

HUG

Multi-Resource Fairness for Correlated and Elastic Demands

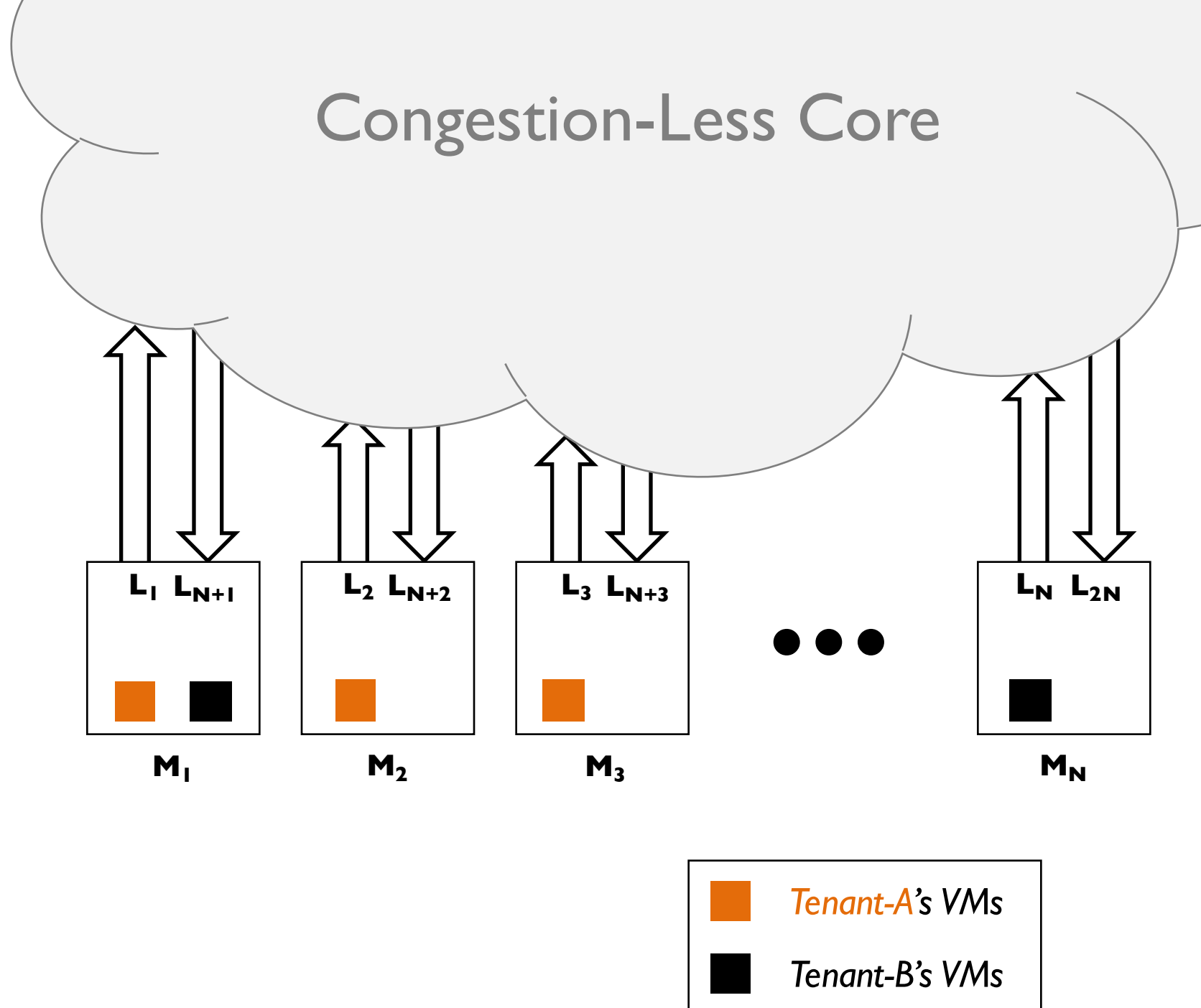
Mosharaf Chowdhury, Zhenhua Liu

Ali Ghodsi, Ion Stoica

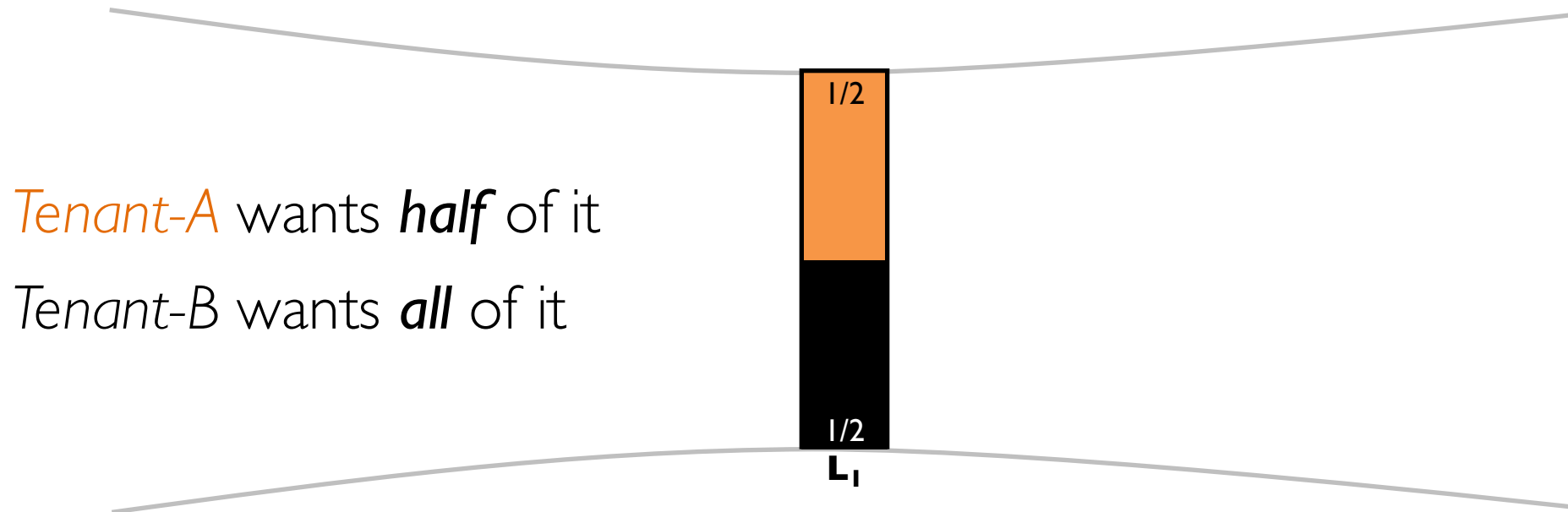


How to share the links between multiple tenants to provide

1. optimal performance guarantees and
2. maximize utilization?



Single-Resource Max-Min Fairness



1. Optimal Isolation Guarantee

Single-Resource Max-Min Fairness

Tenant-A wants *half* of it

Tenant-B wants *all* of it



Progress (M) of a tenant is its demand-normalized allocation

Isolation Guarantee is the minimum progress across all

$$\left. \begin{array}{l} \mathbf{d}_A = 1/2 \quad \mathbf{a}_A = 1/2 \quad \mathbf{M}_A = \frac{\mathbf{a}_A}{\mathbf{d}_A} = 1 \\ \mathbf{d}_B = 1 \quad \mathbf{a}_B = 1/2 \quad \mathbf{M}_B = \frac{\mathbf{a}_B}{\mathbf{d}_B} = 1/2 \end{array} \right\} \text{Min}(\mathbf{M}_A, \mathbf{M}_B) = 1/2$$

Single-Resource Max-Min Fairness

Tenant-A wants *half* of it

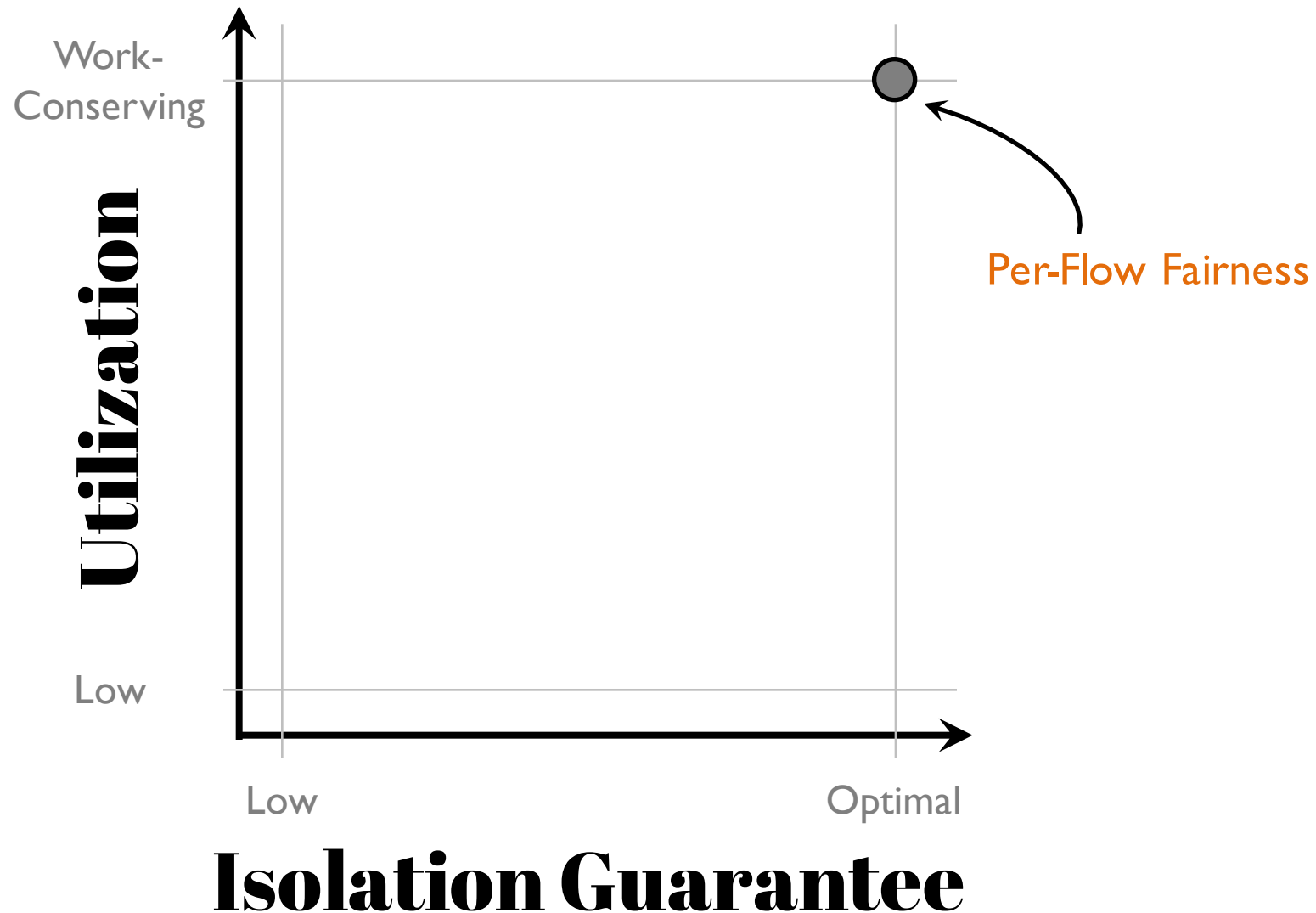
Tenant-B wants *all* of it



- 1. Optimal Isolation Guarantee**
- 2. Work Conservation**

No Tradeoff for Single Resource

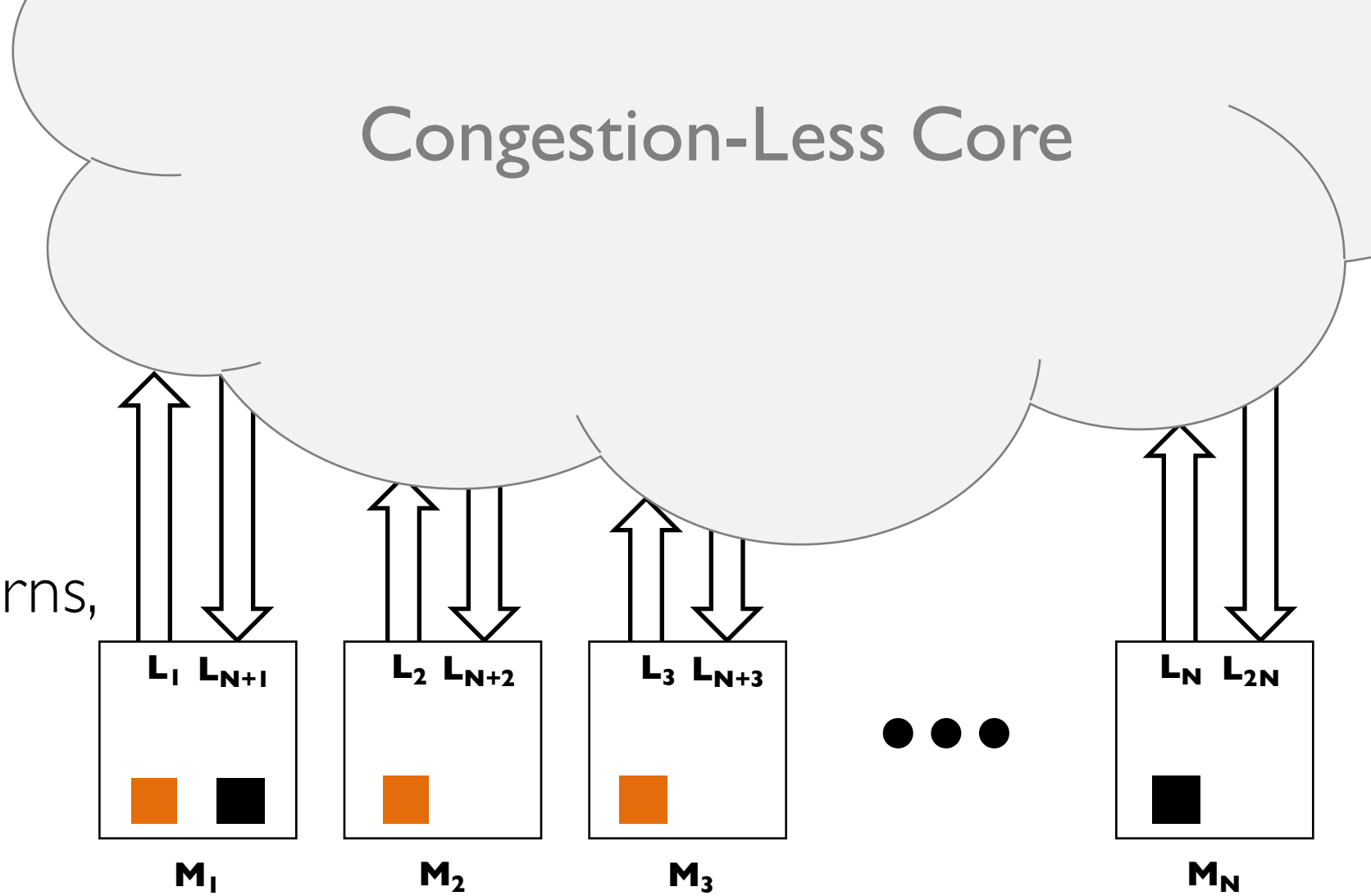
- 1. Optimal Isolation Guarantee**
- 2. Work Conservation**
- 3. Strategyproof**



Congestion-Less Core

Tenants have different

1. placements,
2. communication patterns,
3. demand correlations,
4.

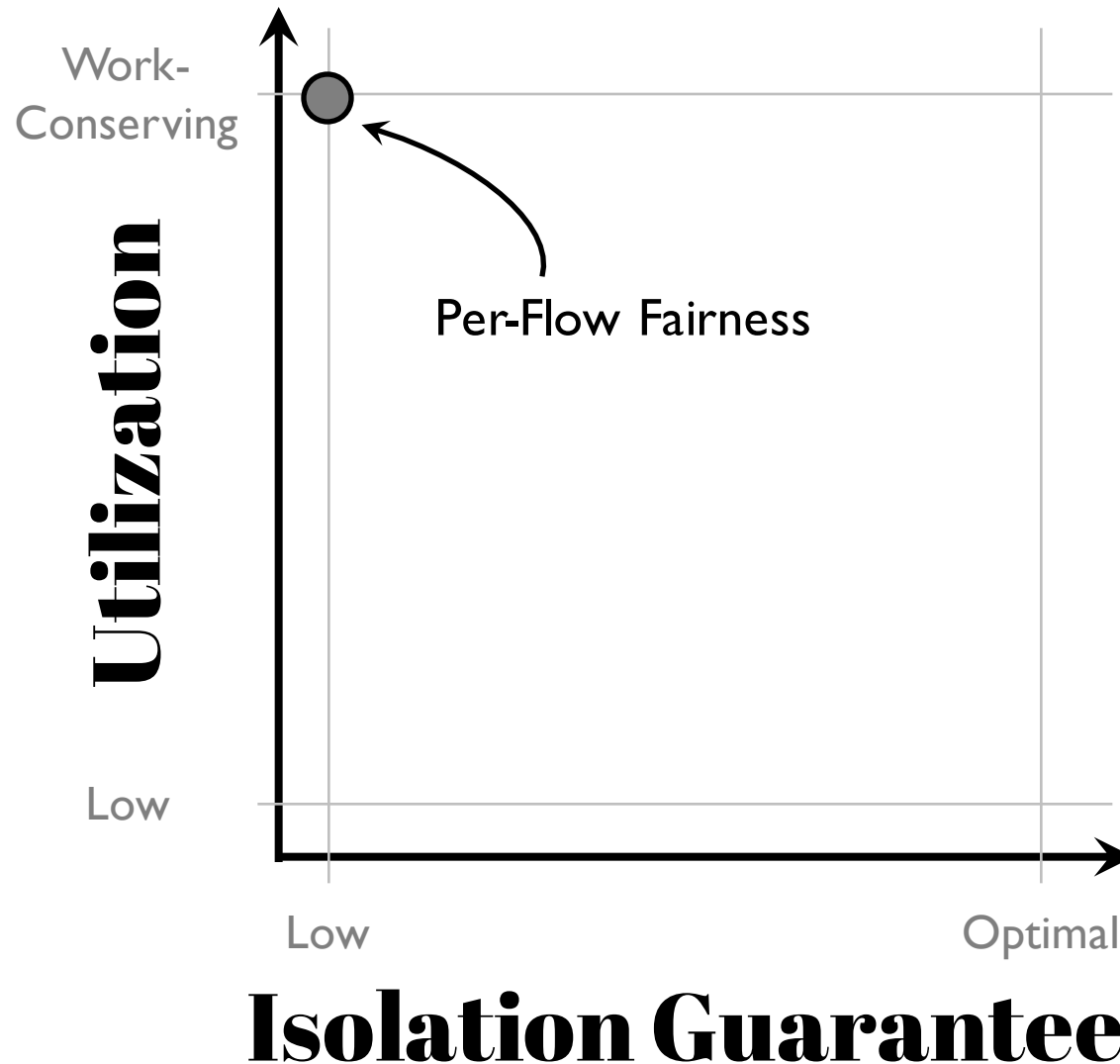


 Tenant-A's VMs
 Tenant-B's VMs

Per-Flow Fairness For Multiple Resources

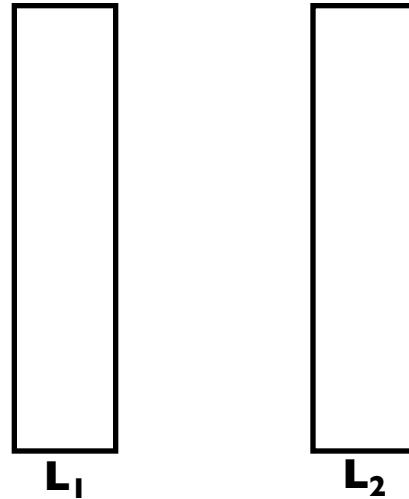
Low

- ~~1. Optimal Isolation Guarantee~~
2. Work Conservation
- ~~3. Strategyproof~~



Elastic Demands¹

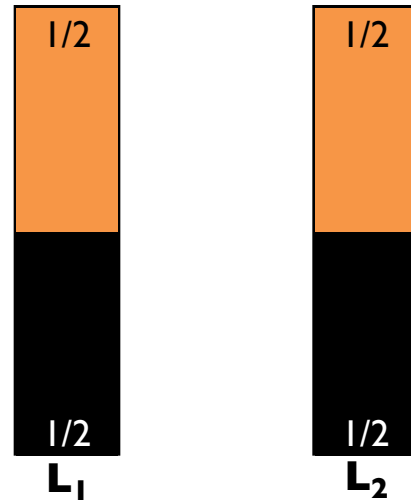
Tenant-A wants **all** of L_1
and **all** of L_2



Tenant-B wants **all** of L_1
and **all** of L_2

Tenant-Level Max-Min Fairness (PS-P)

Tenant-A wants **all** of L_1
and **all** of L_2

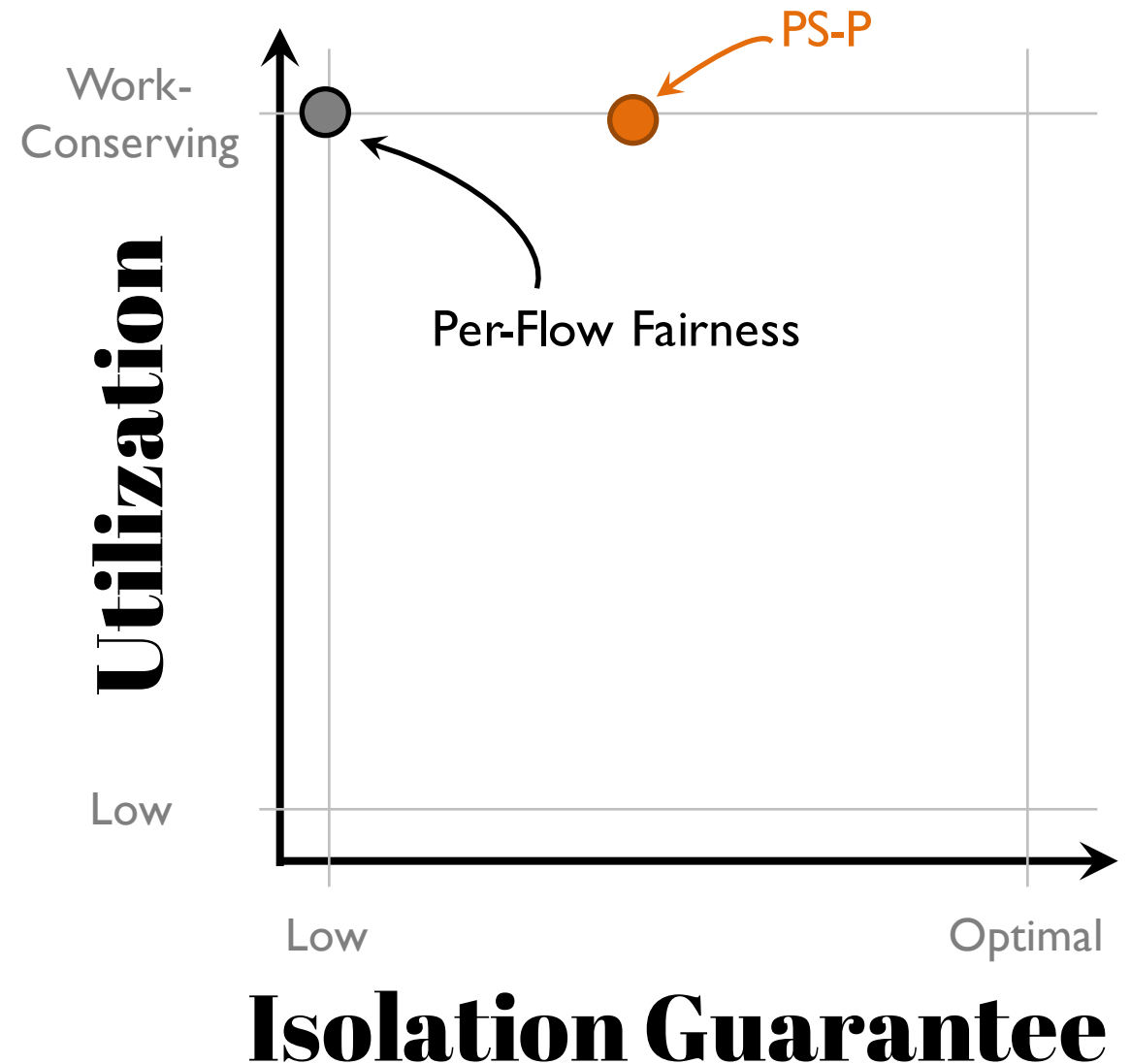


Tenant-B wants **all** of L_1
and **all** of L_2

$$\begin{array}{l} \mathbf{d}_A = \langle 1, 1 \rangle \quad \mathbf{a}_A = \langle 1/2, 1/2 \rangle \quad \mathbf{M}_A = \min\left(\frac{\mathbf{a}_A^i}{\mathbf{d}_A^i}\right) = 1/2 \\ \mathbf{d}_B = \langle 1, 1 \rangle \quad \mathbf{a}_B = \langle 1/2, 1/2 \rangle \quad \mathbf{M}_B = \min\left(\frac{\mathbf{a}_B^i}{\mathbf{d}_B^i}\right) = 1/2 \end{array} \quad \left. \vphantom{\begin{array}{l} \mathbf{d}_A = \langle 1, 1 \rangle \\ \mathbf{d}_B = \langle 1, 1 \rangle \end{array}} \right\} \text{Min}(\mathbf{M}_A, \mathbf{M}_B) = 1/2$$

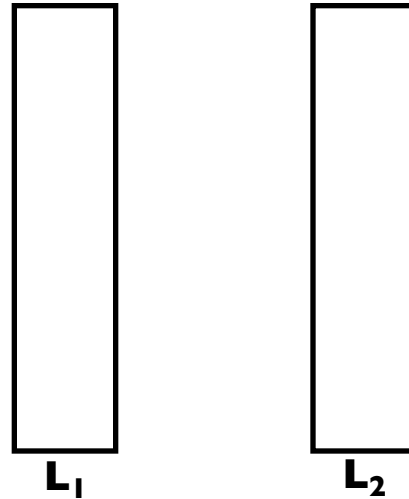
Tenant-Level Max-Min Fairness (PS-P)

1. **Suboptimal Isolation Guarantee**
2. **Work Conservation**



Correlated Demands¹

Tenant-A wants **some** of L_1 and **all** of L_2



Tenant-B wants **some** of L_2 and **all** of L_1

Dominant Resource Fairness (DRF)

Tenant-A wants *exactly half* unit of L_1 for *each* of L_2

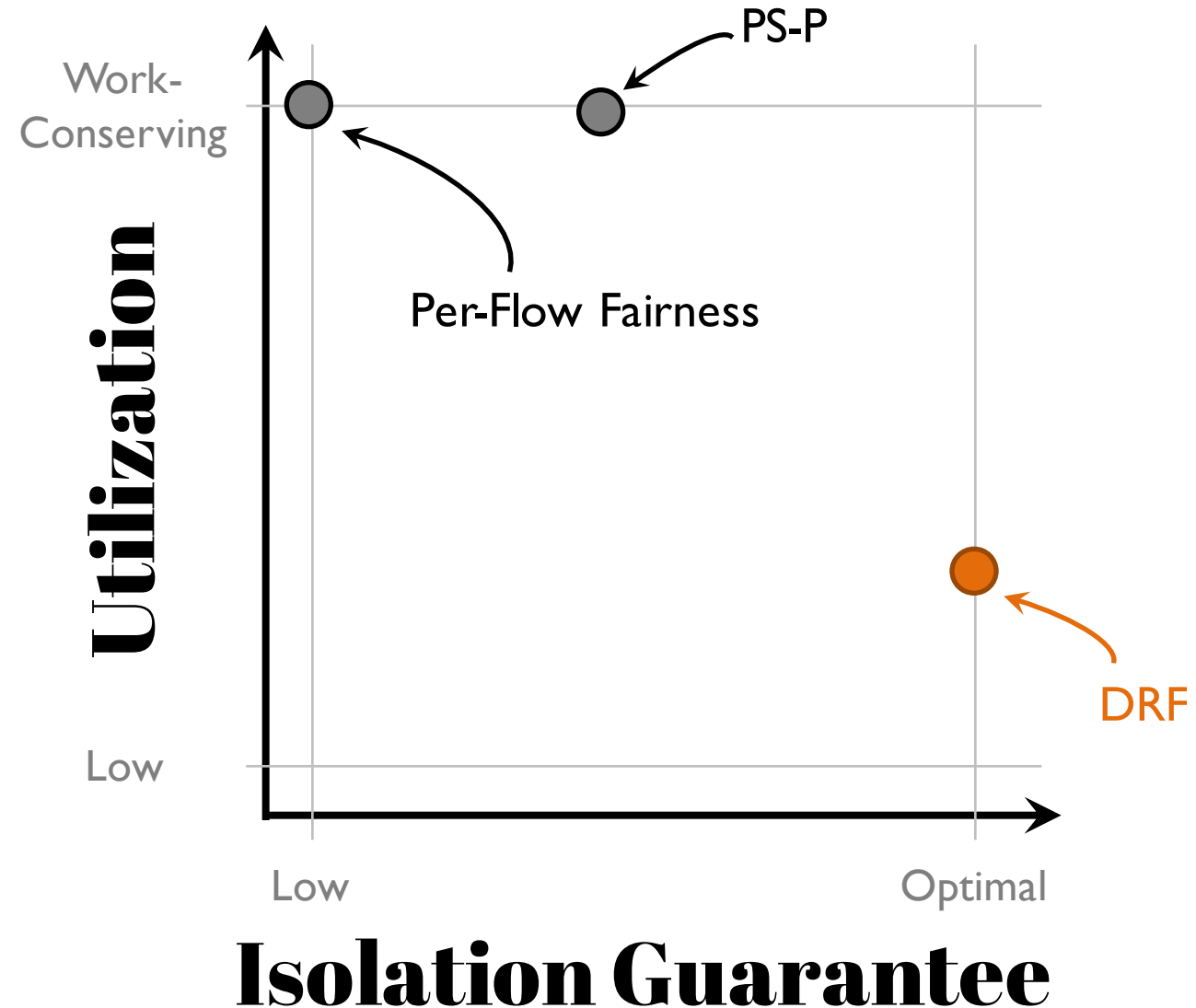


Tenant-B wants *exactly 1/6* unit of L_2 for *each* of L_1

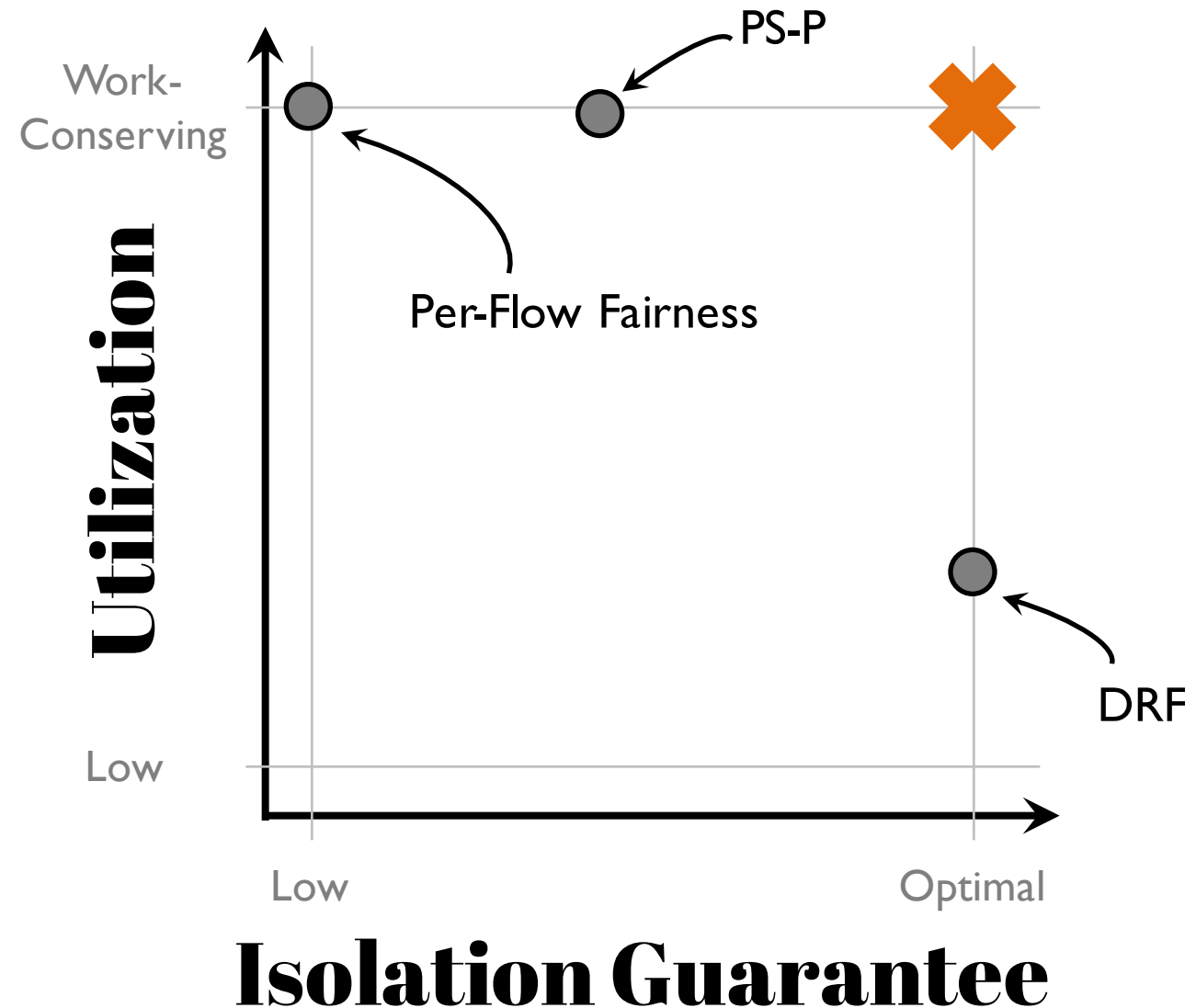
$$\begin{aligned}
 \mathbf{d}_A &= \langle 1/2, 1 \rangle & \mathbf{a}_A &= \langle 1/3, 2/3 \rangle & \mathbf{M}_A &= \min\left(\frac{a_A^i}{d_A^i}\right) = 2/3 \\
 \mathbf{d}_B &= \langle 1, 1/6 \rangle & \mathbf{a}_B &= \langle 2/3, 1/9 \rangle & \mathbf{M}_B &= \min\left(\frac{a_B^i}{d_B^i}\right) = 2/3
 \end{aligned}
 \left. \vphantom{\begin{aligned} \mathbf{d}_A &= \langle 1/2, 1 \rangle \\ \mathbf{d}_B &= \langle 1, 1/6 \rangle \end{aligned}} \right\} \mathbf{Min}(\mathbf{M}_A, \mathbf{M}_B) = 2/3$$

Dominant Resource Fairness (DRF)

- 1. Optimal Isolation Guarantee**
- 2. Arbitrarily Low Utilization**
- 3. Strategyproof**



For elastic and correlated demands, can we *simultaneously* achieve **optimal isolation guarantee** and **maximum utilization**?



For elastic and correlated demands, can we simultaneously achieve
optimal isolation
guarantee *and*
maximum utilization?

NO

1. Why not?

**2. What's the best we
can achieve?**

**3. How can we achieve
that?**

4. Does it matter?

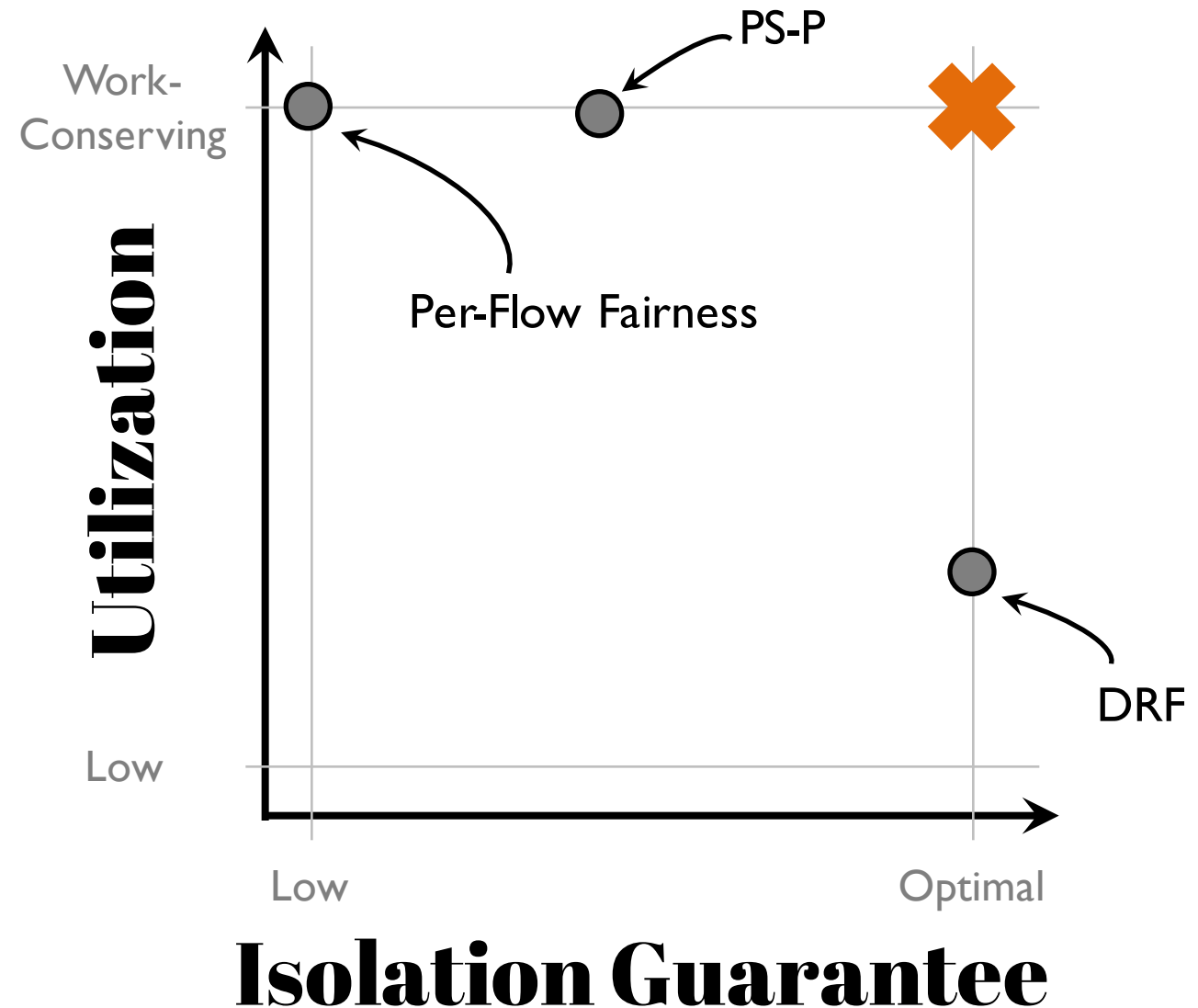
NO

1. **Why not?**

2. **What's the best we can achieve?**

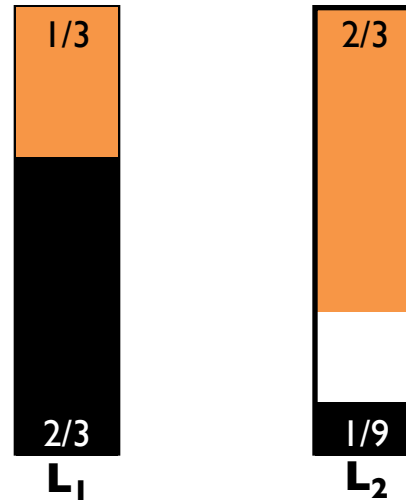
3. **How can we achieve that?**

4. **Does it matter?**



Elastic and Correlated Demands

Tenant-A wants *at least half* unit of L_1 for *each* of L_2

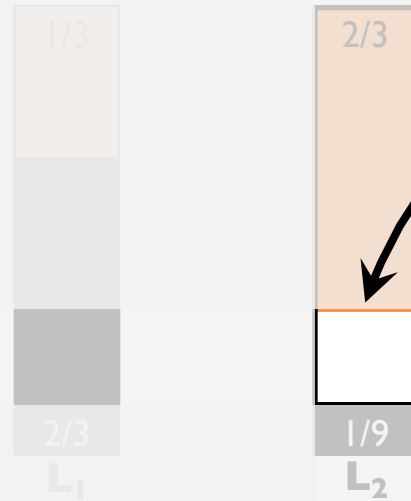


Tenant-B wants *at least 1/6* unit of L_2 for *each* of L_1

$$\left. \begin{array}{l}
 \mathbf{d}_A = \langle 1/2, 1 \rangle \quad \mathbf{a}_A = \langle 1/3, 2/3 \rangle \quad \mathbf{M}_A = 2/3 \\
 \mathbf{d}_B = \langle 1, 1/6 \rangle \quad \mathbf{a}_B = \langle 2/3, 1/9 \rangle \quad \mathbf{M}_B = 2/3
 \end{array} \right\} \text{Min}(\mathbf{M}_A, \mathbf{M}_B) = 2/3$$

Elastic and Correlated Demands

Tenant-A wants *at least half* unit of L_1 for *each* of L_2



Who gets this?

Tenant-B wants *at least 1/6* unit of L_2 for *each* of L_1

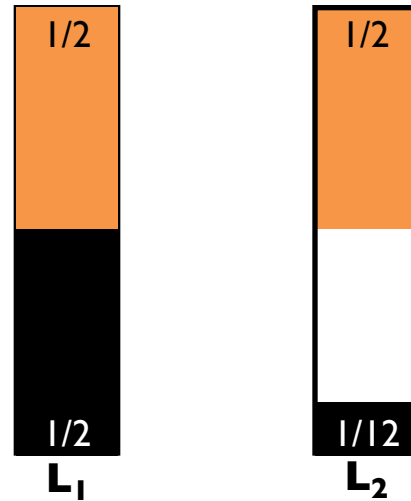
$$\mathbf{d}_A = \langle 1/2, 1 \rangle \quad \mathbf{a}_A = \langle 1/3, 2/3 \rangle \quad \mathbf{M}_A = 2/3$$

$$\mathbf{d}_B = \langle 1, 1/6 \rangle \quad \mathbf{a}_B = \langle 2/3, 1/9 \rangle \quad \mathbf{M}_B = 2/3$$

$$\left. \begin{array}{l} \mathbf{d}_A = \langle 1/2, 1 \rangle \quad \mathbf{a}_A = \langle 1/3, 2/3 \rangle \quad \mathbf{M}_A = 2/3 \\ \mathbf{d}_B = \langle 1, 1/6 \rangle \quad \mathbf{a}_B = \langle 2/3, 1/9 \rangle \quad \mathbf{M}_B = 2/3 \end{array} \right\} \text{Min}(\mathbf{M}_A, \mathbf{M}_B) = 2/3$$

Work Conservation Doesn't Work!

Tenant-A **lies** and asks for **one** unit of L_1 for **each** of L_2



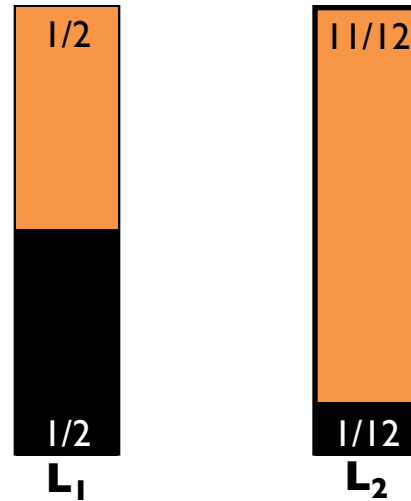
Tenant-B wants **at least** **1/6** unit of L_2 for **each** of L_1

$$\mathbf{d}'_A = \langle 1, 1 \rangle \quad \mathbf{a}'_A = \langle 1/2, 1/2 \rangle$$

$$\mathbf{d}_B = \langle 1, 1/6 \rangle \quad \mathbf{a}_B = \langle 1/2, 1/12 \rangle$$

Work Conservation Doesn't Work!

Tenant-A **lies** and asks for **one** unit of L_1 for **each** of L_2



Tenant-B wants **at least** **1/6** unit of L_2 for **each** of L_1

$$\mathbf{d}'_A = \langle 1, 1 \rangle \quad \mathbf{a}'_A = \langle 1/2, 1/2 \rangle \quad \mathbf{a}'_A = \langle 1/2, 11/12 \rangle \quad \mathbf{M}'_A = 11/12$$

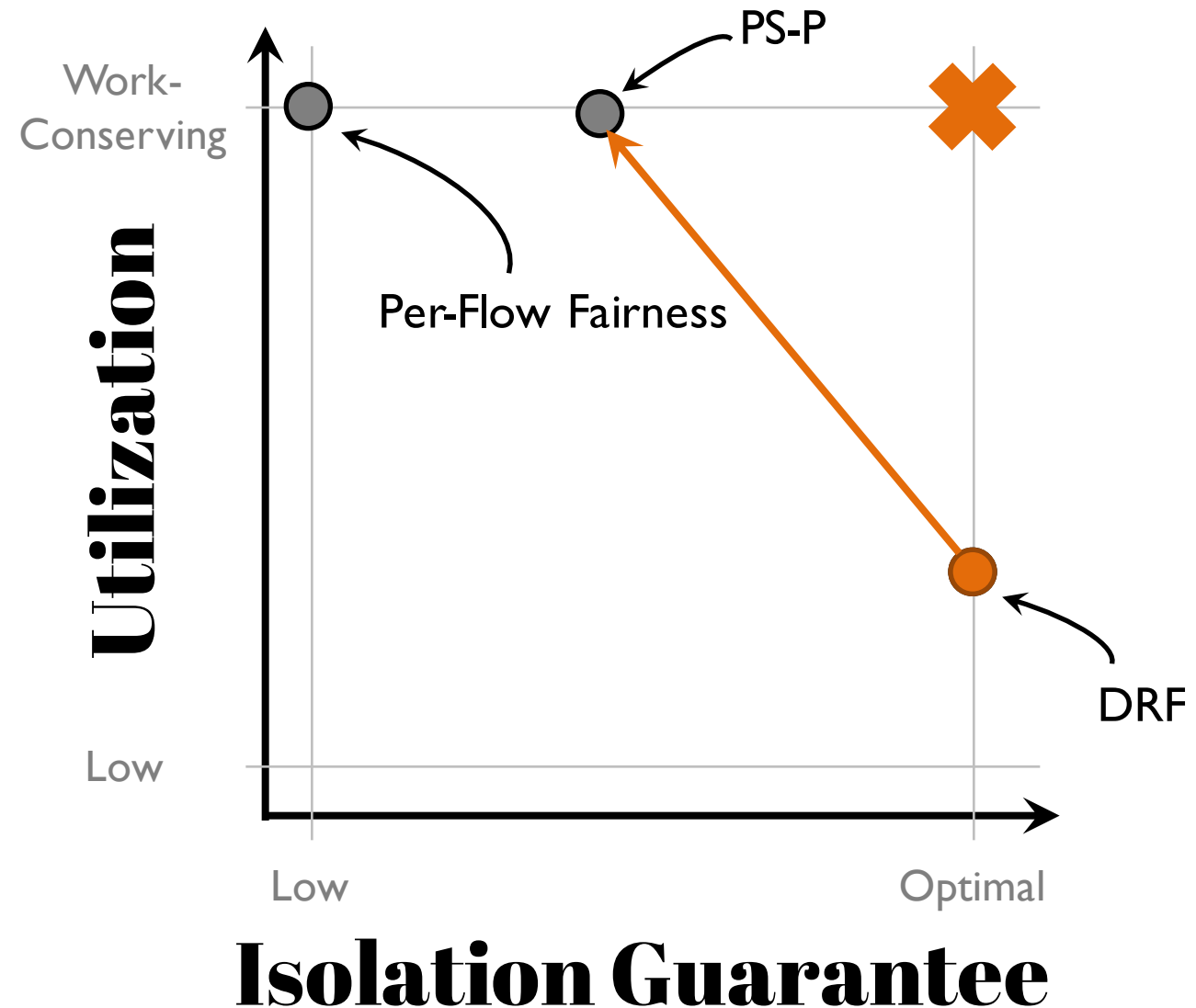
$$\mathbf{d}_B = \langle 1, 1/6 \rangle \quad \mathbf{a}_B = \langle 1/2, 1/12 \rangle \quad \mathbf{M}'_B = 1/2$$

Prisoner's Dilemma

		<i>Tenant-A</i>	
		Doesn't Lie	Lies
<i>Tenant-B</i>	Doesn't Lie	$\frac{2}{3}, \frac{2}{3}$ → $\frac{11}{12}, \frac{1}{2}$	
	Lies	$\frac{1}{2}, \frac{3}{4}$ → $\frac{1}{2}, \frac{1}{2}$	

1. Why not?

Optimal isolation guarantee depends on being strategyproof, but work conservation cannot coexist with strategyproof-ness

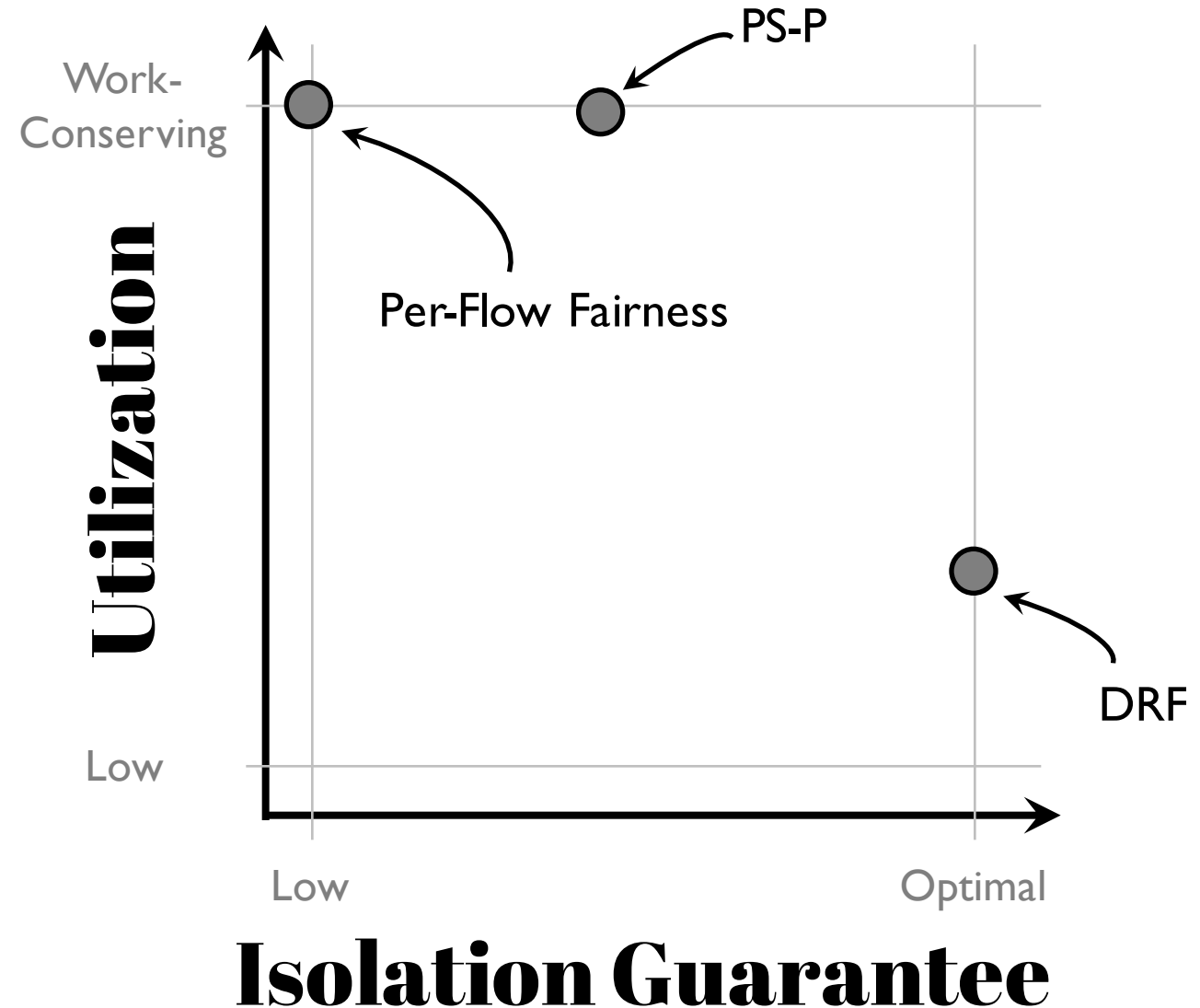


1. Why not?

2. What's the best we can achieve?

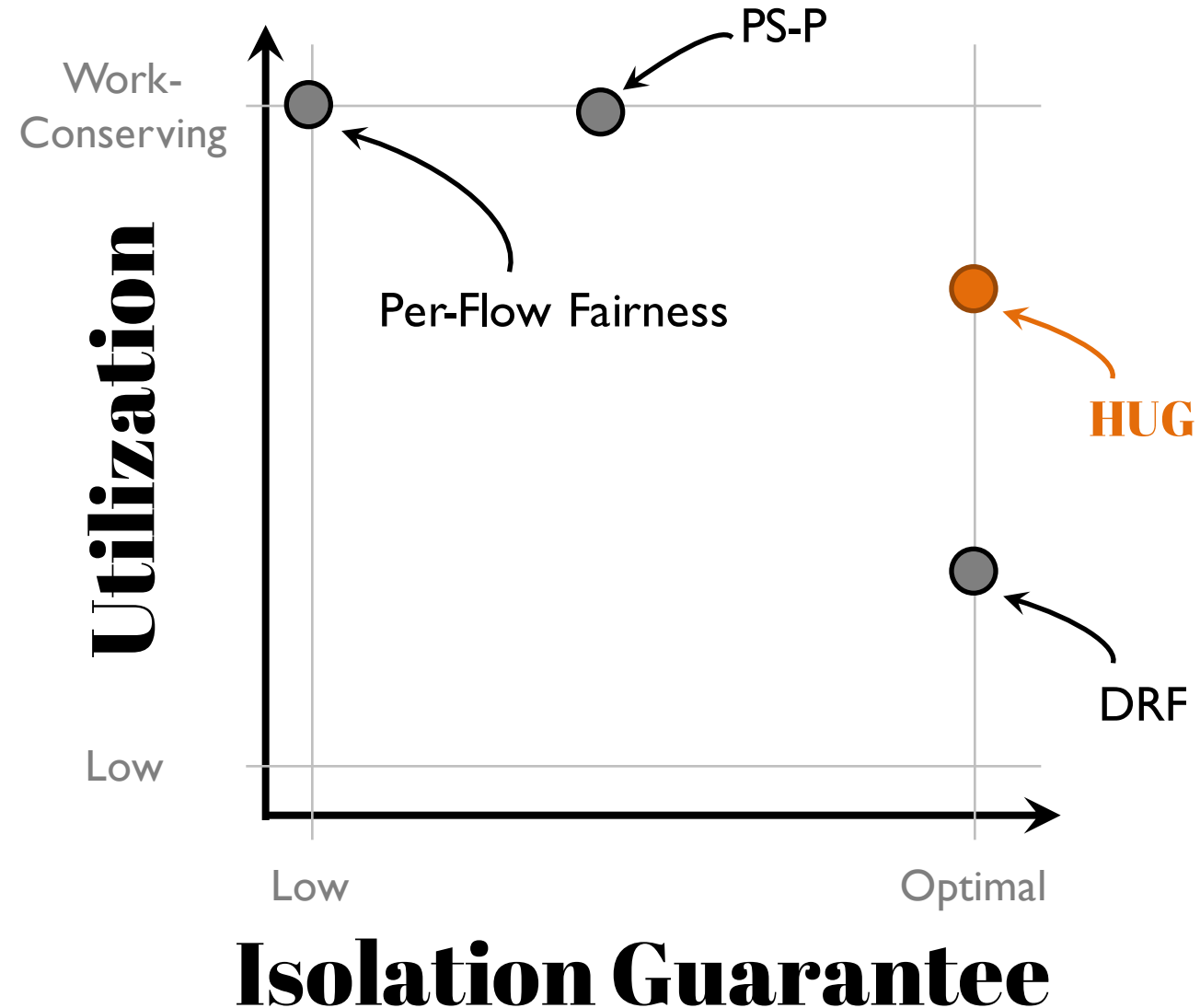
3. How can we achieve that?

4. Does it matter?



HUG in Non-Cooperative Setting

1. Optimal Isolation Guarantee
2. Highest Utilization
3. Strategyproof



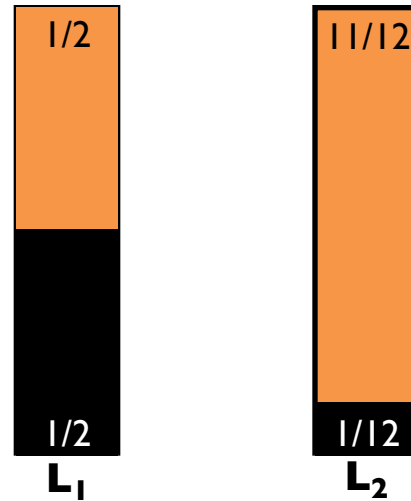
HUG

*Highest Utilization
with the
Optimal Isolation Guarantee*

Restrict a tenant's allocation in any link
to its allocation in the bottleneck link

Tenant-A Lied

Tenant-A **lies** and asks for **one** unit of L_1 for **each** of L_2



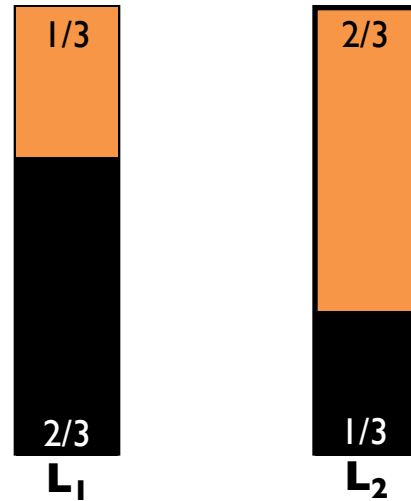
Tenant-B wants **at least** $1/6$ unit of L_2 for **each** of L_1

$$\mathbf{d}'_A = \langle 1, 1 \rangle \quad \mathbf{a}'_A = \langle 1/2, 1/2 \rangle \quad \mathbf{a}'_A = \langle 1/2, 11/12 \rangle \quad \mathbf{M}'_A = 1/2$$

$$\mathbf{d}_B = \langle 1, 1/6 \rangle \quad \mathbf{a}_B = \langle 1/2, 1/12 \rangle \quad \mathbf{M}'_B = 1/2$$

Everyone is Forced to Tell the Truth

Tenant-A wants *at least half* unit of L_1 for *each* of L_2



Tenant-B wants *at least 1/6* unit of L_2 for *each* of L_1

$$\mathbf{d}_A = \langle 1/2, 1 \rangle \quad \mathbf{a}_A = \langle 1/3, 2/3 \rangle$$

$$\mathbf{M}_A = 2/3$$

$$\mathbf{d}_B = \langle 1, 1/6 \rangle \quad \mathbf{a}_B = \langle 2/3, 1/9 \rangle \quad \mathbf{a}'_B = \langle 2/3, 1/3 \rangle$$

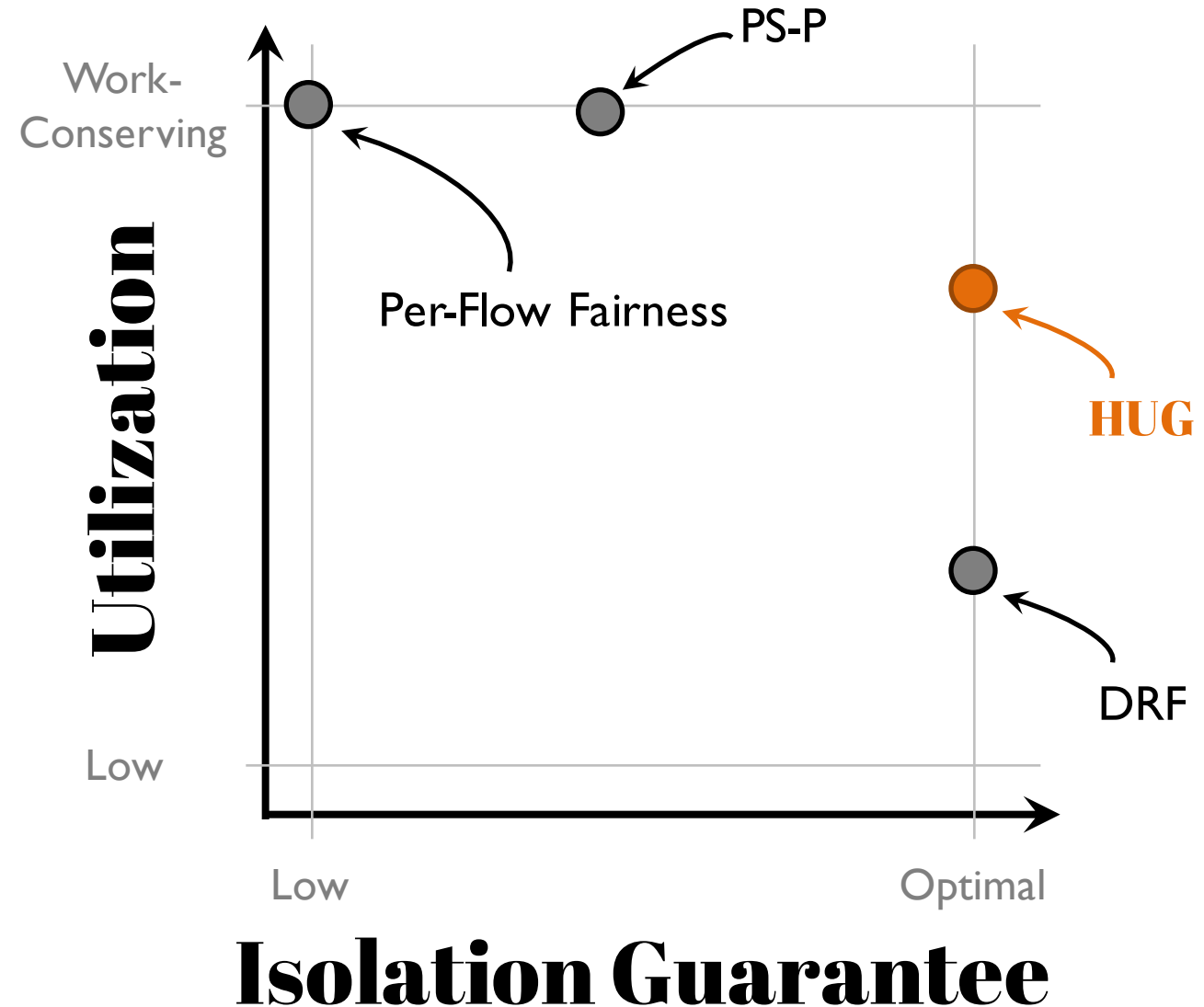
$$\mathbf{M}_B = 2/3$$

HUG

*Highest Utilization
with the
Optimal Isolation Guarantee*

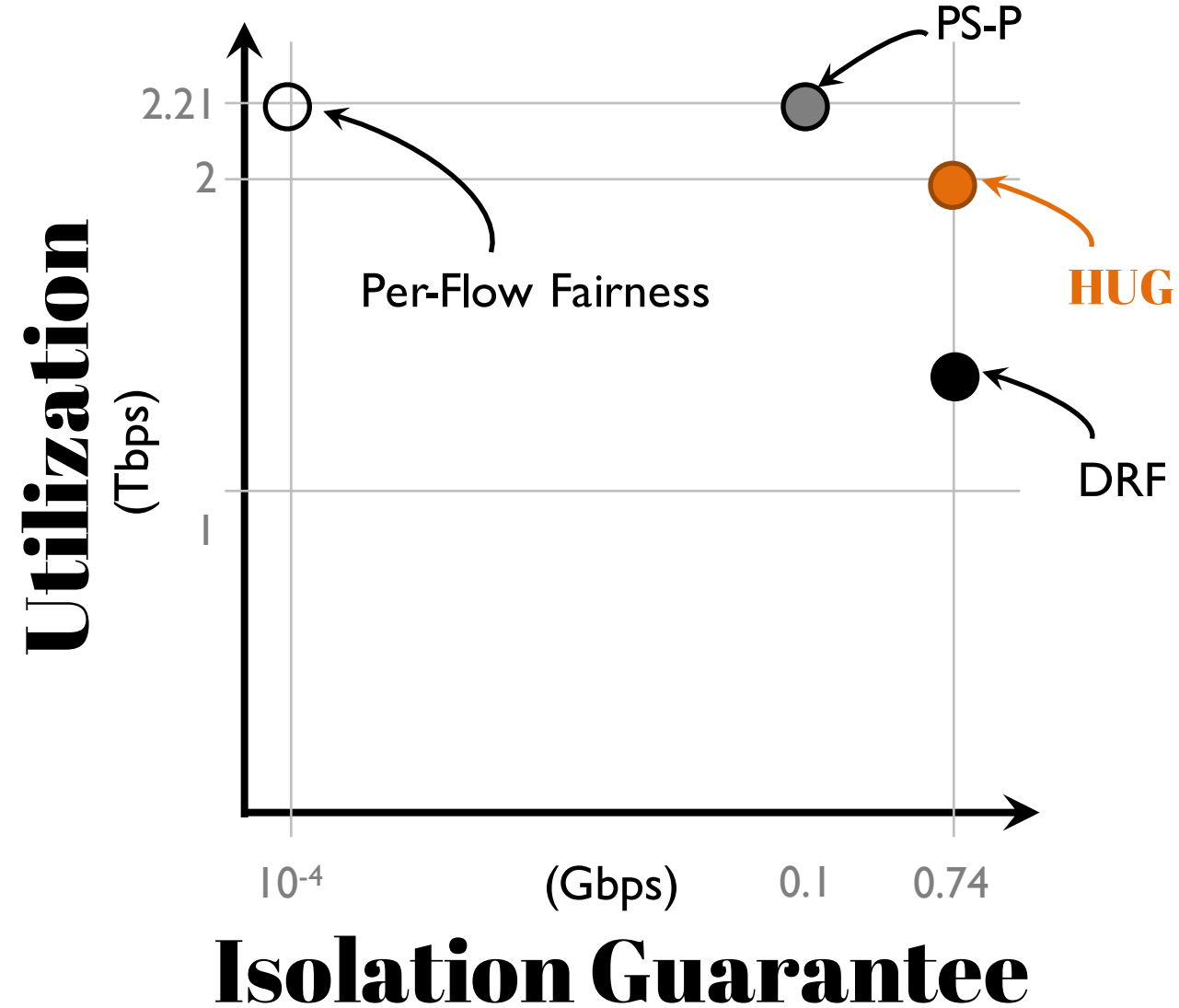
1. Tenants update correlation vectors through an API
2. Operators calculate HUG centrally and enforce it locally

1. Why not?
2. What's the best we can achieve?
3. How can we achieve that?
4. Does it matter?



Evaluation

- 100 concurrent tenants
- 3000 machines with 3Tbps total capacity
- Original placement and communication patterns from the Facebook trace



#1

Bursty Demands

Periodic demand bursts in
Spark streaming

#2

Long-Term Guarantees

Predictable performance
guarantees over time

#3

Decentralized Algorithms

Survive master failures and
enable low response times

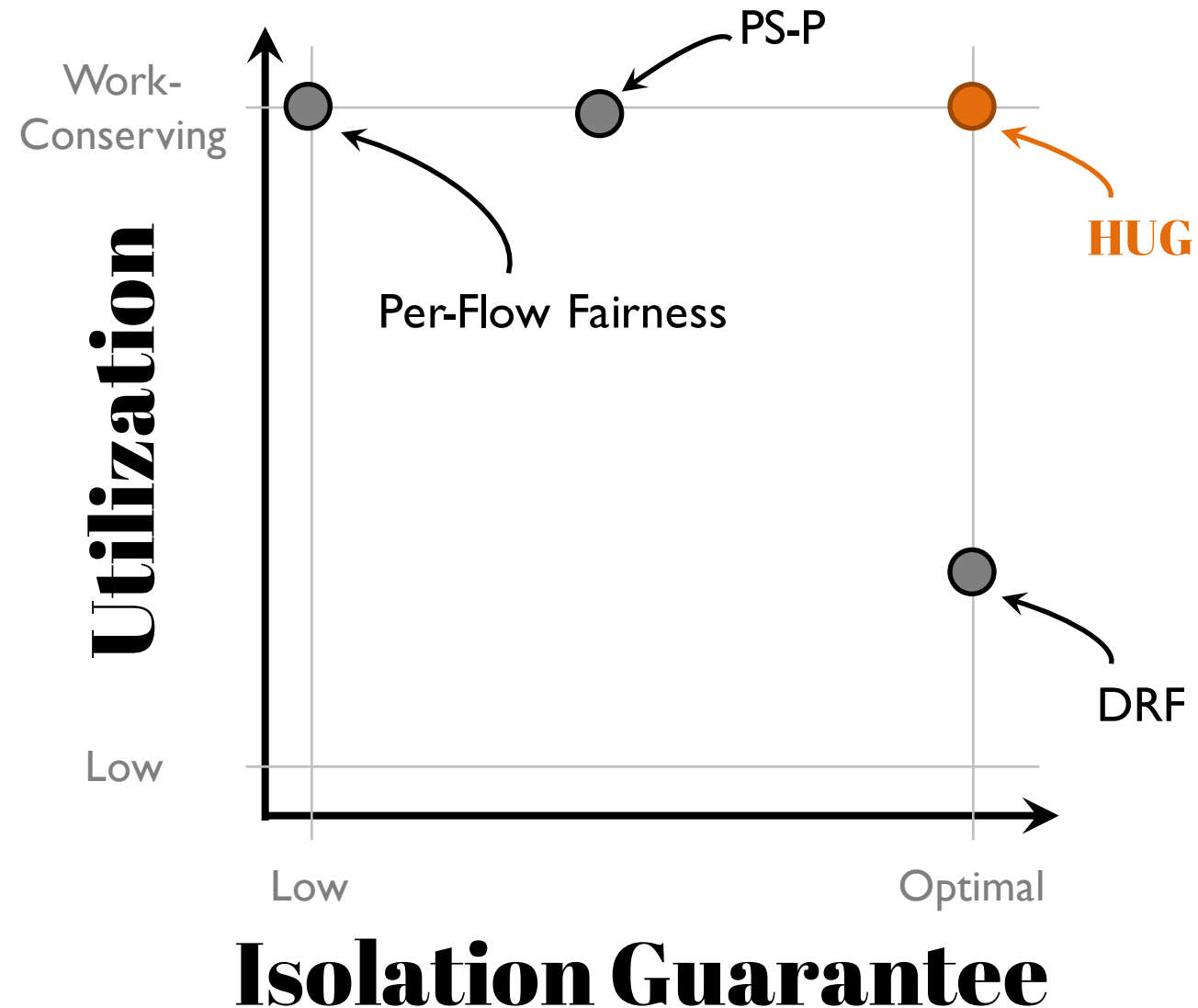
HUG

*Highest Utilization
with the
Optimal Isolation Guarantee*

- Generalizes single- and multi-resource fairness schemes
- Optimal worst-case performance guarantees for tenants
- Highest utilization for operators

HUG in Cooperative Setting

1. Optimal Isolation Guarantee
2. Work Conservation



Evaluation

A 3000-machine trace-driven simulation based on a snapshot of Facebook production trace

1. Does it provide isolation guarantee?
2. Does it improve utilization?
3. Is it practical?

YES

Optimal Progress for ALL

	Per-Flow Fairness	PS-P ²	DRF ³	HUG
Max	1	1	0.74_±	0.74_±
Min	0.0001	0.10	0.74_±	0.74_±
Max-to-Min Progress Ratio	10000X	10X	1X	1X

1. 100 tenants in this particular snapshot. The unit of progress is Gbps.

Higher Network Utilization

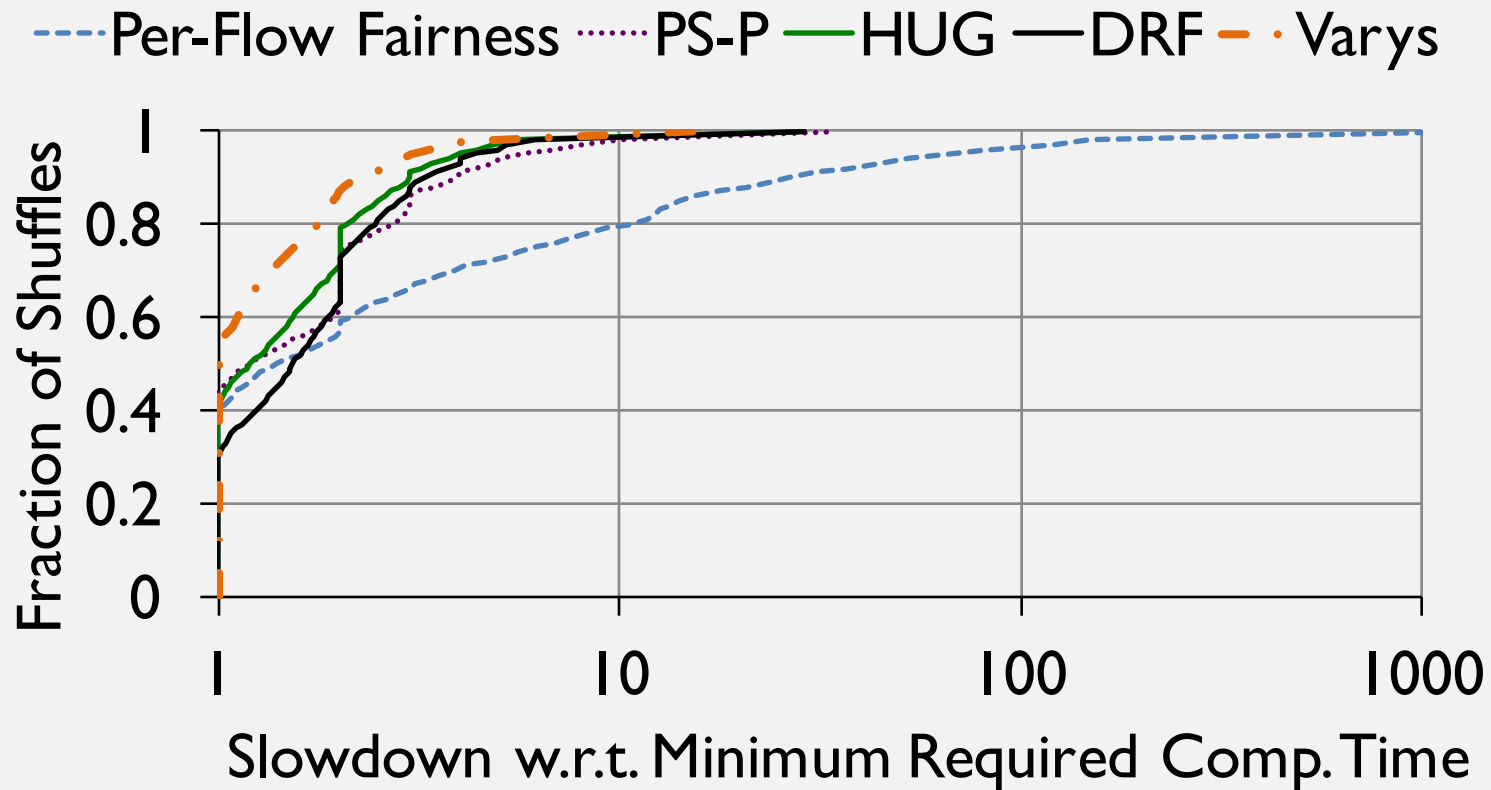
	Per-Flow Fairness	PS-P ²	DRF ³	HUG
Total Utilization (Tbps)	2.21	2.20	1.42	2.00
Max-to-Min Progress Ratio	10000X	10X	1X	1X

1. 100 tenants in this particular snapshot. The unit of progress is Gbps.

2. FairCloud: Sharing the Network in Cloud Computing, SIGCOMM'12

3. Dominant Resource Fairness: Fair Allocation of Multiple Resource Types, NSDI'11

Long-Term Performance



	Average Time
Per-flow Fairness	1.49X
PS-P	1.14X
DRF	1.14X
HUG	1X
Varys ¹	0.69X

Coordination Overheads and Scalability

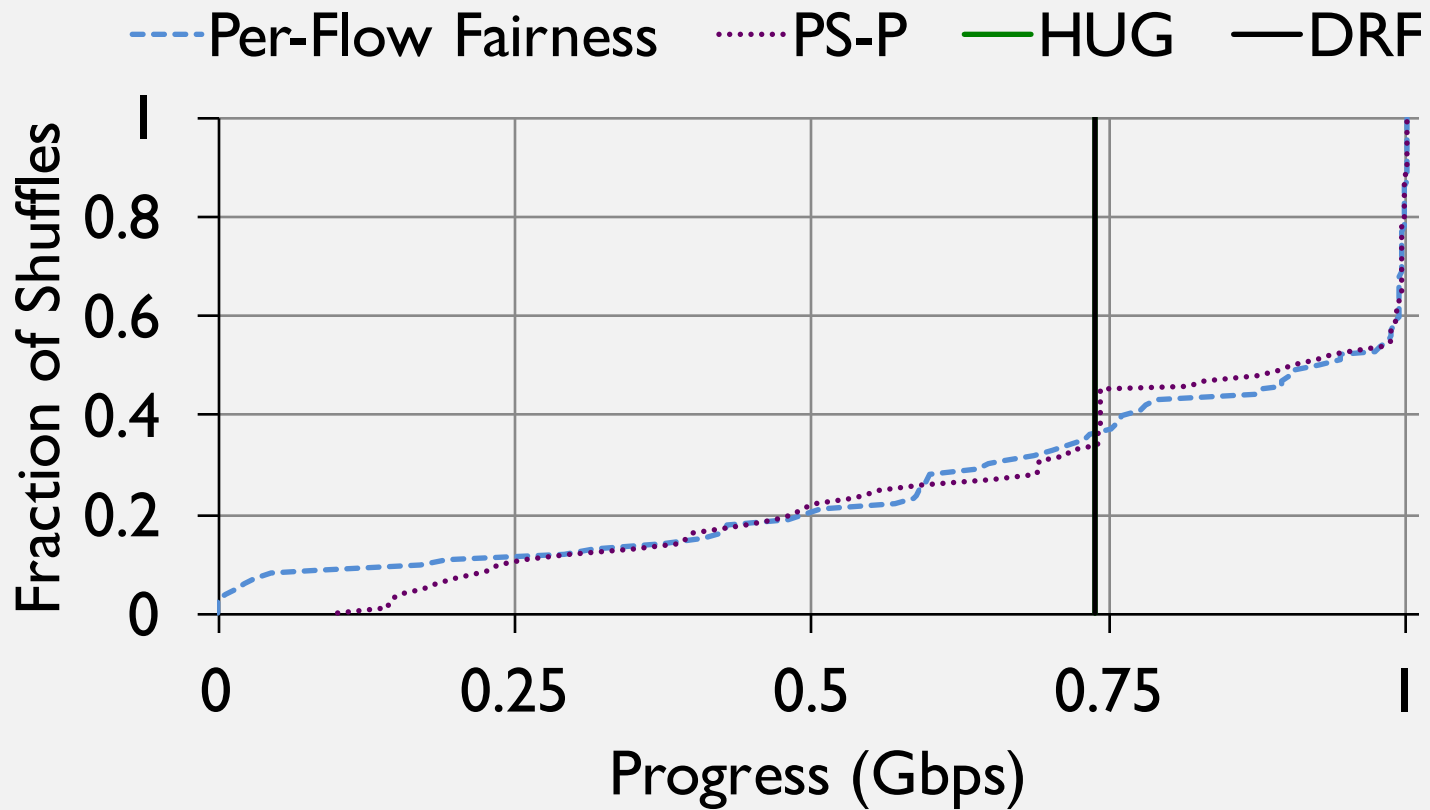
Computation overheads

- Less than 5 μ s for 100-machine cluster
- Less than 10 ms for 100,000 machines

Communication overheads

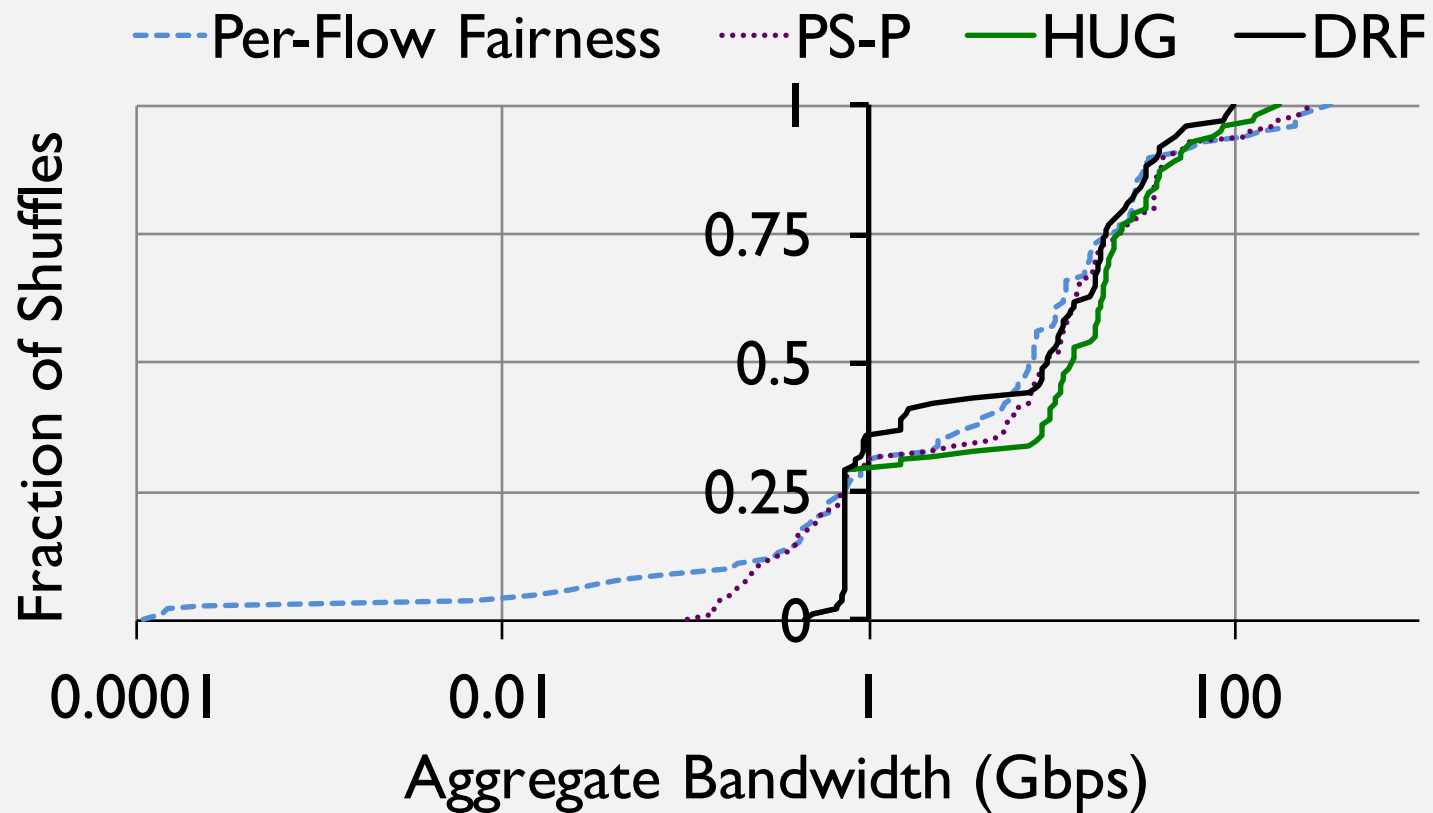
- Less than 10 ms for 100-machine cluster
- Less than 1 second for 100,000 machines

Optimal Progress for ALL



	Max/Min Ratio
Per-flow Fairness	10000X
PS-P	10X
DRF	1X
HUG	1X

Higher Utilization



	Max/Min Ratio
Per-flow Fairness	3240000X
PS-P	2590X
DRF	196X
HUG	340X