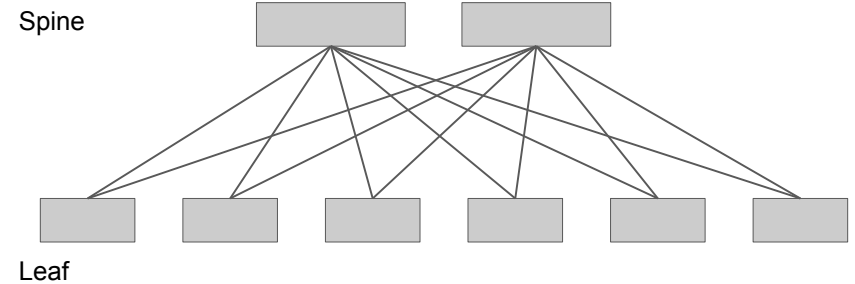# CONGA: Distributed Congestion-Aware Load Balancing for Datacenters

By Alizadeh,M et al.

Presented by Andrew and Jack
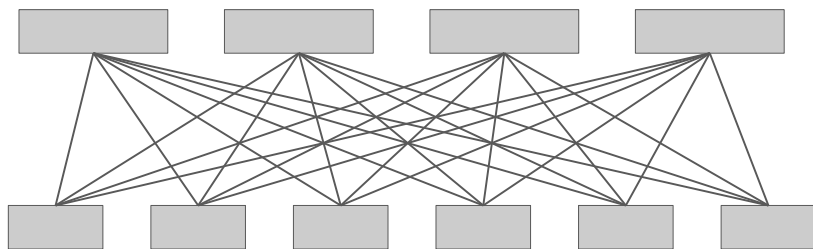
---

## Motivation

Distributed datacenter applications require large bisection bandwidth
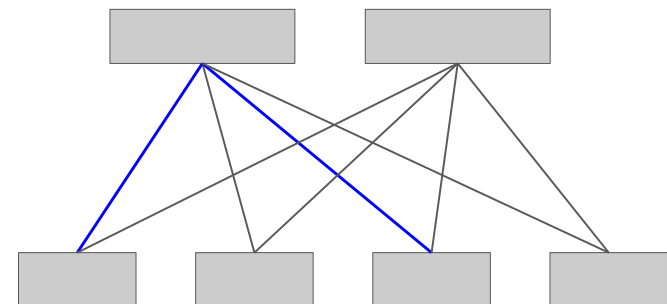
Spine

Leaf



---

## Motivation

Distributed datacenter applications require large bisection bandwidth



---

## ECMP: hash-based hop selection without reordering
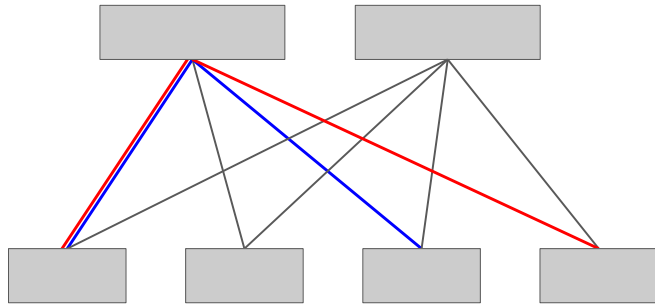
F(sIP, sPort, dIP, dPort, prot) = 0

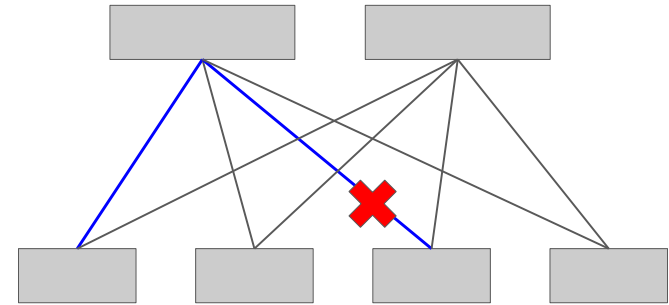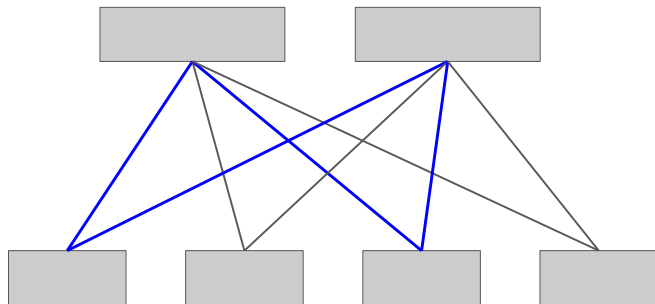## ECMP Problem: hash collisions lead to imbalance

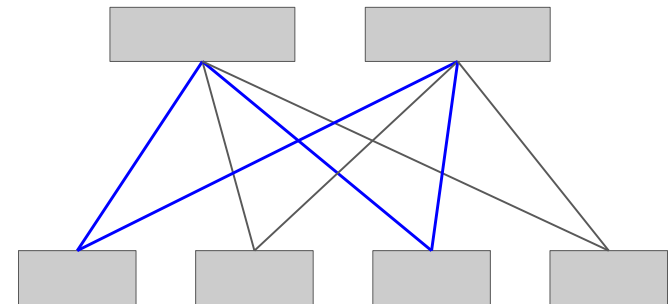F(0, 0, 3, 0, TCP) = 0
F(0, 1, 4, 0, TCP) = 0

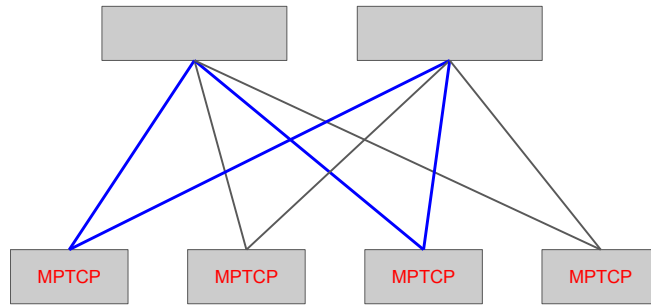## ECMP Problem: local decisions oblivious to downstream asymmetry

## MPTCP: split flows into sub-flows

## MPTCP Problem: higher congestion at edge

# MPTCP Problem: transport layer-specific



# CONGA: Congestion Aware Balancing

Network load-balancing without transport layer interference

Make globally optimal load-balancing decisions

Use common datacenter network features (e.g., overlay networks)
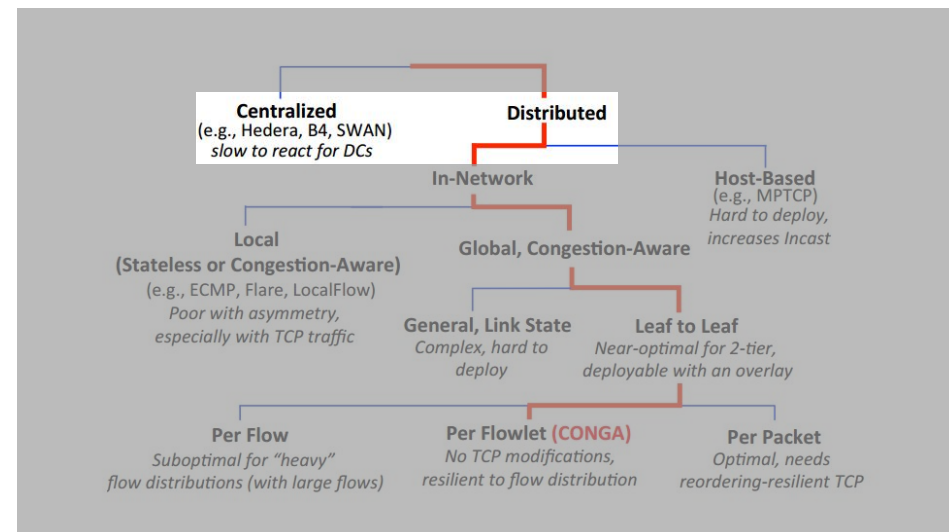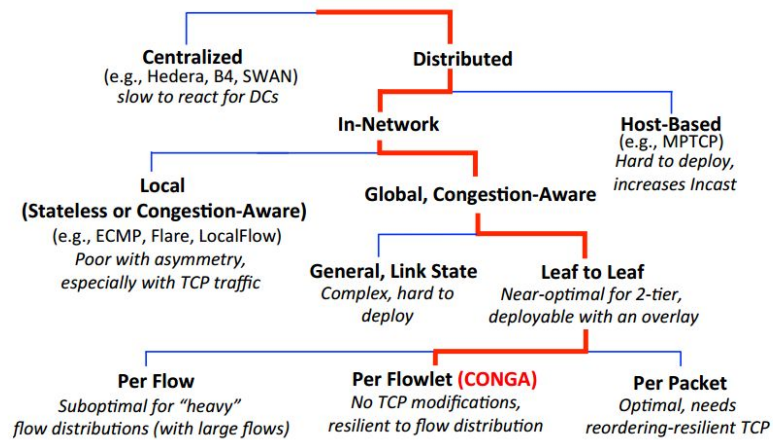
# CONGA Overview

Track end-to-end congestion along path

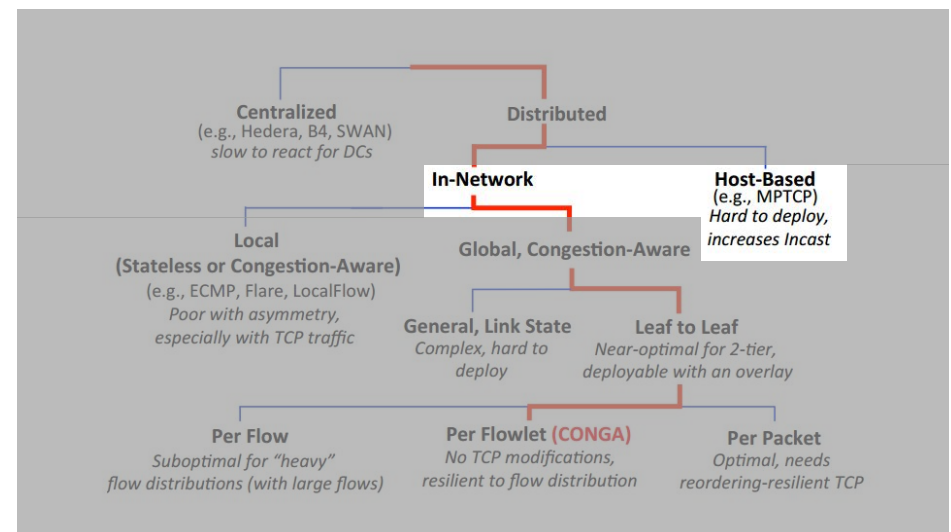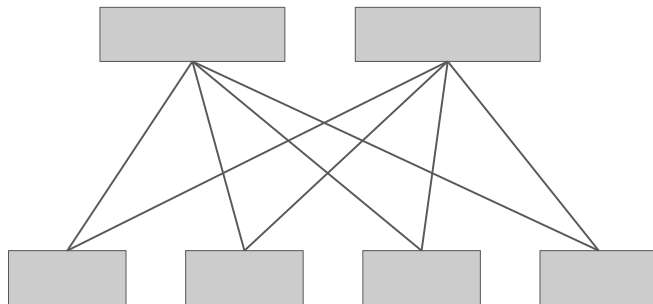Feedback loop between leaf switches: relay congestion information

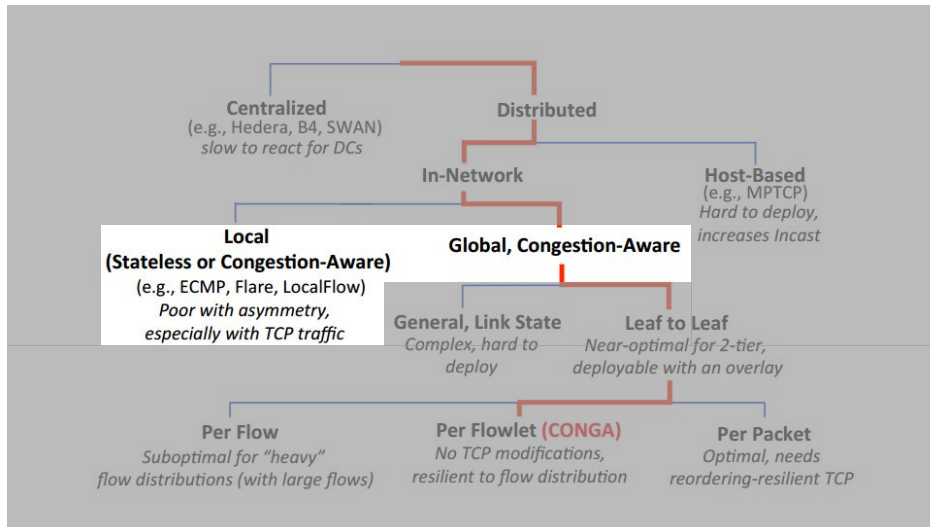Leaf switches send traffic on least congested path

# CONGA Design Goals

1. Responsive
2. Transport independent
3. Robust to asymmetry
4. Incrementally deployable
5. Optimized for Leaf-Spine

**Centralized**
(e.g., Hedera, B4, SWAN)
*slow to react for DCs*

**Distributed**

**In-Network**

**Host-Based**
(e.g., MPTCP)
*Hard to deploy,
increases Incast*

**Local**
**(Stateless or Congestion-Aware)**
(e.g., ECMP, Flare, LocalFlow)
*Poor with asymmetry,
especially with TCP traffic*

**Global, Congestion-Aware**

**General, Link State**
*Complex, hard to
deploy*

**Leaf to Leaf**
*Near-optimal for 2-tier,
deployable with an overlay*

**Per Flow**
*Suboptimal for "heavy"
flow distributions (with large flows)*

**Per Flowlet (CONGA)**
*No TCP modifications,
resilient to flow distribution*

**Per Packet**
*Optimal, needs
reordering-resilient TCP*

---

Distributed load-balancing is highly responsive, near optimal for regular topologies

Centralized
(e.g., Hedera, B4, SWAN)
*slow to react for DCs*

Distributed

In-Network

Host-Based
(e.g., MPTCP)
*Hard to deploy,
increases Incast*

Local
(Stateless or Congestion-Aware)
(e.g., ECMP, Flare, LocalFlow)
*Poor with asymmetry,
especially with TCP traffic*

Global, Congestion-Aware

General, Link State
*Complex, hard to
deploy*

Leaf to Leaf
*Near-optimal for 2-tier,
deployable with an overlay*

Per Flow
*Suboptimal for "heavy"
flow distributions (with large flows)*

Per Flowlet (CONGA)
*No TCP modifications,
resilient to flow distribution*

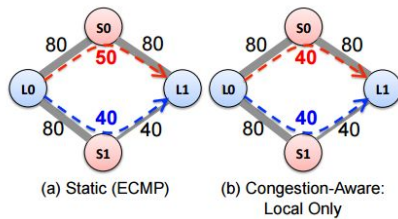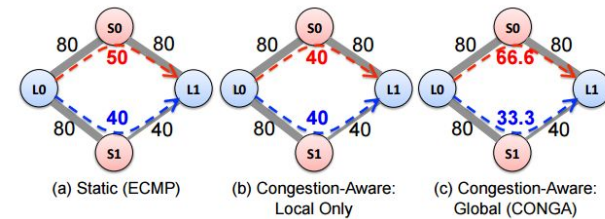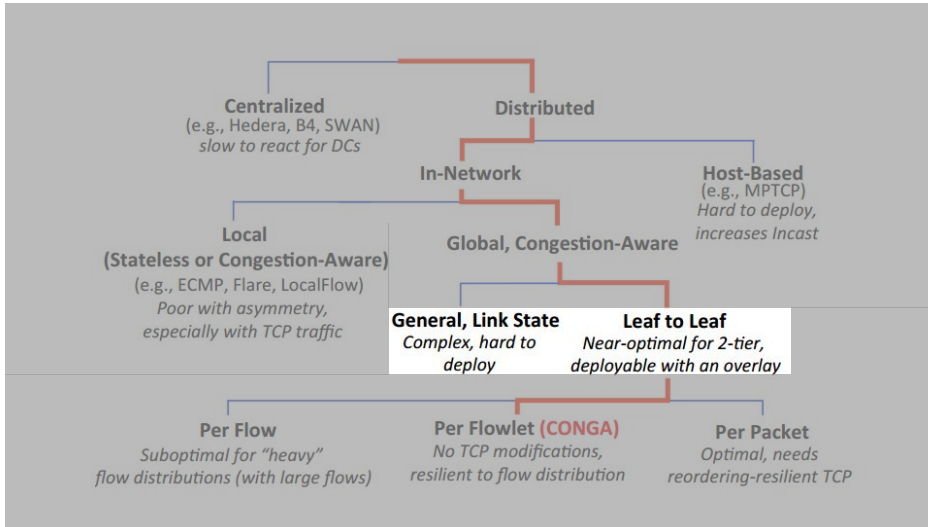Per Packet
*Optimal, needs
reordering-resilient TCP*

Global congestion awareness is necessary to handle network asymmetry



(a) Static (ECMP)

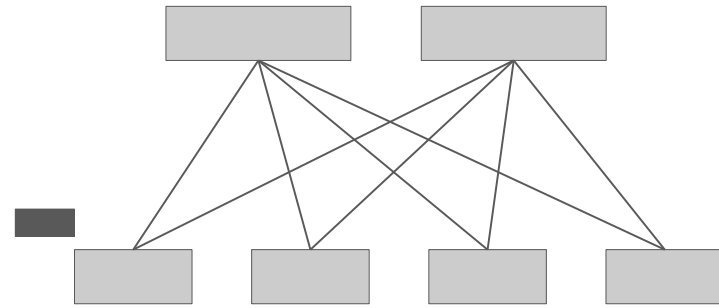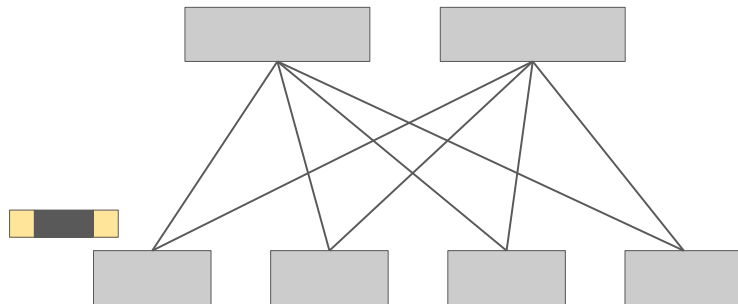Global congestion awareness is necessary to handle network asymmetry



(a) Static (ECMP)    (b) Congestion-Aware:
Local Only

Global congestion awareness is necessary to handle network asymmetry



(a) Static (ECMP)    (b) Congestion-Aware:
Local Only    (c) Congestion-Aware:
Global (CONGA)

Centralized
(e.g., Hedera, B4, SWAN)
*slow to react for DCs*

Distributed

In-Network

Host-Based
(e.g., MPTCP)
*Hard to deploy,
increases Incast*

Local
(Stateless or Congestion-Aware)
(e.g., ECMP, Flare, LocalFlow)
*Poor with asymmetry,
especially with TCP traffic*

Global, Congestion-Aware

**General, Link State**
*Complex, hard to
deploy*

**Leaf to Leaf**
*Near-optimal for 2-tier,
deployable with an overlay*

**Per Flow**
*Suboptimal for "heavy"
flow distributions (with large flows)*

**Per Flowlet (CONGA)**
*No TCP modifications,
resilient to flow distribution*

**Per Packet**
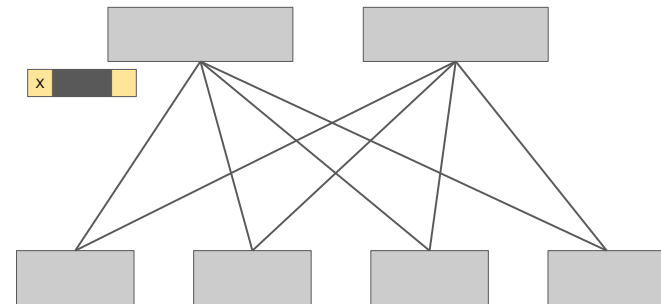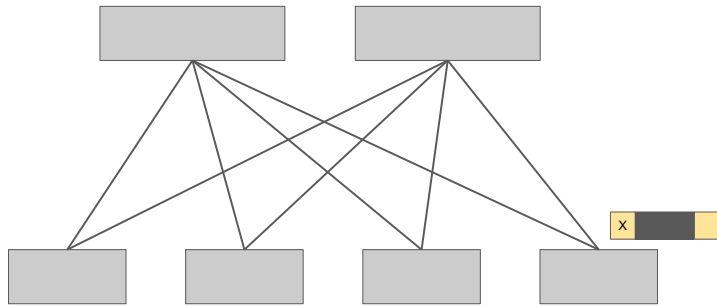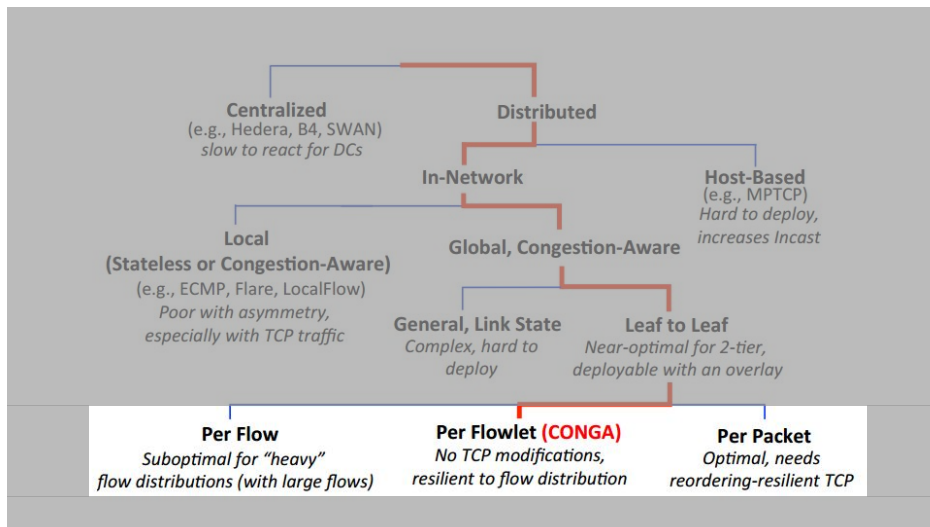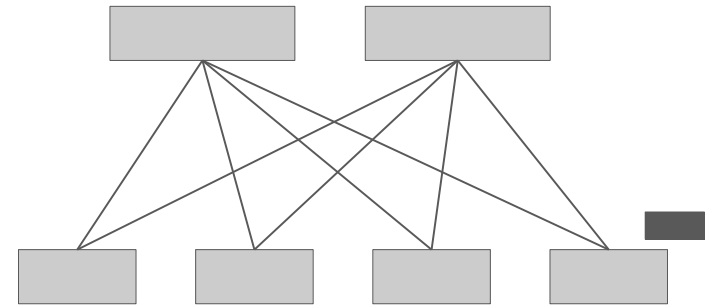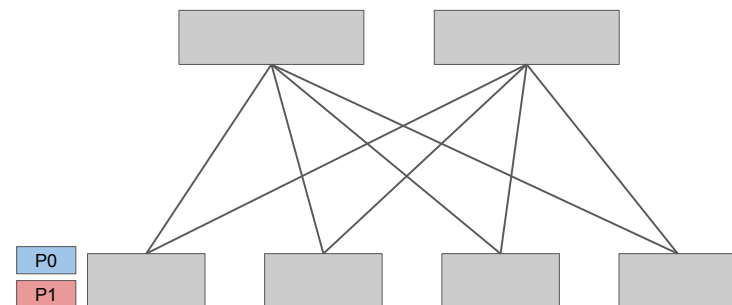*Optimal, needs
reordering-resilient TCP*

Overlay networks allow leaf switches to know destination leaf and carry congestion metrics



Overlay networks allow leaf switches to know destination leaf and carry congestion metrics



Overlay networks allow leaf switches to know destination leaf and carry congestion metrics

Overlay networks allow leaf switches to know destination leaf and carry congestion metrics
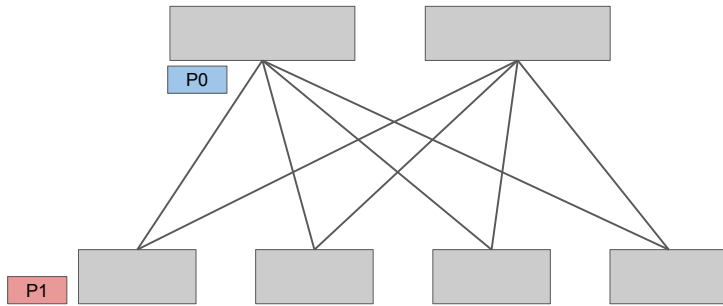


Overlay networks allow leaf switches to know destination leaf and carry congestion metrics
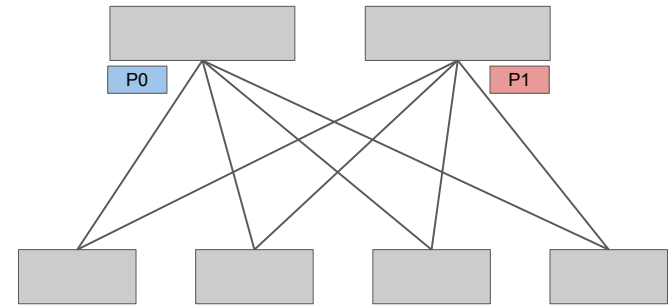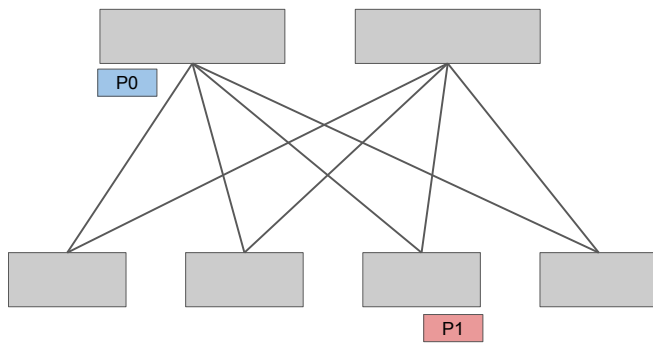




Packet-granularity scheduling can result in reordering → modifications to end-host TCP

Packet-granularity scheduling can result in
reordering → modifications to end-host TCP
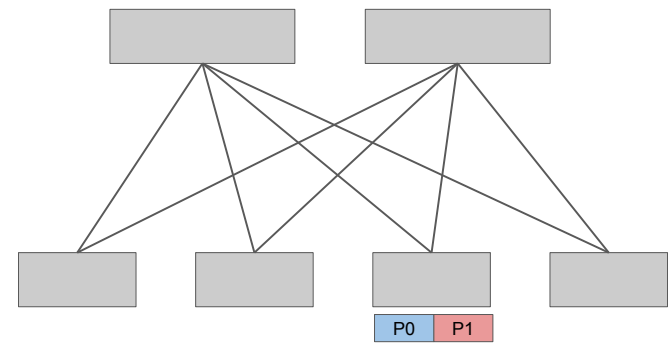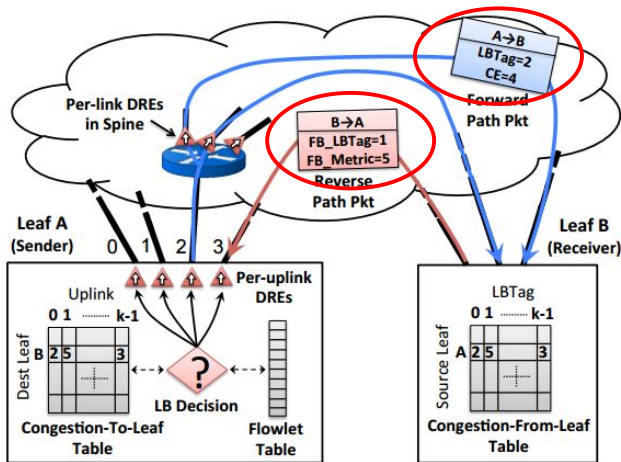
Packet-granularity scheduling can result in
reordering → modifications to end-host TCP

Packet-granularity scheduling can result in
reordering → modifications to end-host TCP

Packet-granularity scheduling can result in
reordering → modifications to end-host TCP
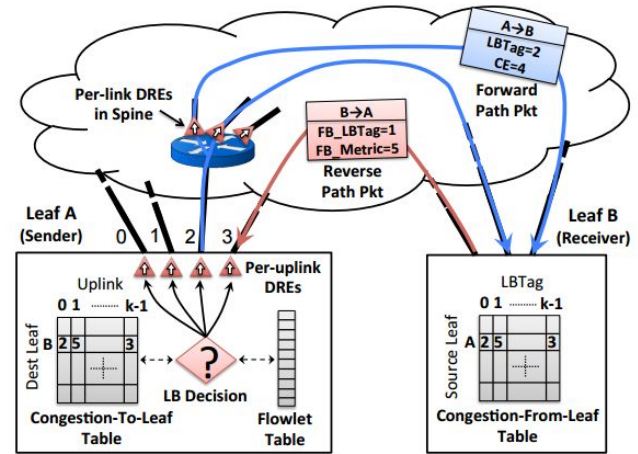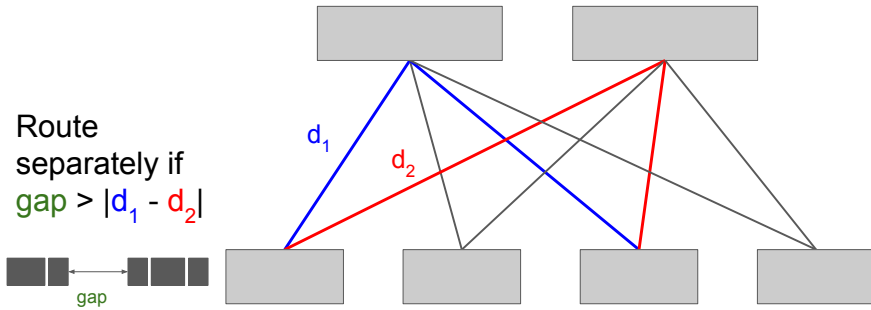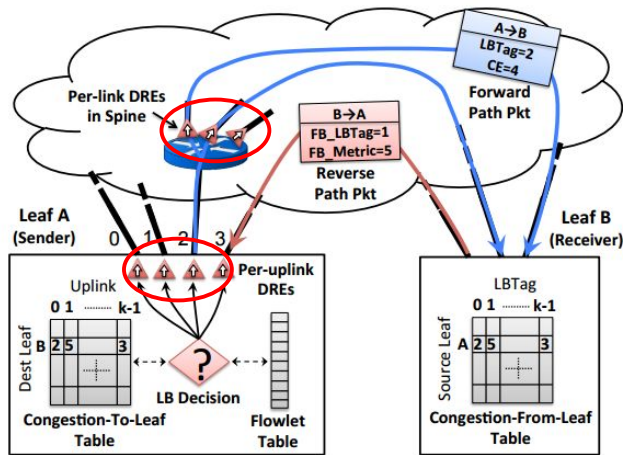
# Flowlet: break apart flow based on delayed bursts
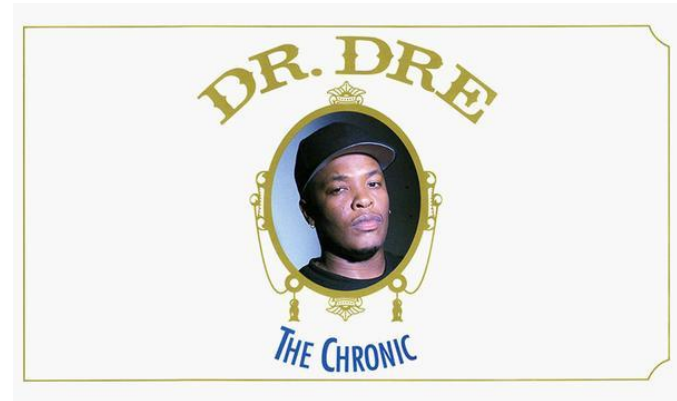
Route
separately if
$gap > |d_1 - d_2|$



$d_1$ $d_2$

gap



A→B
LBTag=2
CE=4

Forward
Path Pkt

Per-link DREs
in Spine

B→A
FB_LBTag=1
FB_Metric=5

Reverse
Path Pkt

Leaf A
(Sender)   0 1 2 3

Leaf B
(Receiver)

Uplink      Per-uplink
0 1 ....... k-1   DREs

Dest Leaf

B 2 5      3

LB Decision

Congestion-To-Leaf
Table

Flowlet
Table

LBTag
0 1 ....... k-1

Source Leaf

A 2 5      3

Congestion-From-Leaf
Table



A→B
LBTag=2
CE=4

Forward
Path Pkt

Per-link DREs
in Spine

B→A
FB_LBTag=1
FB_Metric=5

Reverse
Path Pkt

Leaf A
(Sender)   0 1 2 3

Leaf B
(Receiver)

Uplink      Per-uplink
0 1 ....... k-1   DREs

Dest Leaf

B 2 5      3

LB Decision

Congestion-To-Leaf
Table

Flowlet
Table

LBTag
0 1 ....... k-1

Source Leaf

A 2 5      3

Congestion-From-Leaf
Table

# Overlay packet header contains CONGA metadata

LBTag: src port
CE: path congestion

FB_LBTag: fb dst port
FB_Metric: reported cong.



x

## Discounting Rate Estimator (DRE)



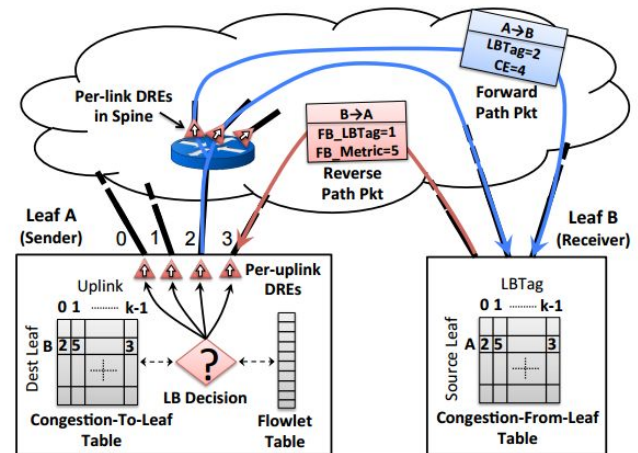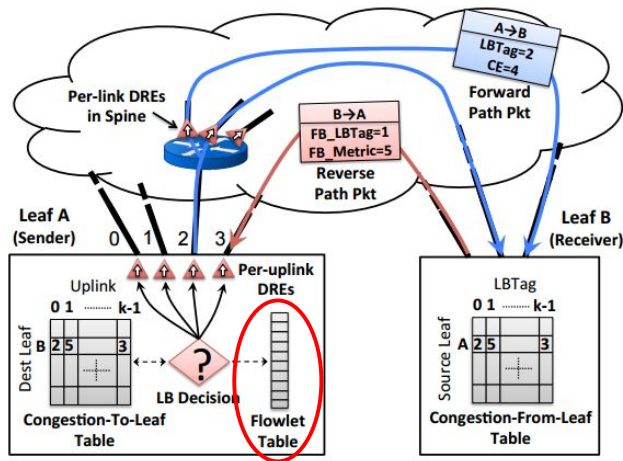## Discounting Rate Estimator (DRE)

*X:* register quantifying load
  Additive increase by bytes sent for each packet
  Multiplicative decrease every $T_{dre}$ by *α*

$$X \leftarrow X * (1 - α)$$

More responsive to traffic bursts than EWMA

Flowlet Detection

$T_{fl}$ : flowlet inactivity gap

Hash flowlets based on 5-tuple
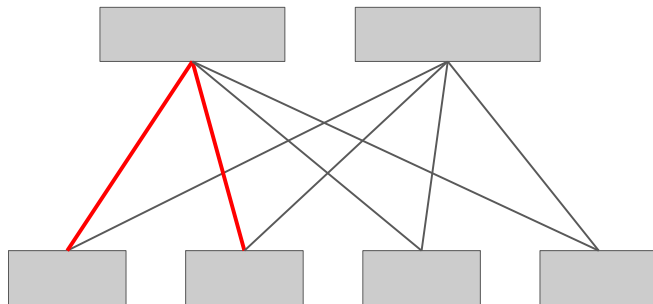    Collision is not a correctness issue

Round-based aging
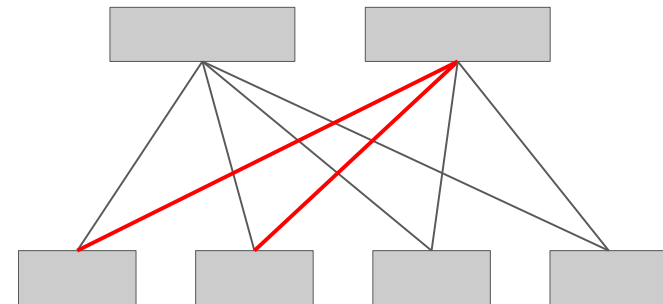
LB decisions made based on first packet

New flowlet: choose uplink minimizing the max local metric

Port:
Valid:
Age:

Implementation: custom ASICs rather than software to reduce overreaction, oscillations



Implementation: custom ASICs rather than software to reduce overreaction, oscillations

# Evaluation

1. How does CONGA impact flow completion times (FCT) vs. state of the art?
2. How does CONGA perform under the impact of failed links?
3. Does CONGA perform well on real-world traffic?

# Experimental Setup

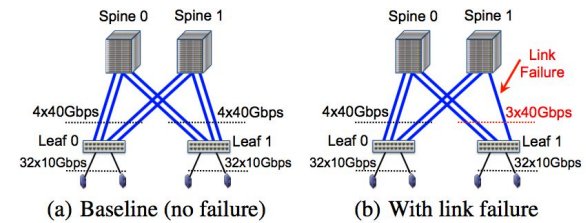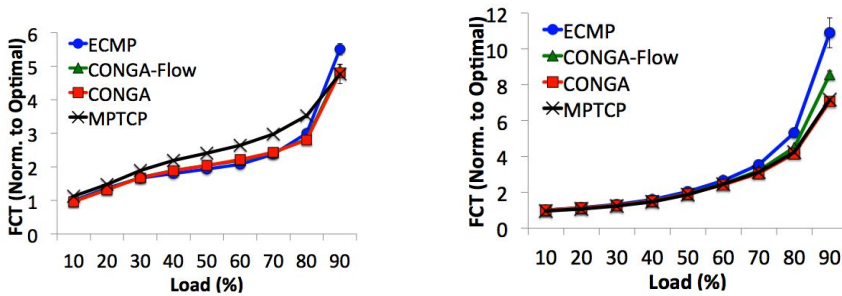Compared CONGA, CONGA-FLOW, ECMP and MPTCP



(a) Baseline (no failure)    (b) With link failure

**Figure 7: Topologies used in testbed experiments.**
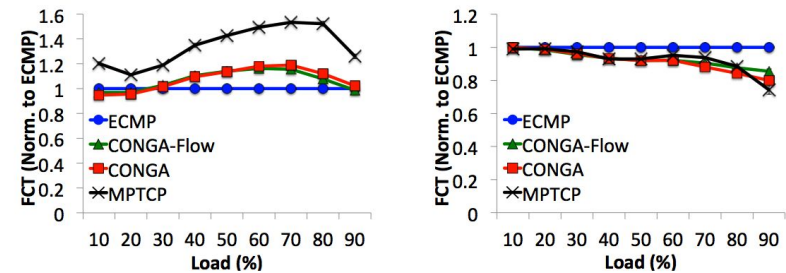
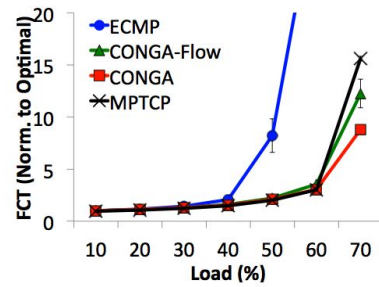# Baseline Performance

Two Workloads: Enterprise and Data-mining



# Baseline Performance

Breakdown: Short Flows and Long Flows
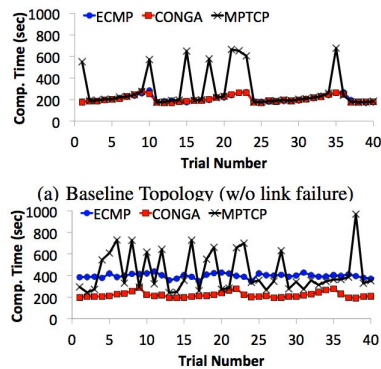
## Link Failure





## Incast



(a) MTU = 1500B



## HDFS Benchmark



(a) Baseline Topology (w/o link failure)



## Analytical Evaluation

Worst-Case performance: The ratio between the most congested link in CONGA and the best possible assignment of flows is 2.

## Analytical Evaluation

What is the expected traffic imbalance?

How does it depend on workload?

$$\mathrm{E}(\chi(t)) \leq \frac{1}{\sqrt{\lambda_e t}} + O(\frac{1}{t}),$$

*where:*

$$\lambda_e = \frac{\lambda}{8n \log n \left(1 + (\frac{\sigma_S}{\mathrm{E}(S)})^2\right)}.$$

## Analytical Evaluation

What is the expected traffic imbalance?

How does it depend on workload?

Less imbalance with many small flows, more imbalance with fewer large flows

$$\mathrm{E}(\chi(t)) \leq \frac{1}{\sqrt{\lambda_e t}} + O($$

*where:*

$$\lambda_e = \frac{\lambda}{8n \log n \left(1 + (\frac{\sigma_S}{\mathrm{E}(S)})^2\right)}.$$

## Conclusion

CONGA: globally aware datacenter load balancing
  - No transport layer intervention

Implemented in custom ASICs

Better flow completion times than ECMP, Incast MPTCP

## Discussion

Leaf-Spine topology has each leaf only two hops apart
  - Significant performance degrade if implemented in software?
  - Extensible to larger, multi-layered topologies?