

Exponential Tail Bounds on Queues

A confluence of non-asymptotic heavy traffic and large deviations

Daniela Hurtado-Lange
Northwestern University

Joint work with

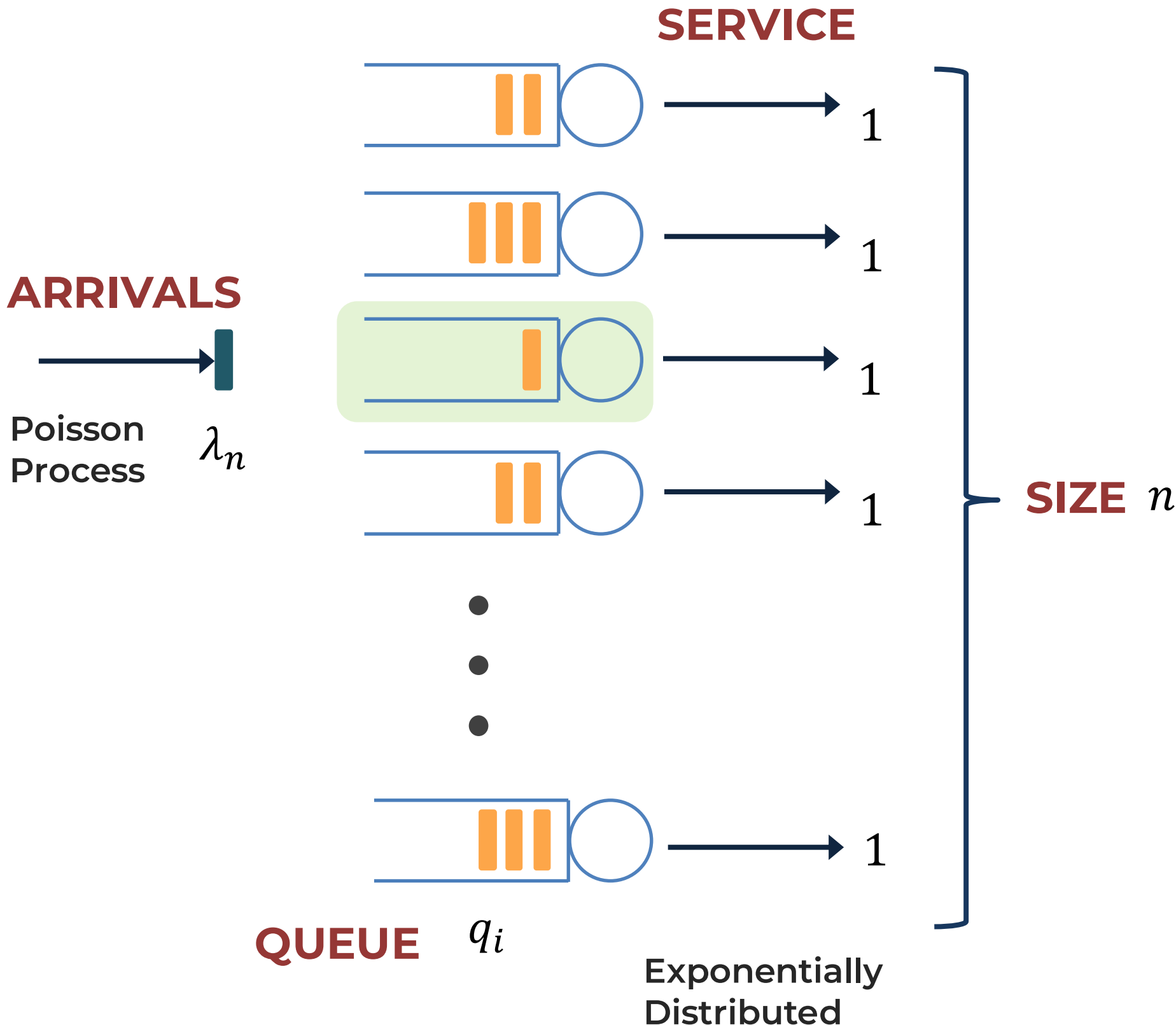


Prakirt Jhunjunwala
On the academic job
market!



Siva Theja Maguluri

Join-the-Shortest Queue (JSQ)



Throughput Optimality:

Stable: $\lambda_n < n$

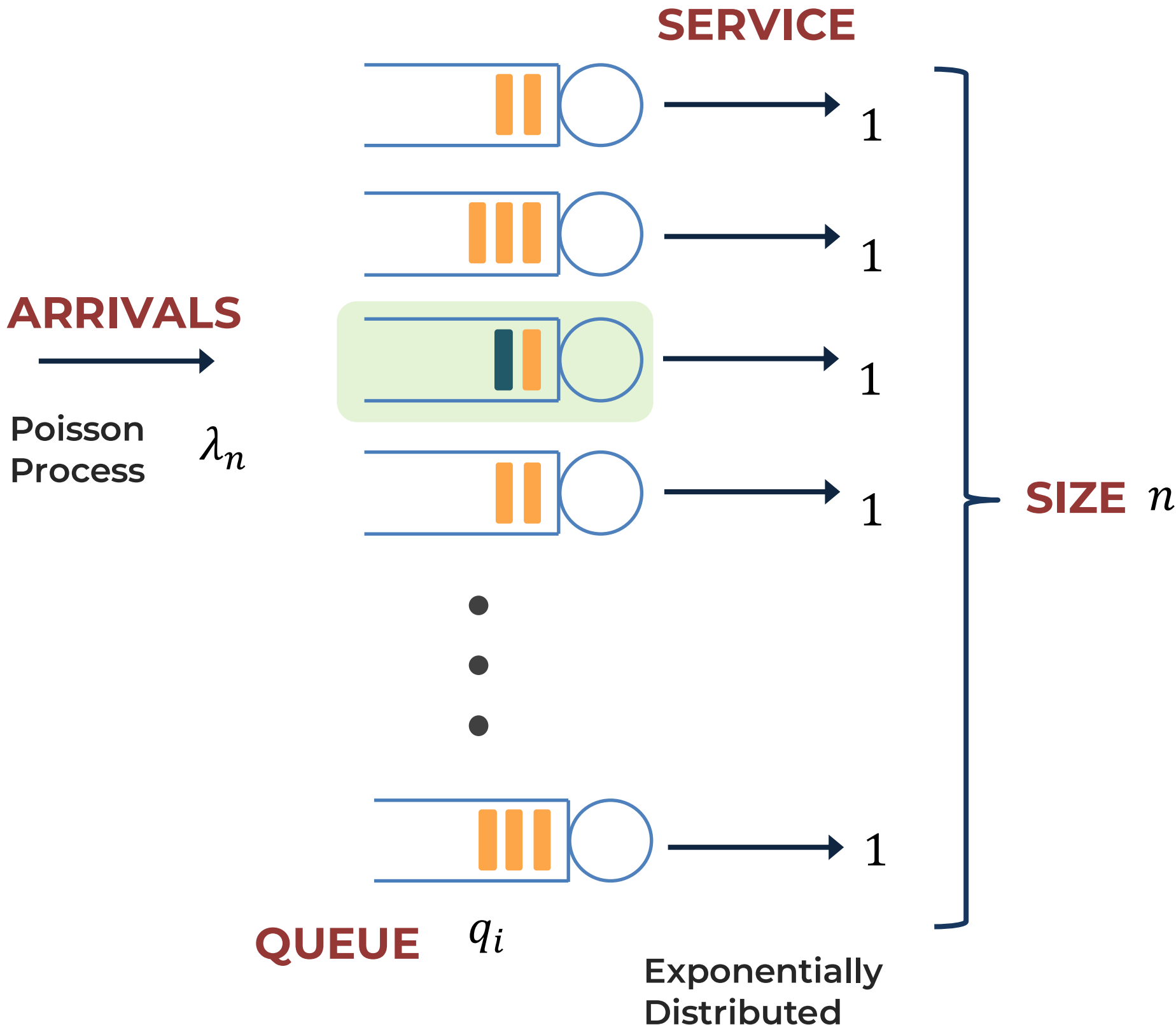
Algorithm Optimality:

Join-the-Shortest Queue (JSQ) **maximizes** the number of customers served

[Winston 1977]

[Weber 1978]

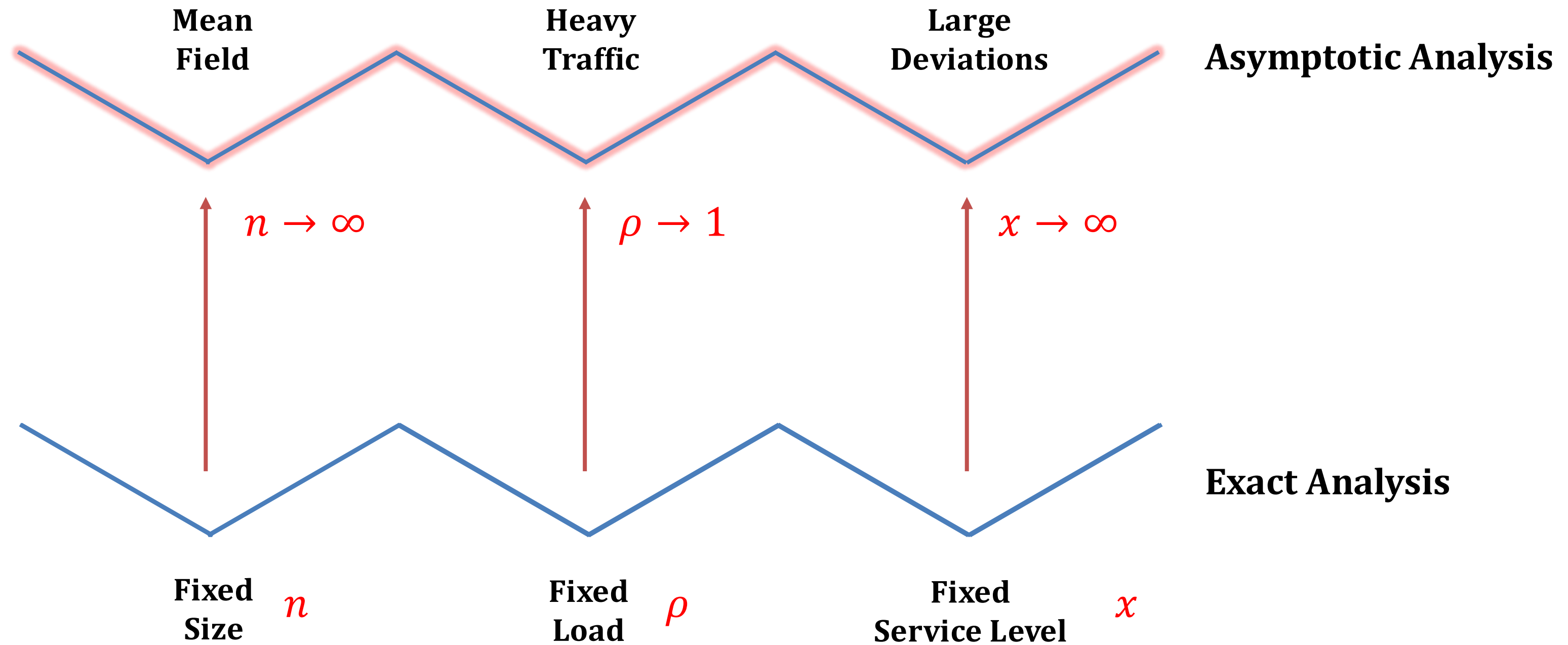
Join-the-Shortest Queue (JSQ)



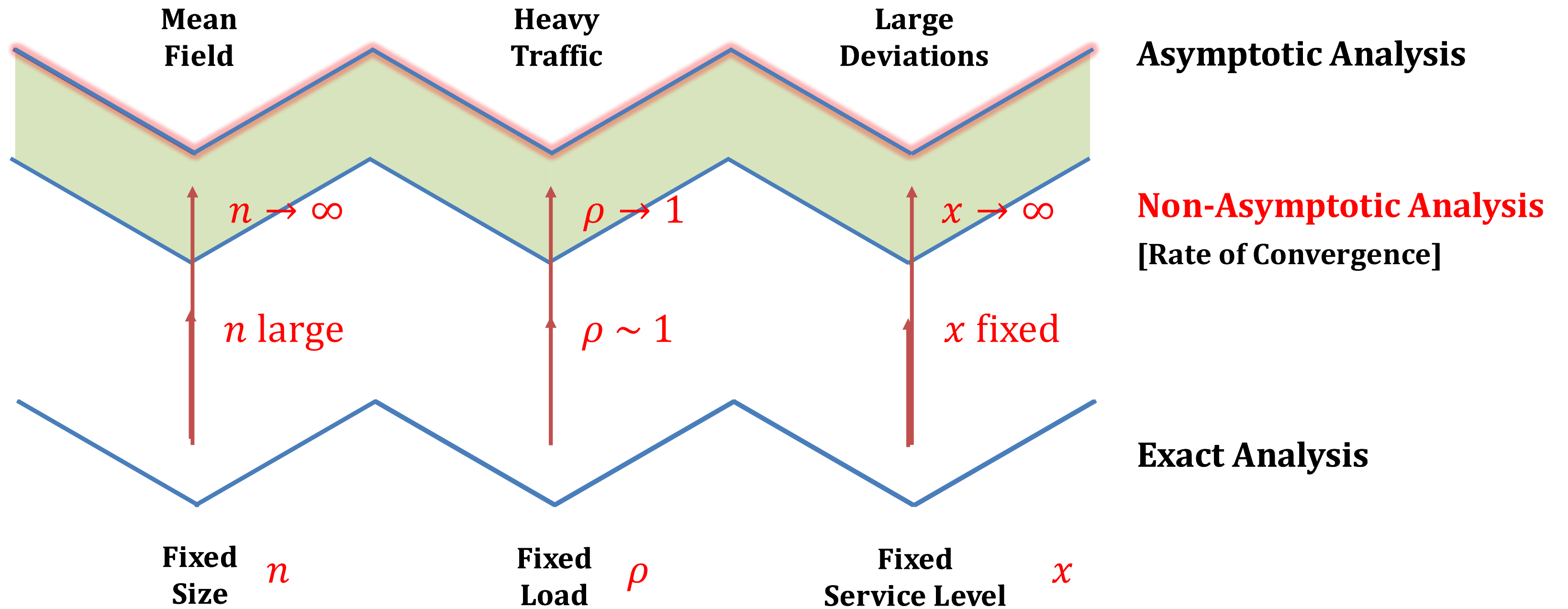
Analysis Questions

1. How does the system behave w.r.t. **SIZE**?
Size = n
2. How does the system behave w.r.t. **LOAD**?
Load = $\rho = \frac{\lambda_n}{n}$
3. How system's **SERVICE LEVEL** change?
 $P(\text{Queue} > x)$

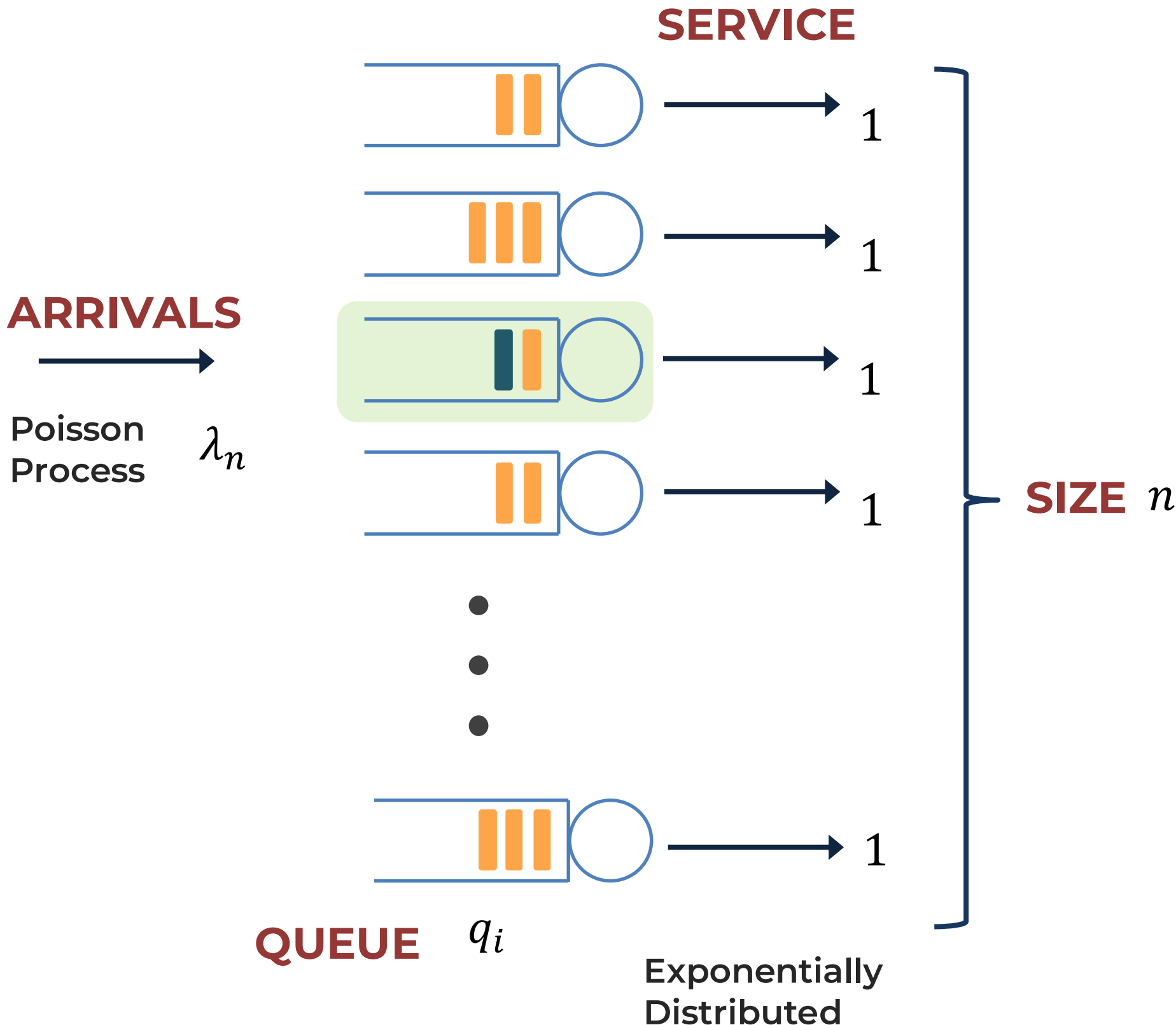
Asymptotic Analysis



Non-Asymptotic Analysis

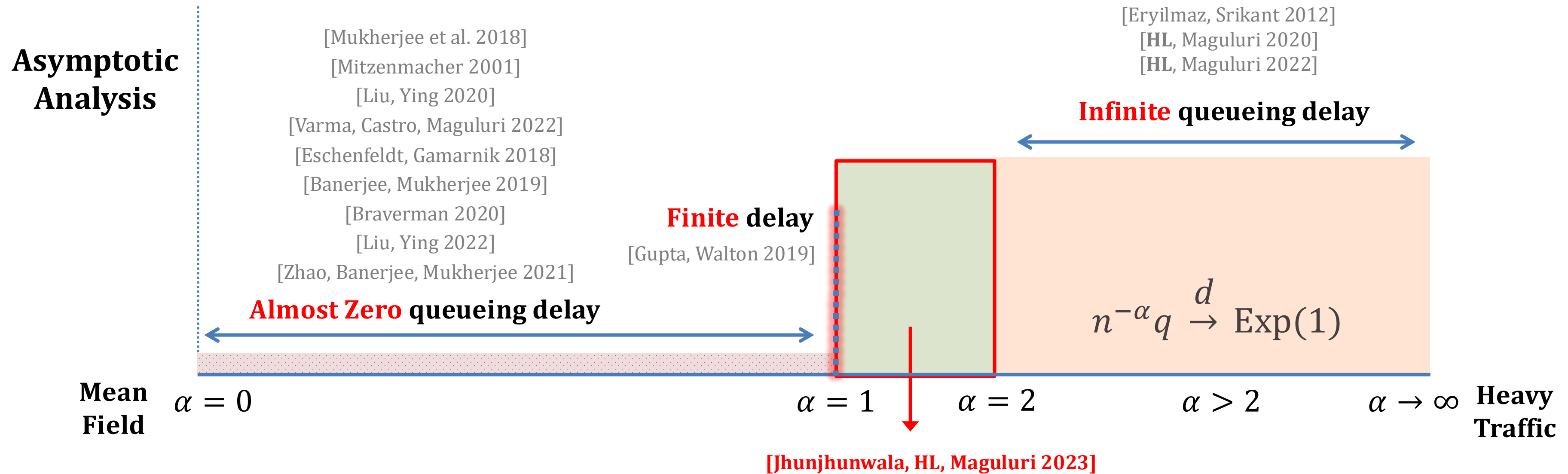


Join-the-Shortest Queue (JSQ) Parametrization



- Heavy Traffic (HT) parameter:**
- Stability:** $\lambda_n < n$
 - Many Server-HT:** $\lambda_n = n - n^{1-\alpha}, \alpha > 0$
 - Super-Slowdown:** $\alpha > 1$
 - Classic-HT:** Fix $n, \lambda = n(1 - \epsilon)$
- $\lambda_n = n - n^{1-\alpha}, \alpha \rightarrow \infty$

Related Work



Non-Asymptotic Analysis

[Braverman 2023],
[HL, Maguluri 2022],
[HL, Varma, Maguluri, 2021],
[Mukherjee 2022]

$$E[f(\epsilon q)] = \text{Limiting Value} + O(\epsilon)$$

$$P(\epsilon q > x) = \text{Limiting Value} + O(\epsilon)$$

Main Contribution: Informal

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime, with large n

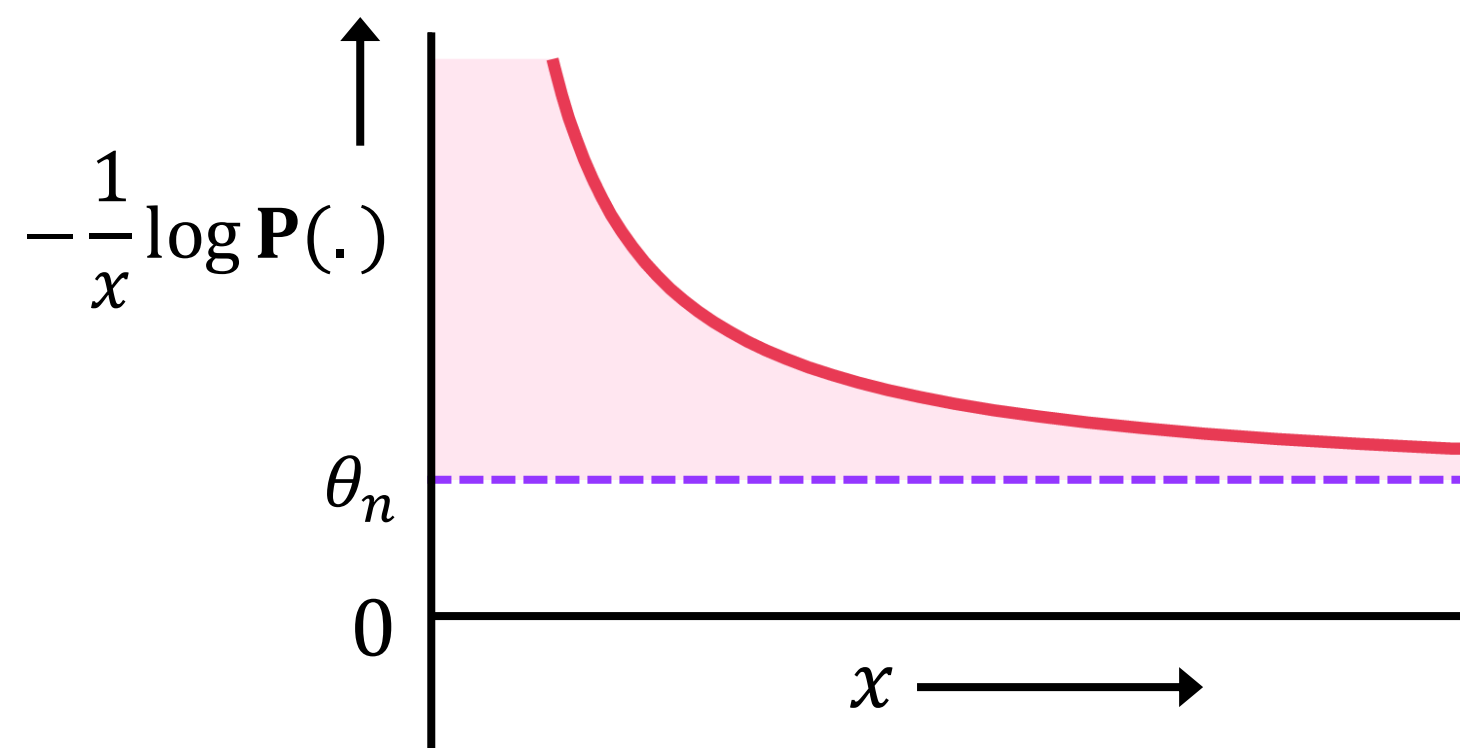
$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(\text{Scaled Queue} > x) \leq [1 + \tilde{O}(n^{1-\alpha})] 2ex \text{Exp}(-\theta_n x)$$

Main Contribution: Informal

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime, with large n

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(\text{Scaled Queue} > x) \leq [1 + \tilde{O}(n^{1-\alpha})] 2ex \text{Exp}(-\theta_n x)$$



Large Deviation:

$$\lim_{x \rightarrow \infty} -\frac{1}{x} \log \mathbf{P}(\cdot) = \theta_n$$

Many-Servers

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in **super slowdown regime** [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(\text{Scaled Queue} > x) \leq [1 + \tilde{O}(n^{1-\alpha})] 2ex \text{Exp}(-\theta_n x)$$

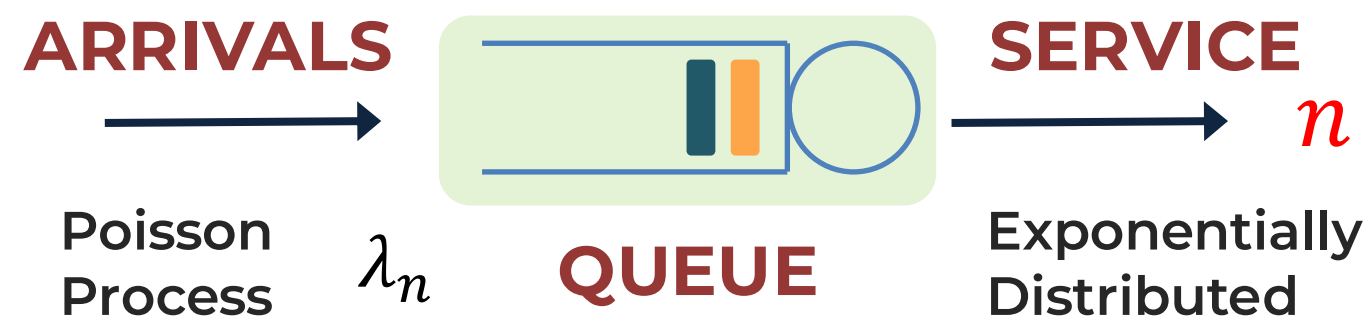
Lower Bound

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(\text{Scaled Queue} > x) \leq [1 + \tilde{O}(n^{1-\alpha})] 2ex \text{Exp}(-\theta_n x)$$

Pool all servers together



Lower bound:

$$\text{For all } n, \quad \mathbf{P}(n^{-\alpha} q > x) \geq e^{-\theta_n x}$$

$$\text{Tail decay rate for SSQ: } \theta_n = n^\alpha \log \frac{1}{1 - n^{-\alpha}}$$

State Space Collapse

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(n^{-\alpha} q > x) \leq [1 + \tilde{O}(n^{1-\alpha})] 2ex \text{Exp}(-\theta_n x)$$



**State Space
Collapse Violation**

State Space Collapse

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(n^{-\alpha} q > x) \leq [1 + K n^{1-\alpha} \log n] 2ex \text{Exp}(-\theta_n x)$$

JSQ with $n = 2$



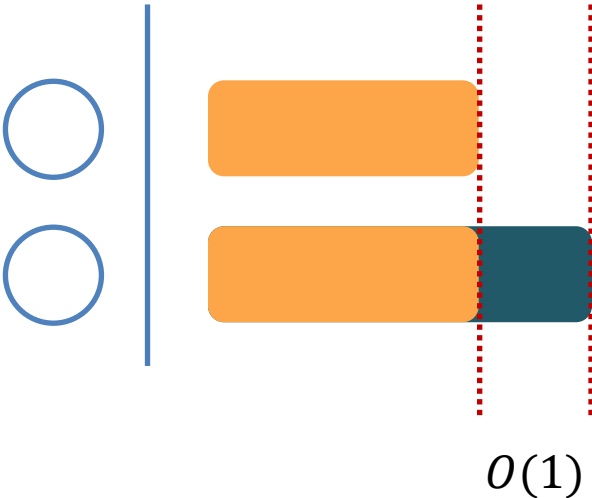
State Space Collapse

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(n^{-\alpha} q > x) \leq [1 + K n^{1-\alpha} \log n] 2ex \text{Exp}(-\theta_n x)$$

JSQ with $n = 2$



State Space Collapse

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(n^{-\alpha} q > x) \leq [1 + K n^{1-\alpha} \log n] 2ex \text{Exp}(-\theta_n x)$$

JSQ with $n = 2$

[Eryilmaz and Srikant 2012]



State Space Collapse

JSQ with $n = 2$

[Eryilmaz and Srikant 2012]



JSQ with n large

[Jhunjunwala, HL and Maguluri 2023] $O(1)$



Improvement: Previously $O(n)$ [Hurtado-Lange and Maguluri 2020]

State Space Collapse

LEMMA

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{MGF}(Q_{SSQ}) \leq \text{MGF}(Q_{JSQ}) \leq [1 + K n^{1-\alpha} \log n] \text{MGF}(Q_{SSQ})$$

$\theta \in (0, \theta_n)$

JSQ with n large

[Jhunjunwala, HL and Maguluri 2023] $O(1)$



Transform Method

$$n^{-\alpha} Q_{\text{JSQ}} \approx n^{-\alpha} (Q_{\text{SSQ}} + O(n))$$



Transform Method

LEMMA

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{MGF}(Q_{\text{SSQ}}) \leq \text{MGF}(Q_{\text{JSQ}}) \leq [1 + K n^{1-\alpha} \log n] \text{MGF}(Q_{\text{SSQ}})$$

$\theta \in (0, \theta_n)$

Transform Method

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(n^{-\alpha} q > x) \leq [1 + K n^{1-\alpha} \log n] 2 \text{ex} \text{Exp}(-\theta_n x)$$

LEMMA

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{MGF}(Q_{SSQ}) \leq \text{MGF}(Q_{JSQ}) \leq [1 + K n^{1-\alpha} \log n] \text{MGF}(Q_{SSQ})$$

$$\theta \in (0, \theta_n)$$

Transform Method

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(n^{-\alpha} q > x) \leq [1 + Kn^{1-\alpha} \log n] \text{ 2ex } \text{Exp}(-\theta_n x)$$

Pre-exponent
term

Pre-limit tail

[Markov's Inequality]

$$\mathbf{P}(n^{-\alpha} q > x) \leq [1 + Kn^{1-\alpha} \log n] \inf_{\theta \in (0, \theta_n)} e^{-\theta x} \text{MGF}(Q_{\text{SSQ}})$$

Main Contribution: Formal

THEOREM: Non-Asymptotic Tail Bound [JHLM23]

JSQ in super slowdown regime [$\lambda_n = n - n^{1-\alpha}$, $\alpha > 1$], with large n [$n^{1-\alpha} \log n$ is small]

$$\text{Exp}(-\theta_n x) \leq \mathbf{P}(n^{-\alpha} q > x) \leq [1 + Kn^{1-\alpha} \log n] 2ex \text{Exp}(-\theta_n x)$$

**Non-Asymptotic
Many Server-HT**

$$[\lambda_n = n - n^{1-\alpha}, \alpha > 1]$$

$n^{1-\alpha} \log n$ is small

**Non-Asymptotic
Large Deviation**

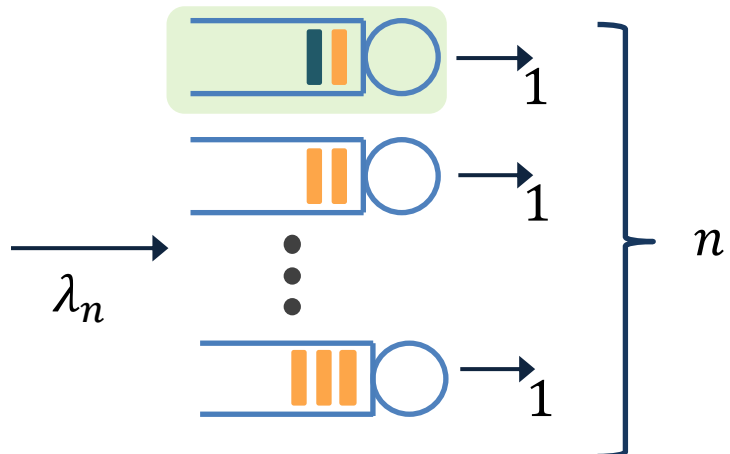

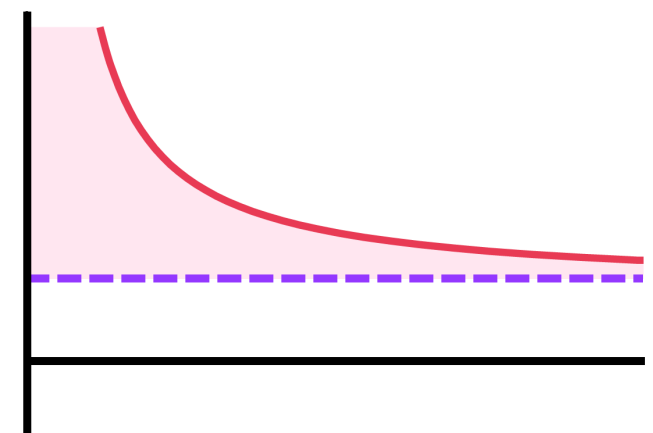
Pre-Exponent Terms

$$[1 + Kn^{1-\alpha} \log n] 2ex$$

**Rate of Convergence of
Tail Decay**

$$\theta_n = n^\alpha \log \frac{1}{1 - n^{-\alpha}}$$

SUMMARY

	<p>JSQ: Many-Server $\xrightarrow{\lambda_n}$ </p>	
Regime	<p>Many-Server-HT: $\lambda_n = n - n^{1-\alpha}, \alpha > 1$</p>	
	<p>Asymptotic: $n \rightarrow \infty$</p>	<p>Non-Asymptotic: n is large</p>
SSC	<p>1-dim</p>	<p>SSC Violation: $[1 + O(n^{1-\alpha} \log n)]$</p>
Transform Eq.	<p>Explicit Linear Eq.</p>	<p>MGF Bound</p>
Limiting Distribution	<p> Exponential</p>	<p>Non-Asymptotic Tail Bounds and Large Deviation </p>