# Heavy-traffic Delay Optimality in Pull-based Load Balancing Systems: Necessary and Sufficient Conditions

Xingyu Zhou



Feb. 7 2019 @ Caltech

### Joint work with...



Jian Tan, Alibaba



Ness Shroff, OSU

### The cloud is everywhere...











- Availability
- Throughput



- Availability
- Throughput
- Responsiveness
  - G: a half second increase in loading time drops traffic by 20%.
  - **3**: every 0.1s of latency costs them 1% in sales.







Discrete-time system, i.e, time-slotted.

<sup>&</sup>lt;sup>1</sup>with all moments bounded



Discrete-time system, i.e, time-slotted.

• Arrival rate at each time slot is  $\lambda_{\Sigma}$ , general distribution <sup>1</sup>.

#### <sup>1</sup>with all moments bounded



- Discrete-time system, i.e, time-slotted.
- Arrival rate at each time slot is  $\lambda_{\Sigma}$ , general distribution <sup>1</sup>.
- Service rate at each server k is  $\mu_k$ , general distribution.

#### <sup>1</sup>with all moments bounded



- Discrete-time system, i.e, time-slotted.
- Arrival rate at each time slot is  $\lambda_{\Sigma}$ , general distribution  $^1$  .
- Service rate at each server k is  $\mu_k$ , general distribution.
- Arrival and service are independent.

<sup>&</sup>lt;sup>1</sup>with all moments bounded



choose the *right* server(s) for each request.



choose the *right* server(s) for each request.



choose the *right* server(s) for each request.

What does *right* mean?

High throughput



choose the *right* server(s) for each request.

- High throughput
- Low latency



choose the *right* server(s) for each request.

- High throughput
- Low latency
  - 1. dispatching time (arrive at dispatcher joining a server).



choose the *right* server(s) for each request.

- High throughput
- Low latency
  - 1. dispatching time (arrive at dispatcher joining a server).
  - 2. response time (joining a server leaving the server).
    - by Little's law, minimize the mean number of requests in system.

Which load balancing policy is the best?

Which load balancing policy is the best?

Maybe the most intuitive one: Join the Shortest Queue

Latency = dispatching time + response time (delay).

#### Latency = dispatching time + response time (delay).

- ► 40 years ago...
  - ► Foschini, et. al, proved that JSQ is delay optimal in heavy traffic.

#### Latency = dispatching time + response time (delay).

#### ► 40 years ago...

- ► Foschini, et. al, proved that JSQ is delay optimal in heavy traffic.
- Initiated by dispatcher (push-based), non-zero dispatching time.
- Often impractical to implement due to high message overhead.

#### Latency = dispatching time + response time (delay).

### ► 40 years ago...

- ► Foschini, et. al, proved that JSQ is delay optimal in heavy traffic.
- Initiated by dispatcher (push-based), non-zero dispatching time.
- Often impractical to implement due to high message overhead.

#### 25 years ago...

- Frank Kelly, et. al, proposed a threshold type policy JBT.
  - ► Join-Below-Threshold (*r*): randomly join a server whose queue length is below *r*, if any; otherwise, randomly select any one to join.

#### Latency = dispatching time + response time (delay).

### ► 40 years ago...

- ► Foschini, et. al, proved that JSQ is delay optimal in heavy traffic.
- Initiated by dispatcher (push-based), non-zero dispatching time.
- Often impractical to implement due to high message overhead.

#### 25 years ago...

- Frank Kelly, et. al, proposed a threshold type policy JBT.
  - ► Join-Below-Threshold (*r*): randomly join a server whose queue length is below *r*, if any; otherwise, randomly select any one to join.
- Initiated by server (pull-based), zero dispatching time.
- The message overhead is at most 1 per arrival.

#### Latency = dispatching time + response time (delay).

#### ► 40 years ago...

- ► Foschini, et. al, proved that JSQ is delay optimal in heavy traffic.
- Initiated by dispatcher (push-based), non-zero dispatching time.
- Often impractical to implement due to high message overhead.

#### 25 years ago...

- Frank Kelly, et. al, proposed a threshold type policy JBT.
  - ▶ Join-Below-Threshold (*r*): randomly join a server whose queue length is below *r*, if any; otherwise, randomly select any one to join.
- Initiated by server (pull-based), zero dispatching time.
- The message overhead is at most 1 per arrival.
- They conjectured that if

 $r \geq K \log(\text{average number of requests in system})$ 

then JBT is delay optimal in heavy traffic.

#### Latency = dispatching time + response time (delay).

#### ► 40 years ago...

- ► Foschini, et. al, proved that JSQ is delay optimal in heavy traffic.
- Initiated by dispatcher (push-based), non-zero dispatching time.
- Often impractical to implement due to high message overhead.

#### 25 years ago...

- Frank Kelly, et. al, proposed a threshold type policy JBT.
  - ▶ Join-Below-Threshold (*r*): randomly join a server whose queue length is below *r*, if any; otherwise, randomly select any one to join.
- Initiated by server (pull-based), zero dispatching time.
- The message overhead is at most 1 per arrival.
- They conjectured that if

 $r \geq K \log(\text{average number of requests in system})$ 

then JBT is delay optimal in heavy traffic.

#### From 1993...

- Logarithmic growth rate of threshold guarantees heavy-traffic delay optimality in several scheduling settings [Harrison'98], [Williams, et al'01].
- However, it is still open in the load balancing setting.

We present necessary and sufficient conditions on the threshold for delay optimality in heavy traffic in **load balancing systems**.



We present necessary and sufficient conditions on the threshold for delay optimality in heavy traffic in **load balancing systems**.

► Theoretical contributions:



We present necessary and sufficient conditions on the threshold for delay optimality in heavy traffic in **load balancing systems**.

#### ► Theoretical contributions:

- Resolves the long-standing open problem mentioned before.
  - > The results are more general compared to Kelly's original conjecture.



We present necessary and sufficient conditions on the threshold for delay optimality in heavy traffic in **load balancing systems**.

#### Theoretical contributions:

- Resolves the long-standing open problem mentioned before.
  - The results are more general compared to Kelly's original conjecture.
- Provides new insights into heavy-traffic delay optimality.
  - The 'King' equation.



We present necessary and sufficient conditions on the threshold for delay optimality in heavy traffic in **load balancing systems**.

#### Theoretical contributions:

- Resolves the long-standing open problem mentioned before.
  - The results are more general compared to Kelly's original conjecture.

47

- Provides new insights into heavy-traffic delay optimality.
  - The 'King' equation.
- Develops new techniques for the analysis of load balancing policies.
  - New type of state-space collapse.

Queueing Systems 13 (1993) 47-86

Dynamic routing in open queueing networks: Brownian models, cut constraints and resource pooling

F.P. Kelly and C.N. Laws\* Statistical Laboratory, University of Cambridge, 16 Mill Lane, Cambridge CB2 1SB, England

The necessary and sufficient conditions on the threshold have...

#### Practical contributions:

- Provides a simple guideline for practical systems, e.g., Netflix Zuul.
- Sheds light on the design of new pull-based algorithms.

M   N THE NETFLIX Follow	
	problems.
ё́ў 2.4К Д	. When possible, instead of configuring static thresholds, use adaptive mechanisms that change based on the current traffic, performance and environment.



Part I: Background

# Low delay...
#### Low delay...

► Closed-form formula in classical regime is very difficult.



#### Low delay...

- Closed-form formula in classical regime is very difficult.
- ► Turn to asymptotic regimes for insights.



#### Low delay...

- Closed-form formula in classical regime is very difficult.
- ► Turn to asymptotic regimes for insights.



#### Low delay ...

- Closed-form formula in classical regime is very difficult.
- Turn to asymptotic regimes for insights.



#### Low delay ...

- Closed-form formula in classical regime is very difficult.
- Turn to asymptotic regimes for insights.



#### Why is the heavy-traffic regime?

#### Why is the heavy-traffic regime?

"In heavy-traffic regime, the important features of good control policies are often displayed in the sharpest relief"

- Frank Kelly



#### Heavy-traffic Delay Optimality



Fact:  $\mathbb{E}\left[\sum Q_n\right] \geq \mathbb{E}\left[q\right]$ , since one service process is stochastically larger.

#### Heavy-traffic Delay Optimality

#### Definition

It can achieve the lower bound on delay when  $\epsilon \to 0$  ( $\epsilon = \sum \mu_n - \lambda_{\Sigma}$ ), that is,  $\lim_{\epsilon \downarrow 0} \epsilon \mathbb{E} \left[ \sum Q_n \right] = \lim_{\epsilon \downarrow 0} \epsilon \mathbb{E} \left[ q \right]$  (the queue length is on the order  $O(1/\epsilon)$ )



Fact:  $\mathbb{E}\left[\sum Q_n\right] \geq \mathbb{E}\left[q\right]$ , since one service process is stochastically larger.

Push algorithm: Join-shortest-queue (JSQ)

- sample each queue length
- ► join the shortest one



Push algorithm: Join-shortest-queue (JSQ)

- sample each queue length
- join the shortest one



Pros:

- > Delay optimal in a stochastic order sense. [Weber'78]
- Heavy-traffic delay optimal. [Foschini and Salz'78], [Eryilmaz and Srikant'12]

Push algorithm: Join-shortest-queue (JSQ)

- sample each queue length
- join the shortest one



Pros:

- > Delay optimal in a stochastic order sense. [Weber'78]
- Heavy-traffic delay optimal. [Foschini and Salz'78], [Eryilmaz and Srikant'12]

Cons:

- Message overhead is undesriable (2N per arrival).
- Non-zero dispatching delay.

Push algorithm: Power-of-d (Pod)

- randomly sample the queue lengths of d servers.
- join the shortest one among them.



Push algorithm: Power-of-d (Pod)

- randomly sample the queue lengths of d servers.
- join the shortest one among them.



#### Pros:

- Double exponential decay in delay tail when N is large. [Mitzenmacher'96]
- Heavy-traffic delay optimal for homogeneous servers. [Chen and Ye'12], [Maguluri, et al'14]
- Improved message overhead (2d per arrival)

Push algorithm: Power-of-d (Pod)

- randomly sample the queue lengths of d servers.
- join the shortest one among them.



#### Pros:

- Double exponential decay in delay tail when N is large. [Mitzenmacher'96]
- Heavy-traffic delay optimal for homogeneous servers. [Chen and Ye'12], [Maguluri, et al'14]
- Improved message overhead (2d per arrival)

#### Cons:

Non-zero dispatching delay.

Pull algorithm: Join-idle-queue (JIQ)

- if possible, join an idle queue randomly.
- > otherwise, join an arbitrary queue randomly.

Pull algorithm: Join-idle-queue (JIQ)

- if possible, join an idle queue randomly.
- otherwise, join an arbitrary queue randomly.

It is a pull algorithm since it behaves like the idle server pulls tasks from the dispatcher.



Pull algorithm: Join-idle-queue (JIQ)

- if possible, join an idle queue randomly.
- otherwise, join an arbitrary queue randomly.

It is a pull algorithm since it behaves like the idle server pulls tasks from the dispatcher.



#### Pros:

- Better delay performance than Pod with a lower message overhead (at most 1 per arrival), when traffic is not heavy. [Lu, et al'11], [Stolyar'15]
- Zero dispatching delay

Pull algorithm: Join-idle-queue (JIQ)

- if possible, join an idle queue randomly.
- otherwise, join an arbitrary queue randomly.

It is a pull algorithm since it behaves like the idle server pulls tasks from the dispatcher.



#### Pros:

- Better delay performance than Pod with a lower message overhead (at most 1 per arrival), when traffic is not heavy. [Lu, et al'11], [Stolyar'15]
- Zero dispatching delay

#### Cons:

Delay performance downgrades substantially under heavy traffic.

▶ Push algorithms, e.g., JSQ and Pod

▶ Pull algorithm, e.g., JIQ

Push algorithms, e.g., JSQ and Pod

elay optimal in heavy traffic.

▶ Pull algorithm, e.g., JIQ

Push algorithms, e.g., JSQ and Pod

elay optimal in heavy traffic.

- Pull algorithm, e.g., JIQ
  - very poor delay in heavy traffic.

- elay optimal in heavy traffic.
- non-zero dispatching delay.

- Pull algorithm, e.g., JIQ
  - very poor delay in heavy traffic.

- elay optimal in heavy traffic.
- inon-zero dispatching delay.

- Pull algorithm, e.g., JIQ
  - very poor delay in heavy traffic.
  - zero dispatching delay.

- Gelay optimal in heavy traffic.
- inon-zero dispatching delay.
- high (relatively) message overhead.
- Pull algorithm, e.g., JIQ
  - very poor delay in heavy traffic.
  - zero dispatching delay.

- Gelay optimal in heavy traffic.
- inon-zero dispatching delay.
- high (relatively) message overhead.
- Pull algorithm, e.g., JIQ
  - very poor delay in heavy traffic.
  - cero dispatching delay.
  - iow message overhead.

- delay optimal in heavy traffic.
- inno-zero dispatching delay. (is by nature of push algorithms)
- high (relatively) message overhead.
- Pull algorithm, e.g., JIQ
  - very poor delay in heavy traffic.
  - cero dispatching delay.
  - iow message overhead.

- e delay optimal in heavy traffic.
- inno-zero dispatching delay. (is by nature of push algorithms)
- high (relatively) message overhead.
- Pull algorithm, e.g., JIQ
  - ▶ 😌 very poor delay in heavy traffic. (⊖ some hope here?)
  - ero dispatching delay.
  - ▶ 🙂 low message overhead.

Can we design a heavy-traffic delay optimal pull-based policy?

Can we design a heavy-traffic delay optimal pull-based policy?

Main idea: A dynamic threshold!

Can we design a heavy-traffic delay optimal pull-based policy?

Main idea: A dynamic threshold!

The hope is that:

- instead of only storing idle servers.
- the dispatcher stores servers with queue lengths being less than a dynamic threshold!

Part II: Necessary Conditions

# (a) Is JIQ delay optimal in heavy traffic?(A). Yes(B). No



(a) Is JIQ delay optimal in heavy traffic?

(A). Yes (B). No

(b) In the heavy-traffic limit: is the delay under JIQ the same as that under Random?

(A). Yes (B). No



(a) Is JIQ delay optimal in heavy traffic?

(A). Yes (B). No

(b) In the heavy-traffic limit: is the delay under JIQ the same as that under Random?

(A). Yes (B). No

Both answers are NO!



(a) Is JIQ delay optimal in heavy traffic?

(A). Yes (B). No

(b) In the heavy-traffic limit: is the delay under JIQ the same as that under Random?

(A). Yes (B). No

Both answers are NO!

In fact, in the heavy-traffic limit:



```
heavy-traffic optimal = JSQ < JIQ < Rand
```
The JBT(r) (Join-Below-Threshold(r)) algorithm:

The JBT(r) (Join-Below-Threshold(r)) algorithm:

1. the memory at the dispatcher stores IDs of servers with queue lengths being less than *r*. (mechanism will be introduced later)

The JBT(r) (Join-Below-Threshold(r)) algorithm:

- 1. the memory at the dispatcher stores IDs of servers with queue lengths being less than *r*. (mechanism will be introduced later)
- 2. if the memory is non-empty, randomly picks a ID and join the server.

The JBT(r) (Join-Below-Threshold(r)) algorithm:

- 1. the memory at the dispatcher stores IDs of servers with queue lengths being less than *r*. (mechanism will be introduced later)
- 2. if the memory is non-empty, randomly picks a ID and join the server.
- 3. otherwise, randomly picks a queue to join.

The JBT(r) (Join-Below-Threshold(r)) algorithm:

- 1. the memory at the dispatcher stores IDs of servers with queue lengths being less than *r*. (mechanism will be introduced later)
- 2. if the memory is non-empty, randomly picks a ID and join the server.
- 3. otherwise, randomly picks a queue to join.

Remark:

▶ JIQ is just a special case of JBT(r) with constant r = 1.

The JBT(r) (Join-Below-Threshold(r)) algorithm:

- 1. the memory at the dispatcher stores IDs of servers with queue lengths being less than *r*. (mechanism will be introduced later)
- 2. if the memory is non-empty, randomly picks a ID and join the server.
- 3. otherwise, randomly picks a queue to join.

#### Remark:

- ▶ JIQ is just a special case of JBT(r) with constant r = 1.
- ► For heterogeneous servers, we can just replace random selection with selection in proportion to the service rate.

## Geometry of JBT...

The JBT(r) (Join-Below-Threshold(r)) algorithm:

- 1. the memory at the dispatcher stores IDs of servers with queue lengths being less than *r*. (mechanism will be introduced later)
- 2. if the memory is non-empty, randomly picks an ID and join the server.
- 3. otherwise, randomly picks a queue to join.



 $\mathbf{Q} \in \mathcal{R}_{I}$ : Random (full memory)  $\mathbf{Q} \in \mathcal{R}_{u}$ : Random (empty memory)

## Geometry of JBT...

The JBT(r) (Join-Below-Threshold(r)) algorithm:

- 1. the memory at the dispatcher stores IDs of servers with queue lengths being less than *r*. (mechanism will be introduced later)
- 2. if the memory is non-empty, randomly picks a ID and join the server.
- 3. otherwise, randomly picks a queue to join.



 $\mathbf{Q} \notin \mathcal{R}_{I} \cup \mathcal{R}_{u}$ : shorter queues are preferred.

Grow, but not too fast...

### Theorem (Necessary Conditions)

Consider the JBT(r) policy.

1. For any constant threshold r, we have the following average delay ordering in heavy traffic:

$$\overline{D}_{JSQ} < \overline{D}_{JBT(r)} < \overline{D}_{Rand}$$

2. For  $r = \Omega((1/\epsilon)^{1+\alpha})$  for any  $\alpha > 0$ , we have that in heavy traffic:

 $\overline{D}_{JSQ} < \overline{D}_{JBT(r)} = \overline{D}_{Rand}$ 

Before the proof, any intuitions?

The sufficient and necessary condition for HT-optimality:

$$\lim_{\epsilon \downarrow 0} \mathbb{E} \left[ \left\| \overline{\mathbf{Q}}^{(\epsilon)}(t+1) \right\|_1 \left\| \overline{\mathbf{U}}^{(\epsilon)}(t) \right\|_1 \right] = 0.$$

where the unused service vector  $\mathbf{U}(t) = \max{\{\mathbf{S}(t) - \mathbf{Q}(t) - \mathbf{A}(t), \mathbf{0}\}}$ .

The sufficient and necessary condition for HT-optimality:

$$\lim_{\epsilon \downarrow 0} \mathbb{E} \left[ \left\| \overline{\mathbf{Q}}^{(\epsilon)}(t+1) \right\|_1 \left\| \overline{\mathbf{U}}^{(\epsilon)}(t) \right\|_1 \right] = 0.$$

where the unused service vector  $\mathbf{U}(t) = \max{\{\mathbf{S}(t) - \mathbf{Q}(t) - \mathbf{A}(t), \mathbf{0}\}}$ .

• Note that  $Q_n(t+1)U_n(t) = 0$  for all n and t.

The sufficient and necessary condition for HT-optimality:

$$\lim_{\epsilon \downarrow 0} \mathbb{E} \left[ \left\| \overline{\mathbf{Q}}^{(\epsilon)}(t+1) \right\|_1 \left\| \overline{\mathbf{U}}^{(\epsilon)}(t) \right\|_1 \right] = 0.$$

where the unused service vector  $\mathbf{U}(t) = \max{\{\mathbf{S}(t) - \mathbf{Q}(t) - \mathbf{A}(t), \mathbf{0}\}}$ .

• Note that  $Q_n(t+1)U_n(t) = 0$  for all *n* and *t*.

IMPLICATIONS: No server is idle while others with high loads.

The sufficient and necessary condition for HT-optimality:

$$\lim_{\epsilon \downarrow 0} \mathbb{E}\left[\left\|\overline{\mathbf{Q}}^{(\epsilon)}(t+1)\right\|_1 \left\|\overline{\mathbf{U}}^{(\epsilon)}(t)\right\|_1\right] = 0.$$

where the unused service vector  $\mathbf{U}(t) = \max{\{\mathbf{S}(t) - \mathbf{Q}(t) - \mathbf{A}(t), \mathbf{0}\}}$ .

Note that  $Q_n(t+1)U_n(t) = 0$  for all *n* and *t*.

IMPLICATIONS: No server is idle while others with high loads.

"Probability theory is nothing but common sense reduced to calculation."

— Pierre Laplace





Completely degenerate to random.

$$\overline{D}_{JSQ} < \overline{D}_{JBT(r)} = \overline{D}_{Rand}$$

 $\mathbf{Q} \in \mathcal{R}_{I}$ : Random (full memory)



Completely degenerate to random.

$$\overline{D}_{JSQ} < \overline{D}_{JBT(r)} = \overline{D}_{Rand}$$

 $\mathbf{Q} \in \mathcal{R}_{I}$ : Random (full memory)



Always has 'nice' things happen.

$$\overline{D}_{JSQ} < \overline{D}_{JBT(r)} < \overline{D}_{Rand}$$
$$\mathbf{Q} \notin \mathcal{R}_{I} \cup \mathcal{R}_{u} : \text{ `nice' things}$$

## The Universal Equality...

$$\mathbb{E}\left[\left\|\overline{\mathbf{Q}}^{(\epsilon)}(t+1)\right\|_{1}\left\|\overline{\mathbf{U}}^{(\epsilon)}(t)\right\|_{1}\right] = \mathcal{T}_{1}^{(\epsilon)} + \mathcal{T}_{2}^{(\epsilon)} - \mathcal{T}_{3}^{(\epsilon)}$$

where

$$\begin{split} \mathcal{T}_{1}^{(\epsilon)} &\triangleq 2\sum_{i=1}^{N}\sum_{j>i}^{N} \mathbb{E}\left[\left(\overline{Q}_{i}^{(\epsilon)} - \overline{Q}_{j}^{(\epsilon)}\right)\left(\overline{A}_{i}^{(\epsilon)} - \overline{A}_{j}^{(\epsilon)} - \overline{S}_{i}^{(\epsilon)} + \overline{S}_{j}^{(\epsilon)}\right)\right] \\ \mathcal{T}_{2}^{(\epsilon)} &\triangleq \sum_{i=1}^{N}\sum_{j>i}^{N} \mathbb{E}\left[\left(\overline{A}_{i}^{(\epsilon)} - \overline{A}_{j}^{(\epsilon)} - \overline{S}_{i}^{(\epsilon)} + \overline{S}_{j}^{(\epsilon)}\right)^{2}\right] \\ \mathcal{T}_{3}^{(\epsilon)} &\triangleq \sum_{i=1}^{N}\sum_{j>i}^{N} \mathbb{E}\left[\left(\overline{U}_{i}^{(\epsilon)} - \overline{U}_{j}^{(\epsilon)}\right)^{2}\right] \\ \overline{\mathbf{Q}}^{+} &\triangleq \overline{\mathbf{Q}}(t+1) \end{split}$$

The first unified method to show a policy is NOT optimal.

Then...how fast should it grow?

Part III: Sufficient Conditions



Question: Which of the following r value guarantees 'optimality'?

(A). 
$$r = 100$$
  
(B).  $r = \theta (\log(1/\epsilon))$   
(C).  $r = \theta (\log^2(1/\epsilon))$   
(D).  $r = \theta ((1/\epsilon)^{1.01})$ 

Hint: The average number of tasks is on the order of  $1/\epsilon$ .



Question: Which of the following r value guarantees 'optimality'?

(A). 
$$r = 100$$
  
(B).  $r = \theta (\log(1/\epsilon))$   
(C).  $r = \theta (\log^2(1/\epsilon))$   
(D).  $r = \theta ((1/\epsilon)^{1.01})$ 

Hint: The average number of tasks is on the order of  $1/\epsilon$ .



Question: Which of the following r value guarantees 'optimality'?

(A). 
$$r = 100$$
  
(B).  $r = \theta (\log(1/\epsilon))$   
(C).  $r = \theta (\log^2(1/\epsilon))$   
(D).  $r = \theta ((1/\epsilon)^{1.01})$ 

Hint: The average number of tasks is on the order of  $1/\epsilon$ .



# 25 years ago...

## 25 years ago...

Queueing Systems 13 (1993) 47-86

47

# Dynamic routing in open queueing networks: Brownian models, cut constraints and resource pooling

#### F.P. Kelly and C.N. Laws\*

Statistical Laboratory, University of Cambridge, 16 Mill Lane, Cambridge CB2 1SB, England

#### Conjecture: 'optimality' is guaranteed if $r \ge K \log(1/\epsilon)$ .<sup>2</sup>

<sup>&</sup>lt;sup>2</sup>Two-server case and diffusion approximation optimality

## Theorem (Sufficient Conditions)

Consider the JBT(r) policy with threshold r. Suppose  $r \ge K \log(1/\epsilon)$ and  $r = o(1/\epsilon)$ , for some constant K, then it is heavy-traffic delay optimal in steady-state.

## Theorem (Sufficient Conditions)

Consider the JBT(r) policy with threshold r. Suppose  $r \ge K \log(1/\epsilon)$ and  $r = o(1/\epsilon)$ , for some constant K, then it is heavy-traffic delay optimal in steady-state.

Remark:

## Theorem (Sufficient Conditions)

Consider the JBT(r) policy with threshold r. Suppose  $r \ge K \log(1/\epsilon)$ and  $r = o(1/\epsilon)$ , for some constant K, then it is heavy-traffic delay optimal in steady-state.

#### Remark:

• This holds for any fixed finite number of servers,  $2 \le N < \infty$ .

## Theorem (Sufficient Conditions)

Consider the JBT(r) policy with threshold r. Suppose  $r \ge K \log(1/\epsilon)$ and  $r = o(1/\epsilon)$ , for some constant K, then it is heavy-traffic delay optimal in steady-state.

#### Remark:

- ▶ This holds for any fixed finite number of servers,  $2 \le N < \infty$ .
- This holds for general arrival and service distributions (with finite moment bounds).

## Theorem (Sufficient Conditions)

Consider the JBT(r) policy with threshold r. Suppose  $r \ge K \log(1/\epsilon)$ and  $r = o(1/\epsilon)$ , for some constant K, then it is heavy-traffic delay optimal in steady-state.

#### Remark:

- ▶ This holds for any fixed finite number of servers,  $2 \le N < \infty$ .
- This holds for general arrival and service distributions (with finite moment bounds).
- The optimality is directly in steady-state.

- (a) Diffusion approximations method.
  - ▶ 😌 multi-dimensional Brownian motion is difficult to analyze.

- (a) Diffusion approximations method.
  - ▶ 😔 multi-dimensional Brownian motion is difficult to analyze.
  - ▶ 😳 optimality often only exists in finite time, not steady-state result.

(a) Diffusion approximations method.

- ▶ 😌 multi-dimensional Brownian motion is difficult to analyze.
- optimality often only exists in finite time, not steady-state result.

(b) Lyapunov drift-based method.

(a) Diffusion approximations method.

- ▶ 😔 multi-dimensional Brownian motion is difficult to analyze.
- optimality often only exists in finite time, not steady-state result.
- (b) Lyapunov drift-based method.
  - ▶ 😌 standard techniques fail in our case.
    - since the state-space collapse is neither a line nor a cone.
    - instead, it is even non-convex.

Our methods:

(a) Diffusion approximations method.

- Witti-dimensional Brownian motion is difficult to analyze.
- optimality often only exists in finite time, not steady-state result.
- (b) Lyapunov drift-based method.
  - Standard techniques fail in our case.
    - since the state-space collapse is neither a line nor a cone.
    - instead, it is even non-convex.

Our methods:

we use drift-based analysis to establish a new type of state-space collapse.
# The challenge to prove it...

(a) Diffusion approximations method.

- Witti-dimensional Brownian motion is difficult to analyze.
- optimality often only exists in finite time, not steady-state result.
- (b) Lyapunov drift-based method.
  - Standard techniques fail in our case.
    - since the state-space collapse is neither a line nor a cone.
    - instead, it is even non-convex.

Our methods:

- we use drift-based analysis to establish a new type of state-space collapse.
- we combine this state-space collapse with Chernoff bound.

## What is *state-space collapse?*

Informally, it means steady state 'concentrates' around a subspace.

- the distance between the steady state and a subspace has bounded moment upper bounds.
- the subspace could be a line or a cone in previous works.



#### Why is *state-space collapse* important? Recall the 'King' equation...

$$\lim_{\epsilon \downarrow 0} \mathbb{E}\left[\left\|\overline{\mathbf{Q}}^{(\epsilon)}(t+1)\right\|_1 \left\|\overline{\mathbf{U}}^{(\epsilon)}(t)\right\|_1\right] = 0.$$

where the unused service vector  $\mathbf{U}(t) = \max{\{\mathbf{S}(t) - \mathbf{Q}(t) - \mathbf{A}(t), \mathbf{0}\}}$ .

Note that 
$$Q_n(t+1)U_n(t) = 0$$
 for all *n* and *t*.

IMPLICATIONS: No server is idle while others with high loads.



In heavy traffic (  $\epsilon \rightarrow$  0), steady state lies far away from boundary.

What would be the *state-space collapse* in our case?



Drift towards the pink region due to preference of shorter queues.



Then, ALL the moments of the distance to pink region are bounded!

$$\mathbb{E}\left[e^{\theta^*d_{\mathcal{R}^{(r)}}\left(\overline{\mathbf{Q}}^{(\epsilon)}\right)}\right] \leq C^*,$$

where both  $\theta^*$  and  $C^*$  are independent of  $\epsilon$ .



$$\mathbb{E}\left[\overline{Q}_{2}^{(\epsilon)}(t+1)\overline{U}_{1}^{(\epsilon)}\right]$$

$$= \mathbb{E}\left[\overline{Q}_{2}(t+1)\overline{U}_{1}\mathcal{I}\left(\overline{Q}_{2}(t+1) \leq 2r, \overline{Q}_{1}(t+1) = 0\right)\right]$$

$$+ \mathbb{E}\left[\overline{Q}_{2}(t+1)\overline{U}_{1}\mathcal{I}\left(\overline{Q}_{2}(t+1) > 2r, \overline{Q}_{1}(t+1) = 0\right)\right]$$

$$(1)$$



$$\mathbb{E}\left[\overline{Q}_{2}^{(\epsilon)}(t+1)\overline{U}_{1}^{(\epsilon)}\right]$$
  
= $\mathbb{E}\left[\overline{Q}_{2}(t+1)\overline{U}_{1}\mathcal{I}\left(\overline{Q}_{2}(t+1) \leq 2r, \overline{Q}_{1}(t+1) = 0\right)\right]$  (1)  
+ $\mathbb{E}\left[\overline{Q}_{2}(t+1)\overline{U}_{1}\mathcal{I}\left(\overline{Q}_{2}(t+1) > 2r, \overline{Q}_{1}(t+1) = 0\right)\right]$  (2)

$$\begin{array}{l} (1) \leq 2r\epsilon, \text{ since } \mathbb{E}\left[\overline{U}_{1}\right] \leq \epsilon. \\ (2) \leq C\frac{1}{\epsilon^{2}}\mathbb{P}\left(\overline{Q}_{2}(t+1) > 2r, \overline{Q}_{1}(t+1) = 0\right) \leq C\frac{1}{\epsilon^{2}}\frac{1}{e^{\theta r}} \end{array}$$

**The key requirement:** set of servers whose queue lengths are below the threshold are known at the dispatcher

**The key requirement:** set of servers whose queue lengths are below the threshold are known at the dispatcher

Definition

**The key requirement:** set of servers whose queue lengths are below the threshold are known at the dispatcher

Definition

JBT(r) is composed of the following components:

Each server is initialized with an empty queue, and a corresponding ID in the local memory of the dispatcher.

**The key requirement:** set of servers whose queue lengths are below the threshold are known at the dispatcher

#### Definition

- Each server is initialized with an empty queue, and a corresponding ID in the local memory of the dispatcher.
- Upon new arrivals at the beginning of each time-slot, the dispatcher checks the available IDs in memory.

**The key requirement:** set of servers whose queue lengths are below the threshold are known at the dispatcher

#### Definition

- Each server is initialized with an empty queue, and a corresponding ID in the local memory of the dispatcher.
- Upon new arrivals at the beginning of each time-slot, the dispatcher checks the available IDs in memory.
  - If one or more IDs exist, it removes one uniformly at random, and sends all the new arrivals to the corresponding server.

**The key requirement:** set of servers whose queue lengths are below the threshold are known at the dispatcher

#### Definition

- Each server is initialized with an empty queue, and a corresponding ID in the local memory of the dispatcher.
- Upon new arrivals at the beginning of each time-slot, the dispatcher checks the available IDs in memory.
  - If one or more IDs exist, it removes one uniformly at random, and sends all the new arrivals to the corresponding server.
  - Otherwise, all the new arrivals are dispatched uniformly at random to one of the servers in the system.

**The key requirement:** set of servers whose queue lengths are below the threshold are known at the dispatcher

#### Definition

- Each server is initialized with an empty queue, and a corresponding ID in the local memory of the dispatcher.
- Upon new arrivals at the beginning of each time-slot, the dispatcher checks the available IDs in memory.
  - If one or more IDs exist, it removes one uniformly at random, and sends all the new arrivals to the corresponding server.
  - Otherwise, all the new arrivals are dispatched uniformly at random to one of the servers in the system.
- Each server reports its ID to the dispatcher at the end of each time-slot if
  - its queue length is below the threshold
  - AND the dispatcher does not contain its ID (how?)

If we can estimate the traffic load, then we can directly apply the sufficient condition.

- If we can estimate the traffic load, then we can directly apply the sufficient condition.
- If not, we can use sampling every T time-slots.
  - In particular, randomly sample d queues and take the minimum as threshold.
  - We can prove the following result.

- If we can estimate the traffic load, then we can directly apply the sufficient condition.
- If not, we can use sampling every T time-slots.
  - ► In particular, randomly sample *d* queues and take the minimum as threshold.
  - We can prove the following result.

#### Theorem

For any finite T and  $d \ge 1$ , the policy is throughput and delay optimal in heavy traffic.

Then...is the Logarithmic growth rate also necessary?

Then...is the Logarithmic growth rate also necessary?

We conjecture so!

# Conclusion...

#### Theorem (Necessary Conditions)

- ▶ The threshold *r* should grow with the traffic load.
- But, it can not grow too fast.
- It provides a sharp characterization of JIQ policy.

#### Theorem (Sufficient Conditions)

- It is sufficient to have a logarithmic growth rate.
- This resolves a long-standing open problem.
- It provides a useful guideline for practical systems.

Thank you! Q & A