# Carousel: Scalable Traffic Shaping at End Hosts

**Ahmed Saeed**, Nandita Dukkipati, Vytautas Valancius, Vinh The Lam, Carlo Contavalli, and Amin Vahdat

Rate limiting and isolation between thousands of flows per machine
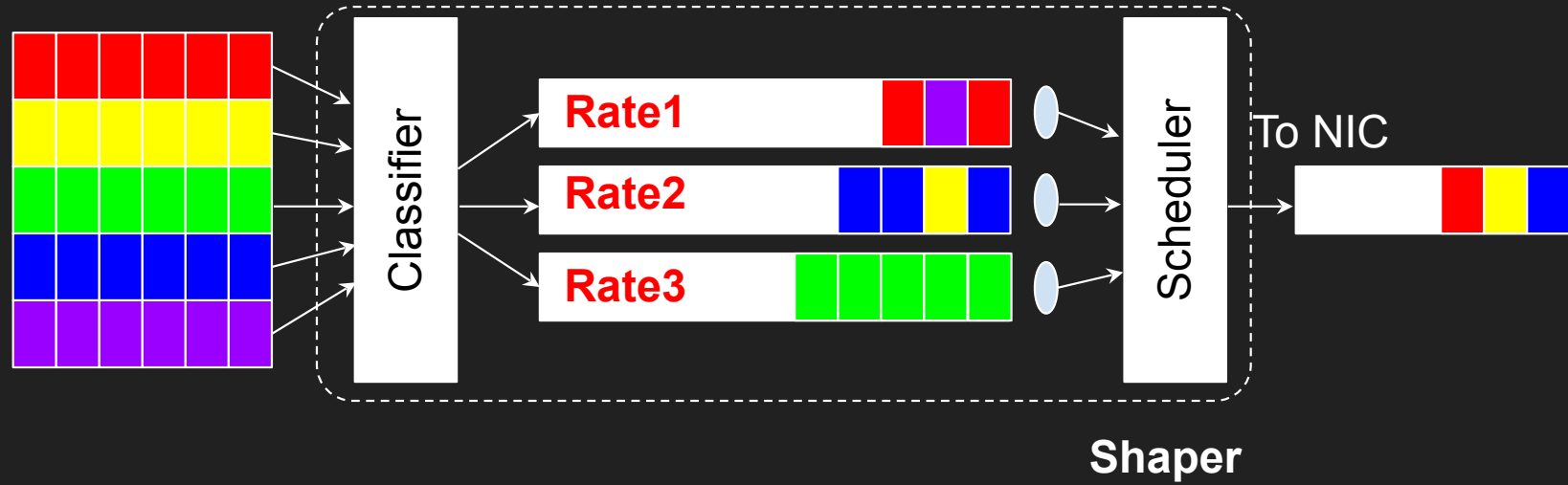[BwE - SIGCOMM '15]

**Rate limiting and isolation between thousands of flows per machine [BwE - SIGCOMM '15]**

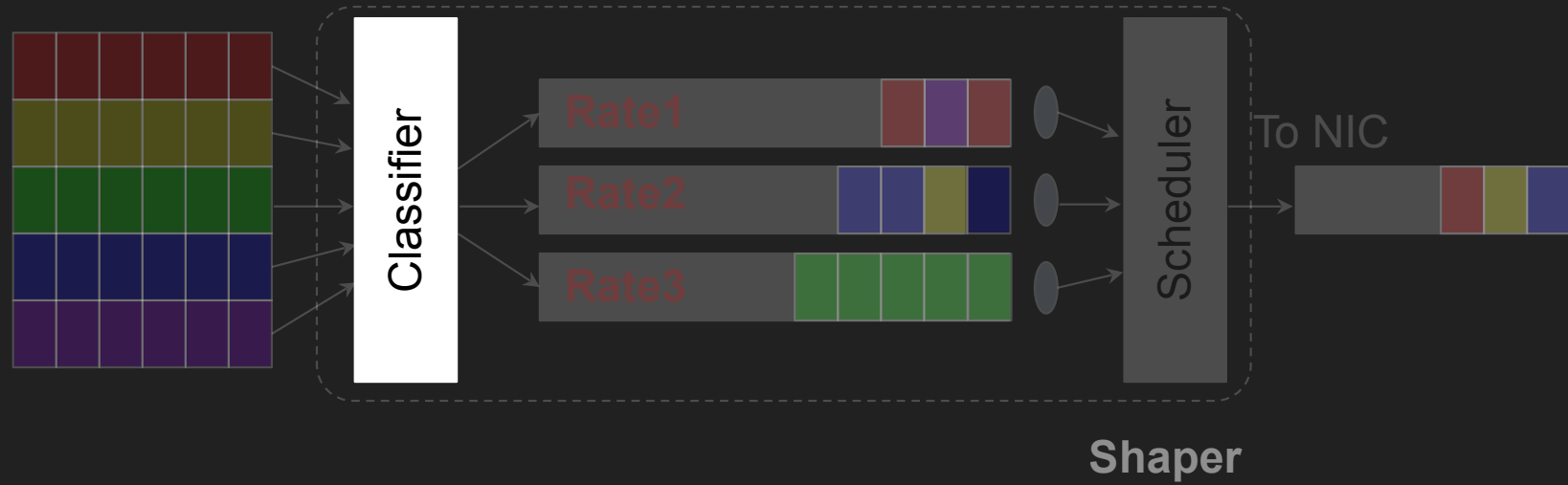**New protocols that require per-flow pacing [TCP BBR and TIMELY - SIGCOMM '15]**
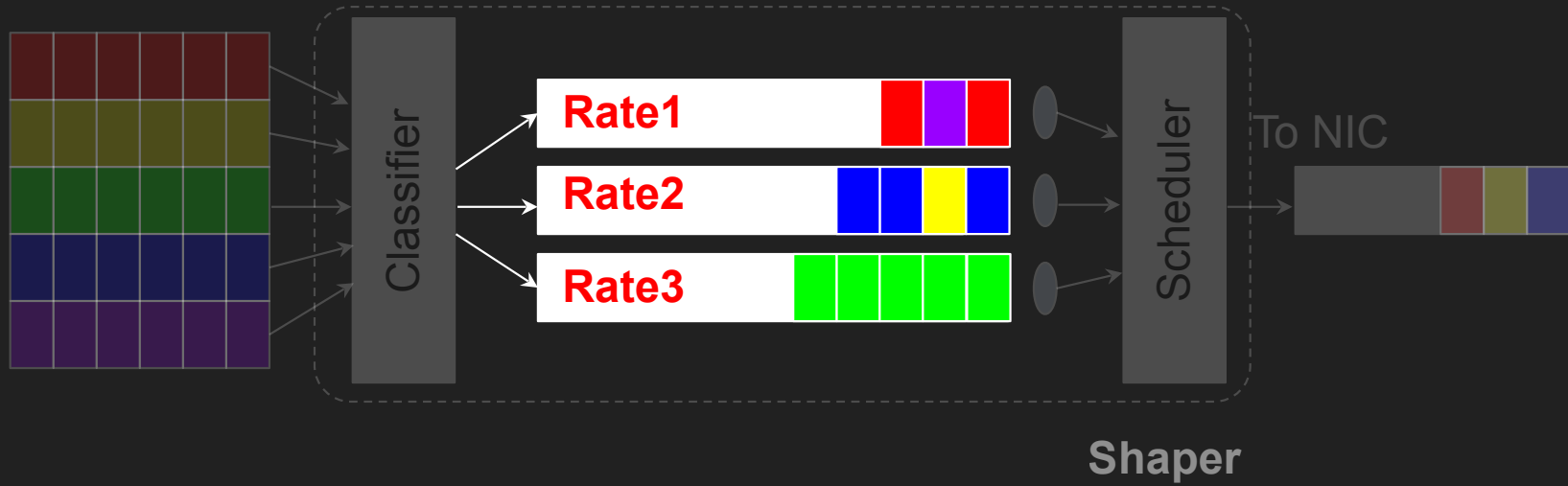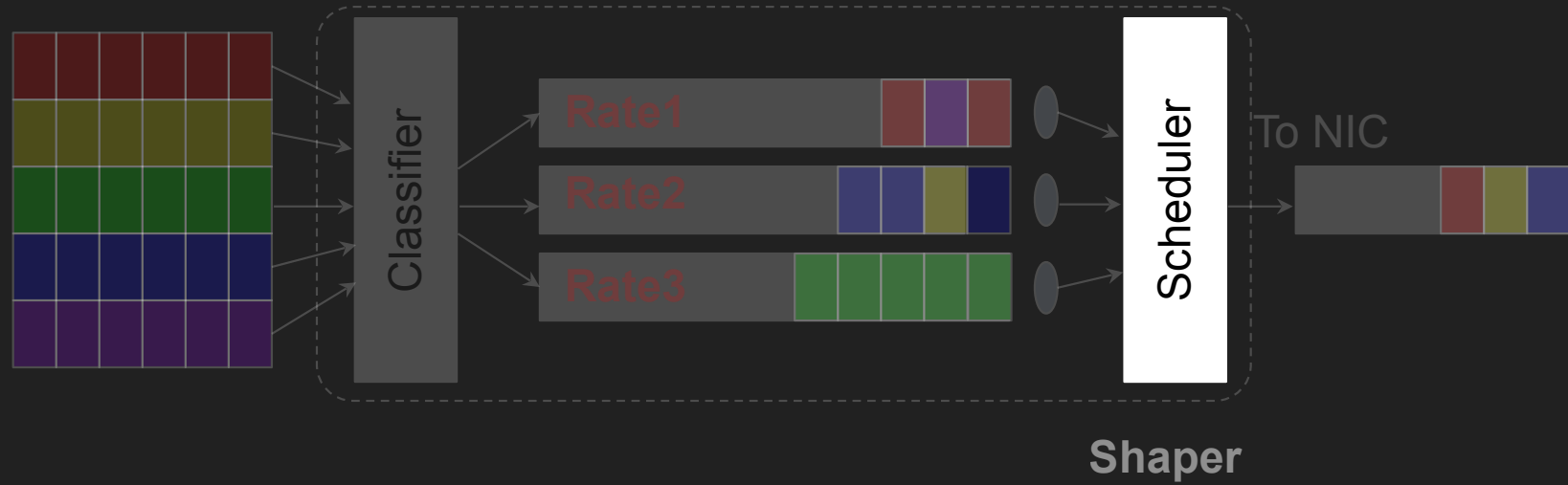
# Traffic Shaping

# Traffic Shaping

# Traffic Shaping

# Traffic Shaping



Packet sources

Classifier

Rate1

Rate2

Rate3

Scheduler

To NIC

Shaper

# Traffic Shaping

# Traffic Shaping

Packet sources

Classifier

Rate1

Rate2

Rate3

Scheduler

To NIC

**Shaper**

**Overhead of managing a queue per configured rate**

# Traffic Shaping

# Traffic Shaping

*We need new traffic shapers that can handle tens of thousands of flows and rates*

Main Idea💡
Replace the **many queues**
with **a single low-overhead queue**

# Contributions

# Contributions

# Contributions



Timestamper

**Shaper**

**Timing Wheel**

To NIC

**Single, O(1), Queue**

# Contributions

# Contributions

# Outline

- Problems with Current Shapers

- Carousel Overview

- Single Queue Shaping

- Backpressure

- Evaluation

# Problems with Current Shapers

# FQ/Pacing


*Source flows*

- Implements per TCP flow pacing

- Requires a queue per flow
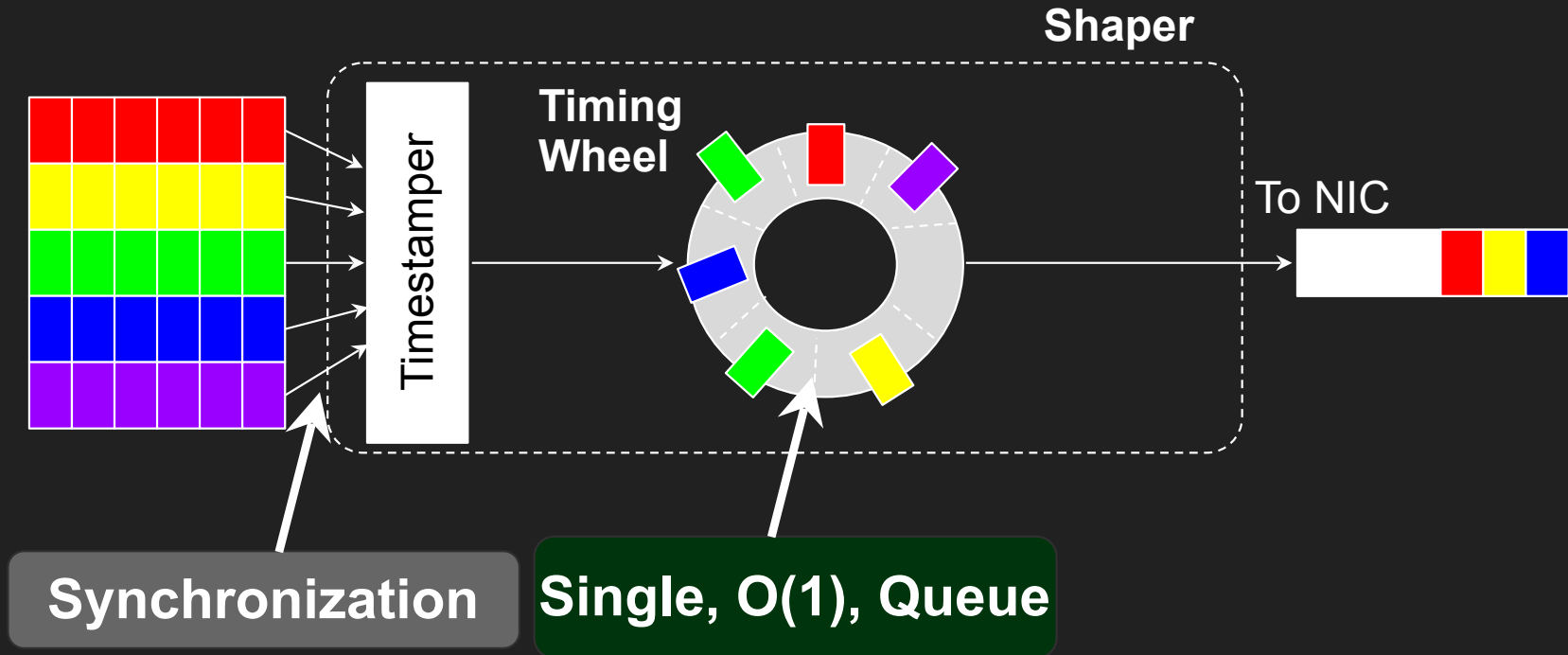
  - Flows are kept in order of their scheduled transmission time

  - Flows are dequeued in order

- O(log n) operations per packet to operate on a sorted list of flows



**Flow state lookup**

*Hashtable of red-black trees on flow ids*

**Flow pacing**

time_next_packet

F1
t=now+5

F2
t=now+3

F3
t=now+8

F7
t=now+2

F6
t=now+4

F5
t=now+6

F4
t=now+9

*Flow delays (red-black tree on timestamps)*

NIC

CPU utilization for FQ/pacing and a NOOP Qdisc for the same load

**FQ/Pacing introduces 10% more
CPU overhead**

# Carousel Overview

# Carousel Overview



- Relies on a single queue for all packets from all flows

- Requires a high frequency timer or busy polling

- Pinned to a single core

# Single Queue Shaping

# Single Queue Shaping

- All packet are sorted by their transmission time in one data structure

- A single queue for all traffic will need to handle tens of thousands of packets

- **Challenge:**
  Enqueue and dequeue in a data structure of sorted elements at line rate

# Single Queue Shaping

- All packet are sorted by their transmission time in one data structure

- A single queue for all traffic will need to handle tens of thousands of packets

- **Challenge:**
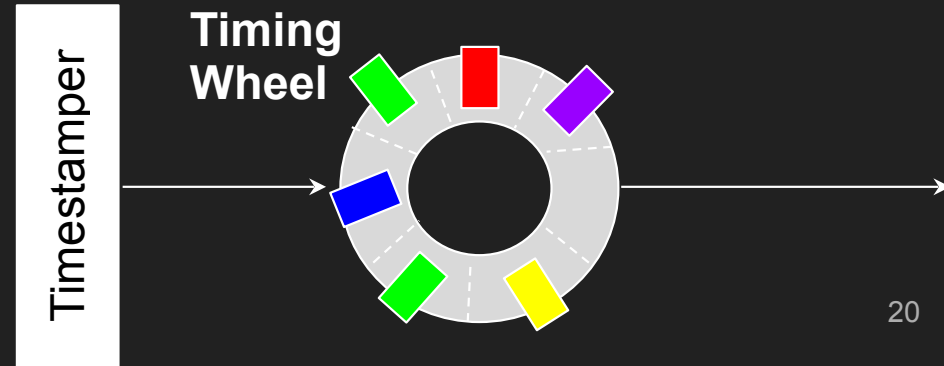  Enqueue and dequeue in a data structure of sorted elements at line rate



Timestamper

**Timing Wheel**

# Timing Wheel [Varghese et al. SOSP '87]

- Bucket sort approach to Calendar Queue covering a time horizon

  - Relies on having a minimum rates

- Implemented as an array of buckets each a linked list of packets

  - Each bucket represents a certain time range



Time Horizon (h)

# Timing Wheel Benchmark

- Measured overhead per enqueue/dequeue pairs

- Overhead per element is between 21-22 nanoseconds

    - Fixed for 2000 to 2 million sorted elements
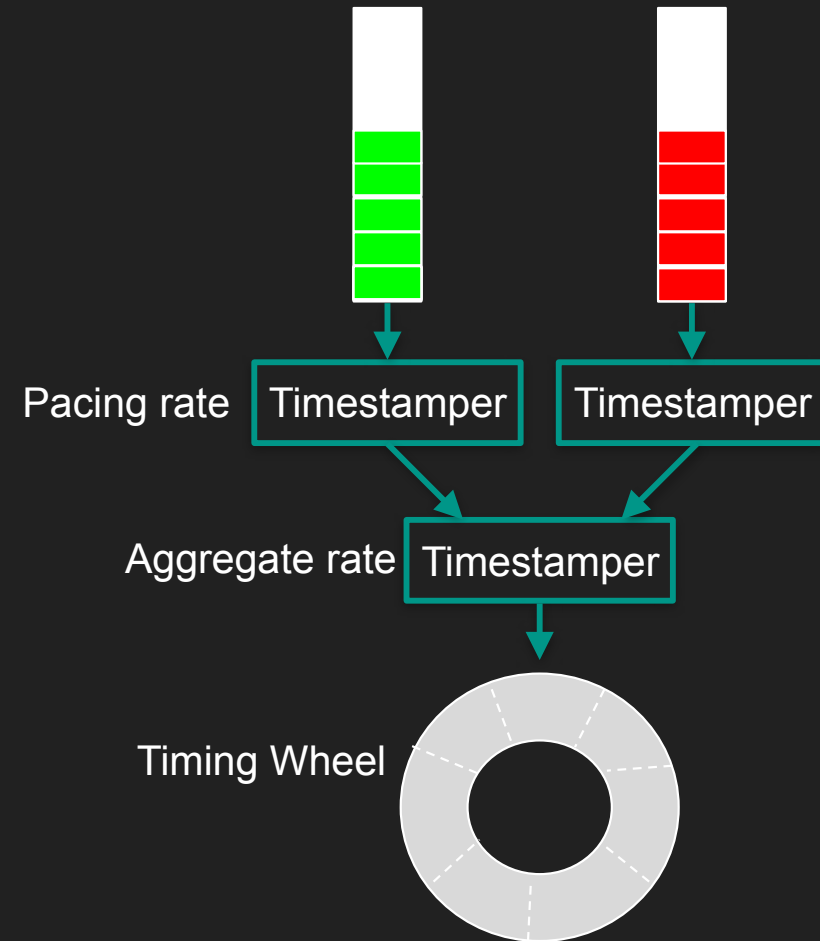
    - 21 nanoseconds per packet = 500 Gbps (for 1500 byte packets)

# Timestampers

- Packets are timestamped by policy enforcers in their transmission path
  - TCP timestamps a packet based on its pacing rate
  - Bandwidth enforcer timestamps a packet based on its policy-based aggregate rate

- Carousel picks the largest timestamp

- $NextTimestamp = LastTimestamp + \dfrac{SizeOfPacket}{ConfiguredRate}$

Pacing rate

Timestamper    Timestamper

Aggregate rate   Timestamper

Timing Wheel

# Example of Shaping using Carousel

# Example of Shaping using Carousel



Rate (1 pps)

Rate (0.5 pps)

Timestamper

T=1

**Timing Wheel**

Each bucket represents 1 second

# A time step **0**

# Example of Shaping using Carousel



Rate (1 pps)

Rate (0.5 pps)

Timestamper

T=2

**Timing Wheel**

Each bucket represents 1 second

## A time step **0**

# Example of Shaping using Carousel



Rate (1 pps)

Rate (0.5 pps)

Timestamper

T=3

**Timing Wheel**

Each bucket represents 1 second

# A time step **0**

# Example of Shaping using Carousel

Rate (1 pps)

Rate (0.5 pps)

Timestamper

T=2

**Timing Wheel**

Each bucket represents 1 second

A time step **0**

# Example of Shaping using Carousel



Rate (1 pps)

Rate (0.5 pps)

Timestamper

T=4

**Timing Wheel**

Each bucket represents 1 second

A time step **0**

# Example of Shaping using Carousel

Rate (1 pps)

Rate (0.5 pps)

Timestamper

**Timing Wheel**

Each bucket represents 1 second

## A time step **0**

# Example of Shaping using Carousel

Rate (1 pps)

Rate (0.5 pps)

Timestamper

**Timing Wheel**

Each bucket represents 1 second

# A time step **1**

# Example of Shaping using Carousel



Rate (1 pps)

Rate (0.5 pps)

Timestamper

**Timing Wheel**

Each bucket represents 1 second

## A time step **2**

# Example of Shaping using Carousel

Rate (1 pps)

Rate (0.5 pps)

Timestamper

**Timing Wheel**

Each bucket represents 1 second

# Backpressure with Deferred Completion

# The Value of Backpressure

- Without backpressure shaper queues get full with small number of flows causing

# The Value of Backpressure

- Without backpressure shaper queues get full with small number of flows causing
  - Unnecessary drops (when the queue is full the queue tail drops)

# The Value of Backpressure

- Without backpressure shaper queues get full with small number of flows causing
    - Unnecessary drops (when the queue is full the queue tail drops)
    - Head of Line Blocking
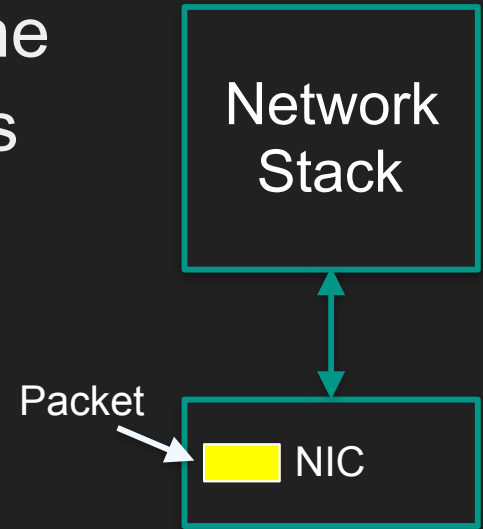
# The Value of Backpressure

- Without backpressure shaper queues get full with small number of flows causing
    - Unnecessary drops (when the queue is full the queue tail drops)
    - Head of Line Blocking

- Backpressure allows shapers to control sender rate and avoid overwhelming the shaper
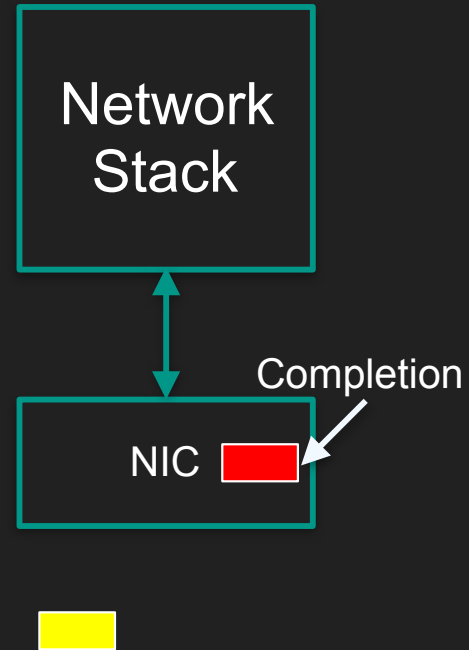
# The Completion Signal

● Completions are signals from the NIC to the network stack to inform it that a packet has been transmitted

Network Stack

Packet

NIC

# The Completion Signal

- Completions are signals from the NIC to the network stack to inform it that a packet has been transmitted

Network Stack
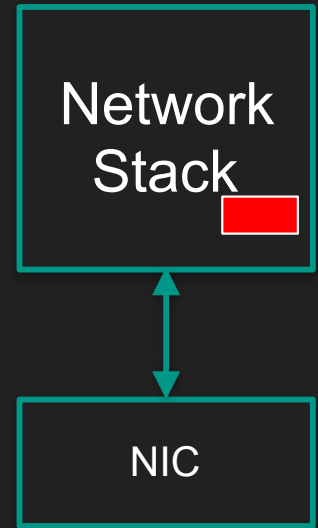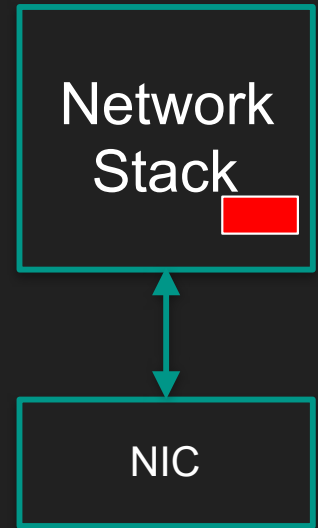
Completion

NIC

27

# The Completion Signal

- Completions are signals from the NIC to the network stack to inform it that a packet has been transmitted

# The Completion Signal

- Completions are signals from the NIC to the network stack to inform it that a packet has been transmitted
  - Completions are typically delivered in order
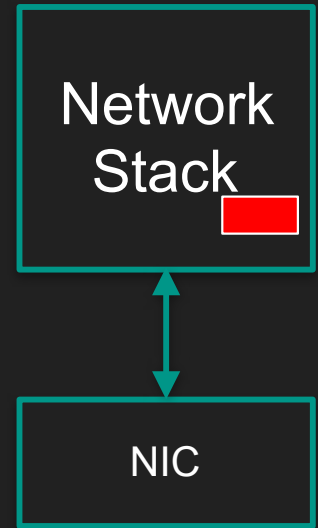


Network Stack

NIC

# The Completion Signal

- Completions are signals from the NIC to the network stack to inform it that a packet has been transmitted
  - Completions are typically delivered in order
  - Completion should be controlled by the hypervisor not the virtual NIC
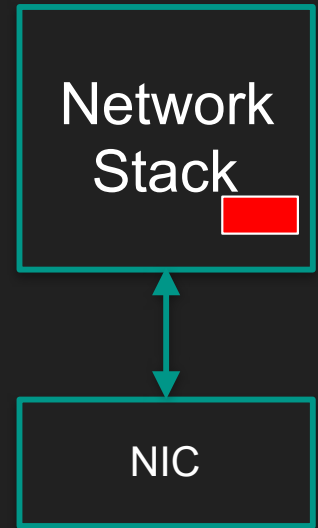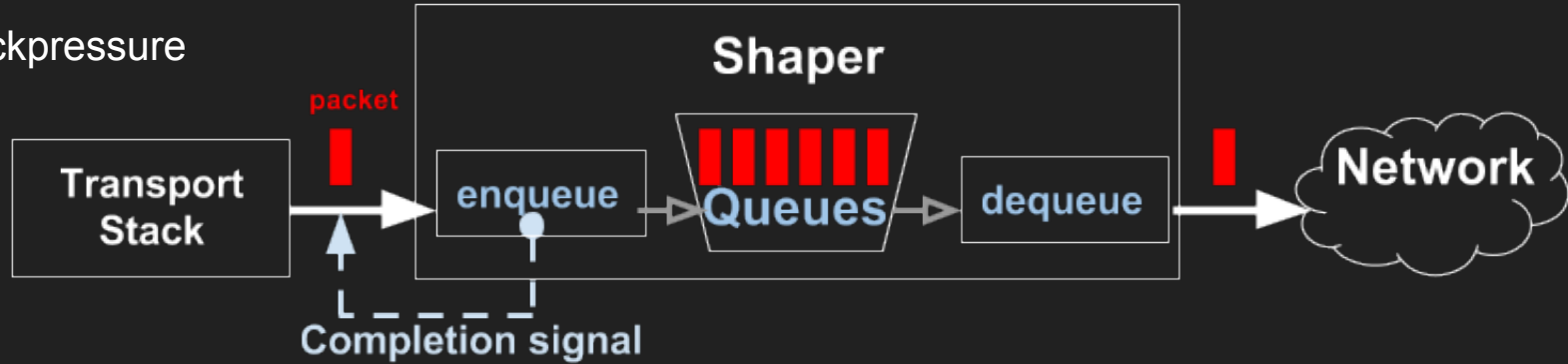
Network
Stack

NIC

# The Completion Signal

- Completions are signals from the NIC to the network stack to inform it that a packet has been transmitted
    - Completions are typically delivered in order
    - Completion should be controlled by the hypervisor not the virtual NIC

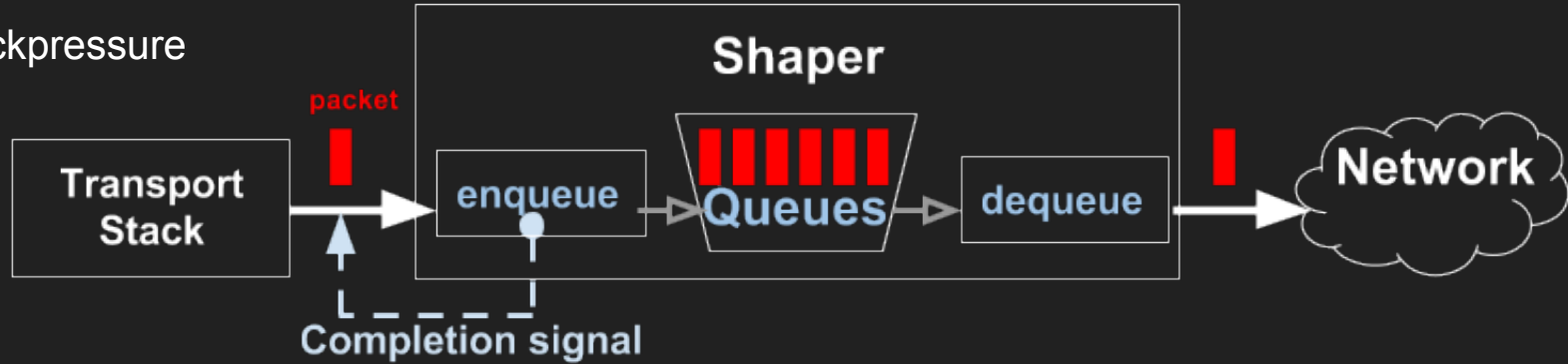- Completions should be delivered out of order and completely controlled by Shapers
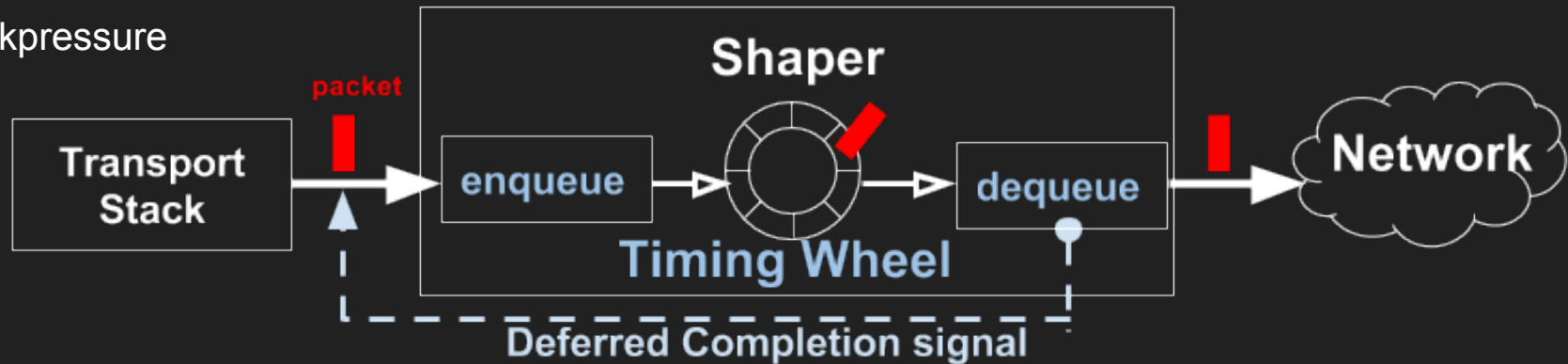
Network
Stack

NIC

# Backpressure with Deferred Completion
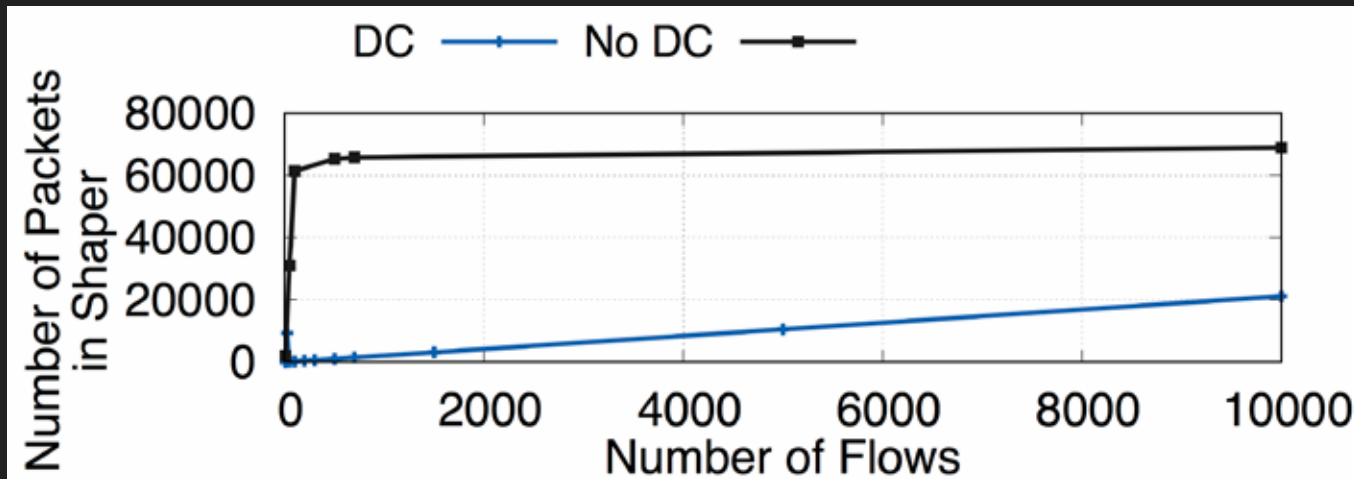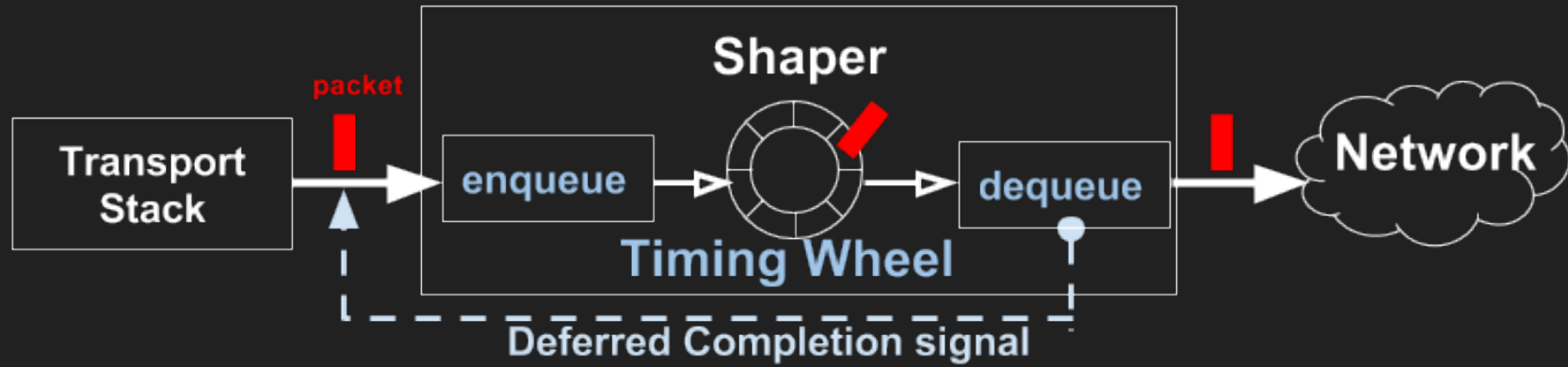
Without Backpressure
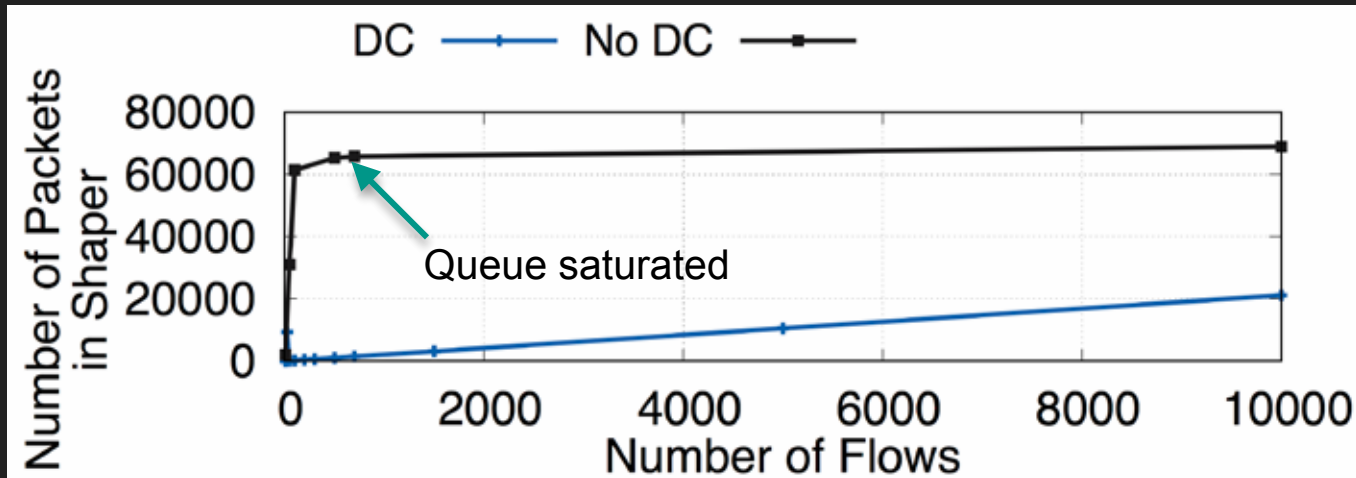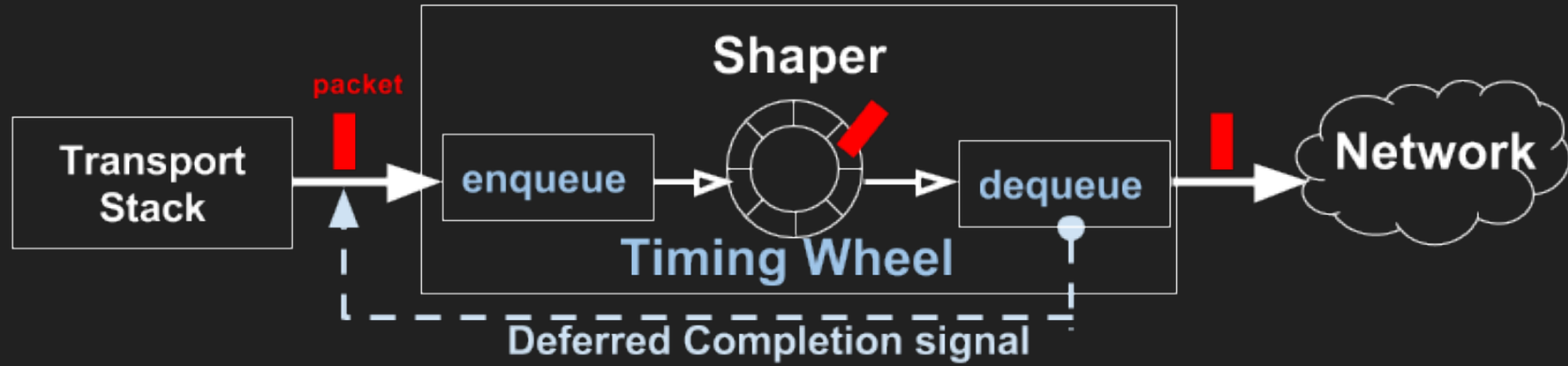
# Backpressure with Deferred Completion

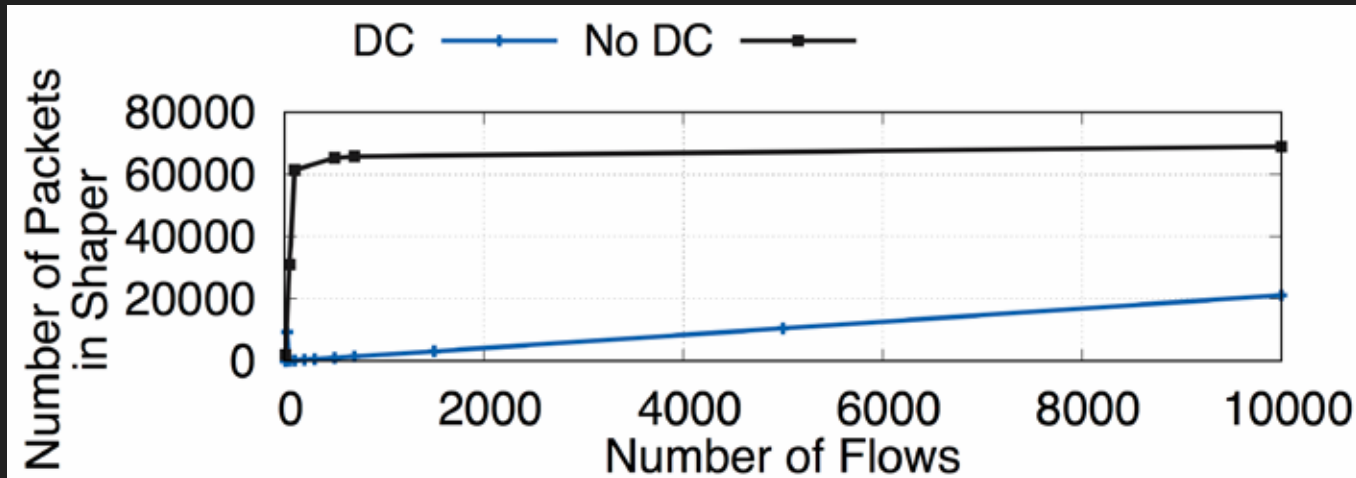# Backpressure with Deferred Completion

# Backpressure with Deferred Completion

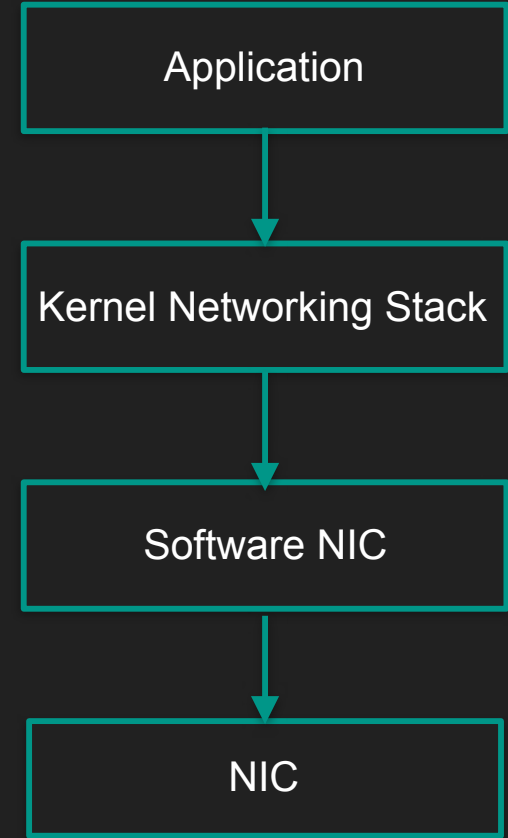# Backpressure with Deferred Completion

**Deferred completions limits the number of packets in shaper reducing its memory footprint**

Shaper

Deferred Completion signal

# Evaluation

# Evaluation Setup

- Carousel deployed within a Software NIC



Stack

C

34

# Evaluation Setup

- Carousel deployed within a Software NIC

- Evaluation on Youtube servers comparing Carousel and FQ/Pacing

```
┌─────────────────────────┐
│       Application       │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│ Kernel Networking Stack │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│      Software NIC       │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│           NIC           │
└─────────────────────────┘
```

34
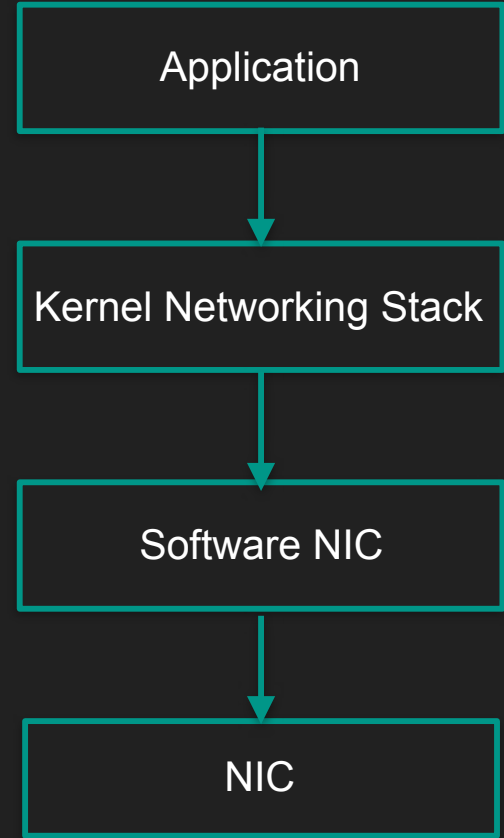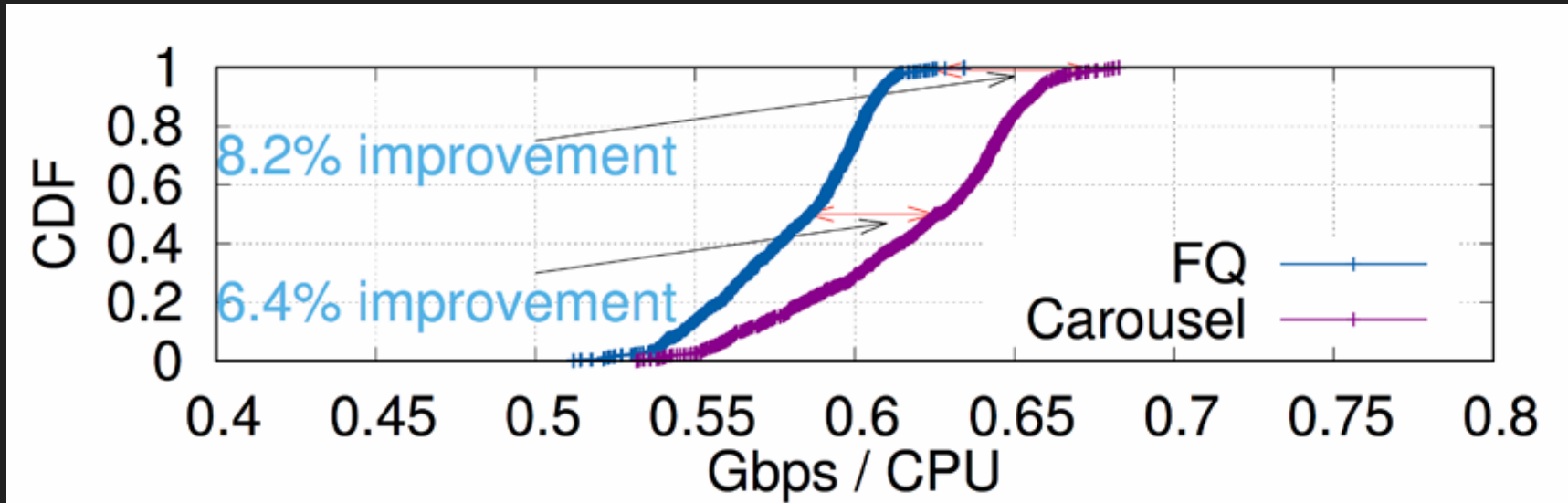
# Evaluation Setup

- Carousel deployed within a Software NIC

- Evaluation on Youtube servers comparing Carousel and FQ/Pacing

- Each server handles up to 50k sessions concurrently

```
┌─────────────────────────┐
│       Application        │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│ Kernel Networking Stack  │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│      Software NIC        │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│          NIC             │
└─────────────────────────┘
```
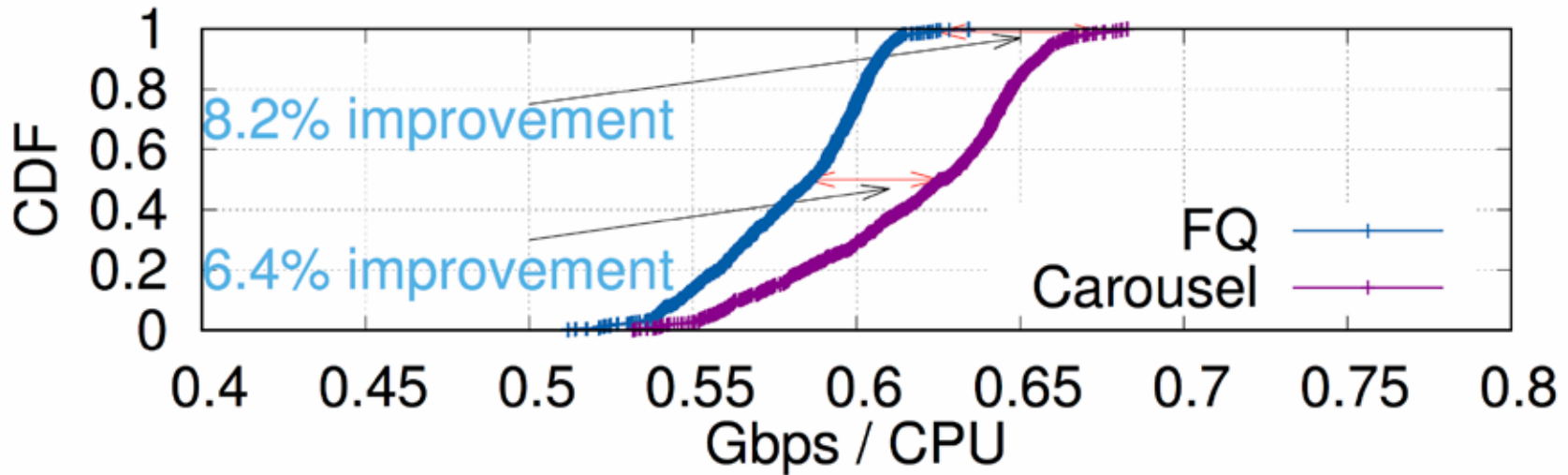
34

34

# Evaluation Metric

- Measures Gbps served per CPU utilization
  - Metric used is Gbps/CPU (higher is better)
  - Compare machines with similar CPU utilization
  - Measurements performed during peak 12-hours per day

- Evaluation is performed for:
  - Overall CPU utilization
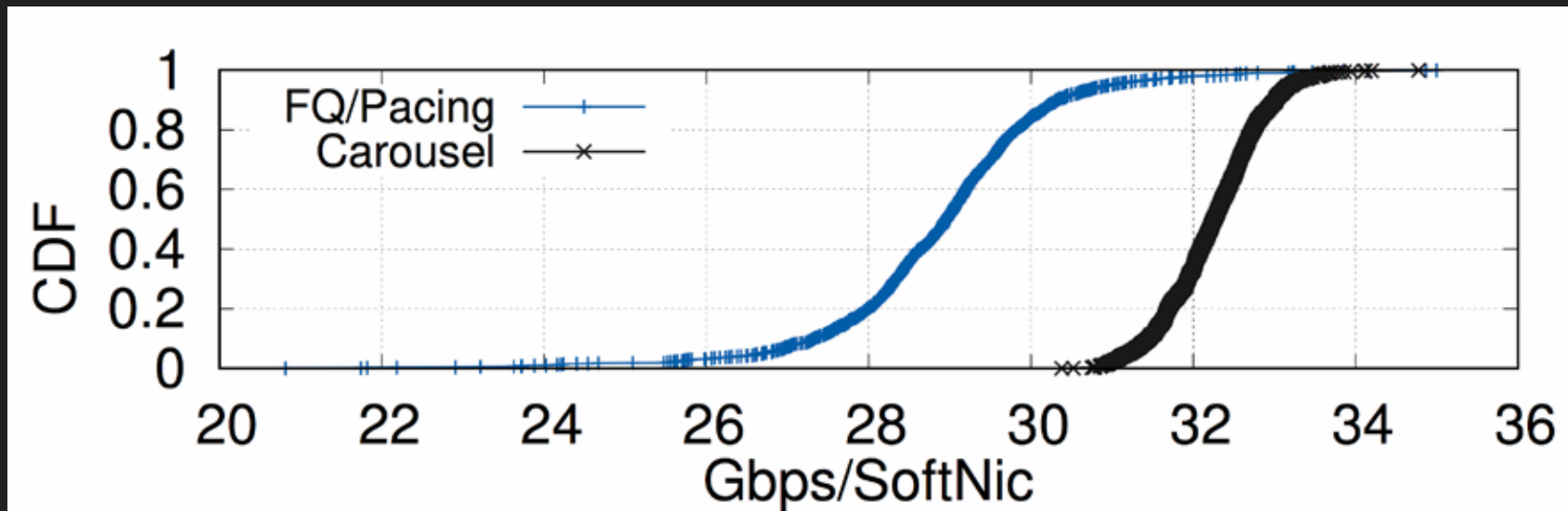  - Software NIC utilization

# Overall CPU Utilization
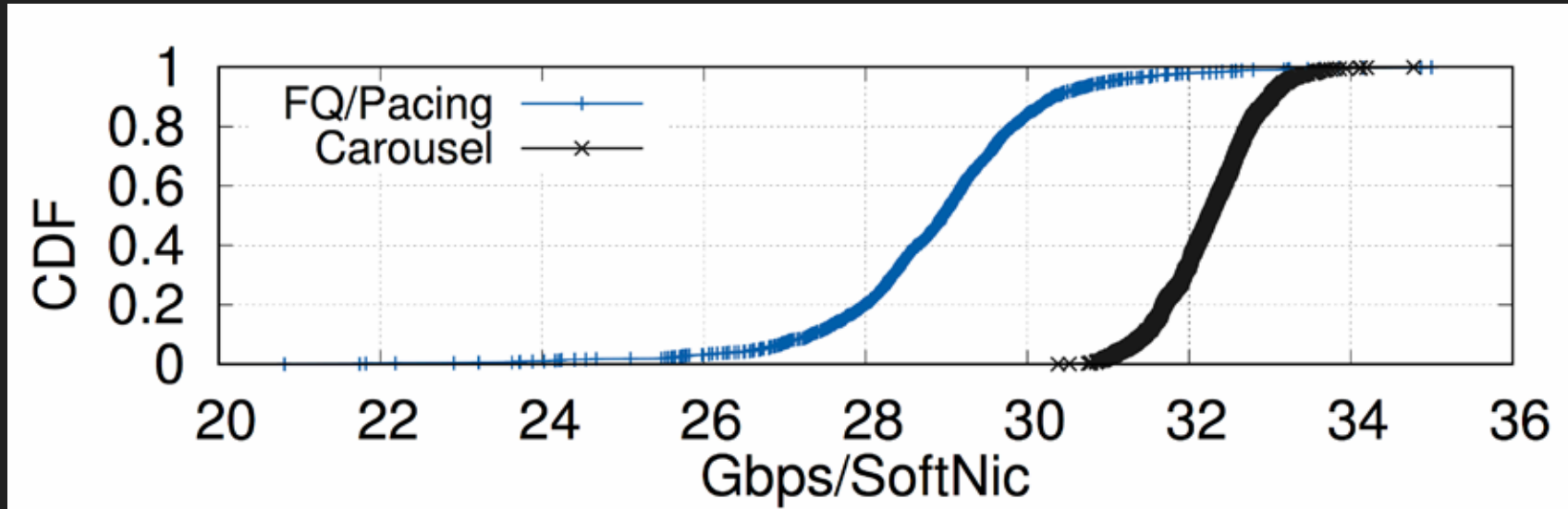
# Overall CPU Utilization



**Carousel saves up to 8.2% of overall CPU utilization
(5.9 cores on a 72 core machine)**
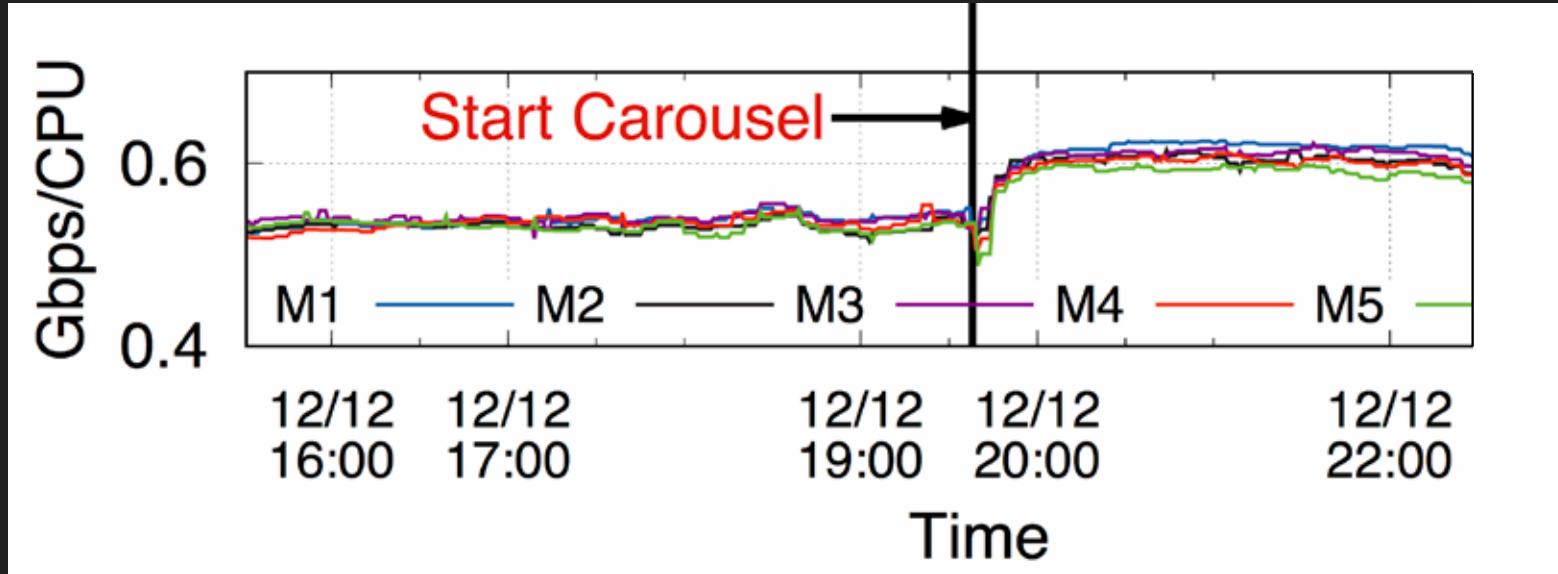
# SoftNIC Utilization

# SoftNIC Utilization



**Carousel improves even Software NIC utilization by 12% by increasing size of batches of packets enqueue in the Software NIC**

# Evaluation Summary



**Performance improvement when Carousel starts on 5 different machines**

# Conclusion

- Carousel allows networks operators for the first time to shape tens of thousands of flows individually

- Carousel advantages make a strong case for providing single-queue shaping and backpressure in kernel, userspace stacks, hypervisors, and hardware

# Questions?