

A Distributed Algorithm to Compute Max-Min Fair Rates Without Per-Flow State

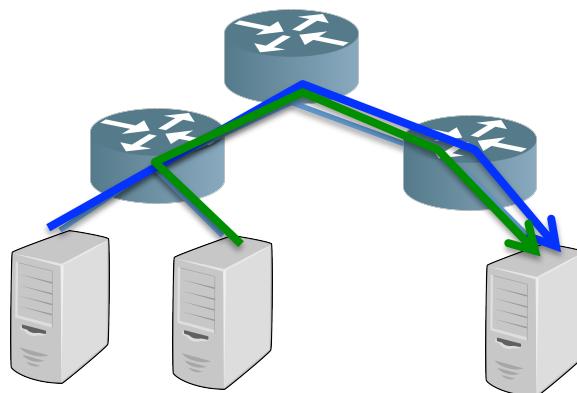
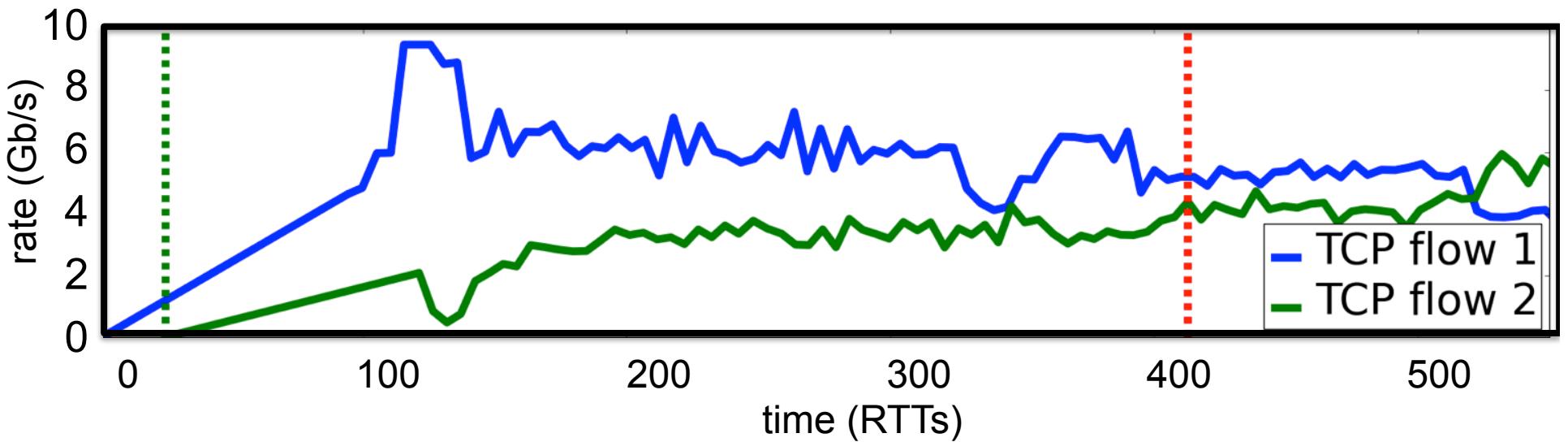
Lavanya Jose, Stephen Ibanez, Nick McKeown
Stanford University

Mohammad Alizadeh
MIT CSAIL

Imagine you're a cloud provider..

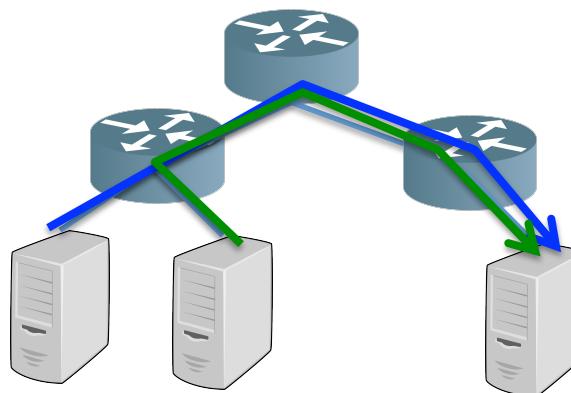
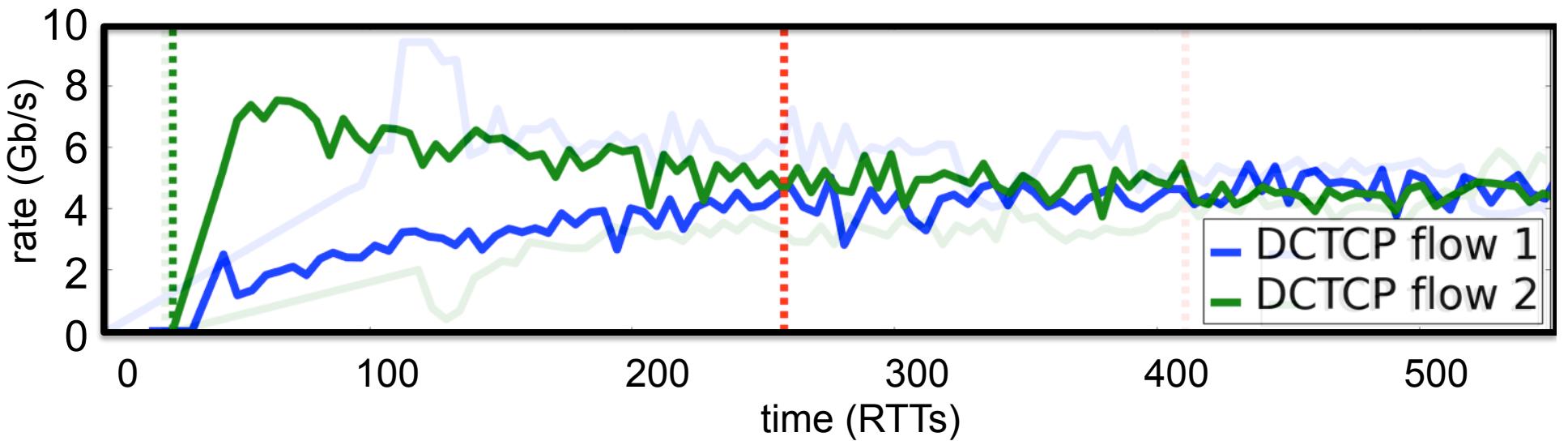
- You want
 - Fair + High Throughput Rates for Long Flows
 - Low Latency for Short Flows
- Many attempts to achieve this, but they all have the same *problem*—they *react* to congestion instead of preventing it.

TCP takes 400 round trips for 2 flows!



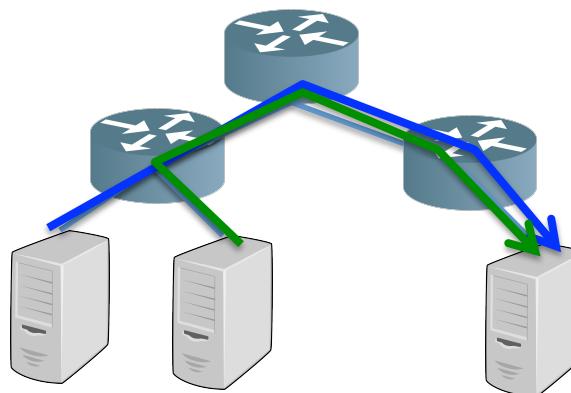
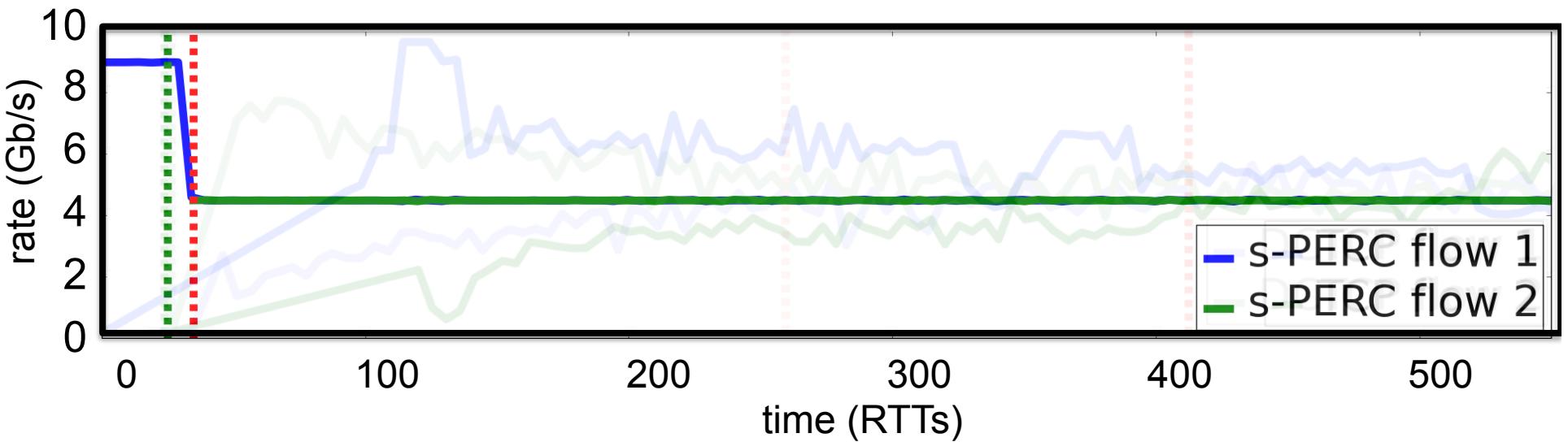
10 Gb/s links
0.2 ms RTT

DCTCP takes 250 round trips for 2 flows!



10 Gb/s links
0.2 ms RTT

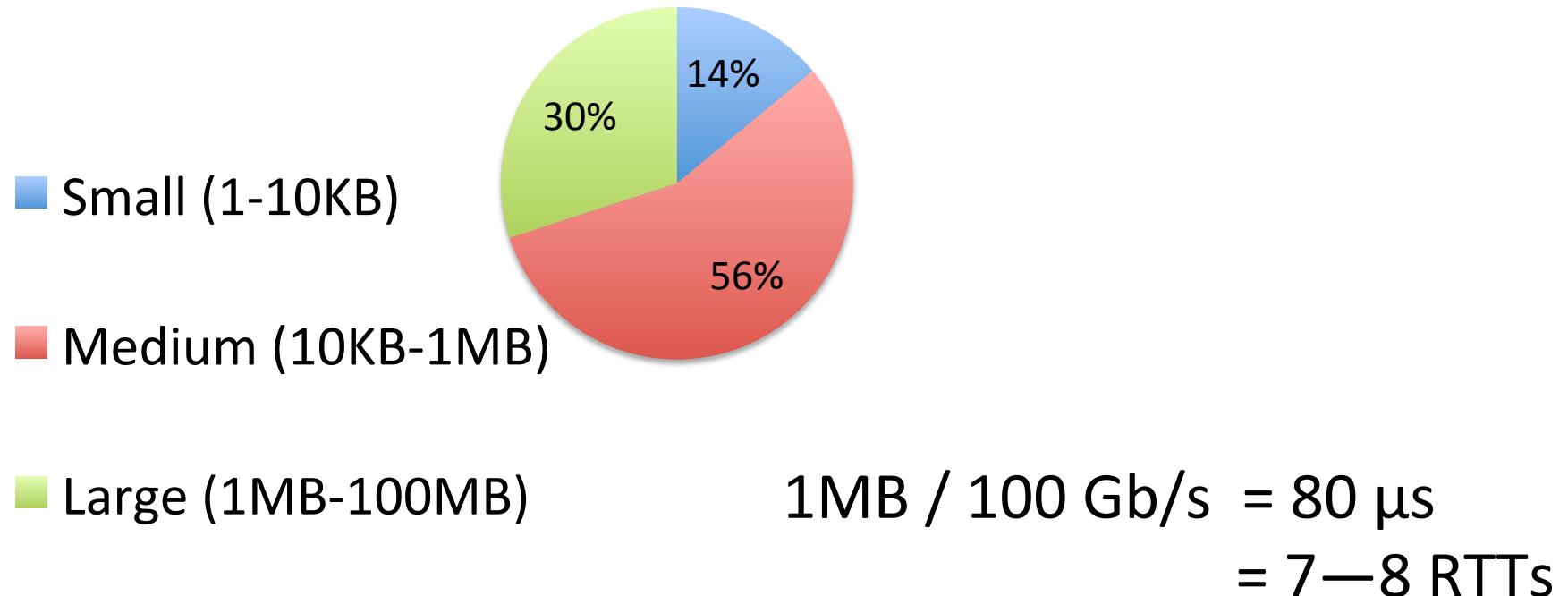
Why not ~1 round trip to find rates?



10 Gb/s links
0.2 ms RTT

Flows can finish much faster

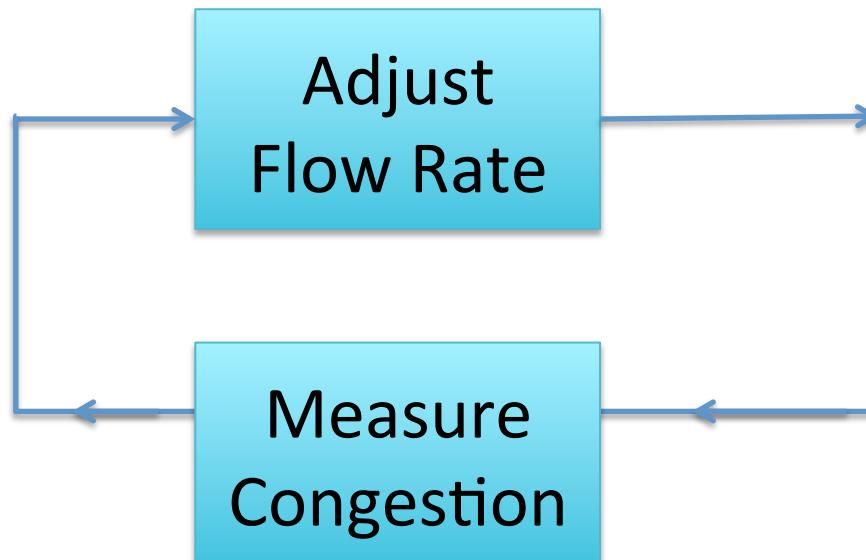
Fraction of Total Flows in Bing Workload



Most flows can finish long before
congestion control finds good rates.

Why “Reactive” Algorithms Are Slow

1. No info. about routing matrix, link speeds.
2. Measure congestion and adjust flow rates.
3. Timid steps to remain stable.



What if instead of reacting to congestion,
we prevent it in the first place?

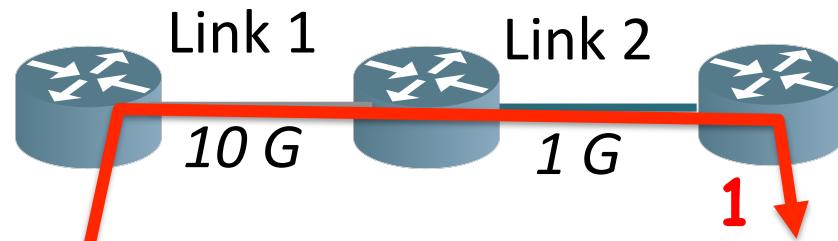
Proactive Algorithms

- Use knowledge of the routing matrix and link speeds to calculate flow rates directly.
- Centralized: can't do this at scale
- Distributed?

PERC Algorithms

Red Flow's Per-Link State	Link	
	1	2
alloc. (a)	10	1

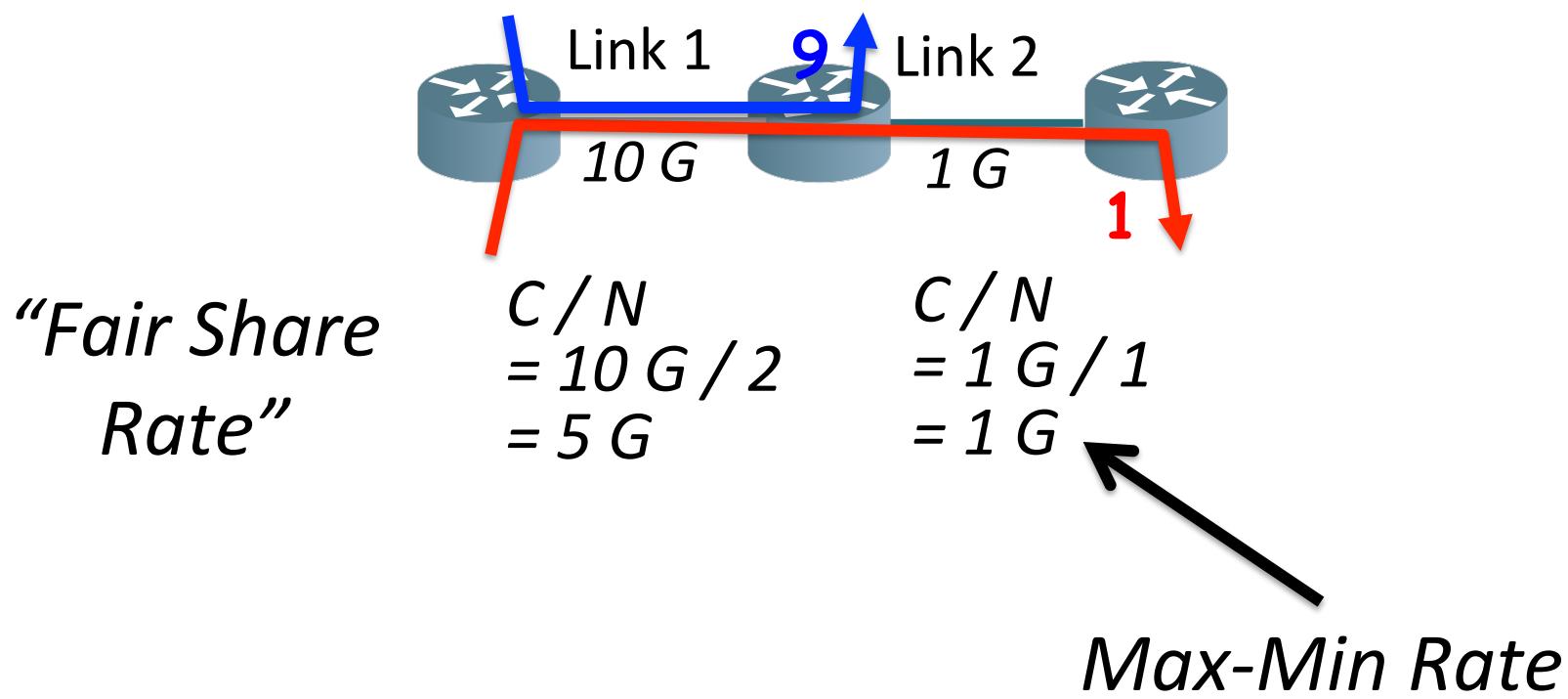
“Proactive Explicit Rate Control (PERC)”



PERC algorithms for the Cloud

- Practical algorithms should
 - update the control packet in a few nanoseconds
 - using very little switch memory, constant state.
- Relevant today: cloud providers own entire infrastructure

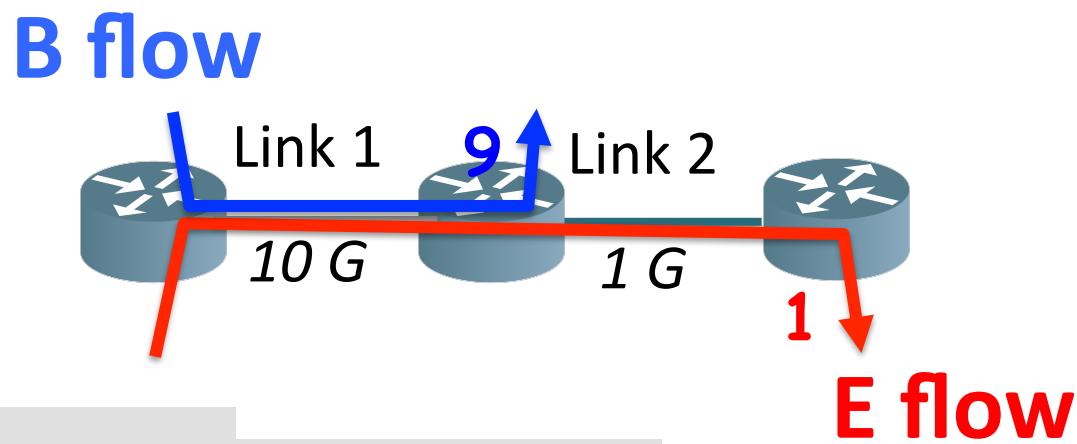
Max-Min Fairness



Outline

- Simple practical “naïve” PERC (n-PERC)
 - Can take arbitrarily long to converge
- Modify n-PERC using constant state to get “stateless” PERC (s-PERC)
 - Converges in known bounded time
 - Without per-flow state at switches

Two Kinds of Flows At a Link



$$\frac{(C - \text{SumE})}{\text{NumB}} = \text{Max-Min Rate}$$

of B flow

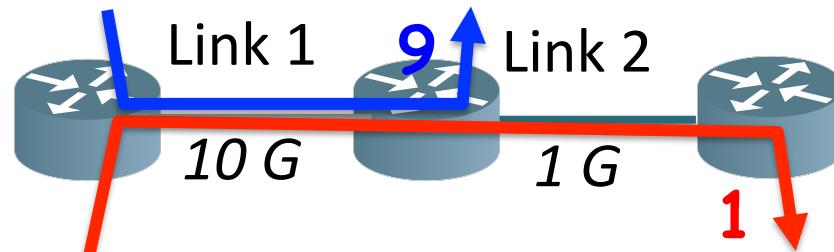
The n-PERC Algorithm

- Calculate *bottleneck* rate assuming flow will be bottlenecked here, $b = (C - \text{SumE})/\text{NumB}$
- Calculate limit rate, e , lowest bottleneck rate elsewhere
- Allocate lower of the two $a = \min(b, e)$

The n-PERC Algorithm

- Calculate bottleneck rate, b ... $b = (1 - 0)/1 = 1 \text{ Gb/s}$
- Calculate limit rate, e ... $e = 10 \text{ Gb/s}$
- Allocate lower of the two ... $a = \min(b, e) = 1 \text{ Gb/s}$

Red Flow's Per-Link State	Link	
	1	2
bn. rate (b)	10	∞

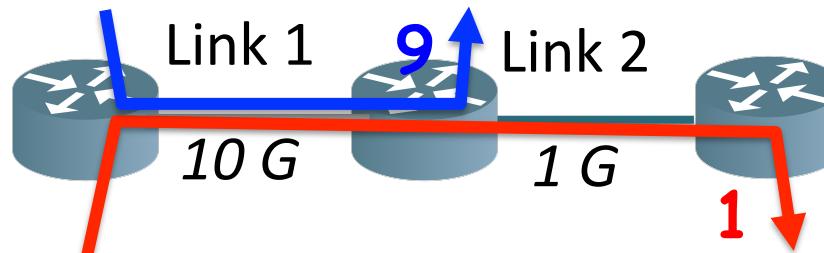


Link 2's State	
NumB	1
SumE	0

The n-PERC Algorithm

- Calculate bottleneck rate, b
 - Calculate limit rate, e ...
 - Allocate lower of the two ..
 - Update packet and link
- $b = (1 - 0)/1 = 1 \text{ Gb/s}$
 $e = 10 \text{ Gb/s}$
 $a = \min(b, e) = 1 \text{ Gb/s}$
 $b < e$: bottlenecked here

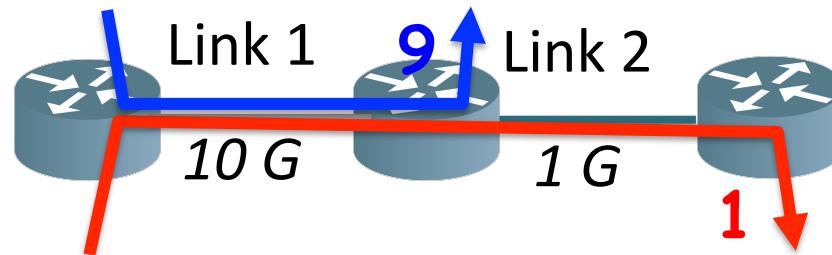
Red Flow's Per-Link State	Link	
	1	2
bn. rate (b)	10	1
state @ alloc. (a)	B @ 10	B @ 1



Link 2's State	
NumB	1
SumE	0

The n-PERC Algorithm

Red Flow's Per-Link State	Link		Flow's Link State	Link	
	1	2		1	2
bn. rate (b)	10	1	b)	10	1
state @ alloc. (a)	E @ 1	B @ 1	alloc. (a)	B @ 10	B @ 1

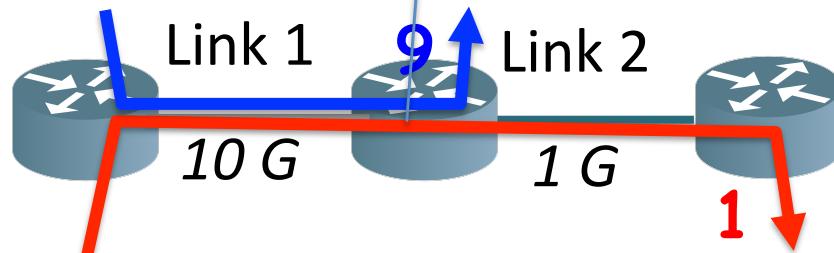


Link 2's State	
NumB	0
SumE	0

The n-PERC Algorithm

$b = 1 \text{ Gb/s}$
 $e = 10 \text{ Gb/s}$
 $a = \min(b, e) = 1 \text{ Gb/s}$
 $b < e$: bottlenecked here.

Red Flow's Per-Link State	Link	
	1	2
bn. rate (b)	10	
state @ alloc. (a)	B @ 10	



Link 2's State	
NumB	1
SumE	0

Propagating a Bad Bottleneck Rate

$$b = 1 \text{ Gb/s}$$

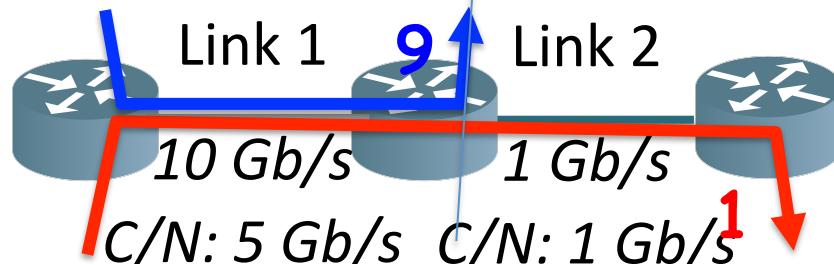
$$e = 0.5 \text{ Gb/s}$$

$$a = \min(b, e) = 0.5 \text{ Gb/s}$$

$e < b$: bottlenecked

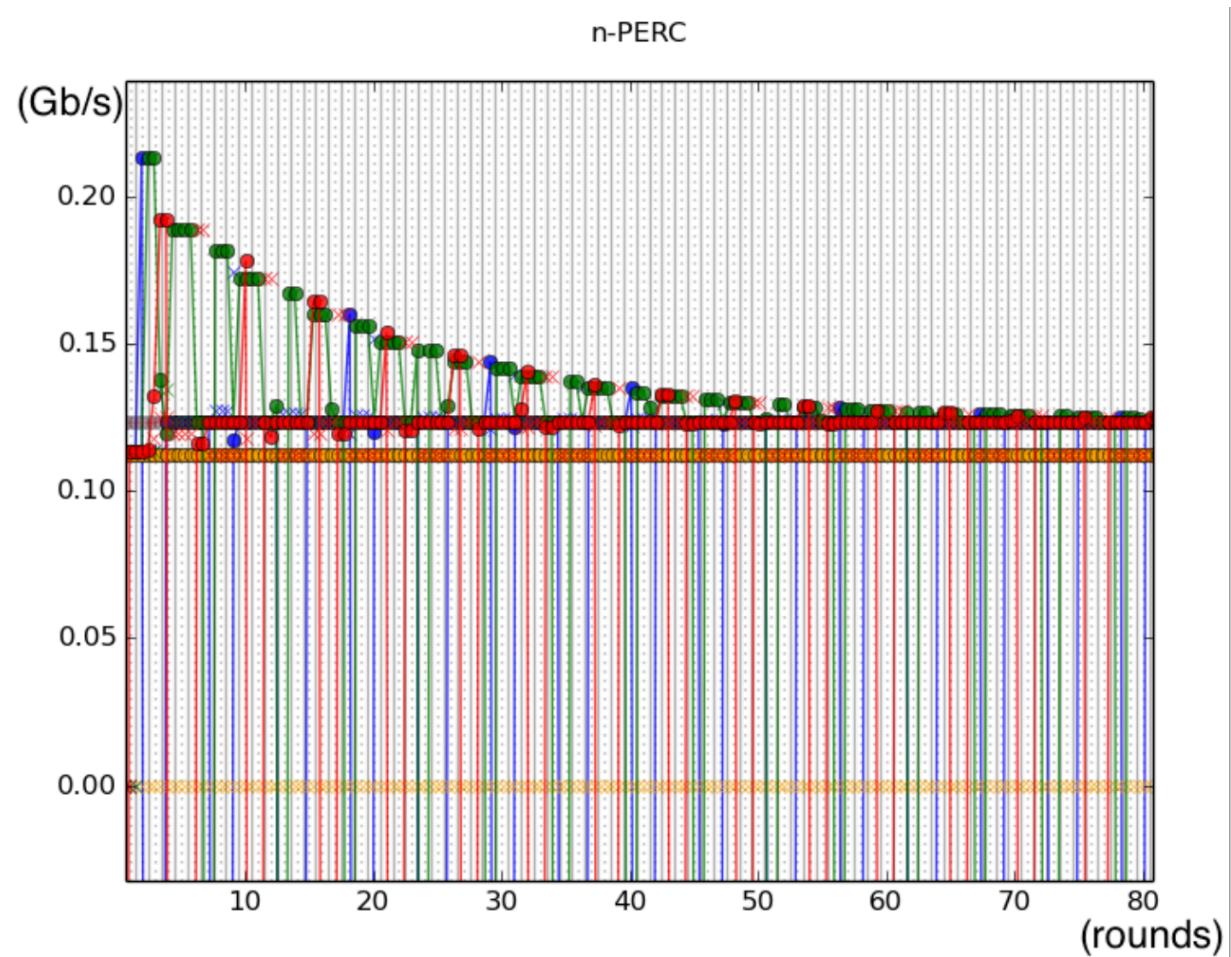
elsewhere at 0.5 Gb/s!?

Red Flow's Per-Link State	Link	
	1	2
bn. rate (b)	0.5	



Link 2's State	
NumB	1
SumE	0

Bad bottleneck rates delay convergence



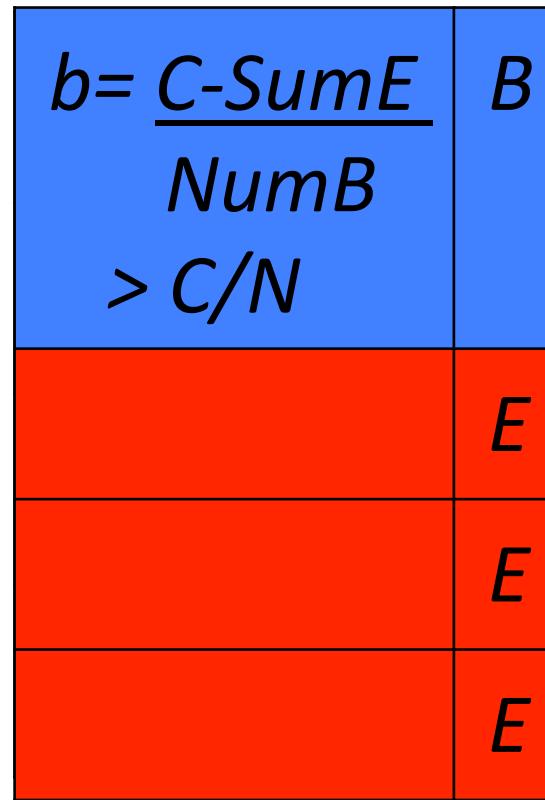
Key Insight

- Recognize when the bottleneck rate is safe to propagate, using constant state.
- MaxE , maximum rate allocated to E flow in the last two RTTs.
- Propagate b only when $b \geq \text{MaxE}$.
- Propagated rate at least fair share rate.
- Link with lowest fair share recognizes B flows.

The s-PERC Algorithm

- Calculate bottleneck rate assuming flow will be bottlenecked here, $b = (C - \text{SumE})/\text{NumB}$
^{propagated}
- Calculate limit rate, e , lowest bottleneck rate elsewhere
- Allocate lower of the two $a = \min(b, e)$
- Update packet and link
- Decide whether to propagate b ($b \geq \text{MaxE?}$)

Link propagates at least its Fair Share Rate.



propagate
 $b > MaxE$.

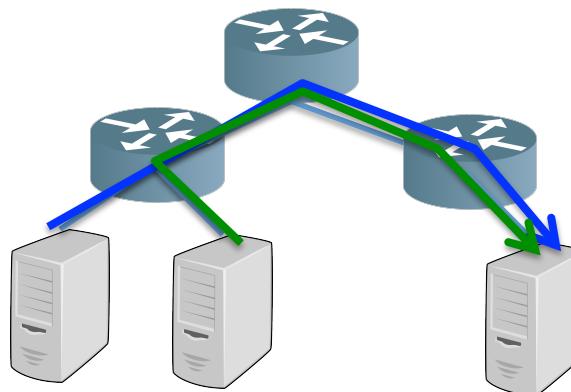
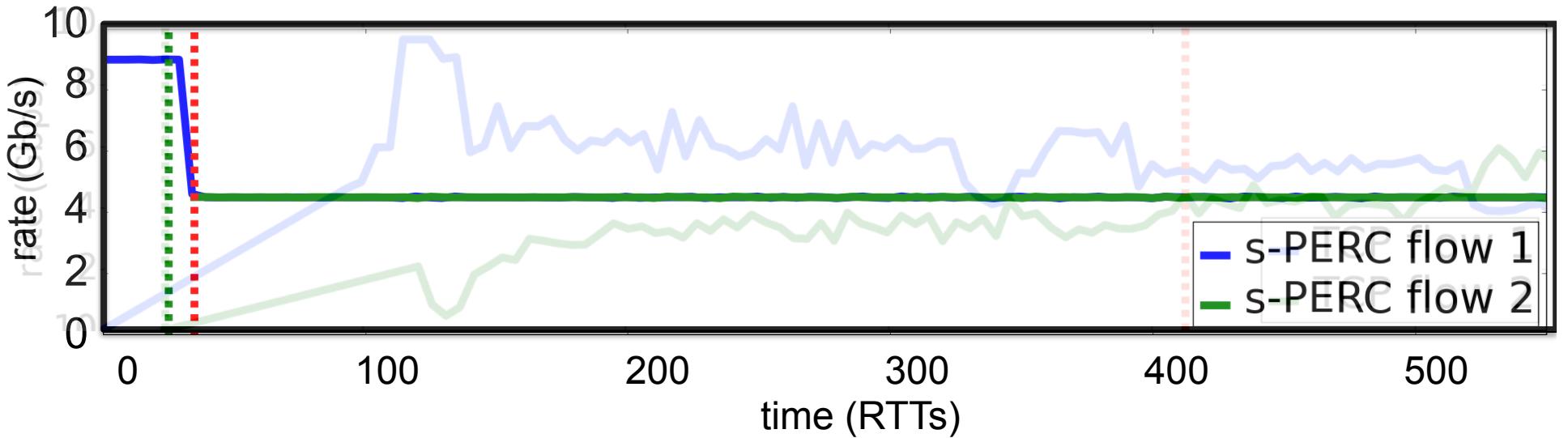
$MaxE < C/N$

Convergence of s-PERC

s-PERC converges in $6N$ RTTs, where N is the number of steps that $k=2$ Waterfilling takes for the same routing matrix and link speeds.

s-PERC is practical

4x10 Gb/s NetFPGA SUME prototype



10 Gb/s links
0.2 ms RTT

Fairness + Throughput for Long Flows, Low Latency for Short Flows

Convergence Times (Fairness)

- 10x faster than RCP in DC networks
- 1.3-6x faster than RCP in WAN networks

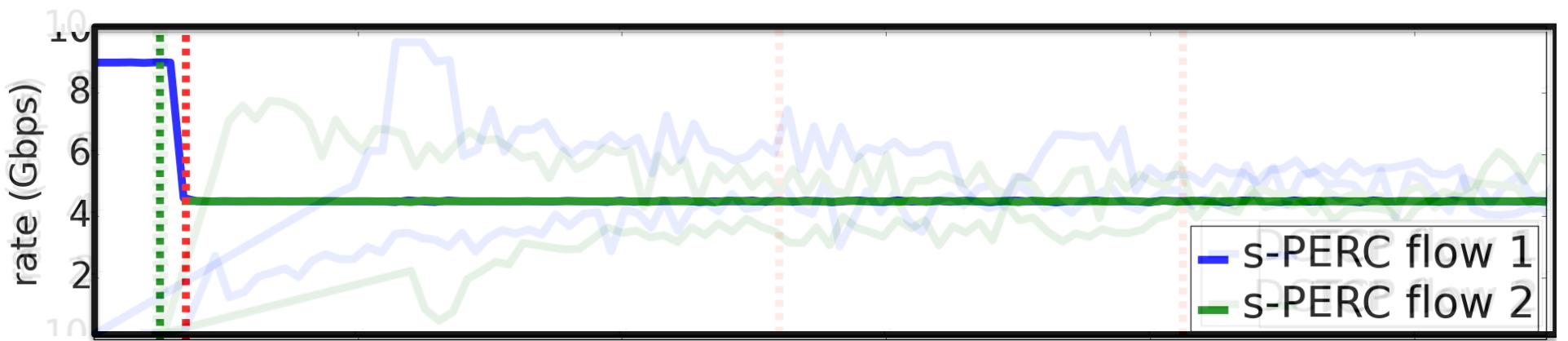
Flow Completion Times (Throughput for Long Flows, Latency for Short Flows)

- Close to Ideal Max-Min For Long Flows
- Close to Minimum For Short Flows

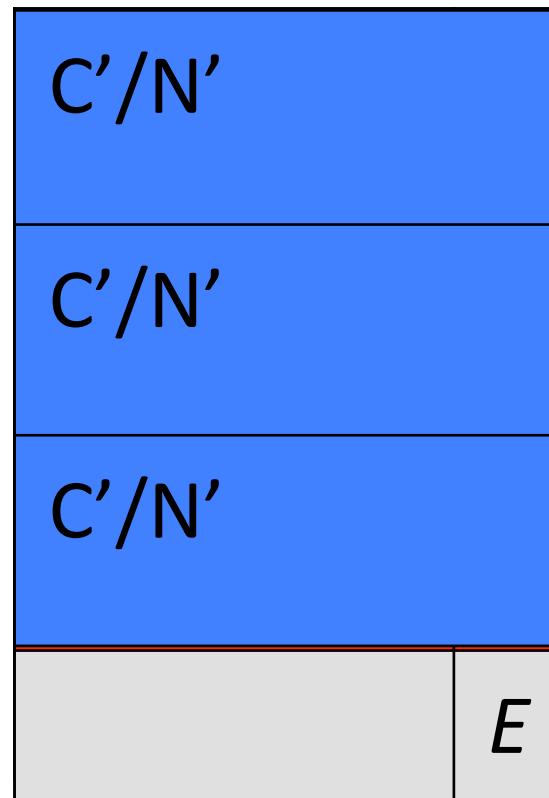
Conclusion

- Reactive algorithms converge slowly relative to flow durations in fast networks.
- Proactive PERC algorithms, use information about the network and flows to compute rates directly, they are much faster.
- PERC algorithms are promising and relevant today and s-PERC, a distributed PERC algorithm for max-min fair rates without per-flow state, is only a first step.

Thank you!



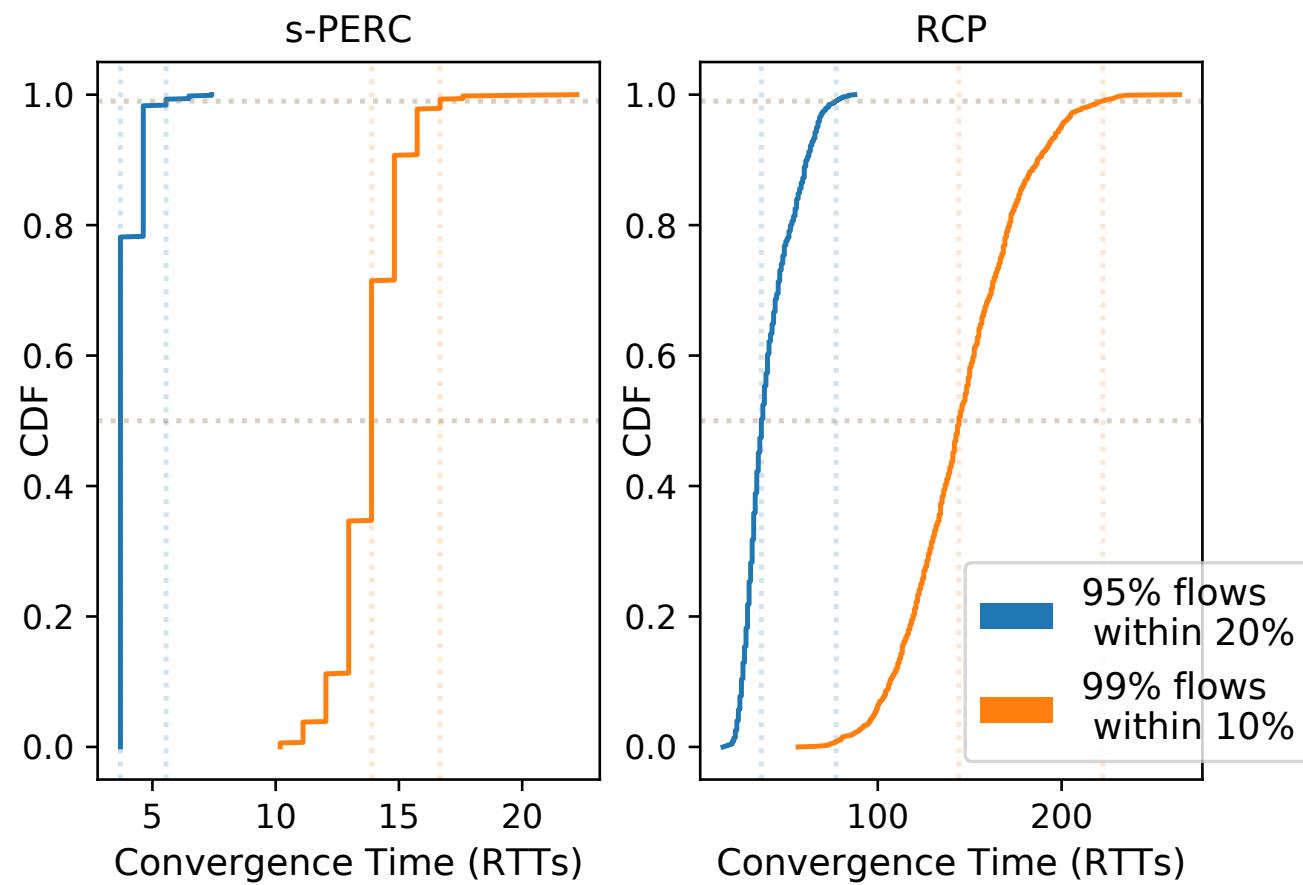
The bottleneck rate propagated by a link is at least its Fair Share Rate.



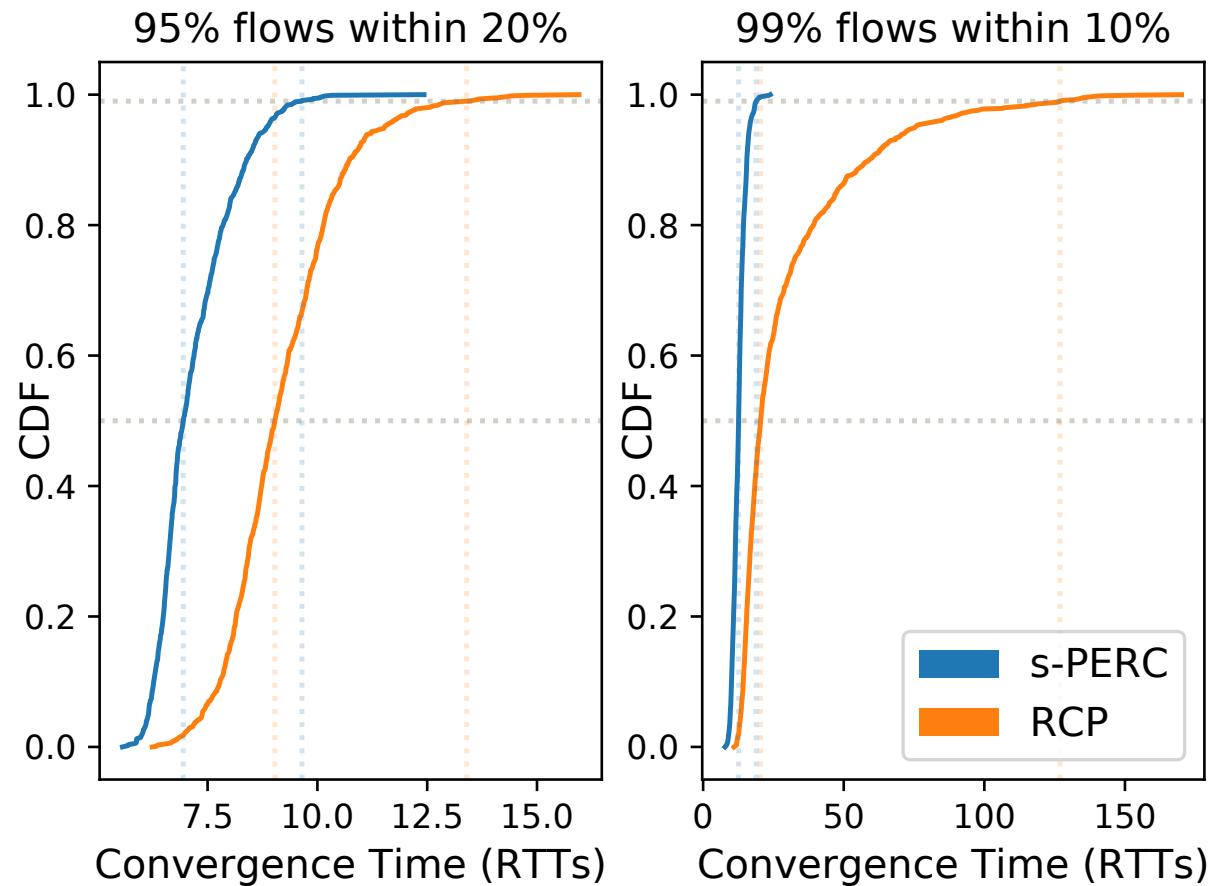
propagate
 $b > \text{Max}E$.

$\text{Max}E < C/N$

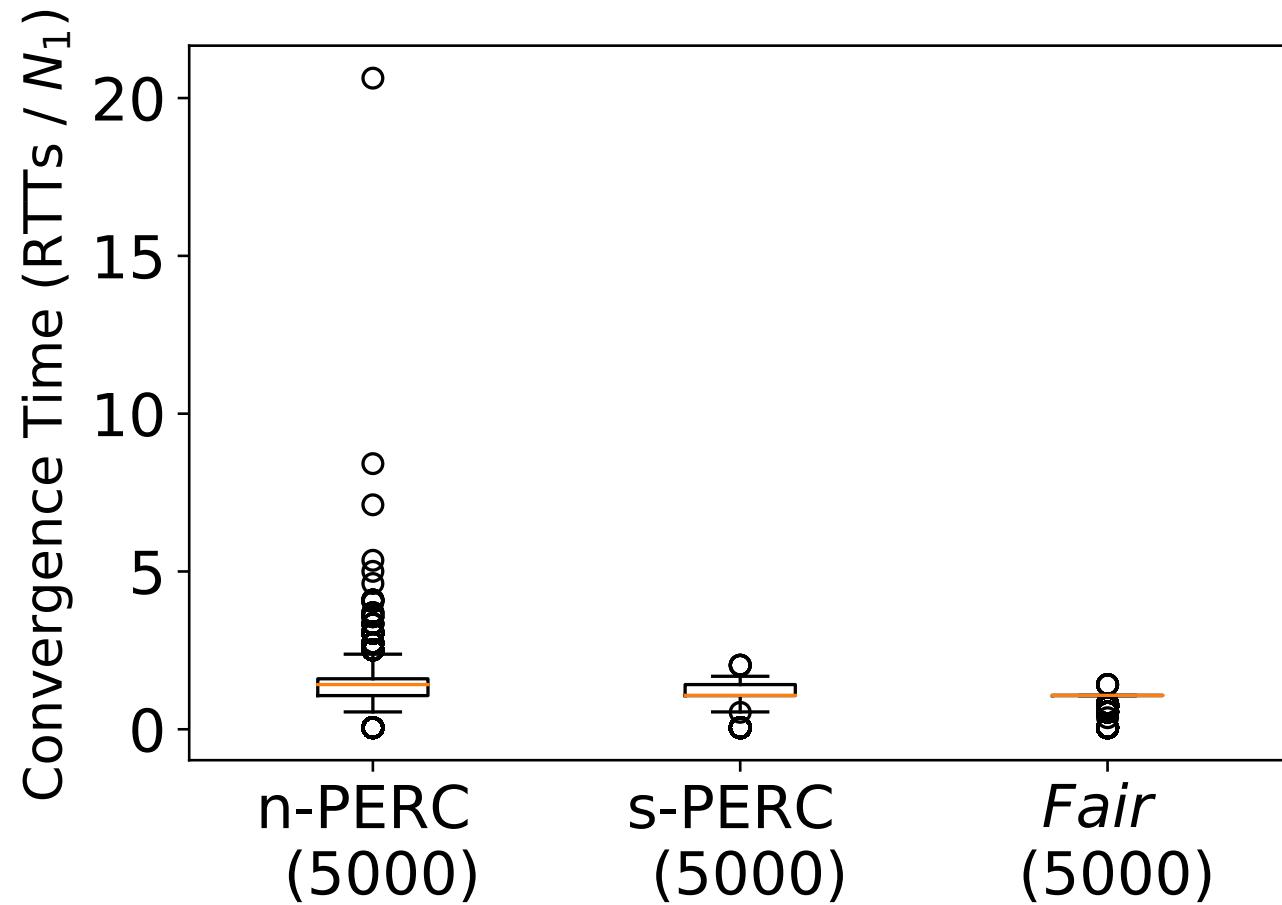
Convergence Time (DC, 100G)



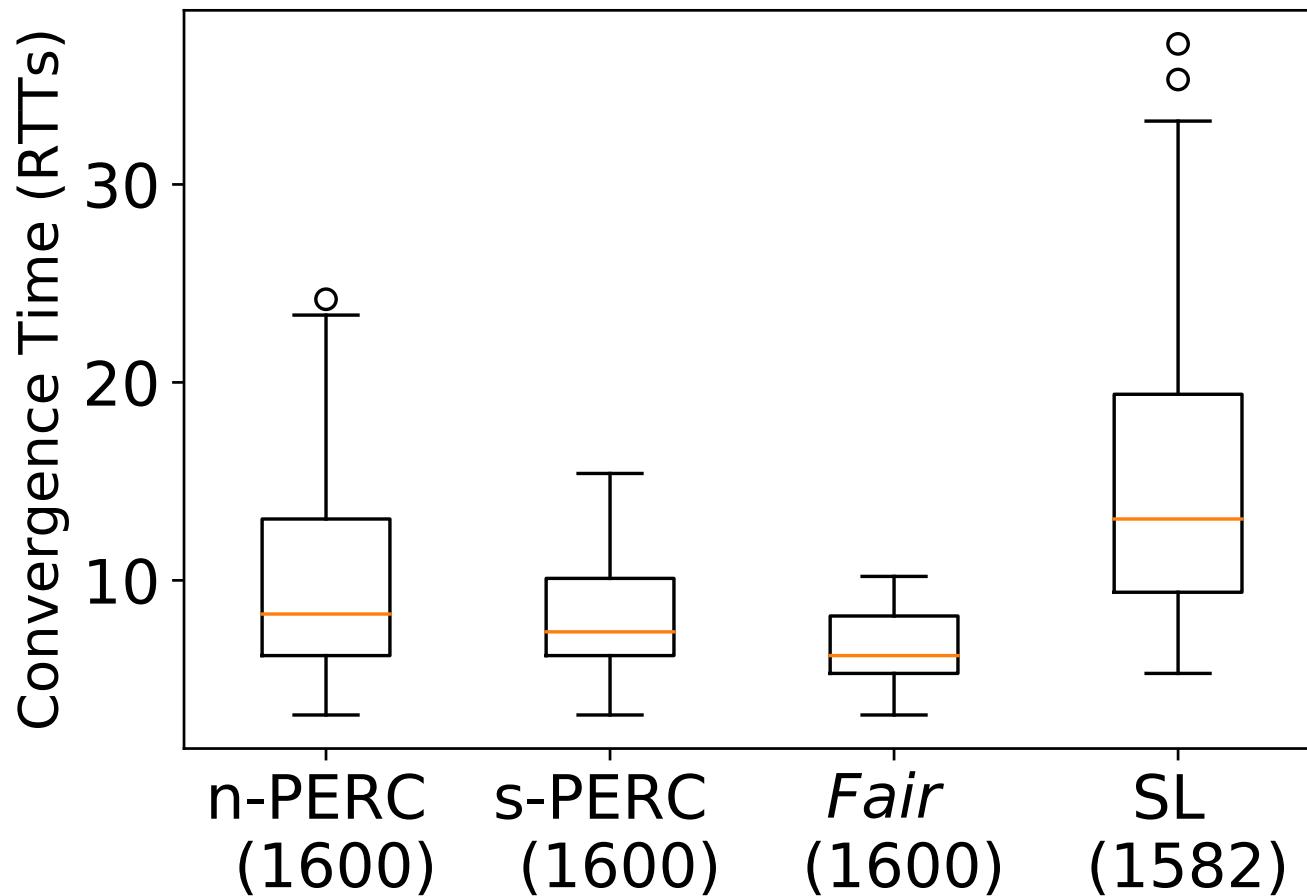
Convergence Times (WAN, 10G)



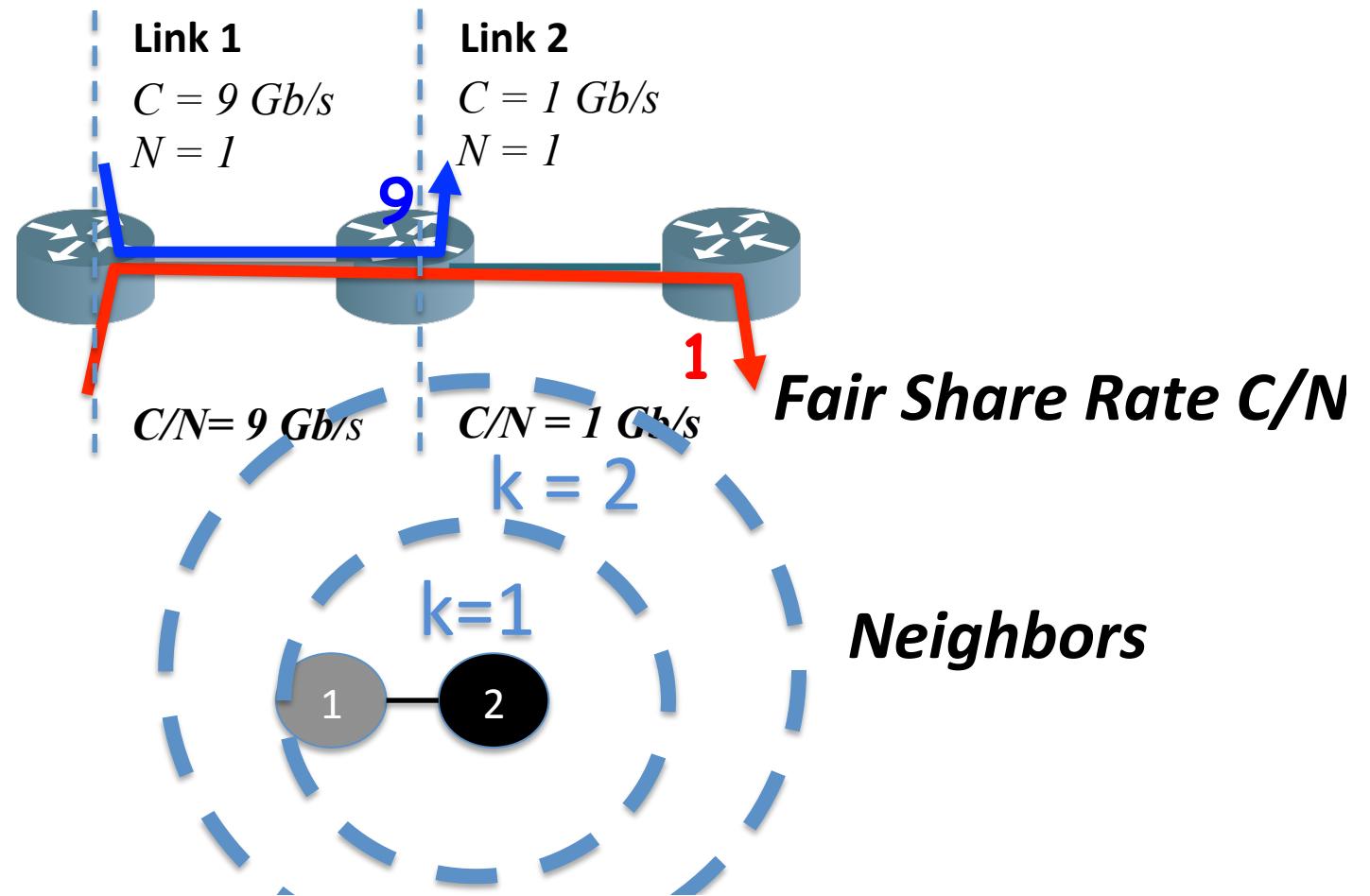
Convergence Times (PERC, worst-case)



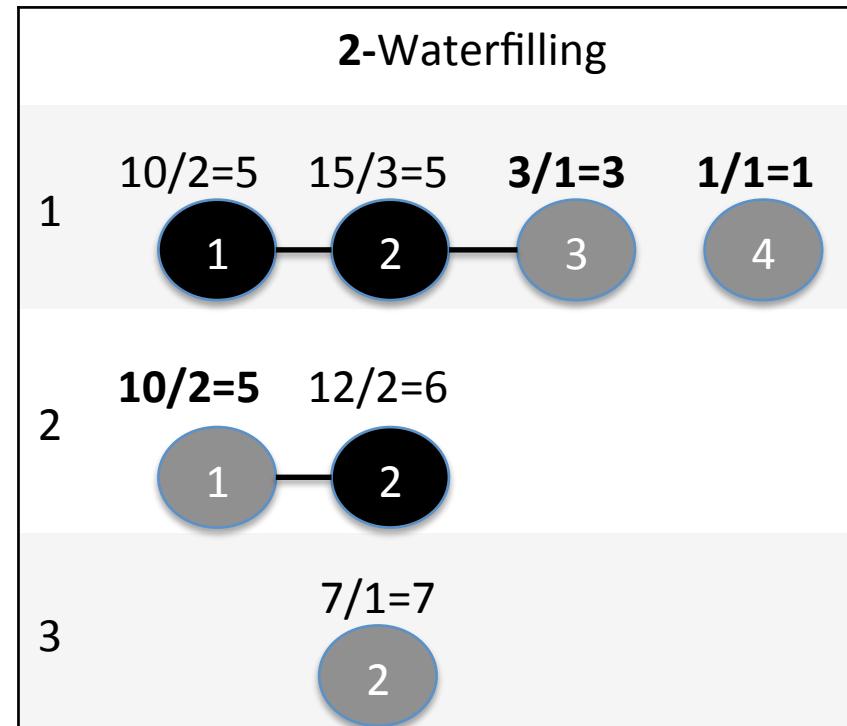
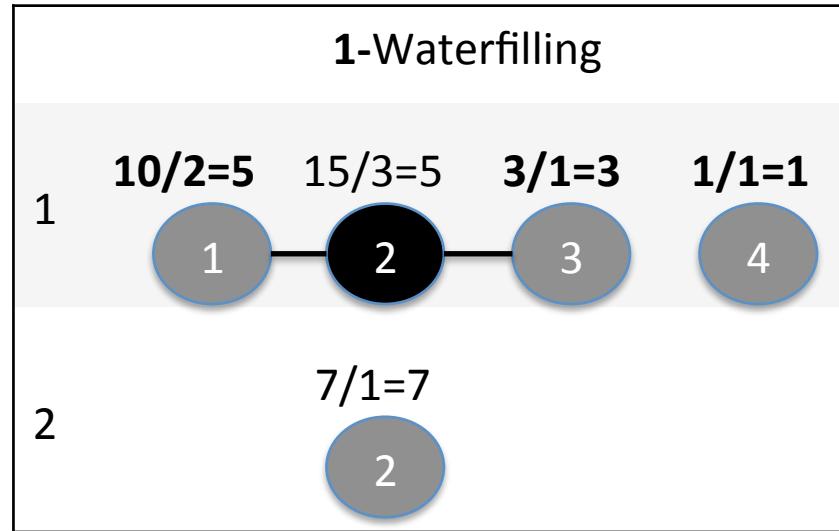
Convergence Times (PERC, common)



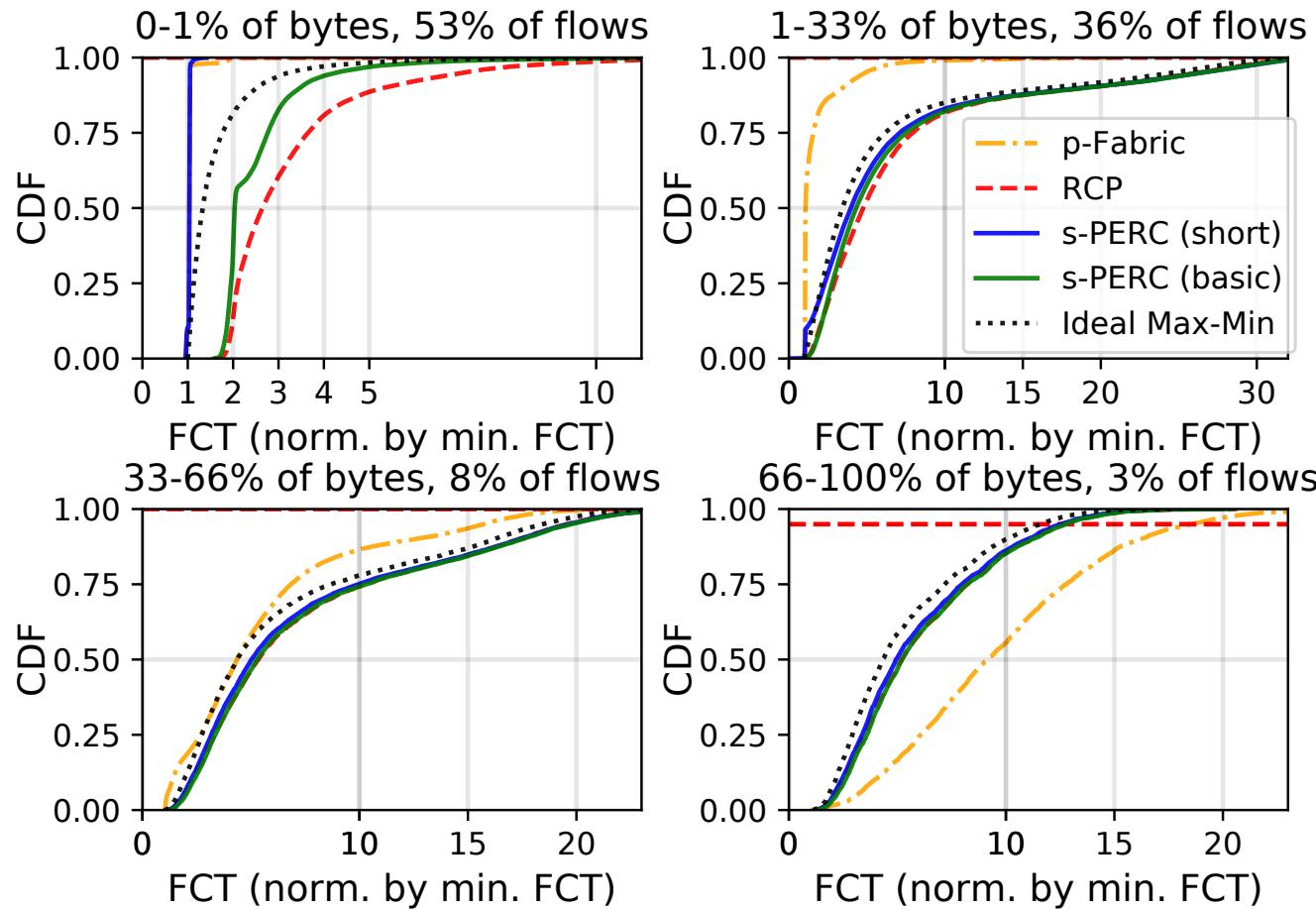
k-Waterfilling: Recursive Algorithms To Compute Max-Min Rates



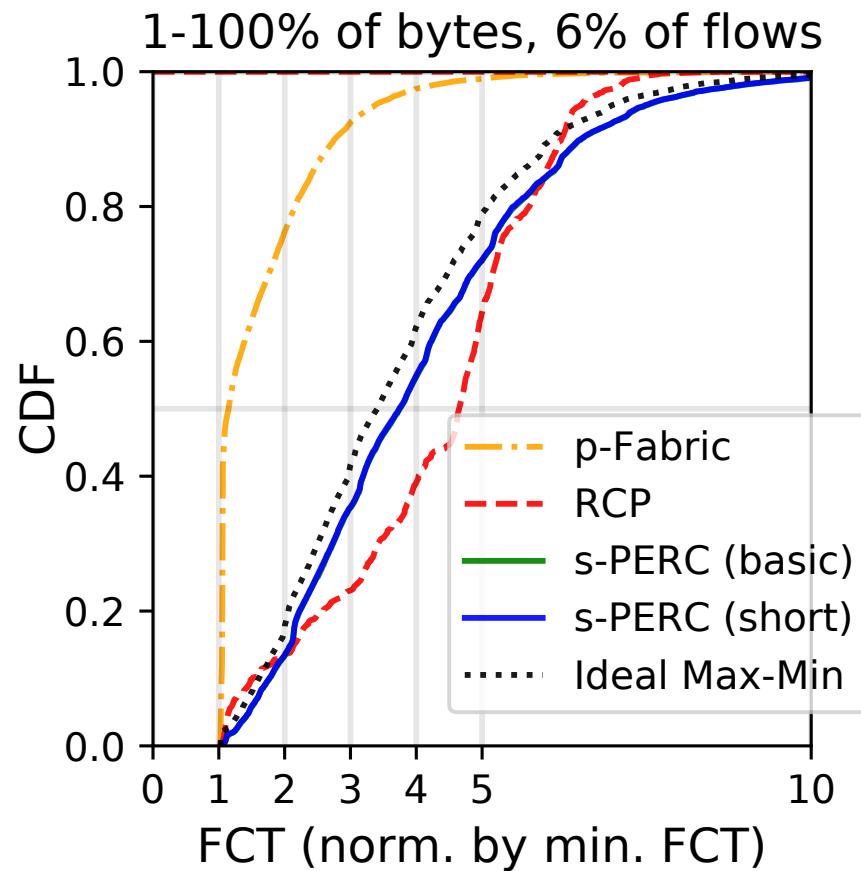
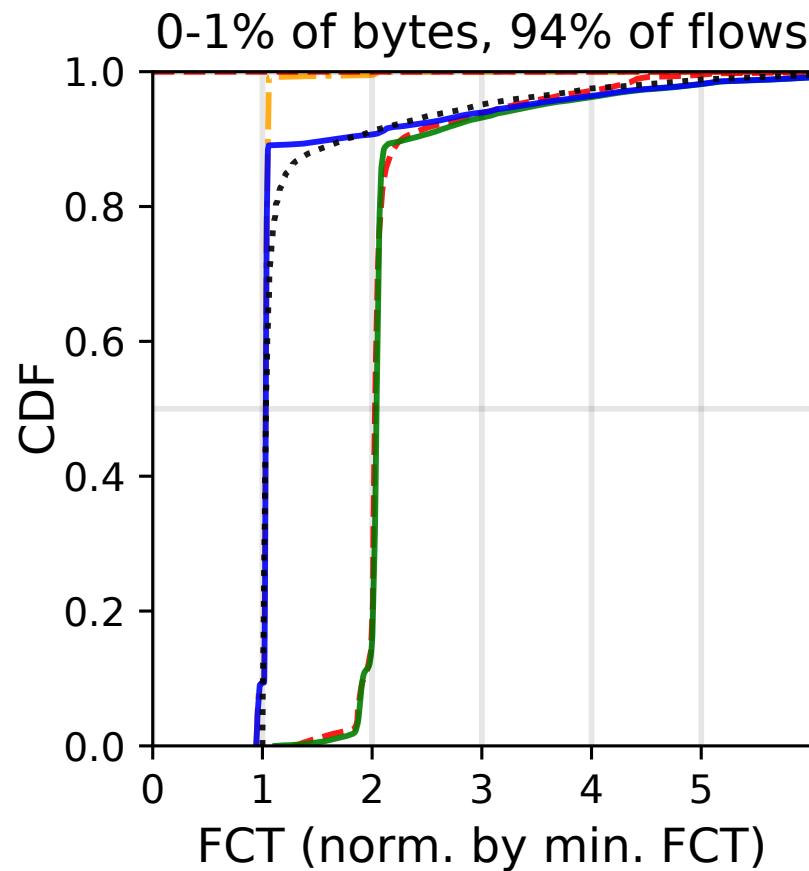
k-Waterfilling (k=1, 2)



FCT (DC, 100G, search, 60%)



FCT (DC, data-mining, 60%)



Testbed, 2-level experiment

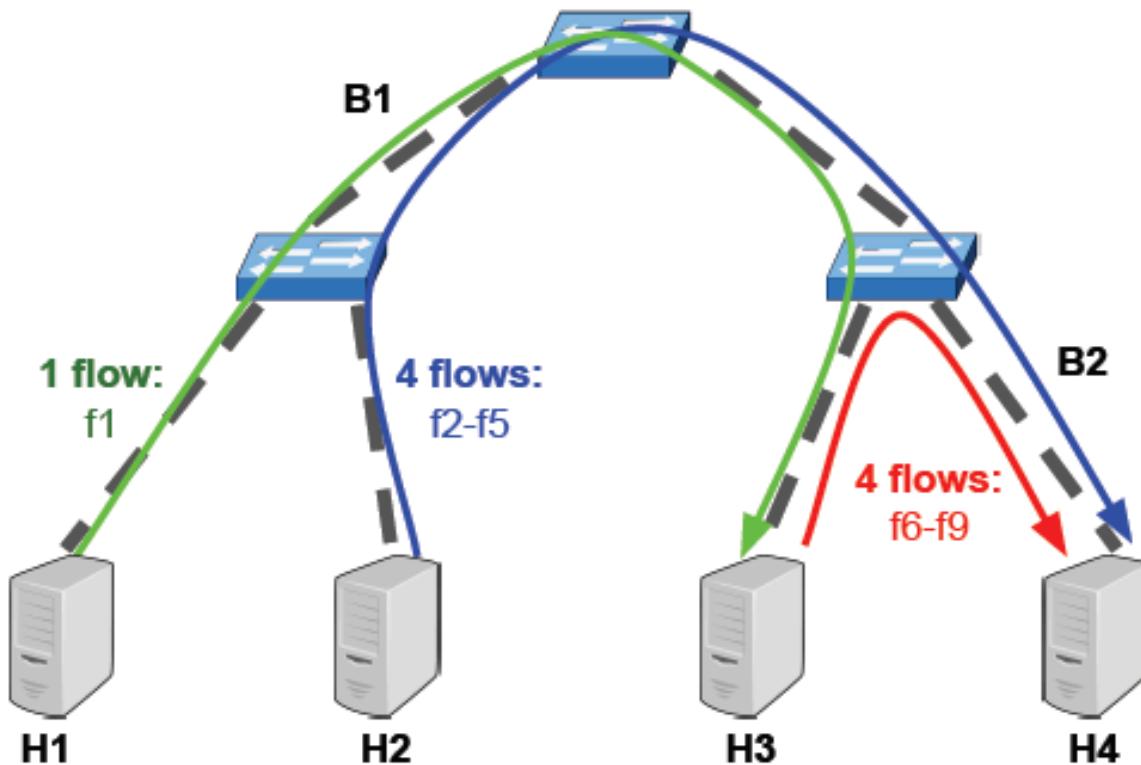


Fig. 11. The topology and traffic pattern used for the two-level dependency chain experiment.

Why 6 RTTs per level?

Time	Updates of flow f at link k	Updates of flows f, g at link l (1st-deg.)	Updates of flow g at link j (2nd-deg.)
T		$b \geq c_k/n_k$ or $i = 1$	$b \geq c_k/n_k$ or $i = 1$
$T + 1 \cdot round$	$e \geq c_k/n_k$	$e \geq c_k/n_k$	
$T + 2 \cdot round$			
$T + 3 \cdot round$	$b = c_k/n_k$ $s = B$ $x = b = c_k/n_k$	$b > c_k/n_k$	
$T + 4 \cdot round$	$i = 0$		
$T + 5 \cdot round$		$e = c_k/n_k$ $s = E$ $x = e = c_k/n_k$ (for flow f)	

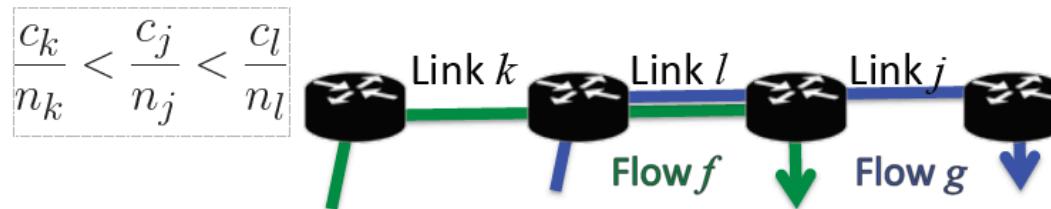


Fig. 5. Links k and j share flows f and g respectively with link l .

Maintaining MaxE with constant state.

- Timeout every RTT (e.g., 10^5 cycles)
- Initially,
 - $\text{MaxE}' = 0$ (“maximum E alloc. in this RTT”)
 - $\text{MaxE} = 0$ (“maximum E alloc. in last two RTTs”)
- On timeout:
 - MaxE reset to MaxE' (so it doesn’t increase forever)
 - MaxE' reset to 0 (reflects max. E since last timeout)