i,2n-1}^r) + m) - \xi_i^r(A_i^r(\tau_{i,2n-1}^r)) \leq s\}, \quad (7.38)$$

$$\check{S}_{i,j}^{r,n}(s) = \sup\{m \geq 0 : \eta_j^r(S_j^r(T_j^r(\tau_{i,2n-1}^r)) + m) - \eta_j^r(S_j^r(T_j^r(\tau_{i,2n-1}^r))) \leq s\}, \quad (7.39)$$

and, for concreteness, on  $\{\tau_{i,2n-1}^r = \infty\}$  we define  $A_{i,i}^{r,n}, S_{i,j}^{r,n}, \check{A}_{i,i}^{r,n}, \check{S}_{i,j}^{r,n}$  to be identically zero. Recall the definition of  $T_{i,j}^{r,n}$  from Section 6.4.

Consider the  $n^{\text{th}}$  up excursion interval for  $R_i^r$ . We have that on  $\{\tau_{i,2n-1}^r < \infty\}$ , for  $0 \leq s \leq \beta_{i,n}^r$ ,

$$R_i^r(\tau_{i,2n-1}^r + s) = 1 + A_{i,i}^{r,n}(s) - \sum_{j \in \mathcal{J}_i} S_{i,j}^{r,n}(T_{i,j}^{r,n}(s)), \quad (7.40)$$

and, taking account of the fact that a new arrival to class  $i$  occurs at  $\tau_{i,2n-1}^r < \infty$  and a job may have been partially served by activity  $j \in \mathcal{J}_i$  at  $\tau_{i,2n-1}^r < \infty$ , we also have that, for  $s \geq 0$ ,

$$A_{i,i}^{r,n}(s) = \check{A}_{i,i}^{r,n}(s), \quad \text{and} \quad S_{i,j}^{r,n}(s) \geq \check{S}_{i,j}^{r,n}(s), \quad j \in \mathcal{J}_i. \quad (7.41)$$

By (6.5), we have that

$$\epsilon_i < \min \left\{ \frac{\mu_j}{2}, \frac{\lambda_i}{2}, \left( 1 - \frac{4\lambda_i + \delta_i}{4\lambda_i + \frac{3}{2}\delta_i} \right) \hat{x}_{i,i}, \left( 1 - \frac{4\lambda_i}{4\lambda_i + \delta_i} \right) x_j^* : j \in \mathcal{J}_i \right\}, \quad (7.42)$$

where  $\delta_i$  and  $\hat{x}_{i,i}$  are given in (7.13) and (6.2), respectively.

In the following, it is assumed that  $r \geq r^*$  and  $t > 0$  satisfies  $r^2 t \geq M^r$ . Then  $r^2 t > 2/\epsilon_i$  by (7.9)–(7.12), and (7.35). Define  $n_i^r = \lfloor (\lambda_i^r + \epsilon_i) r^2 t \rfloor + 1$ . Since each up excursion of  $R_i^r$  is initiated by an arrival to class  $i$ , using the large deviations bounds for renewal processes given in Lemma 6.7 and the choice of  $r^*$  (cf. Assumption 7.2), we have the following estimate of the probability that

at least  $n_i^r$  up excursions of  $R_i^r$  have been initiated in  $[0, r^2t]$  ( $\tau_{i,2n_i^r-1}^r$  is the beginning of the  $(n_i^r)^{\text{th}}$  up excursion interval):

$$\begin{aligned} \mathbf{P}(\tau_{i,2n_i^r-1}^r \leq r^2t) &\leq \mathbf{P}(A_i^r(r^2t) \geq n_i^r) \\ &\leq \mathbf{P}(A_i^r(r^2t) > (\lambda_i^r + \epsilon_i)r^2t) \\ &\leq K_1 \exp(-K_1' r^2t), \end{aligned} \quad (7.43)$$

where  $K_1 = \exp(\Lambda_i^{a,*}((\lambda_i(1 + \epsilon_i/3\lambda_i))^{-1})) > 0$ ,  $K_1' = \lambda_i \Lambda_i^{a,*}((\lambda_i(1 + \epsilon_i/3\lambda_i))^{-1}) > 0$ , depend on the Legendre-Fenchel transform  $\Lambda_i^{a,*}$  of the logarithmic moment generating function  $\Lambda_i^a$  of  $u_i(1)$  (cf. (3.3)),  $\lambda_i$ , and  $\epsilon_i$ , but are independent of  $t$  and  $r$ .

Now, decomposing the probability space according to whether the number of up excursions initiated in  $[0, r^2t]$  is greater than or equal to  $n_i^r$  or less than  $n_i^r$ , we have

$$\begin{aligned} &\mathbf{P}\left(R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \text{ some } s \in [0, r^2t]\right) \\ &\leq \mathbf{P}(\tau_{i,2n_i^r-1}^r \leq r^2t) \\ &\quad + \sum_{n=1}^{n_i^r-1} \mathbf{P}\left(R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \text{ some } s \in (\tau_{i,2n-1}^r, \tau_{i,2n}^r), \tau_{i,2n-1}^r \leq r^2t\right). \end{aligned} \quad (7.44)$$

For each positive integer  $n$ , let

$$\begin{aligned} \Upsilon_i^{r,n} &= \left\{ A_{i,j}^{r,n}(s_i^r) \leq (\lambda_j^r + \epsilon_i)s_i^r; S_{i,j}^{r,n}((\tilde{x}_j^i - \epsilon_i)s_i^r) \geq (\mu_j^r - \epsilon_i)(\tilde{x}_j^i - \epsilon_i)s_i^r, j \in \mathcal{J}_i; \right. \\ &\quad \left. T_{i,j}^{r,n}(s_i^r) > (\tilde{x}_j^i - \epsilon_i)s_i^r, j \in \mathcal{J}_i; \tau_{i,2n-1}^r \leq r^2t \right\}, \end{aligned} \quad (7.45)$$

where  $\tilde{x}_j^i = x_j^*$  for  $j \in \underline{\mathcal{J}}_i$  and  $\tilde{x}_j^i = \hat{x}_{i,i}$  for  $j = a(i)$ . The set  $\Upsilon_i^{r,n}$  is a ‘‘good’’ set, in the sense that, on it, various shifted stochastic processes can be bounded on one side at time  $s_i^r$  by certain linear functions. These bounds will enable us to show that on  $\Upsilon_i^{r,n}$ ,  $R_i^r$  will not reach the level  $L_i^r - |\underline{\mathcal{J}}_i|$  in the  $n^{\text{th}}$  up excursion interval  $(\tau_{i,2n-1}^r, \tau_{i,2n}^r)$  whose length is shorter than  $s_i^r$ .

Let

$$\rho_i^{r,n} = \xi_i^r(A_i^r(\tau_{i,2n-1}^r) + L_i^r - |\underline{\mathcal{J}}_i| - 1) - \xi_i^r(A_i^r(\tau_{i,2n-1}^r)) \text{ on } \{\tau_{i,2n-1}^r < \infty\}, \quad (7.46)$$

and let  $\rho_i^{r,n} \equiv 0$  on  $\{\tau_{i,2n-1}^r = \infty\}$ . Then,

$$\begin{aligned} &\{R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \text{ some } s \in (\tau_{i,2n-1}^r, \tau_{i,2n}^r), \tau_{i,2n-1}^r < \infty\} \\ &= \{R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \text{ some } s \in (\tau_{i,2n-1}^r, \tau_{i,2n}^r), \rho_i^{r,n} \leq \beta_{i,n}^r, \tau_{i,2n-1}^r < \infty\}, \end{aligned} \quad (7.47)$$

since, on  $\{\tau_{i,2n-1}^r < \infty\}$ ,  $\rho_i^{r,n}$  is the minimum possible amount of time required for  $R_i^r$  to reach the level  $L_i^r - |\underline{\mathcal{J}}_i|$  in the  $n^{\text{th}}$  up excursion. Thus,

$$\begin{aligned} &\mathbf{P}\left(R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \text{ some } s \in (\tau_{i,2n-1}^r, \tau_{i,2n}^r), \tau_{i,2n-1}^r \leq r^2t\right) \\ &\leq \mathbf{P}(\tau_{i,2n-1}^r \leq r^2t, (\Upsilon_i^{r,n})^c, \rho_i^{r,n} \leq \beta_{i,n}^r) \\ &\quad + \mathbf{P}\left(R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \text{ some } s \in (\tau_{i,2n-1}^r, \tau_{i,2n}^r), \Upsilon_i^{r,n}\right). \end{aligned} \quad (7.48)$$

Now, on  $\Upsilon_i^{r,n}$ , we have

$$\begin{aligned}
& 1 + A_{i,i}^{r,n}(s_i^r) - \sum_{j \in \mathcal{J}_i} S_{i,j}^{r,n}(T_{i,j}^{r,n}(s_i^r)) \\
& \leq 1 + (\lambda_i^r + \epsilon_i)s_i^r - \sum_{j \in \mathcal{J}_i} (\tilde{x}_j^i - \epsilon_i)(\mu_j^r - \epsilon_i)s_i^r \\
& \leq 1 + \left( \frac{2\lambda_i}{2\lambda_i + \delta_i} \left( \sum_{j \in \mathcal{J}_i} \tilde{x}_j^i \mu_j \right) \right. \\
& \quad \left. - \sum_{j \in \mathcal{J}_i} \left( \frac{4\lambda_i}{4\lambda_i + \delta_i} x_j^* \frac{4\lambda_i + \delta_i}{4\lambda_i + 2\delta_i} \mu_j \right) - \frac{4\lambda_i + \delta_i}{4\lambda_i + \frac{3}{2}\delta_i} \hat{x}_{i,i} \frac{4\lambda_i + \frac{3}{2}\delta_i}{4\lambda_i + 2\delta_i} \mu_{a(i)} \right) s_i^r \\
& = 1 - \frac{\delta_i}{4\lambda_i + 2\delta_i} \hat{x}_{i,i} \mu_{a(i)} s_i^r < 0,
\end{aligned} \tag{7.49}$$

where, in the second inequality we have used (7.16)–(7.18) along with (7.42), and in the last inequality we have used (7.19). It then follows from (7.40) that

$$\beta_{i,n}^r = \tau_{i,2n}^r - \tau_{i,2n-1}^r < s_i^r \text{ on } \Upsilon_i^{r,n}. \tag{7.50}$$

Furthermore, on  $\Upsilon_i^{r,n}$  for  $0 \leq s \leq s_i^r$ ,

$$1 + A_{i,i}^{r,n}(s) - \sum_{j \in \mathcal{J}_i} S_{i,j}^{r,n}(T_{i,j}^{r,n}(s)) \leq 1 + (\lambda_i^r + \epsilon_i)s_i^r = L_i^r - |\underline{\mathcal{J}}_i| - 1, \tag{7.51}$$

by (7.7). Hence by (7.40) we have that on  $\Upsilon_i^{r,n}$ ,  $R_i^r(s) < L_i^r - |\underline{\mathcal{J}}_i|$  for  $s \in (\tau_{i,2n-1}^r, \tau_{i,2n}^r)$ , since  $r \geq r^*$ . Thus the last probability in (7.48) is zero.

Splitting the probability space according to whether  $\rho_i^{r,n} < s_i^r$  or  $\rho_i^{r,n} \geq s_i^r$ , and discarding some qualifiers, we obtain

$$\begin{aligned}
& \mathbf{P}(\tau_{i,2n-1}^r \leq r^2 t, (\Upsilon_i^{r,n})^c, \rho_i^{r,n} \leq \beta_{i,n}^r) \\
& \leq \mathbf{P}(\rho_i^{r,n} < s_i^r, \tau_{i,2n-1}^r \leq r^2 t) \\
& \quad + \mathbf{P}(A_{i,i}^{r,n}(s_i^r) > (\lambda_i^r + \epsilon_i)s_i^r, \tau_{i,2n-1}^r \leq r^2 t) \\
& \quad + \sum_{j \in \mathcal{J}_i} \mathbf{P}(S_{i,j}^{r,n}(\tilde{x}_j^i - \epsilon_i)s_i^r) < (\mu_j^r - \epsilon_i)(\tilde{x}_j^i - \epsilon_i)s_i^r, \tau_{i,2n-1}^r \leq r^2 t) \\
& \quad + \sum_{j \in \mathcal{J}_i} \mathbf{P}(T_{i,j}^{r,n}(s_i^r) \leq (\tilde{x}_j^i - \epsilon_i)s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2 t).
\end{aligned} \tag{7.52}$$

Now, the set  $\{\tau_{i,2n-1}^r < \infty\}$  is contained in the set  $\{\mathcal{T}_{n,i}^r \in \mathbb{N}^{\mathbf{I}} \times \mathbb{N}^{\mathbf{J}}\}$  (cf. Lemma 6.5). Using Lemmas 6.5 and 6.6, we conclude that on  $\{\mathcal{T}_{n,i}^r \in \mathbb{N}^{\mathbf{I}} \times \mathbb{N}^{\mathbf{J}}\}$ , the conditional distribution of  $\{u_i^r(A_i^r(\tau_{i,2n-1}^r) + m), m = 1, 2, 3, \dots\}$  given  $\mathcal{F}_{\mathcal{T}_{n,i}^r}^r$  is equal to that of a sequence of strictly positive independent random variables where the members indexed by  $m = 2, 3, \dots$  are identically distributed with the same distribution as  $u_i^r(1)$  (for the application of these lemmas, observe that  $L_0^r \geq \mathbf{J} + 1$  by (7.35) and (6.5) since  $\epsilon_i \leq \hat{\epsilon} < 1/2(\mathbf{J} + 1)$ , and  $u_i^r(A_i^r(\tau_{i,2n-1}^r) + 1) \in \mathcal{F}_{\mathcal{T}_{n,i}^r}^r$ ). Then,

as in equation (84) of [3], since we have assumed that  $r \geq r^*$  and  $t > 0$  satisfies  $r^2t \geq Mr$ , we have

$$\begin{aligned}
& \mathbf{P} \left( A_{i,i}^{r,n}(s_i^r) > (\lambda_i^r + \epsilon_i) s_i^r, \tau_{i,2n-1}^r \leq r^2t \right) \\
& \leq \mathbf{E} \left( 1_{\{\mathcal{T}_{n,i}^r \in \mathbb{N}^I \times \mathbb{N}^J\}} \mathbf{P} \left( \check{A}_{i,i}^{r,n}(s_i^r) > (\lambda_i^r + \epsilon_i) s_i^r \mid \mathcal{F}_{\mathcal{T}_{n,i}^r}^r \right) \right) \\
& \leq \exp \left( -((\lambda_i^r + \epsilon_i) s_i^r - 1) \Lambda_i^{a,*} \left( \frac{1}{\lambda_i} \left( \frac{1}{1 + \frac{\epsilon_i}{3\lambda_i}} \right) \right) \right) \\
& \leq K_2 \exp(-K_2' L_i^r), \tag{7.53}
\end{aligned}$$

by Lemma 6.7 (since  $s_i^r > 2/\epsilon_i$  by (7.9), (7.12), and (7.35)), and where  $K_2 = \exp((|\mathcal{J}_i| + 3)\Lambda_i^{a,*}((\lambda_i(1 + \epsilon_i/3\lambda_i))^{-1})) > 0$  and  $K_2' = \Lambda_i^{a,*}((\lambda_i(1 + \epsilon_i/3\lambda_i))^{-1}) > 0$  do not depend on  $t, n$ , or  $r$ . Similarly, since  $r \geq r^*$  and  $t > 0$  satisfies  $r^2t \geq Mr$ , we have

$$\begin{aligned}
& \mathbf{P}(\rho_i^{r,n} < s_i^r, \tau_{i,2n-1}^r \leq r^2t) \\
& \leq \mathbf{E} \left( 1_{\{\mathcal{T}_{n,i}^r \in \mathbb{N}^I \times \mathbb{N}^J\}} \mathbf{P} \left( \rho_i^{r,n} < s_i^r \mid \mathcal{F}_{\mathcal{T}_{n,i}^r}^r \right) \right) \\
& \leq \mathbf{E} \left( 1_{\{\mathcal{T}_{n,i}^r \in \mathbb{N}^I \times \mathbb{N}^J\}} \mathbf{P} \left( \check{A}_{i,i}^{r,n}(s_i^r) \geq L_i^r - |\mathcal{J}_i| - 1 \mid \mathcal{F}_{\mathcal{T}_{n,i}^r}^r \right) \right) \\
& \leq \mathbf{E} \left( 1_{\{\mathcal{T}_{n,i}^r \in \mathbb{N}^I \times \mathbb{N}^J\}} \mathbf{P} \left( \check{A}_{i,i}^{r,n}(s_i^r) > (\lambda_i^r + \epsilon_i) s_i^r \mid \mathcal{F}_{\mathcal{T}_{n,i}^r}^r \right) \right) \\
& \leq K_2 \exp(-K_2' L_i^r), \tag{7.54}
\end{aligned}$$

where the third inequality follows by the definition of  $s_i^r$  (cf. (7.7)).

In a similar manner (cf. (85)–(86) of [3]), since  $r \geq r^*$  and  $t > 0$  satisfies  $r^2t \geq Mr$ , we have for all  $j \in \mathcal{J}_i$ ,

$$\begin{aligned}
& \mathbf{P} \left( S_{i,j}^{r,n}((\tilde{x}_j^i - \epsilon_i) s_i^r) < (\mu_j^r - \epsilon_i)(\tilde{x}_j^i - \epsilon_i) s_i^r, \tau_{i,2n-1}^r \leq r^2t \right) \\
& \leq \mathbf{E} \left( 1_{\{\mathcal{T}_{n,i}^r \in \mathbb{N}^I \times \mathbb{N}^J\}} \mathbf{P} \left( \check{S}_{i,j}^{r,n}((\tilde{x}_j^i - \epsilon_i) s_i^r) < (\mu_j^r - \epsilon_i)(\tilde{x}_j^i - \epsilon_i) s_i^r, \tau_{i,2n-1}^r \leq r^2t \mid \mathcal{F}_{\mathcal{T}_{n,i}^r}^r \right) \right) \\
& \leq \exp \left( -(\mu_j^r - 2\epsilon_i)(\tilde{x}_j^i - \epsilon_i) s_i^r \Lambda_j^{s,*} \left( \frac{1}{\mu_j^r} \left( 1 + \frac{\epsilon_i}{2\mu_j^r} \right) \right) \right) \\
& \quad + \mathbf{E} \left( 1_{\{\mathcal{T}_{n,i}^r \in \mathbb{N}^I \times \mathbb{N}^J\}} \mathbf{P} \left( v_{i,j}^{r,n} > \frac{\epsilon_i}{2\mu_j^r}(\tilde{x}_j^i - \epsilon_i) s_i^r, \tau_{i,2n-1}^r \leq r^2t \mid \mathcal{F}_{\mathcal{T}_{n,i}^r}^r \right) \right), \tag{7.55}
\end{aligned}$$

by Lemma 6.7, where  $v_{i,j}^{r,n} = v_j^r(S_j^r(T_j^r(\tau_{i,2n-1}^r)) + 1)$ , and  $\Lambda_j^{s,*}$ ,  $j \in \mathcal{J}_i$  is the Legendre-Fenchel transform of the logarithmic moment generating function  $\Lambda_j^s$  of  $v_j(1)$ ,  $j \in \mathcal{J}_i$  (cf. (3.4)). (Note here that when we use (6.25) of Lemma 6.7, we do not need the condition that  $s > 2/\epsilon$  required in (6.24), i.e., in the case above we do not require that  $(\tilde{x}_j^i - \epsilon_i) s_i^r > 2/\epsilon_i$ .) In a similar manner to that for (87) in [3], using (6.27) and (6.24), since  $r \geq r^*$  and  $t > 0$  satisfies  $r^2t \geq Mr$ , we have for all  $j \in \mathcal{J}_i$ ,

$$\begin{aligned}
& \mathbf{P} \left( v_{i,j}^{r,n} > \frac{\epsilon_i}{2\mu_j^r}(\tilde{x}_j^i - \epsilon_i) s_i^r, \tau_{i,2n-1}^r \leq r^2t \right) \\
& \leq \mathbf{P} \left( \max_{m=1}^{S_j^r(r^2t)+1} v_j^r(m) > \frac{\epsilon_i}{2\mu_j^r}(\tilde{x}_j^i - \epsilon_i) s_i^r \right) \\
& \leq \mathbf{P} \left( \max_{m=1}^{\lfloor (\mu_j^r + \epsilon_i)r^2t \rfloor + 1} v_j^r(m) > \frac{\epsilon_i}{2\mu_j^r}(\tilde{x}_j^i - \epsilon_i) s_i^r \right) \\
& \quad + \mathbf{P} \left( S_j^r(r^2t) > (\mu_j^r + \epsilon_i)r^2t \right) \\
& \leq (\lfloor (\mu_j^r + \epsilon_i)r^2t \rfloor + 1) K_4 \exp(-K_4' s_i^r) + K_5 \exp(-K_5' r^2t), \tag{7.56}
\end{aligned}$$

where, by (6.27),  $K_4 = \max\{\exp(\Lambda_j^s(l_0)) : j \in \mathcal{J}_i\} > 0$ ,  $K'_4 = \min\{(l_0\epsilon_i/2\mu_j)(\tilde{x}_j^i - \epsilon_i) : j \in \mathcal{J}_i\} > 0$  and  $0 < l_0 \in \mathcal{O}_0$  (cf. (3.4)). Also, by (6.24) and since  $r^2t \geq M^r > 2/\epsilon_i$ ,  $K_5 = \max\{\exp(\Lambda_j^{s,*}((\mu_j(1 + \epsilon_i/3\mu_j))^{-1})) : j \in \mathcal{J}_i\} > 0$  and  $K'_5 = \min\{\mu_j\Lambda_j^{s,*}((\mu_j(1 + \epsilon_i/3\mu_j))^{-1}) : j \in \mathcal{J}_i\} > 0$ . Note that  $K_4, K'_4, K_5, K'_5$  do not depend on  $t, n$ , or  $r$ . It follows that the last term in (7.55) is bounded by the expression in (7.56).

Combining all of the above (from (7.43) onwards), we have for all  $r \geq r^*$  and  $t > 0$  satisfying  $r^2t \geq M^r$ ,

$$\begin{aligned} & \mathbf{P}\left(R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \text{ some } s \in [0, r^2t]\right) \\ & \leq K_1 \exp(-K'_1 r^2t) \\ & \quad + (n_i^r - 1) \left\{ 2K_2 \exp(-K'_2 L_i^r) + |\mathcal{J}_i| K_3 \exp(-K'_3 s_i^r) \right. \\ & \quad + \sum_{j \in \mathcal{J}_i} ([(\mu_j^r + \epsilon_i)r^2t] + 1) K_4 \exp(-K'_4 s_i^r) + |\mathcal{J}_i| K_5 \exp(-K'_5 r^2t) \\ & \quad \left. + \sum_{j \in \mathcal{J}_i} \sup_{n < n_i^r} \mathbf{P}(T_{i,j}^{r,n}(s_i^r) \leq (\tilde{x}_j^i - \epsilon_i)s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2t) \right\}, \end{aligned} \quad (7.57)$$

where  $K_3 = 1$ ,  $K'_3 = \min\{(\mu_j - 2\epsilon_i)(\tilde{x}_j^i - \epsilon_i)\Lambda_j^{s,*}(\mu_j^{-1}(1 + \epsilon_i/2\mu_j)) : j \in \mathcal{J}_i\} > 0$  from (7.55).

Now we consider the down excursions of  $R_i^r$ . For this, we assume that  $\underline{\mathcal{J}}_i \neq \emptyset$ , i.e., we assume that class  $i$  is a transition class, since (7.57) suffices for the proof that (I.1) holds when class  $i$  is a non-transition class. The treatment of the down excursions is very similar to that for the up excursions except that we do not need to introduce an analogue of  $\rho_i^{r,n}$  and the effects of the arrivals and services on pushing  $R_i^r$  towards zero are reversed. In particular, one can show that for all  $r \geq r^*$  and  $t > 0$  satisfying  $r^2t \geq M^r$ ,

$$\begin{aligned} & \mathbf{P}\left(R_i^r(s) \leq |\underline{\mathcal{J}}_i| - L_i^r \text{ some } s \in [\tau_{i,0}^r, r^2t]\right) \\ & \leq K_6 \exp(-K'_6 r^2t) + (d_n^r - 1) \left\{ K_7 \exp(-K'_7 d_s^r) \right. \\ & \quad + ([(\lambda_i^r + \epsilon_i)r^2t] + 1) K_8 \exp(-K'_8 d_s^r) + K_9 \exp(-K'_9 r^2t) \\ & \quad + K_{10} \exp(-K'_{10} d_s^r) \\ & \quad \left. + \sum_{j \in \underline{\mathcal{J}}_i} \sup_{n < d_n^r} \mathbf{P}\left(T_{i,j}^{r,n}(d_s^r) \geq (x_j^* + \epsilon_i)d_s^r, d_{\tau_{i,2n-1}^r}^r \leq r^2t\right) \right\}, \end{aligned} \quad (7.58)$$

where  $d_n^r = \left\lceil \sum_{j \in \underline{\mathcal{J}}_i} (\mu_j^r + \epsilon_i)r^2t \right\rceil + |\underline{\mathcal{J}}_i|$ ,  $K_6 = |\underline{\mathcal{J}}_i| \max\{\exp(\Lambda_j^{s,*}((\mu_j(1 + \epsilon_i/3\mu_j))^{-1})) : j \in \underline{\mathcal{J}}_i\} > 0$ ,  $K'_6 = \min\{\mu_j\Lambda_j^{s,*}((\mu_j(1 + \epsilon_i/3\mu_j))^{-1}) : j \in \underline{\mathcal{J}}_i\} > 0$ ,  $K_7 = 1$ ,  $K'_7 = (\lambda_i - 2\epsilon_i)\Lambda_i^{a,*}(\lambda_i^{-1}(1 + \epsilon_i/2\lambda_i)) > 0$ ,  $K_8 = \exp(\Lambda_i^a(l_0)) > 0$ ,  $K'_8 = l_0\epsilon_i/2\lambda_i > 0$ ,  $0 < l_0 \in \mathcal{O}_0$ ,  $K_9 = \exp(\Lambda_i^{a,*}((\lambda_i(1 + \epsilon_i/3\lambda_i))^{-1})) > 0$ ,  $K'_9 = \lambda_i\Lambda_i^{a,*}((\lambda_i(1 + \epsilon_i/3\lambda_i))^{-1}) > 0$ ,  $K_{10} = |\underline{\mathcal{J}}_i| \max\{\exp(\Lambda_j^{s,*}((\mu_j(1 + \epsilon_i/6\mu_j))^{-1})) : j \in \underline{\mathcal{J}}_i\} > 0$ , and  $K'_{10} = \min\{\mu_j(x_j^* + \epsilon_i)\Lambda_j^{s,*}((\mu_j(1 + \epsilon_i/6\mu_j))^{-1}) : j \in \underline{\mathcal{J}}_i\} > 0$ . Details of the argument for this can be found in [2].

On combining the results (7.57) and (7.58) for the up and down excursions, assumptions (i), (ii), and (iv) of Lemma 7.3, the definitions of  $n_i^r, d_n^r, s_i^r, d_s^r, L_i^r$ , and the fact that  $s_i^r \geq L_0^r$  and  $d_s^r \geq L_0^r$  for all  $i \in \mathcal{I}$  (by (7.9) and (7.12)), it follows that for  $r \geq r^*$  and  $t > 0$  satisfying  $r^2t \geq M^r$ , we have that (I.1) holds. Note that since  $i \geq 1$ , terms involving  $(r^2t)^2$  are absorbed in the polynomial term  $p_{1,i}(r^2t)$ .

*Proof of (I.2).* Since  $\underline{\mathcal{K}}_i = \emptyset$ , if  $i$  is a non-transition class, (I.2) trivially holds in this case. So it



suffices to consider the case when  $i$  is a transition class. Note that for  $0 \leq s \leq \tau_{i,0}^r$ ,

$$Q_i^r(s) = A_i^r(s) - \sum_{j \in \underline{\mathcal{J}}_i} S_j^r(T_j^r(s)), \quad (7.59)$$

since activity  $a(i)$  will be turned off for such  $s$ . By (6.5), (6.7), and (7.24), we have

$$0 < \tilde{\epsilon}_i < \min \left\{ 1, \frac{\lambda_i - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j}{4(|\underline{\mathcal{J}}_i| + 2)}, \frac{\lambda_i}{2}, \frac{\mu_j}{2} : j \in \underline{\mathcal{J}}_i \right\}. \quad (7.60)$$

Now,

$$\begin{aligned} \mathbf{P}(I_k^r(\tau_{i,0}^r) \geq t_i^r) &\leq \mathbf{P}(\tau_{i,0}^r \geq t_i^r) \\ &\leq \mathbf{P}\left(A_i^r(t_i^r) - \sum_{j \in \underline{\mathcal{J}}_i} S_j^r(T_j^r(t_i^r)) \leq L_i^r\right), \quad \text{by (7.59),} \\ &\leq \mathbf{P}\left(A_i^r(t_i^r) \geq (\lambda_i^r - \tilde{\epsilon}_i)t_i^r, T_j^r(t_i^r) \leq (x_j^* + \epsilon_i)t_i^r \text{ for all } j \in \underline{\mathcal{J}}_i, \right. \\ &\quad \left. S_j^r((x_j^* + \epsilon_i)t_i^r) \leq (\mu_j^r + \tilde{\epsilon}_i)(x_j^* + \epsilon_i)t_i^r \text{ for all } j \in \underline{\mathcal{J}}_i, \right. \\ &\quad \left. ((\lambda_i^r - \tilde{\epsilon}_i) - \sum_{j \in \underline{\mathcal{J}}_i} (x_j^* + \epsilon_i)(\mu_j^r + \tilde{\epsilon}_i))t_i^r \leq L_i^r\right) \\ &\quad + \mathbf{P}\left(A_i^r(t_i^r) < (\lambda_i^r - \tilde{\epsilon}_i)t_i^r\right) \\ &\quad + \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P}\left(S_j^r((x_j^* + \epsilon_i)t_i^r) > (\mu_j^r + \tilde{\epsilon}_i)(x_j^* + \epsilon_i)t_i^r\right) \\ &\quad + \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P}\left(T_j^r(t_i^r) > (x_j^* + \epsilon_i)t_i^r\right). \end{aligned} \quad (7.61)$$

Now, for  $r \geq r^*$ , using (7.25), (7.26), (7.24), (7.60), (7.8), and the fact that  $x_j^* < 1$  for all  $j \in \underline{\mathcal{J}}_i$ , we have

$$\begin{aligned} &\left( \lambda_i^r - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j^r - \tilde{\epsilon}_i - \tilde{\epsilon}_i \sum_{j \in \underline{\mathcal{J}}_i} x_j^* - \epsilon_i \sum_{j \in \underline{\mathcal{J}}_i} (\mu_j^r + \tilde{\epsilon}_i) \right) t_i^r \\ &\geq \left( \frac{1}{2} \left( \lambda_i - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j \right) - (|\underline{\mathcal{J}}_i| + 1)\tilde{\epsilon}_i - \epsilon_i \sum_{j \in \underline{\mathcal{J}}_i} (\mu_j + 2\tilde{\epsilon}_i) \right) t_i^r \\ &\geq \left( \frac{1}{2} \left( \lambda_i - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j \right) - (|\underline{\mathcal{J}}_i| + 2)\tilde{\epsilon}_i \right) t_i^r > L_i^r. \end{aligned} \quad (7.62)$$

Hence, the first probability in the last expression of (7.61) is zero. From Lemma 6.7 (since  $(x_j^* + \epsilon_i)t_i^r > 2/\tilde{\epsilon}_i$ , for all  $j \in \underline{\mathcal{J}}_i$ , by (7.9), (7.12), and (7.35)), we have for all  $j \in \underline{\mathcal{J}}_i$ ,  $r \geq r^*$ ,

$$\begin{aligned} &\mathbf{P}\left(S_j^r((x_j^* + \epsilon_i)t_i^r) > (\mu_j^r + \tilde{\epsilon}_i)(x_j^* + \epsilon_i)t_i^r\right) \\ &\leq \exp\left(-(\mu_j(x_j^* + \epsilon_i)t_i^r - 1)\Lambda_j^{s,*} \left(\frac{1}{\mu_j} \left(\frac{1}{1 + \frac{\tilde{\epsilon}_i}{3\mu_j}}\right)\right)\right) \\ &\leq K_{11} \exp(-K'_{11} t_i^r), \end{aligned} \quad (7.63)$$

where  $K_{11} = \max\{\exp(\Lambda_j^{s,*}((\mu_j(1+\tilde{\epsilon}_i/3\mu_j))^{-1}) : j \in \underline{\mathcal{J}}_i\} > 0$ , and  $K'_{11} = \min\{\mu_j(x_j^* + \epsilon_i)\Lambda_j^{s,*}((\mu_j(1+\tilde{\epsilon}_i/3\mu_j))^{-1}) : j \in \underline{\mathcal{J}}_i\} > 0$ . Using (6.27) in conjunction with Lemma 6.7, we have for  $r \geq r^*$ ,  $0 < l_0 \in \mathcal{O}_0$ ,

$$\begin{aligned} \mathbf{P}(A_i^r(t_i^r) < (\lambda_i^r - \tilde{\epsilon}_i)t_i^r) &\leq \exp\left(-(\lambda_i - 2\tilde{\epsilon}_i)t_i^r \Lambda_i^{a,*}\left(\frac{1}{\lambda_i}\left(1 + \frac{\tilde{\epsilon}_i}{2\lambda_i}\right)\right)\right) \\ &\quad + \exp\left(-\frac{l_0\tilde{\epsilon}_i t_i^r}{2\lambda_i}\right) \exp(\Lambda_i^a(l_0)) \\ &\leq K_{12} \exp(-K'_{12}t_i^r), \end{aligned} \tag{7.64}$$

where  $K_{12} = \max\{1, \exp(\Lambda_i^a(l_0))\} > 0$ ,  $K'_{12} = \min\{(\lambda_i - 2\tilde{\epsilon}_i)\Lambda_i^{a,*}(\lambda_i^{-1}(1 + \tilde{\epsilon}_i/2\lambda_i)), l_0\tilde{\epsilon}_i/2\lambda_i\} > 0$ .

Then, combining (7.61)–(7.64), we have for each  $k \in \underline{\mathcal{K}}_i$ , for all  $r \geq r^*$ ,

$$\begin{aligned} \mathbf{P}(I_k^r(\tau_{i,0}^r) \geq t_i^r) &\leq |\underline{\mathcal{J}}_i| K_{11} \exp(-K'_{11}t_i^r) + K_{12} \exp(-K'_{12}t_i^r) \\ &\quad + \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P}(T_j^r(t_i^r) > (x_j^* + \epsilon_i)t_i^r). \end{aligned} \tag{7.65}$$

By assumption (iii) of Lemma 7.3, the definitions of  $t_i^r$  and  $L_i^r$ , and the fact that  $t_i^r \geq L_0^r$  (by (7.9) and (7.12)) for  $r \geq r^*$ , it follows that for all  $k \in \underline{\mathcal{K}}_i$ ,  $r \geq r^*$ , and  $t > 0$  satisfying  $r^2t \geq M^r$ , we have

$$\mathbf{P}(I_k^r(\tau_{i,0}^r) \geq t_i^r) \leq p_{2,i}(r^2t) \left( C_{2,i}^{(1)} \exp(-C_{2,i}^{(2)}L_0^r) + C_{2,i}^{(3)} \exp(-C_{2,i}^{(4)}r^2t) \right), \tag{7.66}$$

where  $p_{2,i}$  is a polynomial (of degree at most  $i$ ) with non-negative coefficients, and  $C_{2,i}^{(m)} > 0$ , for  $m = 1, 2, 3, 4$ , and where the constants and the polynomial do not depend on  $t$  or  $r$ .

*Proof of (I.3).* Since  $\underline{\mathcal{K}}_i = \emptyset$  if  $i$  is not a transition class, (I.3) holds trivially in this case. So suppose  $i$  is a transition class. Note that under the threshold policy, since class  $i$  (being above server  $k$ ) is the lowest priority class for server  $k \in \underline{\mathcal{K}}_i$ ,  $I_k^r$  can increase only when  $Q_i^r \leq |\underline{\mathcal{K}}_i| = |\underline{\mathcal{J}}_i|$ . The bound of  $|\underline{\mathcal{J}}_i|$  occurs here because there may be a class  $i$  job in service or in suspension at each of the other  $|\underline{\mathcal{J}}_i|$  servers ( $|\underline{\mathcal{J}}_i| - 1$  servers below and one server above  $i$ ) that can serve class  $i$ . In particular, if server  $k \in \underline{\mathcal{K}}_i$  incurs some idletime in  $[\tau_{i,0}^r, r^2t]$ , i.e.,  $\tau_{i,0}^r < r^2t$  and  $I_k^r(r^2t) - I_k^r(\tau_{i,0}^r) > 0$ , then  $R_i^r(s) \leq -L_i^r + |\underline{\mathcal{J}}_i|$  for some  $s \in [\tau_{i,0}^r, r^2t]$ . Thus, for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ , and  $k \in \underline{\mathcal{K}}_i$ ,

$$\begin{aligned} &\mathbf{P}(I_k^r(r^2t) - I_k^r(\tau_{i,0}^r) > 0, \tau_{i,0}^r < r^2t) \\ &\leq \mathbf{P}\left(\inf_{\tau_{i,0}^r \leq s \leq r^2t} R_i^r(s) \leq -L_i^r + |\underline{\mathcal{J}}_i|\right) \\ &\leq p_{3,i}(r^2t) \left( C_{3,i}^{(1)} \exp(-C_{3,i}^{(2)}L_0^r) + C_{3,i}^{(3)} \exp(-C_{3,i}^{(4)}r^2t) \right), \end{aligned} \tag{7.67}$$

where  $p_{3,i}$  is a polynomial (of degree at most  $i + 1$ ) with non-negative coefficients, and  $C_{3,i}^{(m)} > 0$ , for  $m = 1, 2, 3, 4$ , by the validity of (I.1) proved above, where the constants and the polynomial do not depend on  $t$  or  $r$ .  $\square$

#### 7.4 Estimates on Allocations for Activities Immediately Below Buffers – Proof of Lemma 7.4

**Proof of Lemma 7.4.** Fix  $i \in \mathcal{I} \setminus \{i^*\}$ . Suppose that  $i$  is a transition buffer, so that  $\underline{\mathcal{J}}_i \neq \emptyset$ , and assume that (I) and (II) hold with  $i'$  in place of  $i$ , for all  $i' < i$  (for  $i = 1$  this is a vacuous assumption). In the following, recall (cf. Section 1.1) that a sum over an empty set is defined to equal zero. In particular, the results below hold even if  $\underline{\mathcal{L}}_k = \emptyset$ .

*Proof of (i).* For  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ ,  $j \in \underline{\mathcal{J}}_i$ ,  $k = k(j)$ , and  $n \geq 1$ , we have

$$\begin{aligned}
& \mathbf{P}(T_{i,j}^{r,n}(s_i^r) \leq (x_j^* - \epsilon_i)s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2t) \\
&= \mathbf{P}\left(s_i^r - \sum_{i' \in \underline{\mathcal{I}}_k} T_{i,a(i')}^{r,n}(s_i^r) \leq (x_j^* - \epsilon_i)s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2t\right) \\
&= \mathbf{P}\left(\sum_{i' \in \underline{\mathcal{I}}_k} T_{i,a(i')}^{r,n}(s_i^r) \geq \sum_{i' \in \underline{\mathcal{I}}_k} x_{a(i')}^* s_i^r + \epsilon_i s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2t\right) \\
&\leq \sum_{i' \in \underline{\mathcal{I}}_k} \mathbf{P}\left(T_{i,a(i')}^{r,n}(s_i^r) \geq (x_{a(i')}^* + \epsilon_{i'})s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2t\right). \tag{7.68}
\end{aligned}$$

The first equality holds since server  $k$ ,  $k \in \underline{\mathcal{K}}_i$ , does not incur any idle time in the  $n^{\text{th}}$  (up) excursion interval for  $R_i^r$  which is of length  $\beta_{i,n}^r$ , and the second equality holds since  $\sum_{i' \in \underline{\mathcal{I}}_k} x_{a(i')}^* + x_j^* = 1$ . For the last inequality, we have used the fact that  $\epsilon_i \geq \mathbf{I}\epsilon_{i'} \geq |\underline{\mathcal{I}}_k|\epsilon_{i'}$ , by (6.7) (since  $\gamma_l \leq 1$  for all  $l \in \mathcal{I}$ , and the fact that  $i' \in \underline{\mathcal{I}}_k$  satisfies  $i' < i$  by the ordering assumed for the buffer numbering). Hence, for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ , for each  $j \in \underline{\mathcal{J}}_i$ , (i) of Lemma 7.3 holds, since (II.1) was assumed to hold with  $i'$  in place of  $i$  for all  $i' < i$  and  $|\underline{\mathcal{I}}_k| < \infty$ .

*Proof of (ii).* For  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ ,  $j \in \underline{\mathcal{J}}_i$ ,  $k = k(j)$ , and  $n \geq 1$ , we have

$$\begin{aligned}
& \mathbf{P}(dT_{i,j}^{r,n}(d_s^r) \geq (x_j^* + \epsilon_i)d_s^r, d_{\tau_{i,2n-1}}^r \leq r^2t) \\
&= \mathbf{P}\left(d_s^r - \sum_{i' \in \underline{\mathcal{I}}_k} dT_{i,a(i')}^{r,n}(d_s^r) - (I_k^r(d_{\tau_{i,2n-1}}^r + d_s^r) - I_k^r(d_{\tau_{i,2n-1}}^r)) \geq (x_j^* + \epsilon_i)d_s^r, d_{\tau_{i,2n-1}}^r \leq r^2t\right) \\
&\leq \mathbf{P}\left(\sum_{i' \in \underline{\mathcal{I}}_k} dT_{i,a(i')}^{r,n}(d_s^r) \leq \sum_{i' \in \underline{\mathcal{I}}_k} x_{a(i')}^* d_s^r - \epsilon_i d_s^r, d_{\tau_{i,2n-1}}^r \leq r^2t\right) \\
&\leq \sum_{i' \in \underline{\mathcal{I}}_k} \mathbf{P}\left(dT_{i,a(i')}^{r,n}(d_s^r) \leq (x_{a(i')}^* - \epsilon_{i'})d_s^r, d_{\tau_{i,2n-1}}^r \leq r^2t\right). \tag{7.69}
\end{aligned}$$

In the above, the equality follows from (2.8); the first inequality holds since the idletime process,  $I_k^r$ , is non-decreasing and non-negative, and since  $\sum_{i' \in \underline{\mathcal{I}}_k} x_{a(i')}^* + x_j^* = 1$ . Note also for the last inequality that  $\epsilon_i \geq \mathbf{I}\epsilon_{i'} \geq |\underline{\mathcal{I}}_k|\epsilon_{i'}$ . Hence, for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ , for each  $j \in \underline{\mathcal{J}}_i$ , (ii) of Lemma 7.3 holds, since (II.2) was assumed to hold with  $i'$  in place of  $i$  for all  $i' < i$ , and  $|\underline{\mathcal{I}}_k| < \infty$ .

*Proof of (iii).* For  $j \in \underline{\mathcal{J}}_i$ , and  $k = k(j)$ , we have in a similar manner to that for the proof of (ii) above that for  $r \geq r^*$ ,

$$\begin{aligned}
& \mathbf{P}\left(T_j^r(t_i^r) \geq (x_j^* + \epsilon_i)t_i^r\right) \\
&\leq \mathbf{P}\left(t_i^r - \sum_{i' \in \underline{\mathcal{I}}_k} T_{a(i')}^r(t_i^r) - I_k^r(t_i^r) \geq (x_j^* + \epsilon_i)t_i^r\right) \\
&\leq \mathbf{P}\left(\sum_{i' \in \underline{\mathcal{I}}_k} T_{a(i')}^r(t_i^r) \leq \sum_{i' \in \underline{\mathcal{I}}_k} (x_{a(i')}^* - \epsilon_{i'})t_i^r\right) \\
&\leq \sum_{i' \in \underline{\mathcal{I}}_k} \mathbf{P}\left(T_{a(i')}^r(t_i^r) \leq (x_{a(i')}^* - \epsilon_{i'})t_i^r\right). \tag{7.70}
\end{aligned}$$

Hence, for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ , for each  $j \in \underline{\mathcal{J}}_i$ , (iii) of Lemma 7.3 holds, since (II.3) was assumed to hold with  $i'$  in place of  $i$ , for all  $i' < i$ .  $\square$

## 7.5 Estimates on Allocations for Activities Immediately Above Buffers – Proof of Lemma 7.5

**Proof of Lemma 7.5.** Fix  $i \in \mathcal{I} \setminus \{i^*\}$ , and let  $k = k(a(i))$ . Assume that (I) and (II) hold with  $i'$  in place of  $i$ , for all  $i' < i$ . Note that by the priorities assigned to buffers by server  $k$ , we have that

$$\sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' < i}} T_{i,a(i')}^{r,n}(s) = s, \quad \text{for } 0 \leq s \leq \beta_{i,n}^r, \quad (7.71)$$

since activities that have lower priority than activity  $a(i)$  will not be “on” and server  $k$  will not idle in the  $n^{\text{th}}$  up excursion interval for  $R_i^r$ . We then have for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , and  $n \geq 1$ ,

$$\begin{aligned} & \mathbf{P} \left( T_{i,a(i)}^{r,n}(s_i^r) \leq (\hat{x}_{i,i} - \epsilon_i) s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2 t \right) \\ &= \mathbf{P} \left( s_i^r - \sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' < i}} T_{i,a(i')}^{r,n}(s_i^r) \leq (\hat{x}_{i,i} - \epsilon_i) s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2 t \right) \\ &= \mathbf{P} \left( \sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' < i}} T_{i,a(i')}^{r,n}(s_i^r) \geq \sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' < i}} \hat{x}_{i,i'} s_i^r + \epsilon_i s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2 t \right) \\ &\leq \mathbf{P} \left( \sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' < i}} T_{i,a(i')}^{r,n}(s_i^r) \geq \sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' < i}} x_{a(i')}^* s_i^r + \epsilon_i s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2 t \right) \\ &\leq \sum_{\substack{i' \in \underline{\mathcal{I}}_k \\ i' < i}} \mathbf{P} \left( T_{i,a(i')}^{r,n}(s_i^r) \geq (x_{a(i')}^* + \epsilon_{i'}) s_i^r, s_i^r \leq \beta_{i,n}^r, \tau_{i,2n-1}^r \leq r^2 t \right). \end{aligned} \quad (7.72)$$

The first equality follows by (7.71) and the second uses (6.3). The second to last inequality holds since  $\hat{x}_{i,i'} \geq x_{a(i')}^*$  for all  $i' \in \underline{\mathcal{I}}_k$ ,  $i' < i$ , and the last inequality follows since  $\epsilon_i \geq \mathbf{I}\epsilon_{i'} \geq |\{i'' \in \underline{\mathcal{I}}_k : i'' < i\}| \epsilon_{i'}$  for  $i' < i$ . (Notice that if buffer  $i$  is the highest priority buffer for server  $k$ , i.e.,  $\{i' \in \underline{\mathcal{I}}_k : i' < i\} = \emptyset$ , then the first probability in (7.72) is zero since  $T_{i,a(i)}^{r,n}(s_i^r) = s_i^r$  for  $s_i^r \leq \beta_{i,n}^r$ , by (7.71).) Hence, for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , (iv) of Lemma 7.3 holds, since (II.1) was assumed to hold with  $i'$  in place of  $i$ , for all  $i' < i$ .  $\square$

## 7.6 Transition Between Layers in the Server-Buffer Tree — Proof of Lemma 7.6

We will show below that (II) holds for  $i \in \mathcal{I} \setminus \{i^*\}$  given that (I) and (II) hold with  $i'$  in place of  $i$ , for all  $i' < i$ .

**Proof of Lemma 7.6.** Fix  $i \in \mathcal{I} \setminus \{i^*\}$ . Assume that (I) and (II) hold with  $i'$  in place of  $i$ , for all  $i' < i$ . Then, from Lemmas 7.3–7.5, we have that (I) holds for  $i$  as well. By (6.5) and (6.7), we have, for any  $\iota \in \mathcal{I}$  satisfying  $\iota > i$ ,

$$\begin{aligned} \epsilon_i &< \min \left\{ \frac{\prod_{m=1}^{\mathbf{I}} \gamma_m}{\mathbf{I}^{\mathbf{I}}}, \frac{\mu_{a(i)}}{1024\lambda_\iota}, \frac{\mu_{a(i)}}{1024 \sum_{j \in \mathcal{J}_i} (\mu_j + \epsilon_\iota)}, \right. \\ &\quad \left. \frac{\lambda_i - \sum_{j \in \underline{\mathcal{I}}_\iota} x_j^* \mu_j}{18\lambda_\iota}, \frac{\mu_{a(i)}}{16(\lambda_\iota - \sum_{j \in \underline{\mathcal{I}}_\iota} x_j^* \mu_j)} \right\}, \end{aligned} \quad (7.73)$$

since  $\epsilon_l < \hat{\epsilon}$ , for all  $l \in \mathcal{I}$ . To validate the denominator in the third term in (7.73), we note that  $1024 \sum_{j \in \mathcal{J}_i} (\mu_j + \epsilon_l) \leq 1024(\mu_{\text{sum}} + |\mathcal{J}_i| \hat{\epsilon}) \leq 2048\mu_{\text{sum}}$ . For the fifth term, we note that  $0 < \lambda_l - \sum_{j \in \underline{\mathcal{J}}_l} x_j^* \mu_j \leq \lambda_{\text{max}}$  for all  $l \in \mathcal{I}$ .

*Proof of (II.1).* For  $r \geq 1$ ,  $n \geq 1$ ,  $s \geq 0$ ,  $\iota \in \mathcal{I}$  such that  $\iota > i$ ,  $j \in \mathcal{J}_i$ ,  $k = k(j)$ , on  $\{\tau_{\iota, 2n-1}^r < \infty\}$ , define

$$A_{\iota, i}^{r, n}(s) = A_{\iota}^r(\tau_{\iota, 2n-1}^r + s) - A_{\iota}^r(\tau_{\iota, 2n-1}^r), \quad (7.74)$$

$$S_{\iota, j}^{r, n}(s) = S_j^r(T_j^r(\tau_{\iota, 2n-1}^r) + s) - S_j^r(T_j^r(\tau_{\iota, 2n-1}^r)), \quad (7.75)$$

$$\check{S}_{\iota, j}^{r, n}(s) = \sup\{m \geq 0 : \eta_j^r(S_j^r(T_j^r(\tau_{\iota, 2n-1}^r)) + m) - \eta_j^r(S_j^r(T_j^r(\tau_{\iota, 2n-1}^r))) \leq s\}, \quad (7.76)$$

$$\check{A}_{\iota, i}^{r, n}(s) = \sup\{m \geq 0 : \xi_i^r(A_{\iota}^r(\tau_{\iota, 2n-1}^r) + m + 1) - \xi_i^r(A_{\iota}^r(\tau_{\iota, 2n-1}^r) + 1) \leq s\}, \quad (7.77)$$

$$I_{\iota, k}^{r, n}(s) = I_k^r(\tau_{\iota, 2n-1}^r + s) - I_k^r(\tau_{\iota, 2n-1}^r), \quad (7.78)$$

and for concreteness on  $\{\tau_{\iota, 2n-1}^r = \infty\}$ , we define  $A_{\iota, i}^{r, n}$ ,  $S_{\iota, j}^{r, n}$ ,  $\check{S}_{\iota, j}^{r, n}$ ,  $\check{A}_{\iota, i}^{r, n}$ ,  $I_{\iota, k}^{r, n}$ , to be identically zero. Then we have, for all  $s \geq 0$ ,

$$\check{A}_{\iota, i}^{r, n}(s) \geq A_{\iota, i}^{r, n}(s) - 1, \quad S_{\iota, j}^{r, n}(s) \geq \check{S}_{\iota, j}^{r, n}(s). \quad (7.79)$$

Now, for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,  $\iota \in \mathcal{I}$  satisfying  $\iota > i$ , and  $n \geq 1$ ,

$$\begin{aligned} & \mathbf{P} \left( T_{\iota, a(i)}^{r, n}(s_{\iota}^r) \geq (x_{a(i)}^* + \epsilon_i) s_{\iota}^r, s_{\iota}^r \leq \beta_{\iota, n}^r, \tau_{\iota, 2n-1}^r \leq r^2 t \right) \\ & \leq \mathbf{P} \left( S_{\iota, a(i)}^{r, n}(T_{\iota, a(i)}^{r, n}(s_{\iota}^r)) \geq S_{\iota, a(i)}^{r, n}((x_{a(i)}^* + \epsilon_i) s_{\iota}^r), s_{\iota}^r \leq \beta_{\iota, n}^r, \tau_{\iota, 2n-1}^r \leq r^2 t \right) \\ & \leq \mathbf{P} \left( S_{\iota, a(i)}^{r, n}((x_{a(i)}^* + \epsilon_i) s_{\iota}^r) \leq Q_{\iota}^r(\tau_{\iota, 2n-1}^r) + A_{\iota, i}^{r, n}(s_{\iota}^r) \right. \\ & \quad \left. - \sum_{j \in \underline{\mathcal{J}}_i} S_{\iota, j}^{r, n}(T_{\iota, j}^{r, n}(s_{\iota}^r)), s_{\iota}^r \leq \beta_{\iota, n}^r, \tau_{\iota, 2n-1}^r \leq r^2 t \right) \\ & \leq \mathbf{P} \left( S_{\iota, a(i)}^{r, n}((x_{a(i)}^* + \epsilon_i) s_{\iota}^r) < (\mu_{a(i)}^r - \epsilon_{2, i})(x_{a(i)}^* + \epsilon_i) s_{\iota}^r, \tau_{\iota, 2n-1}^r \leq r^2 t \right) \\ & \quad + \mathbf{P} \left( Q_{\iota}^r(\tau_{\iota, 2n-1}^r) + A_{\iota, i}^{r, n}(s_{\iota}^r) - \sum_{j \in \underline{\mathcal{J}}_i} S_{\iota, j}^{r, n}(T_{\iota, j}^{r, n}(s_{\iota}^r)) \geq (x_{a(i)}^* + \epsilon_i)(\mu_{a(i)}^r - \epsilon_{2, i}) s_{\iota}^r, \right. \\ & \quad \left. s_{\iota}^r \leq \beta_{\iota, n}^r, \tau_{\iota, 2n-1}^r \leq r^2 t \right). \end{aligned} \quad (7.80)$$

For the second inequality, we have used the fact that the number of class  $i$  jobs processed by activity  $a(i)$ , between time  $\tau_{\iota, 2n-1}^r$  and time  $\tau_{\iota, 2n-1}^r + s_{\iota}^r$ , namely,  $S_{\iota, a(i)}^{r, n}(T_{\iota, a(i)}^{r, n}(s_{\iota}^r))$ , is less than or equal to the number of class  $i$  jobs present at time  $\tau_{\iota, 2n-1}^r$ , namely,  $Q_{\iota}^r(\tau_{\iota, 2n-1}^r)$ , plus the number of arrivals to class  $i$  between time  $\tau_{\iota, 2n-1}^r$  and time  $\tau_{\iota, 2n-1}^r + s_{\iota}^r$ , namely  $A_{\iota, i}^{r, n}(s_{\iota}^r)$ , less the number of class  $i$  jobs processed during this time interval by the activities indexed by  $\underline{\mathcal{J}}_i$ , namely,  $\sum_{j \in \underline{\mathcal{J}}_i} S_{\iota, j}^{r, n}(T_{\iota, j}^{r, n}(s_{\iota}^r))$  (cf. (2.9)).

Letting  $v_{\iota, a(i)}^{r, n} = v_{a(i)}^r(S_{a(i)}^r(T_{a(i)}^r(\tau_{\iota, 2n-1}^r)) + 1)$  on  $\{\tau_{\iota, 2n-1}^r < \infty\}$  and  $v_{\iota, a(i)}^{r, n} = 0$  on  $\{\tau_{\iota, 2n-1}^r = \infty\}$ , by Lemmas 6.5–6.7 together with (6.27) and (7.33), in a similar manner to (7.55)–(7.56), for  $r \geq r^*$ ,  $n \geq 1$ , and  $t > 0$  satisfying  $r^2 t \geq M^r$ , we have

$$\begin{aligned} & \mathbf{P} \left( S_{\iota, a(i)}^{r, n}((x_{a(i)}^* + \epsilon_i) s_{\iota}^r) < (\mu_{a(i)}^r - \epsilon_{2, i})(x_{a(i)}^* + \epsilon_i) s_{\iota}^r, \tau_{\iota, 2n-1}^r \leq r^2 t \right) \\ & \leq K_{13} \exp(-K'_{13} s_{\iota}^r) \\ & \quad + \left( \lfloor (\mu_{a(i)}^r + \epsilon_i) r^2 t \rfloor + 1 \right) K_{14} \exp(-K'_{14} s_{\iota}^r) + K_{15} \exp(-K'_{15} r^2 t), \end{aligned} \quad (7.81)$$

where  $K_{13} = 1$ ,  $K'_{13} = (\mu_{a(i)} - 2\epsilon_{2,i})(x_{a(i)}^* + \epsilon_i)\Lambda_{a(i)}^{s,*}(\mu_{a(i)}^{-1}(1 + \epsilon_{2,i}/2\mu_{a(i)})) > 0$ ,  $K_{14} = \exp(\Lambda_{a(i)}^s(l_0)) > 0$ ,  $K'_{14} = (l_0\epsilon_{2,i}/2\mu_{a(i)})(x_{a(i)}^* + \epsilon_i) > 0$ ,  $0 < l_0 \in \mathcal{O}$ ,  $K_{15} = \exp(\Lambda_{a(i)}^{s,*}((\mu_{a(i)}(1 + \epsilon_i/3\mu_{a(i)}))^{-1})) > 0$ , and  $K'_{15} = \mu_{a(i)}\Lambda_{a(i)}^{s,*}((\mu_{a(i)}(1 + \epsilon_i/3\mu_{a(i)}))^{-1}) > 0$ , which are all independent of  $t$ ,  $n$ , and  $r$ .

For  $r \geq 1$ ,  $n \geq 1$ , and  $\iota > i$ , let

$$\begin{aligned} \Upsilon_{\iota,i}^{r,n} &= \left\{ Q_i^r(\tau_{\iota,2n-1}^r) \leq (\mu_{a(i)}^r \epsilon_i / 16) s_\iota^r; A_{\iota,i}^{r,n}(s_\iota^r) \leq (\lambda_i^r + \mu_{a(i)}^r \epsilon_i / 16) s_\iota^r + 1; \right. \\ &\quad S_{\iota,j}^{r,n}((x_j^* - \epsilon_{1,i}) s_\iota^r) \geq \left( \mu_j^r - \frac{\mu_{a(i)}^r \epsilon_i}{16|\underline{\mathcal{J}}_i|} \right) (x_j^* - \epsilon_{1,i}) s_\iota^r, j \in \underline{\mathcal{J}}_i; \\ &\quad \left. T_{\iota,j}^{r,n}(s_\iota^r) > (x_j^* - \epsilon_{1,i}) s_\iota^r, j \in \underline{\mathcal{J}}_i; \tau_{\iota,2n-1}^r \leq r^2 t \right\}, \end{aligned} \quad (7.82)$$

where  $\epsilon_{1,i}$  is defined in (7.32), and if  $\underline{\mathcal{J}}_i = \emptyset$ , we omit the terms involving  $j \in \underline{\mathcal{J}}_i$  from the definition of  $\Upsilon_{\iota,i}^{r,n}$ . Then, for the last probability in (7.80), we have

$$\begin{aligned} &\mathbf{P} \left( Q_i^r(\tau_{\iota,2n-1}^r) + A_{\iota,i}^{r,n}(s_\iota^r) - \sum_{j \in \underline{\mathcal{J}}_i} S_{\iota,j}^{r,n}(T_{\iota,j}^{r,n}(s_\iota^r)) \right. \\ &\quad \left. \geq (x_{a(i)}^* + \epsilon_i)(\mu_{a(i)}^r - \epsilon_{2,i}) s_\iota^r, s_\iota^r \leq \beta_{\iota,n}^r, \tau_{\iota,2n-1}^r \leq r^2 t \right) \\ &\leq \mathbf{P} \left( Q_i^r(\tau_{\iota,2n-1}^r) + A_{\iota,i}^{r,n}(s_\iota^r) - \sum_{j \in \underline{\mathcal{J}}_i} S_{\iota,j}^{r,n}(T_{\iota,j}^{r,n}(s_\iota^r)) \right. \\ &\quad \left. \geq (x_{a(i)}^* + \epsilon_i)(\mu_{a(i)}^r - \epsilon_{2,i}) s_\iota^r, \Upsilon_{\iota,i}^{r,n}, \tau_{\iota,2n-1}^r \leq r^2 t \right) \\ &+ \mathbf{P} \left( Q_i^r(\tau_{\iota,2n-1}^r) > (\mu_{a(i)}^r \epsilon_i / 16) s_\iota^r, \tau_{\iota,2n-1}^r \leq r^2 t \right) \\ &+ \mathbf{P} \left( A_{\iota,i}^{r,n}(s_\iota^r) > \left( \lambda_i^r + \frac{\mu_{a(i)}^r \epsilon_i}{16} \right) s_\iota^r + 1, \tau_{\iota,2n-1}^r \leq r^2 t \right) \\ &+ \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P} \left( S_{\iota,j}^{r,n}((x_j^* - \epsilon_{1,i}) s_\iota^r) < \left( \mu_j^r - \frac{\mu_{a(i)}^r \epsilon_i}{16|\underline{\mathcal{J}}_i|} \right) (x_j^* - \epsilon_{1,i}) s_\iota^r, \tau_{\iota,2n-1}^r \leq r^2 t \right) \\ &+ \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P} \left( T_{\iota,j}^{r,n}(s_\iota^r) \leq (x_j^* - \epsilon_{1,i}) s_\iota^r, s_\iota^r \leq \beta_{\iota,n}^r, \tau_{\iota,2n-1}^r \leq r^2 t \right). \end{aligned} \quad (7.83)$$

For the first term in the right side of the inequality in (7.83) we have on  $\Upsilon_{l,i}^{r,n}$  that for  $r \geq r^*$ ,

$$\begin{aligned}
& Q_i^r(\tau_{l,2n-1}^r) + A_{l,i}^{r,n}(s_l^r) - \sum_{j \in \underline{\mathcal{J}}_i} S_{l,j}^{r,n}(T_{l,j}^{r,n}(s_l^r)) \\
& \leq \frac{\mu_{a(i)}^r \epsilon_i}{16} s_l^r + \lambda_i^r s_l^r + \frac{\mu_{a(i)}^r \epsilon_i}{16} s_l^r + 1 - \sum_{j \in \underline{\mathcal{J}}_i} \left( \mu_j^r - \frac{\mu_{a(i)}^r \epsilon_i}{16 |\underline{\mathcal{J}}_i|} \right) (x_j^* - \epsilon_{1,i}) s_l^r \\
& = \frac{\mu_{a(i)}^r \epsilon_i}{8} s_l^r + \lambda_i^r s_l^r - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j^r s_l^r + 1 \\
& \quad + \sum_{j \in \underline{\mathcal{J}}_i} \frac{\mu_{a(i)}^r \epsilon_i}{16 |\underline{\mathcal{J}}_i|} x_j^* s_l^r + \sum_{j \in \underline{\mathcal{J}}_i} \left( \mu_j^r - \frac{\mu_{a(i)}^r \epsilon_i}{16 |\underline{\mathcal{J}}_i|} \right) \epsilon_{1,i} s_l^r \\
& \leq \frac{\mu_{a(i)}^r \epsilon_i}{8} s_l^r + \lambda_i^r s_l^r - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j^r s_l^r + 1 + \frac{\mu_{a(i)}^r \epsilon_i}{16} s_l^r + \sum_{j \in \underline{\mathcal{J}}_i} \mu_j^r \epsilon_{1,i} s_l^r \\
& \leq \frac{\mu_{a(i)}^r \epsilon_i}{4} s_l^r + \lambda_i^r s_l^r - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j^r s_l^r + 1 \\
& \leq \left( \lambda_i^r - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j^r + \frac{\mu_{a(i)}^r \epsilon_i}{2} \right) s_l^r \\
& \leq \left( x_{a(i)}^* \mu_{a(i)}^r + \frac{\mu_{a(i)}^r \epsilon_i}{32} + \frac{\mu_{a(i)}^r \epsilon_i}{2} \right) s_l^r \\
& < (x_{a(i)}^* + \epsilon_i) (\mu_{a(i)}^r - \epsilon_{2,i}) s_l^r, \tag{7.84}
\end{aligned}$$

where in the second inequality we have used the fact that  $\sum_{j \in \underline{\mathcal{J}}_i} x_j^* \leq 1$ . For the third inequality we have, using (7.32) together with (7.15), if  $\underline{\mathcal{J}}_i \neq \emptyset$ ,

$$\sum_{j \in \underline{\mathcal{J}}_i} \mu_j^r \epsilon_{1,i} = \frac{\sum_{j \in \underline{\mathcal{J}}_i} \mu_j^r}{\sum_{j \in \underline{\mathcal{J}}_i} (\mu_j + \epsilon_i)} \frac{\epsilon_i}{16} \left( \mu_{a(i)} - \frac{\mu_{a(i)} \epsilon_i}{16} \right) \leq \frac{\mu_{a(i)}^r \epsilon_i}{16}. \tag{7.85}$$

For the fourth inequality we have used (7.27), together with (7.15). The fifth inequality follows by (7.28). The final inequality follows by (7.33). Hence, the first probability in the second expression of (7.83) is zero, for all  $n \geq 1$ ,  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ .

For the second term in the right side of the inequality in (7.83), using the result established at the beginning of this proof that (I.1) holds for  $i$ , we have for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,

$$\begin{aligned}
& \mathbf{P} \left( Q_i^r(\tau_{l,2n-1}^r) > (\mu_{a(i)}^r \epsilon_i / 16) s_l^r, \tau_{l,2n-1}^r \leq r^2 t \right) \\
& \leq \mathbf{P} \left( \sup_{\tau_{i,0}^r \leq s \leq r^2 t} R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \right) \\
& \leq p_{1,i}(r^2 t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-C_{1,i}^{(4)} r^2 t) \right), \tag{7.86}
\end{aligned}$$

where  $p_{1,i}$  is a polynomial (of degree at most  $i + 1$ ) with non-negative coefficients, and  $C_{1,i}^{(m)} > 0$ , for  $m = 1, 2, 3, 4$ , and where the polynomial and constants are independent of  $t$  and  $r$ . In (7.86),

the first inequality holds since  $s_\iota^r$  (for  $\iota > i$ ) is considerably larger than  $L_i^r$ . Specifically,

$$\begin{aligned}
\frac{\mu_{a(i)}^r \epsilon_i}{16} s_\iota^r &= \frac{\mu_{a(i)}^r \epsilon_i}{16} \cdot \frac{L_\iota^r - (|\underline{\mathcal{J}}_\iota| + 2)}{\lambda_\iota^r + \epsilon_\iota} \\
&\geq \frac{\mu_{a(i)}^r}{\mu_{a(i)}} \cdot \frac{\lambda_\iota}{\lambda_\iota^r + \epsilon_\iota} \left[ \frac{\mu_{a(i)}}{16\lambda_\iota} \cdot \frac{L_i^r}{\epsilon_i} - \frac{\mu_{a(i)}}{16\lambda_\iota} \epsilon_i (|\underline{\mathcal{J}}_\iota| + 2) \right] \\
&\geq \frac{1}{4} \left[ \frac{\mu_{a(i)}}{16\lambda_\iota} \cdot \frac{L_i^r}{\epsilon_i} - \frac{\mu_{a(i)}}{16\lambda_\iota} \epsilon_i (|\underline{\mathcal{J}}_\iota| + 2) \right] \\
&\geq 2L_i^r + \frac{1}{4} \left[ 8L_i^r - \frac{\mu_{a(i)}}{16\lambda_\iota} \epsilon_i (|\underline{\mathcal{J}}_\iota| + 2) \right] \\
&\geq 2L_i^r,
\end{aligned} \tag{7.87}$$

for  $r \geq r^*$ . In the first inequality of (7.87), we have used (6.4), and the fact that  $\epsilon_\iota \leq 1$  for all  $\iota \in \mathcal{I}$ . In the second inequality, we have used (7.14) and (7.27), together with the fact that  $\epsilon_\iota < \lambda_\iota/2$ , to obtain  $\lambda_\iota/(\lambda_\iota^r + \epsilon_\iota) > 1/2$ . The third inequality follows since  $\epsilon_i < \hat{\epsilon} < \mu_{a(i)}/256\lambda_\iota$  (cf. (7.73)), and the fourth inequality follows from (7.29).

For the third term in the right side of the inequality in (7.83), using (7.15) and (7.31) to obtain  $\mu_{a(i)}^r \epsilon_i/16 \geq \check{\epsilon}_i$ , and using the fact that  $s_\iota^r > 2/\check{\epsilon}_i$  for all  $\iota > i$ , we can proceed in a similar manner to that in (7.53), using Lemmas 6.5–6.7, to obtain for  $r \geq r^*$ ,  $n \geq 1$  and  $t > 0$  satisfying  $r^2 t \geq M^r$  that

$$\mathbf{P} \left( A_{\iota,i}^{r,n}(s_\iota^r) > \left( \lambda_i^r + \frac{\mu_{a(i)}^r \epsilon_i}{16} \right) s_\iota^r + 1, \tau_{\iota,2n-1}^r \leq r^2 t \right) \leq K_{16} \exp(-K'_{16} s_\iota^r), \tag{7.88}$$

where  $K_{16} = \exp(\Lambda_i^{a,*}((\lambda_i(1 + \check{\epsilon}_i/3\lambda_i))^{-1})) > 0$  and  $K'_{16} = \lambda_i \Lambda_i^{a,*}((\lambda_i(1 + \check{\epsilon}_i/3\lambda_i))^{-1}) > 0$  do not depend on  $t$ ,  $n$ , or  $r$ .

Similarly, for the fourth term in the right side of the inequality in (7.83), we have, by a similar argument to that for (7.55)–(7.56), that for  $r \geq r^*$ ,  $n \geq 1$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , and  $j \in \underline{\mathcal{J}}_i \neq \emptyset$ ,

$$\begin{aligned}
&\mathbf{P} \left( S_{\iota,j}^{r,n}((x_j^* - \epsilon_{1,i})s_\iota^r) < \left( \mu_j^r - \frac{\mu_{a(i)}^r \epsilon_i}{16|\underline{\mathcal{J}}_i|} \right) (x_j^* - \epsilon_{1,i})s_\iota^r, \tau_{\iota,2n-1}^r \leq r^2 t \right) \\
&\leq K_{17} \exp(-K'_{17} s_\iota^r) \\
&\quad + (|\mu_j^r + \epsilon_{1,i}|r^2 t + 1) K_{18} \exp(-K'_{18} s_\iota^r) + K_{19} \exp(-K'_{19} r^2 t),
\end{aligned} \tag{7.89}$$

where  $K_{17} = 1$ ,  $K'_{17} = \min\{(\mu_j - 2\check{\epsilon}_i)(x_j^* - \epsilon_{1,i})\Lambda_j^{s,*}(\mu_j^{-1}(1 + \check{\epsilon}_i/2\mu_j)) : j \in \underline{\mathcal{J}}_i\} > 0$ ,  $K_{18} = \max\{\exp(\Lambda_j^s(l_0)) > 0 : j \in \underline{\mathcal{J}}_i\}$ ,  $K'_{18} = \min\{(l_0 \check{\epsilon}_i/2\mu_j)(x_j^* - \epsilon_{1,i}) : j \in \underline{\mathcal{J}}_i\} > 0$ ,  $0 < l_0 \in \mathcal{O}_0$ ,  $K_{19} = \max\{\exp(\Lambda_j^{s,*}((\mu_j(1 + \epsilon_{1,i}/3\mu_j))^{-1})) : j \in \underline{\mathcal{J}}_i\} > 0$ , and  $K'_{19} = \min\{\mu_j \Lambda_j^{s,*}((\mu_j(1 + \epsilon_{1,i}/3\mu_j))^{-1}) : j \in \underline{\mathcal{J}}_i\} > 0$ , which are all independent of  $t$ ,  $n$ , and  $r$ .

For the fifth term in the right side of the inequality in (7.83), we have for  $r \geq r^*$ ,  $n \geq 1$ ,  $t > 0$



satisfying  $r^2t \geq M^r$ ,  $j \in \underline{\mathcal{J}}_i \neq \emptyset$ , and  $k = k(j)$ ,

$$\begin{aligned}
& \mathbf{P}\left(T_{l,j}^{r,n}(s_l^r) \leq (x_j^* - \epsilon_{1,i})s_l^r, s_l^r \leq \beta_{l,n}^r, \tau_{l,2n-1}^r \leq r^2t\right) \\
&= \mathbf{P}\left(s_l^r - \sum_{i' \in \underline{\mathcal{I}}_k} T_{l,a(i')}^{r,n}(s_l^r) - I_{l,k}^{r,n}(s_l^r) \leq (x_j^* - \epsilon_{1,i})s_l^r, s_l^r \leq \beta_{l,n}^r, \tau_{l,2n-1}^r \leq r^2t\right) \\
&\leq \mathbf{P}\left(\sum_{i' \in \underline{\mathcal{I}}_k} T_{l,a(i')}^{r,n}(s_l^r) + I_{l,k}^{r,n}(s_l^r) \geq \sum_{i' \in \underline{\mathcal{I}}_k} x_{a(i')}^* s_l^r + \epsilon_{1,i} s_l^r, s_l^r \leq \beta_{l,n}^r, \tau_{l,2n-1}^r \leq r^2t\right) \\
&\leq \sum_{i' \in \underline{\mathcal{I}}_k} \mathbf{P}\left(T_{l,a(i')}^{r,n}(s_l^r) \geq (x_{a(i')}^* + \epsilon_{i'})s_l^r, s_l^r \leq \beta_{l,n}^r, \tau_{l,2n-1}^r \leq r^2t\right) \\
&\quad + \mathbf{P}\left(I_{l,k}^{r,n}(s_l^r) \geq \frac{\epsilon_{1,i}}{\mathbf{I}} s_l^r, \tau_{l,2n-1}^r \leq r^2t\right). \tag{7.90}
\end{aligned}$$

For the first inequality in (7.90) we use the fact that  $\sum_{i' \in \underline{\mathcal{I}}_k} x_{a(i')}^* + x_j^* = 1$ . The last inequality in (7.90) follows from the fact that when  $\underline{\mathcal{J}}_i \neq \emptyset$ , for all  $i' \in \underline{\mathcal{I}}_k$ ,

$$\begin{aligned}
\epsilon_{1,i} &= \frac{1}{\sum_{j \in \underline{\mathcal{J}}_i} (\mu_j + \epsilon_i)} \frac{\epsilon_i}{16} \left( \mu_{a(i)} - \frac{\mu_{a(i)} \epsilon_i}{16} \right) \\
&\geq \frac{\mu_{a(i)} \epsilon_i}{64 \sum_{j \in \underline{\mathcal{J}}_i} \mu_j} \geq \gamma_i \epsilon_i \geq \mathbf{I} \epsilon_{i'}, \tag{7.91}
\end{aligned}$$

by (6.7) (since  $i' < i$ ), and where  $\mathbf{I} \geq (|\underline{\mathcal{I}}_k| + 1)$ . The first inequality in (7.91) holds since  $\epsilon_i \leq \min\{\mu_{\min}, 1\}$ , by (6.5), and the second inequality follows by the definition of  $\gamma_i$  (cf. (6.6)).

For the second term in the last expression in (7.90), we have that for  $r \geq r^*$ ,  $n \geq 1$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ ,

$$\begin{aligned}
& \sup_{n \geq 1} \mathbf{P}\left(I_{l,k}^{r,n}(s_l^r) \geq \frac{\epsilon_{1,i}}{\mathbf{I}} s_l^r, \tau_{l,2n-1}^r \leq r^2t\right) \\
&\leq \mathbf{P}(I_k^r(2r^2t) - I_k^r(\tau_{i,0}^r) > 0, \tau_{i,0}^r \leq 2r^2t) \\
&\quad + \mathbf{P}(I_k^r(\tau_{i,0}^r) \geq t_i^r) \\
&\leq p_{3,i}(2r^2t) \left( C_{3,i}^{(1)} \exp(-C_{3,i}^{(2)} L_0^r) + C_{3,i}^{(3)} \exp(-2C_{3,i}^{(4)} r^2t) \right) \\
&\quad + p_{2,i}(r^2t) \left( C_{2,i}^{(1)} \exp(-C_{2,i}^{(2)} L_0^r) + C_{2,i}^{(3)} \exp(-C_{2,i}^{(4)} r^2t) \right), \tag{7.92}
\end{aligned}$$

where  $p_{2,i}$  and  $p_{3,i}$  are polynomials (of degree at most  $i$  and  $i + 1$ , respectively) with non-negative coefficients, and  $C_{l,i}^{(m)} > 0$ , for  $l = 2, 3$ ,  $m = 1, 2, 3, 4$ , since (I) holds for  $i$  (with  $2t$  in place of  $t$ ), and  $k \in \underline{\mathcal{K}}_i$ . The polynomials and the constants do not depend on  $t$ ,  $n$ , or  $r$ . The first inequality in (7.92) holds since on  $\{\tau_{l,2n-1}^r \leq r^2t\}$ ,  $I_{l,k}^{r,n}(s_l^r) \leq I_k^r(\tau_{l,2n-1}^r + s_l^r) \leq I_k^r(r^2t + s_l^r) \leq I_k^r(2r^2t)$ , as

$s_l^r \leq M^r \leq r^2 t$ , and since by (7.91),

$$\begin{aligned}
\frac{\epsilon_{1,i} s_l^r}{\mathbf{I}} &\geq \frac{\gamma_i \epsilon_i s_l^r}{\mathbf{I}} = \frac{\gamma_i \epsilon_i (L_l^r - (|\underline{\mathcal{J}}_l| + 2))}{\mathbf{I}(\lambda_l^r + \epsilon_i)} \\
&\geq \frac{\epsilon_i^2 (L_l^r - (|\underline{\mathcal{J}}_l| + 2))}{(\lambda_l^r + \epsilon_i)} \\
&\geq \frac{\lambda_l}{\lambda_l^r + \epsilon_i} \left( \frac{L_l^r}{\lambda_l \epsilon_i} - \frac{\epsilon_i^2 (|\underline{\mathcal{J}}_l| + 2)}{\lambda_l} \right) \\
&\geq \frac{\lambda_l}{\lambda_l^r + \epsilon_i} \left( \frac{18 L_l^r}{\lambda_l - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j} - \frac{\epsilon_i (|\underline{\mathcal{J}}_l| + 2)}{\lambda_l} \right) \\
&\geq \frac{1}{2} \left( \frac{17 L_l^r}{\lambda_l - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j} \right) \geq t_i^r. \tag{7.93}
\end{aligned}$$

Here, in the second inequality in (7.93), we have used the fact that  $\epsilon_i < (\prod_{m=1}^{\mathbf{I}} \gamma_m) / \mathbf{I} \leq \gamma_i / \mathbf{I}$ , as  $\gamma_m \leq 1$  for all  $m$  (cf. (7.73)). In the third inequality of (7.93), we have used (6.4). In the fourth inequality we have used the fact that  $\epsilon_i < (\lambda_l - \sum_{j \in \underline{\mathcal{J}}_i} x_j^* \mu_j) / (18 \lambda_l)$  (cf. (7.73)) and  $\epsilon_i < 1$ . In the fifth inequality we have used (7.29) along with the estimate  $\lambda_l / (\lambda_l^r + \epsilon_i) > 1/2$  used in proving (7.87).

Combining all of the above (from (7.80) onwards), we have for all  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,

$$\begin{aligned}
&\sup_{n \geq 1} \mathbf{P} \left( T_{l,a(i)}^{r,n}(s_l^r) \geq (x_{a(i)}^* + \epsilon_i) s_l^r, s_l^r \leq \beta_{l,n}^r, \tau_{l,2n-1}^r \leq r^2 t \right) \\
&\leq K_{13} \exp(-K'_{13} s_l^r) \\
&\quad + \left( \lfloor (\mu_{a(i)}^r + \epsilon_i) r^2 t \rfloor + 1 \right) K_{14} \exp(-K'_{14} s_l^r) + K_{15} \exp(-K'_{15} r^2 t) \\
&\quad + p_{1,i}(r^2 t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-C_{1,i}^{(4)} r^2 t) \right) \\
&\quad + K_{16} \exp(-K'_{16} s_l^r) + |\underline{\mathcal{J}}_i| K_{17} \exp(-K'_{17} s_l^r) \\
&\quad + \sum_{j \in \underline{\mathcal{J}}_i} \left( \lfloor (\mu_j^r + \epsilon_{1,i}) r^2 t \rfloor + 1 \right) K_{18} \exp(-K'_{18} s_l^r) + |\underline{\mathcal{J}}_i| K_{19} \exp(-K'_{19} r^2 t) \\
&\quad + \sum_{j \in \underline{\mathcal{J}}_i} \left\{ \sum_{i' \in \underline{\mathcal{I}}_k(j)} \sup_{n \geq 1} \mathbf{P} \left( T_{l,a(i')}^{r,n}(s_l^r) \geq (x_{a(i')}^* + \epsilon_{i'}) s_l^r, s_l^r \leq \beta_{l,n}^r, \tau_{l,2n-1}^r \leq r^2 t \right) \right. \\
&\quad \quad \left. + p_{3,i}(2r^2 t) \left( C_{3,i}^{(1)} \exp(-C_{3,i}^{(2)} L_0^r) + C_{3,i}^{(3)} \exp(-2C_{3,i}^{(4)} r^2 t) \right) \right. \\
&\quad \quad \left. + p_{2,i}(r^2 t) \left( C_{2,i}^{(1)} \exp(-C_{2,i}^{(2)} L_0^r) + C_{2,i}^{(3)} \exp(-C_{2,i}^{(4)} r^2 t) \right) \right\}. \tag{7.94}
\end{aligned}$$

By the assumption that (II.1) holds with  $i'$  in place of  $i$ , for all  $i' < i$ , the definition of  $s_l^r$ , and the fact that  $s_l^r \geq L_0^r$ , for all  $l \in \mathcal{I}$ , it follows that for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,

$$\begin{aligned}
&\sup_{n \geq 1} \mathbf{P} \left( T_{l,a(i)}^{r,n}(s_l^r) \geq (x_{a(i)}^* + \epsilon_i) s_l^r, s_l^r \leq \beta_{l,n}^r, \tau_{l,2n-1}^r \leq r^2 t \right) \\
&\leq p_{4,i}(r^2 t) \left( C_{4,i}^{(1)} \exp(-C_{4,i}^{(2)} L_0^r) + C_{4,i}^{(3)} \exp(-C_{4,i}^{(4)} r^2 t) \right), \tag{7.95}
\end{aligned}$$

where  $p_{4,i}$  is a polynomial (of degree at most  $i + 1$ ) with positive coefficients, and  $C_{4,i}^{(m)} > 0$ , for  $m = 1, 2, 3, 4$ . The polynomial and the constants do not depend on  $t$  or  $r$ . This completes the proof that (II.1) holds for  $i$ .

*Proof of (II.2).* Fix a transition class  $\iota > i$ . For each  $r \geq 1$ ,  $n \geq 1$ ,  $s \geq 0$ ,  $j \in \mathcal{J}_i$ , on  $\{d_{\tau_{\iota,2n-1}}^r < \infty\}$  define

$$\begin{aligned} d_{\iota,i}^{r,n}(s) &= A_i^r(d_{\tau_{\iota,2n-1}}^r + s) - A_i^r(d_{\tau_{\iota,2n-1}}^r), \\ d_{\iota,j}^{r,n}(s) &= S_j^r(T_j^r(d_{\tau_{\iota,2n-1}}^r) + s) - S_j^r(T_j^r(d_{\tau_{\iota,2n-1}}^r)), \\ d_{\check{\iota},i}^{r,n}(s) &= \sup\{m \geq 0 : \xi_i^r(A_i^r(d_{\tau_{\iota,2n-1}}^r + m)) - \xi_i^r(A_i^r(d_{\tau_{\iota,2n-1}}^r)) \leq s\}, \\ d_{\check{\iota},j}^{r,n}(s) &= \sup\{m \geq 0 : \eta_j^r(S_j^r(T_j^r(d_{\tau_{\iota,2n-1}}^r)) + m + 1) - \eta_j^r(S_j^r(T_j^r(d_{\tau_{\iota,2n-1}}^r)) + 1) \leq s\}, \end{aligned}$$

and for concreteness on  $\{d_{\tau_{\iota,2n-1}}^r = \infty\}$ , we define  $d_{\iota,i}^{r,n}$ ,  $d_{\iota,j}^{r,n}$ ,  $d_{\check{\iota},i}^{r,n}$ ,  $d_{\check{\iota},j}^{r,n}$ , to be identically zero. Then, for each  $s \geq 0$ ,

$$d_{\iota,i}^{r,n}(s) \geq d_{\check{\iota},i}^{r,n}(s), \quad d_{\check{\iota},j}^{r,n}(s) \geq d_{\iota,j}^{r,n}(s) - 1. \quad (7.96)$$

We have for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , and  $n \geq 1$ ,

$$\begin{aligned} &\mathbf{P}\left(d_{\iota,a(i)}^{r,n}(d_{s_\iota}^r) \leq (x_{a(i)}^* - \epsilon_i) d_{s_\iota}^r, d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right) \\ &\leq \mathbf{P}\left(d_{\iota,a(i)}^{r,n}(d_{\tau_{\iota,a(i)}}^{r,n}(d_{s_\iota}^r)) \leq d_{\iota,a(i)}^{r,n}((x_{a(i)}^* - \epsilon_i) d_{s_\iota}^r), d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right) \\ &\leq \mathbf{P}\left(Q_i^r(d_{\tau_{\iota,2n-1}}^r + d_{s_\iota}^r) - Q_i^r(d_{\tau_{\iota,2n-1}}^r) \geq (\mu_{a(i)}^r \epsilon_i / 32) d_{s_\iota}^r, d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right) \\ &\quad + \mathbf{P}\left(d_{\iota,i}^{r,n}(d_{s_\iota}^r) - \sum_{j \in \underline{\mathcal{J}}_i} d_{\iota,j}^{r,n}(d_{\tau_{\iota,j}}^{r,n}(d_{s_\iota}^r))\right. \\ &\quad \left. - d_{\iota,a(i)}^{r,n}((x_{a(i)}^* - \epsilon_i) d_{s_\iota}^r) < (\mu_{a(i)}^r \epsilon_i / 32) d_{s_\iota}^r, d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right), \end{aligned} \quad (7.97)$$

where the last inequality in (7.97) uses (2.9).

For each  $r \geq 1$  and  $n \geq 1$ , let

$$\begin{aligned} d_{\Upsilon_{\iota,i}}^{r,n} &= \left\{ d_{\iota,i}^{r,n}(d_{s_\iota}^r) \geq (\lambda_i^r - (\mu_{a(i)}^r \epsilon_i / 32)) d_{s_\iota}^r; \right. \\ &\quad d_{\iota,a(i)}^{r,n}((x_{a(i)}^* - \epsilon_i) d_{s_\iota}^r) \leq (\mu_{a(i)}^r + \epsilon_{2,i})(x_{a(i)}^* - \epsilon_i) d_{s_\iota}^r + 1; \\ &\quad d_{\iota,j}^{r,n}((x_j^* + \epsilon_{1,i}) d_{s_\iota}^r) \leq (\mu_j^r + \epsilon_{3,j})(x_j^* + \epsilon_{1,i}) d_{s_\iota}^r + 1, \\ &\quad \left. d_{\tau_{\iota,j}}^{r,n}(d_{s_\iota}^r) < (x_j^* + \epsilon_{1,i}) d_{s_\iota}^r, j \in \underline{\mathcal{J}}_i; d_{\tau_{\iota,2n-1}}^r \leq r^2 t \right\}, \end{aligned} \quad (7.98)$$

where  $\epsilon_{2,i}$  and  $\epsilon_{3,j}$  are defined in (7.33) and (7.34), respectively.

Then, the last term in (7.97) is bounded above by

$$\begin{aligned} &\mathbf{P}\left(d_{\iota,i}^{r,n}(d_{s_\iota}^r) - \sum_{j \in \underline{\mathcal{J}}_i} d_{\iota,j}^{r,n}(d_{\tau_{\iota,j}}^{r,n}(d_{s_\iota}^r))\right. \\ &\quad \left. - d_{\iota,a(i)}^{r,n}((x_{a(i)}^* - \epsilon_i) d_{s_\iota}^r) < (\mu_{a(i)}^r \epsilon_i / 32) d_{s_\iota}^r, d_{\Upsilon_{\iota,i}}^{r,n}, d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right) \\ &+ \mathbf{P}\left(d_{\iota,i}^{r,n}(d_{s_\iota}^r) < (\lambda_i^r - (\mu_{a(i)}^r \epsilon_i / 32)) d_{s_\iota}^r, d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right) \\ &+ \mathbf{P}\left(d_{\iota,a(i)}^{r,n}((x_{a(i)}^* - \epsilon_i) d_{s_\iota}^r) > (\mu_{a(i)}^r + \epsilon_{2,i})(x_{a(i)}^* - \epsilon_i) d_{s_\iota}^r + 1, d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right) \\ &+ \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P}\left(d_{\iota,j}^{r,n}((x_j^* + \epsilon_{1,i}) d_{s_\iota}^r) > (\mu_j^r + \epsilon_{3,j})(x_j^* + \epsilon_{1,i}) d_{s_\iota}^r + 1, d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right) \\ &+ \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P}\left(d_{\tau_{\iota,j}}^{r,n}(d_{s_\iota}^r) \geq (x_j^* + \epsilon_{1,i}) d_{s_\iota}^r, d_{\tau_{\iota,2n-1}}^r \leq r^2 t\right). \end{aligned} \quad (7.99)$$

It can be verified (cf. (8.124) in [2] with  $|\mathcal{J}_i|$  in place of 2 there), using (7.28), (7.30), (7.85), (7.33), (7.34), and the fact that  $\mu_{a(i)} < 2\mu_{a(i)}^r$ , that on  $\mathcal{A}_{l,i}^{r,n}$ , for all  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq Mr$ ,

$$dA_{l,i}^{r,n}(d_{s_l}^r) - \sum_{j \in \underline{\mathcal{J}}_i} dS_{l,j}^{r,n}(dT_{l,j}^{r,n}(d_{s_l}^r)) - dS_{l,a(i)}^{r,n}((x_{a(i)}^* - \epsilon_i) d_{s_l}^r) > \frac{\mu_{a(i)}^r \epsilon_i}{32} d_{s_l}^r, \quad (7.100)$$

and then the first probability in the right side of (7.99) is zero.

For the second, third and fourth terms in the right side of (7.99), in a similar manner to that in (7.88)–(7.89), we have using Lemmas 6.5–6.7, that for  $r \geq r^*$  and  $t > 0$  satisfying  $r^2t \geq Mr$ ,

$$\begin{aligned} & \mathbf{P} \left( dA_{l,i}^{r,n}(d_{s_l}^r) < (\lambda_i^r - (\mu_{a(i)}^r \epsilon_i / 32)) d_{s_l}^r, d_{\tau_{l,2n-1}}^r \leq r^2t \right) \\ & \leq K_{20} \exp(-K'_{20} d_{s_l}^r) + (\lfloor (\lambda_i^r + \epsilon_i) r^2t \rfloor + 1) K_{21} \exp(-K'_{21} d_{s_l}^r) \\ & \quad + K_{22} \exp(-K'_{22} r^2t), \end{aligned} \quad (7.101)$$

$$\begin{aligned} & \mathbf{P} \left( dS_{l,a(i)}^{r,n}((x_{a(i)}^* - \epsilon_i) d_{s_l}^r) > (\mu_{a(i)}^r + \epsilon_{2,i})(x_{a(i)}^* - \epsilon_i) d_{s_l}^r + 1, d_{\tau_{l,2n-1}}^r \leq r^2t \right) \\ & \quad + \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P} \left( dS_{l,j}^{r,n}((x_j^* + \epsilon_{1,i}) d_{s_l}^r) > (\mu_j^r + \epsilon_{3,j})(x_j^* + \epsilon_{1,i}) d_{s_l}^r + 1, d_{\tau_{l,2n-1}}^r \leq r^2t \right) \\ & \leq K_{23} \exp(-K'_{23} d_{s_l}^r), \end{aligned} \quad (7.102)$$

where  $K_{20} = 1$ ,  $K'_{20} = (\lambda_i - \check{\epsilon}_i) \Lambda_i^{a,*}(\lambda_i^{-1}(1 + \check{\epsilon}_i/4\lambda_i)) > 0$ ,  $K_{21} = \exp(\Lambda_i^a(l_0)) > 0$ ,  $K'_{21} = l_0 \check{\epsilon}_i / 4\lambda_i > 0$ ,  $0 < l_0 \in \mathcal{O}_0$ ,  $K_{22} = \exp(\Lambda_i^{a,*}((\lambda_i(1 + \epsilon_i/3\lambda_i))^{-1})) > 0$ ,  $K'_{22} = \lambda_i \Lambda_i^{a,*}((\lambda_i(1 + \epsilon_i/3\lambda_i))^{-1}) > 0$ ,

$K_{23} = |\mathcal{J}_i| \max\{\exp(\Lambda_{a(i)}^{s,*}((\mu_{a(i)}(1 + \epsilon_{2,i}/3\mu_{a(i)}))^{-1})); \exp(\Lambda_j^{s,*}((\mu_j(1 + \epsilon_{3,j}/3\mu_j))^{-1})) : j \in \underline{\mathcal{J}}_i\} > 0$ ,

$K'_{23} = \min\{\mu_{a(i)}(x_{a(i)}^* - \epsilon_i) \Lambda_{a(i)}^{s,*}((\mu_{a(i)}(1 + \epsilon_{2,i}/3\mu_{a(i)}))^{-1}); \mu_j(x_j^* + \epsilon_{1,i}) \Lambda_j^{s,*}((\mu_j(1 + \epsilon_{3,j}/3\mu_j))^{-1}) : j \in \underline{\mathcal{J}}_i\} > 0$ .

For the last probability in (7.99), for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq Mr$ ,  $j \in \underline{\mathcal{J}}_i$ ,  $k = k(j)$ , and  $n \geq 1$ , as in (7.69) and using (7.91), we have

$$\begin{aligned} & \mathbf{P} \left( dT_{l,j}^{r,n}(d_{s_l}^r) \geq (x_j^* + \epsilon_{1,i}) d_{s_l}^r, d_{\tau_{l,2n-1}}^r \leq r^2t \right) \\ & \leq \sum_{i' \in \underline{\mathcal{I}}_k} \mathbf{P} \left( dT_{l,a(i')}^{r,n}(d_{s_l}^r) \leq (x_{a(i')}^* - \epsilon_{i'}) d_{s_l}^r, d_{\tau_{l,2n-1}}^r \leq r^2t \right). \end{aligned} \quad (7.103)$$

Finally, for the first term in the last inequality in (7.97), we have that for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq Mr$ ,

$$\begin{aligned} & \mathbf{P} \left( Q_i^r(d_{\tau_{l,2n-1}}^r + d_{s_l}^r) - Q_i^r(d_{\tau_{l,2n-1}}^r) \geq (\mu_{a(i)}^r \epsilon_i / 32) d_{s_l}^r, d_{\tau_{l,2n-1}}^r \leq r^2t \right) \\ & \leq \mathbf{P} \left( \sup_{0 \leq s \leq 2r^2t} Q_i^r(s) \geq (\mu_{a(i)}^r \epsilon_i / 32) d_{s_l}^r \right) \\ & \leq \mathbf{P} \left( \sup_{\tau_{i,0}^r \leq s \leq 2r^2t} R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i| \right) \\ & \leq p_{1,i}(2r^2t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-2C_{1,i}^{(4)} r^2t) \right), \end{aligned} \quad (7.104)$$

where  $p_{1,i}$  is a polynomial (of degree at most  $i + 1$ ) with non-negative coefficients, and  $C_{1,i}^{(m)} > 0$ , for  $m = 1, 2, 3, 4$ , since (I.1) was already proved to hold for  $i$  (with  $2t$  in place of  $t$ ). The polynomial

and the constants do not depend on  $t$  or  $r$ . The first inequality above uses the fact that  $M^r \geq d_{s_l}^r$ , and the second inequality holds since for  $r \geq r^*$ ,

$$\frac{\epsilon_i \mu_{a(i)}^r}{32} \cdot d_{s_l}^r \geq 2L_i^r, \quad (7.105)$$

which can be verified as in [2], using (7.27), (7.30), (6.4), (7.29),  $\epsilon_i < \min \left\{ \frac{\mu_{a(i)}}{1024 \sum_{j \in \underline{\mathcal{I}}_i} (\mu_j + \epsilon_i)}, 1 \right\}$  (cf. (7.73)).

Combining all of the above from (7.97) onwards, we have for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,

$$\begin{aligned} & \sup_{n \geq 1} \mathbf{P} \left( d_{T_{l,a(i)}^{r,n}}(d_{s_l}^r) \leq (x_{a(i)}^* - \epsilon_i) d_{s_l}^r, d_{\tau_{l,2n-1}}^r \leq r^2 t \right) \\ & \leq p_{1,i}(2r^2 t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-2C_{1,i}^{(4)} r^2 t) \right) \\ & \quad + K_{20} \exp(-K'_{20} d_{s_l}^r) + (\lfloor (\lambda_i^r + \epsilon_i) r^2 t \rfloor + 1) K_{21} \exp(-K'_{21} d_{s_l}^r) \\ & \quad + K_{22} \exp(-K'_{22} r^2 t) + K_{23} \exp(-K'_{23} d_{s_l}^r) \\ & \quad + \sum_{j \in \underline{\mathcal{I}}_i} \sum_{i' \in \underline{\mathcal{I}}_{k(j)}} \sup_{n \geq 1} \mathbf{P} \left( d_{T_{l,a(i')}^{r,n}}(d_{s_l}^r) \leq (x_{a(i')}^* - \epsilon_{i'}) d_{s_l}^r, d_{\tau_{l,2n-1}}^r \leq r^2 t \right). \end{aligned} \quad (7.106)$$

By the induction assumption that (II.2) holds with  $i'$  in place of  $i$  for all  $i' < i$ , we have that for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,

$$\begin{aligned} & \sup_{n \geq 1} \mathbf{P} \left( d_{T_{l,a(i)}^{r,n}}(d_{s_l}^r) \leq (x_{a(i)}^* - \epsilon_i) d_{s_l}^r, d_{\tau_{l,2n-1}}^r \leq r^2 t \right) \\ & \leq p_{5,i}(r^2 t) \left( C_{5,i}^{(1)} \exp(-C_{5,i}^{(2)} L_0^r) + C_{5,i}^{(3)} \exp(-C_{5,i}^{(4)} r^2 t) \right), \end{aligned} \quad (7.107)$$

where  $p_{5,i}$  is a polynomial (of degree at most  $i+1$ ) with non-negative coefficients, and  $C_{5,i}^{(m)} > 0$ , for  $m = 1, 2, 3, 4$ . The polynomial and the constants do not depend on  $r$  or  $t$ . It then follows that (II.2) holds for  $i$ , since  $d_{s_l}^r \geq L_0^r$ .

*Proof of (II.3).* For  $r \geq r^*$  and  $\iota > i$  fixed,

$$\begin{aligned} & \mathbf{P} \left( T_{a(i)}^r(t_\iota^r) \leq (x_{a(i)}^* - \epsilon_i) t_\iota^r \right) \\ & \leq \mathbf{P} \left( A_i^r(t_\iota^r) - \sum_{j \in \underline{\mathcal{I}}_i} S_j^r(T_j^r(t_\iota^r)) - S_{a(i)}^r(T_{a(i)}^r(t_\iota^r)) \right. \\ & \quad \left. \geq A_i^r(t_\iota^r) - \sum_{j \in \underline{\mathcal{I}}_i} S_j^r(T_j^r(t_\iota^r)) - S_{a(i)}^r((x_{a(i)}^* - \epsilon_i) t_\iota^r) \right) \\ & \leq \mathbf{P} \left( Q_i^r(t_\iota^r) \geq (\mu_{a(i)}^r \epsilon_i / 32) t_\iota^r \right) \\ & \quad + \mathbf{P} \left( A_i^r(t_\iota^r) - \sum_{j \in \underline{\mathcal{I}}_i} S_j^r(T_j^r(t_\iota^r)) - S_{a(i)}^r((x_{a(i)}^* - \epsilon_i) t_\iota^r) < (\mu_{a(i)}^r \epsilon_i / 32) t_\iota^r \right). \end{aligned} \quad (7.108)$$

Let,

$$\begin{aligned} \Upsilon_{\iota,i}^r = & \left\{ A_i^r(t_\iota^r) \geq (\lambda_i^r - (\mu_{a(i)}^r \epsilon_i / 32)) t_\iota^r; \right. \\ & S_{a(i)}^r((x_{a(i)}^* - \epsilon_i) t_\iota^r) \leq (\mu_{a(i)}^r + \epsilon_{2,i})(x_{a(i)}^* - \epsilon_i) t_\iota^r; \\ & S_j^r((x_j^* + \epsilon_{1,i}) t_\iota^r) \leq (\mu_j^r + \epsilon_{3,j})(x_j^* + \epsilon_{1,i}) t_\iota^r, j \in \underline{\mathcal{I}}_i; \\ & \left. T_j^r(t_\iota^r) < (x_j^* + \epsilon_{1,i}) t_\iota^r, j \in \underline{\mathcal{I}}_i \right\}. \end{aligned} \quad (7.109)$$

Now, for the last term in (7.108), we have

$$\begin{aligned}
& \mathbf{P}\left(A_i^r(t_l^r) - \sum_{j \in \underline{\mathcal{J}}_i} S_j^r(T_j^r(t_l^r)) - S_{a(i)}^r((x_{a(i)}^* - \epsilon_i)t_l^r) < (\mu_{a(i)}^r \epsilon_i / 32)t_l^r\right) \\
& \leq \mathbf{P}\left(A_i^r(t_l^r) - \sum_{j \in \underline{\mathcal{J}}_i} S_j^r(T_j^r(t_l^r)) - S_{a(i)}^r((x_{a(i)}^* - \epsilon_i)t_l^r) < (\mu_{a(i)}^r \epsilon_i / 32)t_l^r, \Upsilon_{l,i}^r\right) \\
& \quad + \mathbf{P}\left(A_i^r(t_l^r) < (\lambda_i^r - (\mu_{a(i)}^r \epsilon_i / 32))t_l^r\right), \\
& \quad + \mathbf{P}\left(S_{a(i)}^r((x_{a(i)}^* - \epsilon_i)t_l^r) > (\mu_{a(i)}^r + \epsilon_{2,i})(x_{a(i)}^* - \epsilon_i)t_l^r\right) \\
& \quad + \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P}\left(S_j^r((x_j^* + \epsilon_{1,i})t_l^r) > (\mu_j^r + \epsilon_{3,j})(x_j^* + \epsilon_{1,i})t_l^r\right) \\
& \quad + \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P}\left(T_j^r(t_l^r) \geq (x_j^* + \epsilon_{1,i})t_l^r\right). \tag{7.110}
\end{aligned}$$

For the first term in the right side of (7.110) we have that on  $\Upsilon_{l,i}^r$ , for  $r \geq r^*$ ,

$$A_i^r(t_l^r) - \sum_{j \in \underline{\mathcal{J}}_i} S_j^r(T_j^r(t_l^r)) - S_{a(i)}^r((x_{a(i)}^* - \epsilon_i)t_l^r) \geq \frac{\mu_{a(i)}^r \epsilon_i}{32} t_l^r, \tag{7.111}$$

in a similar manner to that in (7.100). Hence the first probability in the right side of (7.110) is zero.

For the second term in the right side of (7.110) we have, using (7.15) and (7.31), for  $r \geq r^*$ ,

$$\begin{aligned}
\mathbf{P}\left(A_i^r(t_l^r) < (\lambda_i^r - (\mu_{a(i)}^r \epsilon_i / 32))t_l^r\right) & \leq \mathbf{P}\left(A_i^r(t_l^r) < (\lambda_i^r - \check{\epsilon}_i / 2)t_l^r\right) \\
& \leq K_{24} \exp(-K'_{24} t_l^r) + K_{25} \exp(-K'_{25} t_l^r), \tag{7.112}
\end{aligned}$$

by Lemma 6.7 and (6.27), where  $K_{24} = 1$ ,  $K'_{24} = (\lambda_i - \check{\epsilon}_i) \Lambda_i^{a,*} (\lambda_i^{-1} (1 + \check{\epsilon}_i / 4 \lambda_i)) > 0$ ,  $K_{25} = \exp(\Lambda_i^a(l_0)) > 0$ ,  $K'_{25} = l_0 \check{\epsilon}_i / 4 \lambda_i > 0$ , and  $0 < l_0 \in \mathcal{O}_0$ .

For the third and fourth terms in the right side of (7.110) we have, using Lemma 6.7 (since  $(x_{a(i)}^* - \epsilon_i)t_l^r > 2/\epsilon_{2,i}$  and  $(x_j^* + \epsilon_{1,i})t_l^r > 2/\epsilon_{3,j}$ , for all  $j \in \underline{\mathcal{J}}_i$ , by (7.9), (7.12), and (7.35)), for  $r \geq r^*$ ,

$$\begin{aligned}
& \mathbf{P}\left(S_{a(i)}^r((x_{a(i)}^* - \epsilon_i)t_l^r) > (\mu_{a(i)}^r + \epsilon_{2,i})(x_{a(i)}^* - \epsilon_i)t_l^r\right) \\
& \quad + \sum_{j \in \underline{\mathcal{J}}_i} \mathbf{P}\left(S_j^r((x_j^* + \epsilon_{1,i})t_l^r) > (\mu_j^r + \epsilon_{3,j})(x_j^* + \epsilon_{1,i})t_l^r\right) \\
& \leq K_{26} \exp(-K'_{26} t_l^r), \tag{7.113}
\end{aligned}$$

where  $K_{26} = |\mathcal{J}_i| \max\{\exp(\Lambda_{a(i)}^{s,*}((\mu_{a(i)}(1 + \epsilon_{2,i}/3\mu_{a(i)}))^{-1})); \exp(\Lambda_j^{s,*}((\mu_j(1 + \epsilon_{3,j}/3\mu_j))^{-1})) : j \in \underline{\mathcal{J}}_i\} > 0$ ,  $K'_{26} = \min\{\mu_{a(i)}(x_{a(i)}^* - \epsilon_i) \Lambda_{a(i)}^{s,*}((\mu_{a(i)}(1 + \epsilon_{2,i}/3\mu_{a(i)}))^{-1}); \mu_j(x_j^* + \epsilon_{1,i}) \Lambda_j^{s,*}((\mu_j(1 + \epsilon_{3,j}/3\mu_j))^{-1}) : j \in \underline{\mathcal{J}}_i\} > 0$ .

For the fifth term in the right side of (7.110) in a similar manner to that in (7.70) and using (7.91), we have for  $j \in \underline{\mathcal{J}}_i$ ,  $k = k(j)$ , and  $r \geq r^*$ ,

$$\mathbf{P}(T_j^r(t_l^r) \geq (x_j^* + \epsilon_{1,i})t_l^r) \leq \sum_{i' \in \underline{\mathcal{I}}_k} \mathbf{P}(T_{a(i')}^r(t_l^r) < (x_{a(i')}^* - \epsilon_{i'})t_l^r). \tag{7.114}$$

For the first term in the last inequality in (7.108), we have that for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,

$$\begin{aligned} & \mathbf{P}\left(Q_i^r(t_\iota^r) \geq (\mu_{a(i)}^r \epsilon_i / 32)t_\iota^r\right) \\ & \leq \mathbf{P}\left(\sup_{\tau_{i,0}^r \leq s \leq r^2 t} R_i^r(s) \geq L_i^r - |\underline{\mathcal{J}}_i|\right) \\ & \leq p_{1,i}(r^2 t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-C_{1,i}^{(4)} r^2 t) \right), \end{aligned} \quad (7.115)$$

where  $p_{1,i}$  is a polynomial (of degree at most  $i + 1$ ) with non-negative coefficients, and  $C_{1,i}^{(m)} > 0$ , for  $m = 1, 2, 3, 4$ , since (I.1) was already proved to hold for  $i$ . The polynomial and the constants do not depend on  $t$  or  $r$ . For the first inequality in (7.115), we have used the fact that

$$\begin{aligned} \frac{\mu_{a(i)}^r \epsilon_i}{32} t_\iota^r &= \frac{\mu_{a(i)}^r \epsilon_i}{32} \cdot \frac{8L_\iota^r}{\lambda_\iota - \sum_{j \in \underline{\mathcal{J}}_\iota} x_j^* \mu_j} \\ &\geq \frac{\mu_{a(i)}^r}{\mu_{a(i)}} \cdot \frac{\mu_{a(i)} \epsilon_i}{4(\lambda_\iota - \sum_{j \in \underline{\mathcal{J}}_\iota} x_j^* \mu_j)} \cdot \frac{L_i^r}{\epsilon_i^3} \\ &\geq \frac{\mu_{a(i)}}{8(\lambda_\iota - \sum_{j \in \underline{\mathcal{J}}_\iota} x_j^* \mu_j)} \cdot \frac{L_i^r}{\epsilon_i^2} \\ &\geq 2L_i^r, \end{aligned} \quad (7.116)$$

by (6.4), (7.27), and the fact that  $\epsilon_i < \min\left\{\frac{\mu_{a(i)}}{16(\lambda_\iota - \sum_{j \in \underline{\mathcal{J}}_\iota} x_j^* \mu_j)}, 1\right\}$  (cf. (7.73)).

Combining all of the above (from (7.108) onwards), we have for all  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,

$$\begin{aligned} & \mathbf{P}\left(T_{a(i)}^r(t_\iota^r) \leq (x_{a(i)}^* - \epsilon_i)t_\iota^r\right) \\ & \leq p_{1,i}(r^2 t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-C_{1,i}^{(4)} r^2 t) \right) \\ & \quad + K_{24} \exp(-K_{24}' t_\iota^r) + K_{25} \exp(-K_{25}' t_\iota^r) + K_{26} \exp(-K_{26}' t_\iota^r) \\ & \quad + \sum_{j \in \underline{\mathcal{J}}_i} \sum_{i' \in \underline{\mathcal{I}}_{k(j)}} \mathbf{P}\left(T_{a(i')}^r(t_\iota^r) < (x_{a(i')}^* - \epsilon_{i'})t_\iota^r\right). \end{aligned} \quad (7.117)$$

By the induction assumption that (II.3) holds with  $i'$  in place of  $i$ , for all  $i' < i$ , and the definition of  $t_\iota^r$ , it follows that for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ ,

$$\begin{aligned} & \mathbf{P}\left(T_{a(i)}^r(t_\iota^r) \leq (x_{a(i)}^* - \epsilon_i)t_\iota^r\right) \\ & \leq p_{6,i}(r^2 t) \left( C_{6,i}^{(1)} \exp(-C_{6,i}^{(2)} L_0^r) + C_{6,i}^{(3)} \exp(-C_{6,i}^{(4)} r^2 t) \right), \end{aligned} \quad (7.118)$$

where  $p_{6,i}$  is a polynomial (of degree at most  $i + 1$ ) with non-negative coefficients, and  $C_{6,i}^{(m)} > 0$  for  $m = 1, 2, 3, 4$ . Thus, (II.3) holds for  $i$ .  $\square$

## 7.7 Proofs of Theorems 6.1, 7.1, and 7.7

**Proof of Theorem 7.7.** Fix  $i \in \mathcal{I} \setminus \{i^*\}$  and assume that (I) and (II) hold for all  $i' < i$ . Then, by Lemmas 7.4–7.5, (i)–(iv) in Lemma 7.3 hold for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2 t \geq M^r$ , for  $i$  and each  $j \in \underline{\mathcal{J}}_i$ , and hence (I) holds for  $i$  by Lemma 7.3. By Lemma 7.6, we have that (II) also holds for  $i$ . The conclusion of the first statement in Theorem 7.7 then follows by the induction principle.

For (III), suppose that  $i^*$  is a transition class, and let  $k \in \underline{\mathcal{K}}_{i^*}$ . By the above, for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ , we have that (II) holds with  $i^*$  in place of  $i$ ,  $i' \in \underline{\mathcal{I}}_k$  in place of  $i$ .

For  $r \geq r^*$ , by the same proof as for Lemma 7.4, (ii) of Lemma 7.3 holds with  $i^*$  in place of  $i$  (cf. (7.69)), and then in the same manner as in the proof of Lemma 7.3 (the part of the proof of (I.1) involving down excursions of  $R_{i^*}^r$ , with  $i^*$  in place of  $i$  there), it can be shown that for  $r \geq r^*$  and  $t > 0$  satisfying  $r^2t \geq M^r$ ,

$$\begin{aligned} & \mathbf{P} \left( \inf_{\tau_{i^*,0}^r \leq s \leq r^2t} R_{i^*}^r(s) \leq -L_{i^*}^r + |\underline{\mathcal{J}}_{i^*}| \right) \\ & \leq p_{1,i^*}(r^2t) \left( C_{1,i^*}^{(1)} \exp(-C_{1,i^*}^{(2)}L_0^r) + C_{1,i^*}^{(3)} \exp(-C_{1,i^*}^{(4)}r^2t) \right), \end{aligned} \quad (7.119)$$

where  $p_{1,i^*}$  is a polynomial (of degree at most  $\mathbf{I} + 1$ ) with non-negative coefficients, and  $C_{1,i^*}^{(m)} > 0$ ,  $m = 1, 2, 3, 4$  are independent of  $r$  and  $t$ . Thus (III.1) holds.

To establish (III.2), we note that by the same proof as for Lemma 7.4, for  $r \geq r^*$  and  $t > 0$  satisfying  $r^2t \geq M^r$ , (iii) of Lemma 7.3 holds with  $i^*$  in place of  $i$  (i.e., for  $j \in \underline{\mathcal{J}}_{i^*}$ ), since (II.3) holds with  $i^*$  in place of  $i$  and  $i' \in \underline{\mathcal{I}}_k$  (where  $k = k(j)$ ) in place of  $i$  (cf. (7.70)). Then, by the same proof as for (I.2) in Lemma 7.3 (with  $i^*$  in place of  $i$  there, cf. (7.59)–(7.66)), we have,

$$\begin{aligned} & \mathbf{P} \left( I_k^r(\tau_{i^*,0}^r) \geq t_{i^*}^r \right) \\ & \leq p_{2,i^*}(r^2t) \left( C_{2,i^*}^{(1)} \exp(-C_{2,i^*}^{(2)}L_0^r) + C_{2,i^*}^{(3)} \exp(-C_{2,i^*}^{(4)}r^2t) \right), \end{aligned} \quad (7.120)$$

where  $p_{2,i^*}$  is a polynomial (of degree at most  $\mathbf{I}$ ) with non-negative coefficients, and  $C_{2,i^*}^{(m)} > 0$ ,  $m = 1, 2, 3, 4$  are independent of  $r$  and  $t$ .

Finally, for (III.3) we argue as in the proof of (I.3) in Lemma 7.3 (cf. (7.67)), that  $I_k^r$ ,  $k \in \underline{\mathcal{K}}_{i^*}$ , can increase only at times  $s \geq 0$  such that  $Q_{i^*}^r(s) \leq |\underline{\mathcal{J}}_{i^*}|$ , so that for  $r \geq r^*$ ,  $t > 0$  satisfying  $r^2t \geq M^r$ ,

$$\begin{aligned} & \mathbf{P} \left( I_k^r(r^2t) - I_k^r(\tau_{i^*,0}^r) > 0, \tau_{i^*,0}^r < r^2t \right) \\ & \leq \mathbf{P} \left( \inf_{\tau_{i^*,0}^r \leq s \leq r^2t} Q_{i^*}^r(s) \leq |\underline{\mathcal{J}}_{i^*}| \right) \\ & \leq p_{1,i^*}(r^2t) \left( C_{1,i^*}^{(1)} \exp(-C_{1,i^*}^{(2)}L_0^r) + C_{1,i^*}^{(3)} \exp(-C_{1,i^*}^{(4)}r^2t) \right), \end{aligned} \quad (7.121)$$

by (7.119). □

**Proof of Theorem 7.1.** Note that (I) and (II) hold for all  $i \in \mathcal{I} \setminus \{i^*\}$  by Theorem 7.7.

Fix  $i \in \mathcal{I} \setminus \{i^*\}$ ,  $k \in \underline{\mathcal{K}}_i$ , and  $\epsilon > 0$ . For  $t = 0$ , (7.1) and (7.3) hold trivially since  $R_i^r(0) = 0$  if  $\tau_{i,0}^r = 0$ . So we assume that  $t > 0$  is fixed. Since  $M^r = O(\log r)$ , there exists an  $r_t \geq r^*$  such that for all  $r \geq r_t$ ,  $r^2t \geq M^r$ . Then for  $r \geq r_t$ , by (I.1),

$$\begin{aligned} & \mathbf{P} \left( \sup_{\tau_{i,0}^r \leq s \leq r^2t} |R_i^r(s)| \geq L_i^r - |\underline{\mathcal{J}}_i| \right) \\ & \leq p_{1,i}(r^2t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)}L_0^r) + C_{1,i}^{(3)} \exp(-C_{1,i}^{(4)}r^2t) \right). \end{aligned} \quad (7.122)$$

Since  $t_i^r = o(r)$ , there is an  $r'_t \geq r_t$  such that  $t_i^r \leq r\epsilon$  for all  $r \geq r'_t$ . Then, by (I.2), for  $r \geq r'_t$

$$\begin{aligned} & \mathbf{P} \left( I_k^r(\tau_{i,0}^r) \geq r\epsilon \right) \\ & \leq \mathbf{P} \left( I_k^r(\tau_{i,0}^r) \geq t_i^r \right) \\ & \leq p_{2,i}(r^2t) \left( C_{2,i}^{(1)} \exp(-C_{2,i}^{(2)}L_0^r) + C_{2,i}^{(3)} \exp(-C_{2,i}^{(4)}r^2t) \right). \end{aligned} \quad (7.123)$$



Finally, we have by (I.3) that for  $r \geq r_t$ ,

$$\begin{aligned} & \mathbf{P} \left( I_k^r(r^2t) - I_k^r(\tau_{i,0}^r) > 0, \tau_{i,0}^r < r^2t \right) \\ & \leq p_{3,i}(r^2t) \left( C_{3,i}^{(1)} \exp(-C_{3,i}^{(2)}L_0^r) + C_{3,i}^{(3)} \exp(-C_{3,i}^{(4)}r^2t) \right). \end{aligned} \quad (7.124)$$

Since  $L_0^r = \lceil c \log r \rceil$  (and hence for  $d > 0$ ,  $\exp(-dL_0^r) \leq r^{-cd}$ ), it follows from (7.122)–(7.124) that there is a constant  $c_0 > 0$  (not depending on  $t$  or  $r$ ) such that if  $c \geq c_0$ , the expressions in (7.122)–(7.124) tend to zero as  $r \rightarrow \infty$  (for each fixed  $t > 0$ ). This verifies (7.2)–(7.3).

For the proof of (7.4)–(7.6), suppose that  $i^*$  is a transition class. Let  $k \in \underline{\mathcal{K}}_{i^*}$ , and  $\epsilon > 0$ . For  $t = 0$ , (7.4) and (7.6) hold trivially since  $\tau_{i^*,0}^r > 0$  for a transition class, so we fix  $t > 0$ . Then, as above, for  $r \geq r_t$ , we have that  $r^2t \geq Mr$ . We then have, for  $r \geq r_t$ , by (III.1),

$$\begin{aligned} & \mathbf{P} \left( \inf_{\tau_{i^*,0}^r \leq s \leq r^2t} Q_{i^*}^r(s) \leq |\underline{\mathcal{J}}_{i^*}| \right) \\ & = \mathbf{P} \left( \inf_{\tau_{i^*,0}^r \leq s \leq r^2t} R_{i^*}^r(s) \leq -L_{i^*}^r + |\underline{\mathcal{J}}_{i^*}| \right) \\ & \leq p_{1,i^*}(r^2t) \left( C_{1,i^*}^{(1)} \exp(-C_{1,i^*}^{(2)}L_0^r) + C_{1,i^*}^{(3)} \exp(-C_{1,i^*}^{(4)}r^2t) \right), \end{aligned} \quad (7.125)$$

Since  $t_{i^*}^r = o(r)$ , there exists  $r_t'' \geq r_t$  such that  $t_{i^*}^r \leq r\epsilon$  for all  $r \geq r_t''$ . Then by (III.2), we have for  $r \geq r_t''$ ,

$$\begin{aligned} & \mathbf{P} \left( I_k^r(\tau_{i^*,0}^r) \geq r\epsilon \right) \\ & \leq \mathbf{P} \left( I_k^r(\tau_{i^*,0}^r) \geq t_{i^*}^r \right) \\ & \leq p_{2,i^*}(r^2t) \left( C_{2,i^*}^{(1)} \exp(-C_{2,i^*}^{(2)}L_0^r) + C_{2,i^*}^{(3)} \exp(-C_{2,i^*}^{(4)}r^2t) \right). \end{aligned} \quad (7.126)$$

Finally, we have by (III.3) that for  $r \geq r_t$ ,

$$\begin{aligned} & \mathbf{P} \left( I_k^r(r^2t) - I_k^r(\tau_{i^*,0}^r) > 0, \tau_{i^*,0}^r < r^2t \right) \\ & \leq p_{3,i^*}(r^2t) \left( C_{3,i^*}^{(1)} \exp(-C_{3,i^*}^{(2)}L_0^r) + C_{3,i^*}^{(3)} \exp(-C_{3,i^*}^{(4)}r^2t) \right). \end{aligned} \quad (7.127)$$

Since  $L_0^r = \lceil c \log r \rceil$ , it follows that there is a constant  $c_1 \geq c_0$  (not depending on  $t$  or  $r$ ), such that if  $c \geq c_1$ , then the expressions in (7.125)–(7.127) tend to zero as  $r \rightarrow \infty$  (for each fixed  $t > 0$ ). This proves (7.4)–(7.6).  $\square$

**Proof of Theorem 6.1.** It suffices to show that as  $r \rightarrow \infty$ ,  $(\hat{Q}_i^r : i \in \mathcal{I} \setminus \{i^*\}; \hat{I}_k^r : k \in \underline{\mathcal{K}}_i, i \in \mathcal{I}) \Rightarrow \mathbf{0}$ . Note for this that  $k \in \mathcal{K} \setminus \{k^*\}$  is either in  $\underline{\mathcal{K}}_i$  for some  $i \in \mathcal{I} \setminus \{i^*\}$  or it is in  $\underline{\mathcal{K}}_{i^*}$  and  $i^*$  is a transition class.

Fix  $t > 0$  and  $\epsilon > 0$ . By Theorem 7.1 and the properties of  $\{L_i^r\}_{i \in \mathcal{I}}$ , there is  $r(\epsilon, t) \geq 1$  such that whenever  $r \geq r(\epsilon, t)$  we have  $2L_i^r/r < \epsilon$  and for all  $i \in \mathcal{I} \setminus \{i^*\}$ ,  $k \in \underline{\mathcal{K}}_i$ ,

$$\mathbf{P} \left( \sup_{\tau_{i,0}^r \leq s \leq r^2t} |R_i^r(s)| \geq L_i^r - |\underline{\mathcal{J}}_i| \right) < \epsilon, \quad (7.128)$$

$$\mathbf{P} \left( I_k^r(\tau_{i,0}^r) \geq r\epsilon \right) < \epsilon, \quad (7.129)$$

$$\mathbf{P} \left( I_k^r(r^2t) - I_k^r(\tau_{i,0}^r) > 0, \tau_{i,0}^r < r^2t \right) < \epsilon, \quad (7.130)$$

and if  $i^*$  is a transition buffer, for all  $k \in \underline{\mathcal{K}}_{i^*}$ ,

$$\mathbf{P} \left( I_k^r(r^2t) - I_k^r(\tau_{i^*,0}^r) > 0, \tau_{i^*,0}^r < r^2t \right) < \epsilon. \quad (7.131)$$

Recalling the definition of  $\|\cdot\|_t$  from Section 1.1, we then have for all  $r \geq r(\varepsilon, t)$ ,

$$\begin{aligned}
& \mathbf{P}\left(\|\hat{Q}_i^r\|_t \geq \varepsilon \text{ for some } i \in \mathcal{I} \setminus \{i^*\}, \text{ or } \|\hat{I}_k^r\|_t \geq \varepsilon \text{ for some } k \in \mathcal{K} \setminus \{k^*\}\right) \\
&= \mathbf{P}\left(\|Q_i^r\|_{r^2t} \geq r\varepsilon \text{ for some } i \in \mathcal{I} \setminus \{i^*\}, \text{ or } \|I_k^r\|_{r^2t} \geq r\varepsilon \text{ for some } k \in \mathcal{K} \setminus \{k^*\}\right) \\
&\leq \mathbf{P}\left(\sup_{\tau_{i,0}^r \leq s \leq r^2t} Q_i^r(s) \geq 2L_i^r \text{ for some } i \in \mathcal{I} \setminus \{i^*\}, \text{ or} \right. \\
&\quad \left. I_k^r(\tau_{i,0}^r) \geq r\varepsilon \text{ for some } k \in \underline{\mathcal{K}}_i \text{ and } i \in \mathcal{I}, \text{ or} \right. \\
&\quad \left. \tau_{i,0}^r < r^2t \text{ and } I_k^r(r^2t) - I_k^r(\tau_{i,0}^r) > 0 \text{ for some } k \in \underline{\mathcal{K}}_i \text{ and } i \in \mathcal{I}\right) \\
&\leq \sum_{i \in \mathcal{I} \setminus \{i^*\}} \mathbf{P}\left(\sup_{\tau_{i,0}^r \leq s \leq r^2t} |R_i^r(s)| \geq L_i^r - |\underline{\mathcal{J}}_i|\right) \\
&\quad + \sum_{i \in \mathcal{I}} \sum_{k \in \underline{\mathcal{K}}_i} \left\{ \mathbf{P}(I_k^r(\tau_{i,0}^r) \geq r\varepsilon) + \mathbf{P}(I_k^r(r^2t) - I_k^r(\tau_{i,0}^r) > 0, \tau_{i,0}^r < r^2t) \right\} \\
&< (\mathbf{I} + 2\mathbf{K} - 3)\varepsilon. \tag{7.132}
\end{aligned}$$

□

## 8 Weak Convergence under the Threshold Policy

This section is devoted to the proof of Theorem 5.3.

### 8.1 Fluid Limits for Allocation Processes

Recall the definitions from Section 2.2 of the functions,

$$i : \mathcal{J} \rightarrow \mathcal{I}, \quad k : \mathcal{J} \rightarrow \mathcal{K}, \tag{8.1}$$

where for  $j \in \mathcal{J}$ ,  $i(j)$  is the buffer processed by activity  $j$  and  $k(j)$  is the server which processes activity  $j$ . Also, recall from Assumption 3.1 that

$$\lambda^r \rightarrow \lambda, \quad \mu^r \rightarrow \mu, \quad \text{as } r \rightarrow \infty. \tag{8.2}$$

Recall from Sections 5 and 6 that  $a(i)$  is the basic activity above buffer  $i$  in the server-buffer tree (where only basic activities are included in that tree), and  $k^*$  is the server at the root of the server-buffer tree in layer  $l = l^*$ . Using the notation of Section 6.1, for each  $r \geq 1$ , for  $i \in \mathcal{I}^{l^*}$ , define

$$z_{k^*}^r = z_{k^*}, \quad y_i^r = \frac{z_{k^*}^r}{\mu_{a(i)}^r}, \tag{8.3}$$

and for layer  $l = l^* - 1, \dots, 1$ , define by backwards induction on  $l$ ,

$$z_k^r = y_{i(b(k))}^r \mu_{b(k)}^r, \quad \text{for each } k \in \mathcal{K}^l, \tag{8.4}$$

$$y_i^r = \frac{z_{k(a(i))}^r}{\mu_{a(i)}^r}, \quad \text{for each } i \in \mathcal{I}^l, \tag{8.5}$$

where  $b(k)$  is the basic activity immediately above server  $k$  which links it to a buffer in the next highest layer (see Figure 7). Here  $z_{k^*}^*$  is the variable, determined from one component of the unique optimal solution  $(y^*, z^*)$  of the dual program (4.8). Then, since each basic activity either links a

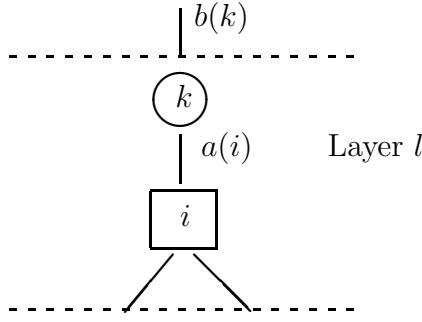


Figure 7: Activity  $b(k)$  above a server  $k$  that is in layer  $l$ .

buffer in some layer  $l + 1$  to a server below it in layer  $l$ , or links a server in some layer  $l$  to a buffer immediately below it in layer  $l$ , we have for each basic activity  $j$ ,

$$y_{i(j)}^r \mu_j^r = z_{k(j)}^r, \quad j = 1, \dots, \mathbf{B}, \quad (8.6)$$

i.e., for any basic activity  $j$ ,

$$((y^r)' \mathbf{R}^r)_j = ((z^r)' \mathbf{A})_j, \quad (8.7)$$

where  $y^r = (y_i^r : i \in \mathcal{I})$ ,  $z^r = (z_k^r : k \in \mathcal{K})$ ,

$$\mathbf{R}^r = \mathbf{C} \text{diag}(\mu^r), \quad \mathbf{C}_{ij} = \delta_{i,i(j)}, \quad \mathbf{A}_{kj} = \delta_{k,k(j)}, \quad \text{for all } i \in \mathcal{I}, k \in \mathcal{K}, j \in \mathcal{J}, \quad (8.8)$$

and  $\delta_{x,y} = 1$  if  $x = y$  and  $\delta_{x,y} = 0$  otherwise.

Lemmas 8.1, 8.4, and 8.3, which are proved below, will be used in Section 9 for the proof of asymptotic optimality (Theorem 5.4) as well as in the proof of Lemma 8.5, which shows that the fluid scaled allocation processes,  $\{\bar{T}^{r,*}\}$ , associated with the threshold policy, converge in distribution (as  $r$  goes to infinity) to the nominal allocation processes  $\bar{T}^*$  defined in (6.10).

**Lemma 8.1** *We have that  $(y^r, z^r) \rightarrow (y^*, z^*)$  as  $r \rightarrow \infty$ , where  $(y^*, z^*)$  is the optimal solution to the dual linear program specified in (4.8).*

**Proof.** We prove that  $z_k^r \rightarrow z_k^*$ ,  $y_k^r \rightarrow y_k^*$ , as  $r \rightarrow \infty$ , for all  $k \in \mathcal{K}^l$ ,  $i \in \mathcal{I}^l$ , for  $l = 1, 2, \dots, l^*$ , by backward induction on  $l$ . The result of the lemma then follows, since each server belongs to some layer and similarly for each buffer. From (4.9) and (8.6), we have that  $y_{i(j)}^* \mu_j = z_{k(j)}^*$  and  $y_{i(j)}^r \mu_j = z_{k(j)}^r$  for all basic activities  $j = 1, 2, \dots, \mathbf{B}$ .

Note that layer  $l^*$  has one server  $k^*$  and the buffers in that layer are indexed by  $\underline{\mathcal{I}}_{k^*}$ . We have that (by definition)  $z_{k^*}^r = z_{k^*}^*$ . Then for  $i \in \underline{\mathcal{I}}_{k^*}$ , we have by (8.2) and (8.3), that

$$y_i^r = \frac{z_{k^*}^r}{\mu_{a(i)}^r} \rightarrow \frac{z_{k^*}^*}{\mu_{a(i)}} = \frac{z_{k^*(a(i))}^*}{\mu_{a(i)}} = y_{i(a(i))}^* = y_i^*, \quad \text{as } r \rightarrow \infty. \quad (8.9)$$

Now, for the induction step, suppose that  $(y_i^r, z_k^r) \rightarrow (y_i^*, z_k^*)$  as  $r \rightarrow \infty$ , for all  $i \in \mathcal{I}^l$ ,  $k \in \mathcal{K}^l$ , some  $l \geq 2$ . For  $k \in \mathcal{K}^{(l-1)}$ ,  $i \in \underline{\mathcal{I}}_k$ , we have by (8.2), (8.6), and (4.9) that as  $r \rightarrow \infty$ ,

$$z_k^r = y_{i(b(k))}^r \mu_{b(k)}^r \rightarrow y_{i(b(k))}^* \mu_{b(k)} = z_{k(b(k))}^* = z_k^*, \quad (8.10)$$

since  $i(b(k)) \in \mathcal{I}^l$ , and

$$y_i^r = \frac{z_{k(a(i))}^r}{\mu_{a(i)}^r} \rightarrow \frac{z_{k(a(i))}^*}{\mu_{a(i)}} = y_{a(i)}^* = y_i^*, \quad (8.11)$$

since  $k(a(i)) = k$ . This completes the induction step and the conclusion of the lemma then follows.  $\square$

**Definition 8.2** *A sequence of processes with paths in  $D^m$  for some  $m \geq 1$  is called C-tight if it is tight in  $D^m$  and any weak limit point of the sequence (obtained as a weak limit along a subsequence) has continuous paths almost surely.*

For Lemma 8.3 below, note that the result holds for *any* sequence of scheduling controls, not just for those associated with our threshold policy. In addition to the scaled processes defined in Section 3.3, for any  $r \geq 1$  and control  $T^r$  for the  $r^{\text{th}}$  parallel server system, we define the following fluid scaled processes. For each  $t \geq 0$ , let

$$\bar{A}^r(t) = r^{-2}A^r(r^2t), \quad (8.12)$$

$$\bar{S}^r(t) = r^{-2}S^r(r^2t), \quad (8.13)$$

$$\bar{I}^r(t) = r^{-2}I^r(r^2t), \quad (8.14)$$

$$\bar{Q}^r(t) = r^{-2}Q^r(r^2t). \quad (8.15)$$

**Lemma 8.3** *Let  $\{T^r\}$  be any sequence of scheduling controls (one for each member of the sequence of parallel server systems). Then*

$$\left\{ (\bar{Q}^r(\cdot), \bar{A}^r(\cdot), \bar{S}^r(\cdot), \bar{T}^r(\cdot), \bar{I}^r(\cdot)) \right\} \text{ is C-tight.}$$

**Proof.** It follows from (3.21) that

$$(\bar{A}^r(\cdot), \bar{S}^r(\cdot)) \Rightarrow (\lambda(\cdot), \mu(\cdot)) \quad \text{as } r \rightarrow \infty, \quad (8.16)$$

where  $\lambda(t) = \lambda t$  and  $\mu(t) = \mu t$  for all  $t \geq 0$ . In addition, since they correspond to cumulative allocations of time, each of the components of  $T^r$  is Lipschitz continuous with a Lipschitz constant of one and this property is preserved by the fluid scaled processes  $\bar{T}^r$ . It follows immediately from this and (8.16) that  $\{(\bar{A}^r(\cdot), \bar{S}^r(\cdot), \bar{T}^r(\cdot))\}$  is C-tight, cf. Theorem 15.5 in [4]. From the equations (2.15)–(2.16) for queue length and idletime we have that for each  $t \geq 0$ ,

$$\bar{Q}^r(t) = \bar{A}^r(t) - \mathbf{C}\bar{S}^r(\bar{T}^r(t)), \quad (8.17)$$

$$\bar{I}^r(t) = \mathbf{1}t - \mathbf{A}\bar{T}^r(t). \quad (8.18)$$

Combining these with the C-tightness established above, the continuous mapping theorem, and a random time change theorem (cf. [4], p. 145), we obtain the desired result.  $\square$

Recall that  $T^{r,*}$  denotes the allocation process when our threshold policy is used in the  $r^{\text{th}}$  parallel server system. For the non-basic activities  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ ,  $T_j^{r,*} \equiv \mathbf{0}$ . We indicate the fact that this threshold policy is being used by appending a superscript of  $*$  to the associated processes. By (3.15), (3.17), (3.18), Assumption 3.3, and (8.8), for any  $r \geq 1$ ,

$$\bar{Q}^{r,*}(t) = r^{-1}\hat{A}^r(t) - r^{-1}\mathbf{C}\hat{S}^r(\bar{T}^{r,*}(t)) + (\lambda^r - \mathbf{R}^r x^*)t + \mathbf{R}^r(x^*t - \bar{T}^{r,*}(t)), \quad (8.19)$$

$$\bar{I}^{r,*}(t) = \mathbf{1}t - \mathbf{A}\bar{T}^{r,*}(t), \quad (8.20)$$

where  $\mathbf{1}$  is the  $\mathbf{K}$ -dimensional vector of all ones. Recall from Assumption 3.6 that

$$r(\lambda^r - \mathbf{R}^r x^*) \rightarrow \theta \quad \text{as } r \rightarrow \infty, \quad (8.21)$$

where  $\theta \in \mathbb{R}^{\mathbf{I}}$ .

**Lemma 8.4** For  $(y^r, z^r)$  as defined in (8.3)–(8.5), we have that

$$(y^r)' \mathbf{R}^r (x^* t - \bar{T}^{r,*}(t)) = (z^r)' \bar{I}^{r,*}(t), \quad \text{for all } t \geq 0.$$

**Proof.** Since  $x_j^* = 0$  and  $\bar{T}_j^{r,*} \equiv \mathbf{0}$  for all non-basic activities  $j$ , we have from Assumption 3.3, (8.7) and (8.20) that for all  $t \geq 0$ ,

$$\begin{aligned} (y^r)' \mathbf{R}^r (x^* t - \bar{T}^{r,*}(t)) &= (z^r)' (\mathbf{A} x^* t - \mathbf{A} \bar{T}^{r,*}(t)) \\ &= (z^r)' (\mathbf{1} t - \mathbf{A} \bar{T}^{r,*}(t)) \\ &= (z^r)' \bar{I}^{r,*}(t). \end{aligned} \tag{8.22}$$

□

Lemma 8.5 below is needed in the proof of Theorem 5.3, in combining the functional central limit result (3.21), with a random time change theorem.

**Lemma 8.5** For the fluid scaled allocation processes,  $\bar{T}_j^{r,*}$ ,  $j \in \mathcal{J}$ , we have,

$$\bar{T}^{r,*} \Rightarrow \bar{T}^* \quad \text{as } r \rightarrow \infty, \tag{8.23}$$

where  $\bar{T}^*(t) = x^* t$ , for all  $t \geq 0$ .

**Proof.** We note first that since  $\bar{T}_j^{r,*} \equiv \mathbf{0}$  and  $\bar{T}_j^* \equiv \mathbf{0}$  for  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ , for all  $r \geq 1$ , we have that (trivially)

$$\bar{T}_j^{r,*}(\cdot) \Rightarrow \bar{T}_j^*(\cdot) \quad \text{as } r \rightarrow \infty, \quad \text{for } j = \mathbf{B} + 1, \dots, \mathbf{J}. \tag{8.24}$$

Now from (3.21) and the fact that  $\bar{T}_j^{r,*}(t) \leq t$ ,  $j \in \mathcal{J}$ , for all  $t \geq 0$ , we have that

$$(r^{-1} \hat{A}_i^r(\cdot), r^{-1} \hat{S}_j^r(\bar{T}_j^{r,*}(\cdot))) : i \in \mathcal{I}, j \in \mathcal{J} \Rightarrow \mathbf{0}, \quad \text{as } r \rightarrow \infty. \tag{8.25}$$

Using (8.19) and Lemma 8.4 we have that for each  $t \geq 0$ ,

$$y^r \cdot \bar{Q}^{r,*}(t) = y^r \cdot \bar{X}^{r,*}(t) + z^r \cdot \bar{I}^{r,*}(t), \tag{8.26}$$

where

$$\bar{X}^{r,*}(t) = r^{-1} \hat{A}^r(t) - r^{-1} \mathbf{C} \hat{S}^r(\bar{T}^{r,*}(t)) + (\lambda^r - \mathbf{R}^r x^*) t, \tag{8.27}$$

and by (8.25), (8.2) and (8.21),

$$\bar{X}^{r,*}(\cdot) \Rightarrow \mathbf{0} \quad \text{as } r \rightarrow \infty. \tag{8.28}$$

Thus, for  $t \geq 0$ ,

$$y_i^r \bar{Q}_i^{r,*}(t) = \bar{\zeta}^{r,*}(t) + z_k^r \bar{I}_k^{r,*}(t), \tag{8.29}$$

where, by (8.26), (8.28), Lemma 8.1, and Theorem 6.1 we have

$$\bar{\zeta}^{r,*}(\cdot) \equiv y^r \cdot \bar{X}^{r,*}(\cdot) - \sum_{i \in \mathcal{I} \setminus \{i^*\}} y_i^r \bar{Q}_i^{r,*}(\cdot) + \sum_{k \in \mathcal{K} \setminus \{k^*\}} z_k^r \bar{I}_k^{r,*}(\cdot) \Rightarrow \mathbf{0} \quad \text{as } r \rightarrow \infty. \tag{8.30}$$

Here  $y^r > 0$ ,  $z^r > 0$ ,  $\bar{Q}_i^{r,*} \geq 0$ ,  $\bar{\zeta}^{r,*}(0) = 0$ ,  $\bar{I}_k^{r,*}(0) = 0$ ,  $\bar{I}_k^{r,*}(\cdot)$  is continuous and non-decreasing. Furthermore, for  $r$  sufficiently large, in the  $r^{\text{th}}$  system operating under the threshold policy, the idletime at server  $k^*$  can increase only if the queue length at buffer  $i^*$  is at or below the level  $L_{i^*}^r$ .

Hence  $\bar{I}_{k^*}^{r,*}$  can only increase if  $\bar{Q}_{i^*}^{r,*}$  is at or below the level  $L_{i^*}^r/r^2$ . (Here, when  $i^*$  is a transition buffer,  $r$  needs to be large enough that  $L_{i^*}^r \geq |\underline{\mathcal{J}}_{i^*}|$ . For if  $L_{i^*}^r < Q_{i^*}^{r,*}(t) \leq |\underline{\mathcal{J}}_{i^*}|$ , then the idletime at server  $k^*$  might increase at  $t$  if the  $Q_{i^*}^{r,*}(t)$  jobs are in service or in suspension at the servers in  $\underline{\mathcal{K}}_{i^*}$ . If  $i^*$  is a non-transition buffer, server  $k^*$  is busy whenever buffer  $i^*$  is nonempty, and in particular whenever  $Q_{i^*}^{r,*} > L_{i^*}^r > 0$ .)

It then follows from an oscillation inequality for solutions of a perturbed Skorokhod problem (cf. the proof of Theorem 5.1 in [40]) that

$$z_{k^*}^r \bar{I}_{k^*}^{r,*}(t) \leq - \inf_{0 \leq s \leq t} (\bar{\zeta}^{r,*}(s)) + y_{i^*}^r L_{i^*}^r r^{-2}, \quad \text{for all } t \geq 0. \quad (8.31)$$

Hence, it follows by (8.30), the continuous mapping theorem, Lemma 8.1, and the facts that  $L_{i^*}^r = O(\log r)$  and  $z_{k^*} > 0$ , that

$$\bar{I}_{k^*}^{r,*} \Rightarrow \mathbf{0} \quad \text{as } r \rightarrow \infty. \quad (8.32)$$

By (8.29) and (8.30), we then have that

$$\bar{Q}_{i^*}^{r,*} \Rightarrow \mathbf{0} \quad \text{as } r \rightarrow \infty, \quad (8.33)$$

since  $y_{i^*}^r > 0$ . To obtain the conclusion of the lemma, we appeal to Lemma 8.3 and assume that  $(\bar{Q}(\cdot), \bar{A}(\cdot), \bar{S}(\cdot), \bar{T}(\cdot), \bar{I}(\cdot))$  is obtained as a weak limit point of  $\{(\bar{Q}^{r,*}(\cdot), \bar{A}^r(\cdot), \bar{S}^r(\cdot), \bar{T}^{r,*}(\cdot), \bar{I}^{r,*}(\cdot))\}$  along a subsequence. Since  $(\bar{Q}(\cdot), \bar{I}(\cdot)) \equiv \mathbf{0}$  (by Theorem 6.1, (8.32), and (8.33)), and  $(\bar{A}(\cdot), \bar{S}(\cdot)) = (\lambda(\cdot), \mu(\cdot))$  (by (8.16)) we have by passing to the limit in (8.17)–(8.18) that  $\bar{T}$  satisfies

$$\mathbf{0} = \lambda t - \mathbf{R}\bar{T}(t), \quad (8.34)$$

$$\mathbf{0} = \mathbf{1}t - \mathbf{A}\bar{T}(t), \quad (8.35)$$

for all  $t \geq 0$ , where  $\bar{T}$  inherits the properties (3.8) from  $\bar{T}^{r,*}$ . Thus, the fluid system is balanced and incurs no idleness (cf. Section 3) under  $\bar{T}$ . Hence, by Definition 3.4, Assumption 3.3, and Lemma 3.5, we have that  $\bar{T} = \bar{T}^*$ . Thus, the weak limit point is the same along all convergent subsequences. It follows that  $(\bar{Q}^{r,*}(\cdot), \bar{A}^r(\cdot), \bar{S}^r(\cdot), \bar{T}^{r,*}(\cdot), \bar{I}^{r,*}(\cdot)) \Rightarrow (\mathbf{0}, \lambda(\cdot), \mu(\cdot), \bar{T}^*(\cdot), \mathbf{0})$ , as  $r \rightarrow \infty$ . The conclusion of the lemma then follows.  $\square$

## 8.2 Convergence of Diffusion Scaled Performance Measures under the Threshold Policy—Proof of Theorem 5.3

We now prove that the diffusion scaled performance measures for our sequence of parallel server systems operating under the allocation processes  $\{T^{r,*}\}$  converge in distribution to the processes given in (4.17)–(4.18). We indicate the fact that this threshold policy is being used by appending a superscript of  $* r$  to the associated processes.

**Proof of Theorem 5.3.** For each  $t \geq 0$ , we have by multiplying (8.26) through by  $r$  that

$$\hat{W}^{r,*}(t) \equiv y^r \cdot \hat{Q}^{r,*}(t) = y^r \cdot \hat{X}^{r,*}(t) + z^r \cdot \hat{I}^{r,*}(t) \quad (8.36)$$

where  $\hat{X}^{r,*}$  is defined by (4.7) (with  $\bar{T}^r$  replaced by  $\bar{T}^{r,*}$  there). By the functional central limit result (3.21), Lemma 8.5, and a random time change theorem (cf. [4], §17), we have that as  $r \rightarrow \infty$ ,

$$(\hat{A}^r(\cdot), \hat{S}^r(\bar{T}^{r,*}(\cdot))) \Rightarrow (\tilde{A}(\cdot), \tilde{S}(\bar{T}^*(\cdot))). \quad (8.37)$$

It then follows from the definition of  $\hat{X}^{r,*}$ , (8.2), and (8.21) that

$$\hat{X}^{r,*} \Rightarrow \tilde{X}, \quad \text{as } r \rightarrow \infty, \quad (8.38)$$

where

$$\tilde{X}(t) = \tilde{A}(t) - \mathbf{C}\tilde{S}(\bar{T}^*(t)) + \theta t, \quad t \geq 0, \quad (8.39)$$

is an  $\mathbf{I}$ -dimensional Brownian motion with drift  $\theta$  and a diagonal covariance matrix whose  $i^{\text{th}}$  diagonal entry is  $\lambda_i a_i^2 + \sum_{j=1}^{\mathbf{J}} \mathbf{C}_{ij} \mu_j b_j^2 x_j^*$ .

Rearranging (8.36), we have that for  $t \geq 0$ ,

$$y_i^r \hat{Q}_i^{r,*}(t) = \hat{\zeta}^{r,*}(t) + z_k^* \hat{I}_k^{r,*}(t), \quad (8.40)$$

where

$$\hat{\zeta}^{r,*}(\cdot) \equiv y^r \cdot \hat{X}^{r,*}(\cdot) - \sum_{i \in \mathcal{I} \setminus \{i^*\}} y_i^r \hat{Q}_i^{r,*}(\cdot) + \sum_{k \in \mathcal{K} \setminus \{k^*\}} z_k^r \hat{I}_k^{r,*}(\cdot) \Rightarrow y^* \cdot \tilde{X}(\cdot), \quad (8.41)$$

as  $r \rightarrow \infty$ , by (8.38), Lemma 8.1, and Theorem 6.1. Here  $y^r > 0$ ,  $z^r > 0$ ,  $\hat{Q}_i^{r,*} \geq 0$ ,  $\hat{\zeta}^{r,*}(0) = 0$ ,  $\hat{I}_k^{r,*}(0) = 0$ ,  $\hat{I}_k^{r,*}(\cdot)$  is continuous and non-decreasing. Furthermore, for sufficiently large  $r$ ,  $\hat{I}_k^{r,*}$  can increase only if  $\hat{Q}_i^{r,*}$  is at or below level  $L_{i^*}^r/r$  (since  $L_{i^*}^r > |\underline{\mathcal{J}}_{i^*}|$ , for sufficiently large  $r$ ). Then, by Corollary 4.3 in [40], since  $L_{i^*}^r/r \rightarrow 0$  as  $r \rightarrow \infty$ , it follows that

$$(\hat{Q}_i^{r,*}, \hat{I}_k^{r,*}) \Rightarrow (\tilde{Q}_i^*, \tilde{I}_k^*), \quad \text{as } r \rightarrow \infty, \quad (8.42)$$

where  $\tilde{Q}_i^*, \tilde{I}_k^*$  are given by (4.17)-(4.18). Combining this with Theorem 6.1 and Lemma 8.1 yields,

$$(\hat{Q}^{r,*}, \hat{I}^{r,*}, \hat{W}^{r,*}) \Rightarrow (\tilde{Q}^*, \tilde{I}^*, \tilde{W}^*), \quad \text{as } r \rightarrow \infty, \quad (8.43)$$

where  $\tilde{W}^*$  is defined in (4.15)-(4.16), and  $\tilde{Q}^*, \tilde{I}^*$  are defined by (4.17)-(4.18).

## 9 Asymptotic Optimality of the Threshold Policy

In this section we prove Theorem 5.4. We follow a similar development to that in Section 9 of [3]. Before proceeding with the proof, we first establish some preliminary results concerning fluid scaled processes under a sequence of scheduling controls,  $T = \{T^r\}$  (one for each member of the sequence of parallel server systems). The associated queue length and idletime processes will be denoted by  $Q^r, I^r$ , and the fluid and diffusion scaled versions of these processes will be denoted by  $\bar{Q}^r, \bar{I}^r$  and  $\hat{Q}^r, \hat{I}^r$ , respectively. We also let

$$\underline{J}(T) = \liminf_{r \rightarrow \infty} \hat{J}^r(T^r), \quad (9.1)$$

where  $\hat{J}^r(T^r)$  is defined by (3.13). When our sequence of threshold controls  $\{T^{r,*}\}$  is used, we append a superscript  $*$  to the queue length, idletime etc. processes, i.e., we use  $Q^{r,*}, I^{r,*}$ , etc.

The next lemma implies that, when searching for an asymptotically optimal policy, we may restrict to those policies whose associated fluid scaled allocation processes converge (along a subsequence) to those given by  $\bar{T}^*$ .

**Lemma 9.1** *Let  $T = \{T^r\}$  be a sequence of scheduling controls such that  $\underline{J}(T) < \infty$ . Consider a subsequence  $\{T^{r'}\}$  of  $\{T^r\}$  along which the liminf in the definition of  $\underline{J}(T)$  is achieved, i.e.,*

$$\lim_{r' \rightarrow \infty} \hat{J}^{r'}(T^{r'}) = \underline{J}(T). \quad (9.2)$$

Then,

$$(\bar{Q}^{r'}(\cdot), \bar{A}^{r'}(\cdot), \bar{S}^{r'}(\cdot), \bar{T}^{r'}(\cdot), \bar{I}^{r'}(\cdot)) \Rightarrow (\mathbf{0}, \lambda(\cdot), \mu(\cdot), \bar{T}^*(\cdot), \mathbf{0}) \quad \text{as } r' \rightarrow \infty, \quad (9.3)$$

where  $\bar{T}^r(t) = T^r(r^2 t)/r^2$  and  $\bar{T}^*(t) = x^* t$ , for all  $r \geq 1$  and  $t \geq 0$ ,  $x^*$  is given by the heavy traffic Assumption 3.3,  $\mathbf{0}$  denotes the constant process that stays at the origin (in the appropriate dimension) for all time, and  $\lambda(t) = \lambda t$ ,  $\mu(t) = \mu t$  for all  $t \geq 0$ .

**Proof.** It follows from Lemma 8.3 that

$$\left\{ (\bar{Q}^{r'}(\cdot), \bar{A}^{r'}(\cdot), \bar{S}^{r'}(\cdot), \bar{T}^{r'}(\cdot), \bar{I}^{r'}(\cdot)) \right\} \quad (9.4)$$

is C-tight. Thus, it suffices to show that all weak limit points of this sequence are given by the right member of (9.3). For this, suppose that

$$(\bar{Q}(\cdot), \bar{A}(\cdot), \bar{S}(\cdot), \bar{T}(\cdot), \bar{I}(\cdot)), \quad (9.5)$$

is obtained as a weak limit of (9.4) along a subsequence indexed by  $r''$ . Without loss of generality, by appealing to the Skorokhod representation theorem (cf. [10], Theorem 3.1.8), we may choose an equivalent distributional representation (for which we use the same symbols) such that all of the random processes in (9.4) indexed by  $r''$  in place of  $r'$ , as well as the limit (9.5), are defined on the same probability space and the convergence in distribution is replaced by almost sure convergence uniformly on compact time intervals, so that a.s., as  $r'' \rightarrow \infty$ ,

$$(\bar{Q}^{r''}(\cdot), \bar{A}^{r''}(\cdot), \bar{S}^{r''}(\cdot), \bar{T}^{r''}(\cdot), \bar{I}^{r''}(\cdot)) \rightarrow (\bar{Q}(\cdot), \bar{A}(\cdot), \bar{S}(\cdot), \bar{T}(\cdot), \bar{I}(\cdot)) \quad \text{u.o.c.} \quad (9.6)$$

From (8.16) we have that a.s.,  $\bar{A}(\cdot) = \lambda(\cdot)$  and  $\bar{S}(\cdot) = \mu(\cdot)$ . We next show that a.s.,  $\bar{Q}(\cdot) \equiv \mathbf{0}$ . Combining the fact that  $\lim_{r'' \rightarrow \infty} \hat{J}^{r''}(T^{r''}) = \underline{J}(T) < \infty$  with (9.6) and Fatou's lemma, we have

$$\begin{aligned} 0 = \lim_{r'' \rightarrow \infty} \frac{1}{r''} \hat{J}^{r''}(T^{r''}) &\geq \mathbf{E} \left( \int_0^\infty e^{-\gamma t} \liminf_{r'' \rightarrow \infty} (h \cdot \bar{Q}^{r''}(t)) dt \right) \\ &= \mathbf{E} \left( \int_0^\infty e^{-\gamma t} h \cdot \bar{Q}(t) dt \right). \end{aligned} \quad (9.7)$$

Since  $h_i > 0$  for all  $i \in \mathcal{I}$ , and a.s.,  $\bar{Q}$  has continuous paths in  $\mathbb{R}_+^{\mathbf{I}}$ , it follows from the above that a.s.,  $\bar{Q}(\cdot) \equiv \mathbf{0}$ . Then, by letting  $r = r'' \rightarrow \infty$  in (8.17)–(8.18) and using (8.16), (9.6), and the definition of  $\mathbf{R}$ , we have a.s., for each  $t \geq 0$ ,

$$\mathbf{0} = \lambda t - \mathbf{R}\bar{T}(t), \quad (9.8)$$

$$\bar{I}(t) = \mathbf{1}t - \mathbf{A}\bar{T}(t). \quad (9.9)$$

Multiplying (9.8) by  $(y^*)'$  while recalling that  $y^* \cdot \lambda = 1 = z^* \cdot \mathbf{1}$  and  $(y^*)'\mathbf{R} = (z^*)'\mathbf{A} - [0' (u^*)']$ , where  $u^* > 0$  (cf. Lemma 4.5), we obtain

$$0 = z^* \cdot \bar{I}(t) + [0' (u^*)']\bar{T}(t). \quad (9.10)$$

Since a.s., the components of  $\bar{I}(\cdot)$  and  $\bar{T}(\cdot)$  inherit the property from  $\bar{I}^{r''}(\cdot)$ ,  $\bar{T}^{r''}(\cdot)$  that they are all non-negative for all time, it follows from (9.10), and the fact that  $z_k^* > 0$  for all  $k \in \mathcal{K}$ ,  $u_j^* > 0$  for all  $j = 1, 2, \dots, \mathbf{J} - \mathbf{B}$ , that a.s.,  $\bar{I}(\cdot) = \mathbf{0}$ , and  $\bar{T}_j(\cdot) = \mathbf{0}$  for all  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ . We then observe that  $\bar{T}$  is a fluid control under which the fluid system in (3.6)–(3.7) is balanced and incurs no idleness (cf. Section 3). Hence, by Definition 3.4, Assumption 3.3, and Lemma 3.5 we have that  $\bar{T}(\cdot) = \bar{T}^*(\cdot)$ .  $\square$

**Proof of Theorem 5.4.** We first concentrate on proving the inequality on the left side of (5.2). For this, let  $T \equiv \{T^r\}$  be a sequence of scheduling controls. If  $\underline{J}(T) = \infty$ , then the inequality holds trivially and so we assume that  $\underline{J}(T) < \infty$ . Recall the definitions of  $(y^r, z^r)$  from (8.3)–(8.5). For each  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ , let  $u_{j-\mathbf{B}}^r = ((z^r)'\mathbf{A} - (y^r)'\mathbf{R}^r)_j$ . Then by (8.7), we have that

$$(y^r)'\mathbf{R}^r = (z^r)'\mathbf{A} - [0' (u^r)'], \quad (9.11)$$



where for sufficiently large  $r$ ,  $u^r > 0$  by Lemma 4.5, (8.2), and Lemma 8.1.

For each  $i \in \mathcal{I}$ , let

$$h_i^r = \frac{h_i y_i^r}{y_i^*}, \quad (9.12)$$

where  $h$  is the holding cost vector appearing in (3.13). Since  $i^*$  is the ‘‘cheapest’’ buffer (cf. Section 4.2), we have that

$$\frac{h_{i^*}}{y_{i^*}^*} \leq \frac{h_i}{y_i^*}, \quad \text{for all } i \in \mathcal{I}. \quad (9.13)$$

Then, using (9.11), (3.17)–(3.18) and (9.12)–(9.13), we have for all  $t \geq 0$ ,

$$\begin{aligned} h^r \cdot \hat{Q}^r(t) &= \sum_{i=1}^{\mathbf{I}} h_i^r \hat{Q}_i^r(t) \\ &\geq \frac{h_{i^*}}{y_{i^*}^*} \sum_{i=1}^{\mathbf{I}} y_i^r \hat{Q}_i^r(t) \\ &= \frac{h_{i^*}}{y_{i^*}^*} y^r \cdot \hat{Q}^r(t) \\ &= \frac{h_{i^*}}{y_{i^*}^*} \left( y^r \cdot \hat{X}^r(t) + \hat{V}^r(t) \right), \end{aligned} \quad (9.14)$$

where  $\hat{X}^r$  is given in (4.7),

$$\hat{V}^r(t) = ((z^r)' \mathbf{A} - [0' (u^r)']) \hat{Y}^r(t) = z^r \cdot \hat{T}^r(t) - u^r \cdot \hat{Y}_N^r(t), \quad t \geq 0, \quad (9.15)$$

$\hat{Y}^r$  is defined by (3.15) and  $\hat{Y}_N^r$  is the  $(\mathbf{J} - \mathbf{B})$ -dimensional process whose components are  $\hat{Y}_j^r$ ,  $j = \mathbf{B} + 1, \dots, \mathbf{J}$ .

Now, since  $h^r \cdot \hat{Q}^r(t) \geq 0$  for all  $t \geq 0$ ,  $y^r \cdot \hat{X}^r$  starts from zero, and  $\hat{V}^r$  is non-decreasing (for sufficiently large  $r$ ) and starts from zero, it follows from the well known minimality of the solution of the Skorokhod problem (cf. Appendix B in [3]) that for all  $r$  sufficiently large,

$$\hat{V}^r(t) \geq \sup_{0 \leq s \leq t} \left( -y^r \cdot \hat{X}^r(s) \right) \quad \text{for all } t \geq 0, \quad (9.16)$$

and hence

$$h^r \cdot \hat{Q}^r(t) \geq \frac{h_{i^*}}{y_{i^*}^*} \varphi \left( y^r \cdot \hat{X}^r \right) (t) \quad \text{for all } t \geq 0, \quad (9.17)$$

where  $\varphi(x)(t) \equiv x(t) + \sup_{0 \leq s \leq t} (-x(s))$  for all  $t \geq 0$  and  $x \in \mathbf{D}$  satisfying  $x(0) = 0$ .

Now, let  $\{T^{r'}\}$  be a subsequence of  $\{T^r\}$  such that  $\lim_{r' \rightarrow \infty} \hat{J}^{r'}(T^{r'}) = \underline{J}(T)$ . By Lemma 9.1, the fact that the limit in (9.3) is deterministic, and (3.21), we have that as  $r' \rightarrow \infty$ ,

$$\left( \hat{A}^{r'}(\cdot), \hat{S}^{r'}(\cdot), \hat{T}^{r'}(\cdot) \right) \Rightarrow \left( \tilde{A}(\cdot), \tilde{S}(\cdot), \tilde{T}^*(\cdot) \right). \quad (9.18)$$

By invoking the Skorokhod representation theorem, we may assume without loss of generality that the convergence above is almost surely uniform on compact time intervals and then for  $\hat{X}^r$  given by (4.7), using Assumption 3.6, we have that a.s. as  $r' \rightarrow \infty$ ,

$$\left( \hat{A}^{r'}(\cdot), \hat{S}^{r'}(\cdot), \hat{T}^{r'}(\cdot), \hat{X}^{r'}(\cdot) \right) \rightarrow \left( \tilde{A}(\cdot), \tilde{S}(\cdot), \tilde{T}^*(\cdot), \tilde{X}(\cdot) \right) \quad \text{u.o.c.}, \quad (9.19)$$

where  $\tilde{X}(t) = \tilde{A}(t) - \mathbf{C}\tilde{S}(\bar{T}^*(t)) + \theta t$ , for  $t \geq 0$ , defines a Brownian motion as described in Definition 4.1. By Fatou's lemma, we have

$$\underline{J}(T) = \lim_{r' \rightarrow \infty} \hat{J}^{r'}(T^{r'}) \geq \mathbf{E} \left( \int_0^\infty e^{-\gamma t} \liminf_{r' \rightarrow \infty} (h \cdot \hat{Q}^{r'}(t)) dt \right). \quad (9.20)$$

Now we claim that a.s., for each  $t \geq 0$ ,

$$\liminf_{r' \rightarrow \infty} (h \cdot \hat{Q}^{r'}(t)) \geq h \cdot \tilde{Q}^*(t), \quad (9.21)$$

where  $\tilde{Q}^*$  is given by (4.15)–(4.18). To see this, fix  $\omega \in \Omega$  such that  $\omega$  is in the set of probability one where the convergence in (9.19) holds u.o.c., and the limits have continuous paths. Fix  $t \geq 0$ . If the left member of the inequality (9.21) is infinite at  $\omega$ , then the inequality clearly holds. On the other hand, if the left member is finite at  $\omega$ , then there is a further subsequence indexed by  $r''$  (possibly depending on  $\omega$  and  $t$ ) such that

$$\lim_{r'' \rightarrow \infty} (h \cdot \hat{Q}^{r''}(t, \omega)) = \liminf_{r'' \rightarrow \infty} (h \cdot \hat{Q}^{r''}(t, \omega)) < \infty. \quad (9.22)$$

Since  $h_i > 0$  and  $\hat{Q}_i^{r''}(t, \omega) \geq 0$ , for all  $i \in \mathcal{I}$ , it follows that  $\hat{Q}_i^{r''}(t, \omega)$  is bounded as  $r'' \rightarrow \infty$ , for all  $i \in \mathcal{I}$ , and then using the fact that  $h_i^r \rightarrow h_i$ , for all  $i \in \mathcal{I}$  (cf. Lemma 8.1), we have

$$\lim_{r'' \rightarrow \infty} (h - h^{r''}) \cdot \hat{Q}^{r''}(t, \omega) = 0. \quad (9.23)$$

Then, using (9.17), (9.19), the continuity of  $\varphi$  on  $\mathbf{D}$ , the fact that  $\varphi(y^* \cdot \tilde{X})(\cdot, \omega)$  is continuous, and (4.16)–(4.17), we have

$$\begin{aligned} \lim_{r'' \rightarrow \infty} h \cdot \hat{Q}^{r''}(t, \omega) &= \lim_{r'' \rightarrow \infty} \left( h^{r''} \cdot \hat{Q}^{r''}(t, \omega) + (h - h^{r''}) \cdot \hat{Q}^{r''}(t, \omega) \right) \\ &\geq \liminf_{r'' \rightarrow \infty} \frac{h_{i^*}}{y_{i^*}^*} \varphi \left( y^{r''} \cdot \hat{X}^{r''} \right) (t, \omega) \\ &= \frac{h_{i^*}}{y_{i^*}^*} \varphi \left( y^* \cdot \tilde{X} \right) (t, \omega) = \frac{h_{i^*}}{y_{i^*}^*} \tilde{W}^*(t, \omega) = h \cdot \tilde{Q}^*(t, \omega). \end{aligned}$$

Thus, (9.21) holds. Now, substituting this in (9.20), we conclude that

$$\underline{J}(T) \geq \mathbf{E} \left( \int_0^\infty e^{-\gamma t} h \cdot \tilde{Q}^*(t) dt \right) \equiv J^*. \quad (9.24)$$

This completes the proof of the inequality in the left side of (5.2).

Suppose now that the threshold control  $T^{r,*}$  is used in the  $r^{\text{th}}$  parallel server system. For the purpose of establishing the finiteness of  $J^*$  and the equality in the right side of (5.2), by appealing to Theorem 5.3 and the Skorokhod representation theorem, we may assume that a.s.,

$$\hat{Q}^{r,*} \rightarrow \tilde{Q}^* \text{ u.o.c. as } r \rightarrow \infty, \quad (9.25)$$

where  $\tilde{Q}^*$  is given by (4.15)–(4.17). Then for

$$\hat{H}^{r,*} \equiv h \cdot \hat{Q}^{r,*} \quad \text{and} \quad \tilde{H}^* \equiv h \cdot \tilde{Q}^* \quad (9.26)$$

we have

$$\hat{H}^{r,*} \rightarrow \tilde{H}^* \quad (m \times \mathbf{P})\text{-a.e. on } \mathbb{R}_+ \times \Omega, \quad (9.27)$$

where  $dm = \gamma e^{-\gamma t} dt$  on  $(\mathbb{R}_+, \mathcal{B}_+)$  and  $\mathcal{B}_+$  denotes the Borel  $\sigma$ -algebra on  $\mathbb{R}_+$ . Then, since  $(\mathbb{R}_+ \times \Omega, \mathcal{B}_+ \times \mathcal{F}, m \times \mathbf{P})$  is a probability space, to establish

$$\hat{J}^r(T^{r,*}) \equiv \mathbf{E} \left( \int_0^\infty e^{-\gamma t} \hat{H}^{r,*}(t) dt \right) \rightarrow J^* < \infty \text{ as } r \rightarrow \infty, \quad (9.28)$$

it suffices to show that

$$\limsup_{r \rightarrow \infty} \mathbf{E} \left( \int_0^\infty e^{-\gamma t} (\hat{H}^{r,*}(t))^2 dt \right) < \infty \quad (9.29)$$

which implies the required uniform integrability. From (8.36) we have

$$\hat{H}^{r,*} = h \cdot \hat{Q}^{r,*} \leq \left( \sum_{i \in \mathcal{I}} \frac{h_i}{y_i^r} \right) \hat{W}^{r,*}. \quad (9.30)$$

For each  $r \geq 1$  and  $t > 0$ , let

$$G^{r,t} = \left\{ Q_i^{r,*}(s) \leq 2L_i^r \text{ for all } s \in [0, r^2 t], i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*} \right\}. \quad (9.31)$$

By the definition of  $T^{r,*}$ , on  $G^{r,t}$  we have that (for  $r$  large enough that  $L_i^r \geq |\underline{\mathcal{J}}_i|$  for all  $i \in \underline{\mathcal{I}}_{k^*}$ ),  $\hat{I}_{k^*}^{r,*}$  can have a point of increase at  $s \in [0, t]$  only if  $\hat{Q}_i^r(s)$  is at or below the level  $L_i^r/r$  for all  $i \in \underline{\mathcal{I}}_{k^*}$ , which occurs only if

$$\hat{W}^{r,*}(s) \leq \left( \sum_{i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}} 2y_i^r L_i^r + \sum_{i \in \underline{\mathcal{I}}_{k^*}} y_i^r L_i^r \right) / r. \quad (9.32)$$

Thus, on  $G^{r,t}$ , it follows from (8.36) and an oscillation inequality for solutions of a perturbed Skorokhod problem (cf. the proof of Theorem 5.1 in [40]) that

$$\begin{aligned} z_{k^*}^r \hat{I}_{k^*}^{r,*}(t) &\leq \sup_{0 \leq s \leq t} \left( -y^r \cdot \hat{X}^{r,*}(s) - \sum_{k \neq k^*} z_k^r \hat{I}_k^{r,*}(s) \right) \\ &\quad + \left( \sum_{i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}} 2y_i^r L_i^r + \sum_{i \in \underline{\mathcal{I}}_{k^*}} y_i^r L_i^r \right) r^{-1} \\ &\leq \sup_{0 \leq s \leq t} |y^r \cdot \hat{X}^{r,*}(s)| + \sum_{k \neq k^*} z_k^r \hat{I}_k^{r,*}(t) \\ &\quad + \left( \sum_{i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}} 2y_i^r L_i^r + \sum_{i \in \underline{\mathcal{I}}_{k^*}} y_i^r L_i^r \right) r^{-1}, \end{aligned} \quad (9.33)$$

where we have used the fact that  $\hat{I}_k^{r,*}$ ,  $k \in \mathcal{K} \setminus \{k^*\}$ , is non-decreasing, to obtain the last inequality.

Since  $I_k^r(r^2 t) \leq r^2 t$  for all  $k \in \mathcal{K}$ , we have for  $t$  satisfying  $r^2 t < M^r$ ,

$$\mathbf{E} \left( (\hat{I}_k^{r,*}(t))^2 \right) \leq \left( \frac{M^r}{r} \right)^2 \text{ for all } k \in \mathcal{K}. \quad (9.34)$$

On the other hand, for  $r \geq r^*$ , for each  $t > 0$  satisfying  $r^2 t \geq Mr$ ,

$$\begin{aligned}
\mathbf{E}\left(\left(\hat{I}_{k^*}^{r,*}(t)\right)^2; \Omega \setminus G^{r,t}\right) &\leq r^2 t^2 \mathbf{P}\left(\Omega \setminus G^{r,t}\right) \\
&\leq r^2 t^2 \sum_{i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}} \mathbf{P}\left(Q_i^{r,*}(s) > 2L_i^r \text{ for some } s \in [0, r^2 t]\right) \\
&\leq r^2 t^2 \sum_{i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}} \mathbf{P}\left(\sup_{\tau_{i,0}^r \leq s \leq r^2 t} |R_i^r(s)| \geq L_i^r - |\underline{\mathcal{J}}_i|\right) \\
&\leq r^2 t^2 \sum_{i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}} p_{1,i}(r^2 t) \left(C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-C_{3,i}^{(4)} r^2 t)\right) \\
&\leq t p(r^2 t) \left(C^{(1)} \exp(-C^{(2)} L_0^r) + C^{(3)} \exp(-C^{(4)} r^2 t)\right), \tag{9.35}
\end{aligned}$$

by Theorem 7.7, where the  $p_{1,i}$ ,  $C_{1,i}^{(m)}$ ,  $m = 1, 2, 3, 4$ , are as in (I.1),  $p(s) = s \sum_{i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}} p_{1,i}(s)$ , for all  $s \geq 0$ ,  $C^{(1)} = \max\{C_{1,i}^{(1)} : i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}\}$ ,  $C^{(3)} = \max\{C_{1,i}^{(3)} : i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}\}$ ,  $C^{(2)} = \min\{C_{1,i}^{(2)} : i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}\}$ ,  $C^{(4)} = \min\{C_{1,i}^{(4)} : i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}\}$ . Here  $p$  is a polynomial (of degree at most  $\mathbf{I} + 1$ ) with non-negative coefficients and  $C^{(m)} > 0$ ,  $m = 1, 2, 3, 4$ ; the polynomial and constants do not depend on  $t$  or  $r$ .

By (8.21) and the fact that  $L_i^r$  is of order  $\log r$  for all  $i \in \mathcal{I}$ , there is  $r' \geq r^*$ , such that for all  $r \geq r'$ ,  $y^r \cdot (\lambda^r - \mathbf{R}^r x^*) r \leq y^* \cdot |\theta| + 1$  and  $\left(\sum_{i \in \mathcal{I} \setminus \underline{\mathcal{I}}_{k^*}} 2y_i^r L_i^r + \sum_{i \in \underline{\mathcal{I}}_{k^*}} y_i^r L_i^r\right) r^{-1} \leq 1$ . Then we have for  $r \geq r'$ , and  $t > 0$  satisfying  $r^2 t \geq Mr$ , using the inequality  $\left(\sum_{i=1}^n |x_i|\right)^2 \leq n \sum_{i=1}^n x_i^2$  repeatedly, that

$$\begin{aligned}
\mathbf{E}\left(\left(\hat{I}_{k^*}^{r,*}(t)\right)^2\right) &= \mathbf{E}\left(\left(\hat{I}_{k^*}^{r,*}(t)\right)^2; G^{r,t}\right) + \mathbf{E}\left(\left(\hat{I}_{k^*}^{r,*}(t)\right)^2; \Omega \setminus G^{r,t}\right) \\
&\leq \frac{5}{(z_{k^*}^r)^2} \left\{ \mathbf{I} \sum_{i \in \mathcal{I}} (y_i^r)^2 \mathbf{E}\left(\sup_{0 \leq s \leq t} (\hat{A}_i^r(s))^2\right) \right. \\
&\quad \left. + \mathbf{J} \sum_{j \in \mathcal{J}} (y_{i(j)}^r)^2 \mathbf{E}\left(\sup_{0 \leq s \leq t} (\hat{S}_j^r(T_j^{r,*}(s)))^2\right) + ((y^* \cdot |\theta| + 1)t)^2 \right. \\
&\quad \left. + (\mathbf{K} - 1) \sum_{k \neq k^*} (z_k^r)^2 \mathbf{E}\left(\left(\hat{I}_k^{r,*}(t)\right)^2\right) + 1 \right\} \\
&\quad + t p(r^2 t) \left(C^{(1)} \exp(-C^{(2)} L_0^r) + C^{(3)} \exp(-C^{(4)} r^2 t)\right), \tag{9.36}
\end{aligned}$$

where the first term to the right of the equality sign is controlled via (9.33), and the second term is controlled via (9.35). Using the fact that  $\exp(-C^{(2)} L_0^r) \leq r^{-C^{(2)} c}$  (since  $L_0^r = \lceil c \log r \rceil$ ), and the fact that for any polynomial  $q$  and  $d > 0$ ,  $q(x)e^{-dx}$  is bounded for  $x \in \mathbb{R}_+$ , we have that there is a constant  $c_2 \geq c_1$  (independent of  $t$  and  $r$ ) and  $r'' \geq r'$  such that for each fixed  $c \geq c_2$  and all  $r \geq r''$ , the last term in (9.36) is bounded by a polynomial in  $t$ , independent of  $r$ . (The constant  $c_1$  was used in the proof of Theorem 7.1.)

Then, using (9.30), (8.36), (4.7), (3.11), (9.34), (9.36), and the fact that  $Mr = o(r)$ , for  $c$  sufficiently large (chosen independently of  $t$  and  $r$ ), we see that to prove (9.29) it suffices to show that, as functions of  $t$ , the following are all in a bounded subset of  $L^1(m) \equiv L^1(\mathbb{R}_+, \mathcal{B}_+, m)$  for  $r$  sufficiently large:

$$\mathbf{E}\left(\sup_{0 \leq s \leq t} \left(\hat{A}_i^r(s)\right)^2\right), \quad \mathbf{E}\left(\sup_{0 \leq s \leq t} \left(\hat{S}_j^r(\bar{T}_j^{r,*}(s))\right)^2\right), \quad \mathbf{E}\left(\left(\hat{I}_k^{r,*}(t)\right)^2\right), \tag{9.37}$$

for all  $i \in \mathcal{I}$ ,  $j \in \mathcal{J}$ ,  $k \in \mathcal{K} \setminus \{k^*\}$ . The bounds for the first two expectations can be obtained as in Section 9 in [3]. For the third bound, fix  $k \in \mathcal{K} \setminus \{k^*\}$ , and let  $i$  be the transition buffer that is immediately above  $k$ . Let  $\epsilon > 0$  and choose  $r''' \geq r''$  such that for all  $i \in \mathcal{I}$  and  $r \geq r'''$ ,  $r\epsilon \geq t_i^r$  (cf. (7.8)). Fix  $r \geq r'''$ . By (9.34), we need only consider  $t > 0$  satisfying  $r^2 t \geq M^r$ . Since  $I_k^{r,*}(r^2 t) \leq r^2 t$ ,

$$\begin{aligned}
\mathbf{E}\left(\left(\hat{I}_k^{r,*}(t)\right)^2\right) &= \int_0^\infty \mathbf{P}\left(\left(\hat{I}_k^{r,*}(t)\right)^2 > s\right) ds \\
&= \int_0^{r^2 t^2} \mathbf{P}\left(I_k^{r,*}(r^2 t) > r\sqrt{s}\right) ds \\
&\leq \int_0^{r^2 t^2} \left\{ \mathbf{P}\left(I_k^{r,*}(\tau_{i,0}^r) > r\sqrt{s}\right) \right. \\
&\quad \left. + \mathbf{P}\left(\inf_{\tau_{i,0}^r \leq u \leq r^2 t} R_i^r(u) \leq -L_i^r + |\underline{\mathcal{J}}_i|\right) \right\} ds \\
&\leq \epsilon^2 + \int_{\epsilon^2}^{r^2 t^2} \mathbf{P}\left(I_k^{r,*}(\tau_{i,0}^r) > r\epsilon\right) ds \\
&\quad + r^2 t^2 \mathbf{P}\left(\inf_{\tau_{i,0}^r \leq u \leq r^2 t} R_i^r(u) \leq -L_i^r + |\underline{\mathcal{J}}_i|\right) \\
&\leq \epsilon^2 + t \left\{ r^2 t \mathbf{P}\left(I_k^{r,*}(\tau_{i,0}^r) \geq t_i^r\right) \right. \\
&\quad \left. + r^2 t \mathbf{P}\left(\inf_{\tau_{i,0}^r \leq u \leq r^2 t} R_i^r(u) \leq -L_i^r + |\underline{\mathcal{J}}_i|\right) \right\} \\
&\leq \epsilon^2 + t \left\{ r^2 t p_{2,i}(r^2 t) \left( C_{2,i}^{(1)} \exp(-C_{2,i}^{(2)} L_0^r) + C_{2,i}^{(3)} \exp(-C_{2,i}^{(4)} r^2 t) \right) \right. \\
&\quad \left. + r^2 t p_{1,i}(r^2 t) \left( C_{1,i}^{(1)} \exp(-C_{1,i}^{(2)} L_0^r) + C_{1,i}^{(3)} \exp(-C_{1,i}^{(4)} r^2 t) \right) \right\}, \quad (9.38)
\end{aligned}$$

by (I.1)–(I.2), if  $i \neq i^*$ , or by (III.1)–(III.2), if  $k \in \underline{\mathcal{K}}_{i^*}$ , which all hold by Theorem 7.7. In the first inequality in (9.38), we have used that fact that  $I_k^{r,*}(r^2 t) - I_k^{r,*}(\tau_{i,0}^r) = 0$  if  $\inf_{\tau_{i,0}^r \leq u \leq r^2 t} R_i^r(u) > -L_i^r + |\underline{\mathcal{J}}_i|$ .

Since  $\exp(-C_{m,i}^{(2)} L_0^r) \leq r^{-C_{m,i}^{(2)} c}$ ,  $m = 1, 2$ , and for any polynomial  $q$  and  $d > 0$ ,  $q(x)e^{-dx}$  is bounded for  $x \in \mathbb{R}_+$ , it follows that there is a constant  $c_3 \geq c_2$  (independent of  $t$  and  $r$ ),  $r^{**} \geq r'''$ , such that for each fixed  $c \geq c_3$ , and all  $r \geq r^{**}$ , the right member above can be bounded by a polynomial in  $t$  (not depending on  $r$ ). Combining (9.34) with the above, we conclude that  $\{\mathbf{E}((\hat{I}_k^{r,*}(\cdot))^2) : r \geq r^{**}\}$  is a bounded sequence of functions in  $L^1(m)$ , provided that  $c$  is fixed and sufficiently large.  $\square$

## References

- [1] B. Ata and S. Kumar. Heavy traffic analysis of open processing networks with complete resource pooling: asymptotic optimality of discrete review policies. *Annals of Applied Probability*, **15** (2005), 331–391.
- [2] S. L. Bell. *Dynamic Scheduling of a Parallel Server System in Heavy Traffic with Complete Resource Pooling: Asymptotic Optimality of a Threshold Policy*. Ph.D. dissertation, Department of Mathematics, University of California, San Diego, 2003.
- [3] S. L. Bell and R. J. Williams. Dynamic scheduling of a system with two parallel servers in heavy traffic with resource pooling: asymptotic optimality of a threshold policy. *Annals of Applied Probability*, **11** (2001), 608–649.

- [4] P. Billingsley. *Convergence of Probability Measures*. John Wiley and Sons, New York, 1968.
- [5] M. Bramson. State space collapse with applications to heavy traffic limits for multiclass queueing networks. *Queueing Systems*, **30** (1998), 89–148.
- [6] M. Bramson and R. J. Williams. Two workload properties for Brownian networks. *Queueing Systems*, **45** (2003), 191–221.
- [7] A. Budhiraja and A. P. Ghosh, A large deviations approach to asymptotically optimal control of crisscross network in heavy traffic. *Annals of Applied Probability*, **15** (2005), 1887–1935.
- [8] P. B. Chevalier and L. Wein. Scheduling networks of queues: heavy traffic analysis of a multi-station closed network. *Operations Research*, **41** (1993), 743–758.
- [9] K. L. Chung and R. J. Williams. *Introduction to Stochastic Integration*. 2nd edition, Birkhäuser, Boston, 1990.
- [10] S. N. Ethier and T. G. Kurtz. *Markov Processes: Characterization and Convergence*. Wiley, New York, 1986.
- [11] J. M. Harrison. *Brownian Motion and Stochastic Flow Systems*. John Wiley and Sons, New York, 1985.
- [12] J. M. Harrison. Brownian models of queueing networks with heterogeneous customer populations. In *Stochastic Differential Systems, Stochastic Control Theory and Their Applications, IMA Volume 10*, W. Fleming and P. L. Lions (eds.), Springer Verlag, New York, 1988, pp. 147–186.
- [13] J. M. Harrison. Balanced fluid models of multiclass queueing networks: a heavy traffic conjecture. In *Stochastic Networks*, F. P. Kelly and R. J. Williams (eds.), Springer-Verlag, 1995, pp. 1–20.
- [14] J. M. Harrison. The BIGSTEP approach to flow management in stochastic processing networks. In *Stochastic Networks: Theory and Applications*, F. P. Kelly, S. Zachary and I. Ziedins (eds.), Oxford University Press, 1996, pp. 57–90.
- [15] J. M. Harrison. Heavy traffic analysis of a system with parallel servers: asymptotic optimality of discrete-review policies. *Annals of Applied Probability*, **8** (1998), 822–848.
- [16] J. M. Harrison. Brownian models of open processing networks: canonical representation of workload. *Annals of Applied Probability*, **10** (2000), 75–103. Correction: **13** (2003), 390–393.
- [17] J. M. Harrison. A broader view of Brownian networks. *Annals of Applied Probability*, **13** (2003), 1119–1150.
- [18] J. M. Harrison and M. J. López. Heavy traffic resource pooling in parallel-server systems. *Queueing Systems*, **33** (1999), 339–368.
- [19] J. M. Harrison and J. A. Van Mieghem. Dynamic control of Brownian networks: state space collapse and equivalent workload formulations. *Annals of Applied Probability*, **7** (1997), 747–771.
- [20] J. M. Harrison and L. Wein. Scheduling networks of queues: heavy traffic analysis of a simple open network. *Queueing Systems*, **5** (1989), 265–280.

- [21] J. M. Harrison and L. Wein. Scheduling networks of queues: heavy traffic analysis of a two-station closed network. *Operations Research*, **38** (1990), 1052–1064.
- [22] D. L. Iglehart and W. Whitt. The equivalence of functional central limit theorems for counting processes and associated partial sums. *Ann. Math. Statist.*, **42** (1971), 1372–1378.
- [23] W. C. Jordan and C. Graves. Principles on the benefits of manufacturing process flexibility. *Management Science*, **41** (1995), 577–594.
- [24] F. P. Kelly and C. N. Laws. Dynamic routing in open queueing networks: Brownian models, cut constraints and resource pooling. *Queueing Systems*, **13** (1993), 47–86.
- [25] H. J. Kushner and Y. N. Chen. Optimal control of assignment of jobs to processors under heavy traffic. *Stochastics and Stochastics Rep.*, **68** (2000), no. 3-4, pp. 177–228.
- [26] H. J. Kushner and P. Dupuis. *Numerical Methods for Stochastic Control Problems in Continuous Time*. Springer-Verlag, New York, 1992.
- [27] H. J. Kushner and L. F. Martins. Numerical methods for stochastic singular control problems. *SIAM J. Control and Optimization* **29** (1991), 1443–1475.
- [28] C. N. Laws. *Dynamic Routing in Queueing Networks*. Ph.D. dissertation, Statistical Laboratory, University of Cambridge, U.K., 1990.
- [29] C. N. Laws. Resource pooling in queueing networks with dynamic routing. *Advances in Applied Probability*, **24** (1992), 699–726.
- [30] C. N. Laws and G. M. Louth. Dynamic scheduling of a four-station queueing network. *Probab. Engrg. Inform. Sci.*, **4** (1990), 131–156.
- [31] A. Mandelbaum and A. L. Stolyar. Scheduling flexible servers with convex delay costs: heavy-traffic optimality of the generalized  $c\mu$ -rule. *Operations Research*, **52** (2004), 836–855.
- [32] S. P. Meyn. Dynamic safety-stocks for asymptotic optimality in stochastic networks. *Queueing Systems*, **50** (2005), 255–297.
- [33] T. Osogami, M. Harchol-Balter, A. Scheller-Wolf, and L. Zhang. Exploring threshold-based policies for load sharing. *Proceedings of 42nd Annual Allerton Conference on Communication, Control and Computing*, University of Illinois, Urbana-Champaign, October 2004.
- [34] W. P. Peterson. A heavy traffic limit theorem for networks of queues with multiple customer types. *Mathematics of Operations Research*, **16** (1991), 90–118.
- [35] M. Squillante, C. H. Xia, D. Yao, and L. Zhang. Threshold-based priority policies for parallel-server systems with affinity scheduling. *Proc. IEEE American Control Conference* (2001), 2992–2999.
- [36] A. L. Stolyar. Maxweight scheduling in a generalized switch: state space collapse and workload minimization in heavy traffic. *Annals of Applied Probability*, **14** (2004), 1–53.
- [37] Y. C. Teh. *Threshold Routing Strategies for Queueing Networks*. D. Phil. thesis, University of Oxford, 1999.
- [38] Y. C. Teh and A. R. Ward. Critical thresholds for dynamic routing in queueing networks. *Queueing Systems*, **42** (2002), 297–316.

- [39] L. Wein. Scheduling networks of queues: heavy traffic analysis of a two-station network with controllable inputs. *Operations Research*, **38** (1990), 1065–1078.
- [40] R. J. Williams. An invariance principle for semimartingale reflecting Brownian motions in an orthant. *Queueing Systems*, **30** (1998), 5–25.
- [41] R. J. Williams. Diffusion approximations for open multiclass queueing networks: sufficient conditions involving state space collapse. *Queueing Systems*, **30** (1998), 27–88.
- [42] R. J. Williams. On dynamic scheduling of a parallel server system with complete resource pooling. In *Analysis of Communication Networks: Call Centres, Traffic and Performance*, D. R. McDonald and S. R. E. Turner (eds.), American Mathematical Society, Providence, RI, 2000, 49–71.