

MANY-SERVER ASYMPTOTICS FOR JSQ IN SUPER-HALFIN-WHITT REGIME

Zhisheng Zhao[✉]

Georgia Institute of Technology

Model: Parallel Queueing System

- Single dispatcher: tasks arriving as a **Poisson process** of rate $\lambda(N)$,

$$\lambda(N) = N - \beta N^{-1/2+\varepsilon},$$

where $\beta > 0$ and $\varepsilon \in (0, 1/2)$;

- Incoming tasks must immediately be sent to one of the queues under the **Joint-the-Shortest Queue** (JSQ) policy;
- N servers working at **unit rate**, service requirements are **exponential**.

State Description and Notation

- $S^{(N)}(t)$: the total number of tasks at time t in the N -th system;
- $Q_i^{(N)}(t)$: the number of servers with at least $i \in \mathbb{N}_0$ tasks at time t in the N -th system;
- $I^{(N)}(\cdot) = N - Q_1^{(N)}(\cdot)$: the idle process of the N -th system;
- $A(\cdot)$ and $D(\cdot)$ are independent Poisson processes with unit rate;
- $W(\cdot)$: the standard Brownian motion;
- Define a **centered** and **scaled** process

$$X^{(N)}(t) := \frac{S^{(N)}(N^{2\varepsilon}t) - N}{N^{1/2+\varepsilon}}.$$

Literature Review

Value of α	Regime	Asymptotic behavior	References
0	Meanfield	$Q_1^{(N)} = N\lambda_N \pm \Theta_p(\sqrt{N\lambda_N})$, $Q_i^{(N)} = o_p(1)$ for $i \geq 2$	[6]
$(0, \frac{1}{2})$	Sub-Halfin-Whitt	$\sum_{i=1}^b Q_i^{(N)} = N\lambda_N + O_p(\sqrt{N \log N})$	[4]
$\frac{1}{2}$	Halfin-Whitt	$Q_1^{(N)} = N - \Theta_p(\sqrt{N})$, $Q_2^{(N)} = \Theta(\sqrt{N})$, $Q_i^{(N)} = o_p(1)$ for $i \geq 3$	[1]
$(\frac{1}{2}, 1)$	Super-Halfin-Whitt	$Q_1^{(N)} = N - \Theta(N^{1-\alpha})$, $Q_2^{(N)} = \Theta_p(N^\alpha)$, $Q_i^{(N)} = o_p(1)$ for $i \geq 3$	[5], current paper
1	NDS	$Q_i = \Theta_p(N)$ for all $i \geq 1$	[2]
$(1, \infty)$	Super Slowdown	Unknown for $\alpha \in (1, 2]$. For $\alpha > 2$, $\sum_{i=1}^\infty Q_i^{(N)} = \Theta_p(N^\alpha)$	[3]

Tab. 1: Analysis of JSQ in various regimes ($\alpha = 1/2 - \varepsilon$)

Main Results

Process-level convergence

With appropriate assumptions on $S^{(N)}(0)$, $Q_i^{(N)}(0)$, $\forall N, i \in \mathbb{N}_0$, for any finite $T > 0$, $X^{(N)}(\cdot)$ weakly converges to $X(\cdot)$ uniformly on $[0, T]$, where $X(\cdot)$ is the solution of the SDE:

$$dX(t) = \left(\frac{1}{X(t)} - \beta \right) dt + \sqrt{2} dW(t). \quad (1)$$

Remark: The SDE in (1) is a Langevin diffusion so it is ergodic and has a unique stationary distribution $\pi \sim \text{Gamma}(2, \beta)$, having p -th moment $\Gamma(p+2)/\beta^p$.

Stationary distribution of the N -system

There exist constants C_1, C_2 and B such that for large enough N ,

$$\mathbb{P}(X^{(N)}(\infty) \geq x) \leq \begin{cases} C_1 \exp\{-C_2 x^{1/5}\}, & 4B \leq x \leq 2N^{1/2-\varepsilon}, \\ C_1 \exp\{-C_2 x^{1/44}\}, & x \geq 2N^{1/2-\varepsilon}. \end{cases} \quad (2)$$

Moreover, $\sup_{N \geq 1} \mathbb{E}[N^{-1/2-\varepsilon} Q_2^{(N)}(\infty)] < \infty$, $\mathbb{E}[N^{-1/2+\varepsilon} I^{(N)}(\infty)] = \beta$ for large enough N , and $\sum_{i=3}^\infty Q_i^{(N)}(\infty) \xrightarrow{P} 0$ as $N \rightarrow \infty$.

Interchange of limits

Let $X^{(N)}(\infty)$ be the stationary distribution of the scaled process $X^{(N)}(\cdot)$ in the N -th system. The sequence of random variables $\{X^{(N)}(\infty)\}_{N \geq 1}$ converges weakly to the $\text{Gamma}(2, \beta)$ distribution as $N \rightarrow \infty$.

Remark: The interchange of limits holds:

$$\lim_{t \rightarrow \infty} \lim_{N \rightarrow \infty} X^N(t) = \lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} X^N(t) \sim \text{Gamma}(2, \beta). \quad (3)$$

Remark: The centered and scaled total number of tasks in steady state is distributed as the sum of two independent exponential random variables for the JSQ policy, as opposed to a single exponential random variable in the M/M/N case.

Proof Scheme

Martingale representation

$$\begin{aligned} X^{(N)}(t) - X^{(N)}(0) &= N^{-1/2-\varepsilon} \left[A(N^{1+2\varepsilon}\lambda_N t) - D\left(\int_0^{N^{2\varepsilon}t} (N - I^{(N)}(s)) ds\right) \right] \\ &= \mathcal{M}_A^{(N)}(\lambda_N t) - \mathcal{M}_D^{(N)}\left(t - \frac{1}{N^{1+2\varepsilon}} \int_0^{N^{2\varepsilon}t} I^{(N)}(s) ds\right) \\ &\quad + \frac{1}{N^{1/2+\varepsilon}} \int_0^{N^{2\varepsilon}t} I^{(N)}(s) ds - \int_0^t \frac{1}{X^{(N)}(s)} ds \\ &\quad - \beta t + \int_0^t \frac{1}{X^{(N)}(s)} ds \end{aligned} \quad (4)$$

where $\mathcal{M}_H(t) = \frac{H(N^{1+2\varepsilon}t) - N^{1+2\varepsilon}t}{N^{1/2+\varepsilon}}$, $H = A, D$.

Analysis of the process $I^{(N)}$

For the proof of (4) \Rightarrow 0, the main idea is to approximate each excursion of $I^{(N)}$ by M/M/1 queues. Consider an excursion during $[\sigma_1, \sigma_2] \subseteq [0, T]$ (i.e., $I^{(N)}(t) > 0$, $t \in (\sigma_1, \sigma_2)$, and $I^{(N)}(\sigma_i) = 0$, $i = 1, 2$). We have

$$\sup_{t \in [\sigma_1, \sigma_2]} |S^{(N)}(t) - S^{(N)}(\sigma_1)| = o(N^{1/2-\varepsilon}) \quad \text{and} \quad \sup_{t \in [\sigma_1, \sigma_2]} I^{(N)}(t) = o(N^{1/2-\varepsilon})$$

Hence, each excursion of $I^{(N)}$ can be bounded by two M/M/1 queues $\bar{I}_l^{(N)}$ and $\bar{I}_u^{(N)}$ such that with natural coupling, $\bar{I}_l^{(N)} \leq I^{(N)} \leq \bar{I}_u^{(N)}$, and

$$\lim_{N \rightarrow \infty} \frac{1}{N^{1/2+\varepsilon}} \int_{\sigma_1}^{\sigma_2} |\bar{I}_u^{(N)}(s) - \bar{I}_l^{(N)}(s)| ds = 0.$$

Renewal representation of stationary measure

Let the initial state of the N -th system be

$$\{I^{(N)}(0) = 0, Q_2^{(N)}(0) = \lfloor 2BN^{1/2+\varepsilon} \rfloor, Q_3^{(N)}(0) = 0\},$$

where $B > 0$ is appropriately selected. Let $\Theta^{(N)}$ be the next renewal time point, i.e. at time $\Theta^{(N)}$, the system backs to the initial state. Define

$$\pi(X^{(N)}(\infty) \in A) = \frac{\mathbb{E}_{(0, \lfloor 2BN^{1/2+\varepsilon} \rfloor, 0)} \left(\int_0^{\Theta^{(N)}} \mathbf{1}(X^{(N)}(\infty) \in A) du \right)}{\mathbb{E}_{(0, \lfloor 2BN^{1/2+\varepsilon} \rfloor, 0)} (\Theta^{(N)})}.$$

$\Theta^{(N)}$ can be analyzed by two parts: *down-crossing* and *up-crossing*.

From (4) and (5), we have a drift term of $X^{(N)}$:

$$\frac{1}{N^{1/2+\varepsilon}} \int_0^{N^{2\varepsilon}t} I^{(N)}(s) ds - \beta t. \quad (6)$$

Down-crossing. When $Q_2^{(N)} : 2BN^{1/2+\varepsilon} \rightarrow BN^{1/2+\varepsilon}$, $I^{(N)} \leq \bar{I}_B^{(N)}$ where $\bar{I}_B^{(N)}$ is an M/M/1 queue with increase rate $N - BN^{1/2+\varepsilon}$ and $\frac{1}{N^{1/2+\varepsilon}} \int_0^{N^{2\varepsilon}t} \bar{I}_B^{(N)}(s) ds - \beta t < 0$ w.h.p. so the drift (6) would be negative w.h.p..

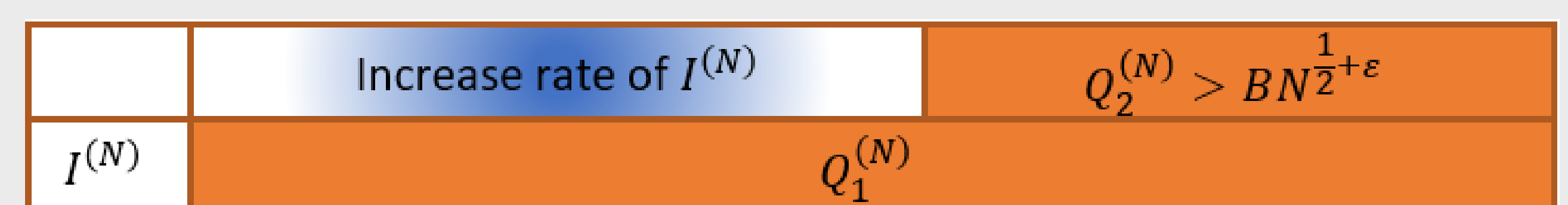


Fig. 1: Down-crossing

Up-crossing. The key observation is that if the system starts with the state in Fig 2, then the probability that $Q_2^{(N)}$ hits $2BN^{1/2+\varepsilon}$ within $N^{2\varepsilon}$ is a constant independent on N . This leads to a geometric number of such excursions required for $Q_2^{(N)}$ to hit the level $2BN^{1/2+\varepsilon}$.

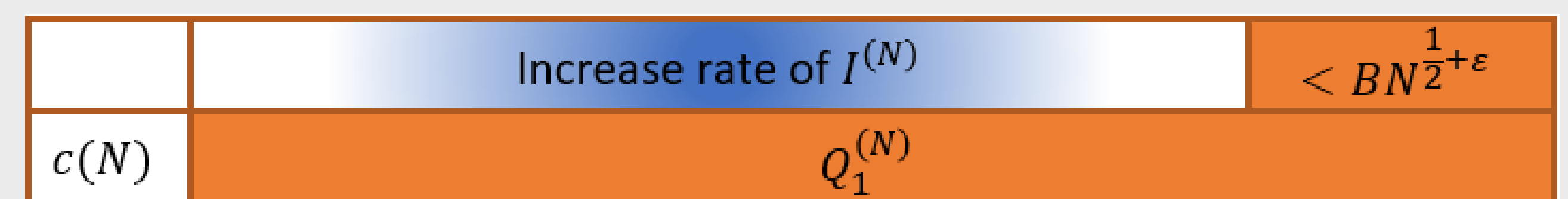


Fig. 2: Up-crossing ($c(N) < 2\beta N^{1/2-\varepsilon}$)

References

- [1] Patrick Eschenfeldt and David Gamarnik. "Join the shortest queue with many servers. The heavy traffic-asymptotics". In: *Math. Oper. Res.* 43.3 (2018), pp. 867–886.
- [2] Varun Gupta and Neil Walton. "Load balancing in the nondegenerate slowdown regime". In: *Oper. Res.* 67.1 (2019), pp. 281–294.
- [3] Daniela Hurtado-Lange and Siva Theja Maguluri. "Load balancing system under Join the Shortest Queue: Many-server-heavy-traffic asymptotics". In: *arXiv:2004.04826* (2020).
- [4] Xin Liu and Lei Ying. "A simple steady-state analysis of load balancing algorithms in the sub-Halfin-Whitt regime". In: *ACM SIGMETRICS Perform. Eval. Rev.* 46.2 (2019), pp. 15–17.
- [5] Xin Liu and Lei Ying. "Universal scaling of distributed queues under load balancing in the super-Halfin-Whitt regime". In: *IEEE/ACM Trans. Netw.* 30.1 (2022), pp. 190–201.
- [6] Debankur Mukherjee et al. "Universality of power-of-d load balancing in many-server systems". In: *Stoch. Syst.* 8.4 (2018), pp. 265–292.