

---

# Let It Flow

## Resilient Asymmetric Load Balancing with Flowlet Switching

---

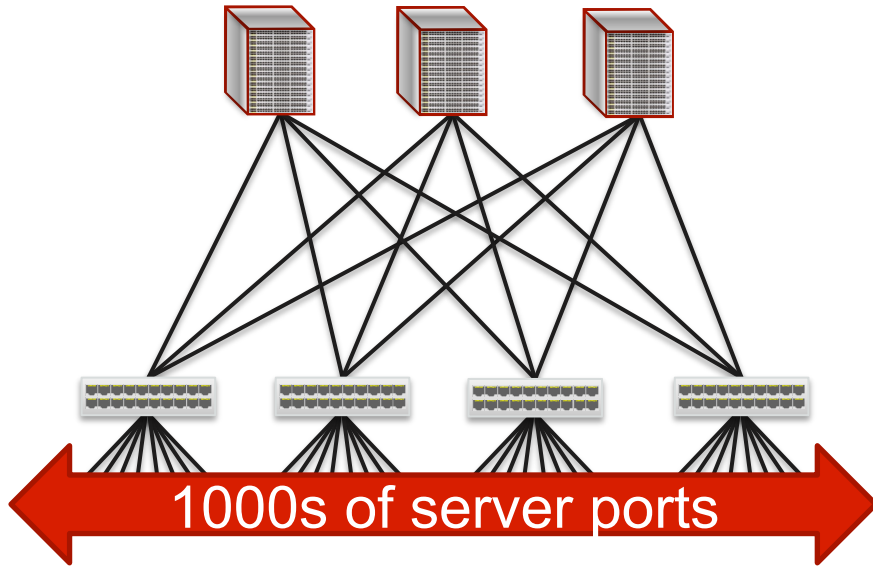
Erico Vanini\*, Rong Pan\*, Mohammad Alizadeh†, Parvin Taheri\*, Tom Edsall\*



# Load Balancing in Data Centers

---

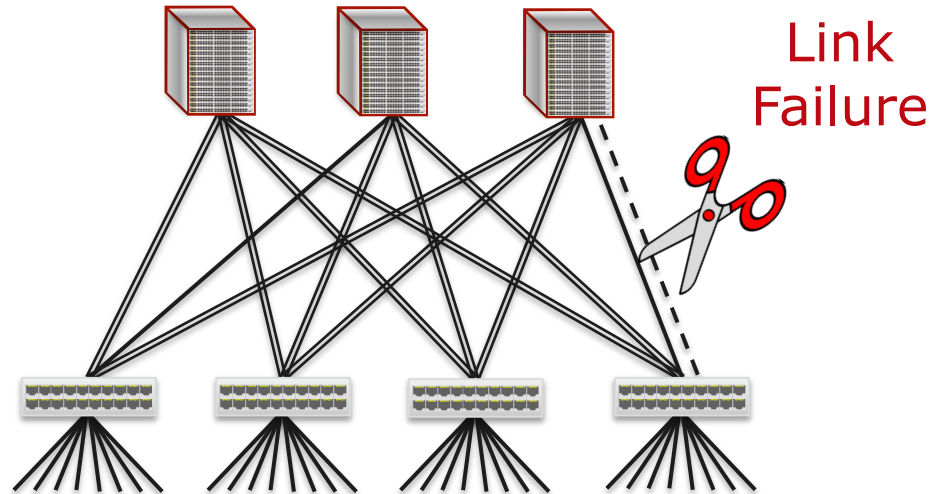
Multi-rooted tree



- **Goal:** Avoid congestion hotspots
- Active research area
- Solved for symmetric topologies
- Still open question in asymmetric scenarios

# Asymmetry Is Common in Practice

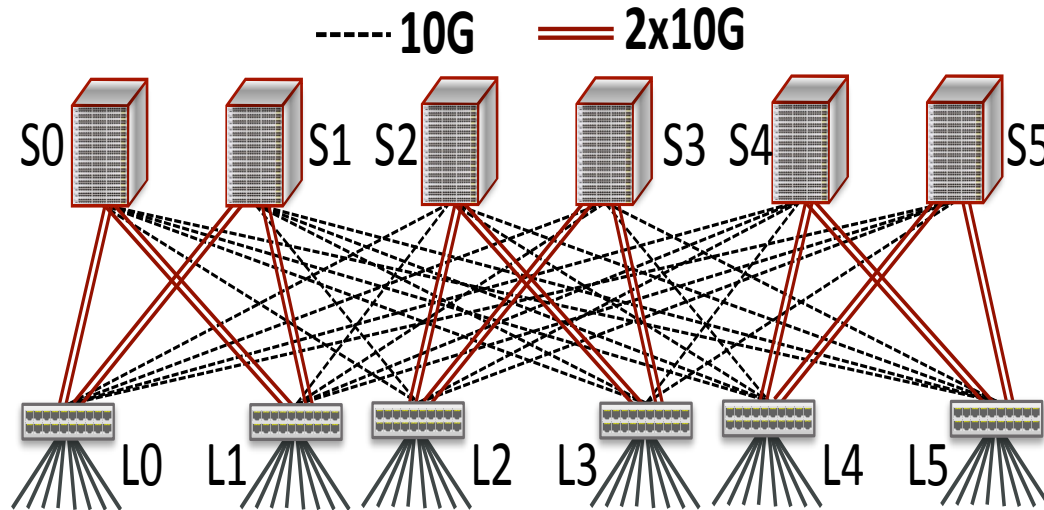
---



# Asymmetry Is Common in Practice

---

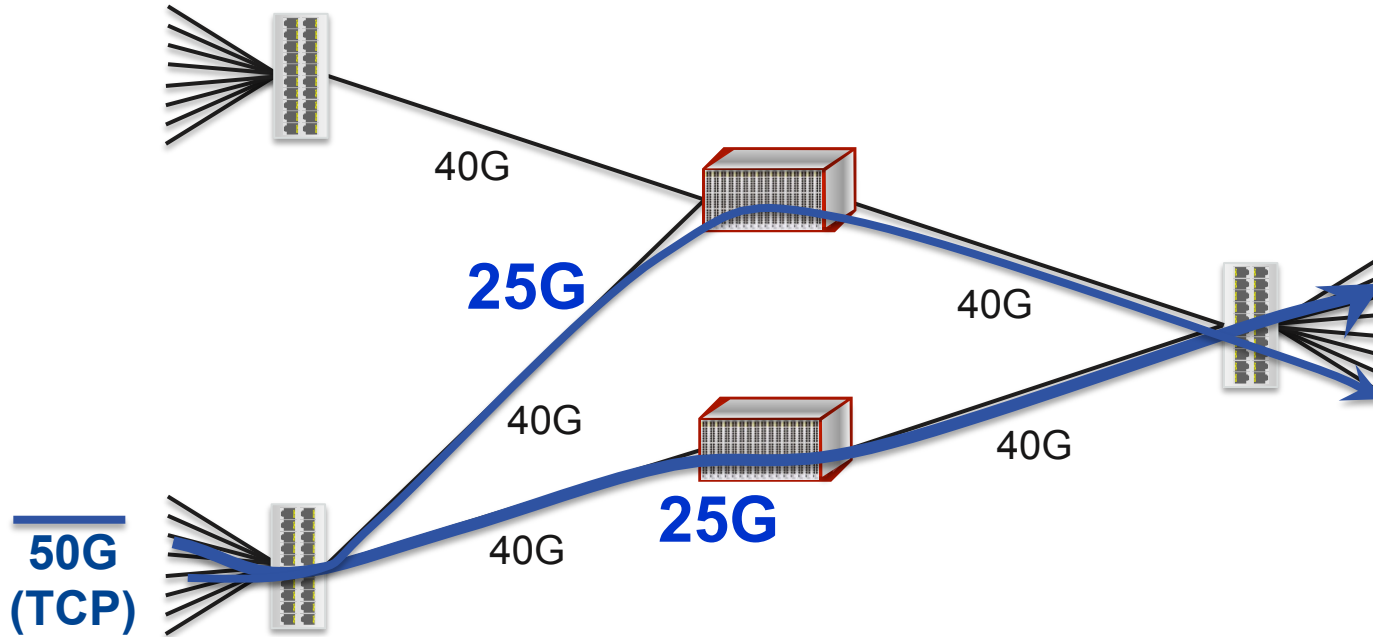
**Imbalanced striping:** # of ports indivisible by # of switches in other tier



Zhou et al. "[WCMP: Weighted cost multipathing for improved fairness in data centers](#)," CoNEXT 2014.

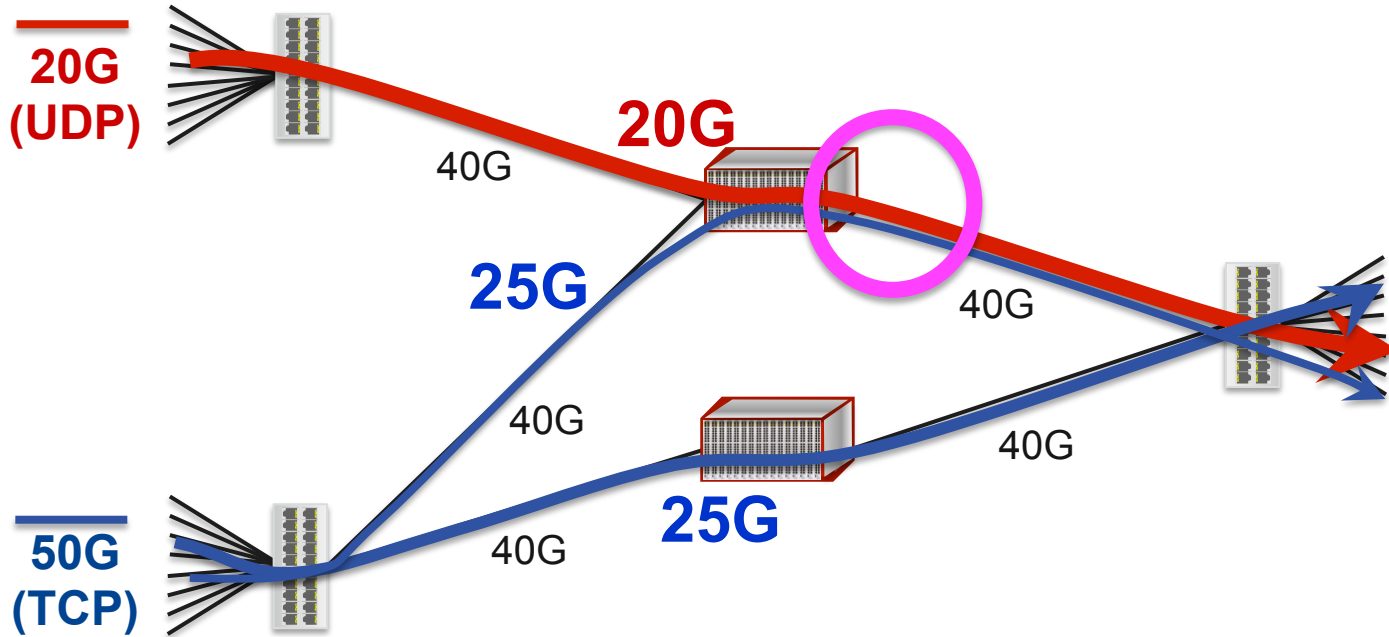
# Dealing with Asymmetry

---

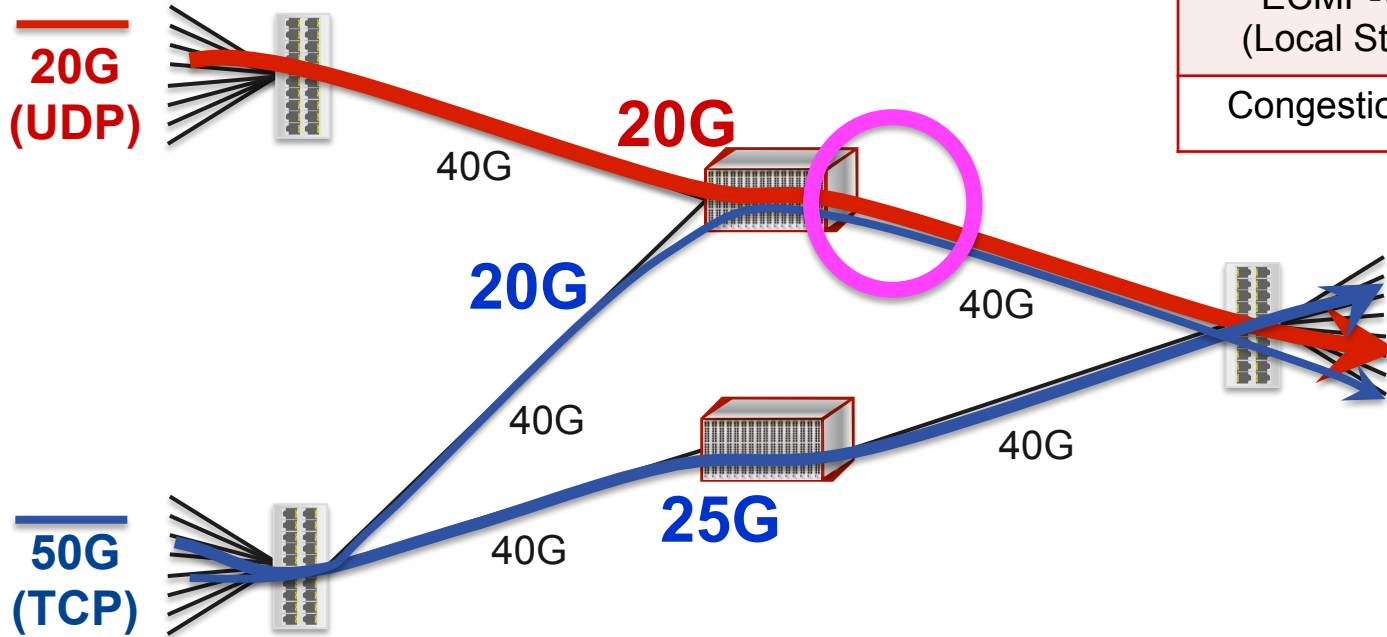


# Dealing with Asymmetry

---

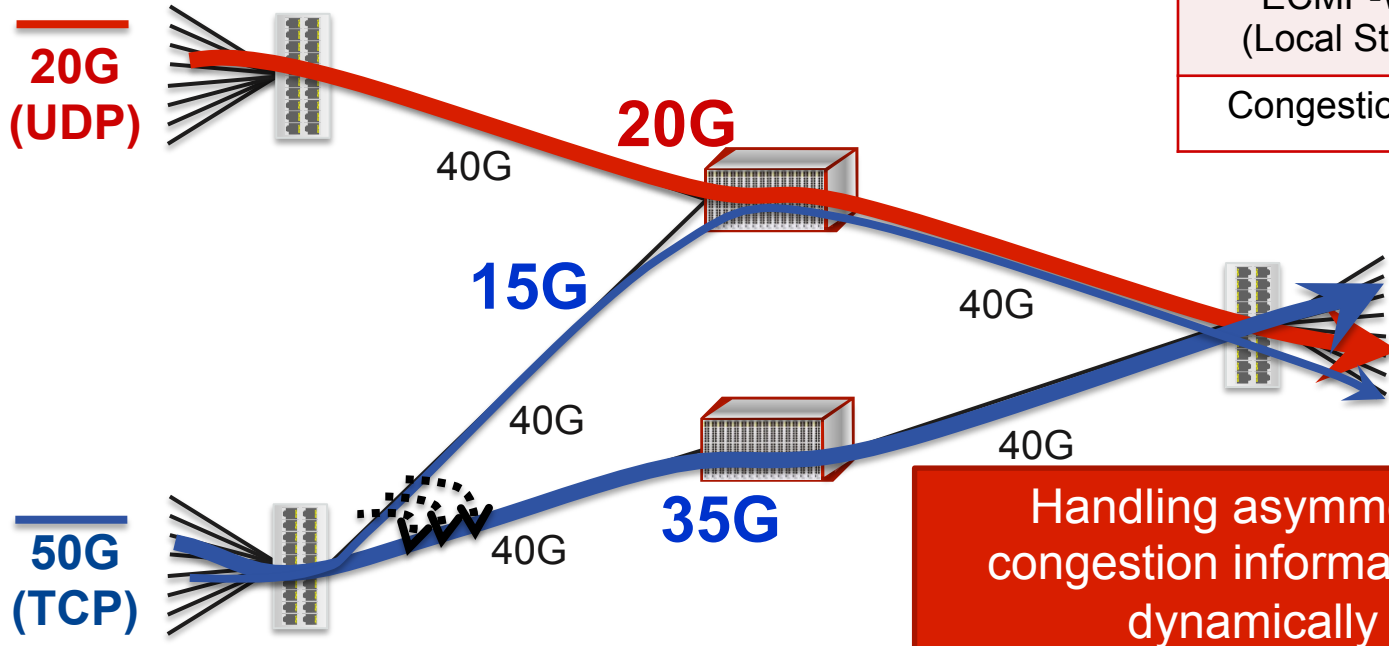


# Dealing with Asymmetry



Scheme	Thrput
ECMP-WCMP (Local Stateless)	<b>65G</b>
Congestion-Aware	

# Dealing with Asymmetry

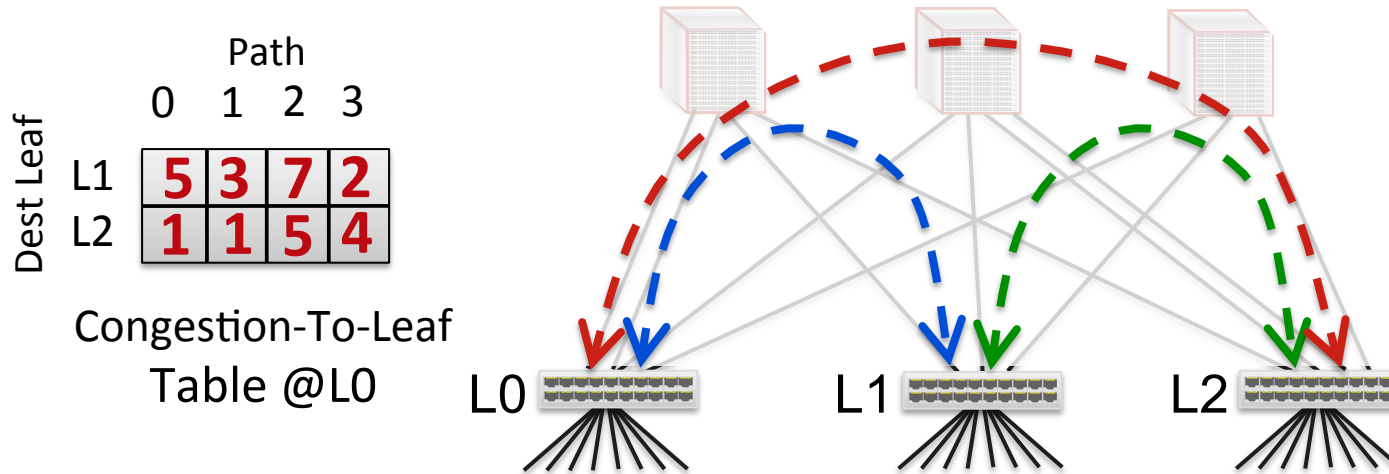


Scheme	Thrput
ECMP-WCMP (Local Stateless)	<b>65G</b>
Congestion-Aware	<b>70G</b>



# Example: CONGA

1. Leaf switches (top-of-rack) track congestion to other leaves on different paths
2. Use this information to minimize bottleneck utilization



# Existing Load Balancing Schemes

---

## **Congestion-aware decisions: complex**

- Measure and feed back congestion in real time
- CONGA, Hedera, HULA, MPTCP, FlowBender,...

## **Congestion-oblivious decisions: simple**

- Random, round robin, hashing decision process
- ECMP, WCMP, Packet-Spray, Presto,...

**Is there a simple load balancing scheme  
(with congestion-oblivious decisions)  
that can handle asymmetry?**

# LetFlow

---

## Simple:

Randomly assign Flowlets to available paths

## Flowlets:



“Flowlets are bursts from same flow separated by at least  $\Delta$ ”

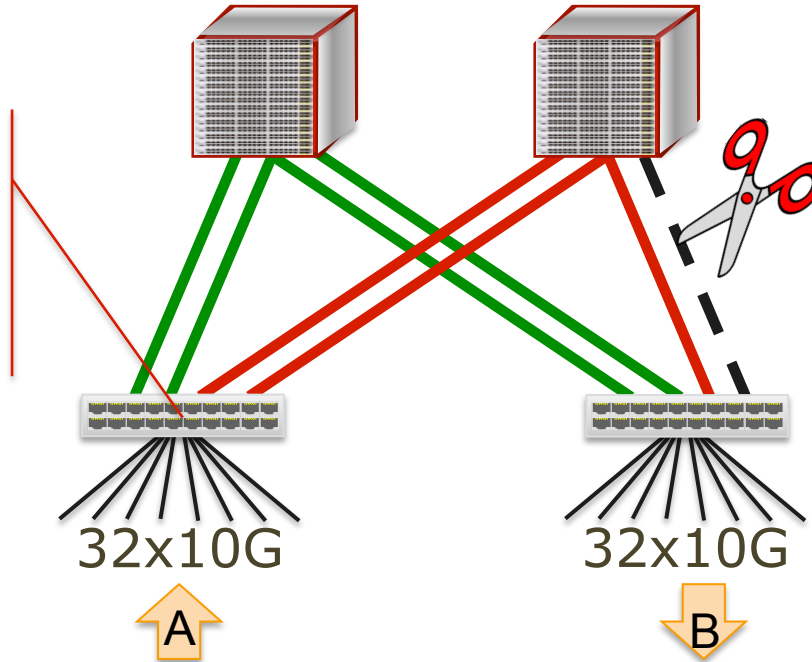
*“the main origin of flowlets is the burstiness of TCP at RTT and sub-RTT scales.”*

Kandula et al, [“Dynamic load balancing without packet reordering”](#), (2007)

# Simple Asymmetric Scenario

---

Detect and  
randomly assign  
Flowlets to  
available paths

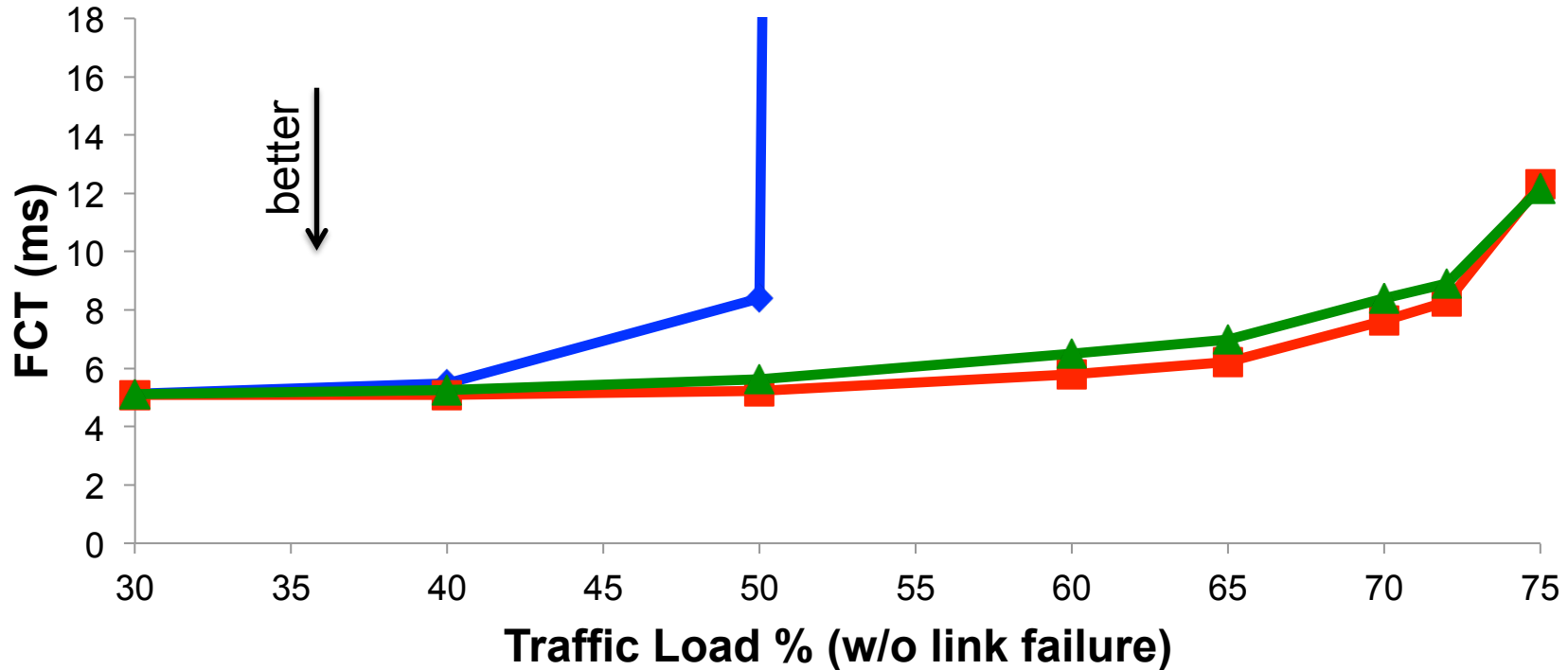
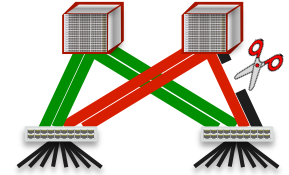


Link  
Failure

Extremely simple!

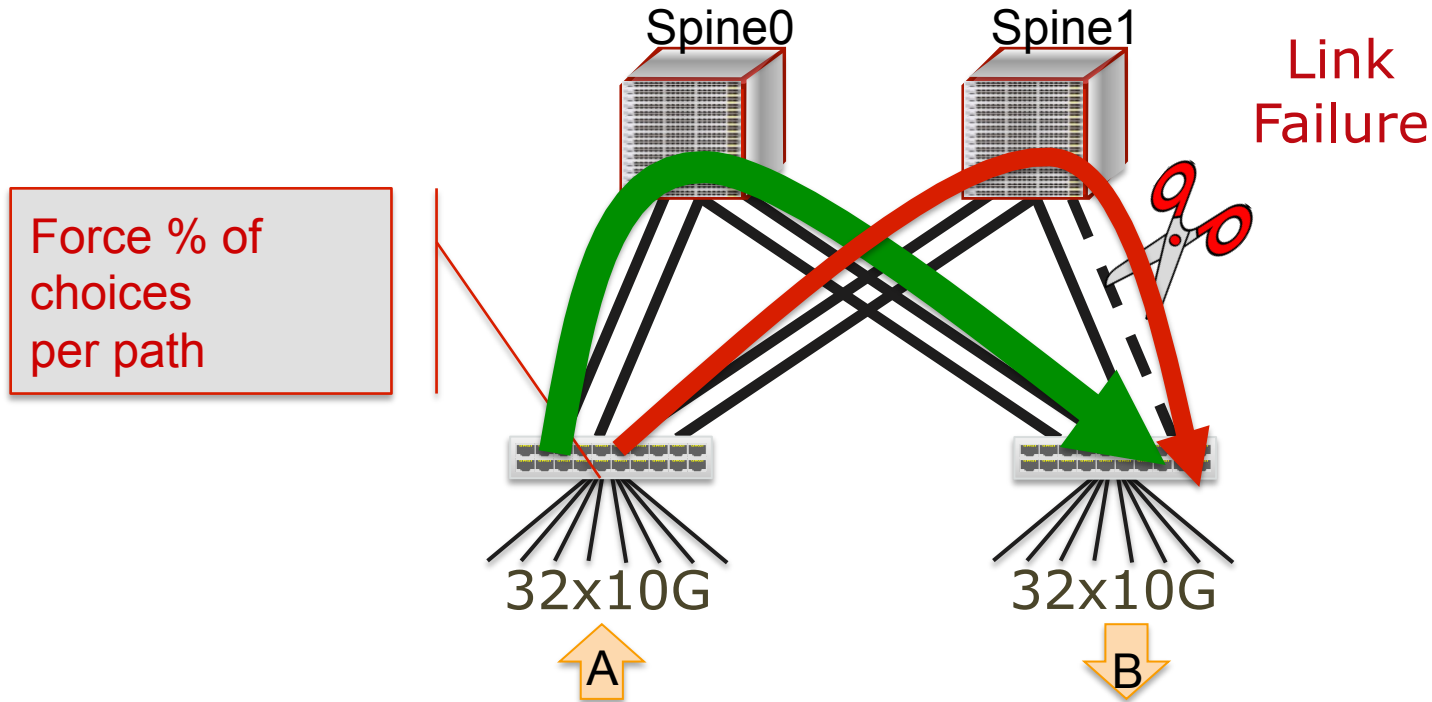
- No measurements
- No feedback
- No congestion state

# LetFlow

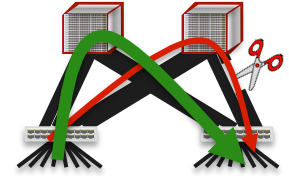


# What's Going On?

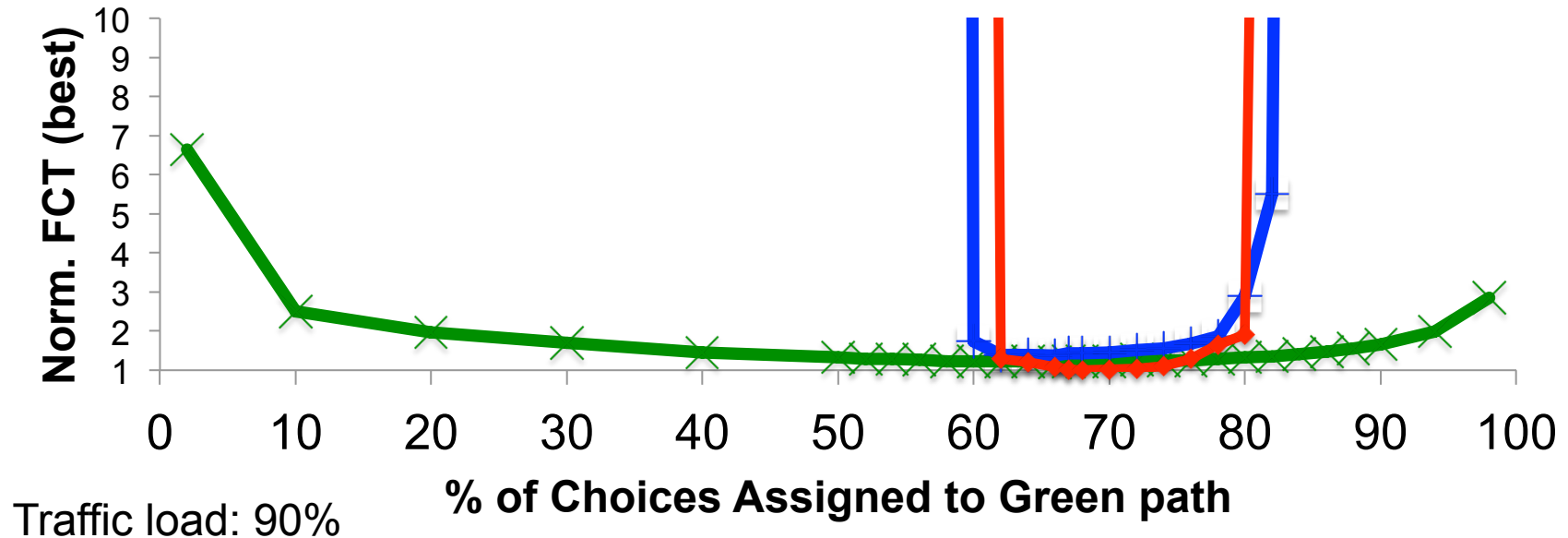
---



# Flowlets are Robust

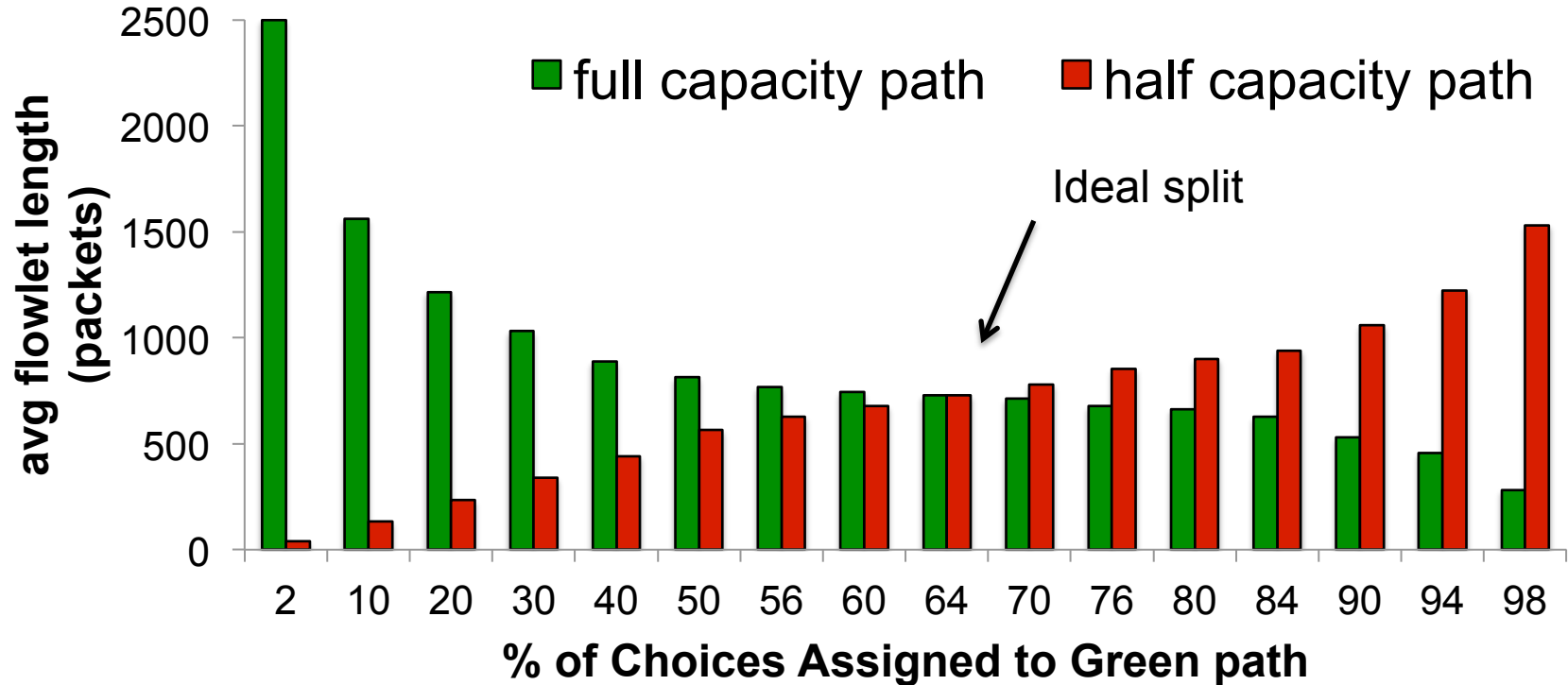
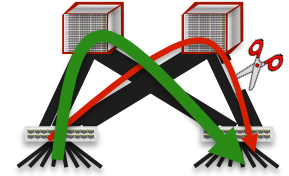


Performance is not sensitive to load balancing decisions





# Flowlet Length



# Flowlets are Elastic

---

- Flowlets change size based on congestion on the path
    - Uncongested path → larger flowlets
    - Congested path → smaller flowlets
- Flowlet sizes **implicitly** encode path congestion information  
... this determines the amount of traffic on each path – not just load balancing decisions

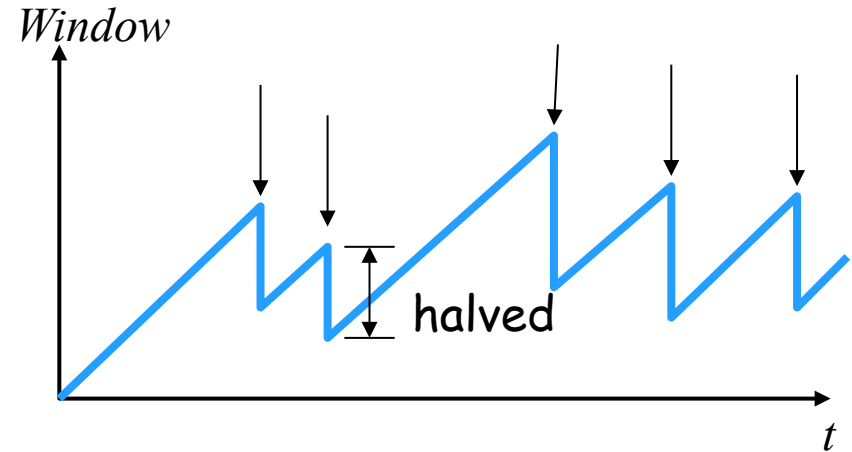
LetFlow *is* congestion-aware,  
despite simple random decisions

# Why Are Flowlets Elastic?

---

- Because of congestion control (e.g., TCP)

- A flowlet gap occurs on
  - Window cuts (Loss/ECN)
  - Latency spikes (ACK clocking)

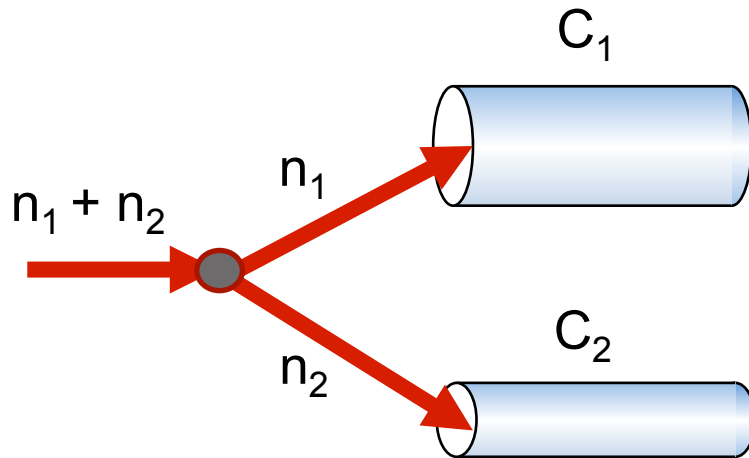


- But, there's a more basic reason, applicable to any congestion control protocol ...

# LetFlow Analysis

---

- Assume flows transmit as Poisson processes

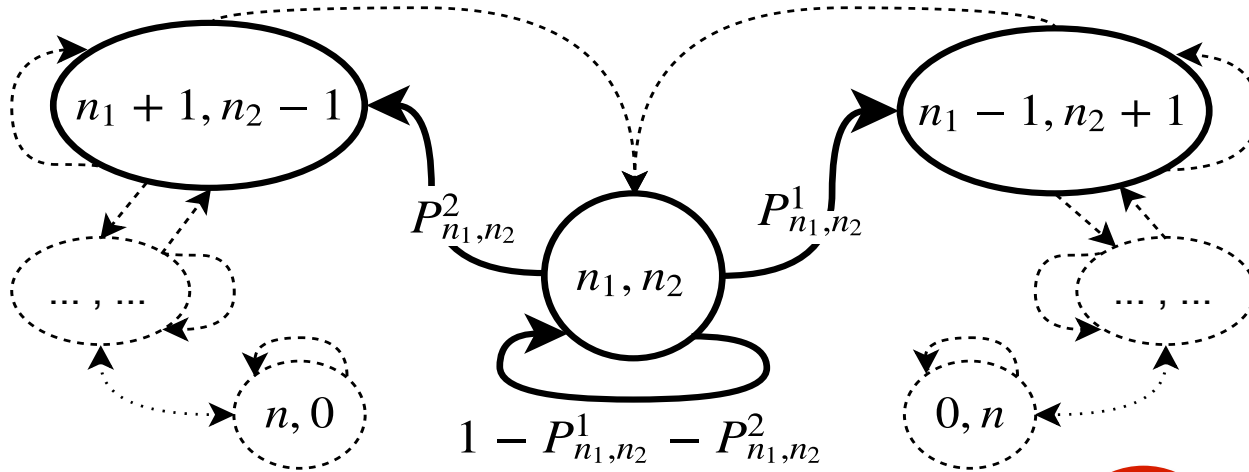
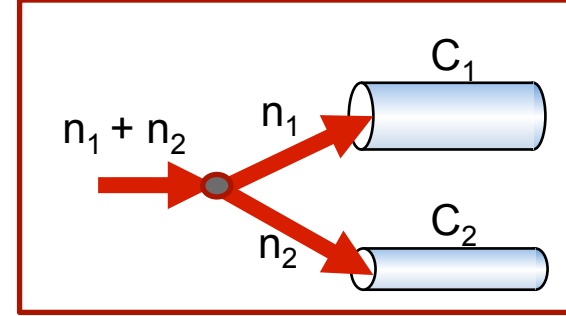


Avg. rate of each flow:

$$\lambda_1 = C_1 / n_1$$

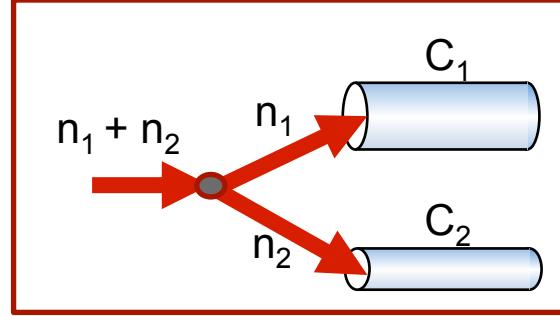
$$\lambda_2 = C_2 / n_2$$

# LetFlow Analysis



State transition probability  $P_{n_1, n_2}^j \approx \frac{c_j}{2(c_1 + c_2)} e^{-\lambda_j \Delta}$ ,  $j \in \{1, 2\}$

# LetFlow Analysis



## Takeaways

1. Flows move from low rate paths (small  $\lambda$ ) to high rate paths (large  $\lambda$ )
2. The flowlet timeout ( $\Delta$ ) is important
  - Shouldn't be too small or large

State transition probability  $P_{n_1, n_2}^j \approx \frac{c_j}{2(c_1 + c_2)} e^{-\lambda_j \Delta}$ ,  $j \in \{1, 2\}$

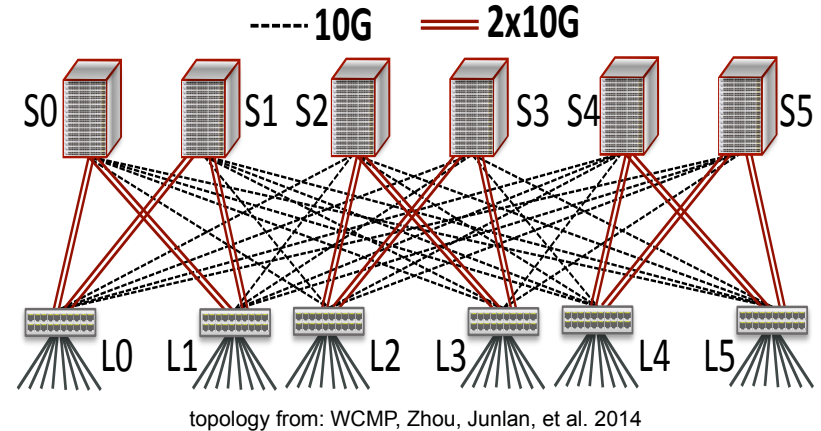
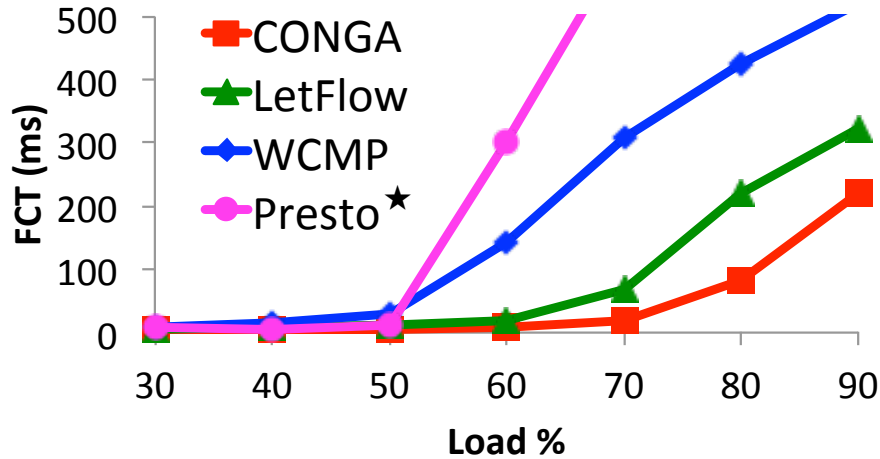
# Experiments Summary

---

Different workloads: web search, data mining, enterprise

- Testbed experiments: ECMP, CONGA, LetFlow
  - 2 leaves 2 spines, 64 servers: symmetric & asymmetric topologies
- Simulations: ECMP, WCMP, Presto★, CONGA, LetFlow
  - Large topology: 6 leaves 6 spines, 288 servers
  - Complex asymmetric topologies: speed mismatch, combined workloads, multitier
  - Different protocols: TCP, DCTCP, DCQCN

# Large Scale Simulations

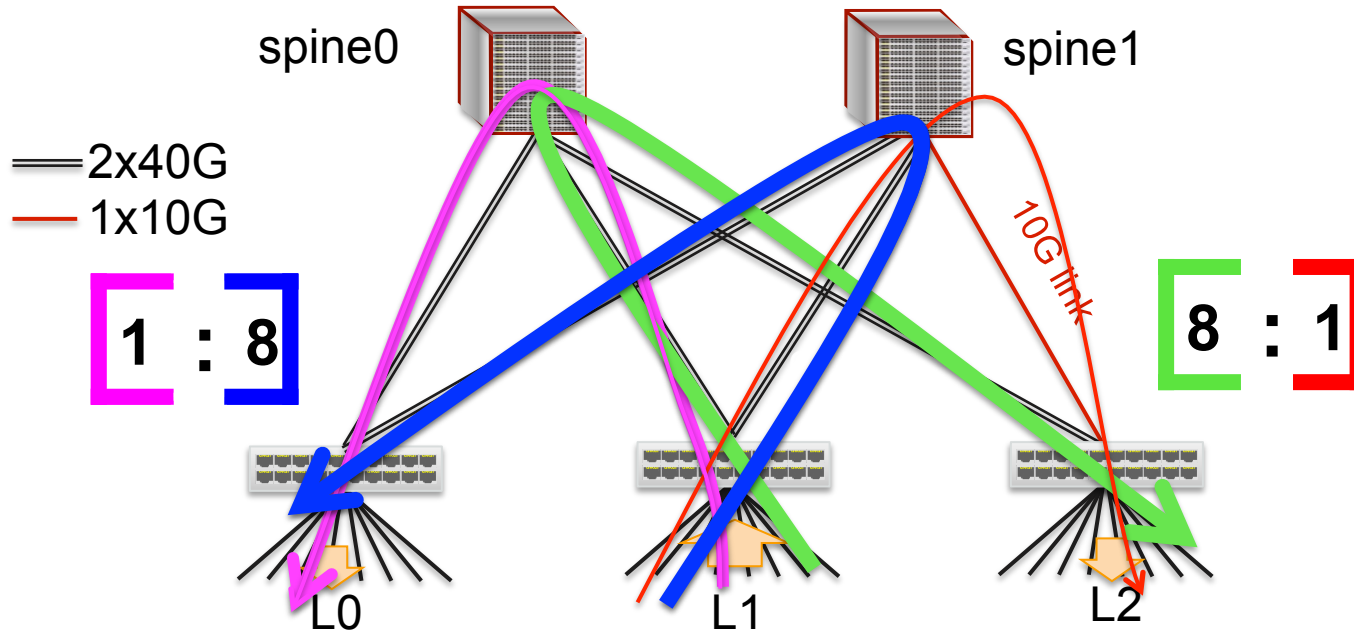


LetFlow within 2X of CONGA;  
Both are much better than other schemes



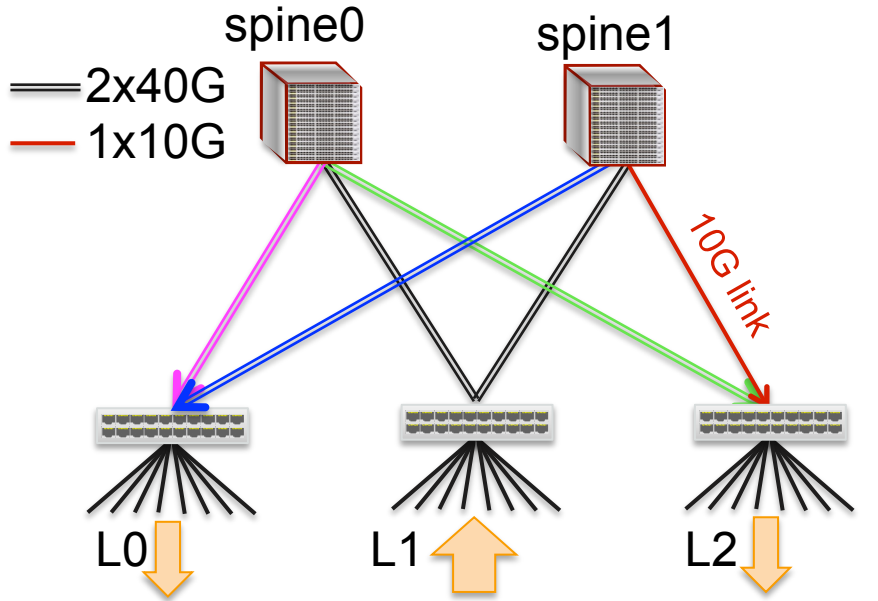
# Multi Destination Scenario

---



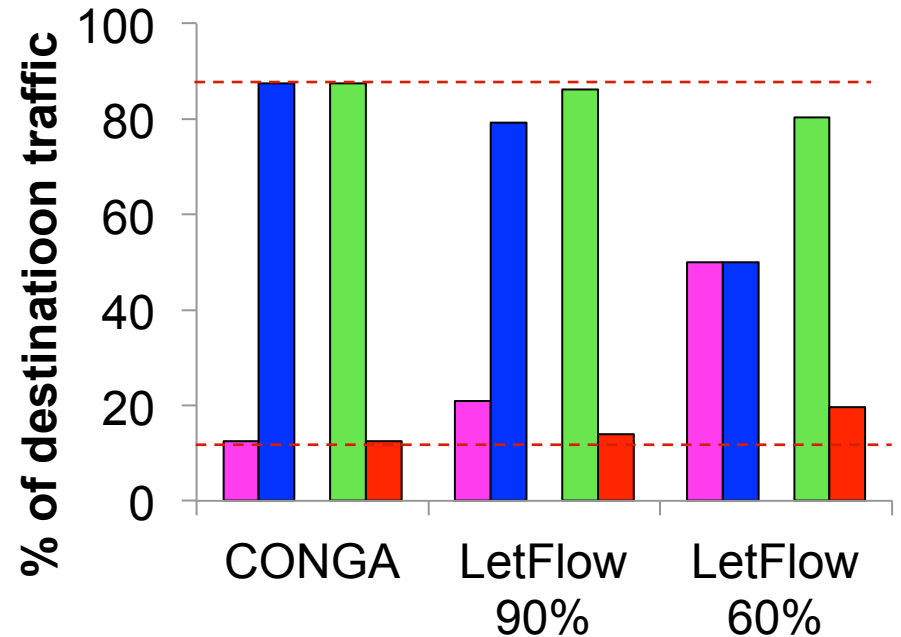
Traffic Load uniform to the 2 destinations

# Multi Destination Scenario

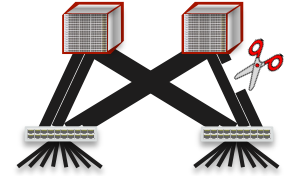


Traffic Load uniform to the 2 destinations

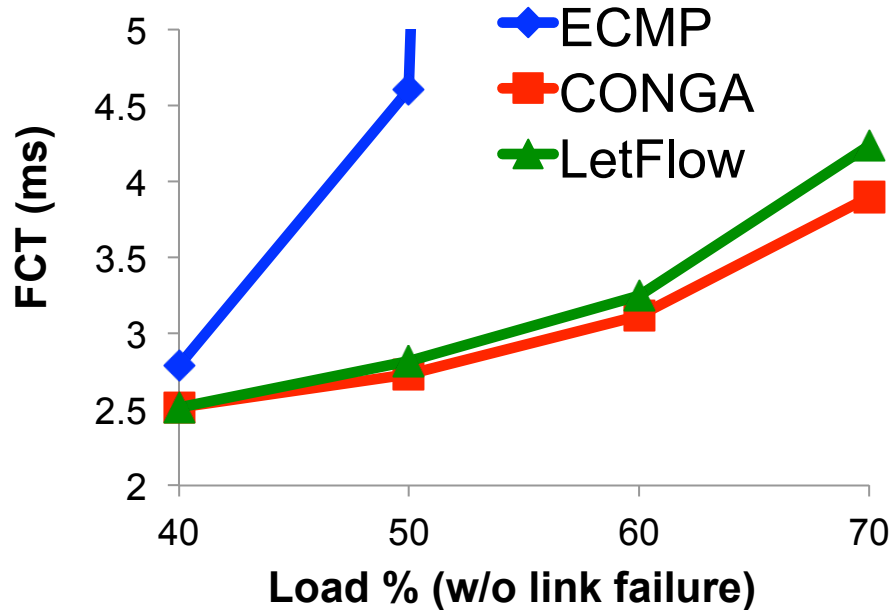
FCT is similar between Conga and LetFlow



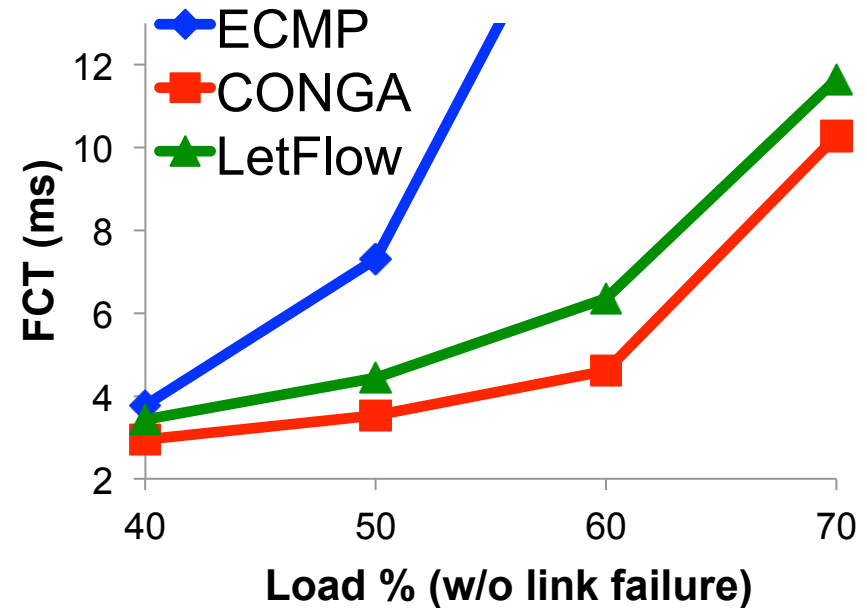
# Other Transport Protocols



## DCTCP



## DCQCN



# Conclusion

---

- Flowlet switching is a powerful technique for asymmetric load balancing
- LetFlow: a simple LB mechanism that handles asymmetry
  - Random decisions but implicitly congestion-aware
  - Suitable for standalone switches – does not need feedback
- Letflow is stochastic and reactive in nature
  - Cannot proactively prevent congestion / queue buildup like more sophisticated schemes

**LET it FLOW !**