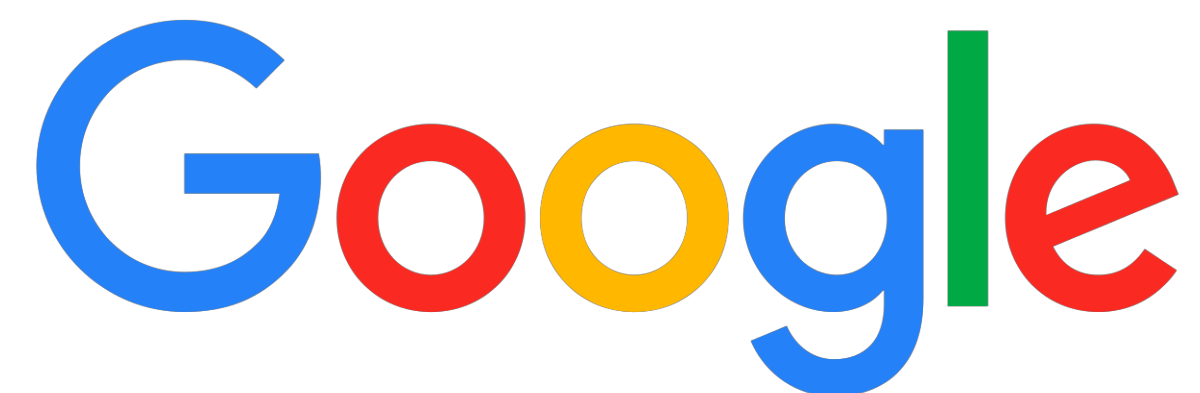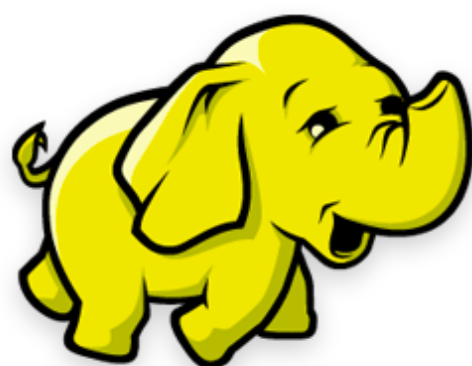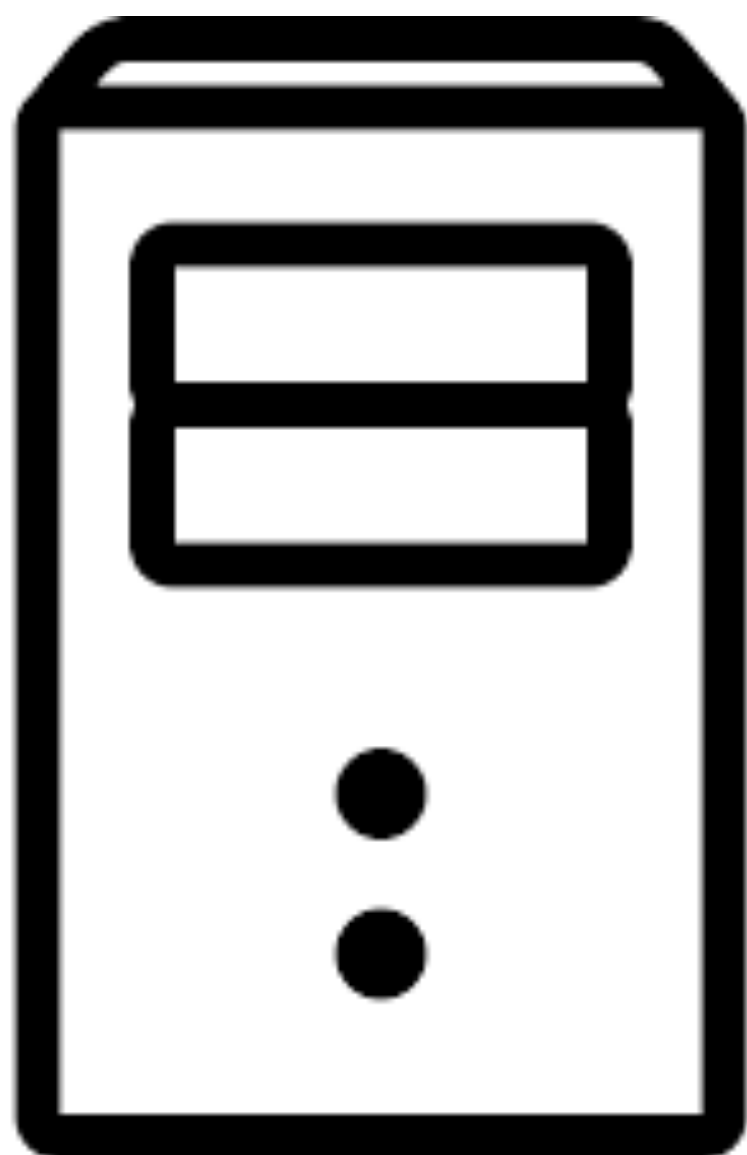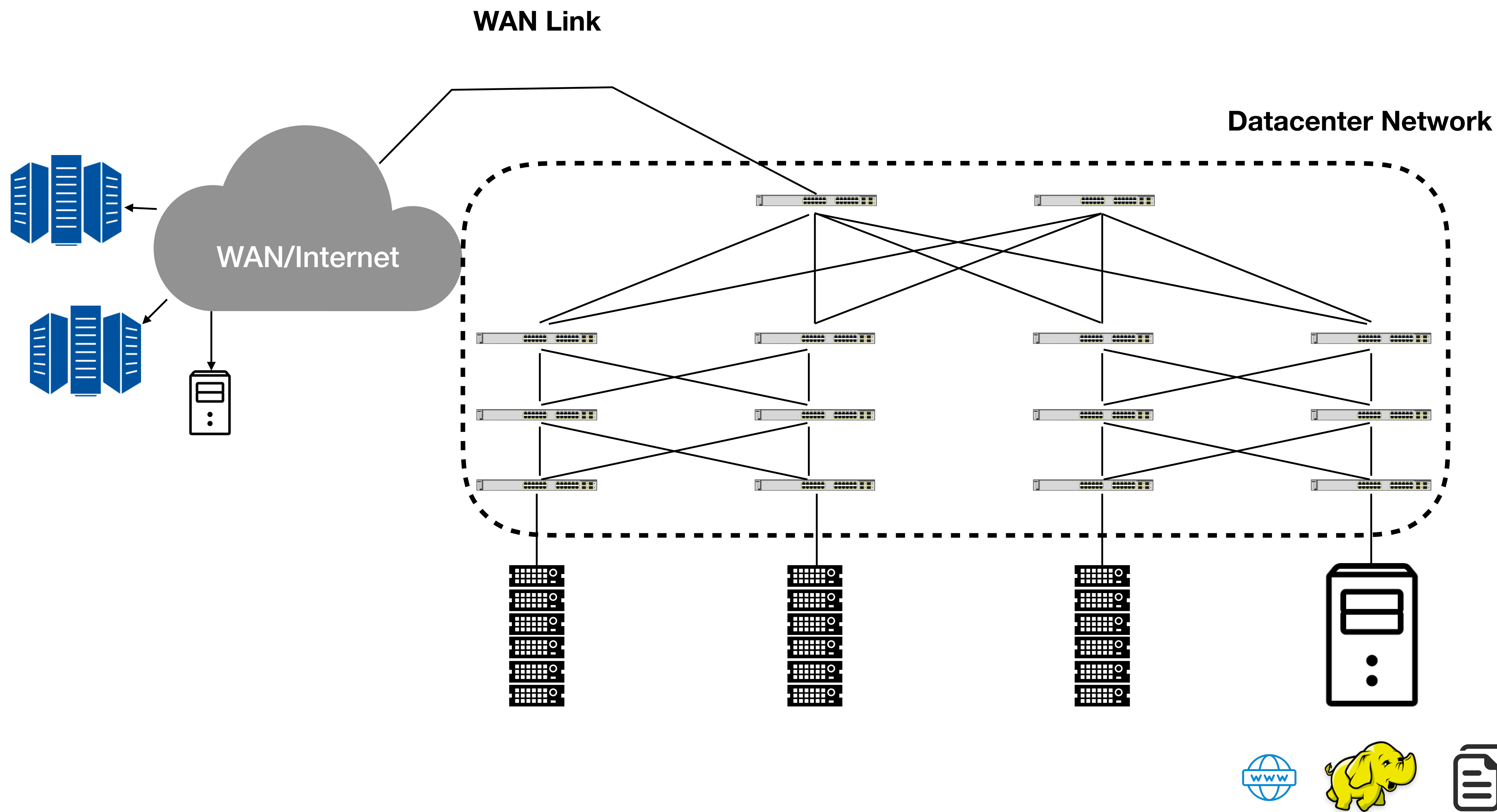# Eiffel: Efficient and Flexible Software Packet Scheduling

**Ahmed Saeed,** *Yimeng Zhao, Nandita Dukkipati,*
*Mostafa Ammar, Ellen Zegura, Khaled Harras, and Amin Vahdat*
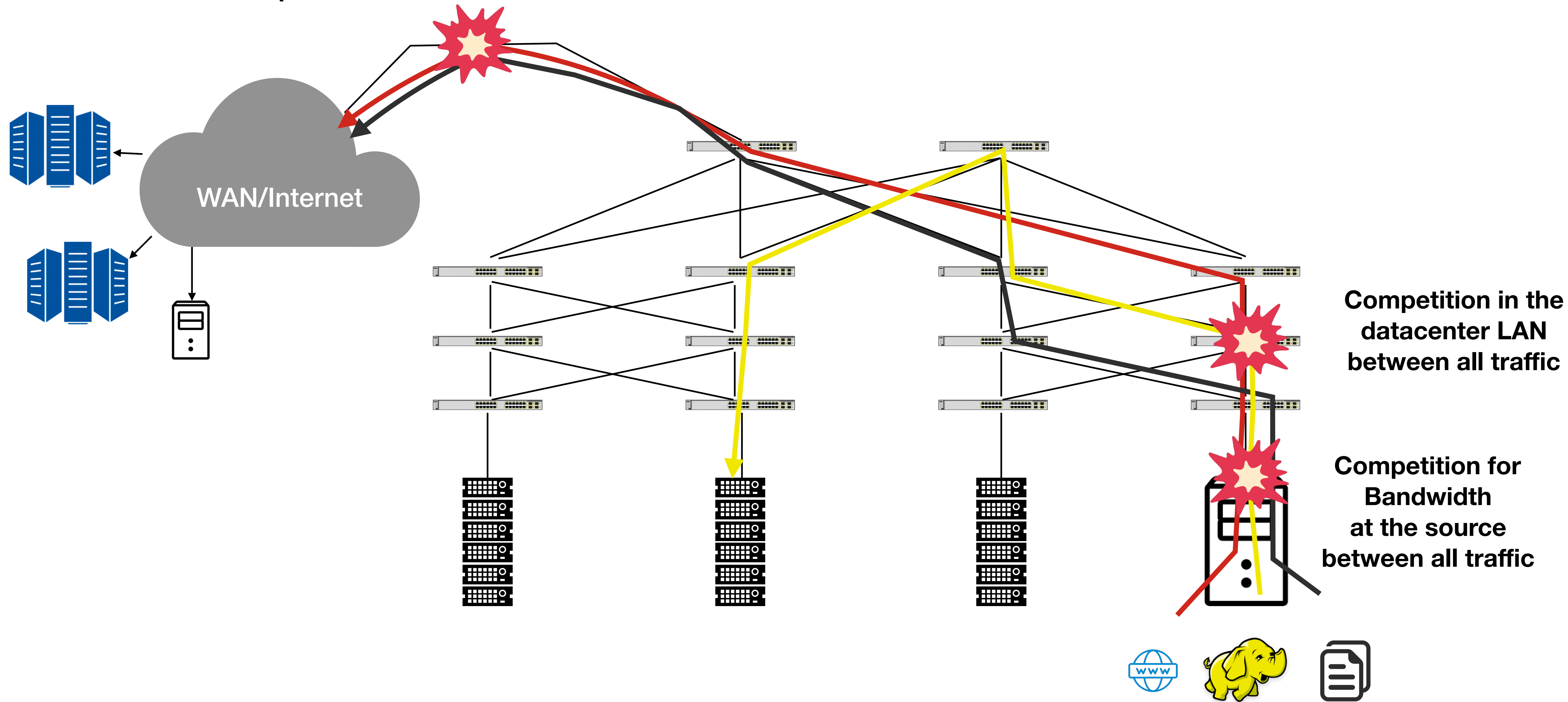
WAN Link
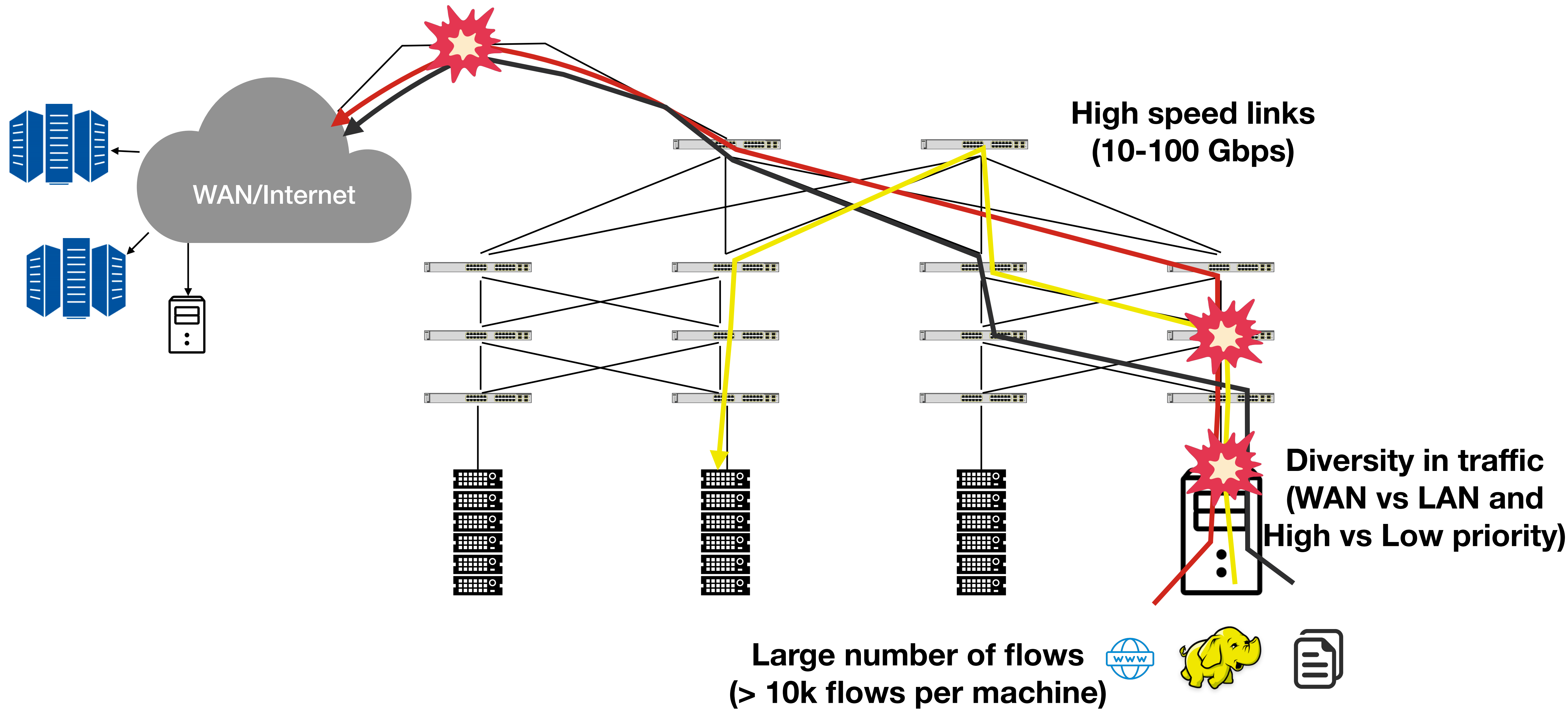
Datacenter Network

WAN/Internet

**Competition for Bandwidth
at premium links between WAN traffic**

**WAN/Internet**

**Competition in the
datacenter LAN
between all traffic**

**Competition for
Bandwidth
at the source
between all traffic**

High speed links
(10-100 Gbps)

Diversity in traffic
(WAN vs LAN and
High vs Low priority)

Large number of flows
(> 10k flows per machine)

7

# Packet Scheduling

- Scheduling determines the relative ordering as well as transmission time of packets in a queuing data structure with respect to some ranking function

- Packet scheduling implements policies to solve such problems

  - Traffic Isolation

  - Flow Completion Time Optimization

  - Congestion Control

# Scheduler Implementation

- Hardware schedulers in ASICs, FPGAs, or NPUs

  - Preprogrammed policies in switches or NICs

  - Programmable schedulers

- Software schedulers at end hosts or middleboxes

  - Kernel Queuing Disciplines (Qdiscs)

  - Userspace networking stacks

# Software vs Hardware

- Hardware lags behind network needs

- Software serves as a good experimental environment before hardware deployment

- Software provides a "build once, deploy many"

**Middlebox**

**End host**

# Challenges of Network Scheduling

- Accurate scheduling



- Efficient CPU and memory implementation   *O(log(n))* ⇒ *O(1)*

- Diversity of requirements

*Hierarchical Weighted Fair Queuing*

*Strict Priority*       *Rate Limiting*

*Shortest Remaining Time First*

***Objective: Design an accurate, efficient, and programmable software scheduler***

# Outline

- Eiffel Overview

- Characteristics of Packet Ranks

- Efficient Packet Ordering: Integer Priority Queues

- Scheduler Programmability

- Evaluation

# Eiffel Overview

# Eiffel Overview



- Efficient building block for packet sorting operating at line rate

- Expressive abstraction that can capture a wide range of policies

15

# Characteristics of Packet Ranks

# Ranks are Integers

- Packet carry limited precision integer
  priorities of width $w$ bits

  - QoS-based priority

  - Time-based priority

  - Flow size-based priority



TCP header format

# Ranks have Known Ranges

**Values of Interest**

$0$          $2^w$

- Semantics of priority values typical have limited ranges within the whole range of integer representation

  - Time-based priorities: from now to a few seconds in the future

  - Flow size-based priorities: values are known from typical application behavior

  - Strict priority ranges: policy/network operator defined

# Packets are Processed in Batches



Application

Application
Processing Speed

Network Processing

60B packet every 24 ns
1500B packet every 600 ns

NIC (20Gbps)

# Packets are Processed in Batches



**Application**

Application
Processing Speed

**Processing delay
is in 100s of ns**

**Network Processing**

60B packet every 24 ns
1500B packet every 600 ns

**NIC (20Gbps)**

# Packets are Processed in Batches



Application

**Application**

**Processing delay**

**Packets have to be processed in batches, rendering all packets in a batch to have virtually the same rank**

66B packet every 24 ns
**1500B packet every 600 ns**

NIC (20Gbps)

# Eiffel Building Block

**Bucketed Data
Structure**
+
**Limited number
of buckets**
+
**Algorithm to find min/max
non-empty bucket**

**= Integer Priority Queues**

# Efficient Packet Ordering: Integer Priority Queues

# Priority Queues 101

- Binary trees, Binomial Heap, Fibonacci Heap

- Support ExtractMin/ExtractMax

- Overhead of O(log n) on insertion or extraction

- Requires definition of a comparison operator: **Comparison-based Priority Queues**

# Integer Priority Queue

- Bucketed queues of N buckets

- Bucket index is the priority of elements in the bucket

- O(1) insertion and change priority

- O($Log_w$ N) ExtractMin/ExtractMax

| $P_0$ | $P_1$ | ... | ... | $P_N$ |
|-------|-------|-----|-----|-------|

# Integer Priority Queue

**Packets have known priority range and can be grouped into coarse granularity buckets**

- Bucketed queues of $\big(\!$ N buckets $\!\big)$

- Bucket index is the priority of elements in the bucket

- O(1) insertion and change priority

- O(Log $_w$ N) ExtractMin/ExtractMax

| $P_0$ | $P_1$ | ... | ... | $P_N$ |
|-------|-------|-----|-----|-------|

**Packets have integer priority are captured in limited precision integers**

# FFS-based Integer Priority Queue

# FFS-based Integer Priority Queue

- FindFirstSet (FFS) in a 64-bit word in 3 CPU cycles
  - Every bucket is represented by a bit
  - Bit is set iff bucket is not empty

- O(1) Integer Priority Queue in for N=64
  - Linux Real Time Process Scheduler
  - Quick Fair Queuing (QFQ)
    [F. Checconi et al. INFOCOM '13]

Count Leading Zeros          Count Trailing Zeros

0001 1111 … 1111 1100 0000

Find Highest Set          Find First Set

Count Leading Ones          Count Trailing Ones

1110000000…00000 1111111

Find Highest Zero          Find First Zero

# Hierarchical FFS-based Queue

Packets

| 0 | 1 | 2 | 3 | 4 | 5 |

Queue

Bitmap
Meta Data

| 1 | 0 | | 0 | 0 | | 1 | 1 |

Leaf

| 1 | 0 | | 1 | 0 |

| 1 | 1 |

Root

# Circular Hierarchical FFS-based Queue

# Circular Hierarchical FFS-based Queue

Primary                                          Secondary                          Packets

| 242 | 243 | 244 | 245 | 246 | 247 | | 248 | 249 | 250 | 251 | 252 | 253 | Queue

**cFFS-based queues has a small memory footprint and requires $O(\log_w N)$ steps for ExtractMin operating over a small N**

# Scheduler Programmability

# PIFO Programming Model

- Eiffel extends Push In First Out (PIFO) model

- PIFO model capture hierarchical policies using tress of priority queues [Sivaraman et. al SIGCOMM '16]

  - Packet ranking is performed on enqueue

  - Scheduling and shaping are tightly coupled in a single transaction

  - Implemented in hardware through parallel comparisons

# Eiffel Programming Model

- Eiffel model extends the PIFO model

  - Packets can be ordered based on flow ranking

  - Flows and packets can be ranked on enqueue and dequeue

  - Shaping and scheduling are decoupled for efficiency

# Eiffel Example: pFabric

- Each packet is tagged with Remaining Processing Time

- Packets are transmitted with *Shortest Remaining Processing Time First (SRPTF)*

- To avoid starvation, earliest packet from the highest priority flow is transmitted



- pFabric requires prioritizing flows based on ranks of packets

# Eiffel Example: pFabric

- Each packet is tagged with Remaining Processing Time

- Packets are transmitted with *Shortest Remaining Processing Time First (SRPTF)*

- To avoid starvation, earliest packet from the highest priority flow is transmitted

- pFabric requires prioritizing flows based on ranks of packets

# Eiffel Example: pFabric

- Each packet is tagged with Remaining Processing Time

- Packets are transmitted with *Shortest Remaining Processing Time First (SRPTF)*
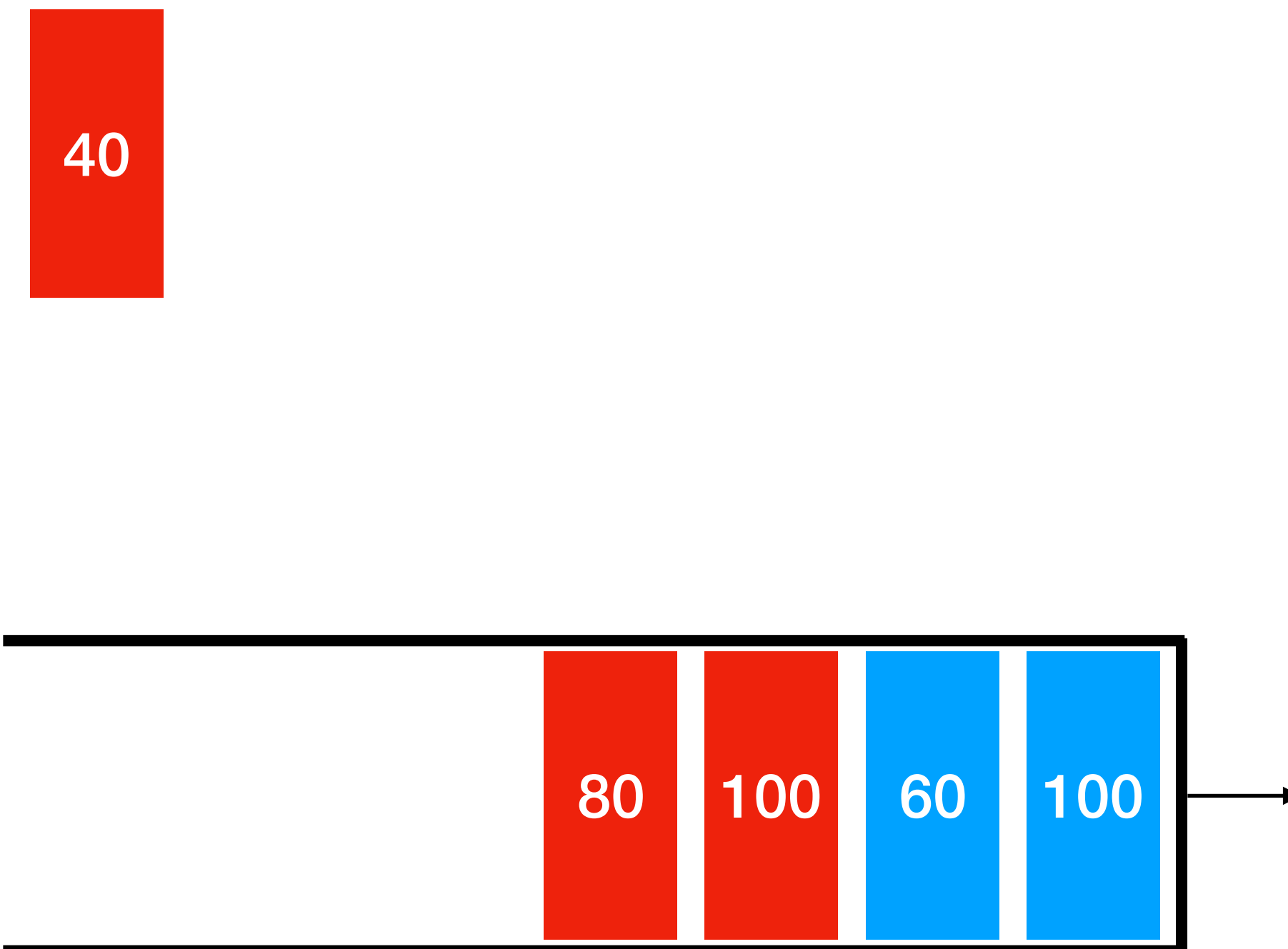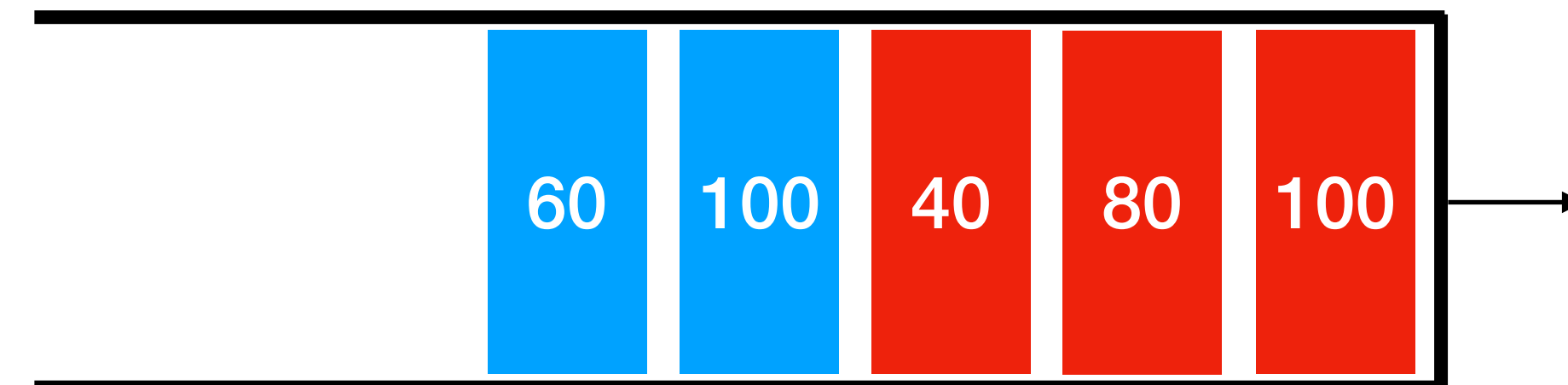
- To avoid starvation, earliest packet from the highest priority flow is transmitted

- pFabric requires prioritizing flows based on ranks of packets

| 60 | 100 | 40 | 80 | 100 |

# Eiffel Example: Implementation

**Rank 80**

**Rank 60**

- Data structures

  - Priority Queue per policy that ranks flows

  - FIFO queue per-flow

- On packet enqueue

  - Check packet tag and update flow rank

  - Update flow position in priority queue

| 80 |
|----|
| 100 |

| 60 |
|----|
| 100 |

| 20 | 40 | 60 | 80 | 100 |
|----|----|----|----|-----|

# Eiffel Example: Implementation

**Rank 40**   **Rank 60**

- Data structures
  - Priority Queue per policy that ranks flows
  - FIFO queue per-flow

- On packet enqueue
  - Check packet tag and update flow rank
  - Update flow position in priority queue

| 40 |
|----|
| 80 |
| 100 |

| 60 |
|----|
| 100 |

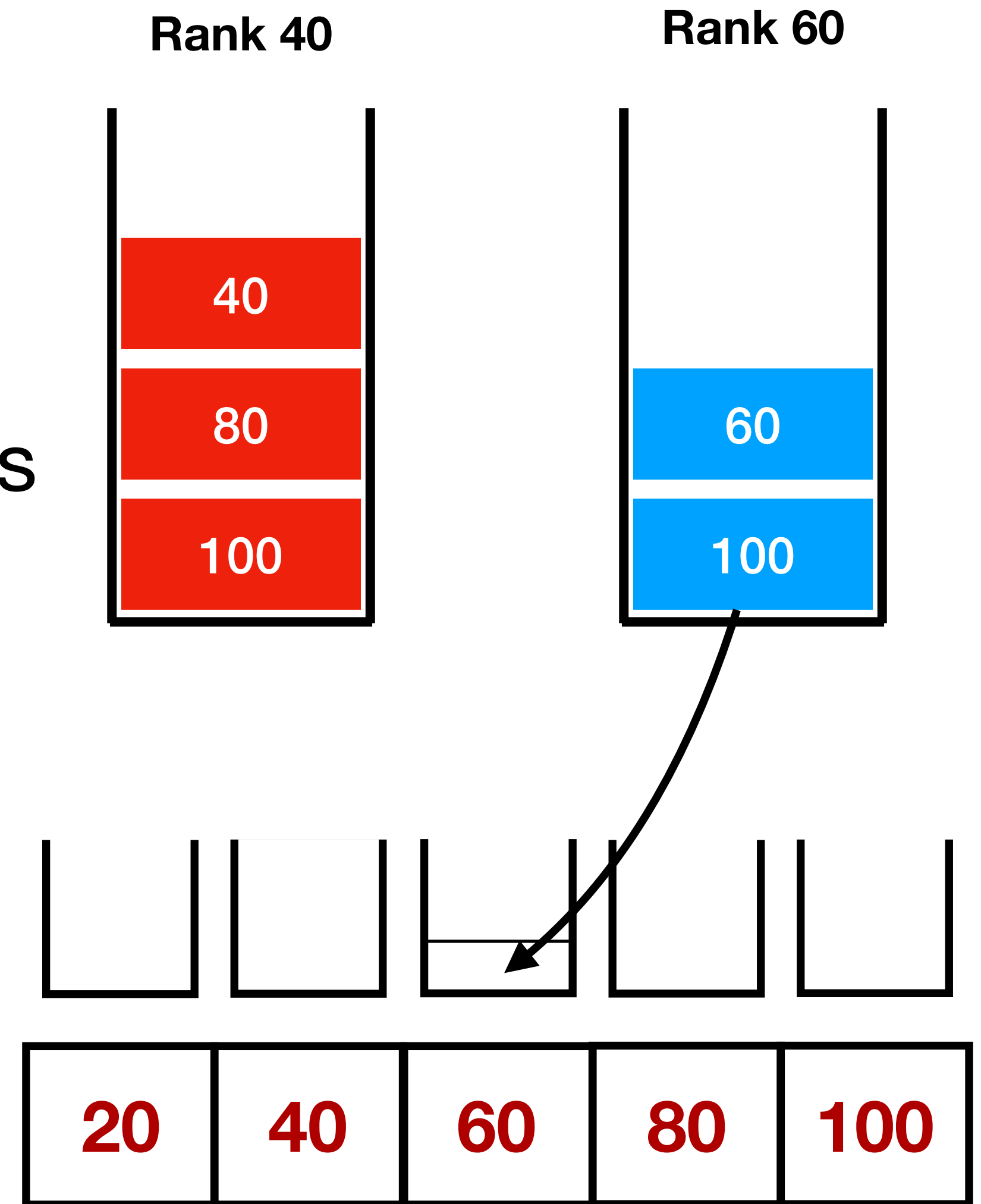| 20 | 40 | 60 | 80 | 100 |
|----|----|----|----|-----|

# Eiffel Example: Implementation

- ## Data structures

  - Priority Queue per policy that ranks flows

  - FIFO queue per-flow

- ## On packet enqueue

  - Check packet tag and update flow rank

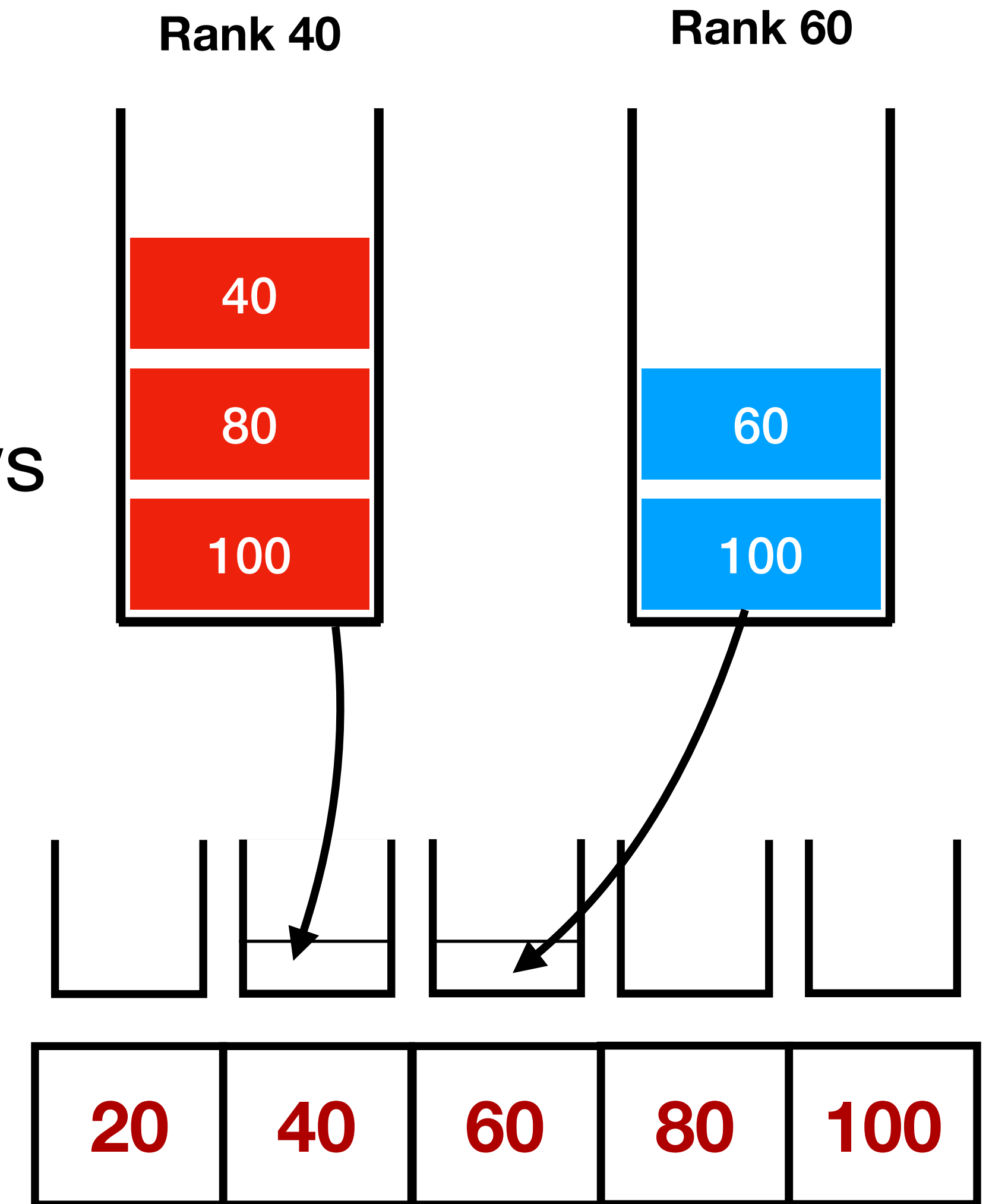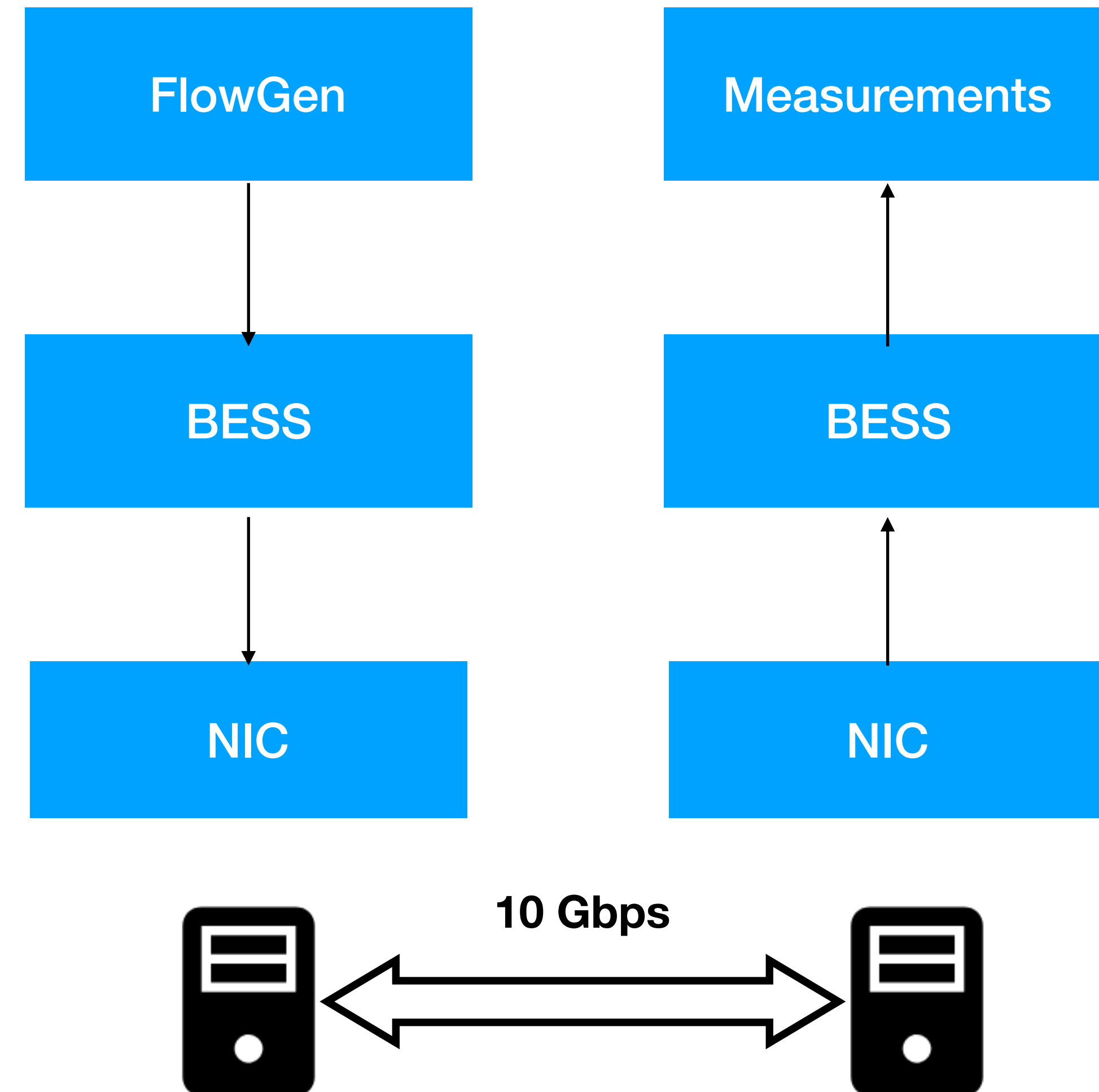  - Update flow position in priority queue

**Rank 40**     **Rank 60**

| 40 |
| 80 |
| 100 |

| 60 |
| 100 |

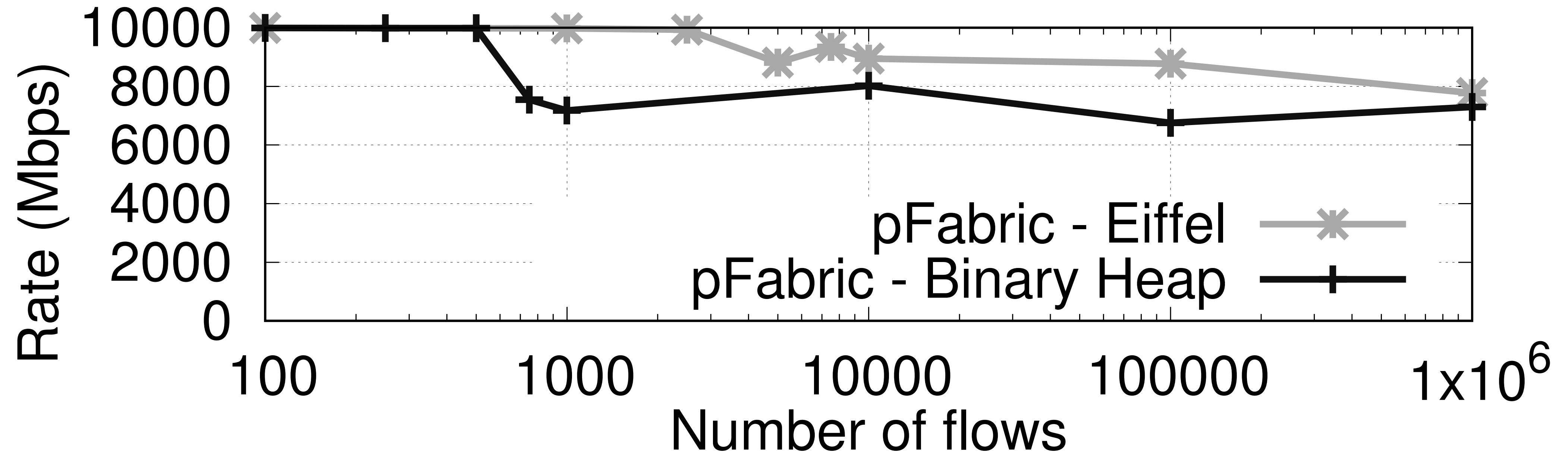| **20** | **40** | **60** | **80** | **100** |

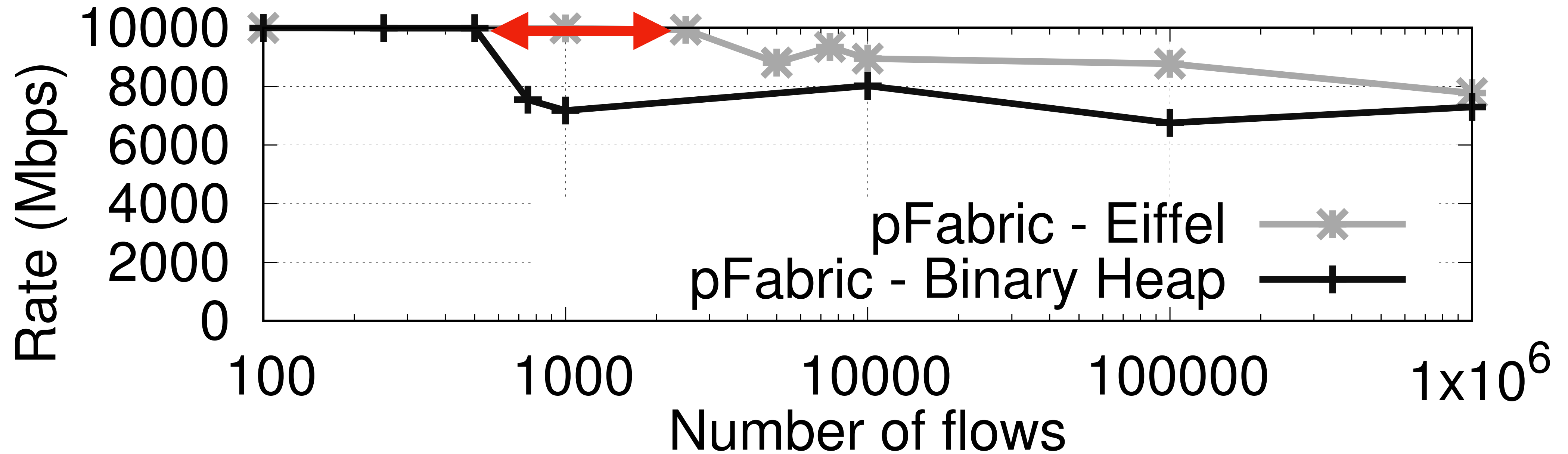# Evaluation

# Evaluation Setup

- Two servers with Intel X520-SR2 dual port NICs

- Eiffel implemented in Berkeley Extensible Software Switch (BESS)

- BESS runs on a single dedicated core

- Traffic generated using BESS FlowGen with varying number of flows and fixed 1500B packets

FlowGen

BESS

NIC

Measurements

BESS

NIC

**10 Gbps**

# Evaluation

# Evaluation



**Eiffel improves capacity by 5x in terms of number of flows that can be handled at line rate**

# Conclusion

- Eiffel network operators to deploy complex scheduling policies at end hosts and middle boxes

- Eiffel advantages make a strong case for rethinking the building blocks of packet in scheduling in hardware

# Questions?