# Burst-tolerant Datacenter Networks with VERTIGO

**Sepehr Abdous\*, Erfan Sharafzadeh\*, Soudeh Ghorbani**

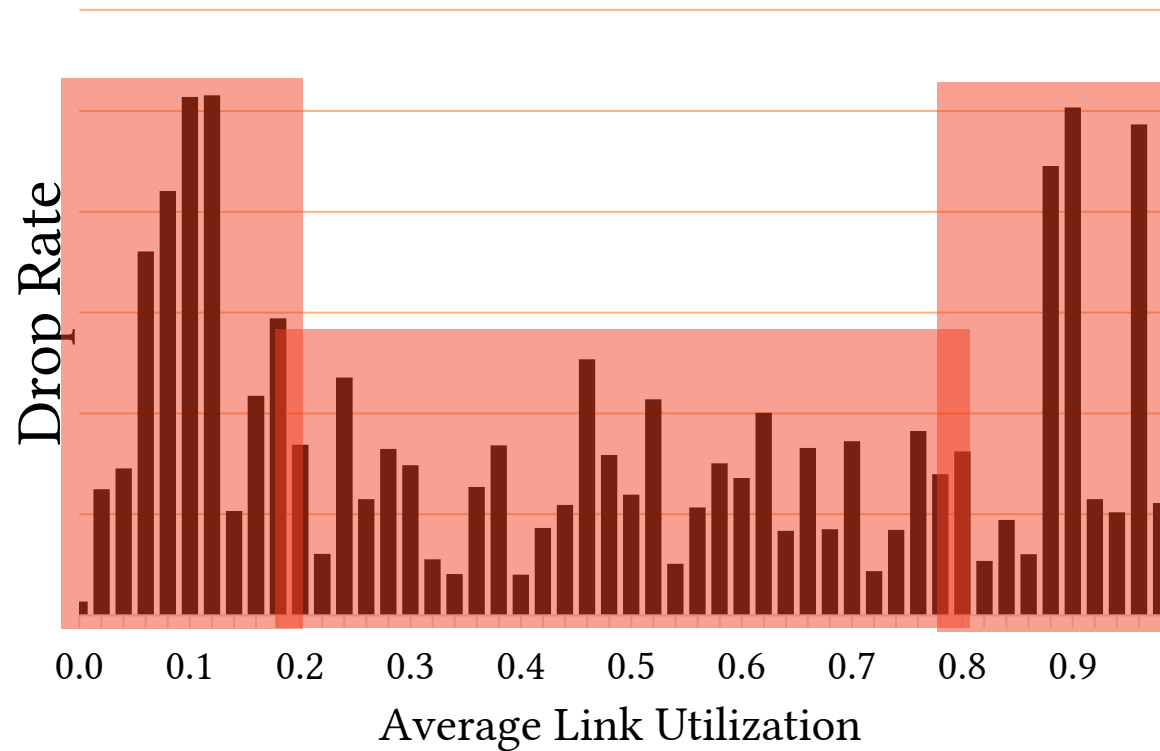**\*Co-first Authors**

JOHNS HOPKINS UNIVERSITY

# Datacenter traffic is **bursty** in short timescales

2

# Majority of drops are due to microbursts



[Zhang et al., "High-Resolution Measurement of Data Center Microbursts.", IMC '17]

!Microbursts!

High utilization periods in switch buffers that lasting 10s of μseconds

# Edge-centric congestion control: slow for microbursts

Congestion control using queue occupancy data

- HPCC [SIGCOMM '19]

Congestion control using round-trip time variations

- Swift [SIGCOMM '20]

**Deployed at the edge**
Require at least **1 RTT** to identify and recover from packet loss
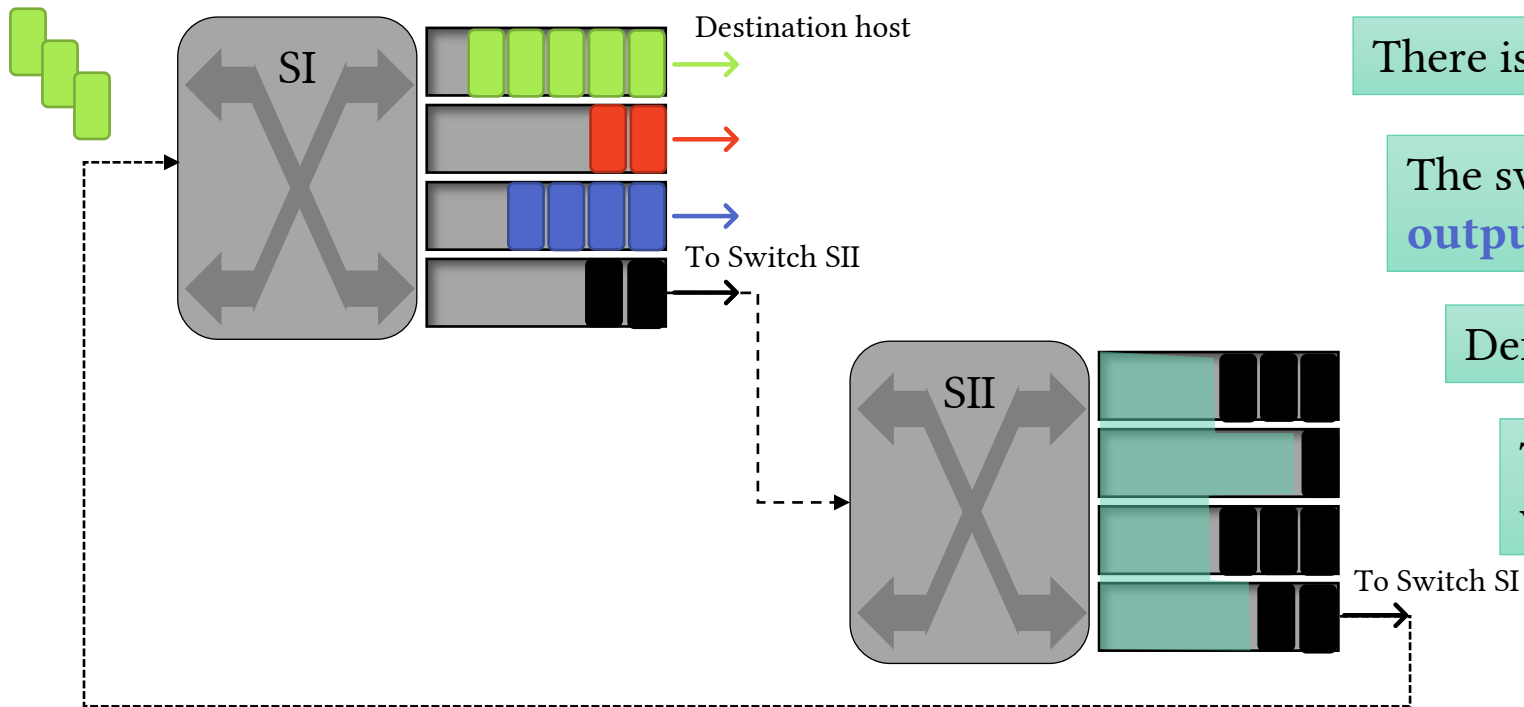
**Edge** is still slow for microbursts.
Why not react to them in the **network core**?

# Goal: Managing microbursts in the network, in real-time

5

# Deflection: a realization of in-network reaction

Randomly re-routing packets that arrive at a full buffer



Destination host

To Switch SII

SI

SII

To Switch SI

There is plenty of free buffer in **neighbors**

The switch pushes the packet into **another output buffer** instead of **dropping** it
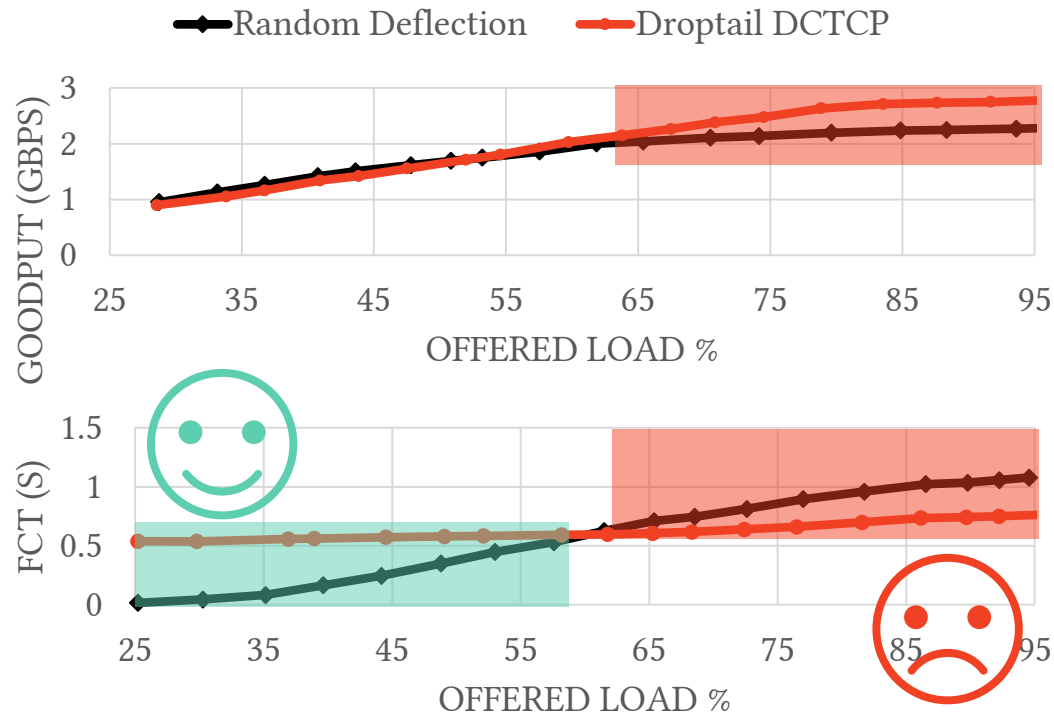
Deflected packets shortly linger in the network

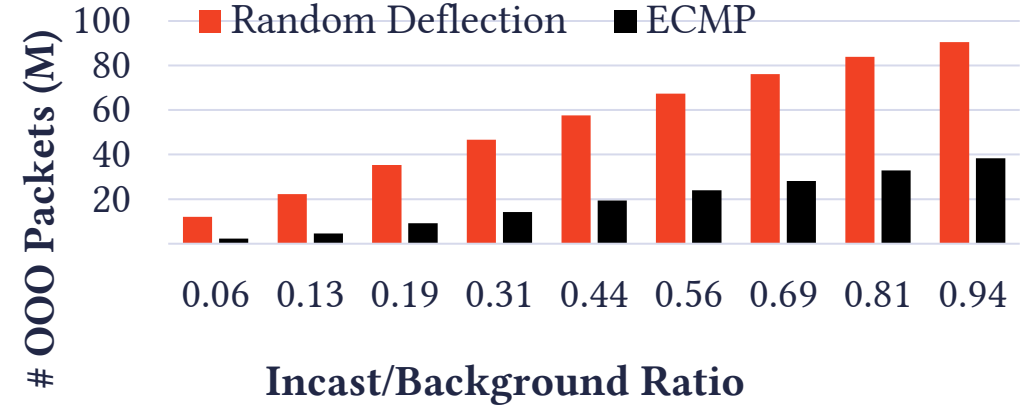There is enough room in the destination buffer when they arrive back at ToR

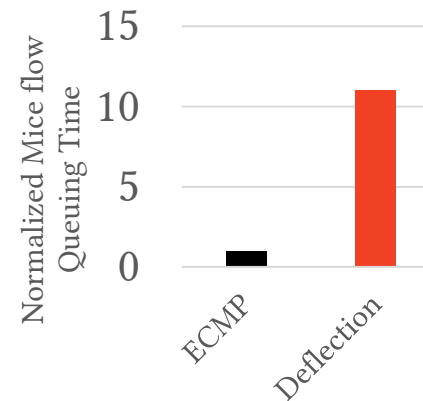# Challenges of random deflection

Setup
- 4-8-40 Two-tiered leaf-spine
- 10GB server-to-ToR, 40GB aggregate links
- DCTCP transport
- Workload: FB cache, fixed background + variable Incast



── Random Deflection   ── Droptail DCTCP



**1.** Deflection **collapses** under high loads.

■ Random Deflection   ■ ECMP



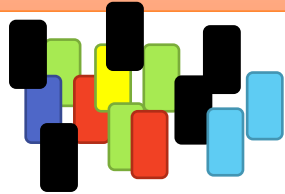**Incast/Background Ratio**

**2.** **Deflection causes heavy reordering**
Up to **10x** more out-of-order packets
**~17%** Goodput reduction



**3.** **Deflection leads to head of line blocking & starvation**
**111%** longer waits for **mice** flows (<100KB)

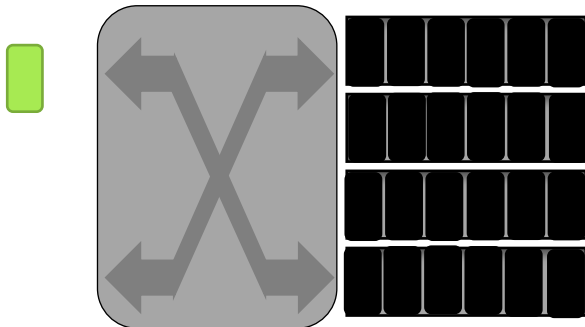# Random deflection causes head-of-the-line blocking

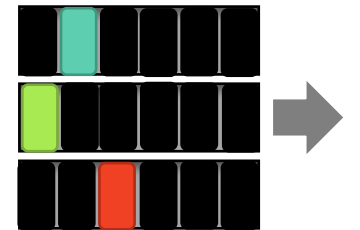Random deflection saves flows regardless of their **size**

A Large flow **continues** to send traffic instead of **backing off**

Neighbor buffers **fill up**, innocent flows are victimized

Short flows are **stuck** in congested buffers

# Random deflection breaks under load

## Problem
Random deflection treats the flows contributing to **long lasting congestion** similar to **short-lived microbursts**
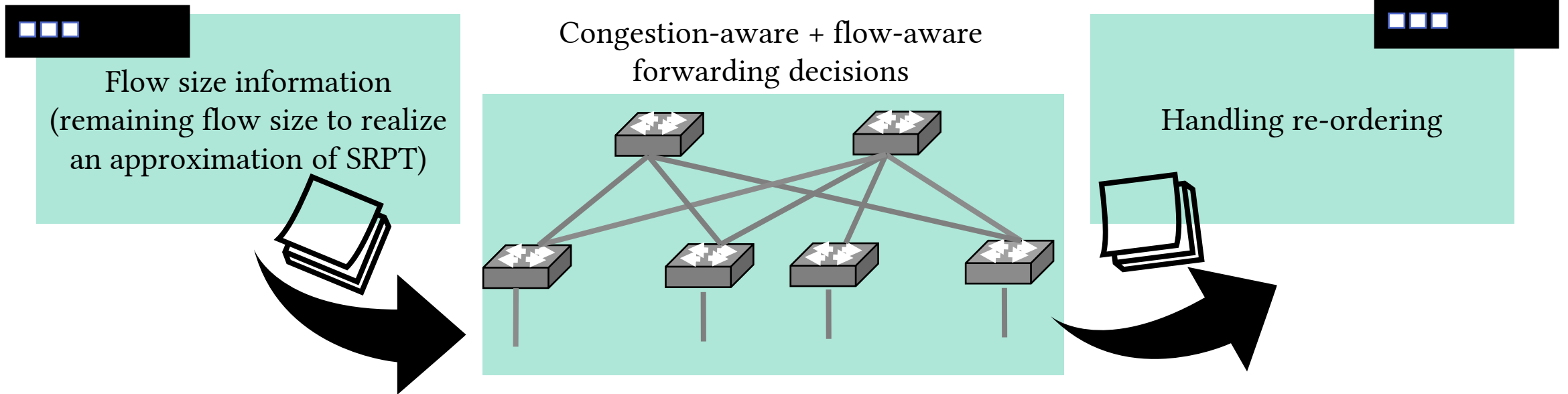
## Solution
**Detecting** the flows that are more likely to contribute to **lasting congestion** and **prioritizing** their packets for:
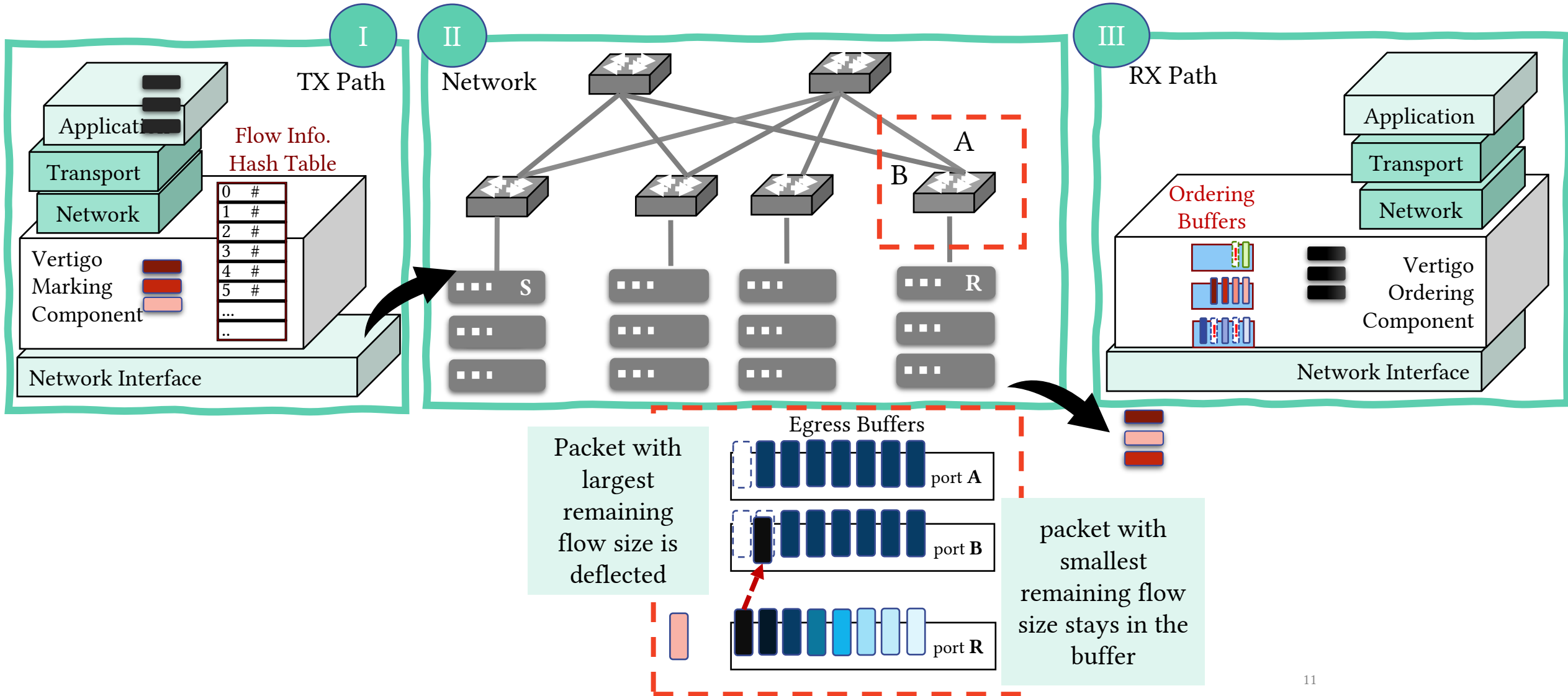(a) **deflection** under light load
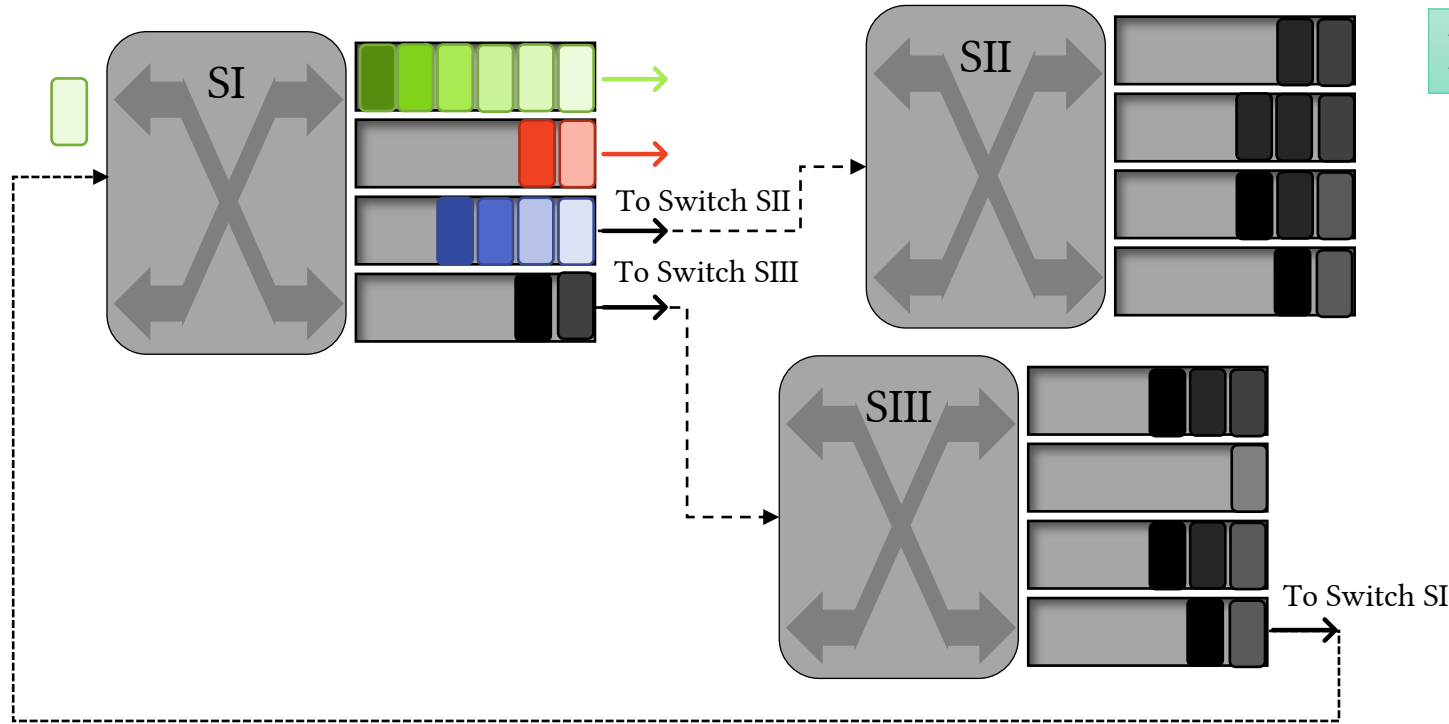(b) **drop** under high load

# Host-assisted deflection

**Remaining flow size**, a good indicator for lasting congestion.

Flow size information (remaining flow size to realize an approximation of SRPT)

Congestion-aware + flow-aware forwarding decisions

Handling re-ordering

# Vertigo: the big picture



I

TX Path

Flow Info. Hash Table

Application

Transport

Network

Vertigo Marking Component

Network Interface

| 0 | # |
| 1 | # |
| 2 | # |
| 3 | # |
| 4 | # |
| 5 | # |
| ... | |
| .. | |

II

Network

A

B

S

R

III

RX Path

Application

Transport

Ordering Buffers

Network

Vertigo Ordering Component

Network Interface

Egress Buffers

Packet with largest remaining flow size is deflected

port **A**

port **B**

port **R**

packet with smallest remaining flow size stays in the buffer

11

# Preventing collapse using flow length information



**Vertigo Fabric**
I. **Forwarding:** Least remaining flow size
II. **Congestion:** Deflect instead of Drop
III. **Deflection:** Highest remaining flow size
IV. **Load-balancing:** Power of 2 choices

Packet from a short flow arrives at a full buffer

Vertigo identifies the packet with highest remaining flow size from a full buffer

Randomly chooses two destination buffers, selects the one with least queue occupancy

Deflects the selected packet to chosen buffer

Inserts the arrived flow to its correct position w.r.t. its remaining flow size

12

# Vertigo components at the host

- **Marking** the packets based on remaining flow size
- Detecting re-transmissions to ensure **consistency**
- Boosting re-transmissions to avoid **starvation**
  - Re-transmitted packets appear as packets of a small flow
- **Ordering** shim layer at the destination
- Detailed design can be found in paper

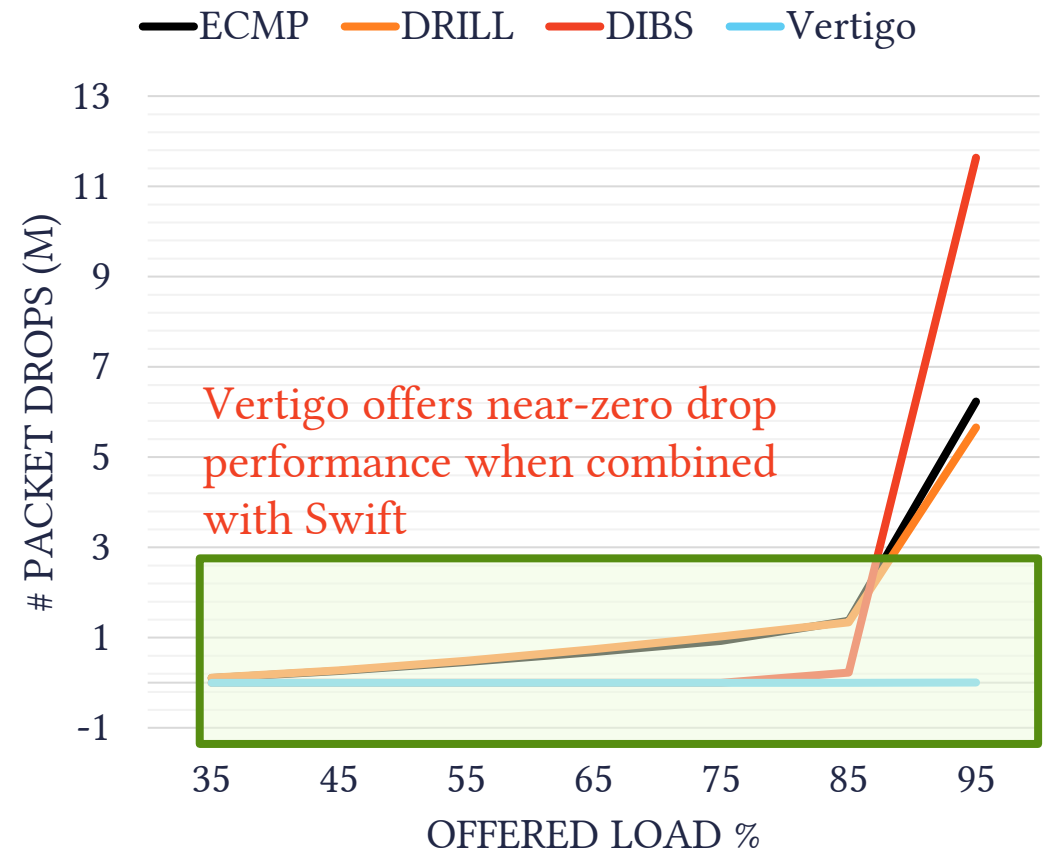# Simulation results: Vertigo's superior performance

Setup

- 4-8-40 Two-tiered leaf-spine
- 10GB server-to-ToR, 40GB aggregate links
- DCTCP transport
- Workload: FB cache, fixed background + var. Incast

| ECMP | DRILL | DIBS | Vertigo |
|------|-------|------|---------|

Improvements in tail query completion times at the highest load

57%   54%   76%

ECMP — DRILL — DIBS — Vertigo

**P99 QCT (S)** vs **OFFERED LOAD %**

ECMP — DRILL — DIBS — Vertigo

FCT of all flows (large and small) grow gracefully under Vertigo

**MEAN FCT (S)** vs **OFFERED LOAD %**

ECMP — DRILL — DIBS — Vertigo

Vertigo & DIBS, same drop rate, different outcomes!

**# DROPS (M)** vs **OFFERED LOAD %**

14

# Vertigo achieves near-0 drops with Swift



Swift cuts the tail QCT for all systems

Vertigo offers near-zero drop performance when combined with Swift

# Vertigo component analysis



Scheduling prevents collapse under high loads

Deflection helps complete incasts faster

VERTIGO — No Sched — No Deflection — No Ordering — ECMP Baseline



Reduced probability of finding free buffers

Impact of power of two forwarding vs random

Mean QCT (s) / Offered Load %

■ Power-of-two Vertigo  ■ Random Vertigo



Boosting prevents starvation, helps complete flows

>50%

Query Completion Ratio % / No Boosting / x2 Boosting / x4 Boosting / x8 Boosting



Ordering layer improves application throughput

~10%

Aggregate App. throughput (Gbps) / Offered Load %

■ Vertigo  ■ No-ordering

16

**Deflection**: Cuts the completion time tail

**Scheduling**: Prevents the collapse

**Ordering**: Preserves app throughput

**Boosting**: Prevents starvation

17

# Vertigo Conclusions

**Key Takeaway:**

To properly react to microbursts, network-centric **real-time action** and end-host's **advance knowledge of flow sizes** are vital!

**Vertigo:**

A hybrid solution to tolerate micro-scale bursty traffic by changing the forwarding decisions upon facing imminent packet loss

**Challenges:**
- Both host and network must be changed
- Existing queue management abstractions are not enough

Check out Vertigo artifacts!
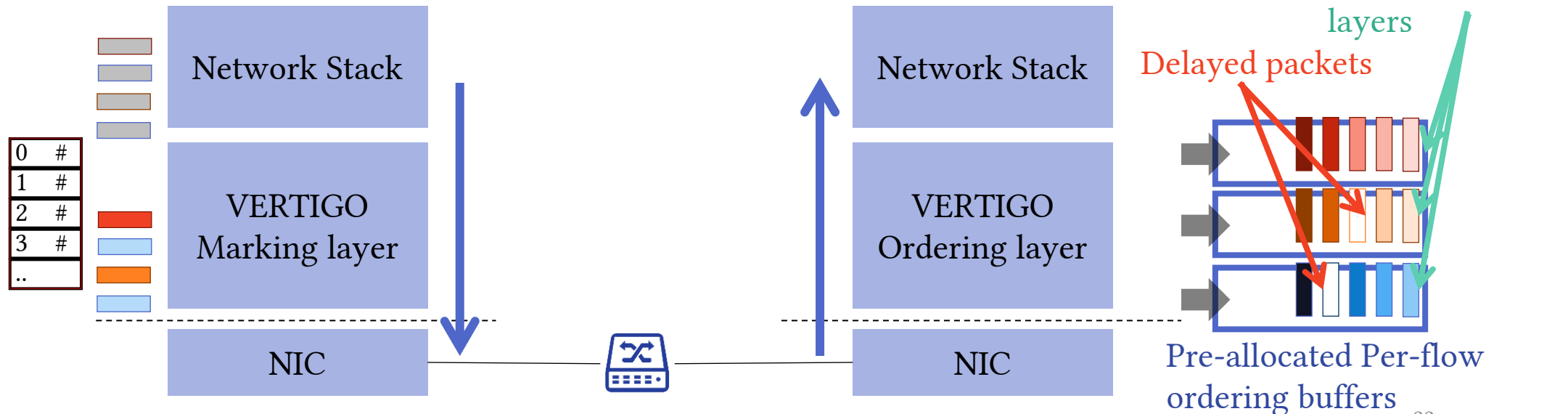https://github.com/hopnets/vertigo-artifacts

# Thank you!

Contact us

- [sabdous1@jhu.edu](mailto:sabdous1@jhu.edu)
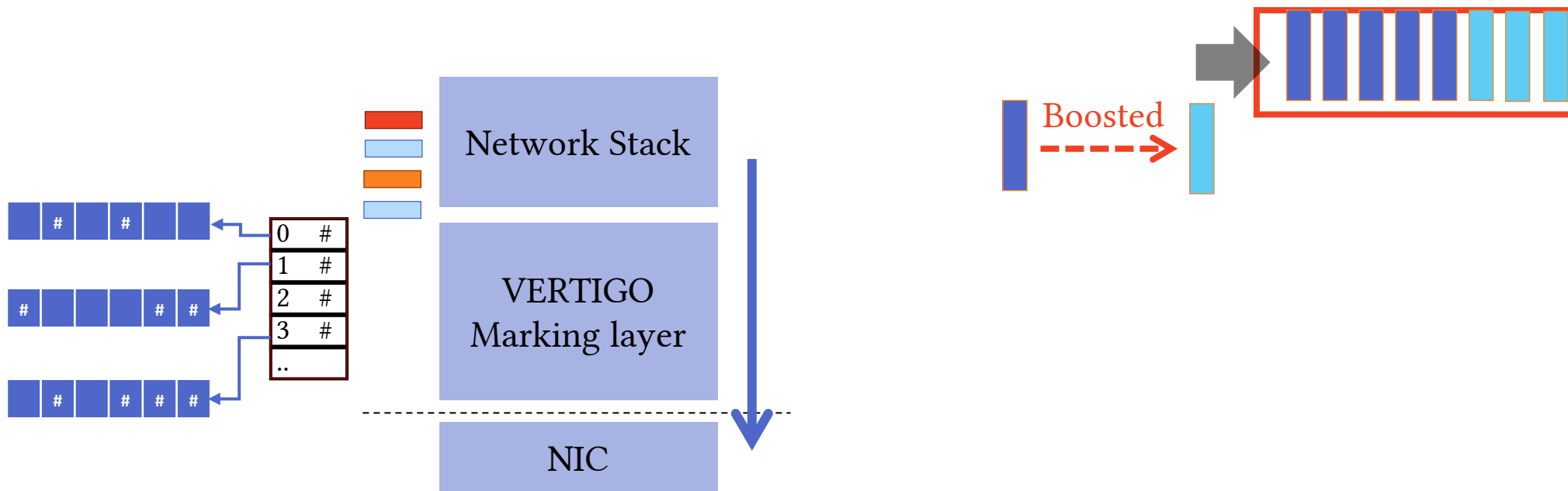- [erfan@cs.jhu.edu](mailto:erfan@cs.jhu.edu)

# Backup slides

# Handling packet reordering

- Mark packets with **remaining flow size (RFS)** @sender
- Flow size tracking is **transport-independent**
- RFS must be **unique** per-flow
- RFS used @destination to **order** the packets



In-order packets immediately flushed to upper layers

Delayed packets
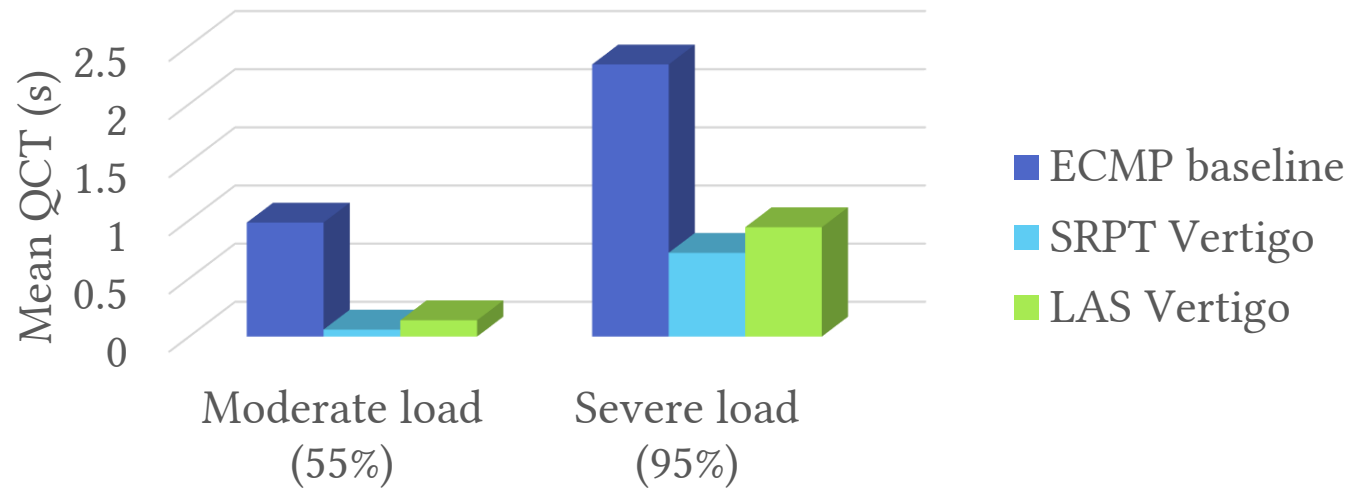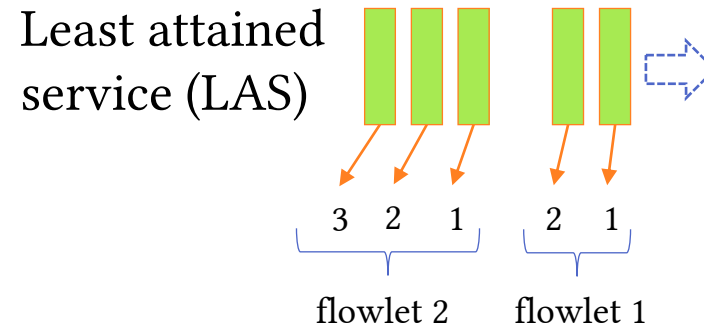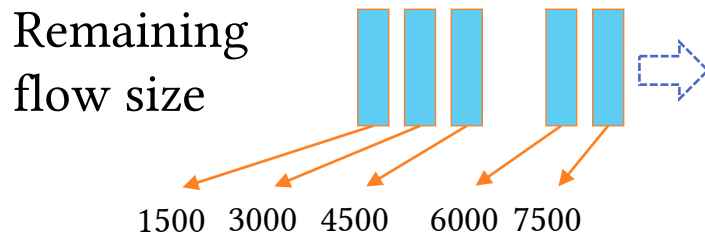
Pre-allocated Per-flow ordering buffers

# Saying no to starvation

- Keeping track of **re-transmissions** to ensure RFS **consistency**
- **Boost** the re-transmitted packet by cutting its RFS

# Simple marking by counting upwards

What if flow size information is not available?