

Utility maximizing load balancing policies

Online companion

Diego Goldsztajn, Sem C. Borst

Eindhoven University of Technology, 5612 AZ Eindhoven, Netherlands, d.e.goldsztajn@tue.nl, s.c.borst@tue.nl

Johan S.H. van Leeuwen

Tilburg University, 5037 AB Tilburg, Netherlands, j.s.h.vanleeuwen@uvt.nl

This online supplement contains the proofs of some of the technical results stated in the body of the paper. Specifically, various auxiliary results are proved in Appendix A, a relative compactness result is established in Appendix B and several results pertaining to the limiting behavior of SLTA are proved in Appendix C.

Appendix A: Auxiliary results

Construction of Section 5.1.2. First we define the functions η_k . For this purpose, suppose that $q \in Q_n$ and $r \geq 1$ are the current values of q_n and r_n , respectively. Then

$$\mathcal{G}(q, r) := \begin{cases} \{(i, j) \in \mathcal{I}_+ : (i, j) \triangleright (i_r, j_r)\} & \text{if } r > 1, \\ \emptyset & \text{if } r = 1, \end{cases}$$

is the set of the coordinates that could belong to a server pool with a green token; i.e., a server pool of class i with exactly $j - 1$ tasks has a green token if and only if $(i, j) \in \mathcal{G}(q, r)$. We partition the latter set of coordinates as follows:

$$\begin{aligned} \mathcal{G}_1(q, r) &:= \mathcal{G}(q, r) \cap \{(i, j) \in \mathcal{I}_+ : i \neq i_{r-1}\}, \\ \mathcal{G}_2(q, r) &:= \mathcal{G}(q, r) \cap \{(i, j) \in \mathcal{I}_+ : i = i_{r-1}\}. \end{aligned}$$

Let $\delta_q(i, j) := q(i, j - 1) - q(i, j)$ denote the fraction of server pools which are of class i and have exactly $j - 1$ tasks when the occupancy state is q , and define

$$I(q, i, j) := \left[1 - \sum_{(k, l) \preceq (i, j)} \delta_q(k, l), 1 - \sum_{(k, l) \triangleleft (i, j)} \delta_q(k, l) \right) \quad \text{for all } (i, j) \in \mathcal{I}_+.$$

These intervals form a partition of $[0, 1)$ such that the length of $I(q, i, j)$ is the fraction of server pools which are of class i and have exactly $j - 1$ tasks. The length of $I(q, i, j)$ is equal to the probability of picking a server pool with coordinates (i, j) uniformly at random.

Note that $G_1(q, r) := \sum_{(k, l) \in \mathcal{G}_1(q, r)} \delta_q(k, l)$ is the fraction of server pools of class $i \neq i_{r-1}$ with a green token. If the latter quantity is positive, then we let

$$G_1(q, r, i, j) := \left[\frac{1 - \sum_{(k, l) \preceq (i, j), (k, l) \in \mathcal{G}_1(q, r)} \delta_q(k, l)}{G_1(q, r)}, \frac{1 - \sum_{(k, l) \triangleleft (i, j), (k, l) \in \mathcal{G}_1(q, r)} \delta_q(k, l)}{G_1(q, r)} \right).$$

These intervals yield another partition of $[0, 1]$, where the length of $G_1(q, r, i, j)$ is the fraction of server pools which are of class i and have exactly $j - 1$ tasks, but only among those server pools which are of class $i \neq i_{r-1}$ and have a green token. If $G_1(q, r) = 0$, then we define $G_1(q, r, i, j) := \emptyset$ for all $(i, j) \in \mathcal{I}_+$. Sets $G_2(q, r, i, j)$ are defined similarly.

The functions η_k are defined as follows:

$$\eta_k(q, r, i, j) := \begin{cases} \mathbb{1}_{\{U_k \in G_1(q, r, i, j)\}} & \text{if } (i, j) \in \mathcal{G}_1(q, r), \\ \mathbb{1}_{\{G_1(q, r) = 0, U_k \in G_2(q, r, i, j)\}} & \text{if } (i, j) \in \mathcal{G}_2(q, r), \\ \mathbb{1}_{\{G_1(q, r) = G_2(q, r) = 0, I(q, i_r, j_r) \neq \emptyset\}} & \text{if } (i, j) = (i_r, j_r), \\ \mathbb{1}_{\{G_1(q, r) = G_2(q, r) = 0, I(q, i_r, j_r) = \emptyset, U_k \in I(q, i, j)\}} & \text{otherwise.} \end{cases}$$

Observe that $\eta_k(q, r, i, j) \in \{0, 1\}$ for all k and $(i, j) \in \mathcal{I}_+$, and that for a fixed k there exists a unique $(i, j) \in \mathcal{I}_+$ such that $\eta_k(q, r, i, j) = 1$. Also, if (q, r) is the value of $(\mathbf{q}_n, \mathbf{r}_n)$ when the k^{th} task arrives to the system, then $\eta_k(q, r, i, j) = 1$ if and only if this task has to be sent to a server pool of class i with exactly $j - 1$ tasks.

It only remains to define the sets $I_{n,k}$ and $D_{n,k}$ that appear in (20b). The former is the set of those $\omega \in \Omega$ that satisfy:

$$\begin{aligned} \mathbf{q}(\tau_{n,k}^-(\omega), i, j) &= \alpha_n(i) \quad \text{for all } (i, j) \triangleright (i_{\mathbf{r}(\tau_{n,k}^-(\omega))}, j_{\mathbf{r}(\tau_{n,k}^-(\omega))}), \\ n\mathbf{q}(\tau_{n,k}^-(\omega), i_{\mathbf{r}(\tau_{n,k}^-(\omega))}, j_{\mathbf{r}(\tau_{n,k}^-(\omega))}) &\geq n\alpha_n(i_{\mathbf{r}(\tau_{n,k}^-(\omega))}) - 1. \end{aligned}$$

If the state of the system is (\mathbf{q}, \mathbf{r}) , then the first condition means that there are no green tokens right before the k^{th} task arrives, and the second condition means that the number of yellow tokens is at most one. Finally, $D_{n,k}$ is the set of those $\omega \in \Omega$ such that

$$\begin{aligned} \mathbf{r}(\tau_{n,k}^-(\omega)) &> 1, \\ n - \sum_{i=1}^m n\mathbf{q}(\tau_{n,k}^-(\omega), i, \ell_i(\mathbf{r}(\tau_{n,k}^-(\omega)))) &\geq n\beta_n, \\ \mathbf{q}(\tau_{n,k}^-(\omega), i_{\mathbf{r}(\tau_{n,k}^-(\omega))-1}, j_{\mathbf{r}(\tau_{n,k}^-(\omega))-1}) &< \alpha_n(i_{\mathbf{r}(\tau_{n,k}^-(\omega))-1}). \end{aligned}$$

The second condition means that the number of green tokens is larger than or equal to $n\beta_n$ right before the k^{th} task arrives. The last condition means that at least one of these tokens is of class i_{r-1} . \square

Proof of Lemma 2. Fix $T \geq 0$ and let $M_T := \lambda T(e-1) + 1$. The set

$$E_T := \{\omega \in \Omega : \mathcal{N}_n^\lambda(\omega, T) \leq nM_T \text{ for all large enough } n\}$$

has probability one as a consequence of the Borel-Cantelli lemma and the Chernoff bound

$$\mathbb{P}(\mathcal{N}_n^\lambda(T) > nM_T) \leq \frac{e^{\lambda T(e-1)n}}{e^{M_T n}} = e^{-n}.$$

As a result, the following set also has probability one:

$$\Gamma := \bigcap_{T \in \mathbb{N}} (E_T \cap \Gamma_\infty).$$

Note that for each $\omega \in \Gamma$ and $T \geq 0$ there exists $n_T^1(\omega)$ such that $\mathcal{N}_n^\lambda(\omega, T) \leq nM_T$ for all $n \geq n_T^1(\omega)$. Fix $T \geq 0$, $\omega \in \Gamma$ and $0 < \varepsilon < \min\{\alpha(i) : 1 \leq i \leq m\}$. Next we establish that there exist $j_T(\omega)$ and $n_T(\omega)$ such that

$$\mathbf{q}_n(\omega, t, i, j_T(\omega)) \leq \varepsilon < \alpha_n(i) \quad \text{for all } t \in [0, T], \quad 1 \leq i \leq m \quad \text{and} \quad n \geq n_T(\omega). \quad (51)$$

If the load balancing policy is JLMU, then (8) and the above statement imply that $\sigma_j(\mathbf{q}_n(\omega, t)) \leq j_T(\omega)$ for all $t \in [0, T]$ and $n \geq n_T(\omega)$, which in turn implies that the claim of the lemma holds. Suppose instead that SLTA is used and let

$$r_T(\omega) := \min \{k \geq 1 : (i_k, j_k) \preceq (i, j_T(\omega)) \text{ for all } i\}.$$

It follows from (23) that $\mathbf{r}_n(\omega, t) \leq \max\{r_T(\omega), R(\omega)\} =: R_T(\omega)$ for all $t \in [0, T]$ and $n \geq n_T(\omega)$. Furthermore, the latter property and Remark 1 imply that $\mathcal{A}_n(\omega, t, i, j) = 0$ for all $t \in [0, T]$, $n \geq n_T(\omega)$ and $(i, j) \prec (i_{R_T(\omega)+1}, j_{R_T(\omega)+1})$. Therefore, the claim of the lemma also holds for SLTA.

The following arguments apply both for JLMU and SLTA. Below we omit ω from the notation for brevity. Since $q_0 \in \ell_1$, there exists j_0 such that $q_0(i, j_0) \leq \varepsilon/4$ for all i , and by (9) and (22a), there exists $n_T \geq n_T^1$ such that

$$\mathbf{q}_n(0, i, j_0) \leq \frac{\varepsilon}{2} < \varepsilon < \alpha_n(i) \quad \text{for all } 1 \leq i \leq m \quad \text{and} \quad n \geq n_T.$$

Define $j_T := j_0 + \lceil 2M_T/\varepsilon \rceil - 1$. Also, fix $t \in [0, T]$, $1 \leq i \leq m$ and $n \geq n_T$. We have

$$\sum_{j=j_0}^{j_T} n \left[\mathbf{q}_n(t, i, j) - \frac{\varepsilon}{2} \right] \leq \sum_{j=j_0}^{j_T} n [\mathbf{q}_n(t, i, j) - \mathbf{q}_n(0, i, j)] \leq \sum_{j=j_0}^{j_T} n \mathcal{A}_n(t, i, j) \leq nM_T.$$

For the first inequality, observe that $\mathbf{q}_n(0, i, j) \leq \mathbf{q}_n(0, i, j_0) \leq \varepsilon/2$ for all $j \geq j_0$, and for the last inequality, note that $\mathcal{N}_n^\lambda(t) \leq \mathcal{N}_n^\lambda(T) \leq nM_T$ because $n \geq n_T^1$. Since $\mathbf{q}_n(t, i, j)$ is non-increasing in j , we conclude that

$$\mathbf{q}_n(t, i, j_T) \leq \frac{1}{j_T - j_0 + 1} \sum_{j=j_0}^{j_T} \mathbf{q}_n(t, i, j) \leq \frac{M_T + \frac{\varepsilon}{2}(j_T - j_0 + 1)}{j_T - j_0 + 1} < \varepsilon < \alpha_n(i).$$

This completes the proof. \square

Proof of Lemma 3. Consider a function $\mathbf{x} \in D[0, \infty)$ such that $\mathbf{x}(0) \geq 0$. Then there exist unique $\mathbf{y}, \mathbf{z} \in D[0, \infty)$ such that the following statements hold.

- (a') $\mathbf{z}(t) = \mathbf{x}(t) + \mathbf{y}(t) \geq 0$ for all $t \geq 0$.
- (b') \mathbf{y} is non-decreasing and $\mathbf{y}(0) = 0$.
- (c') \mathbf{y} is flat off $\{t \geq 0 : \mathbf{z}(t) = 0\}$, i.e., $\dot{\mathbf{y}}(t) \mathbb{1}_{\{\mathbf{z}(t) > 0\}} = 0$ almost everywhere.

The map (Ψ, Φ) such that $\Psi(\mathbf{x}) = \mathbf{y}$ and $\Phi(\mathbf{x}) = \mathbf{z}$ is called the one-dimensional Skorokhod mapping with lower reflecting barrier at zero and satisfies

$$\Psi(\mathbf{x})(t) = \sup_{s \in [0, t]} [-\mathbf{x}(s)]^+ \quad \text{and} \quad \Phi(\mathbf{x})(t) = \mathbf{x}(t) + \Psi(\mathbf{x})(t).$$

Also, if $\mathbf{x}, \mathbf{y} \in D[0, \infty)$ satisfy $\mathbf{x}(0), \mathbf{y}(0) \geq 0$, then for each $T \geq 0$, we have

$$\|\Psi(\mathbf{x}) - \Psi(\mathbf{y})\|_T \leq \|\mathbf{x} - \mathbf{y}\|_T \quad \text{and} \quad \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\|_T \leq 2\|\mathbf{x} - \mathbf{y}\|_T.$$

The latter definition and properties of the Skorokhod mapping with lower reflecting barrier at zero can be found in (Chen and Yao 2013, Theorem 6.1) and (Karatzas and Shreve 2014, Lemma 6.14).

Suppose \mathbf{x} is as in the statement of the lemma. Define \mathbf{y} and \mathbf{z} as in (24); i.e.,

$$\mathbf{y}(t) := \sup_{s \in [0, T]} [\mathbf{x}(s) - \alpha]^+ = \Psi(\alpha - \mathbf{x})(t) \quad \text{and} \quad \mathbf{z}(t) := \mathbf{x}(t) - \mathbf{y}(t) = \alpha - \Phi(\alpha - \mathbf{x})(t).$$

Properties (b') and (c') imply (b) and (c). Also, (a) holds since $\mathbf{z} = \alpha - \Phi(\alpha - \mathbf{x}) \leq \alpha$ by (a') and $\mathbf{z} = \mathbf{x} - \mathbf{y}$ by definition.

Conversely, suppose that \mathbf{y} and \mathbf{z} satisfy (a)-(c). Then $\alpha - \mathbf{z} = \alpha - \mathbf{x} + \mathbf{y} \geq 0$ and \mathbf{y} satisfies (b') and (c') with \mathbf{z} replaced by $\alpha - \mathbf{z}$, hence we conclude that $\mathbf{y} = \Psi(\alpha - \mathbf{x})$ and $\alpha - \mathbf{z} = \Phi(\alpha - \mathbf{x})$. Therefore, \mathbf{y} and \mathbf{z} are as in (24).

The Lipschitz properties of Ψ_α and Φ_α follow from the Lipschitz properties of Ψ and Φ and the identities $\Psi_\alpha(\mathbf{x}) = \Psi(\alpha - \mathbf{x})$ and $\Phi_\alpha(\mathbf{x}) = \alpha - \Phi(\alpha - \mathbf{x})$. \square

Proof of Lemma 8. By assumption and Tonelli's theorem,

$$\begin{aligned} \sum_{x, y \in S_n} \pi_n(x) |A_n(x, y)| f(y) &= \sum_{x \in S_n} \pi_n(x) \sum_{y \in S_n} |A_n(x, y)| f(y) \\ &= E \left[\sum_{y \in S_n} |A_n(x_n, y)| f(y) \right] < \infty. \end{aligned}$$

Therefore, we may use Fubini's theorem, which yields

$$0 = \sum_{y \in S_n} f(y) \sum_{x \in S_n} \pi_n(x) A_n(x, y) = \sum_{x \in S_n} \pi_n(x) \sum_{y \in S_n} A_n(x, y) f(y) = E[A_n f(x_n)].$$

The first equality follows from the stationarity of π_n . \square

Appendix B: Relative compactness

Consider the metric

$$d(x, y) := \sum_{(i, j) \in \mathcal{I}} \frac{\min\{|x(i, j) - y(i, j)|, 1\}}{2^j} \quad \text{for all } x, y \in \mathbb{R}^{\mathcal{I}},$$

which induces the product topology on $\mathbb{R}^{\mathcal{I}}$. For each $T \geq 0$, let $D_{\mathbb{R}^{\mathcal{I}}}[0, T]$ denote the space of all càdlàg functions on $[0, T]$ with values in $\mathbb{R}^{\mathcal{I}}$ and equip this space with the uniform metric, defined by

$$\varrho_T(\mathbf{x}, \mathbf{y}) := \sup_{t \in [0, T]} d(\mathbf{x}(t), \mathbf{y}(t)) \quad \text{for all } \mathbf{x}, \mathbf{y} \in D_{\mathbb{R}^{\mathcal{I}}}[0, T].$$

The topology of uniform convergence over compact sets on $D_{\mathbb{R}^{\mathcal{I}}}[0, \infty)$ is compatible with the metric

$$\varrho_\infty(\mathbf{x}, \mathbf{y}) := \sum_{k=0}^{\infty} \frac{\min\{\varrho_k(\mathbf{x}, \mathbf{y}), 1\}}{2^k} \quad \text{for all } \mathbf{x}, \mathbf{y} \in D_{\mathbb{R}^{\mathcal{I}}}[0, \infty).$$

Since $D_{\mathbb{R}^{\mathcal{I}}}[0, \infty)$ is metrizable, a subset of $D_{\mathbb{R}^{\mathcal{I}}}[0, \infty)$ is relatively compact if and only if every sequence within the subset has a convergent subsequence. The same holds for the metric spaces $D_{\mathbb{R}^{\mathcal{I}}}[0, T]$ defined above.

The following arguments apply both for JLMU and SLTA. Fix some arbitrary time horizon $T \geq 0$. In order to prove Proposition 2, we first demonstrate that the sequences $\{\mathcal{A}_n : n \geq 1\}$, $\{\mathcal{D}_n : n \geq 1\}$ and $\{\mathbf{q}_n : n \geq 1\}$ are relatively compact in the space $D_{\mathbb{R}^{\mathcal{I}}}[0, T]$ with probability one, and that every convergent subsequence has a componentwise Lipschitz limit. Here we are actually referring to the restrictions to the interval $[0, T]$ of the processes \mathcal{A}_n , \mathcal{D}_n and \mathbf{q}_n , but the symbols $|_{[0, T]}$ have been omitted in order to maintain a lighter notation; this is also done in the sequel when there is no risk of confusion. The latter properties are established using the methodological framework Bramson (1998). The following result defines a set of probability one where the relative compactness holds.

LEMMA 1. *There exists a set of probability one $\Gamma_T \subset \Gamma_0$ where:*

$$\lim_{n \rightarrow \infty} \|\mathbf{q}_n(0) - q_0\|_1 = 0, \quad (52a)$$

$$\lim_{n \rightarrow \infty} \sup_{t \in [0, T]} n^\gamma \left| \frac{1}{n} \mathcal{N}_n^\lambda(t) - \lambda t \right| = 0, \quad (52b)$$

$$\lim_{n \rightarrow \infty} \sup_{t \in [0, \mu_j T]} n^\gamma \left| \frac{1}{n} \mathcal{N}_{(i,j)}(nt) - t \right| = 0 \quad \text{for all } (i, j) \in \mathcal{I}_+, \quad (52c)$$

for all $\gamma \in [0, 1/2)$. Also, if the load balancing policy is SLTA, then apart from the latter properties, there exists a random variable $R \geq 1$ such that

$$\mathbf{r}_n(0) \leq R \quad \text{and} \quad \mathbf{q}_n(0, i, j) < \alpha_n(i) \quad \text{for all } (i, j) \in (i_{\mathbf{r}_n(0)}, j_{\mathbf{r}_n(0)}) \quad \text{and } n.$$

Proof. This result is a straightforward consequence of (9), (15) and the refined strong law of large numbers for the Poisson process proven in (Goldsztajn et al. 2021, Lemma 4). \square

We fix some $\omega \in \Gamma_T$, which we omit from the notation for brevity. Next we prove that the sequences $\{\mathcal{A}_n : n \geq 1\}$ and $\{\mathcal{D}_n : n \geq 1\}$ are relatively compact subsets of $D_{\mathbb{R}^{\mathcal{I}}}[0, T]$ and that the limit of every convergent subsequence is a function with Lipschitz coordinates. If JLMU is used, it then follows from (19) and (52a) that the sequence $\{\mathbf{q}_n : n \geq 1\}$ has the same properties. If SLTA is used, then we need to invoke (21a) instead of (19).

The following characterization of relative compactness with respect to ϱ_T will be useful. Let $D[0, T]$ denote the space of real càdlàg functions on $[0, T]$, with the uniform norm:

$$\|\mathbf{x}\|_T = \sup_{t \in [0, T]} |\mathbf{x}(t)| \quad \text{for all } \mathbf{x} \in D[0, T].$$

Observe that a sequence of functions $\mathbf{x}_n \in D_{\mathbb{R}^{\mathcal{I}}}[0, T]$ converges to $\mathbf{x} \in D_{\mathbb{R}^{\mathcal{I}}}[0, T]$ with respect to ϱ_T if and only if $\mathbf{x}_n(i, j)$ converges to $\mathbf{x}(i, j)$ with respect to the uniform norm for all $(i, j) \in \mathcal{I}$. Moreover, we have the following proposition.

PROPOSITION 1. *The sequence $\{\mathbf{x}_n : n \geq 1\}$ is relatively compact in $D_{\mathbb{R}^{\mathcal{I}}}[0, T]$ if and only if $\{\mathbf{x}_n(i, j) : n \geq 1\}$ is a relatively compact subset of $D[0, T]$ for all $(i, j) \in \mathcal{I}$.*

Proof. We only need to prove the converse, so assume that $\{\mathbf{x}_n(i, j) : n \geq 1\}$ is relatively compact for all $(i, j) \in \mathcal{I}$. Given an increasing sequence $\mathcal{K} \subset \mathbb{N}$, we need to establish that there exists a subsequence of $\{\mathbf{x}_k : k \in \mathcal{K}\}$ that converges with respect to ϱ_T . For this purpose, we may construct a family of sequences $\{\mathcal{K}_j : j \geq 0\}$ with the following properties.

- (a) $\mathcal{K}_{j+1} \subset \mathcal{K}_j \subset \mathcal{K}$ for all $j \geq 0$.
- (b) $\{\mathbf{x}_k(i, j) : k \in \mathcal{K}_j\}$ has a limit $\mathbf{x}(i, j) \in D[0, T]$ for each $(i, j) \in \mathcal{I}$.

Define $\{k_l : l \geq 1\} \subset \mathcal{K}$ such that k_l is the l^{th} element of \mathcal{K}_l . Then

$$\lim_{l \rightarrow \infty} \|\mathbf{x}_{k_l}(i, j) - \mathbf{x}(i, j)\|_T = 0 \quad \text{for all } (i, j) \in \mathcal{I}.$$

Let $\mathbf{x} \in D_{\mathbb{R}^{\mathcal{I}}}[0, T]$ be the function with coordinates the limiting functions $\mathbf{x}(i, j)$ introduced in (b). Then \mathbf{x}_{k_l} converges to \mathbf{x} with respect to ϱ_T . \square

As a result, it suffices to establish that $\{\mathcal{A}_n(i, j) : n \geq 1\}$ and $\{\mathcal{D}_n(i, j) : n \geq 1\}$ are relatively compact in $D[0, T]$ for all $(i, j) \in \mathcal{I}$. Consider the sets

$$L_M := \{\mathbf{x} \in D[0, T] : \mathbf{x}(0) = 0 \text{ and } |\mathbf{x}(t) - \mathbf{x}(s)| \leq M|t - s| \text{ for all } s, t \in [0, T]\}$$

which are compact for each $M > 0$ by the Arzelá-Ascoli theorem. For each $(i, j) \in \mathcal{I}$, we prove that there exists M_j such that $\mathcal{A}_n(i, j)$ and $\mathcal{D}_n(i, j)$ approach L_{M_j} as n grows to large. Then we use the compactness of L_{M_j} to show that $\{\mathcal{A}_n(i, j) : n \geq 1\}$ and $\{\mathcal{D}_n(i, j) : n \geq 1\}$ are relatively compact subsets of $D[0, T]$. To this end, we introduce the spaces

$$L_M^\varepsilon := \{\mathbf{x} \in D[0, T] : \mathbf{x}(0) = 0 \text{ and } |\mathbf{x}(t) - \mathbf{x}(s)| \leq M|t - s| + \varepsilon \text{ for all } s, t \in [0, T]\}.$$

LEMMA 2. *If $\mathbf{x} \in L_M^\varepsilon$, then there exists $\mathbf{y} \in L_M$ such that $\|\mathbf{x} - \mathbf{y}\|_T \leq 4\varepsilon$.*

The above lemma is a restatement of (Bramson 1998, Lemma 4.2), and together with the next lemma, it implies that for each $(i, j) \in \mathcal{I}$ there exists a constant M_j such that $\mathcal{A}_n(i, j)$ and $\mathcal{D}_n(i, j)$ approach L_{M_j} as n grows large.

LEMMA 3. *For each $(i, j) \in \mathcal{I}$ there exist $M_j > 0$ and $\{\varepsilon_n(i, j) > 0 : n \geq 1\}$ such that $\mathcal{A}_n(i, j)$ and $\mathcal{D}_n(i, j)$ lie in $L_{M_j}^{\varepsilon_n(i, j)}$ for all n and $\varepsilon_n(i, j) \rightarrow 0$ as n grows large.*

Proof. If $j = 0$, then $\mathcal{A}_n(i, j)$ and $\mathcal{D}_n(i, j)$ are identically zero for all i and n , so the claim holds trivially. Therefore, let us fix $(i, j) \in \mathcal{I}_+$, in this case we have

$$|\mathcal{A}_n(t, i, j) - \mathcal{A}_n(s, i, j)| \leq \frac{1}{n} |\mathcal{N}_n^\lambda(t) - \mathcal{N}_n^\lambda(s)| \leq \lambda|t - s| + 2 \sup_{r \in [0, T]} \left| \frac{1}{n} \mathcal{N}_n^\lambda(r) - \lambda r \right|$$

for all $s, t \in [0, T]$. It follows from Lemma 1 that there exists a vanishing sequence of positive real numbers $\{\delta_n^1 > 0 : n \geq 1\}$ such that

$$|\mathcal{A}_n(t, i, j) - \mathcal{A}_n(s, i, j)| \leq \lambda|t - s| + \delta_n^1 \quad \text{for all } s, t \in [0, T].$$

Consider now the non-decreasing function defined as

$$\mathbf{f}_n(t) := \int_0^t \mu j [\mathbf{q}_n(s, i, j) - \mathbf{q}_n(s, i, j + 1)] ds \quad \text{for all } t \in [0, T].$$

Note that $\mathbf{f}_n(0) = 0$ and $|\mathbf{f}_n(t) - \mathbf{f}_n(s)| \leq \mu j |t - s|$ for all $s, t \in [0, T]$, so in particular, $\mathbf{f}_n(T) \leq \mu j T$. For all $s, t \in [0, T]$, we have

$$\begin{aligned} |\mathcal{D}_n(t, i, j) - \mathcal{D}_n(s, i, j)| &= \frac{1}{n} |\mathcal{N}_{(i, j)}(n\mathbf{f}_n(t)) - \mathcal{N}_{(i, j)}(n\mathbf{f}_n(s))| \\ &\leq |\mathbf{f}_n(t) - \mathbf{f}_n(s)| + 2 \sup_{r \in [0, T]} \left| \frac{1}{n} \mathcal{N}_{(i, j)}(n\mathbf{f}_n(r)) - \mathbf{f}_n(r) \right| \\ &\leq \mu j |t - s| + 2 \sup_{r \in [0, \mu j T]} \left| \frac{1}{n} \mathcal{N}_{(i, j)}(nr) - r \right|. \end{aligned}$$

By (52c), there exists a vanishing sequence $\{\delta_n^2(i, j) > 0 : n \geq 1\}$ such that

$$|\mathcal{D}_n(t, i, j) - \mathcal{D}_n(s, i, j)| \leq \mu j |t - s| + \delta_n^2(i, j) \quad \text{for all } s, t \in [0, T].$$

The result follows setting $M_j := \max\{\lambda, \mu j\}$ and $\varepsilon_n(i, j) := \max\{\delta_n^1, \delta_n^2(i, j)\}$. \square

We now demonstrate that $\{\mathcal{A}_n : n \geq 1\}$, $\{\mathcal{D}_n : n \geq 1\}$ and $\{\mathbf{q}_n : n \geq 1\}$ are relatively compact subsets of $D_{\mathbb{R}^x}[0, T]$ with probability one, and that the limit of every convergent subsequence is componentwise Lipschitz.

LEMMA 4. *The sequences $\{\mathcal{A}_n(\omega) : n \geq 1\}$, $\{\mathcal{D}_n(\omega) : n \geq 1\}$ and $\{\mathbf{q}_n(\omega) : n \geq 1\}$ are all relatively compact in $D_{\mathbb{R}^x}[0, T]$ for all $\omega \in \Gamma_T$. Also, they have the property that the limit of every convergent subsequence is a function with Lipschitz coordinates.*

Proof. We fix an arbitrary $\omega \in \Gamma_T$ and we omit it from the notation for brevity. As noted above, it suffices to prove that $\{\mathcal{A}_n(i, j) : n \geq 1\}$ and $\{\mathcal{D}_n(i, j) : n \geq 1\}$ are relatively compact subsets of $D[0, T]$ for all $(i, j) \in \mathcal{I}$. We fix some $(i, j) \in \mathcal{I}$ and demonstrate that $\{\mathcal{A}_n(i, j) : n \geq 1\}$ is relatively compact in $D[0, T]$, and that the limit of every convergent subsequence is Lipschitz. The same arguments can be used if \mathcal{A}_n is replaced by \mathcal{D}_n .

Let M_j and $\{\varepsilon_n(i, j) : n \geq 1\}$ be as in the statement of Lemma 3. It follows from Lemma 2 that for each n there exists $\mathbf{x}_n(i, j) \in L_{M_j}$ such that

$$\|\mathcal{A}_n(i, j) - \mathbf{x}_n(i, j)\|_T \leq 4\varepsilon_n(i, j).$$

Recall that L_{M_j} is compact, thus every increasing sequence of natural numbers has a subsequence \mathcal{K} such that $\{\mathbf{x}_k(i, j) : k \in \mathcal{K}\}$ converges to a function $\mathbf{x} \in L_{M_j}$. Furthermore,

$$\limsup_{k \rightarrow \infty} \|\mathcal{A}_k(i, j) - \mathbf{x}(i, j)\|_T \leq \lim_{k \rightarrow \infty} 4\varepsilon_k(i, j) + \lim_{k \rightarrow \infty} \|\mathbf{x}_k(i, j) - \mathbf{x}(i, j)\|_T = 0,$$

where the limits are taken along \mathcal{K} . This shows that every subsequence of $\{\mathcal{A}_n(i, j) : n \geq 1\}$ has a further subsequence which converges to a Lipschitz function. \square

We are now ready to prove Proposition 2

Proof of Proposition 2. Recall that Γ_T has probability one for all $T \geq 0$. Hence,

$$\Gamma_\infty := \bigcap_{T \in \mathbb{N}} \Gamma_T$$

has probability one as well. Also, (22) and (23) hold within Γ_∞ by Lemma 1.

Next we fix an arbitrary $\omega \in \Gamma_\infty$, which we omit from the notation for brevity, and we prove that $\{\mathbf{q}_n : n \geq 1\}$ is a relatively compact subset of $D_{\mathbb{R}^x}[0, \infty)$ such that the limit of every convergent subsequence has locally Lipschitz coordinate functions. Exactly the same arguments can be used if \mathbf{q}_n is replaced by \mathcal{A}_n or \mathcal{D}_n .

Fix a sequence $\mathcal{K} \subset \mathbb{N}$. We must show that $\{\mathbf{q}_k : k \in \mathcal{K}\}$ has a subsequence which converges uniformly over compact sets to a function with locally Lipschitz coordinates. To this end, we may construct a family of sequences $\{\mathcal{K}_T : T \in \mathbb{N}\}$ such that:

- (a) $\mathcal{K}_{T+1} \subset \mathcal{K}_T \subset \mathcal{K}$ for all $T \in \mathbb{N}$;
- (b) for each $T \in \mathbb{N}$, there exists $\mathbf{q}_T \in D_{\mathbb{R}^x}[0, T]$ such that the coordinates of \mathbf{q}_T are Lipschitz and $\varrho_T(\mathbf{q}_k|_{[0, T]}, \mathbf{q}_T) \rightarrow 0$ as $k \rightarrow \infty$ with $k \in \mathcal{K}_T$.

Let k_l denote the l^{th} element of \mathcal{K}_l . It follows from (a) and (b) that

$$\lim_{l \rightarrow \infty} \varrho_T(\mathbf{q}_{k_l}|_{[0, T]}, \mathbf{q}_T) = 0 \quad \text{for all } T \in \mathbb{N}. \tag{53}$$

Note that $\mathbf{q}_{k_l}(t) \rightarrow \mathbf{q}_T(t)$ with l for all $t \leq T$, so $\mathbf{q}_S(t) = \mathbf{q}_T(t)$ for all $t \leq S, T$. Thus, we may define a function $\mathbf{q} \in D_{\mathbb{R}^x}[0, \infty)$ such that $\mathbf{q}(t) = \mathbf{q}_T(t)$ for all $t \leq T$. Moreover, (b) implies that \mathbf{q} has locally Lipschitz coordinates and (53) means that $\{\mathbf{q}_{k_l} : l \geq 1\}$ converges uniformly over compact sets to \mathbf{q} . \square

Appendix C: Limiting behavior of SLTA

Proof of Proposition 4. We fix an arbitrary $\omega \in \Gamma$, which is omitted from the notation for brevity. It follows from (21a) that

$$\mathbf{s}_n = \mathbf{s}_n(0) + \sum_{(i,j) \in \mathcal{I}_+} [\mathcal{A}_n(i,j) - \mathcal{D}_n(i,j)] = \mathbf{s}_n(0) + \frac{1}{n} \mathcal{N}_n^\lambda(t) - \sum_{(i,j) \in \mathcal{I}_+} \mathcal{D}_n(i,j) \quad \text{for all } n.$$

By Proposition 3, every increasing sequence of natural numbers has a subsequence \mathcal{K} such that the sequences $\{\mathcal{D}_k : k \in \mathcal{K}\}$ and $\{\mathbf{q}_k : k \in \mathcal{K}\}$ converge in $D_{\ell_1}[0, \infty)$ to certain functions \mathbf{d} and \mathbf{q} , respectively. It follows from (22c) that

$$\mathbf{d}(t, i, j) = \int_0^t \mu_j [\mathbf{q}(s, i, j) - \mathbf{q}(s, i, j + 1)] ds \quad \text{for all } t \geq 0 \quad \text{and } (i, j) \in \mathcal{I}_+.$$

Moreover, $\{\mathbf{s}_k : k \in \mathcal{K}\}$ converges uniformly over compact sets to

$$\tilde{\mathbf{s}} := \sum_{(i,j) \in \mathcal{I}_+} \mathbf{q}(i, j).$$

It follows from (22a) and (22b) that

$$\lim_{k \rightarrow \infty} \mathbf{s}_k(0) = \mathbf{s}_0 \quad \text{and} \quad \lim_{k \rightarrow \infty} \sup_{t \in [0, T]} \left| \frac{1}{k} \mathcal{N}_k^\lambda(t) - \lambda t \right| = 0 \quad \text{for all } T \geq 0.$$

Furthermore, the fact that $\{\mathcal{D}_k : k \in \mathcal{K}\}$ converges to \mathbf{d} in $D_{\ell_1}[0, \infty)$ implies that

$$\begin{aligned} & \lim_{k \rightarrow \infty} \sup_{t \in [0, T]} \left| \sum_{(i,j) \in \mathcal{I}_+} \mathcal{D}_k(t, i, j) - \int_0^t \mu \tilde{\mathbf{s}}(s) ds \right| = \\ & \lim_{k \rightarrow \infty} \sup_{t \in [0, T]} \left| \sum_{i=1}^m \sum_{j=1}^{\infty} \left[\mathcal{D}_k(t, i, j) - \int_0^t \mu_j [\mathbf{q}(s, i, j) - \mathbf{q}(s, i, j + 1)] ds \right] \right| = \\ & \lim_{k \rightarrow \infty} \sup_{t \in [0, T]} \left| \sum_{(i,j) \in \mathcal{I}_+} [\mathcal{D}_k(t, i, j) - \mathbf{d}(t, i, j)] \right| = 0 \quad \text{for all } T \geq 0. \end{aligned}$$

We conclude that $\tilde{\mathbf{s}}$ satisfies

$$\tilde{\mathbf{s}}(t) = \mathbf{s}_0 + \lambda t - \int_0^t \mu \tilde{\mathbf{s}}(s) ds \quad \text{for all } t \geq 0.$$

Equivalently, $\tilde{\mathbf{s}}$ solves (35), and therefore $\tilde{\mathbf{s}} = \mathbf{s}$. This demonstrates that every subsequence of $\{\mathbf{s}_n : n \geq 1\}$ has a further subsequence that converges to \mathbf{s} uniformly over compact sets. The claim of the proposition follows from this fact. \square

Proof of Proposition 5. Let us omit ω from the notation for brevity. Recall that tasks are always assigned to server pools with coordinates $(i, j) \succeq (i_{r_n}, j_{r_n})$, as explained in Remark 1. By (b), $(i_{r_k}, j_{r_k}) \succeq (i_r, j_r)$ along the interval $[t_0, t_1]$ for each $k \in \mathcal{K}$. Moreover, (b) implies that the system has green tokens when the equality holds. Therefore, tasks are only assigned to server pools with coordinates $(i, j) \triangleright (i_r, j_r)$ along the interval $[t_0, t_1]$ in all the systems comprised in \mathcal{K} . Thus, $\mathcal{A}_k(t, i, j) = \mathcal{A}_k(t_0, i, j)$ for all $t \in [t_0, t_1]$ and $(i, j) \preceq (i_r, j_r)$. We conclude from (21a) that

$$\mathbf{q}_k(t, i, j) = \mathbf{q}_k(t_0, i, j) - [\mathcal{D}_k(t, i, j) - \mathcal{D}_k(t_0, i, j)]$$

for all $t \in [t_0, t_1]$, $(i, j) \preceq (i_r, j_r)$ and $k \in \mathcal{K}$. As in the proof of Proposition 4, we obtain

$$\mathbf{q}(t, i, j) = \mathbf{q}(t_0, i, j) - \int_{t_0}^t \mu j [\mathbf{q}(s, i, j) - \mathbf{q}(s, i, j+1)] ds$$

for all $t \in [t_0, t_1]$, $(i, j) \preceq (i_r, j_r)$ and $k \in \mathcal{K}$. It follows that $\mathbf{q}(i, j)$ is differentiable for all $(i, j) \preceq (i_r, j_r)$ because \mathbf{q} has continuous coordinates by Proposition 2. Furthermore, the system of differential equations in the statement of the proposition holds.

Note that $\{\mathbf{v}_k(r) : k \geq 1\}$ converges uniformly over compact sets to

$$\mathbf{v}(r) := \sum_{(i,j) \preceq (i_r, j_r)} \mathbf{q}(i, j)$$

by (a) and the definition of \mathbf{v}_k . Let $J_i := \min \{j \geq 1 : (i, j) \preceq (i_r, j_r)\}$ and note that

$$\begin{aligned} \dot{\mathbf{v}}(t, r) &= \sum_{i=1}^m \sum_{j=J_i}^{\infty} -\mu j [\mathbf{q}(t, i, j) - \mathbf{q}(t, i, j+1)] \\ &= \sum_{i=1}^m \left[-\mu (J_i - 1) \mathbf{q}(t, i, J_i) - \mu \sum_{j=J_i}^{\infty} \mathbf{q}(t, i, j) \right] = -\mu \sum_{i=1}^m (J_i - 1) \mathbf{q}(t, i, J_i) - \mu \mathbf{v}(t, r). \end{aligned}$$

Since the first term on the right-hand side is non-positive, we have

$$\mathbf{v}(t, r) \leq \mathbf{v}(t_0, r) e^{-\mu(t-t_0)} < \mathbf{s}(t_0) e^{-\mu(t-t_0)} \quad \text{for all } t \in [t_0, t_1].$$

This completes the proof. \square

Proof of Proposition 6. Define

$$\alpha_{\min} := \min_{1 \leq i \leq m} \alpha(i), \quad J := \left\lfloor \frac{\rho+1}{\alpha_{\min}} \right\rfloor \quad \text{and} \quad r_{\text{bd}} := \max \{r \geq 2 : j_{r-1} \leq J\}.$$

The function τ_{bd} mentioned in the statement of the proposition is defined by

$$\tau_{\text{bd}}(s) = \begin{cases} \frac{1}{\mu} \left[\log \left(\frac{s-\rho}{\sum_{(i,j) \succeq \sigma_*} -\rho} \right) \right]^+ + \frac{J-j_{r_*}}{\mu} (r_{\text{bd}} - r_*) & \text{if } s > \rho, \\ \frac{J-j_{r_*}}{\mu} (r_{\text{bd}} - r_*) & \text{if } s \leq \rho. \end{cases}$$

For brevity, we omit ω from the notation, and we let

$$\Sigma(r) := \sum_{(i,j) \succ (i_r, j_r)} \alpha(i) \quad \text{and} \quad \Sigma_n(r) := \sum_{(i,j) \succ (i_r, j_r)} \alpha_n(i) \quad \text{for all } n, r \geq 1.$$

It follows from (10) that we can write $\Sigma(r_*+1) = \rho + \delta + \varepsilon$ with $\delta, \varepsilon > 0$. We regard ε as a function of δ , determined by the latter identity, and for each

$$0 < \delta < \min \{\Sigma(r_*+1) - \rho, \alpha_{\min}\} \quad \text{and} \quad \theta \in (0, 1)$$

we define $\tau(\delta, \theta) := \min \{t \geq 0 : (s_0 - \rho) e^{-\mu t} \leq \theta \varepsilon\}$. Explicitly,

$$\tau(\delta, \theta) = \begin{cases} \frac{1}{\mu} \left[\log \left(\frac{s_0 - \rho}{\theta [\Sigma(r_*+1) - \rho - \delta]} \right) \right]^+ & \text{if } s_0 > \rho, \\ 0 & \text{if } s_0 \leq \rho. \end{cases}$$

Thus, the function \mathbf{s} defined by (35) satisfies

$$\mathbf{s}(t) \leq \rho + \theta \varepsilon \quad \text{for all } t \geq \tau(\delta, \theta). \quad (54)$$

Note that $\tau_{\text{bd}}(s_0)$ is the infimum over δ and θ of

$$\tau(\delta, \theta) + \delta + \frac{J-j_{r_*}}{\mu \theta} (r_{\text{bd}} - r_*).$$

We fix δ and θ such that the above expression is smaller than τ and we choose n_0 such that the following properties hold for all $n \geq n_0$.

- (i) $\mathbf{r}_n(t) \leq R_T$ for all $t \in [0, T]$.
- (ii) $|\mathbf{s}_n(t) - \mathbf{s}(t)| \leq (1 - \theta)\varepsilon$ for all $t \in [0, T]$.
- (iii) $1/n \leq \beta_n \leq \min\{\alpha_n(i) : 1 \leq i \leq m\}$ and $\Sigma_n(r_* + 1) - L(r_* + 1)\beta_n \geq \rho + \varepsilon$, where we used the notation $L(r) := \max\{\ell_i(r) : 1 \leq i \leq m\}$.

Properties (i) and (ii) are possible by Lemma 2 and Proposition 4, respectively, and (iii) can be enforced by assumptions (9) and (14).

Fix $n \geq n_0$ and note that for each $r > r_*$ and $t \in [\tau(\delta, \theta), T]$ we have

$$\Sigma_n(r) - L(r)\beta_n \geq \Sigma_n(r_* + 1) - L(r_* + 1)\beta_n \geq \rho + \varepsilon \geq \mathbf{s}_n(t).$$

The first inequality follows from the fact that $\Sigma_n(\cdot) - L(\cdot)\beta_n$ is non-decreasing provided that $\beta_n \leq \min\{\alpha_n(i) : 1 \leq i \leq m\}$, and the last inequality follows from (ii) and (54). We conclude from the above inequalities that

$$\sum_{i=1}^m \mathbf{q}_n(t, i, \ell_i(r)) \leq 1 - \beta_n \leq 1 - \frac{1}{n} \quad \text{for all } t \in [\tau(\delta, \theta), T], \quad (55)$$

since otherwise we would have $\mathbf{s}_n(t) > \Sigma_n(r) - L(r)\beta_n$. It follows from (55) with $r = r_* + 1$ that the total number of tokens is at least one if $t \in [\tau(\delta, \theta), T]$ and $\mathbf{r}_n(t) \geq r_*$, so we conclude that \mathbf{r}_n does not increase beyond r_* along the interval $[\tau(\delta, \theta), T]$. Hence, it only remains to prove that the amount of time after $\tau(\delta, \theta)$ until $\mathbf{r}_n \leq r_*$ is upper bounded by

$$\delta + \frac{J - j_{r_*}}{\mu\theta} (r_{\text{bd}} - r_*).$$

It follows from (55) that the number of green tokens is

$$n - \sum_{i=1}^m n \mathbf{q}_n(t, i, \ell_i(\mathbf{r}_n(t))) \geq n\beta_n \quad \text{if } \mathbf{r}_n(t) > r_*, \quad t \in [\tau(\delta, \theta), T] \quad \text{and } n \geq n_0.$$

Therefore, \mathbf{r}_n decreases with the next arrival if the latter conditions hold and

$$\mathbf{q}_n(t, i_{\mathbf{r}_n(t)-1}, j_{\mathbf{r}_n(t)-1}) < \alpha_n(i_{\mathbf{r}_n(t)-1}).$$

Choose $n_{\text{bd}}^{\tau, T} \geq n_0$ such that the following conditions hold.

- (iv) $\min\{\alpha_n(i) : 1 \leq i \leq m\} > \alpha(\delta) := \alpha_{\min} - \delta$.
- (v) Let $d := \delta [2(r_{\text{bd}} - r_*) + 1]^{-1}$. The following two properties hold on every subinterval of $[0, T]$ of length at least d . First, $\mathcal{N}_n^\lambda(\cdot)$ has at least $R_T - r_{\text{bd}}$ jumps. Second, for all $1 \leq i \leq m$ and $1 \leq j \leq J$, the process $\mathcal{N}_{(i,j)}(n \cdot)$ has at least $n\theta d$ jumps followed by at least one jump of $\mathcal{N}_n^\lambda(\cdot)$.

- (vi) $\mu n [\alpha_n(i) - \alpha(\delta)] \theta d \geq 1$ for all $1 \leq i \leq m$.

Properties (iv) and (v) can be enforced by (9) and Proposition 2, respectively.

We may assume without loss of generality that

$$J \geq \left\lceil \frac{\rho + \varepsilon}{\alpha(\delta)} \right\rceil,$$

choosing a smaller δ at the start of the proof if needed. By the above inequality,

$$\mathbf{q}_n(t, i, j) \leq \frac{\mathbf{s}_n(t)}{j} \leq \frac{\rho + \varepsilon}{J + 1} < \alpha(\delta) \leq \alpha_n(i) \quad \text{if } j > J, \quad t \in [\tau(\delta, \theta), T] \quad \text{and } n \geq n_{\text{bd}}^{\tau, T}.$$

Therefore, if $n \geq n_{\text{bd}}^{\tau, T}$, $t \in [\tau(\delta, \theta), T]$ and $\mathbf{r}_n(t) > r_{\text{bd}}$, then \mathbf{r}_n decreases with each arrival until it reaches r_{bd} . We conclude from (i) and (v) that

$$\mathbf{r}_n(t) \leq r_{\text{bd}} \quad \text{for all } t \in [\tau(\delta, \theta) + d, T] \quad \text{and} \quad n \geq n_{\text{bd}}^{\tau, T}.$$

Suppose now that $\mathbf{r}_n(t) = r \in (r_*, r_{\text{bd}}]$ for some $t \in [\tau(\delta, \theta) + d, T]$ and $n \geq n_{\text{bd}}^{\tau, T}$. Then \mathbf{r}_n decreases as soon as a server pool of class i_{r-1} has a green token and a task arrives. Moreover, since the number of green tokens remains positive, no tasks are sent to server pools of class i_{r-1} before \mathbf{r}_n decreases. At least $n[\alpha_n(i) - \alpha(\delta)]$ server pools of class i_{r-1} have strictly less than $J + 1$ tasks, thus it follows from (v) that the time until $\mathbf{q}_n(i_{r-1}, j_{r-1})$ drops below $\alpha_n(i)$ is at most

$$\left\lceil \frac{n(J - j_{r-1})[\alpha_n(i) - \alpha(\delta)] + 1}{\mu n [\alpha_n(i) - \alpha(\delta)] j_{r-1} \theta d} \right\rceil d \leq \frac{J - j_{r-1}}{\mu \theta} + \frac{1}{\mu n [\alpha_n(i) - \alpha(\delta)] \theta} + d \leq \frac{J - j_{r_*}}{\mu \theta} + 2d.$$

The numerator inside the ceiling function is the maximum number of tasks that need to leave the $n[\alpha_n(i) - \alpha(\delta)]$ server pools of class i_{r-1} with strictly fewer than $J + 1$ tasks to ensure that one of these server pools has strictly less than j_{r-1} tasks, and the quantity $n[\alpha_n(i) - \alpha(\delta)] j_{r-1}$ in the denominator is a lower bound for the cumulative number of tasks in these server pools while all the server pools have at least j_{r-1} tasks. Also, note that the first inequality uses $j_{r-1} \geq 1$ and the last inequality uses (vi). Since \mathbf{r}_n decreases at most $r_{\text{bd}} - r_*$ times before it reaches r_* , we have

$$\mathbf{r}_n(t) \leq r_* \quad \text{for all } t \in \left[\tau(\delta, \theta) + \delta + \frac{J - j_{r_*}}{\mu \theta} (r_{\text{bd}} - r_*), T \right] \quad \text{and} \quad n \geq n_{\text{bd}}^{\tau, T}.$$

This completes the proof. \square

Proof of Lemma 6. We omit ω from the notation. It follows from (21a) that

$$\begin{aligned} \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) - \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(\tau_{k,1}, i, j) &= \sum_{(i,j) \triangleright (i_r, j_r)} [\mathcal{A}_k(t, i, j) - \mathcal{A}_k(\tau_{k,1}, i, j)] \\ &\quad - \sum_{(i,j) \triangleright (i_r, j_r)} [\mathcal{D}_k(t, i, j) - \mathcal{D}_k(\tau_{k,1}, i, j)]. \end{aligned}$$

Note that all tasks are assigned to server pools with coordinates $(i, j) \triangleright (i_r, j_r)$ during the interval $(\tau_{k,1}, t]$, so the first term on the right-hand side can be expressed as follows:

$$\sum_{(i,j) \triangleright (i_r, j_r)} [\mathcal{A}_k(t, i, j) - \mathcal{A}_k(\tau_{k,1}, i, j)] = \frac{1}{k} [\mathcal{N}_k^\lambda(t) - \mathcal{N}_k^\lambda(\tau_{k,1})].$$

Using the process $\delta_k(r)$ defined in (36), we may write

$$\begin{aligned} \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) - \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(\tau_{k,1}, i, j) &= (t - \tau_{k,1}) \lambda - \sum_{(i,j) \triangleright (i_r, j_r)} \int_{\tau_{k,1}}^t \mu j [\mathbf{q}_k(s, i, j) - \mathbf{q}_k(s, i, j + 1)] ds \\ &\quad + \delta_k(t, r) - \delta_k(\tau_{k,1}, r) \\ &\geq (t - \tau_{k,1}) \left[\lambda - \mu \sum_{(i,j) \triangleright (i_r, j_r)} \alpha_k(i) \right] - 2 \sup_{s \in [0, T]} |\delta_k(s, r)|. \end{aligned}$$

For the last inequality, observe that

$$\begin{aligned}
\sum_{(i,j) \triangleright (i_r, j_r)} j [\mathbf{q}_k(i, j) - \mathbf{q}_k(i, j+1)] &= \sum_{i=1}^m \sum_{j=1}^{\ell_i(r)} j [\mathbf{q}_k(i, j) - \mathbf{q}_k(i, j+1)] \\
&= \sum_{i=1}^m \left[\sum_{j=1}^{\ell_i(r)} \mathbf{q}_k(i, j) - \ell_i(r) \mathbf{q}_k(i, \ell_i(r)+1) \right] \\
&\leq \sum_{(i,j) \triangleright (i_r, j_r)} \alpha_k(i) \quad \text{for all } k \in \mathcal{K}.
\end{aligned}$$

This completes the proof. \square

Proof of Lemma 7. We omit ω from the notation and we assume that $r > 1$, otherwise the statement holds trivially. Fix an arbitrary $\gamma \in [0, 1/2)$ and let

$$\begin{aligned}
\xi_{k,2} &:= \inf \left\{ t \geq \zeta_{k,1} : \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) \leq \sum_{(i,j) \triangleright (i_r, j_r)} \alpha_k(i) - \min\{\beta_k, k^{-\gamma}\} \right\}, \\
\xi_{k,1} &:= \sup \left\{ t \leq \xi_{k,2} : \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) = \sum_{(i,j) \triangleright (i_r, j_r)} \alpha_k(i) \right\}.
\end{aligned}$$

By assumption, $\mathbf{r}_k \leq r$ along the interval $[\zeta_{k,1}, \zeta_{k,2}]$ and $\mathbf{r}_k(\zeta_{k,1}) = r$, since otherwise the number of tokens at time $\zeta_{k,1}$ would be zero, and this cannot happen by Remark 1. Moreover, the number of green tokens is upper bounded by

$$\sum_{(i,j) \triangleright (i_r, j_r)} k \alpha_k(i) - \sum_{(i,j) \triangleright (i_r, j_r)} k \mathbf{q}_k(i, j) < k \beta_k$$

along the interval $[\zeta_{k,1}, \xi_{k,2})$, which implies that \mathbf{r}_k cannot decrease along $[\zeta_{k,1}, \xi_{k,2}]$. We conclude that $\mathbf{r}_k(t) = r$ for all $t \in [\zeta_{k,1}, \zeta_{k,2} \wedge \xi_{k,2}]$ and $k \in \mathcal{K}$. Moreover,

$$\max \{ \alpha_k(i) - \mathbf{q}_k(i, j) : (i, j) \triangleright (i_r, j_r) \} \leq \sum_{(i,j) \triangleright (i_r, j_r)} \alpha_k(i) - \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(i, j) \leq k^{-\gamma}$$

along $[\zeta_{k,1}, \xi_{k,2}]$. Thus, it suffices to prove that $\xi_{k,2} \geq \zeta_{k,2}$ for all large enough $k \in \mathcal{K}$.

In order to prove that $\xi_{k,2} \geq \zeta_{k,2}$, we show that

$$\begin{aligned}
&\sum_{(i,j) \triangleright (i_r, j_r)} \alpha_k(i) - \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) = \\
&\sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(\xi_{k,1}^-, i, j) - \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) < \min\{\beta_k, k^{-\gamma}\}
\end{aligned} \tag{56}$$

for all $t \in [\xi_{k,1}, \zeta_{k,2} \wedge \xi_{k,2}]$ and all large enough $k \in \mathcal{K}$. Note that, by definition of $\xi_{k,2}$, the latter statement implies that $\xi_{k,2} > \zeta_{k,2}$.

Let $\tau_{k,1} := \zeta_{k,2} \wedge \xi_{k,1}$ and $\tau_{k,2} := \zeta_{k,2} \wedge \xi_{k,2}$. In addition, let M_r be the cardinality of the finite set $\{(i, j) \in \mathcal{I}_+ : (i, j) \triangleright (i_r, j_r)\}$. For all $t \in [\tau_{k,1}, \tau_{k,2}]$, we have

$$\begin{aligned}
&\sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) - \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(\tau_{k,1}^-, i, j) \\
&\geq \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) - \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(\tau_{k,1}, i, j) - \frac{M_r}{k}
\end{aligned}$$

$$\begin{aligned}
&\geq (t - \tau_{k,1}) \left[\lambda - \mu \sum_{(i,j) \triangleright (i_r, j_r)} \alpha_k(i) \right] - 2 \sup_{s \in [0, T]} |\delta_k(s, r)| - \frac{M_r}{k} \\
&\geq -2 \sup_{s \in [0, T]} |\delta_k(s, r)| - \frac{M_r}{k}.
\end{aligned}$$

For the first inequality, note that the processes $\mathbf{q}_k(i, j)$ have jumps of size $1/k$. The second inequality follows from Lemma 6, since

$$\mathbf{r}_k(t) = r \quad \text{and} \quad \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) < \sum_{(i,j) \triangleright (i_r, j_r)} \alpha_k(i)$$

for all $t \in [\tau_{k,1}, \tau_{k,2})$ and $k \in \mathcal{K}$. The third inequality follows from (10) and $r \leq r_*$.

We conclude from (14) and (37) that

$$\lim_{k \rightarrow \infty} \frac{1}{\beta_k} \left[2 \sup_{s \in [0, T]} |\delta_k(s, r)| + \frac{M_r}{k} \right] = \lim_{k \rightarrow \infty} \frac{1}{k^{\gamma_0} \beta_k} \left[2 \sup_{s \in [0, T]} k^{\gamma_0} |\delta_k(s, r)| + \frac{M_r}{k^{1-\gamma_0}} \right] = 0.$$

If β_k is replaced by $k^{-\gamma}$ on the left-hand side, then the limit is also equal to zero by (37). Consequently, for all sufficiently large $k \in \mathcal{K}$, we have

$$\begin{aligned}
\sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(t, i, j) - \sum_{(i,j) \triangleright (i_r, j_r)} \mathbf{q}_k(\tau_{k,1}^-, i, j) &\geq -2 \sup_{s \in [0, T]} |\delta_k(s, r)| - \frac{M_r}{k} \\
&> -\min\{\beta_k, k^{-\gamma}\}
\end{aligned}$$

for all $t \in [\tau_{k,1}, \tau_{k,2}]$. This implies that (56) holds for all sufficiently large $k \in \mathcal{K}$; observe that (56) holds trivially when $\xi_{k,1} > \zeta_{k,2}$. \square

References

- Bramson M (1998) State space collapse with application to heavy traffic limits for multiclass queueing networks. *Queueing Systems* 30(1-2):89–140.
- Chen H, Yao DD (2013) *Fundamentals of queueing networks: Performance, asymptotics, and optimization* (Springer Science & Business Media).
- Goldsztajn D, Borst SC, Van Leeuwaarden JS (2021) Learning and balancing unknown loads in large-scale systems. *arXiv preprint arXiv:2012.10142* .
- Karatzas I, Shreve S (2014) *Brownian motion and stochastic calculus*, volume 113 (Springer).