

# In-band Network Telemetry: Network Load Balancing Use-cases

Presenter: Mukesh Hira, VMware

Slide content credits: Co-authors of HULA and CLOVE papers

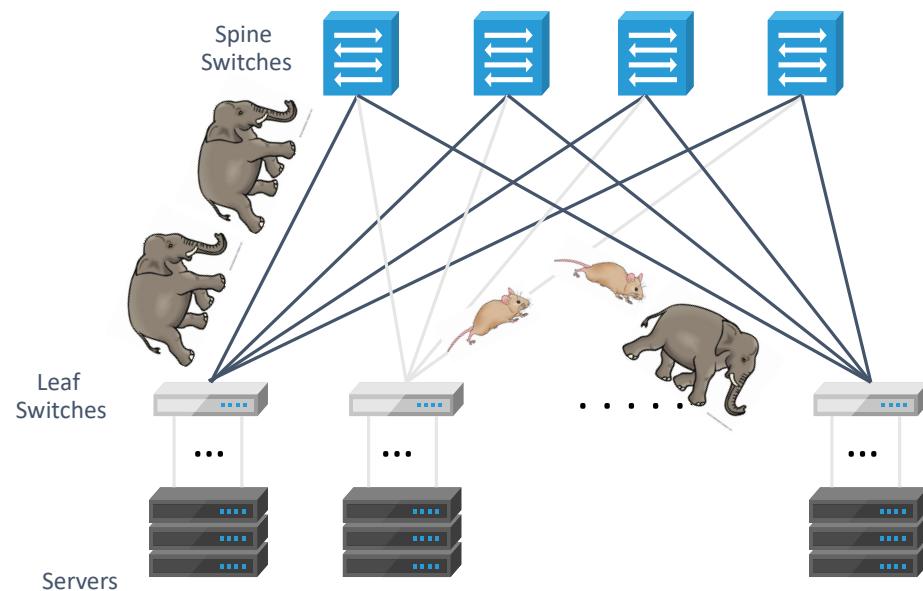
Naga Katta (Salesforce.com), Aditi Ghag (Vmware), Aran Bergmann (Technion), Isaac Keslassy (Technion),  
Changhoon Kim (Barefoot Networks), Anirudh Sivaraman (NYU), Jennifer Rexford (Princeton University)

# Predominant Data center Network Load Balancing

**Equal-Cost Multi-Path (ECMP) routing:**

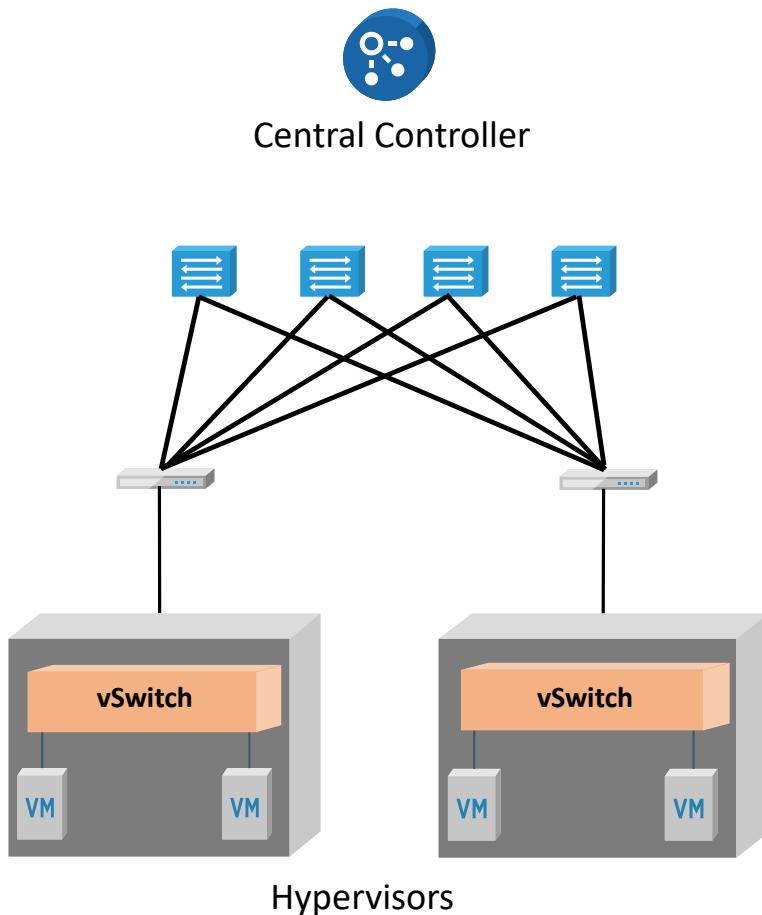
- $\text{Path(packet)} = \text{hash(packet's 5-tuple)}$

src,dst IP + src,dst port +  
protocol



- Coarse-grained
- Elephant Hash collisions
- Congestion-oblivious

# Alternative Proposals

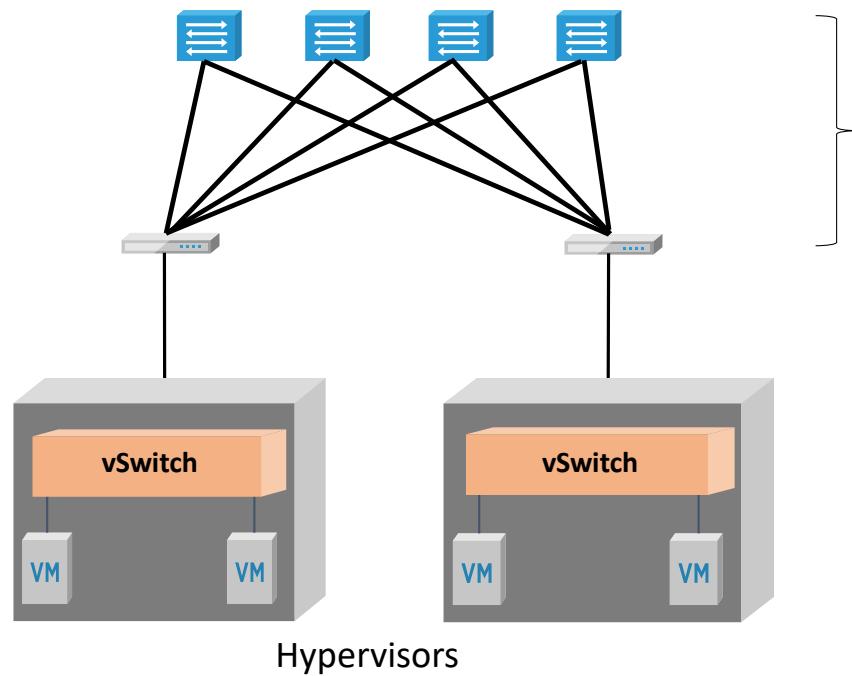


## Centralized Load Balancing

Hedera, Fastpass, MicroTE, SWAN

- Control-driven feedback
- Slow reaction time
- Scalability Issues

# Alternative Proposals

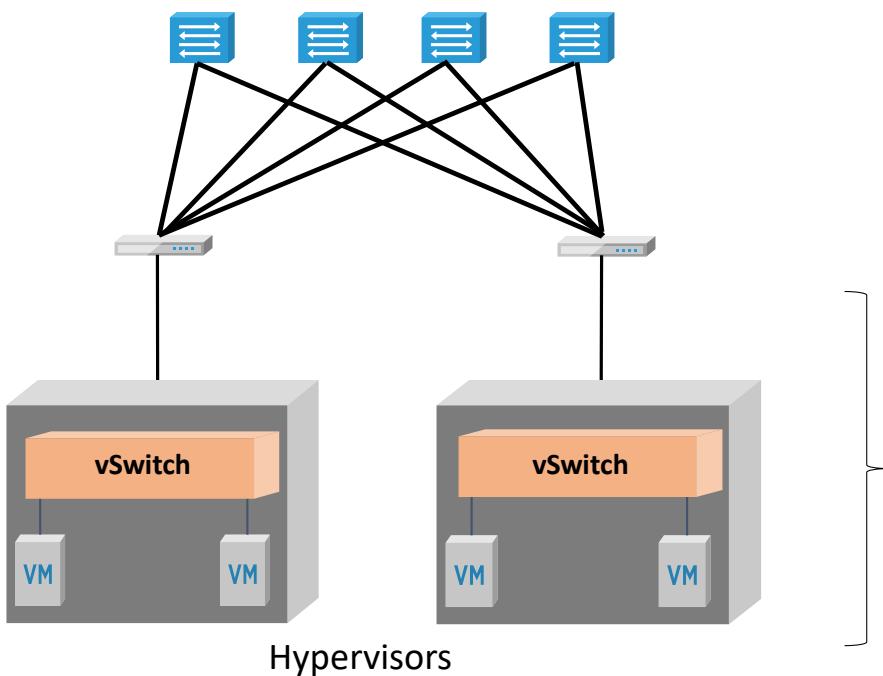


## In-network load balancing

CONGA, LetFlow, DRILL

- Needs custom ASIC data center fabric
- High capital cost
- Controller Involvement may still be required

# Alternative Proposals



## **End-host load balancing**

Presto

- Congestion Oblivious
- Controller intervention in case of topology asymmetry

MPTCP

- Incast collapse
- Guest VM network stack changes

# Congestion Aware Network Load Balancing based on INT

- Two approaches
  - In-Network Network Load Balancing at each Hop  
HULA: Scalable Load Balancing using Programmable Data Planes  
(Hop-by-Hop Link Utilization Aware Load Balancing Architecture)
  - Congestion Aware Network Load Balancing from the Edge  
CLOVE: Congestion-aware Load Balancing at the Virtual Edge

# HULA: Scalable Load Balancing using Programmable Data Planes (Hop-by-Hop Link Utilization Aware Load Balancing Architecture)

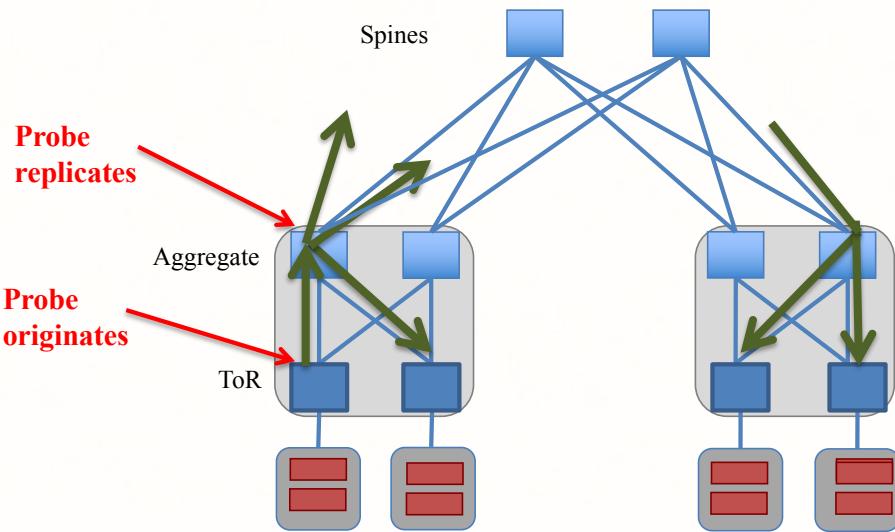
Naga Katta (Princeton University), Mukesh Hira (VMware), Changhoon Kim (Barefoot Networks),  
Anirudh Sivaraman (MIT), Jennifer Rexford (Princeton University)

ACM SOSR 2016

# HULA: Key Benefits

- Works in any network topology
  - Scalability via summarization of state
  - Only propagate network state for best next-hop upstream
  - Scales to any number of network tiers, any number of destinations
    - Subject to forwarding table size limitations
- Leverages INT to gather network state
- No vendor-proprietary state exchange protocol

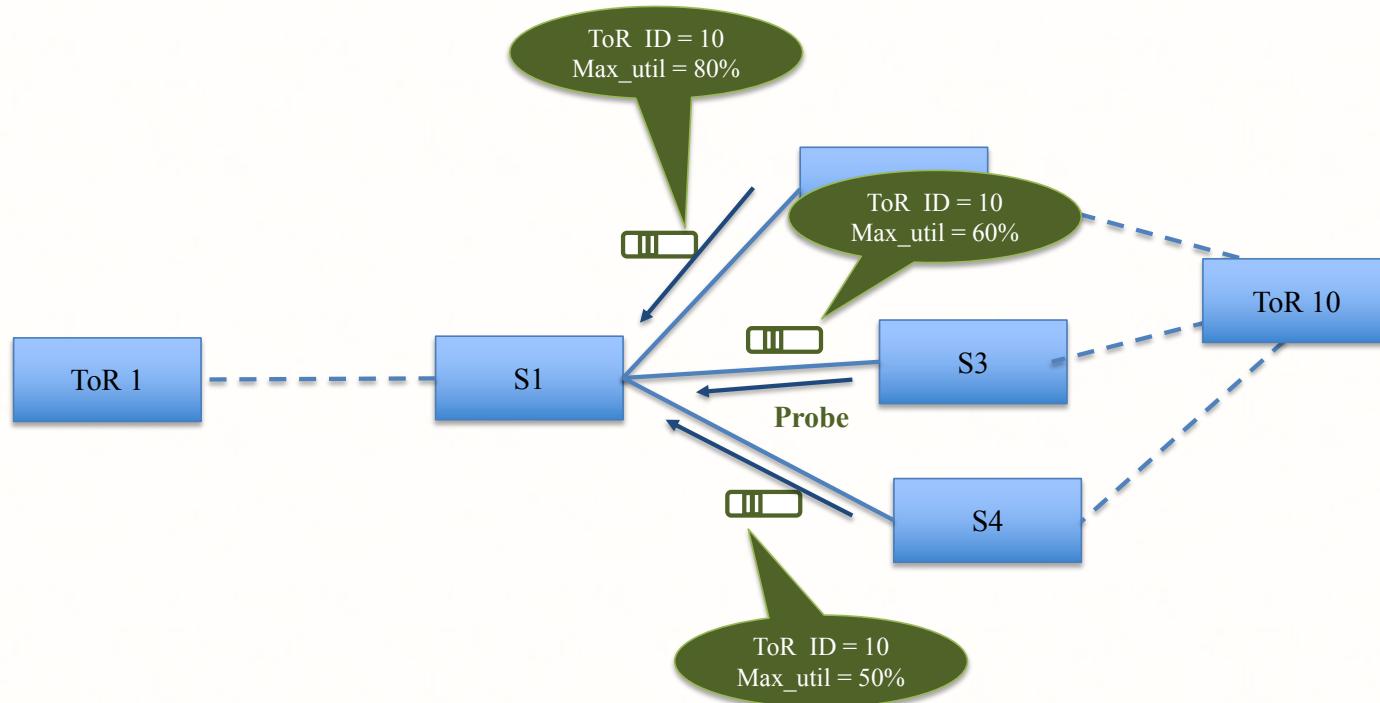
# 1. Probes accumulate path utilization in forward direction



P4 primitives
New header format
Programmable Parsing
Switch metadata

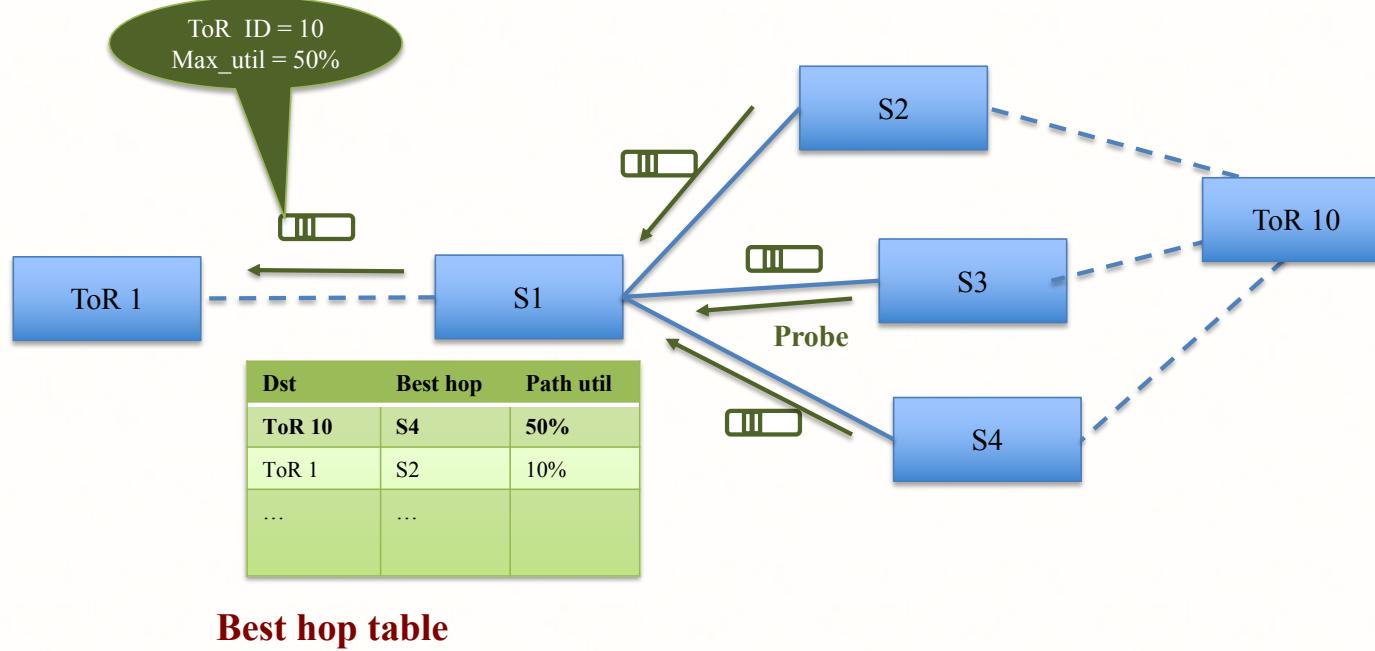
- Probes originate at each ToR
- Replicated at intermediate hops
- Path Utilization collected using INT

## 1. Forward direction path information is piggybacked onto reverse direction probes



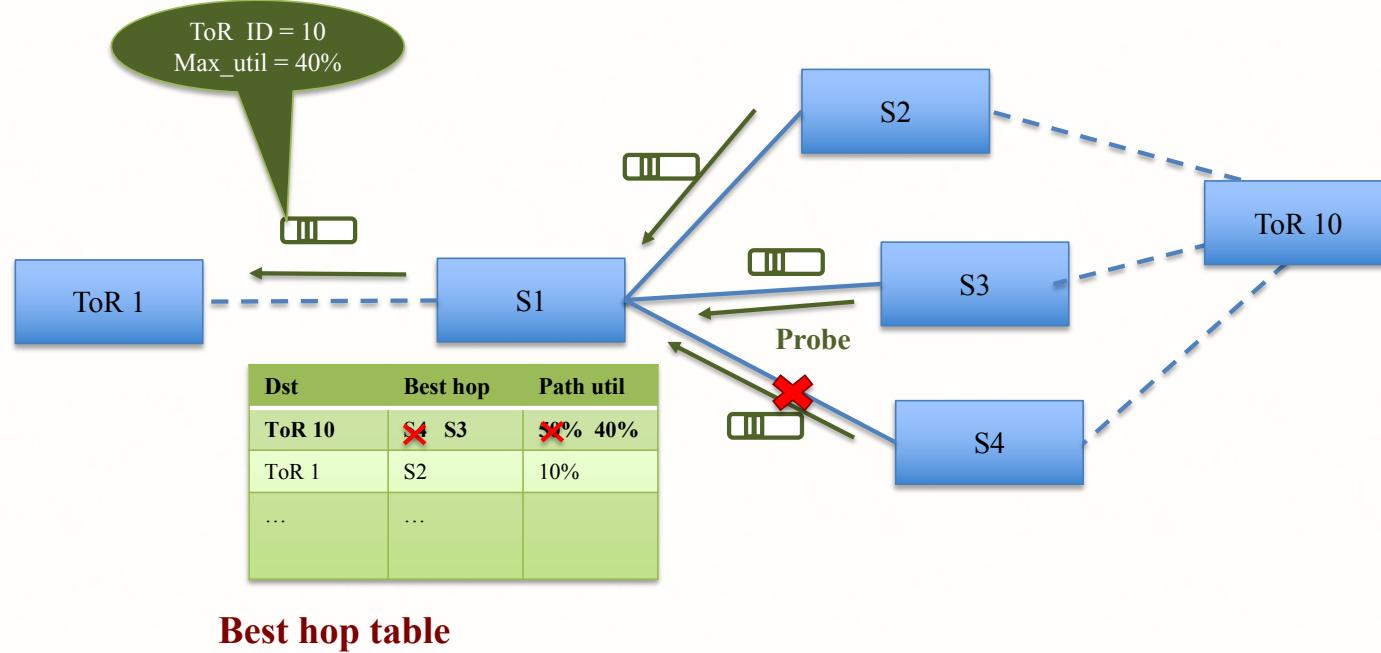
Probe carries INT instructions to collect path state in its direction and previously accumulated state for opposite direction from probes received at probe-originating ToR

## 2. Switch identifies best next-hop to each ToR

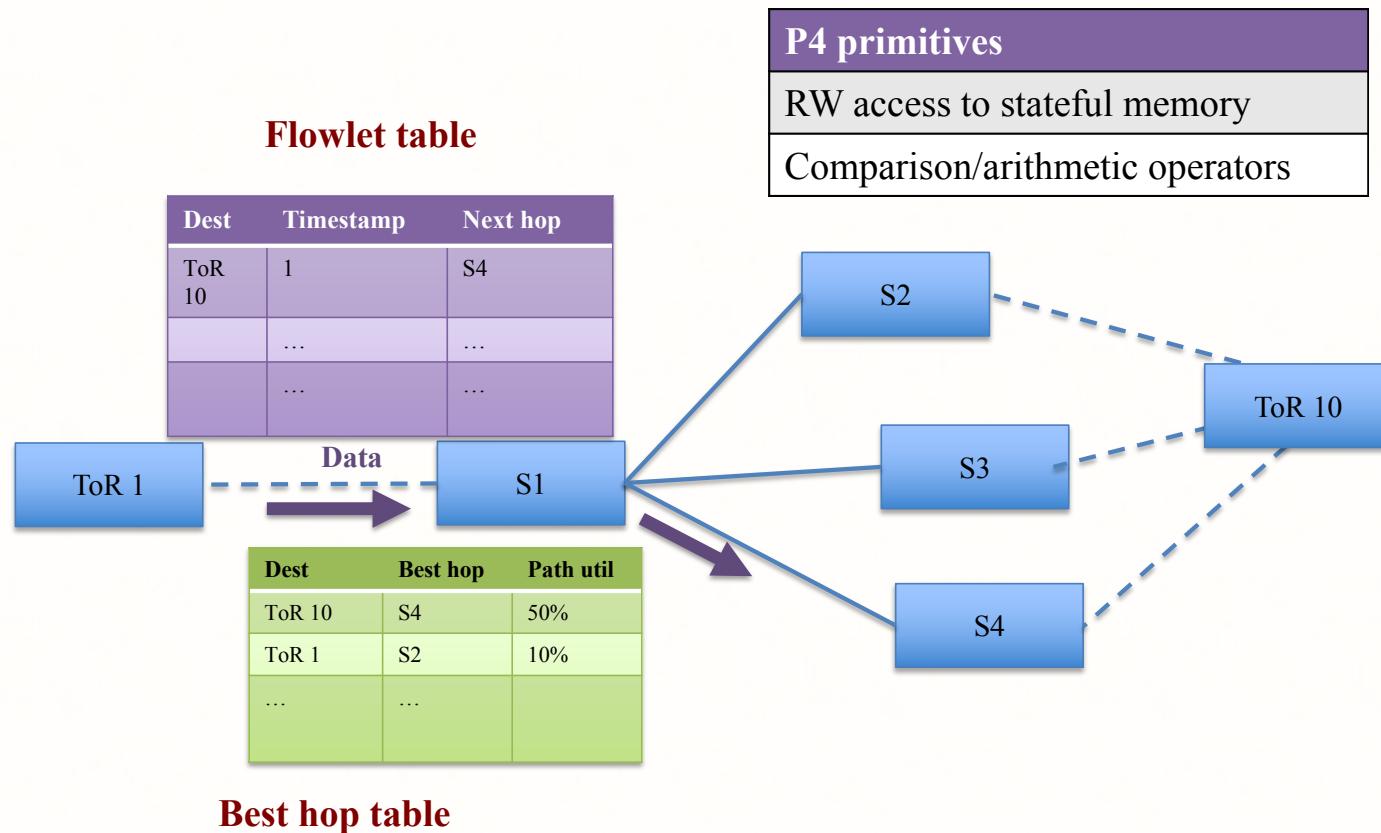


End-host to ToR mapping known via out of band mechanism

## 2. Switch identifies best next-hop to each ToR



### 3. Switches load balance flowlets based on real-time path utilization



## Evaluation Setup

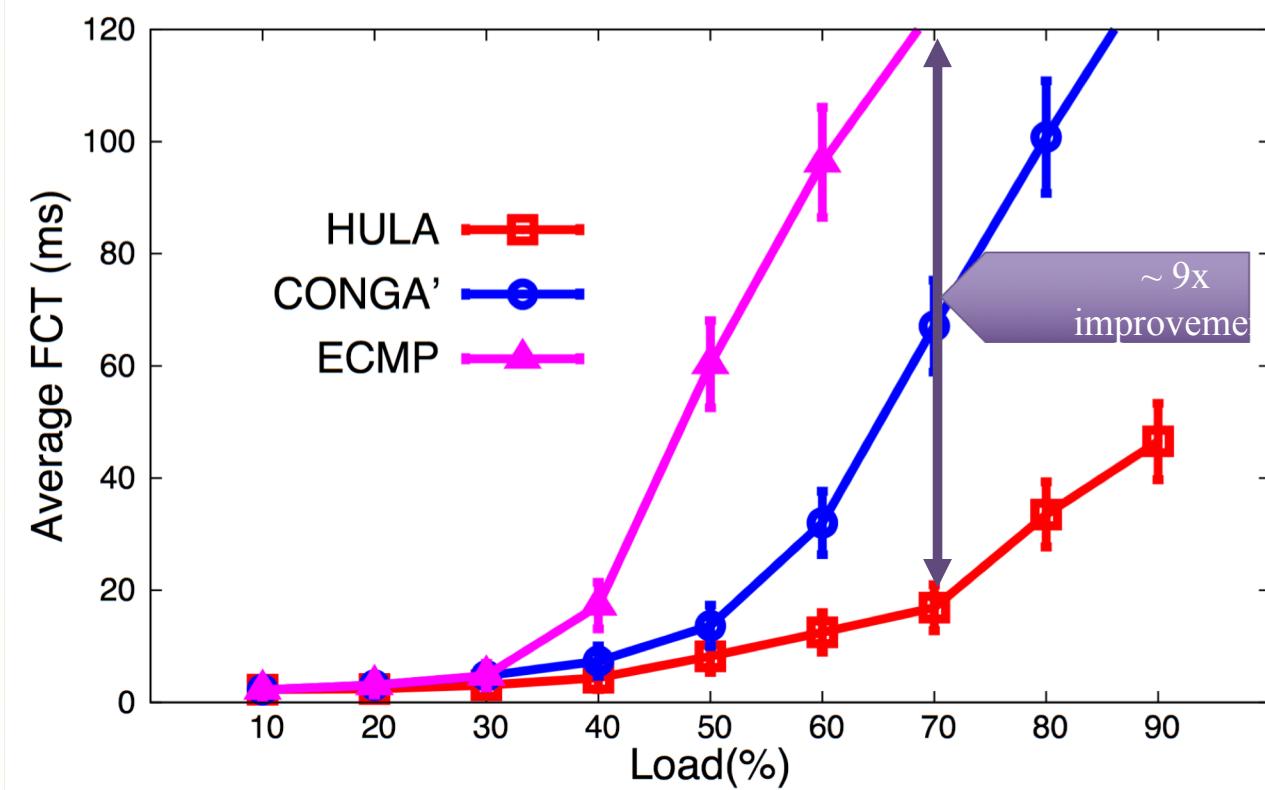
- NS2 packet-level simulator
- RPC-based workload generator
  - Empirical flow size distributions
  - Websearch and Datamining
- End-to-end metric
  - Average Flow Completion Time (FCT)

## Compared with

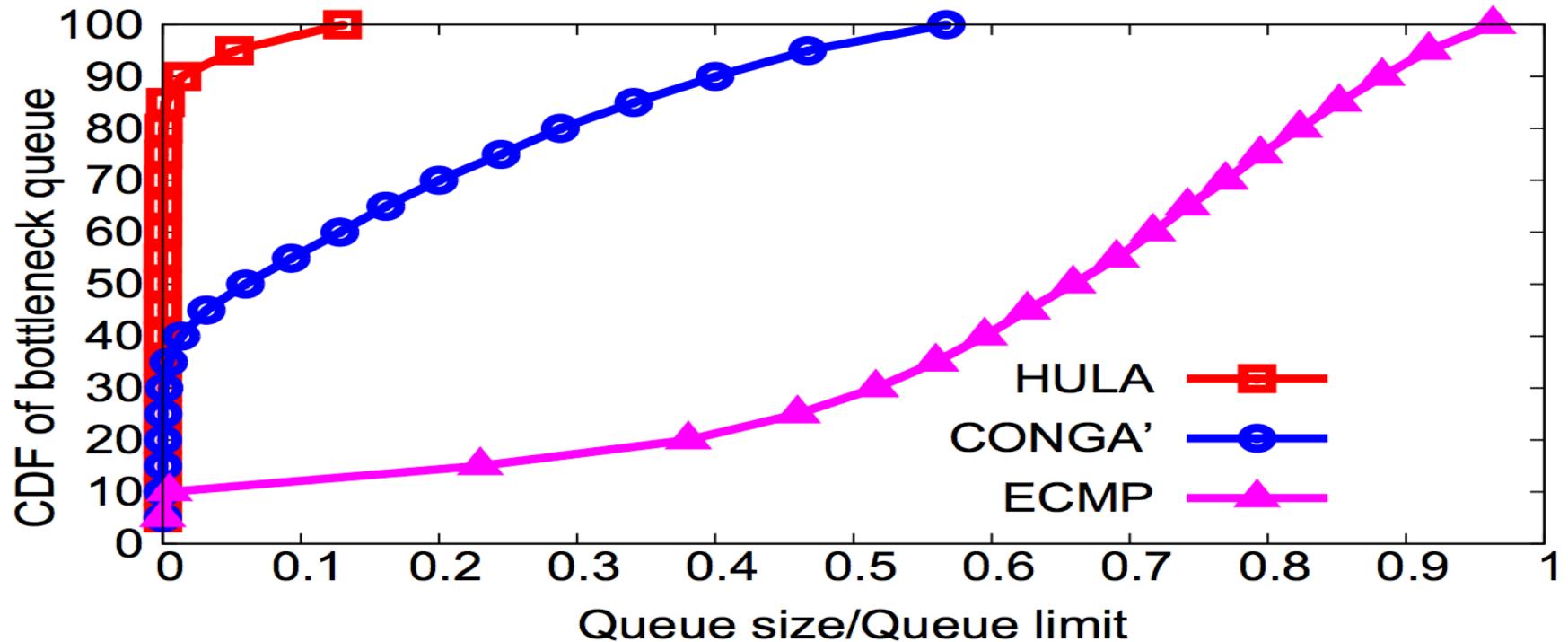
- ECMP
  - Flow level hashing at each switch
- CONGA'
  - CONGA within each leaf-spine pod
  - ECMP on flowlets for traffic across pods<sup>1</sup>

1. Based on communication with the authors

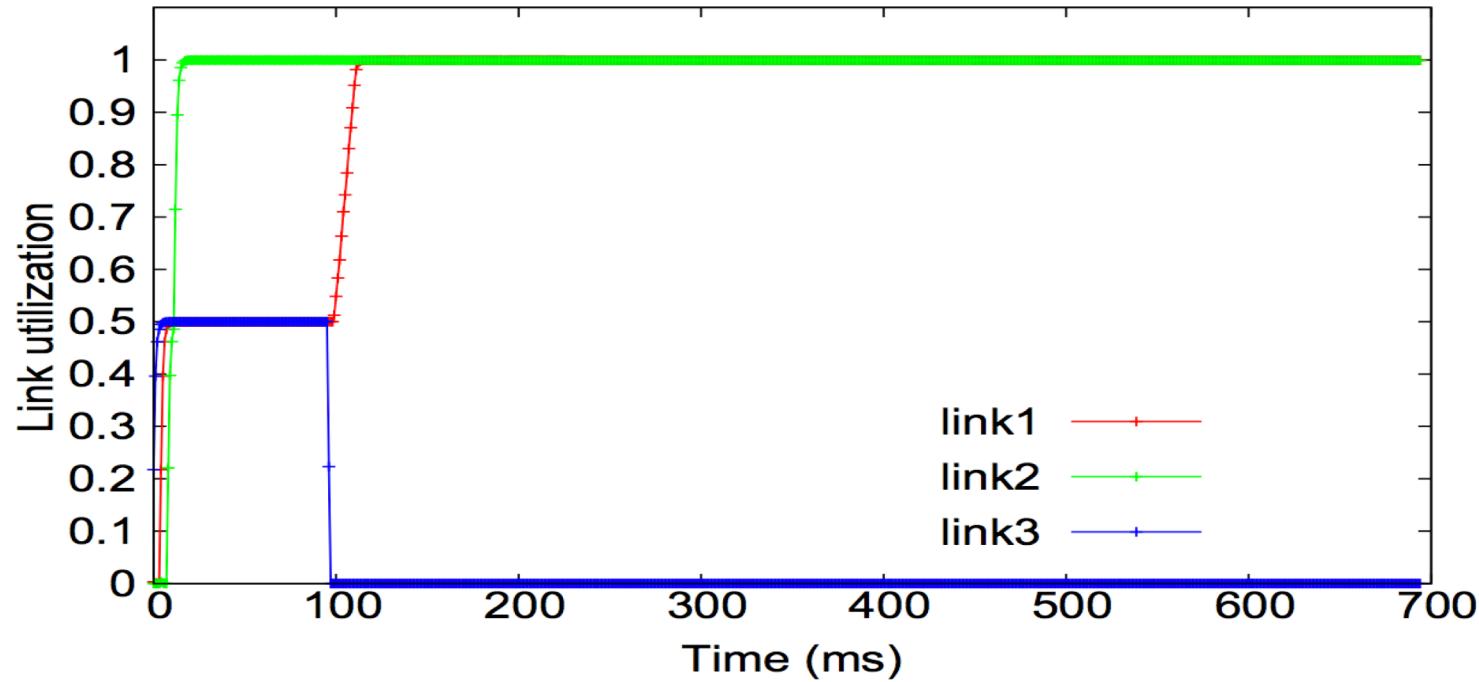
## HULA handles high load much better



## HULA keeps queue occupancy low



## HULA is stable on link failure



# **Clove: Congestion-Aware Load Balancing at the Virtual Edge**

Naga Katta<sup>1,2</sup>, Aditi Ghag<sup>3</sup>, Mukesh Hira<sup>3</sup>, Isaac Keslassy<sup>3,4</sup>,

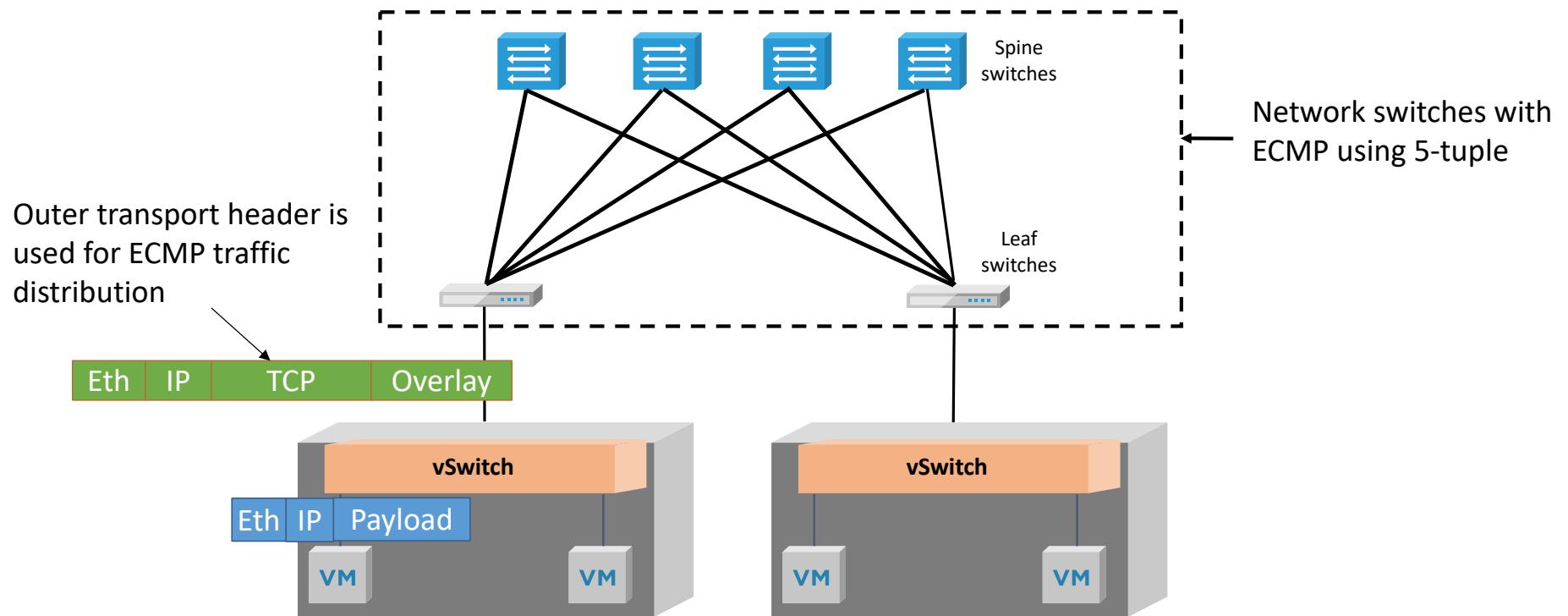
Aran Bergman<sup>3,4</sup>, Changhoon Kim<sup>5</sup>, Jennifer Rexford<sup>2</sup>

<sup>1</sup> *Salesforce.com*    <sup>2</sup> *Princeton University*    <sup>3</sup> *VMware*    <sup>4</sup> *Technion*    <sup>5</sup> *Barefoot Networks*

HotNets 2016, CoNEXT 2017

# CLOVE assumptions

- Clove operates over a DC Overlay – e.g., Stateless Transport Tunneling (STT)



## CLOVE in 1 slide

- Path discovery using traceroute probes
- Load-balancing of flowlets [FLARE '05]
- vSwitch path selection based on RTT-scale feedback
  - Explicit Congestion Notification - ECN
  - In-band Network Telemetry - INT

## Path Discovery

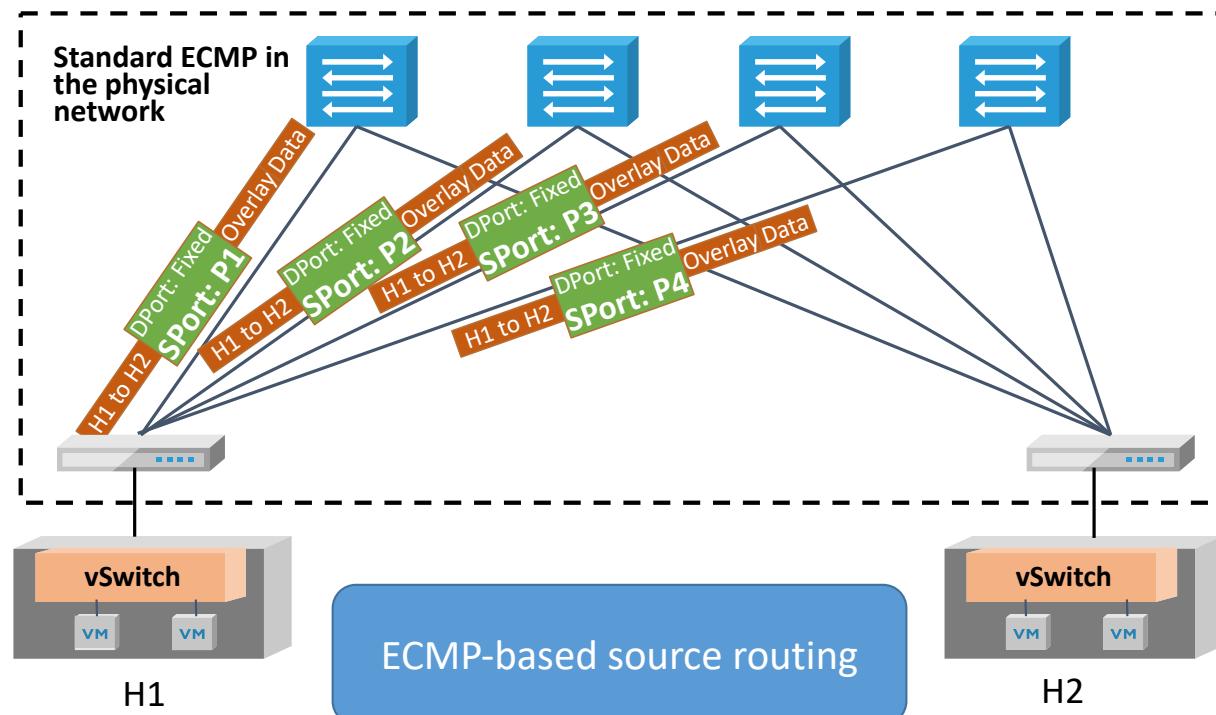
Load balancing flowlets

vSwitch Load balancing

- Outer transport source port (with ECMP) maps to network path

Hypervisor learns source port to path mapping

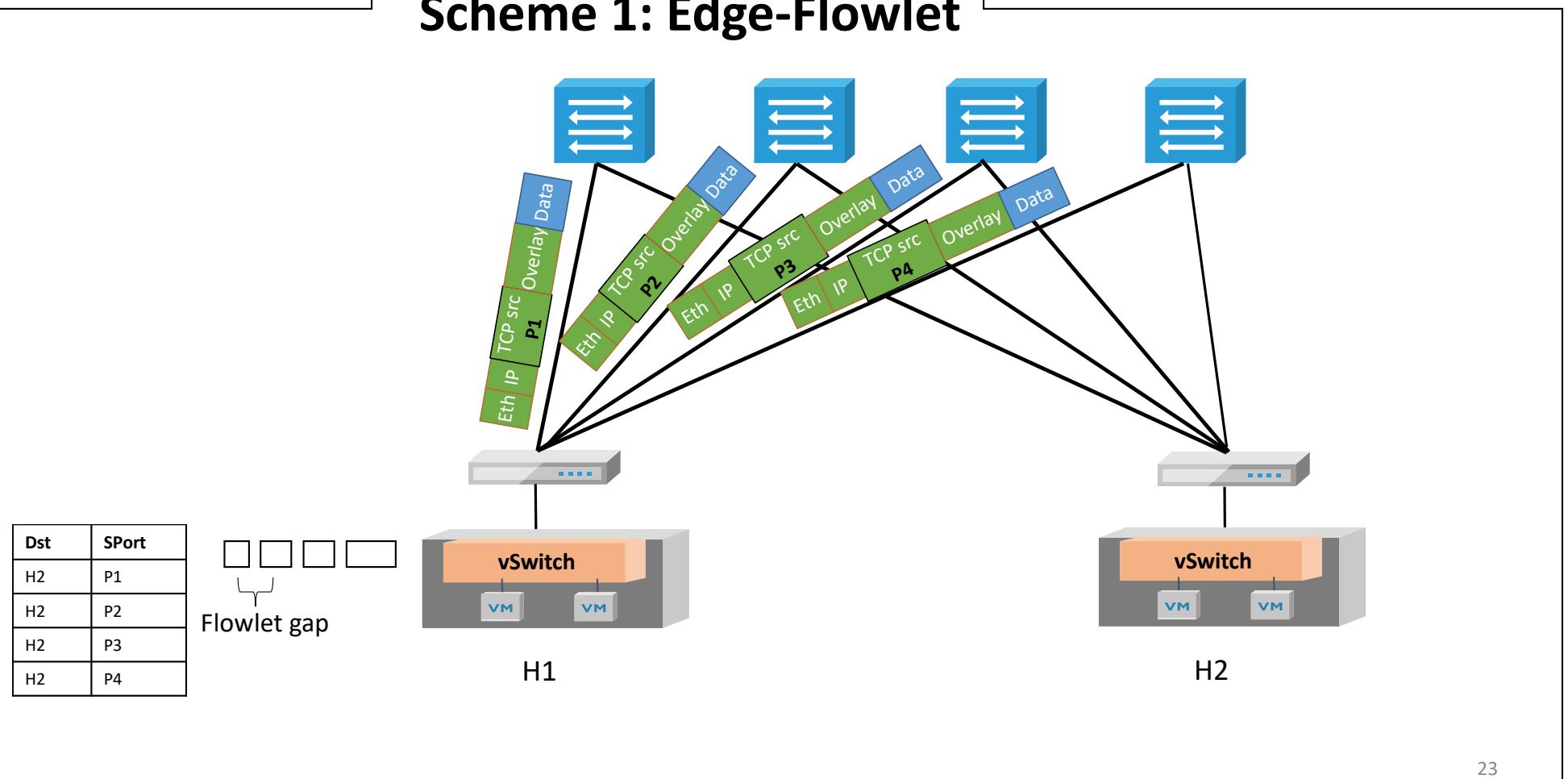
Dst	SPort
H2	P1
H2	P2
H2	P3
H2	P4



Path Discovery

## Load balancing flowlets Scheme 1: Edge-Flowlet

vSwitch Load balancing

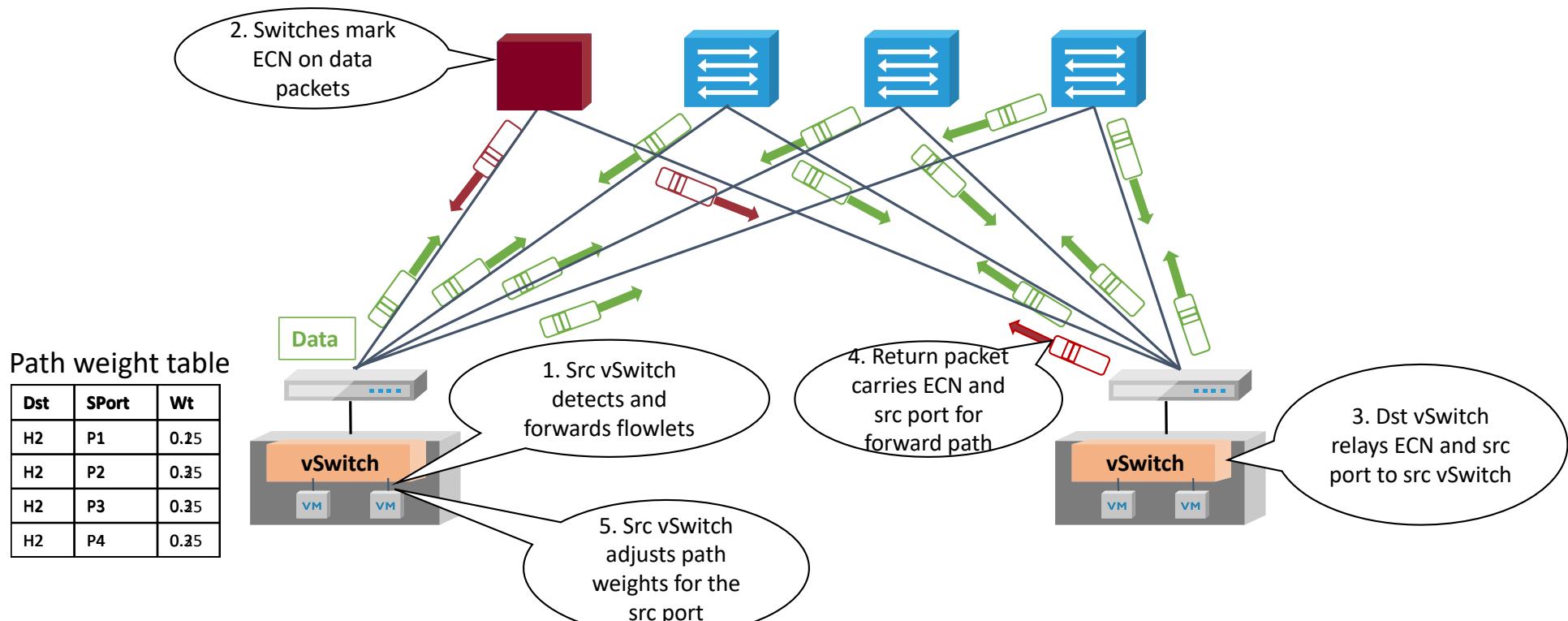


Path Discovery

Load balancing flowlets

## vSwitch Load balancing Scheme2: CLOVE-ECN

- Congestion-aware balancing based on ECN feedback

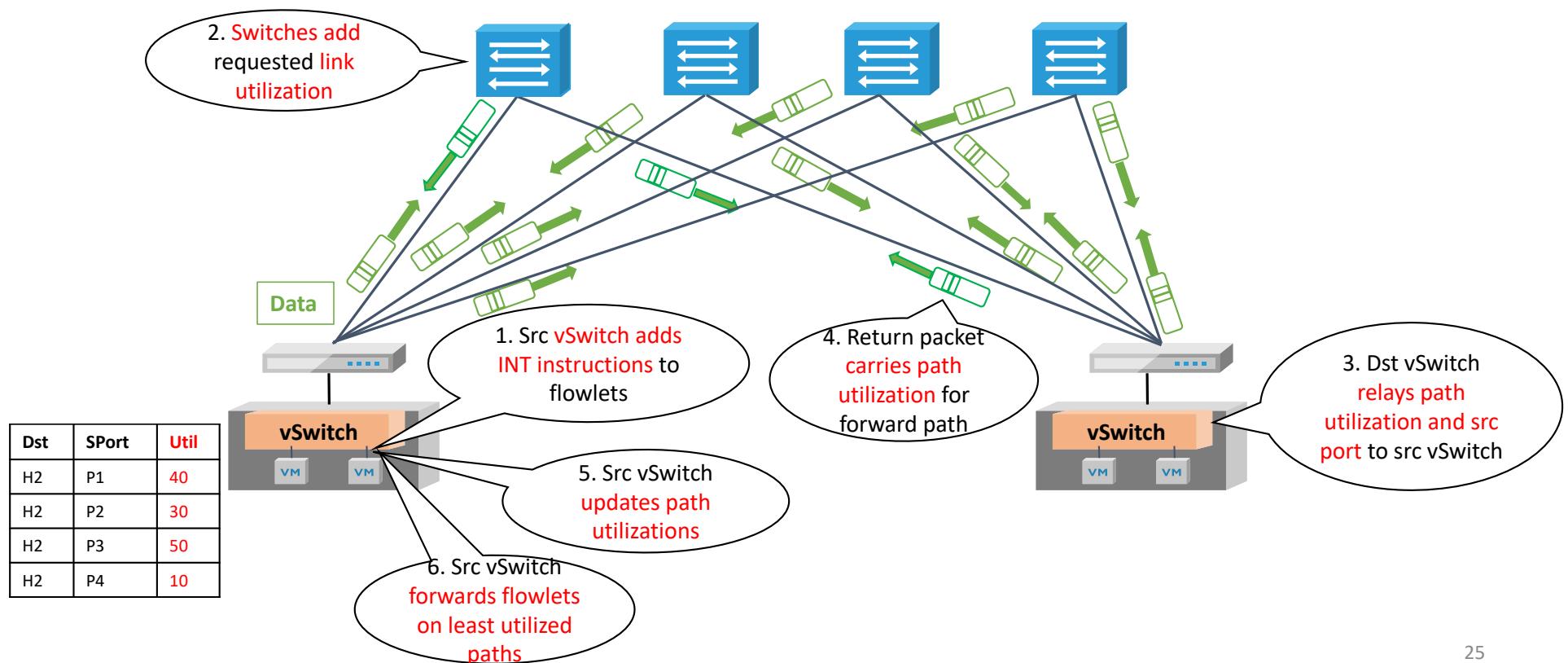


## Path Discovery

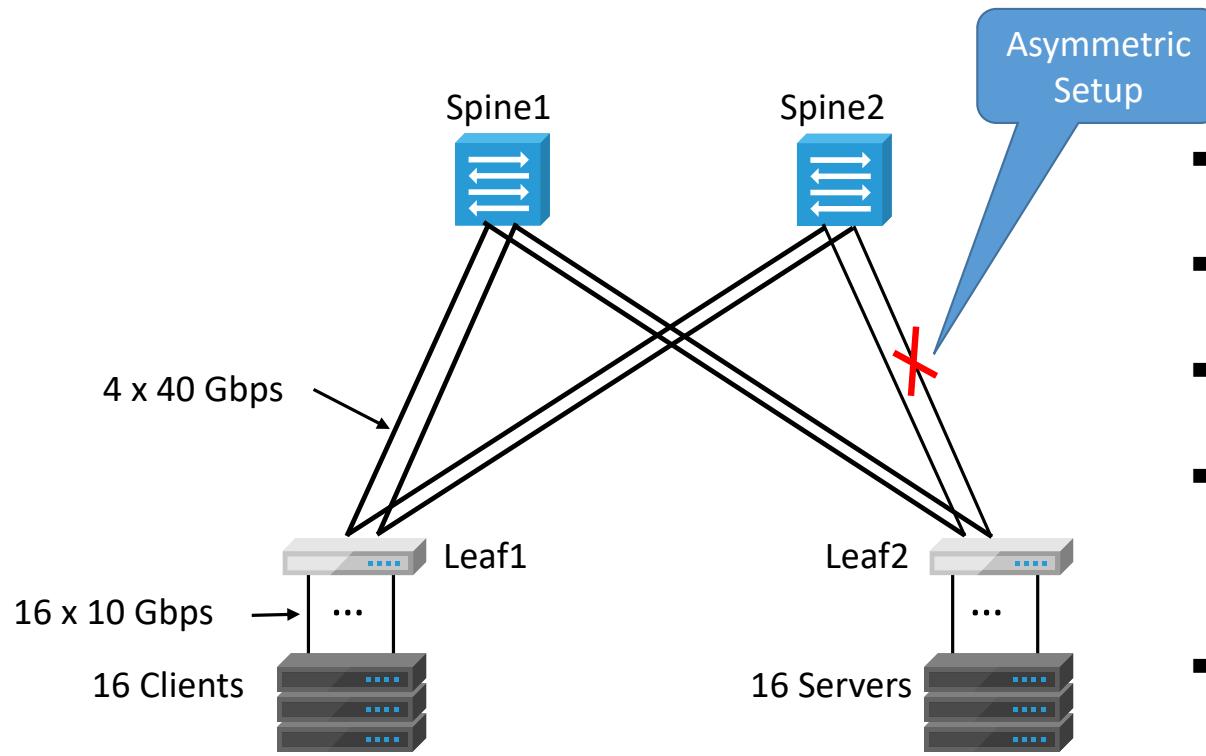
## Load balancing flowlets

# vSwitch Load balancing Scheme 3: CLOVE-INT

- Utilization-aware balancing based on INT feedback



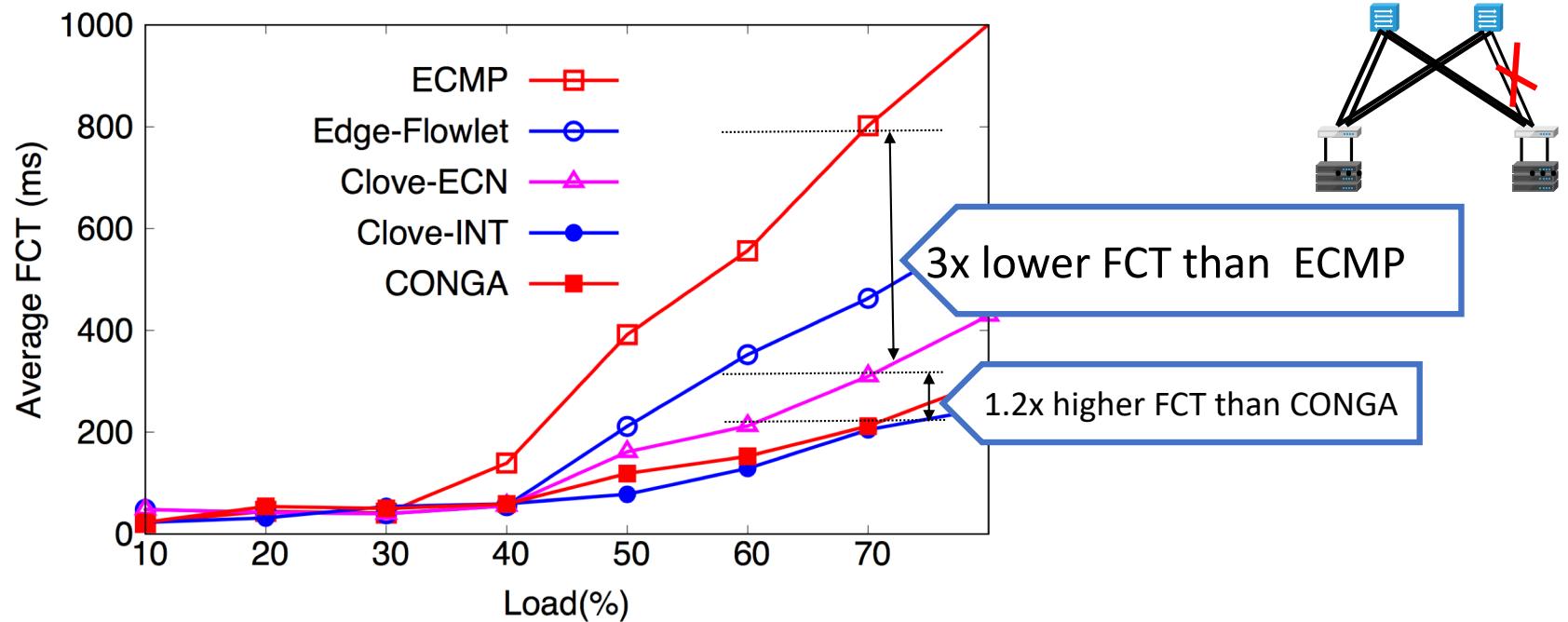
# Performance evaluation setup



- 2-tier leaf-spine symmetric topology
- Web Search Workload
- Client on Leaf1 <-> server on Leaf2
- Measure Average Flow Completion Time (FCT)
- Compare Edge-Flowlet and Clove-ECN to ECMP, MPTCP and Presto

# NS2 Simulation with CONGA – Asymmetric

Better



CLOVE-ECN captures 80% of the performance gain between ECMP and CONGA

## CLOVE highlights

- CLOVE-ECN captures 80% of the performance gain of CONGA
- No changes to network hardware, VMs, applications
- Adapts to asymmetry within the data plane
- Scalable by virtue of distributed state