

Scalable Load Balancing in Networked Systems

Sem Borst

Eindhoven University of Technology (TU/e) & Nokia Bell Labs

UFRJ, 2 August 2018

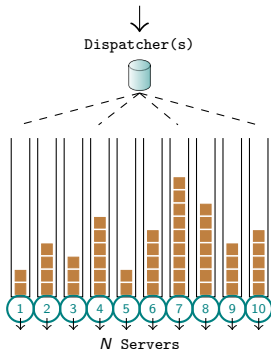
Based on joint work with Mark van der Boor, Johan van Leeuwen, Debankur Mukherjee & Phil Whiting



TWO FUNDAMENTAL PROBLEMS IN NETWORKED SYSTEMS

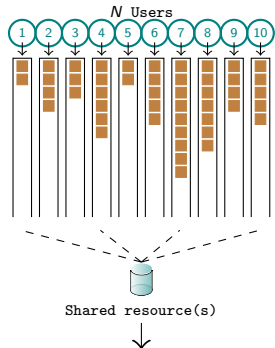
Centralized queue impractical/undesirable in large-scale networked systems

Achieve **low delay** and maintain **low overhead** as system grows large
(with highly distributed queues)



Load balancing/routing

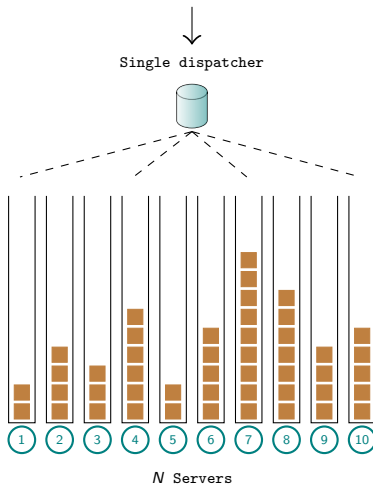
Assign tasks from single arrival stream to be performed by one of several servers, each with their own queue



Resource allocation/scheduling

Allocate shared resource to execute tasks from one of several users, each with their own arrival stream and queue

LOAD BALANCING/ROUTING IN PARALLEL-SERVER SYSTEMS

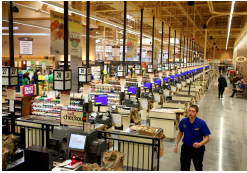


LARGE-SCALE PARALLEL-SERVER SYSTEMS: SOME EXAMPLES



Supermarket checkout line

LARGE-SCALE PARALLEL-SERVER SYSTEMS: SOME EXAMPLES

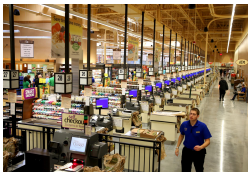


Supermarket checkout line

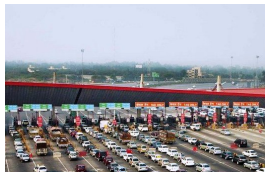


Road toll plaza

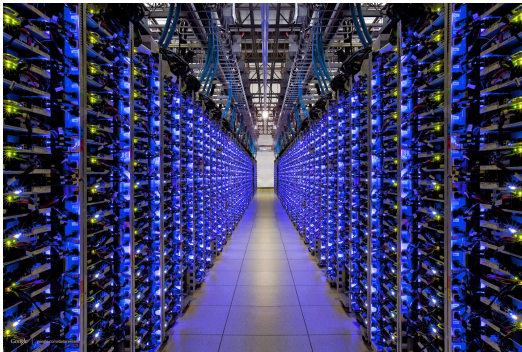
LARGE-SCALE PARALLEL-SERVER SYSTEMS: SOME EXAMPLES



Supermarket checkout line



Road toll plaza



Data center

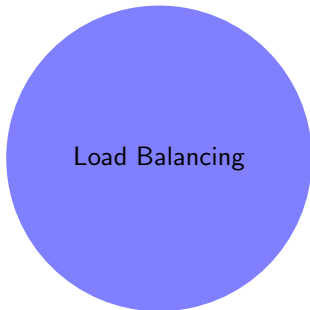
HIGH-LEVEL OUTLINE

- I. Scalability challenges and classical results
- II. Asymptotic optimality and universality (no memory)
- III. Reduction in communication overhead (memory)
- IV. Heterogeneity issues and network scenarios

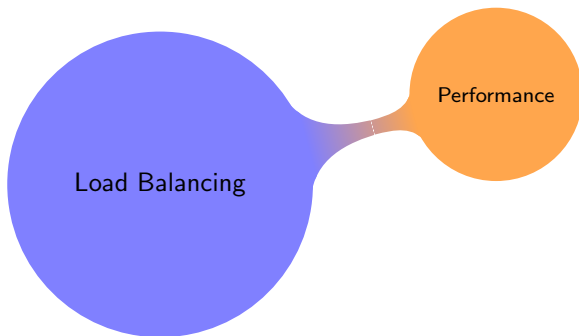
I. Scalability challenges and classical results

- II. Asymptotic optimality and universality (no memory)
- III. Reduction in communication overhead (memory)
- IV. Heterogeneity issues and network scenarios

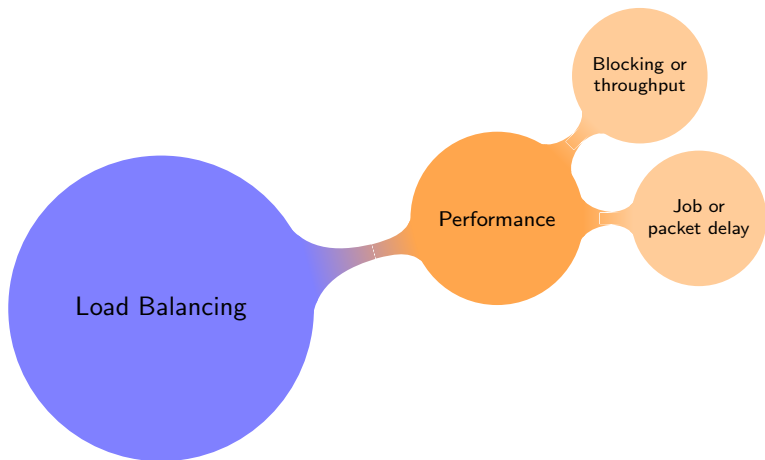
LOAD BALANCING: SCALABILITY ISSUES



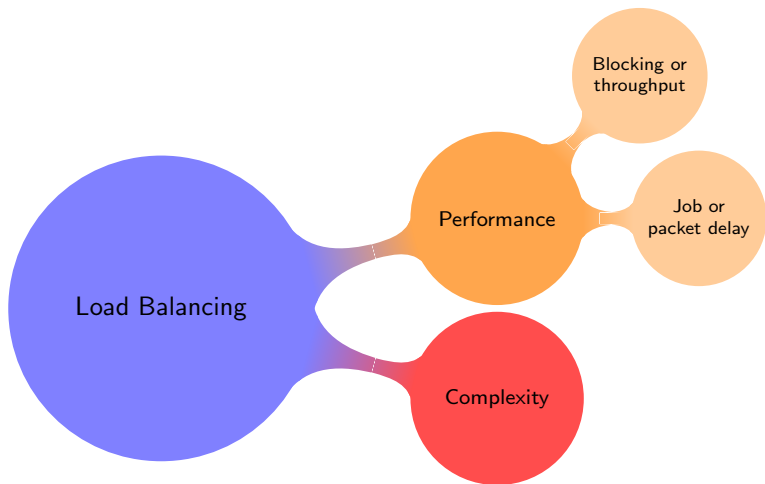
LOAD BALANCING: SCALABILITY ISSUES



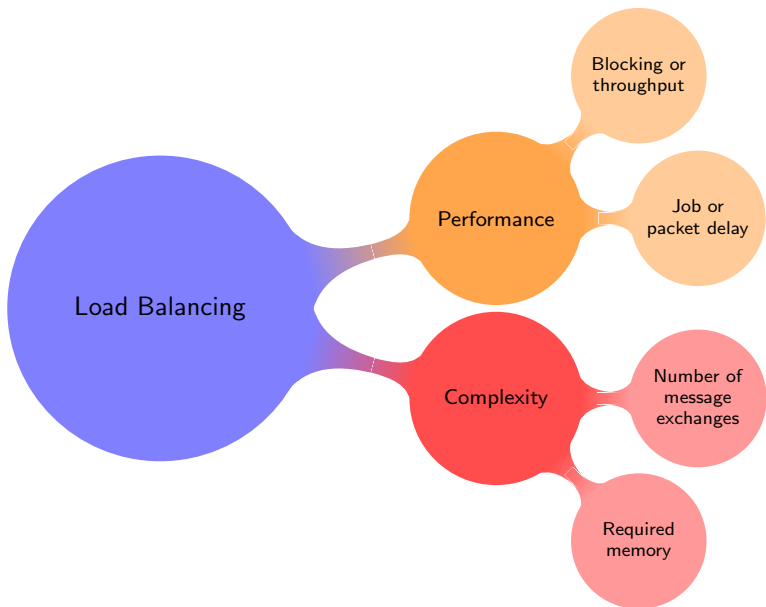
LOAD BALANCING: SCALABILITY ISSUES



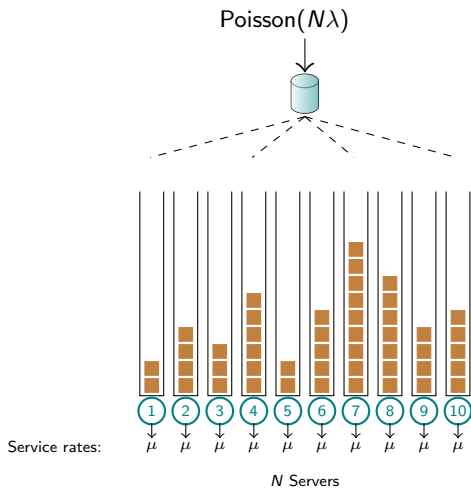
LOAD BALANCING: SCALABILITY ISSUES



LOAD BALANCING: SCALABILITY ISSUES

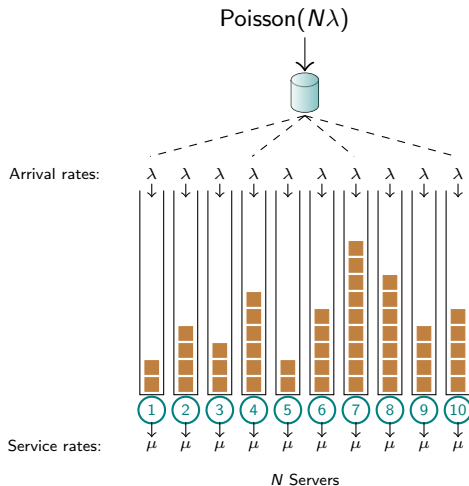


PURELY RANDOM ASSIGNMENT



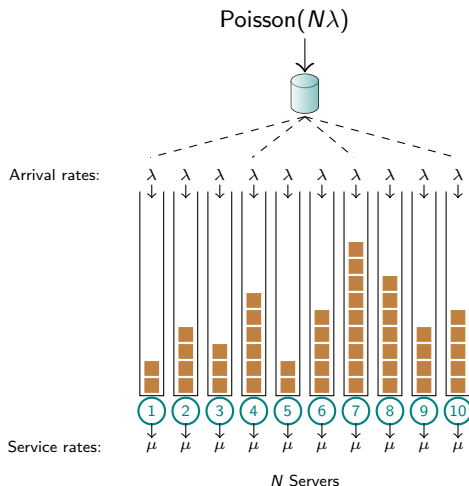
PURELY RANDOM ASSIGNMENT

Assign each task to server selected uniformly at random



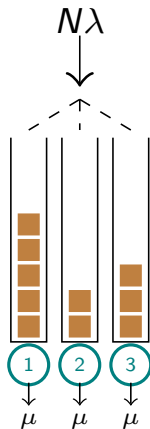
PURELY RANDOM ASSIGNMENT

Assign each task to server selected uniformly at random



N independent M/M/1 queues with arrival rate λ and service rate μ

LOAD BALANCING SCENARIOS

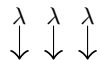


LOAD BALANCING SCENARIOS

Separate Queues

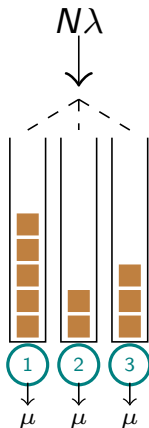
Strictly Random Routing

$N \times M/M/1$



queue length $N \times \frac{\lambda^2}{\mu(\mu-\lambda)}$

waiting time $\frac{\lambda}{\mu(\mu-\lambda)}$

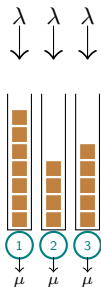


LOAD BALANCING SCENARIOS

Separate Queues

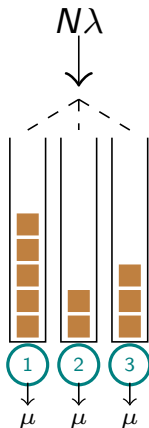
Strictly Random Routing

$N \times M/M/1$



queue length $N \times \frac{\lambda^2}{\mu(\mu-\lambda)}$

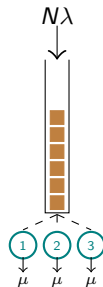
waiting time $\frac{\lambda}{\mu(\mu-\lambda)}$



Centralized Queue

Complete Resource Pooling

$M/M/N$

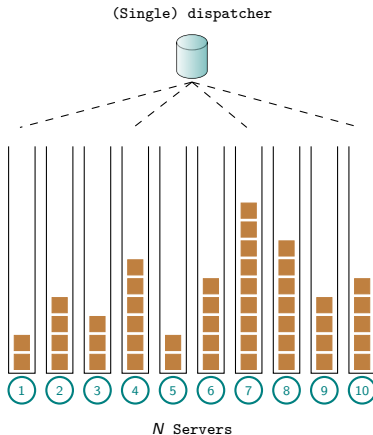


queue length $\Pi_W \frac{\lambda}{\mu-\lambda}$

waiting time $\frac{1}{N} \Pi_W \frac{1}{\mu-\lambda}$

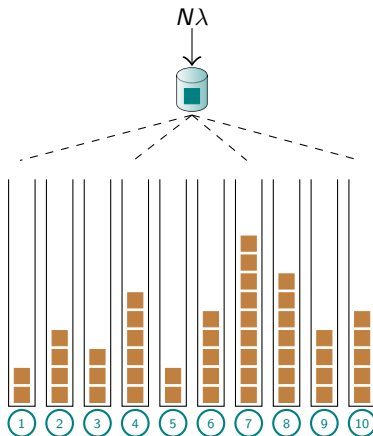
JOIN THE SHORTEST QUEUE (JSQ) POLICY

Classical natural policy: Join the Shortest Queue (JSQ)



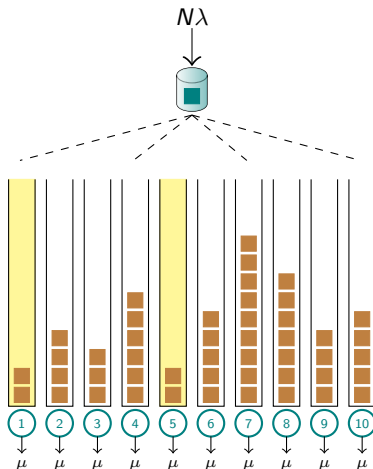
JOIN THE SHORTEST QUEUE (JSQ) POLICY

Classical natural policy: Join the Shortest Queue (JSQ)



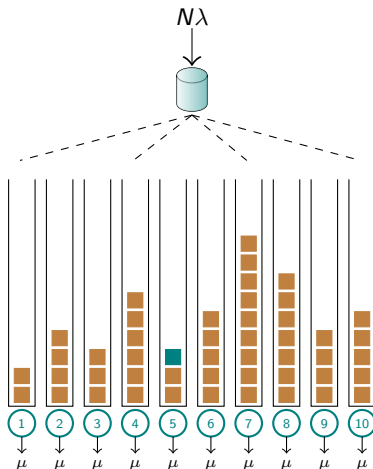
JOIN THE SHORTEST QUEUE (JSQ) POLICY

Classical natural policy: Join the Shortest Queue (JSQ)



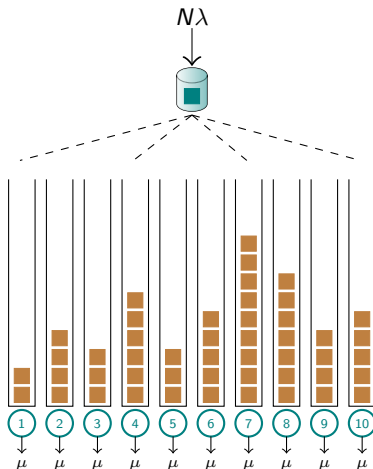
JOIN THE SHORTEST QUEUE (JSQ) POLICY

Classical natural policy: Join the Shortest Queue (JSQ)



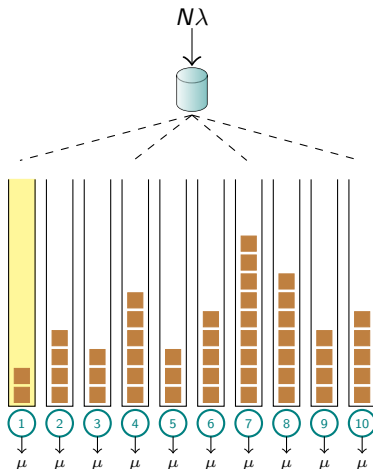
JOIN THE SHORTEST QUEUE (JSQ) POLICY

Classical natural policy: Join the Shortest Queue (JSQ)



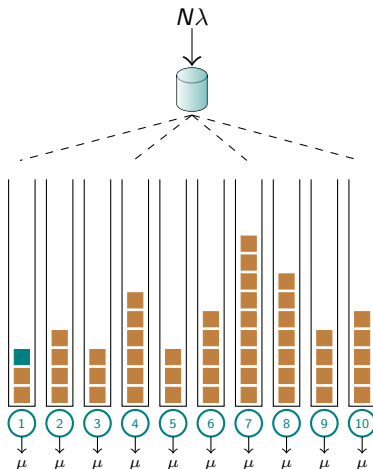
JOIN THE SHORTEST QUEUE (JSQ) POLICY

Classical natural policy: Join the Shortest Queue (JSQ)



JOIN THE SHORTEST QUEUE (JSQ) POLICY

Classical natural policy: Join the Shortest Queue (JSQ)



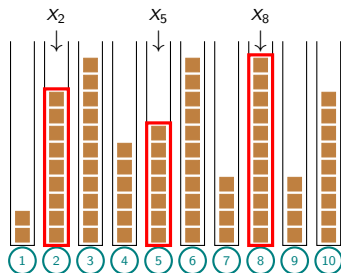
EQUIVALENT STATE DESCRIPTION

Symmetry among servers allows equivalent state description

EQUIVALENT STATE DESCRIPTION

Symmetry among servers allows equivalent state description

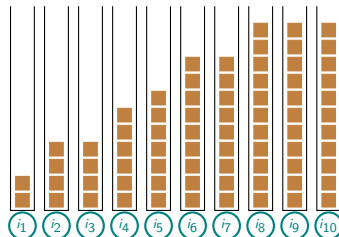
- Denote by X_k queue length at server k , $k = 1, 2, \dots, N$



EQUIVALENT STATE DESCRIPTION

Symmetry among servers allows equivalent state description

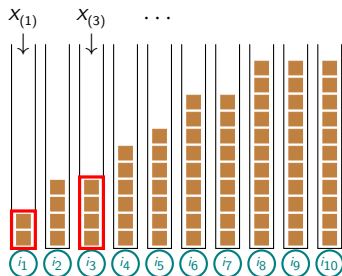
- ▶ Denote by X_k queue length at server k , $k = 1, 2, \dots, N$
- ▶ Rearrange servers from left to right in non-decreasing order of X_k 's



EQUIVALENT STATE DESCRIPTION

Symmetry among servers allows equivalent state description

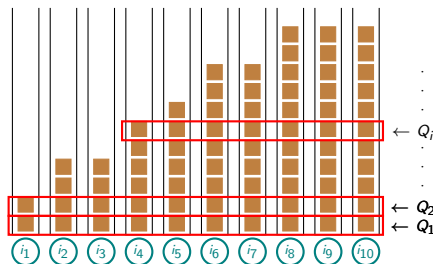
- ▶ Denote by X_k queue length at server k , $k = 1, 2, \dots, N$
- ▶ Rearrange servers from left to right in non-decreasing order of X_k 's
- ▶ $X_{(k)}$ is size of k^{th} shortest queue, i.e., height of k^{th} left-most stack



EQUIVALENT STATE DESCRIPTION

Symmetry among servers allows equivalent state description

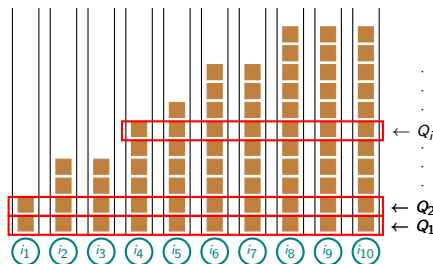
- ▶ Denote by X_k queue length at server k , $k = 1, 2, \dots, N$
- ▶ Rearrange servers from left to right in non-decreasing order of X_k 's
- ▶ $X_{(k)}$ is size of k^{th} shortest queue, i.e., height of k^{th} left-most stack
- ▶ Denote by Q_i number of servers with queue length i or larger, $i = 1, 2, \dots$, i.e., width of i^{th} horizontal bar



EQUIVALENT STATE DESCRIPTION

Symmetry among servers allows equivalent state description

- ▶ Denote by X_k queue length at server k , $k = 1, 2, \dots, N$
- ▶ Rearrange servers from left to right in non-decreasing order of X_k 's
- ▶ $X_{(k)}$ is size of k^{th} shortest queue, i.e., height of k^{th} left-most stack
- ▶ Denote by Q_i number of servers with queue length i or larger, $i = 1, 2, \dots$, i.e., width of i^{th} horizontal bar
- ▶ Fluid-scaled state variables $q_i^N(t) = Q_i^N(t)/N$ represent fraction of servers with queue length i or larger



STOCHASTIC OPTIMALITY OF JSQ POLICY

- ▶ JSQ stochastically minimizes aggregate size of l right-most stacks, i.e., total number of tasks in l largest queues

$$\sum_{k=N-l+1}^N X_{(k)}^{\text{JSQ}} \leq_{st} \sum_{k=N-l+1}^N X_{(k)}^{\Pi}$$

for all $l = 1, \dots, N$, for any non-anticipating policy Π

[Towsley-Sparaggis-Cassandras 1992; Sparaggis-Towsley-Cassandra 1994]

STOCHASTIC OPTIMALITY OF JSQ POLICY

- ▶ JSQ stochastically minimizes aggregate size of l right-most stacks, i.e., total number of tasks in l largest queues

$$\sum_{k=N-l+1}^N X_{(k)}^{\text{JSQ}} \leq_{st} \sum_{k=N-l+1}^N X_{(k)}^{\Pi}$$

for all $l = 1, \dots, N$, for any non-anticipating policy Π

[Towsley-Sparaggis-Cassandras 1992; Sparaggis-Towsley-Cassandra 1994]

- ▶ JSQ stochastically minimizes aggregate size of bars at level j or higher, i.e., total number of tasks in queue position j or higher

$$\sum_{i=j}^{\infty} Q_i^{\text{JSQ}} \leq_{st} \sum_{i=j}^{\infty} Q_i^{\Pi}$$

for all $j = 1, 2, \dots$, for any non-anticipating policy Π

[Mukherjee, B, Van Leeuwen, Whiting 2016]

STOCHASTIC OPTIMALITY OF JSQ POLICY

- ▶ JSQ stochastically minimizes aggregate size of l right-most stacks, i.e., total number of tasks in l largest queues

$$\sum_{k=N-l+1}^N X_{(k)}^{\text{JSQ}} \leq_{st} \sum_{k=N-l+1}^N X_{(k)}^{\Pi}$$

for all $l = 1, \dots, N$, for any non-anticipating policy Π

[Towsley-Sparagis-Cassandras 1992; Sparagis-Towsley-Cassandra 1994]

- ▶ JSQ stochastically minimizes aggregate size of bars at level j or higher, i.e., total number of tasks in queue position j or higher

$$\sum_{i=j}^{\infty} Q_i^{\text{JSQ}} \leq_{st} \sum_{i=j}^{\infty} Q_i^{\Pi}$$

for all $j = 1, 2, \dots$, for any non-anticipating policy Π

[Mukherjee, B, Van Leeuwen, Whiting 2016]

In particular (take $l = N$ or $j = 1$), JSQ stochastically minimizes total number of tasks in system, and hence overall mean delay

JSQ POLICY IN MANY-SERVER REGIME

JSQ yields dramatic performance improvements as $N \rightarrow \infty$

JSQ POLICY IN MANY-SERVER REGIME

JSQ yields dramatic performance improvements as $N \rightarrow \infty$

- Eliminates queues

Scheme	Queue length	
Random	$p_i^N = \lambda^i$	
JSQ	$p_1^N \rightarrow \lambda,$ $p_2^N \rightarrow 0$	

p_i^N is stationary probability that queue length at server is $\geq i$ in N^{th} system

JSQ POLICY IN MANY-SERVER REGIME

JSQ yields dramatic performance improvements as $N \rightarrow \infty$

- ▶ Eliminates queues
- ▶ Achieves zero wait

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	
JSQ	$p_1^N \rightarrow \lambda,$ $p_2^N \rightarrow 0$	$o(1)$	$\Theta(1/\sqrt{N})$	

p_i^N is stationary probability that queue length at server is $\geq i$ in N^{th} system

JSQ POLICY IN MANY-SERVER REGIME

JSQ yields dramatic performance improvements as $N \rightarrow \infty$

- ▶ Eliminates queues
- ▶ Achieves zero wait

But...

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	
JSQ	$p_1^N \rightarrow \lambda,$ $p_2^N \rightarrow 0$	$o(1)$	$\Theta(1/\sqrt{N})$	

p_i^N is stationary probability that queue length at server is $\geq i$ in N^{th} system

JSQ POLICY IN MANY-SERVER REGIME

JSQ yields dramatic performance improvements as $N \rightarrow \infty$

- ▶ Eliminates queues
- ▶ Achieves zero wait

But...

- ▶ JSQ involves high communication overhead as $N \rightarrow \infty$

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	Overhead
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	0
JSQ	$p_1^N \rightarrow \lambda,$ $p_2^N \rightarrow 0$	$o(1)$	$\Theta(1/\sqrt{N})$	$2N$

p_i^N is stationary probability that queue length at server is $\geq i$ in N^{th} system

JSQ POLICY IN MANY-SERVER REGIME

JSQ yields dramatic performance improvements as $N \rightarrow \infty$

- ▶ Eliminates queues
- ▶ Achieves zero wait

But...

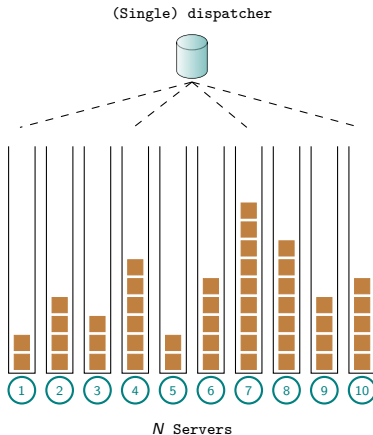
- ▶ JSQ involves high communication overhead as $N \rightarrow \infty$
- ▶ Straightforward implementation of JSQ (no memory at dispatcher) requires queue lengths at all servers to be checked at each arrival, which may be prohibitive in large-scale systems

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	Overhead
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	0
JSQ	$p_1^N \rightarrow \lambda,$ $p_2^N \rightarrow 0$	$o(1)$	$\Theta(1/\sqrt{N})$	$2N$

p_i^N is stationary probability that queue length at server is $\geq i$ in N^{th} system

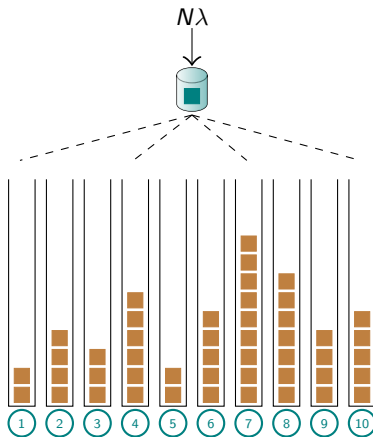
JSQ(d) SCHEME: REDUCED COMMUNICATION OVERHEAD

At each arrival, select d servers (e.g. $d = 2$) uniformly at random, and assign task to shortest queue among selected servers



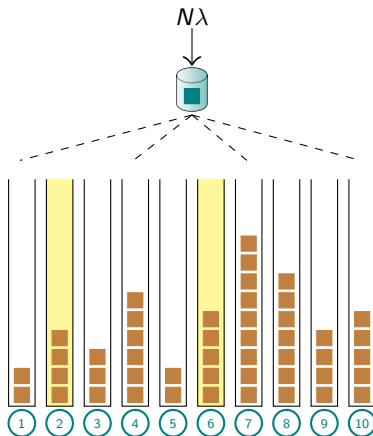
JSQ(d) SCHEME: REDUCED COMMUNICATION OVERHEAD

At each arrival, select d servers (e.g. $d = 2$) uniformly at random, and assign task to shortest queue among selected servers



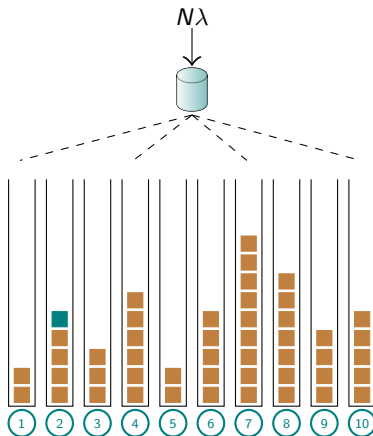
JSQ(d) SCHEME: REDUCED COMMUNICATION OVERHEAD

At each arrival, select d servers (e.g. $d = 2$) uniformly at random, and assign task to shortest queue among selected servers



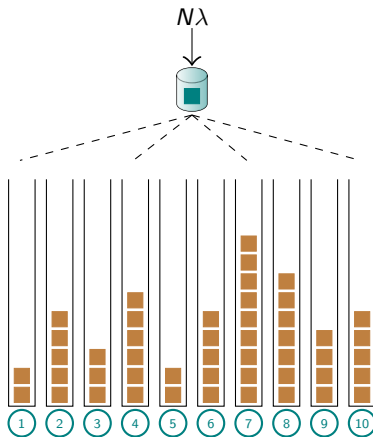
JSQ(d) SCHEME: REDUCED COMMUNICATION OVERHEAD

At each arrival, select d servers (e.g. $d = 2$) uniformly at random, and assign task to shortest queue among selected servers



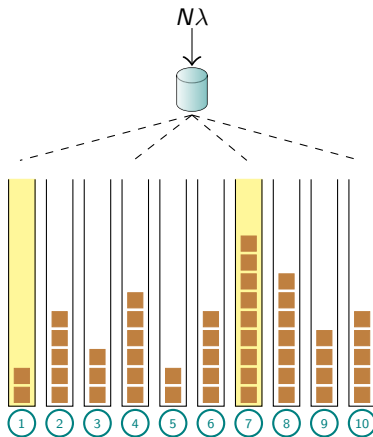
JSQ(d) SCHEME: REDUCED COMMUNICATION OVERHEAD

At each arrival, select d servers (e.g. $d = 2$) uniformly at random, and assign task to shortest queue among selected servers



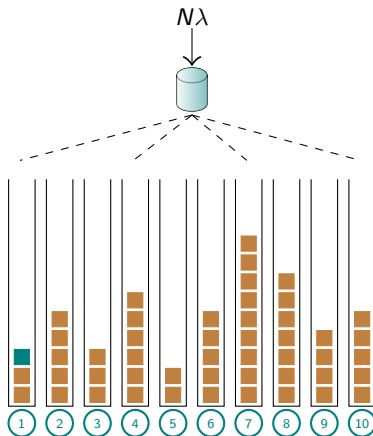
JSQ(d) SCHEME: REDUCED COMMUNICATION OVERHEAD

At each arrival, select d servers (e.g. $d = 2$) uniformly at random, and assign task to shortest queue among selected servers



JSQ(d) SCHEME: REDUCED COMMUNICATION OVERHEAD

At each arrival, select d servers (e.g. $d = 2$) uniformly at random, and assign task to shortest queue among selected servers



YET LOWER COMMUNICATION OVERHEAD PER TASK FOR BATCH ARRIVALS

For batch arrivals d queue samples can be amortized over several tasks

L. Ying, R. Srikant, X. Kang (2015). The power of slightly more than one sample in randomized load balancing. In: *Proc. IEEE Infocom 2015*.

CONNECTION BETWEEN JSQ(d) AND REDUNDANCY- d POLICIES

Redundancy- d policy assigns replicas of each arriving task to d servers sampled uniformly at random, and kills $d - 1$ replicas as soon as execution of first replica is started/completed: Wait before see rather than see before wait

K.S. Gardner, M. Harchol-Balter, A. Scheller-Wolf, M. Velednitsky, S. Zbarsky (2017). Redundancy- d : The power of d choices for redundancy. *Oper. Res.* **65** (4), 1078–1094.

K.S. Gardner, S. Zbarsky, S. Doroudi, M. Harchol-Balter, E. Hyytiä, A. Scheller-Wolf (2015). Reducing latency via redundant requests: Exact analysis. In: *ACM SIGMETRICS Perf. Eval. Rev.* **43** (1), 347–360.

K.S. Gardner, S. Zbarsky, S. Doroudi, M. Harchol-Balter, E. Hyytiä, A. Scheller-Wolf (2016). Queueing with redundant requests: Exact analysis. *Queueing Systems* **83** (3–4), 227–259.

K.S. Gardner, S. Zbarsky, M. Harchol-Balter, A. Scheller-Wolf (2016). The power of d choices for redundancy. In: *ACM SIGMETRICS Perf. Eval. Rev.* **44** (1), 409–410.

K.S. Gardner, S. Zbarsky, M. Velednitsky, M. Harchol-Balter, A. Scheller-Wolf (2016). Understanding response time in the redundancy- d system. In: *ACM SIGMETRICS Perf. Eval. Rev.* **44** (2), 33–35.

FLUID LIMIT FOR JSQ(d) SCHEME

Fluid limit for JSQ(d) [Mitzenmacher 1996, 2001; Vvedenskaya *et al.* 1996]

If $\mathbf{q}^{JSQ(d)}(0) \rightarrow \mathbf{q}^\infty$ as $N \rightarrow \infty$, then $\{\mathbf{q}^{JSQ(d)}(t)\}_{t \geq 0}$ weakly converges to $\{\mathbf{q}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$, with

$$\frac{dq_i(t)}{dt} = \lambda[(q_{i-1}(t))^d - (q_i(t))^d] - \mu[q_i(t) - q_{i+1}(t)],$$

where $\mathbf{q}(0) = \mathbf{q}^\infty$

FLUID LIMIT FOR JSQ(d) SCHEME

Fluid limit for JSQ(d) [Mitzenmacher 1996, 2001; Vvedenskaya *et al.* 1996]

If $\mathbf{q}^{JSQ(d)}(0) \rightarrow \mathbf{q}^\infty$ as $N \rightarrow \infty$, then $\{\mathbf{q}^{JSQ(d)}(t)\}_{t \geq 0}$ weakly converges to $\{\mathbf{q}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$, with

$$\frac{dq_i(t)}{dt} = \lambda[(q_{i-1}(t))^d - (q_i(t))^d] - \mu[q_i(t) - q_{i+1}(t)],$$

where $\mathbf{q}(0) = \mathbf{q}^\infty$

$[(q_{i-1}(t))^d - (q_i(t))^d]$ is (instantaneous) fraction of arriving tasks assigned to servers with queue length $i - 1$ in fluid-level state $\mathbf{q}(t)$

POWER-OF- d EFFECT FOR JSQ(d) SCHEME

Assuming $\mu = 1$ as before, fixed point of fluid limit is

$$q_i^* := \lim_{t \rightarrow \infty} q_i(t) = \lambda^{\frac{d^i - 1}{d - 1}}, \quad i = 1, 2, \dots$$

POWER-OF- d EFFECT FOR JSQ(d) SCHEME

Assuming $\mu = 1$ as before, fixed point of fluid limit is

$$q_i^* := \lim_{t \rightarrow \infty} q_i(t) = \lambda^{\frac{d^i - 1}{d - 1}}, \quad i = 1, 2, \dots$$

Denote $p_i^* := \lim_{N \rightarrow \infty} p_i^N$,

with p_i^N stationary probability that queue length at server is $\geq i$ in N^{th} system

POWER-OF- d EFFECT FOR JSQ(d) SCHEME

Assuming $\mu = 1$ as before, fixed point of fluid limit is

$$q_i^* := \lim_{t \rightarrow \infty} q_i(t) = \lambda^{\frac{d^i - 1}{d - 1}}, \quad i = 1, 2, \dots$$

Denote $p_i^* := \lim_{N \rightarrow \infty} p_i^N$,

with p_i^N stationary probability that queue length at server is $\geq i$ in N^{th} system

Interchange of large- N and large- t limits: $p_i^* = q_i^*$

POWER-OF- d EFFECT FOR JSQ(d) SCHEME

Assuming $\mu = 1$ as before, fixed point of fluid limit is

$$q_i^* := \lim_{t \rightarrow \infty} q_i(t) = \lambda^{\frac{d^i - 1}{d - 1}}, \quad i = 1, 2, \dots$$

Denote $p_i^* := \lim_{N \rightarrow \infty} p_i^N$,

with p_i^N stationary probability that queue length at server is $\geq i$ in N^{th} system

Interchange of large- N and large- t limits: $p_i^* = q_i^*$

Tail of stationary queue length distribution falls off much faster for $d \geq 2$ than for purely random assignment ($d = 1$)

POWER-OF- d EFFECT FOR JSQ(d) SCHEME

Assuming $\mu = 1$ as before, fixed point of fluid limit is

$$q_i^* := \lim_{t \rightarrow \infty} q_i(t) = \lambda^{\frac{d^i - 1}{d - 1}}, \quad i = 1, 2, \dots$$

Denote $p_i^* := \lim_{N \rightarrow \infty} p_i^N$,

with p_i^N stationary probability that queue length at server is $\geq i$ in N^{th} system

Interchange of large- N and large- t limits: $p_i^* = q_i^*$

Tail of stationary queue length distribution falls off much faster for $d \geq 2$ than for purely random assignment ($d = 1$)

“power-of-two” effect: Even value as small as $d = 2$ yields significant performance improvements over purely random assignment, while drastically reducing communication overhead compared to JSQ ($d = N$)

JSQ(d) PROVIDES STRONG BENEFITS OVER RANDOM ASSIGNMENT

JSQ(d) provides doubly-exponential rather than (singly-)exponential decay of stationary queue length tail

Scheme	Queue length		Overhead
Random	$p_i^N = \lambda^i$		0
JSQ(d)	$p_i^* = \lambda^{\frac{d^i-1}{d-1}}$		$2d$
JSQ	$p_1^* = \lambda,$ $p_2^* = 0$		$2N$

JSQ(d) PROVIDES STRONG BENEFITS OVER RANDOM ASSIGNMENT

JSQ(d) provides doubly-exponential rather than (singly-)exponential decay of stationary queue length tail

But...

Scheme	Queue length		Overhead
Random	$p_i^N = \lambda^i$		0
JSQ(d)	$p_i^* = \lambda^{\frac{d^i-1}{d-1}}$		$2d$
JSQ	$p_1^* = \lambda,$ $p_2^* = 0$		$2N$

JSQ(d) PROVIDES STRONG BENEFITS OVER RANDOM ASSIGNMENT

JSQ(d) provides doubly-exponential rather than (singly-)exponential decay of stationary queue length tail

But...

- ▶ JSQ(d) scheme is **suboptimal** for any fixed value of d

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	Overhead
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	0
JSQ(d)	$p_i^* = \lambda^{\frac{d^i-1}{d-1}}$	$\Theta(1)$	$\Omega(\log N)$	$2d$
JSQ	$p_1^* = \lambda,$ $p_2^* = 0$	$o(1)$	$\Theta(1/\sqrt{N})$	$2N$

JSQ(d) PROVIDES STRONG BENEFITS OVER RANDOM ASSIGNMENT

JSQ(d) provides doubly-exponential rather than (singly-)exponential decay of stationary queue length tail

But...

- ▶ JSQ(d) scheme is **suboptimal** for any fixed value of d
- ▶ Waiting time **does not** vanish as $N \rightarrow \infty$

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	Overhead
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	0
JSQ(d)	$p_i^* = \lambda^{\frac{d^i-1}{d-1}}$	$\Theta(1)$	$\Omega(\log N)$	$2d$
JSQ	$p_1^* = \lambda,$ $p_2^* = 0$	$o(1)$	$\Theta(1/\sqrt{N})$	$2N$

JSQ(d) PROVIDES STRONG BENEFITS OVER RANDOM ASSIGNMENT

JSQ(d) provides doubly-exponential rather than (singly-)exponential decay of stationary queue length tail

But...

- ▶ JSQ(d) scheme is **suboptimal** for any fixed value of d
- ▶ Waiting time **does not** vanish as $N \rightarrow \infty$

In absence of any memory at dispatcher, communication overhead **must grow** with N in order for zero delay to be achievable [Gamarnik-Tsitsiklis-Zubeldia 2016]

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	Overhead
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	0
JSQ(d)	$p_i^* = \lambda^{\frac{d^i-1}{d-1}}$	$\Theta(1)$	$\Omega(\log N)$	$2d$
JSQ	$p_1^* = \lambda,$ $p_2^* = 0$	$o(1)$	$\Theta(1/\sqrt{N})$	$2N$

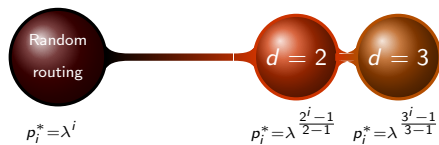
- I. Scalability challenges and classical results
- II. Asymptotic optimality and universality (no memory)
- III. Reduction in communication overhead (memory)
- IV. Heterogeneity issues and network scenarios

PERFORMANCE VERSUS COMMUNICATION OVERHEAD AT A GLANCE

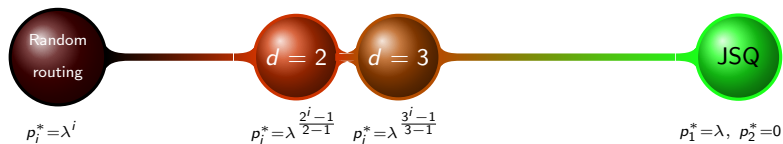


$$p_i^* = \lambda^i$$

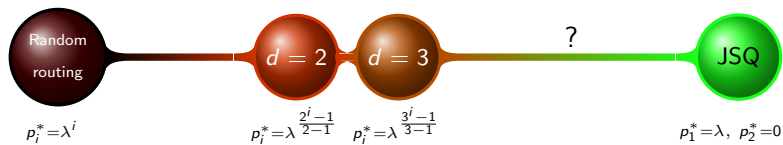
PERFORMANCE VERSUS COMMUNICATION OVERHEAD AT A GLANCE



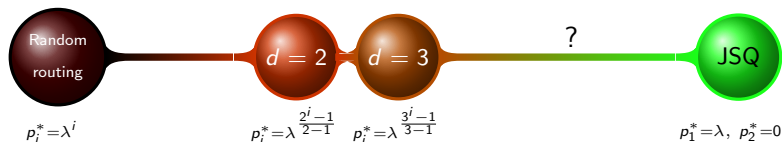
PERFORMANCE VERSUS COMMUNICATION OVERHEAD AT A GLANCE



PERFORMANCE VERSUS COMMUNICATION OVERHEAD AT A GLANCE

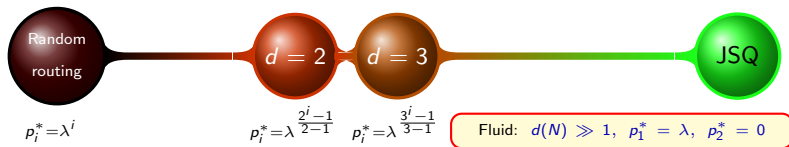


PERFORMANCE VERSUS COMMUNICATION OVERHEAD AT A GLANCE



In order to examine performance versus communication trade-off, we allow parameter d to depend on N and write $d(N)$ to reflect that

PERFORMANCE VERSUS COMMUNICATION OVERHEAD AT A GLANCE



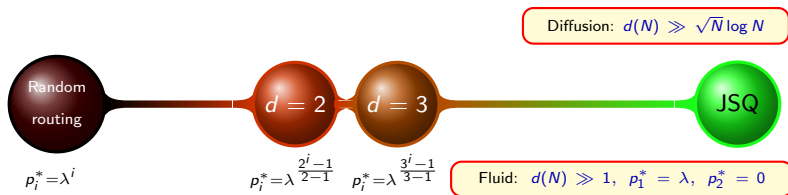
In order to examine performance versus communication trade-off, we allow parameter d to depend on N and write $d(N)$ to reflect that

Universality and Asymptotic Optimality Properties

- ▶ JSQ($d(N)$) has same **fluid limit** as JSQ, and achieves **fluid-level optimality** when

$$d(N) \rightarrow \infty$$

PERFORMANCE VERSUS COMMUNICATION OVERHEAD AT A GLANCE



In order to examine performance versus communication trade-off, we allow parameter d to depend on N and write $d(N)$ to reflect that

Universality and Asymptotic Optimality Properties

- ▶ JSQ($d(N)$) has same **fluid limit** as JSQ, and achieves **fluid-level optimality** when

$$d(N) \rightarrow \infty$$

- ▶ JSQ($d(N)$) has same **diffusion limit** as JSQ, and achieves **diffusion-level optimality** when

$$\frac{d(N)}{\sqrt{N} \log N} \rightarrow \infty$$

UNIVERSALITY OF FLUID LIMIT FOR JSQ($d(N)$) SCHEME

- ▶ Fluid-scaled state variables $q_i^{d(N)}(t) = Q_i^{d(N)}(t)/N$ represent fraction of servers with queue length i or larger in N^{th} system with JSQ($d(N)$) scheme

UNIVERSALITY OF FLUID LIMIT FOR JSQ($d(N)$) SCHEME

- Fluid-scaled state variables $q_i^{d(N)}(t) = Q_i^{d(N)}(t)/N$ represent fraction of servers with queue length i or larger in N^{th} system with JSQ($d(N)$) scheme

Universality of fluid limit for JSQ($d(N)$) [Mukherjee, B, Van L, W 2016]

If $\mathbf{q}^{d(N)}(0) \rightarrow \mathbf{q}^\infty$ and $d(N) \rightarrow \infty$ as $N \rightarrow \infty$, then $\{\mathbf{q}^{d(N)}(t)\}_{t \geq 0}$ has same weak limit $\{\mathbf{q}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$ as ordinary JSQ policy, with

$$\frac{dq_i(t)}{dt} = \lambda p_{i-1}(\mathbf{q}(t)) - \mu[q_i(t) - q_{i+1}(t)],$$

where $\mathbf{q}(0) = \mathbf{q}^\infty$

UNIVERSALITY OF FLUID LIMIT FOR JSQ($d(N)$) SCHEME

- Fluid-scaled state variables $q_i^{d(N)}(t) = Q_i^{d(N)}(t)/N$ represent fraction of servers with queue length i or larger in N^{th} system with JSQ($d(N)$) scheme

Universality of fluid limit for JSQ($d(N)$) [Mukherjee, B, Van L, W 2016]

If $\mathbf{q}^{d(N)}(0) \rightarrow \mathbf{q}^\infty$ and $d(N) \rightarrow \infty$ as $N \rightarrow \infty$, then $\{\mathbf{q}^{d(N)}(t)\}_{t \geq 0}$ has same weak limit $\{\mathbf{q}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$ as ordinary JSQ policy, with

$$\frac{dq_i(t)}{dt} = \lambda p_{i-1}(\mathbf{q}(t)) - \mu[q_i(t) - q_{i+1}(t)],$$

where $\mathbf{q}(0) = \mathbf{q}^\infty$

$p_{i-1}(\mathbf{q}(t))$ is (instantaneous) fraction of arriving tasks assigned to servers with queue length $i - 1$ in fluid-level state $\mathbf{q}(t)$

UNIVERSALITY OF FLUID LIMIT FOR JSQ($d(N)$) SCHEME

Observations

- ▶ Fluid-level behavior coincides with that of ordinary JSQ policy as long as $d(N) \rightarrow \infty$ as $N \rightarrow \infty$

UNIVERSALITY OF FLUID LIMIT FOR JSQ($d(N)$) SCHEME

Observations

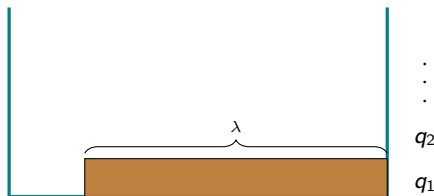
- ▶ Fluid-level behavior coincides with that of ordinary JSQ policy as long as $d(N) \rightarrow \infty$ as $N \rightarrow \infty$
- ▶ JSQ($d(N)$) scheme achieves fluid-level optimality while reducing communication overhead by nearly a factor $O(N)$

UNIVERSALITY OF FLUID LIMIT FOR JSQ($d(N)$) SCHEME

Observations

- ▶ Fluid-level behavior coincides with that of ordinary JSQ policy as long as $d(N) \rightarrow \infty$ as $N \rightarrow \infty$
- ▶ JSQ($d(N)$) scheme achieves fluid-level optimality while reducing communication overhead by nearly a factor $O(N)$
- ▶ Fixed point of fluid limit

$$q_1^* = \lambda, \quad q_i^* = 0, \quad i \geq 2$$



UNIVERSALITY OF DIFFUSION LIMIT FOR JSQ($d(N)$) SCHEME

- ▶ Fluid limit is rather crude, and diffusion limit provides more refined characterization when $\lambda \sim 1$

UNIVERSALITY OF DIFFUSION LIMIT FOR JSQ($d(N)$) SCHEME

- ▶ Fluid limit is rather crude, and diffusion limit provides more refined characterization when $\lambda \sim 1$
- ▶ Suppose total arrival rate $\lambda(N)$ satisfies $N - \lambda(N) \sim \beta\sqrt{N}$ as $N \rightarrow \infty$, so relative slack capacity is β/\sqrt{N} (Halfin-Whitt regime)

UNIVERSALITY OF DIFFUSION LIMIT FOR JSQ($d(N)$) SCHEME

- ▶ Fluid limit is rather crude, and diffusion limit provides more refined characterization when $\lambda \sim 1$
- ▶ Suppose total arrival rate $\lambda(N)$ satisfies $N - \lambda(N) \sim \beta\sqrt{N}$ as $N \rightarrow \infty$, so relative slack capacity is β/\sqrt{N} (Halfin-Whitt regime)
- ▶ Introduce **diffusion-scaled state variables**:

$$\bar{Q}_i^{d(N)}(t) = \begin{cases} \frac{Q_i^{d(N)-N}(t)}{\sqrt{N}} & \text{for } i = 1 \\ \frac{Q_i^{d(N)}(t)}{\sqrt{N}} & \text{for } i \geq 2 \end{cases}$$

UNIVERSALITY OF DIFFUSION LIMIT FOR JSQ($d(N)$) SCHEME

- ▶ Fluid limit is rather crude, and diffusion limit provides more refined characterization when $\lambda \sim 1$
- ▶ Suppose total arrival rate $\lambda(N)$ satisfies $N - \lambda(N) \sim \beta\sqrt{N}$ as $N \rightarrow \infty$, so relative slack capacity is β/\sqrt{N} (Halfin-Whitt regime)
- ▶ Introduce **diffusion-scaled state variables**:

$$\bar{Q}_i^{d(N)}(t) = \begin{cases} \frac{Q_i^{d(N)-N}(t)}{\sqrt{N}} & \text{for } i = 1 \\ \frac{Q_i^{d(N)}(t)}{\sqrt{N}} & \text{for } i \geq 2 \end{cases}$$

Universality of diffusion limit for JSQ($d(N)$) scheme [Mukherjee, B, Van L, W 2016]

Assume $d(N)/(\sqrt{N} \log(N)) \rightarrow \infty$ as $N \rightarrow \infty$. Then, under suitable initial conditions, $\{\bar{\mathbf{Q}}^{d(N)}(t)\}_{t \geq 0}$ has same weak limit $\{\bar{\mathbf{Q}}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$ as ordinary JSQ policy as established by Eschenfeldt & Gamarnik 2015

UNIVERSALITY OF DIFFUSION LIMIT FOR JSQ($d(N)$) SCHEME

Observations

- Diffusion-level behavior coincides with that of ordinary JSQ policy provided $d(N)/(\sqrt{N} \log(N)) \rightarrow \infty$ as $N \rightarrow \infty$

UNIVERSALITY OF DIFFUSION LIMIT FOR JSQ($d(N)$) SCHEME

Observations

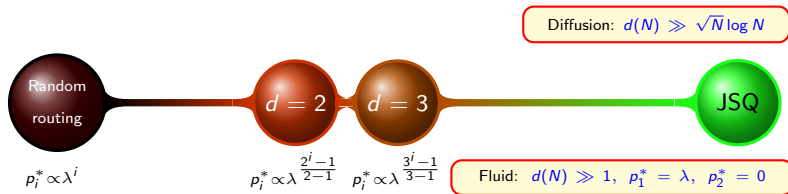
- ▶ Diffusion-level behavior coincides with that of ordinary JSQ policy provided $d(N)/(\sqrt{N} \log(N)) \rightarrow \infty$ as $N \rightarrow \infty$
- ▶ JSQ($d(N)$) scheme achieves diffusion-level optimality while reducing communication overhead by nearly a factor $O(\sqrt{N}/\log(N))$

UNIVERSALITY OF DIFFUSION LIMIT FOR JSQ($d(N)$) SCHEME

Observations

- ▶ Diffusion-level behavior coincides with that of ordinary JSQ policy provided $d(N)/(\sqrt{N} \log(N)) \rightarrow \infty$ as $N \rightarrow \infty$
- ▶ JSQ($d(N)$) scheme achieves diffusion-level optimality while reducing communication overhead by nearly a factor $O(\sqrt{N}/\log(N))$
- ▶ Latter condition is nearly necessary: if $d(N)/(\sqrt{N} \log(N)) \rightarrow 0$ as $N \rightarrow \infty$, then diffusion limit of JSQ($d(N)$) scheme differs from that of JSQ policy

HIGH-LEVEL SUMMARY



Universality and Asymptotic Optimality Properties

- ▶ $\text{JSQ}(d(N))$ has same **fluid limit** as JSQ, and achieves **fluid-level optimality** when

$$d(N) \rightarrow \infty$$

- ▶ $\text{JSQ}(d(N))$ has same **diffusion limit** as JSQ, and achieves **diffusion-level optimality** when

$$\frac{d(N)}{\sqrt{N} \log N} \rightarrow \infty$$

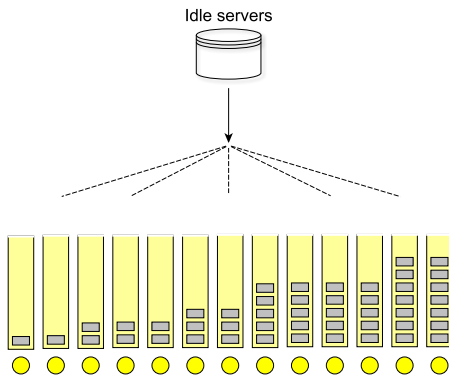
PERFORMANCE VERSUS COMMUNICATION OVERHEAD TRADE-OFF

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	Overhead
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	0
JSQ(d)	$p_i^* = \lambda^{\frac{d^i-1}{d-1}}$	$\Theta(1)$	$\Omega(\log N)$	$2d$
$d(N)$ $\rightarrow \infty$	same as JSQ	same as JSQ	??	$2d(N)$
$\frac{d(N)}{\sqrt{N \log(N)}}$ $\rightarrow \infty$	same as JSQ	same as JSQ	same as JSQ	$2d(N)$
JSQ	$p_1^* = \lambda,$ $p_2^* = 0$	$o(1)$	$\Theta(1/\sqrt{N})$	$2N$

- I. Scalability challenges and classical results
- II. Asymptotic optimality and universality (no memory)
- III. Reduction in communication overhead (memory)**
- IV. Heterogeneity issues and network scenarios

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

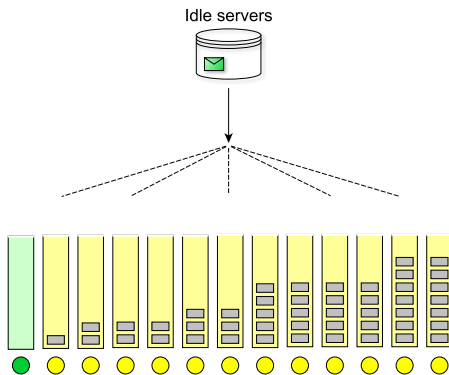
Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

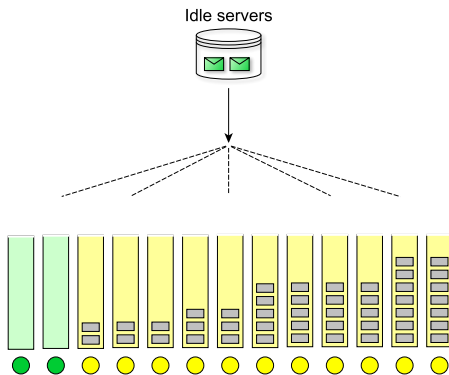
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

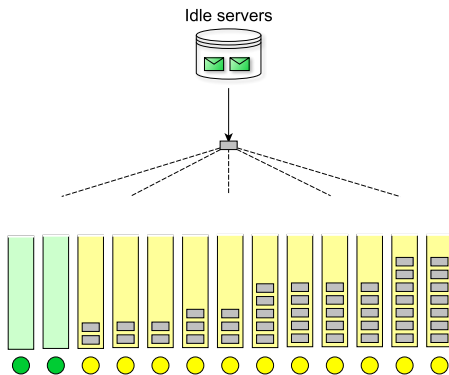
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

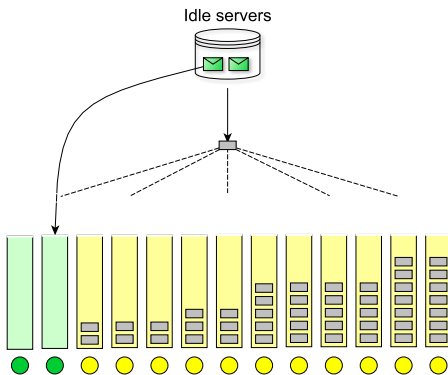
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

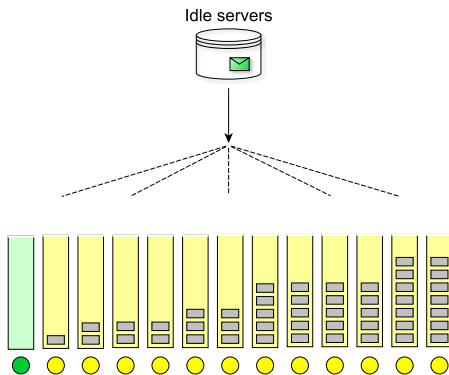
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

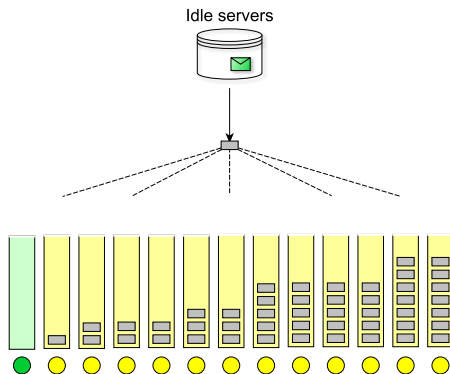
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

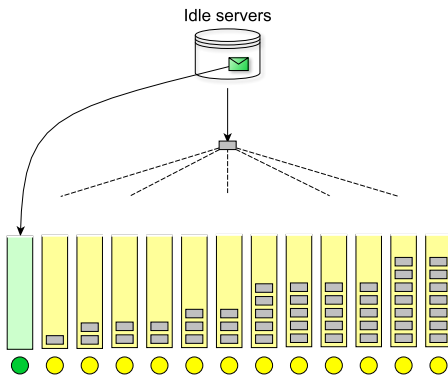
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

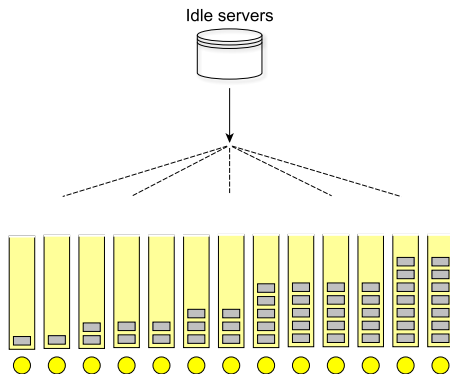
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

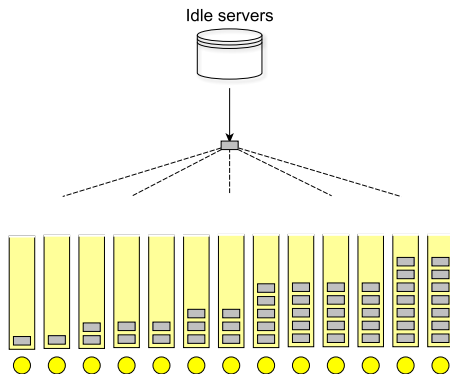
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

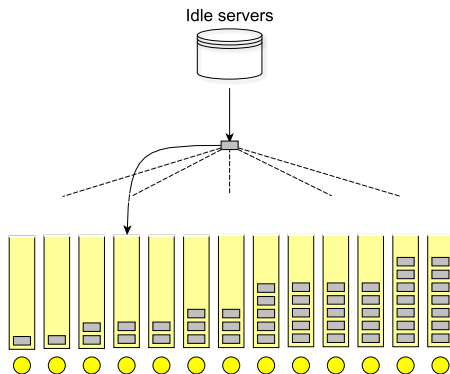
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)
- ▶ When no tokens are available arriving task is forwarded to randomly selected server



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

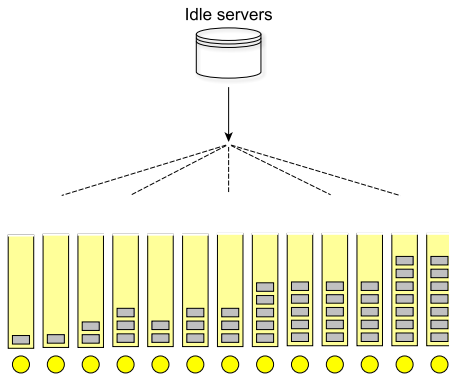
- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)
- ▶ When no tokens are available arriving task is forwarded to randomly selected server



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Server-driven rather than dispatcher-driven [Badonell & Burgess 2008, Lu *et al.* 2011]

- ▶ When server becomes idle, it sends token to dispatcher to advertise its availability
- ▶ When dispatcher has outstanding tokens at time of arrival, it assigns task to one of corresponding servers (and disposes of token)
- ▶ When no tokens are available arriving task is forwarded to randomly selected server

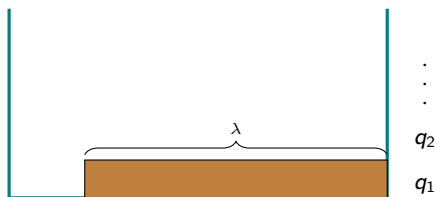


Server only sends token when service completion leaves its queue empty, implying that at most one token is generated per task

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Fixed point of fluid limit for JIQ strategy [Stolyar 2015]

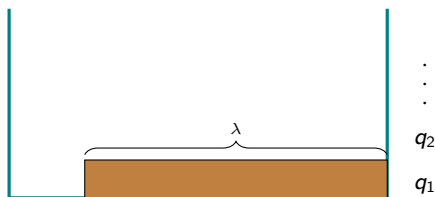
$$q_1^* = \lambda, \quad q_i^* = 0, \quad i \geq 2$$



JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Fixed point of fluid limit for JIQ strategy [Stolyar 2015]

$$q_1^* = \lambda, \quad q_i^* = 0, \quad i \geq 2$$



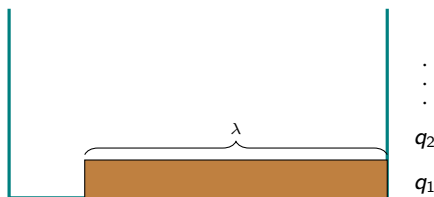
Observations

- Fluid-level behavior of JIQ strategy coincides with that of JSQ policy

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

Fixed point of fluid limit for JIQ strategy [Stolyar 2015]

$$q_1^* = \lambda, \quad q_i^* = 0, \quad i \geq 2$$



Observations

- ▶ Fluid-level behavior of JIQ strategy coincides with that of JSQ policy
- ▶ JIQ strategy achieves **fluid-level optimality** with **$O(1)$ communication overhead** per task

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

In fact, JIQ strategy provides **diffusion-level optimality** as well with **$O(1)$ communication overhead** per task

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

In fact, JIQ strategy provides **diffusion-level optimality** as well with **$O(1)$ communication overhead** per task

Diffusion limit for JIQ strategy [Mukherjee, B, Van L, W 2016]

Under suitable initial conditions, $\{\bar{\mathbf{Q}}^{JIQ}(t)\}_{t \geq 0}$ has same weak limit $\{\bar{\mathbf{Q}}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$ as ordinary JSQ policy

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

In fact, JIQ strategy provides **diffusion-level optimality** as well with **$O(1)$ communication overhead** per task

Diffusion limit for JIQ strategy [Mukherjee, B, Van L, W 2016]

Under suitable initial conditions, $\{\bar{\mathbf{Q}}^{JIQ}(t)\}_{t \geq 0}$ has same weak limit $\{\bar{\mathbf{Q}}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$ as ordinary JSQ policy

- ▶ JIQ strategy outperforms JSQ(d) schemes in terms of performance and communication overhead

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

In fact, JIQ strategy provides **diffusion-level optimality** as well with **$O(1)$ communication overhead** per task

Diffusion limit for JIQ strategy [Mukherjee, B, Van L, W 2016]

Under suitable initial conditions, $\{\bar{\mathbf{Q}}^{JIQ}(t)\}_{t \geq 0}$ has same weak limit $\{\bar{\mathbf{Q}}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$ as ordinary JSQ policy

- ▶ JIQ strategy outperforms $JSQ(d)$ schemes in terms of performance and communication overhead
- ▶ For no single value of d , $JSQ(d)$ scheme can offer both zero delay and constant overhead as $N \rightarrow \infty$

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

In fact, JIQ strategy provides **diffusion-level optimality** as well with **$O(1)$ communication overhead** per task

Diffusion limit for JIQ strategy [Mukherjee, B, Van L, W 2016]

Under suitable initial conditions, $\{\bar{\mathbf{Q}}^{JIQ}(t)\}_{t \geq 0}$ has same weak limit $\{\bar{\mathbf{Q}}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$ as ordinary JSQ policy

- ▶ JIQ strategy outperforms $JSQ(d)$ schemes in terms of performance and communication overhead
- ▶ For no single value of d , $JSQ(d)$ scheme can offer both zero delay and constant overhead as $N \rightarrow \infty$
- ▶ Note that JIQ strategy uses $O(N)$ amount of memory at dispatcher

JOIN THE IDLE QUEUE (JIQ) STRATEGY: YET LOWER OVERHEAD

In fact, JIQ strategy provides **diffusion-level optimality** as well with **$O(1)$ communication overhead** per task

Diffusion limit for JIQ strategy [Mukherjee, B, Van L, W 2016]

Under suitable initial conditions, $\{\bar{\mathbf{Q}}^{JIQ}(t)\}_{t \geq 0}$ has same weak limit $\{\bar{\mathbf{Q}}(t)\}_{t \geq 0}$ as $N \rightarrow \infty$ as ordinary JSQ policy

- ▶ **JIQ strategy outperforms JSQ(d) schemes in terms of performance and communication overhead**
- ▶ For no single value of d , JSQ(d) scheme can offer both zero delay and constant overhead as $N \rightarrow \infty$
- ▶ Note that JIQ strategy uses $O(N)$ amount of memory at dispatcher
- ▶ In order for zero delay to be achievable [Gamarnik-Tsitsiklis-Zubeldia 2016]
 - ▶ either overhead per task $\rightarrow \infty$ as $N \rightarrow \infty$,
so overhead per time unit $\gg N$
 - ▶ or amount of memory at dispatcher $\rightarrow \infty$ as $N \rightarrow \infty$

PERFORMANCE VERSUS COMMUNICATION OVERHEAD TRADE-OFF

Scheme	Queue length	Waiting time (fixed $\lambda < 1$)	Waiting time ($1 - \lambda \sim 1/\sqrt{N}$)	Overhead
Random	$p_i^N = \lambda^i$	$\frac{\lambda}{1-\lambda}$	$\Theta(\sqrt{N})$	0
JSQ(d)	$p_i^* = \lambda^{\frac{d^i-1}{d-1}}$	$\Theta(1)$	$\Omega(\log N)$	$2d$
$d(N)$ $\rightarrow \infty$	same as JSQ	same as JSQ	??	$2d(N)$
$\frac{d(N)}{\sqrt{N \log(N)}}$ $\rightarrow \infty$	same as JSQ	same as JSQ	same as JSQ	$2d(N)$
JSQ	$p_1^* = \lambda,$ $p_2^* = 0$	$o(1)$	$\Theta(1/\sqrt{N})$	$2N$
JIQ	same as JSQ	same as JSQ	same as JSQ	≤ 1

EXTENSIONS

- ▶ Multiple dispatchers
- ▶ Joint **auto-scaling** and **load balancing**:
TABS scheme simultaneously achieves **zero delay** and **zero relative energy wastage** in many-server regime

Each server provides periodic queue status updates to dispatcher
[Van der Boor, B, Van Leeuwen 2018]

- ▶ Either in synchronized or asynchronous manner
- ▶ At either exponentially distributed or constant time intervals with (mean) duration T and frequency $\delta = 1/T$, so overhead per task is δ/λ

Each server provides periodic queue status updates to dispatcher
[Van der Boor, B, Van Leeuwen 2018]

- ▶ Either in synchronized or asynchronous manner
- ▶ At either exponentially distributed or constant time intervals with (mean) duration T and frequency $\delta = 1/T$, so overhead per task is δ/λ
- ▶ SUJSQ^{det}: Synchronized Constant Update Intervals

Each server provides periodic queue status updates to dispatcher
[Van der Boor, B, Van Leeuwen 2018]

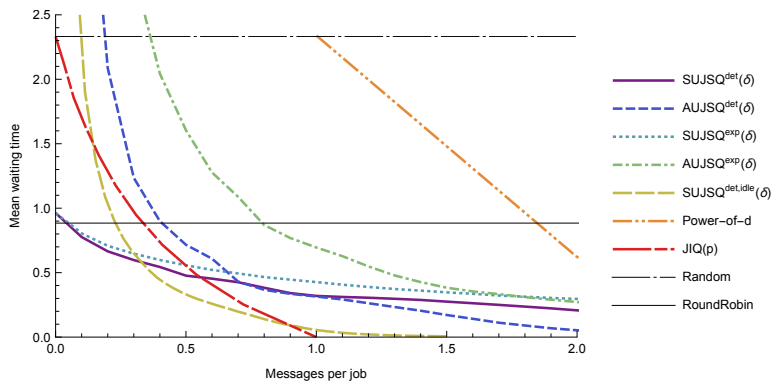
- ▶ Either in synchronized or asynchronous manner
- ▶ At either exponentially distributed or constant time intervals with (mean) duration T and frequency $\delta = 1/T$, so overhead per task is δ/λ
- ▶ SUJSQ^{det}: Synchronized Constant Update Intervals
- ▶ AUJSQ^{exp}: Asynchronous Exponentially Distributed Update Intervals

Each server provides periodic queue status updates to dispatcher
[Van der Boor, B, Van Leeuwen 2018]

- ▶ Either in synchronized or asynchronous manner
- ▶ At either exponentially distributed or constant time intervals with (mean) duration T and frequency $\delta = 1/T$, so overhead per task is δ/λ
- ▶ SUJSQ^{det}: Synchronized Constant Update Intervals
- ▶ AUJSQ^{exp}: Asynchronous Exponentially Distributed Update Intervals

When task arrives, dispatcher assigns it server with minimum value of queue estimate (reset at updates, and incremented for each assigned task)

HYPER-SCALABLE JSQ



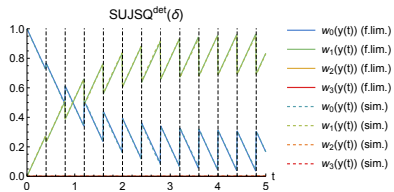
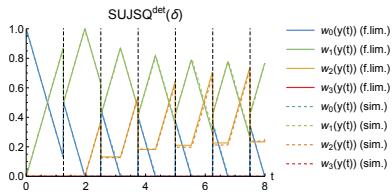
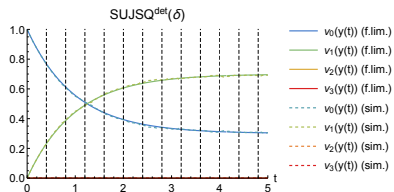
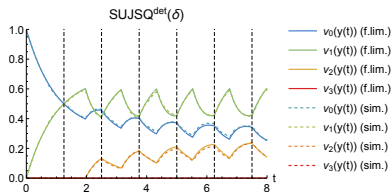
FLUID LIMITS FOR HYPER-SCALABLE JSQ

- ▶ $Y_{i,j}^N(t)$: number of servers with actual queue length i and queue estimate $j \geq i$ at time t

FLUID LIMITS FOR HYPER-SCALABLE JSQ

- ▶ $Y_{i,j}^N(t)$: number of servers with actual queue length i and queue estimate $j \geq i$ at time t
- ▶ $y_{i,j}(t) = \lim_{N \rightarrow \infty} Y_{i,j}^N(t)/N$: fluid-limit variables
- ▶ $v_i(t) = \sum_{j=i}^{\infty} y_{i,j}(t)$: fraction of servers with queue length i
- ▶ $w_j(t) = \sum_{i=0}^j y_{i,j}(t)$: fraction of servers with queue estimate j

SUJSQ^{det}: SYNCHRONIZED CONSTANT UPDATE INTERVALS



Update frequency $\delta = 0.85$

Update frequency $\delta = 2.5$

SUJSQ^{det}: SYNCHRONIZED CONSTANT UPDATE INTERVALS

Zero-delay threshold SUJSQ^{det}(δ)

If update frequency $\delta > \frac{\lambda}{1-\lambda}$, i.e., at least $\frac{1}{1-\lambda}$ messages per task, then no queues and zero delay in stationarity at fluid level

SUJSQ^{det}: SYNCHRONIZED CONSTANT UPDATE INTERVALS

Zero-delay threshold SUJSQ^{det}(δ)

If update frequency $\delta > \frac{\lambda}{1-\lambda}$, i.e., at least $\frac{1}{1-\lambda}$ messages per task, then no queues and zero delay in stationarity at fluid level

Finite maximum queue length SUJSQ^{det}($1/T$)

For any $\lambda < 1$, maximum queue length in stationarity at fluid level is bounded from above by

$$M(\lambda, T) = \max\{L : \lambda T \geq H(L; \lambda, T)\} < \infty,$$

with

$$H(L; \lambda, T) = \left(1 - \frac{\lambda T}{L}\right) \left[\sum_{l=0}^L l e^{-T} \frac{T^l}{l!} + L \sum_{l=L+1}^{\infty} e^{-T} \frac{T^l}{l!} \right]$$

SUJSQ^{det}: SYNCHRONIZED CONSTANT UPDATE INTERVALS

Zero-delay threshold SUJSQ^{det}(δ)

If update frequency $\delta > \frac{\lambda}{1-\lambda}$, i.e., at least $\frac{1}{1-\lambda}$ messages per task, then no queues and zero delay in stationarity at fluid level

Finite maximum queue length SUJSQ^{det}($1/T$)

For any $\lambda < 1$, maximum queue length in stationarity at fluid level is bounded from above by

$$M(\lambda, T) = \max\{L : \lambda T \geq H(L; \lambda, T)\} < \infty,$$

with

$$H(L; \lambda, T) = \left(1 - \frac{\lambda T}{L}\right) \left[\sum_{l=0}^L l e^{-T} \frac{T^l}{l!} + L \sum_{l=L+1}^{\infty} e^{-T} \frac{T^l}{l!} \right]$$

- Upper bound $M(\lambda, T) \rightarrow \infty$ as $T \rightarrow \infty$

SUJSQ^{det}: SYNCHRONIZED CONSTANT UPDATE INTERVALS

Zero-delay threshold SUJSQ^{det}(δ)

If update frequency $\delta > \frac{\lambda}{1-\lambda}$, i.e., at least $\frac{1}{1-\lambda}$ messages per task, then no queues and zero delay in stationarity at fluid level

Finite maximum queue length SUJSQ^{det}($1/T$)

For any $\lambda < 1$, maximum queue length in stationarity at fluid level is bounded from above by

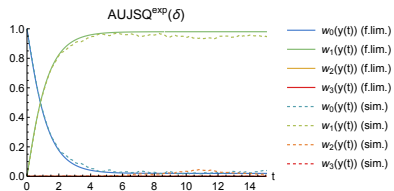
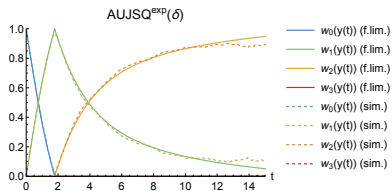
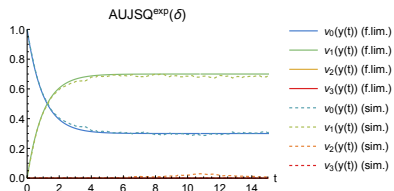
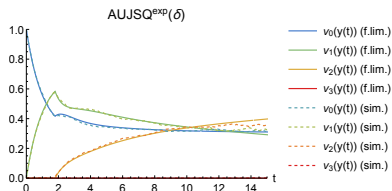
$$M(\lambda, T) = \max\{L : \lambda T \geq H(L; \lambda, T)\} < \infty,$$

with

$$H(L; \lambda, T) = \left(1 - \frac{\lambda T}{L}\right) \left[\sum_{l=0}^L l e^{-T} \frac{T^l}{l!} + L \sum_{l=L+1}^{\infty} e^{-T} \frac{T^l}{l!} \right]$$

- ▶ Upper bound $M(\lambda, T) \rightarrow \infty$ as $T \rightarrow \infty$
- ▶ SUJSQ^{det}($1/T$) reduces to Round Robin policy as $T \rightarrow \infty$, and each individual server behaves as D/M/1 queue

AUJSQ^{exp}: ASYNCHRONOUS EXPONENTIALLY DISTRIBUTED UPDATE INTERVALS



Update frequency $\delta = 0.85$

Update frequency $\delta = 2.5$

AUJSQ^{exp}: ASYNCHRONOUS EXPONENTIALLY DISTRIBUTED UPDATE INTERVALS

Zero-delay threshold AUJSQ^{exp}(δ)

If update frequency $\delta > \frac{\lambda}{1-\lambda}$, i.e., at least $\frac{1}{1-\lambda}$ messages per task, then no queues and zero delay in stationarity at fluid level

AUJSQ^{exp}: ASYNCHRONOUS EXPONENTIALLY DISTRIBUTED UPDATE INTERVALS

Zero-delay threshold AUJSQ^{exp}(δ)

If update frequency $\delta > \frac{\lambda}{1-\lambda}$, i.e., at least $\frac{1}{1-\lambda}$ messages per task, then no queues and zero delay in stationarity at fluid level

Finite maximum queue length AUJSQ^{exp}(δ)

For any $\lambda < 1$, maximum queue length in stationarity at fluid level is

$$M(\lambda, \delta) = \left\lfloor -\frac{\log(1-\lambda)}{\log(1+\delta)} \right\rfloor < \infty$$

AUJSQ^{exp}: ASYNCHRONOUS EXPONENTIALLY DISTRIBUTED UPDATE INTERVALS

Zero-delay threshold AUJSQ^{exp}(δ)

If update frequency $\delta > \frac{\lambda}{1-\lambda}$, i.e., at least $\frac{1}{1-\lambda}$ messages per task, then no queues and zero delay in stationarity at fluid level

Finite maximum queue length AUJSQ^{exp}(δ)

For any $\lambda < 1$, maximum queue length in stationarity at fluid level is

$$M(\lambda, \delta) = \left\lfloor -\frac{\log(1-\lambda)}{\log(1+\delta)} \right\rfloor < \infty$$

- Maximum queue length $M(\lambda, \delta) \rightarrow \infty$ as $\delta \downarrow 0$

AUJSQ^{exp}: ASYNCHRONOUS EXPONENTIALLY DISTRIBUTED UPDATE INTERVALS

Zero-delay threshold AUJSQ^{exp}(δ)

If update frequency $\delta > \frac{\lambda}{1-\lambda}$, i.e., at least $\frac{1}{1-\lambda}$ messages per task, then no queues and zero delay in stationarity at fluid level

Finite maximum queue length AUJSQ^{exp}(δ)

For any $\lambda < 1$, maximum queue length in stationarity at fluid level is

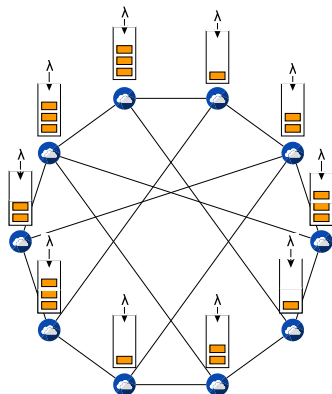
$$M(\lambda, \delta) = \left\lfloor -\frac{\log(1-\lambda)}{\log(1+\delta)} \right\rfloor < \infty$$

- ▶ Maximum queue length $M(\lambda, \delta) \rightarrow \infty$ as $\delta \downarrow 0$
- ▶ Queue length at individual server under SUJSQ^{det}($1/T$) wildly oscillates as $\delta \downarrow 0$,
rising to $M(\lambda, \delta)$ shortly after each update and then gradually dropping to low value before next update

- I. Scalability challenges and classical results
- II. Asymptotic optimality and universality (no memory)
- III. Reduction in communication overhead (memory)
- IV. Heterogeneity issues and network scenarios**

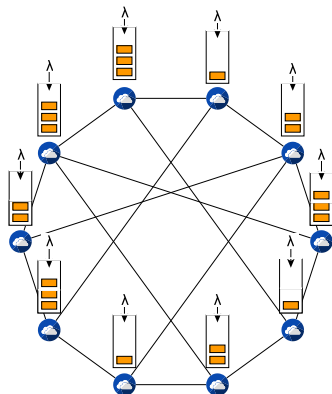
LOAD BALANCING IN NETWORK SCENARIOS

- Suppose N servers are interconnected by underlying graph topology G_N



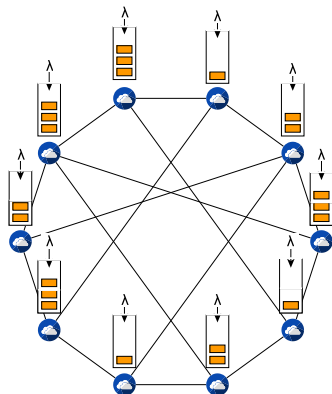
LOAD BALANCING IN NETWORK SCENARIOS

- ▶ Suppose N servers are interconnected by underlying graph topology G_N
- ▶ Tasks arrive at each server as Poisson process of rate λ and can either be executed locally or forwarded to neighbor in G_N with shorter queue, if any



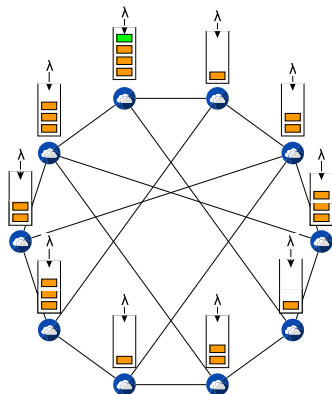
LOAD BALANCING IN NETWORK SCENARIOS

- ▶ Suppose N servers are interconnected by underlying graph topology G_N
- ▶ Tasks arrive at each server as Poisson process of rate λ and can either be executed locally or forwarded to neighbor in G_N with shorter queue, if any
- ▶ Data 'locality', server-task affinity relations, compatibility constraints



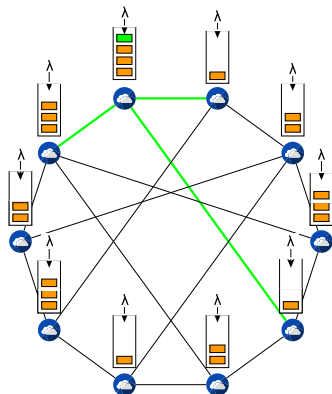
LOAD BALANCING IN NETWORK SCENARIOS

- ▶ Suppose N servers are interconnected by underlying graph topology G_N
- ▶ Tasks arrive at each server as Poisson process of rate λ and can either be executed locally or forwarded to neighbor in G_N with shorter queue, if any
- ▶ Data 'locality', server-task affinity relations, compatibility constraints



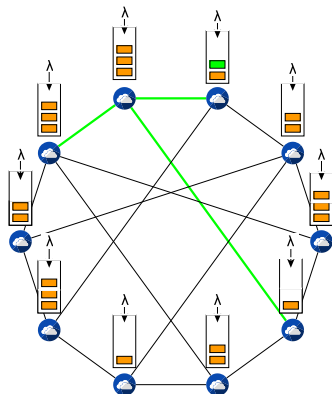
LOAD BALANCING IN NETWORK SCENARIOS

- ▶ Suppose N servers are interconnected by underlying graph topology G_N
- ▶ Tasks arrive at each server as Poisson process of rate λ and can either be executed locally or forwarded to neighbor in G_N with shorter queue, if any
- ▶ Data 'locality', server-task affinity relations, compatibility constraints



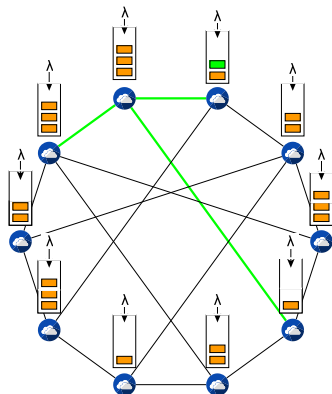
LOAD BALANCING IN NETWORK SCENARIOS

- ▶ Suppose N servers are interconnected by underlying graph topology G_N
- ▶ Tasks arrive at each server as Poisson process of rate λ and can either be executed locally or forwarded to neighbor in G_N with shorter queue, if any
- ▶ Data 'locality', server-task affinity relations, compatibility constraints



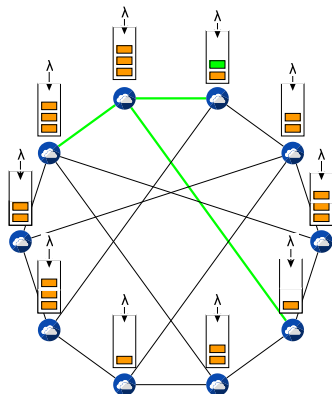
LOAD BALANCING IN NETWORK SCENARIOS

- ▶ Suppose N servers are interconnected by underlying graph topology G_N
- ▶ Tasks arrive at each server as Poisson process of rate λ and can either be executed locally or forwarded to neighbor in G_N with shorter queue, if any
- ▶ Data 'locality', server-task affinity relations, compatibility constraints
- ▶ How does graph structure affect performance?



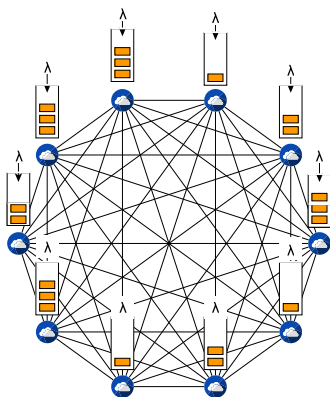
LOAD BALANCING IN NETWORK SCENARIOS

- ▶ Suppose N servers are interconnected by underlying graph topology G_N
- ▶ Tasks arrive at each server as Poisson process of rate λ and can either be executed locally or forwarded to neighbor in G_N with shorter queue, if any
- ▶ Data 'locality', server-task affinity relations, compatibility constraints
- ▶ How does graph structure affect performance?
- ▶ Lack of exchangeability among servers breaks underpinning for stochastic coupling and fluid and diffusion limits



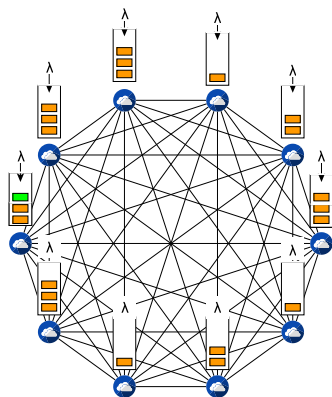
LOAD BALANCING IN NETWORK SCENARIOS: COMPLETE GRAPH

Case of complete graph (clique) corresponds to ordinary supermarket model



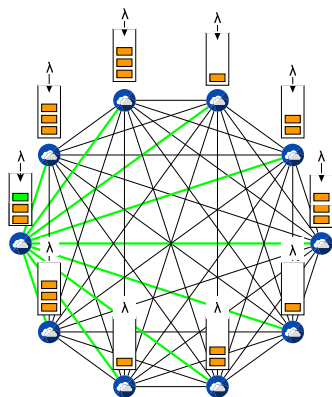
LOAD BALANCING IN NETWORK SCENARIOS: COMPLETE GRAPH

Case of complete graph (clique) corresponds to ordinary supermarket model



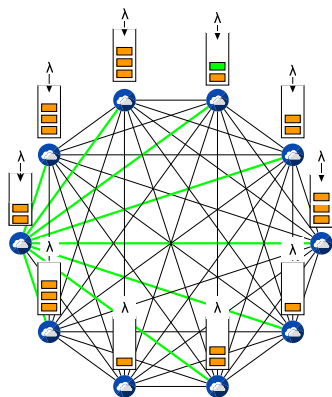
LOAD BALANCING IN NETWORK SCENARIOS: COMPLETE GRAPH

Case of complete graph (clique) corresponds to ordinary supermarket model



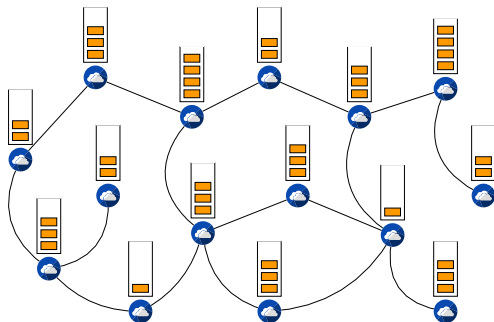
LOAD BALANCING IN NETWORK SCENARIOS: COMPLETE GRAPH

Case of complete graph (clique) corresponds to ordinary supermarket model



OUR PERSPECTIVE: ASYMPTOTIC OPTIMALITY AND UNIVERSALITY

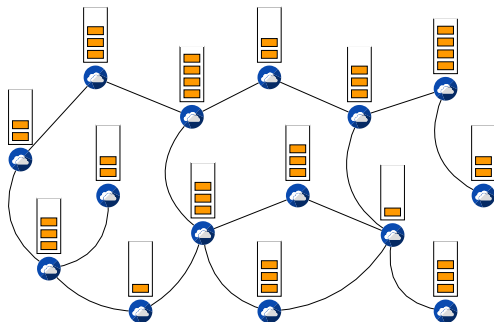
How much connectivity is required in order for JSQ to achieve asymptotically similar performance in G_N as in complete graph as $N \rightarrow \infty$?



OUR PERSPECTIVE: ASYMPTOTIC OPTIMALITY AND UNIVERSALITY

How much connectivity is required in order for JSQ to achieve asymptotically similar performance in G_N as in complete graph as $N \rightarrow \infty$?

Slightly different criterion: what degree of connectivity is required for JSQ(d) to yield asymptotically similar performance in G_N as in complete graph



JSQ ON DETERMINISTIC GRAPHS: INFORMAL STATEMENTS

Let $\{G_N\}_{N \geq 1}$ be sequence of graphs

Condition 1

Neighborhood size of any $\Theta(N)$ vertices is $N - o(N)$

Condition 2

Neighborhood size of any $\Theta(\sqrt{N})$ vertices is $N - o(\sqrt{N})$

JSQ ON DETERMINISTIC GRAPHS: INFORMAL STATEMENTS

Let $\{G_N\}_{N \geq 1}$ be sequence of graphs

Condition 1

Neighborhood size of any $\Theta(N)$ vertices is $N - o(N)$

Condition 2

Neighborhood size of any $\Theta(\sqrt{N})$ vertices is $N - o(\sqrt{N})$

Graph sequence $\{G_N\}_{N \geq 1}$ is said to be **fluid-optimal** or **diffusion-optimal** if JSQ in this graph sequence yields same **fluid limit** and **diffusion limit** as in classical setup, respectively

JSQ ON DETERMINISTIC GRAPHS: INFORMAL STATEMENTS

Let $\{G_N\}_{N \geq 1}$ be sequence of graphs

Condition 1

Neighborhood size of any $\Theta(N)$ vertices is $N - o(N)$

Condition 2

Neighborhood size of any $\Theta(\sqrt{N})$ vertices is $N - o(\sqrt{N})$

Graph sequence $\{G_N\}_{N \geq 1}$ is said to be **fluid-optimal** or **diffusion-optimal** if JSQ in this graph sequence yields same **fluid limit** and **diffusion limit** as in classical setup, respectively

Theorem JSQ on deterministic graphs

Graph sequence $\{G_N\}_{N \geq 1}$ is

- (i) **fluid-optimal** if **Condition 1** is satisfied
- (ii) **diffusion-optimal** if **Condition 2** is satisfied

Theorem JSQ on random graphs

Sequence of (Erdős-Rényi or random regular) graphs with avg degree $c(N)$ is

- ▶ fluid-optimal if $c(N) \rightarrow \infty$
- ▶ diffusion-optimal if $c(N)/(\sqrt{N} \log N) \rightarrow \infty$

Theorem Worst-case scenario

For any graph sequence $\{G_N\}_{N \geq 1}$, if

- ▶ $d_{\min}(G_N)/N \rightarrow 1$, then sequence **must be** fluid-optimal
- ▶ $d_{\min}(G_N)/N \rightarrow c < 1/2$, then sequence **may not be** fluid-optimal
- ▶ $\#$ bounded degree vertices is $\Theta(N)$, then sequence **is not** fluid-optimal

CONCLUSION

- ▶ Much sparser graphs can asymptotically match optimal performance of complete graph, provided they are suitably random
- ▶ In worst-case scenario, performance can be sub-optimal even when graph is sufficiently dense

SOME REFERENCES

M. van der Boor, S.C. Borst, J.S.H. van Leeuwen, D. Mukherjee (2018). Scalable load balancing in networked systems: Universality properties and stochastic coupling methods In: *Proc. ICM 2018*, Rio de Janeiro.

M. van der Boor, S.C. Borst, J.S.H. van Leeuwen, D. Mukherjee (2018). Scalable load balancing in networked systems: A survey of recent advances.

D. Mukherjee (2018). Scalable load balancing in networked systems. PhD thesis, Eindhoven University of Technology.

Thank You for your Attention!!