

Scalable Load Balancing in the Presence of Heterogeneous Servers

Kristen Gardner¹, Jazeem Abdul Jaleel², Alexander Wickeham², Sherwin Doroudi²

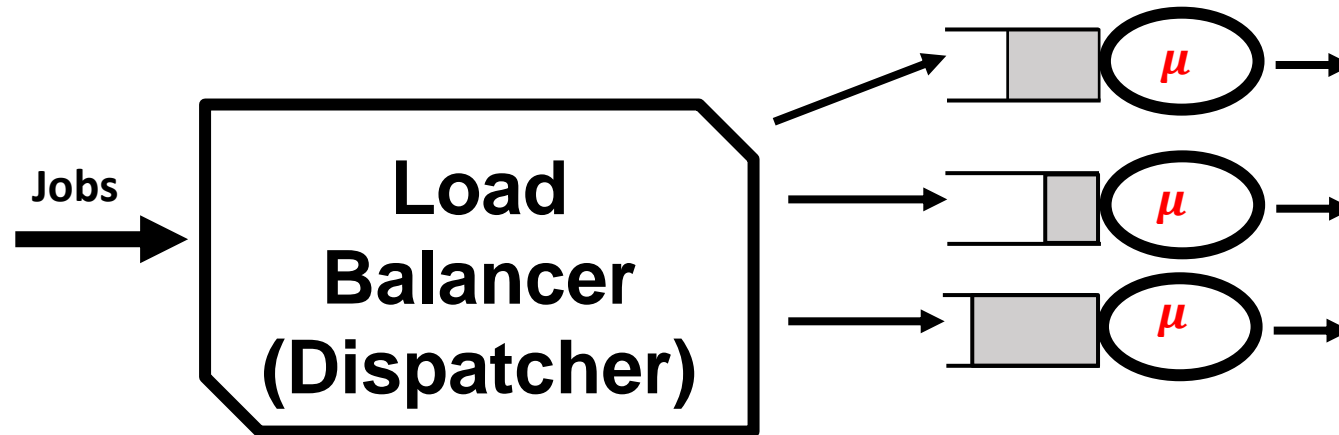
Performance 2020, November 3rd 2020

¹ Computer Science Department, Amherst College,
² Industrial and Systems Engineering, University of Minnesota



Introduction

Scalable **Load Balancing** in the Presence of Heterogeneous Servers



- Load Balancing – homogeneous servers

Load Balancing Policies	
Random	Round Robin
Join Shortest Queue (JSQ)	Join Idle Queue (JIQ)

Examples

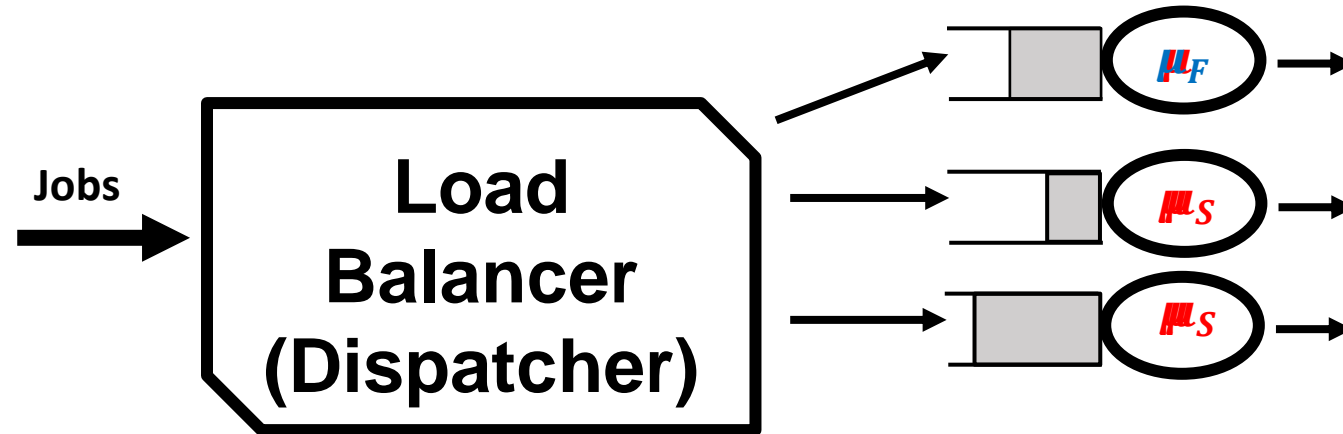
- AWS
- IBM network dispatcher
- Citrix ADC
- Cisco Local director

Important

Weber 1978, Winston 1977, Lu *et al* 2011, Wang *et al* 2018

Introduction

Scalable Load Balancing in the Presence of Heterogeneous Servers



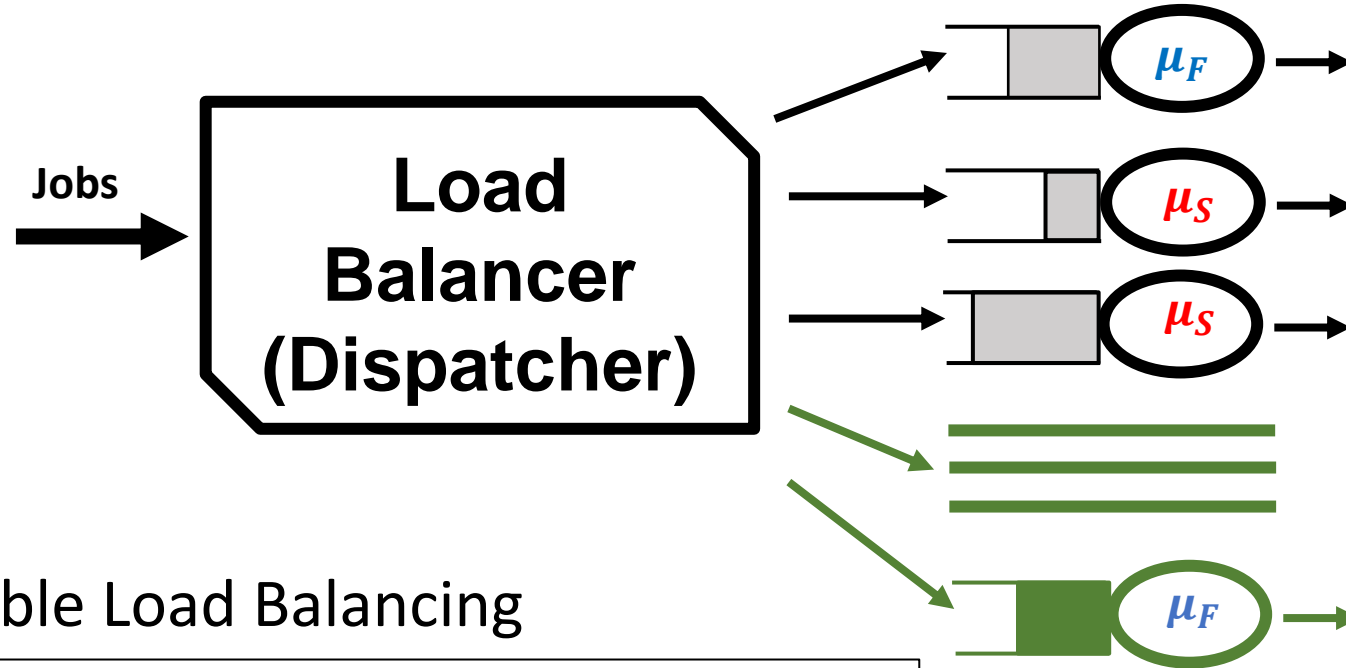
- Load Balancing of heterogeneous Server Systems (Slow server problems)

Load balancing policies

- Threshold type
Hyytiä 2013, Luh et al 2002

Introduction

Scalable Load Balancing in the Presence of Heterogeneous Servers



Difficulties

- Complex
- Error prone
- Costly

- Scalable Load Balancing

Power of d policies (query d)

- **Join Shortest Queue – query d (JSQ- d)**
- **Join Idle Queue – query d (JIQ- d)**

Mitzenmacher 2001

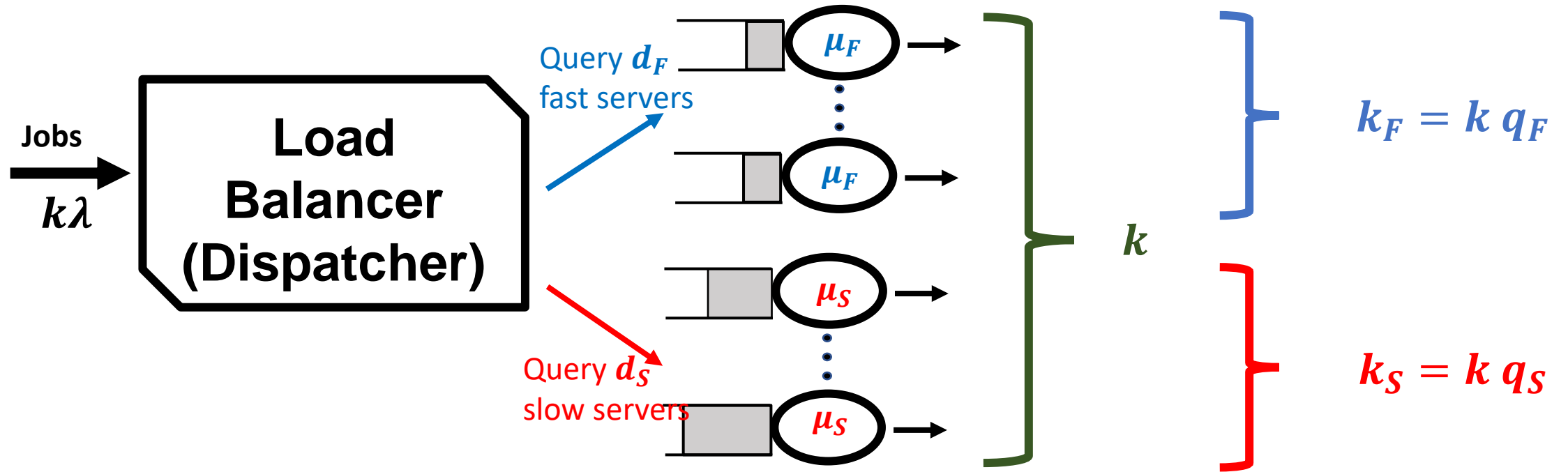
Bramson et al 2012

Li et al. 2014

Outline

- System model
- Explain the policies:
 - Join Idle Queue – (d_F, d_S)
 - Join Shortest Queue – (d_F, d_S)
- Stability of the proposed policies
- Techniques used and the optimization problem
- Performance comparison

Model



Policies:

1. Join Idle Queue – (d_F, d_S)
2. Join Shortest Queue – (d_F, d_S)

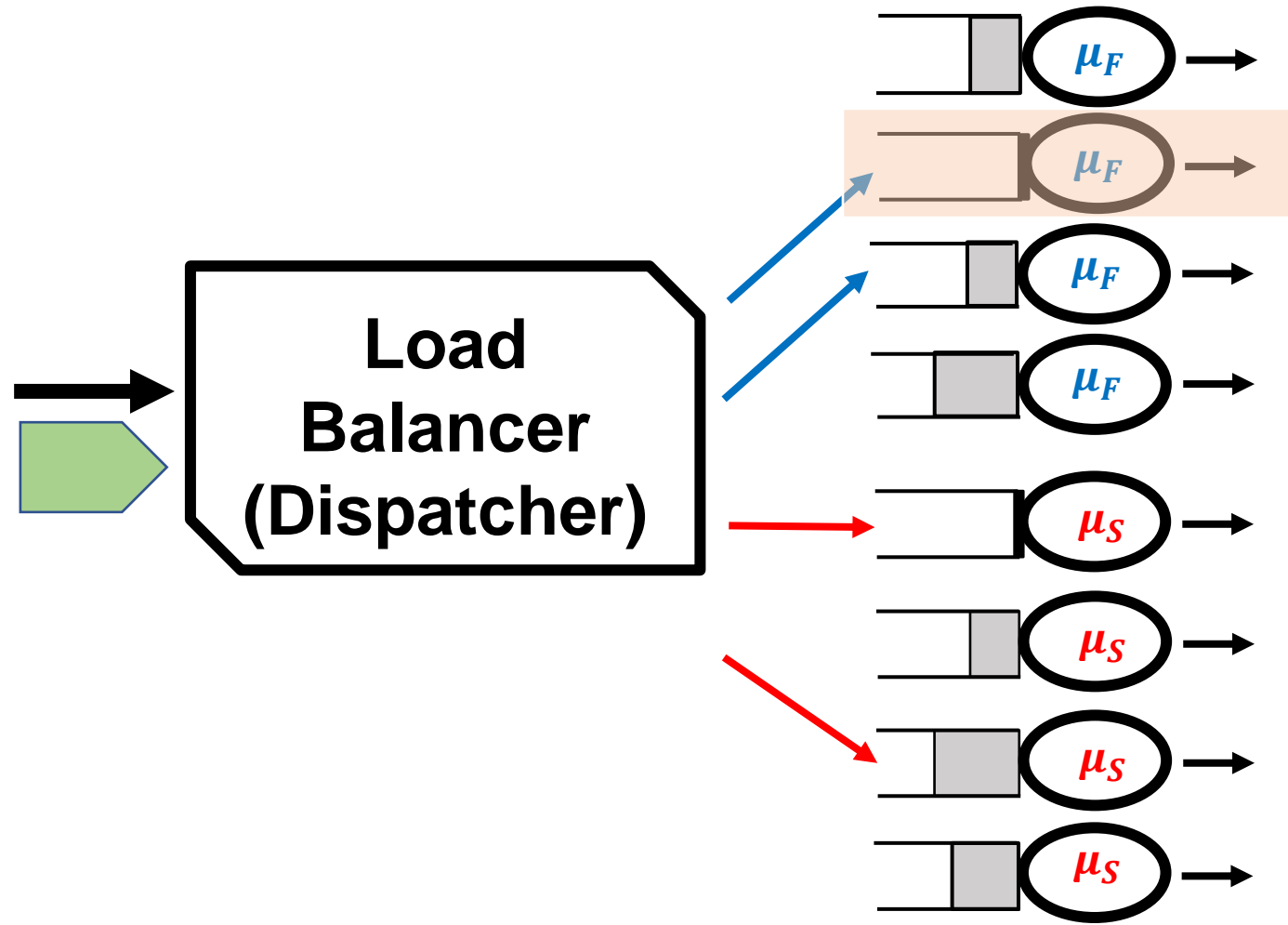
JIQ – (d_F, d_S)

JSQ – (d_F, d_S)

JIQ – (d_F, d_S)

Query randomly d_F fast servers and d_S slow servers.

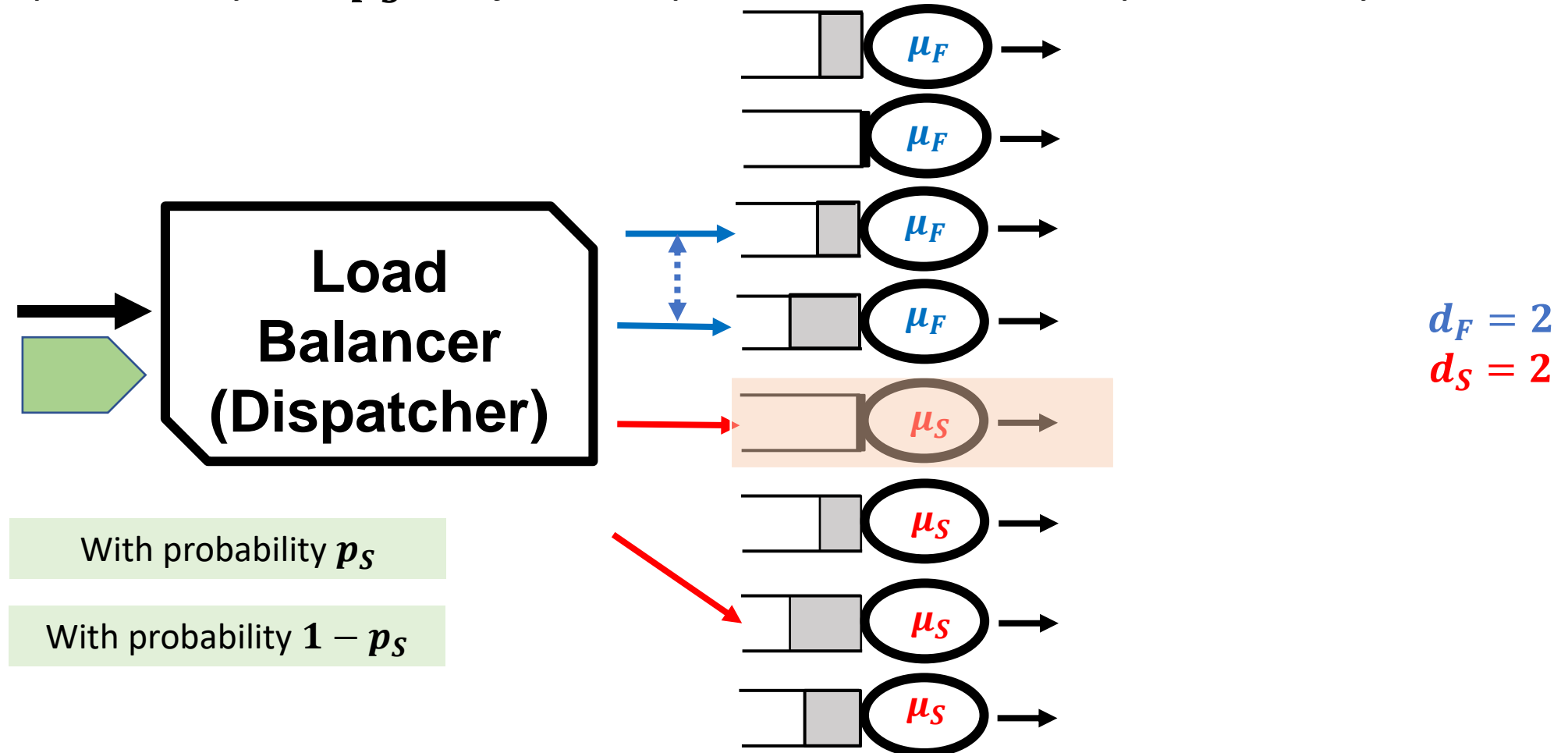
1. If any of the d_F fast servers are idle, the job begins service on one of them.



$$d_F = 2$$
$$d_S = 2$$

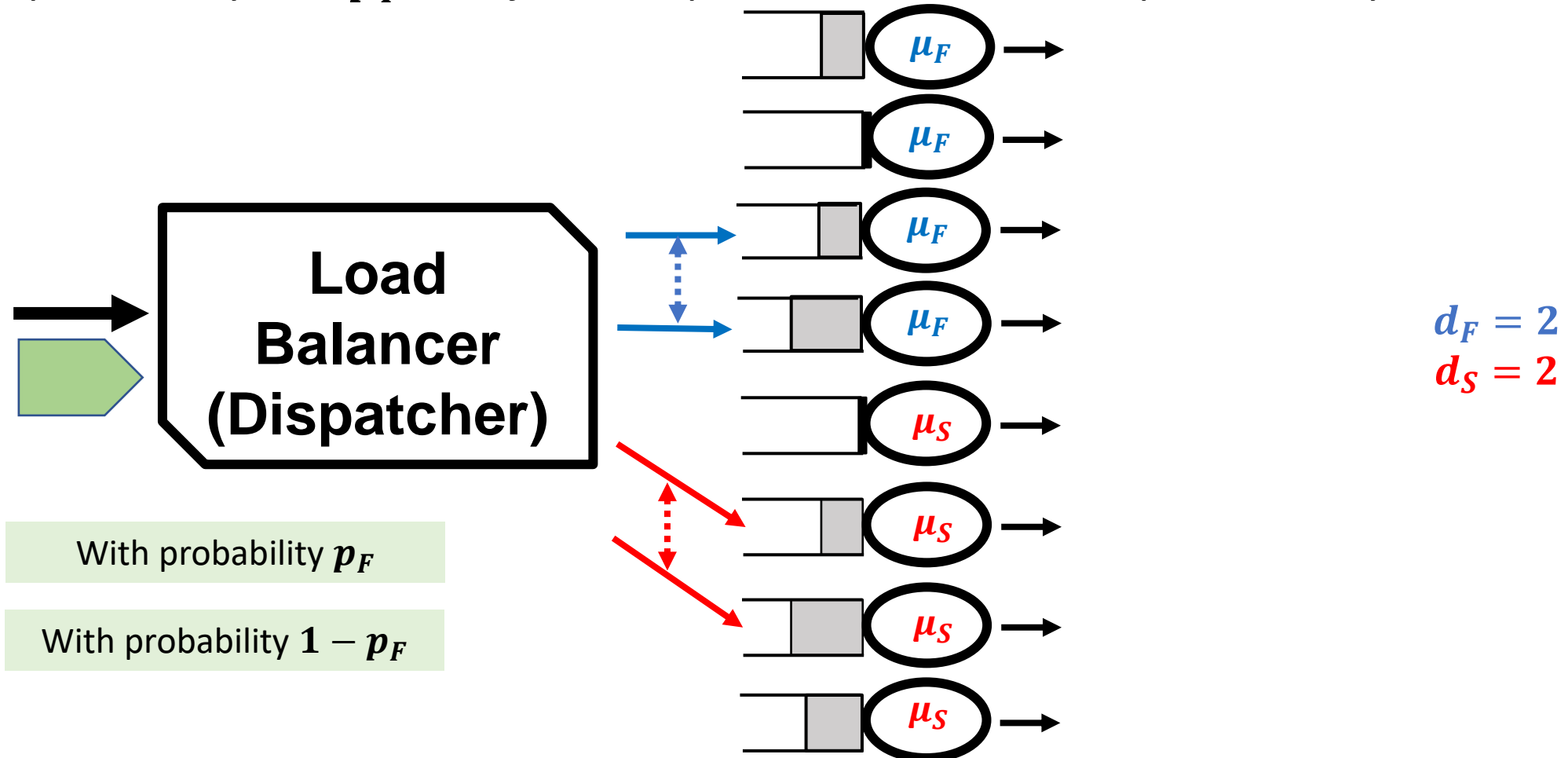
JIQ – (d_F, d_S)

2. If all d_F fast servers are busy and any of the d_S slow servers are idle,
- with probability p_S , the job begins service on a random queried idle slow server.
 - with probability $1 - p_S$ the job is dispatched to a random queried busy fast server



JIQ – (d_F, d_S)

3. If all $d_F + d_S$ queried servers are busy,
- with probability p_F , the job is dispatched to a random queried busy fast server.
 - with probability $1 - p_F$, the job is dispatched to a random queried busy slow server.



JIQ – (d_F, d_S)

Query randomly d_F fast servers and d_S slow servers.

1. If any of the d_F fast servers are idle, the job begins service on one of them.
2. If all d_F fast servers are busy and any of the d_S slow servers are idle,
 - with probability p_S , the job begins service on a random queried idle slow server.
 - with probability $1 - p_S$ the job is dispatched to a random queried busy fast server
3. If all $d_F + d_S$ queried servers are busy,
 - with probability p_F , the job is dispatched to a random queried busy fast server.
 - with probability $1 - p_F$, the job is dispatched to a random queried busy slow server.



JSQ – (d_F, d_S)

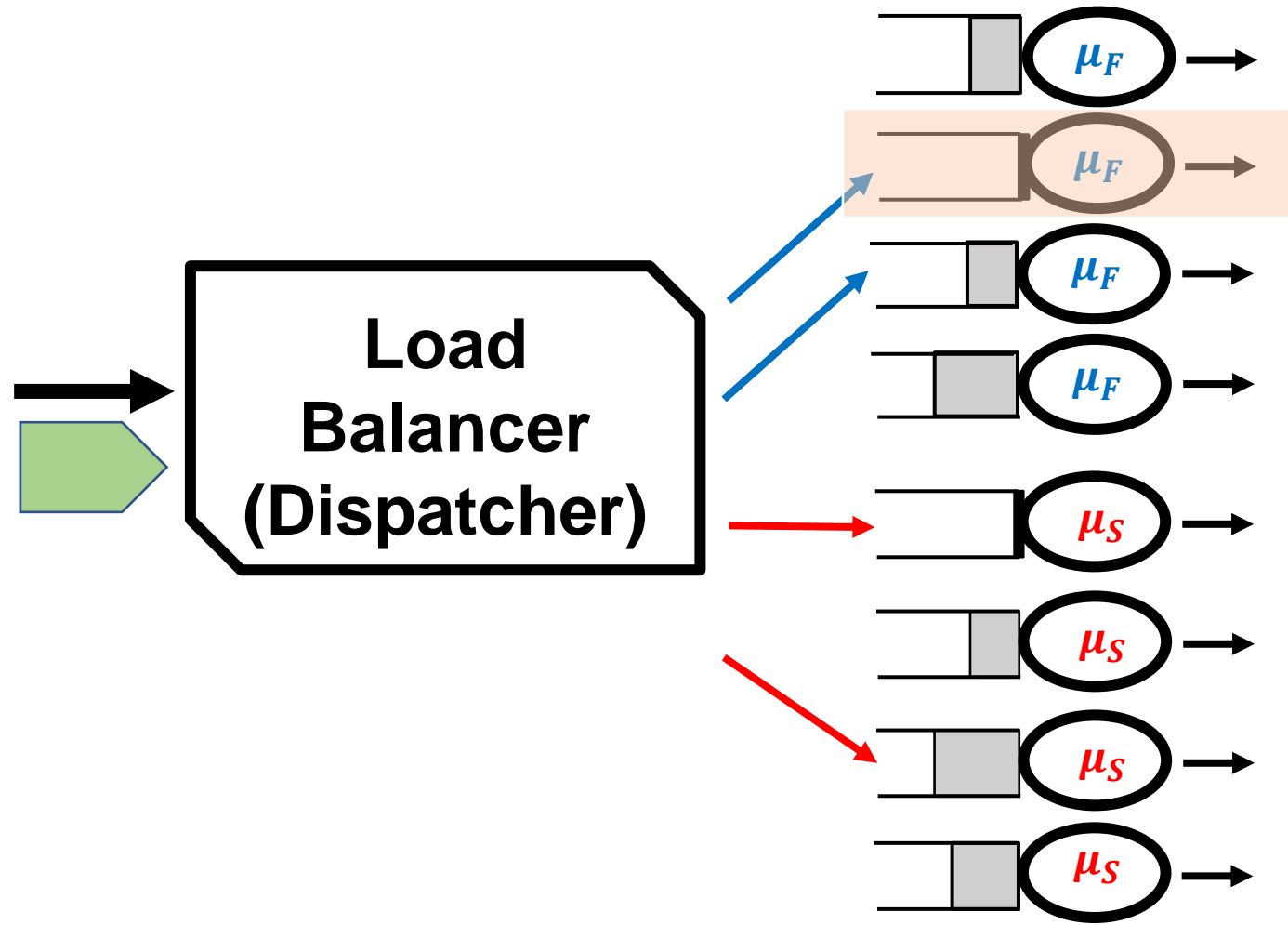
When job arrives, query at random d_F fast servers and d_S slow servers.

1. If any of the d_F fast servers are idle, the job begins service on one of them.
2. If all d_F fast servers are busy and any of the d_S slow servers are idle,
 - with probability p_S , the job begins service on a random queried idle slow server.
 - with probability $1 - p_S$, the job is dispatched to the fast server with shortest queue among the d_F queried fast servers
3. If all $d_F + d_S$ queried servers are busy,
 - with probability p_F , the job is dispatched to the fast server with shortest queue among the d_F queried fast servers
 - with probability $1 - p_F$, the job is dispatched to the slow server with shortest queue among the d_S queried slow servers

JSQ – (d_F, d_S)

Query randomly d_F fast servers d_S slow servers.

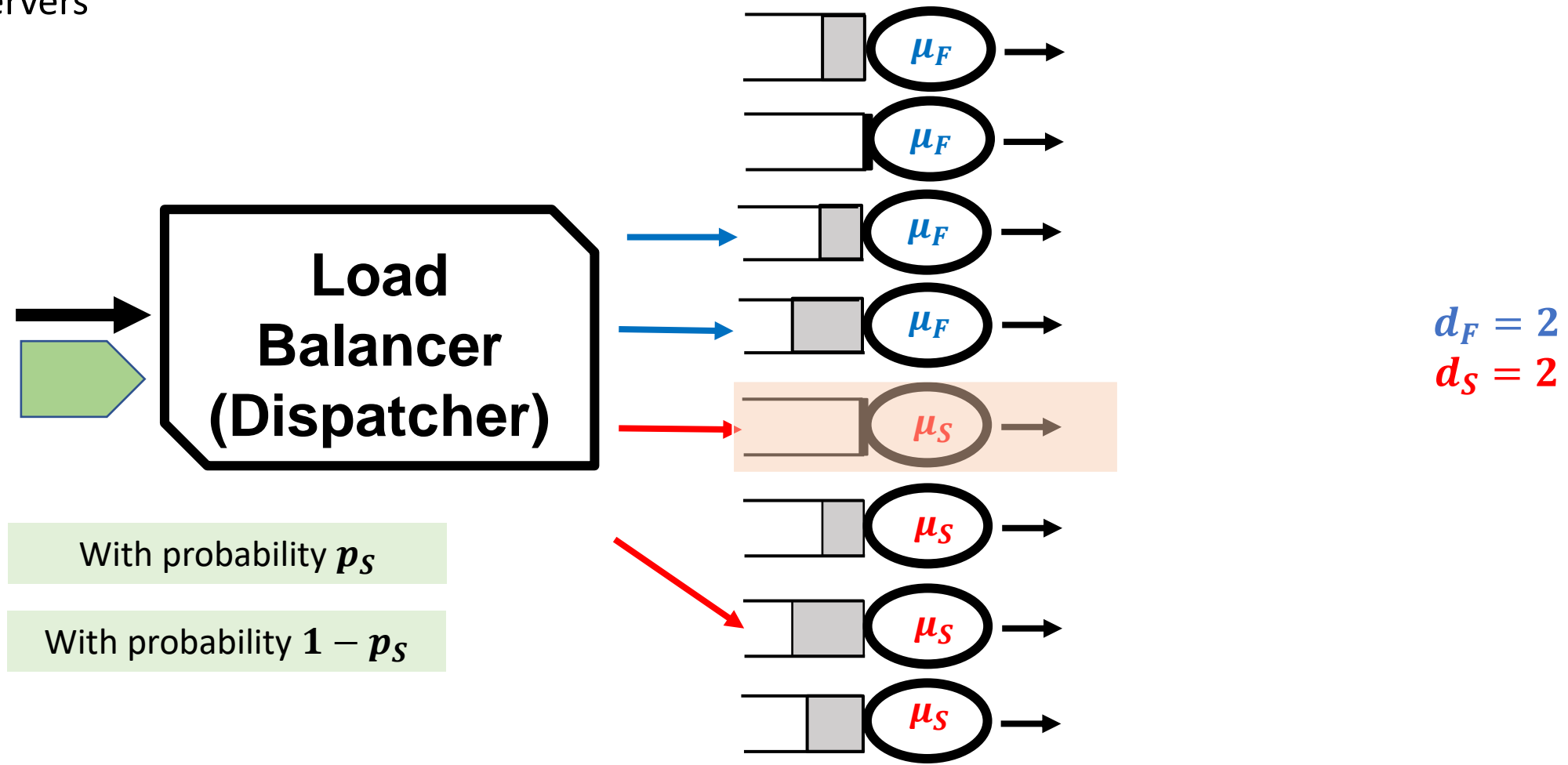
1. If any of the d_F fast servers are idle, the job begins service on one of them.



$$d_F = 2$$
$$d_S = 2$$

JSQ – (d_F, d_S)

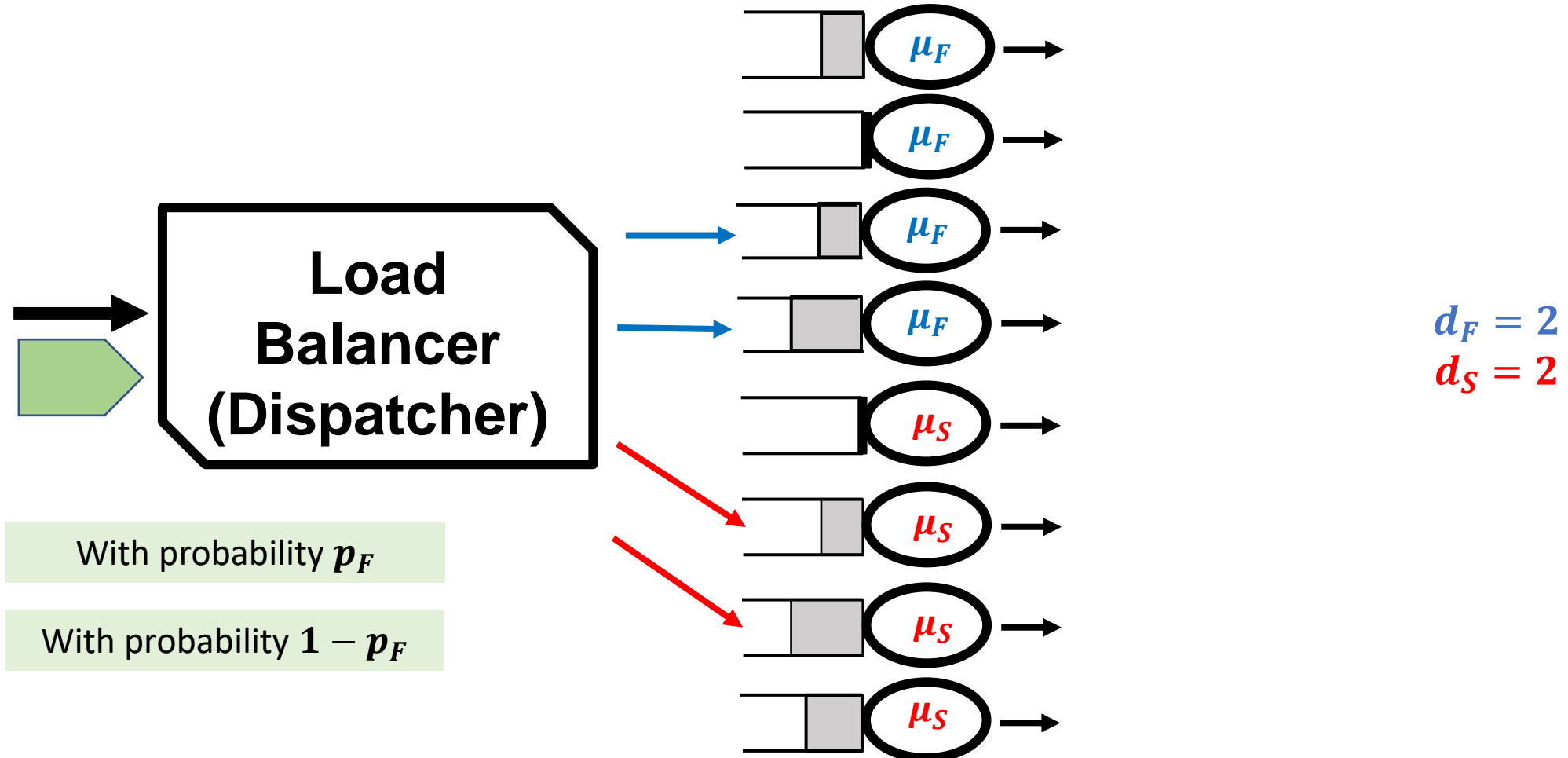
2. If all d_F fast servers are busy and any of the d_S slow servers are idle,
- with probability p_S , the job begins service on a random queried idle slow server.
 - with probability $1 - p_S$, the job is dispatched to the fast server with shortest queue among the d_F queried fast servers



JSQ – (d_F, d_S)

3. If all $d_F + d_S$ queried servers are busy,

- with probability p_F , the job is dispatched to the fast server with shortest queue among the d_F fast servers
- with probability $1 - p_F$, the job is dispatched to the slow server with shortest queue among the d_S slow servers



Theorem 1 - Stability

Theorem

Under both **JIQ** – (d_F, d_S) and **JSQ** – (d_F, d_S) , for any values of $d_F, d_S \geq 1$, there exist choices of p_F and p_S such that the system is stable for $\lambda < \mu_F q_F + \mu_S q_S$.

Calculating Mean Response Time

- **Mean Field Analysis:** asymptotic independence assumption
- **JIQ** – (d_F, d_S) :
 - Tagged server approach
 - General Service time distributions
- **JSQ** – (d_F, d_S) :
 - Solved differential equations to obtain the limiting probabilities of the servers
 - Exponential Service time distribution

Optimization Problem

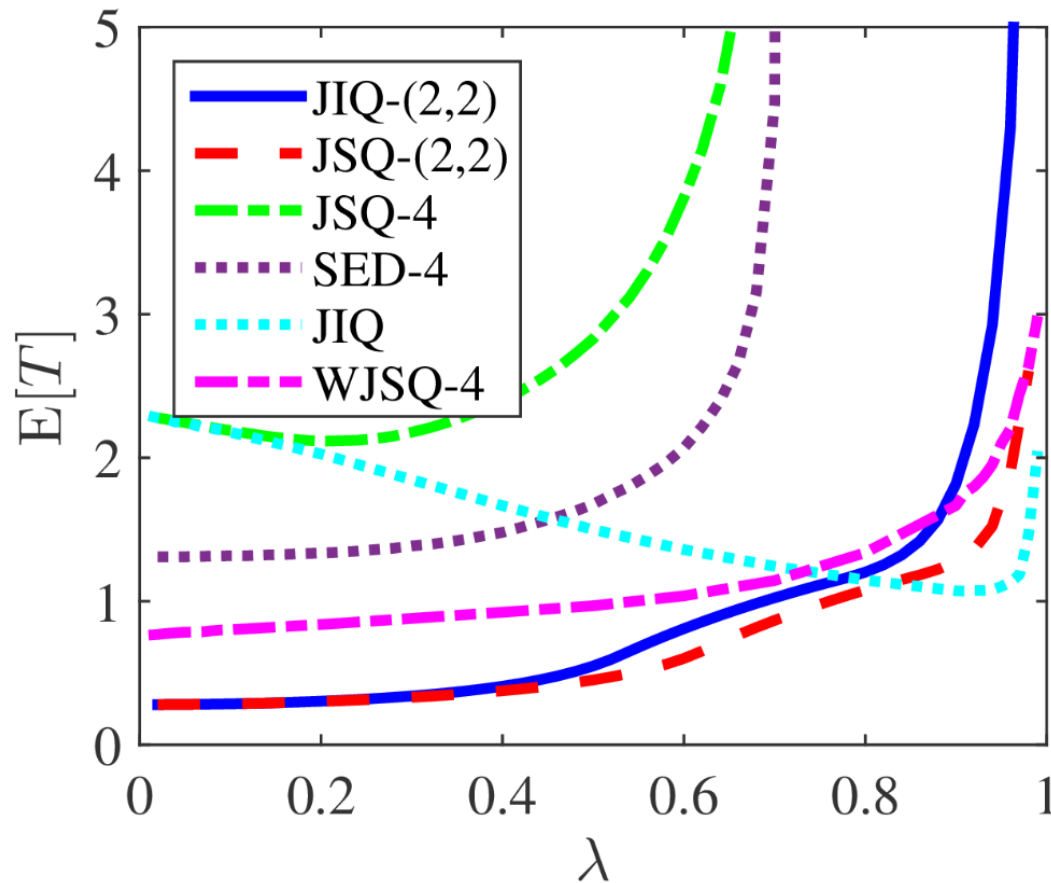
- Minimize Expected Response Time
- Problem Parameters $\lambda, \mu_F, \mu_S, q_F, q_S, d_F, d_S$
- Decision Variables p_F, p_S

Performance comparison

JIQ – (d_F, d_S) and **JSQ** – (d_F, d_S) against:

- **Join Shortest Queue** – (d)
- **Shortest Expected Delay** – (d)
- **Weighted Join Shortest Queue** – (d)
- **Join Idle Queue**

Performance comparison

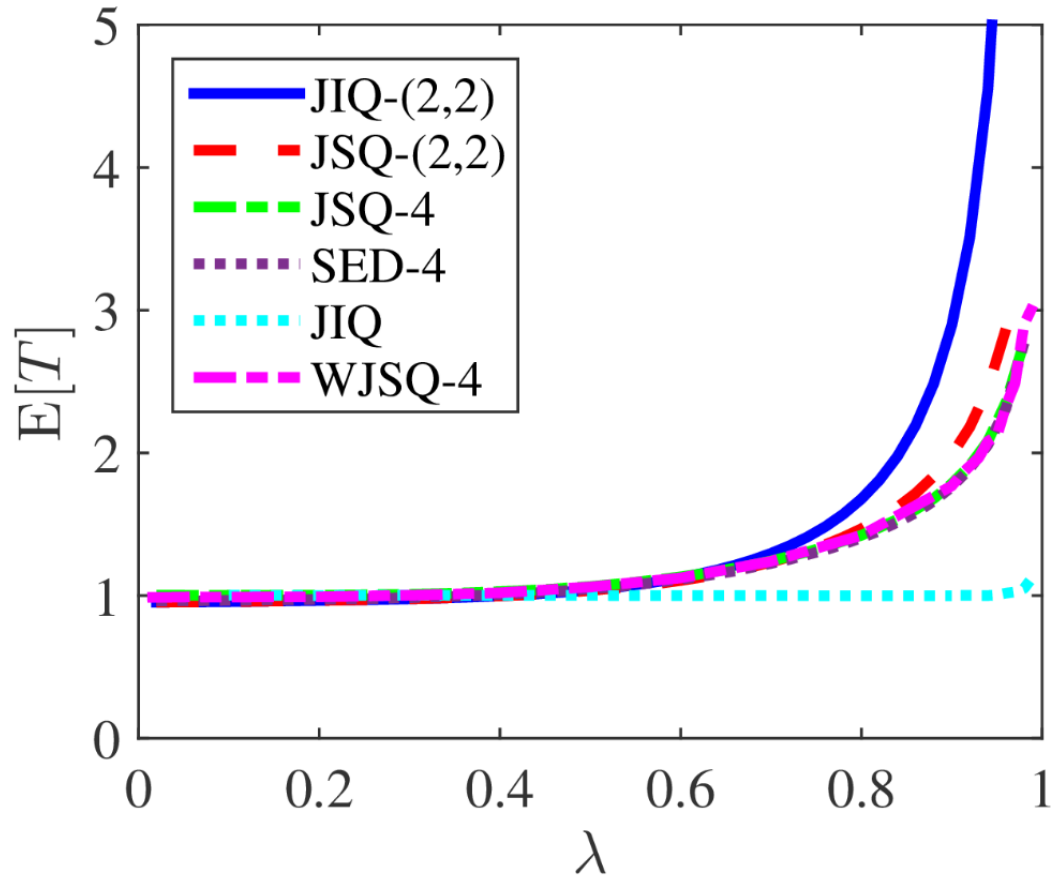


Percentage of fast servers= 20%

$$\frac{\mu_F}{\mu_S} = 10$$

- Apparent instability of **JSQ – (4)** and **SED – (4)**
- Stability of **JIQ – (2, 2)** and **JSQ – (2, 2)** by Theorem 1
- At low loads, **JIQ** does not efficiently utilize the fast servers

Performance comparison



Percentage of fast servers= 50%

$$\frac{\mu_F}{\mu_S} = 1.1$$

Main Takeaways

- Stability of **JIQ** – (d_F, d_S) and **JSQ** – (d_F, d_S) by Theorem 1
- The best dispatching policy will depend on the problem parameters.

Conclusion

- Introduced and analyzed new scalable power of d load balancing policies for heterogeneous systems
- Established the stability of the proposed policies
- Identified parameter settings where the policies outperform others
- Future Work: A more general framework
 - see [MAMA 2020](#)

Thank you