# BBQ: A Fast and Scalable Integer Priority Queue for Hardware Packet Scheduling
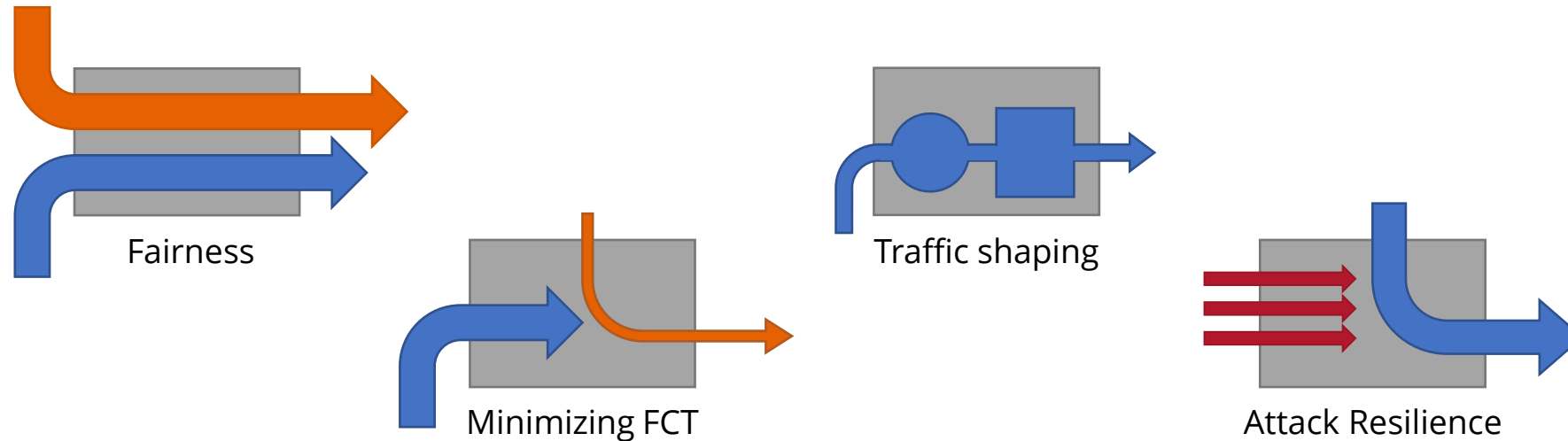
**Nirav Atre**, Hugo Sadok, Justine Sherry

*Carnegie Mellon University*

# Packet Scheduling in the Wild

Rich literature on packet scheduling algorithms optimizing for different performance objectives in various network settings
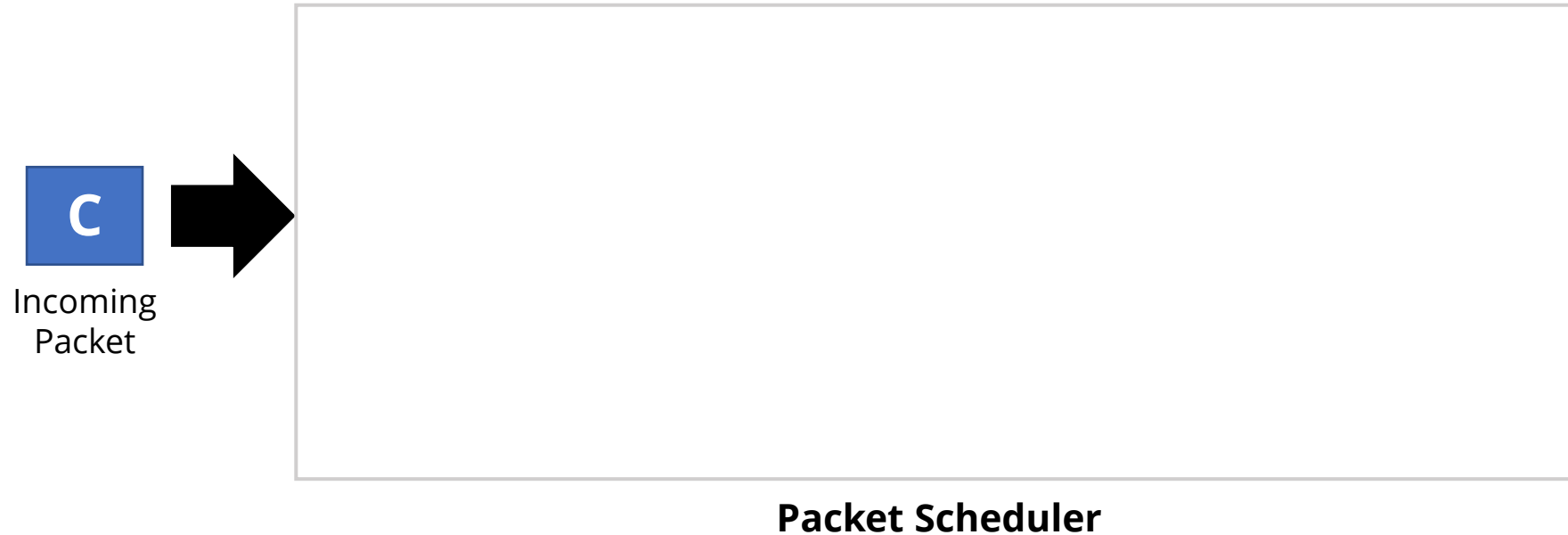
Fairness

Minimizing FCT

Traffic shaping

Attack Resilience

The key to deployment is **programmable hardware packet scheduling**

# *Programmable* Packet Scheduling

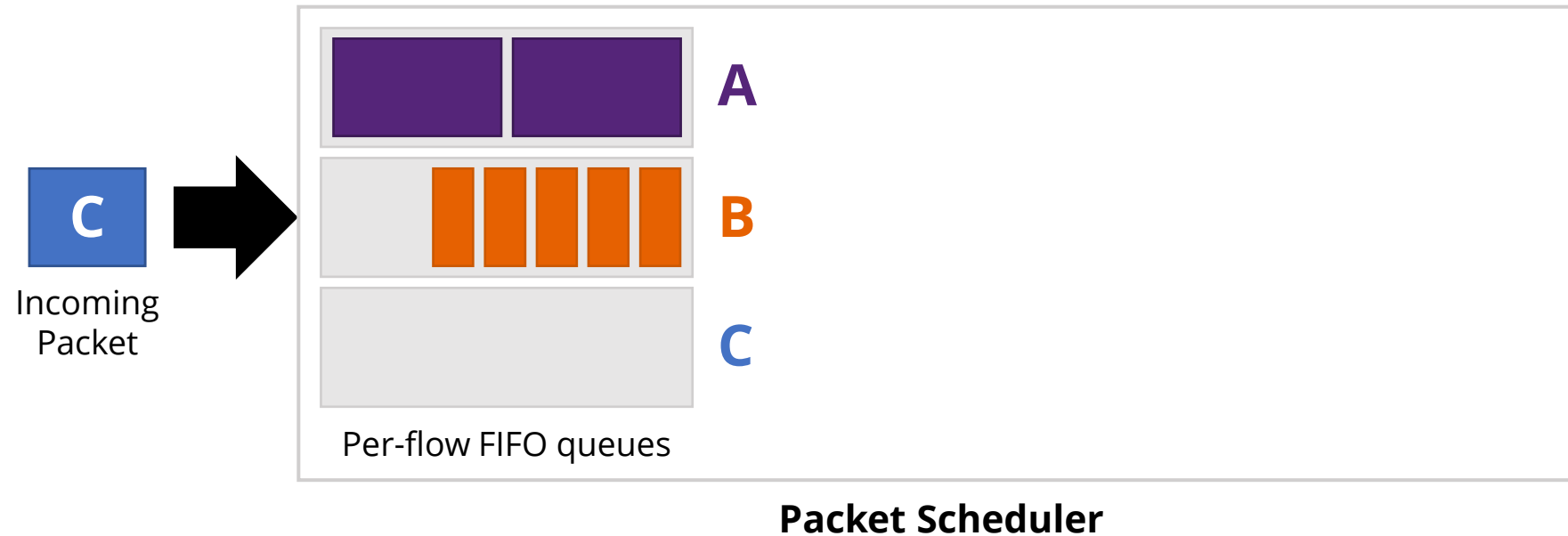Programmable Packet Scheduling at Line Rate [SIGCOMM '16]
→ **Push-In First-Out (PIFO)**



**Packet Scheduler**

# *Programmable* Packet Scheduling

Programmable Packet Scheduling at Line Rate [SIGCOMM '16]
→ **Push-In First-Out (PIFO)**



Packet Scheduler

# *Programmable* Packet Scheduling

Programmable Packet Scheduling at Line Rate [SIGCOMM '16]
→ **Push-In First-Out (PIFO)**



**Packet Scheduler**
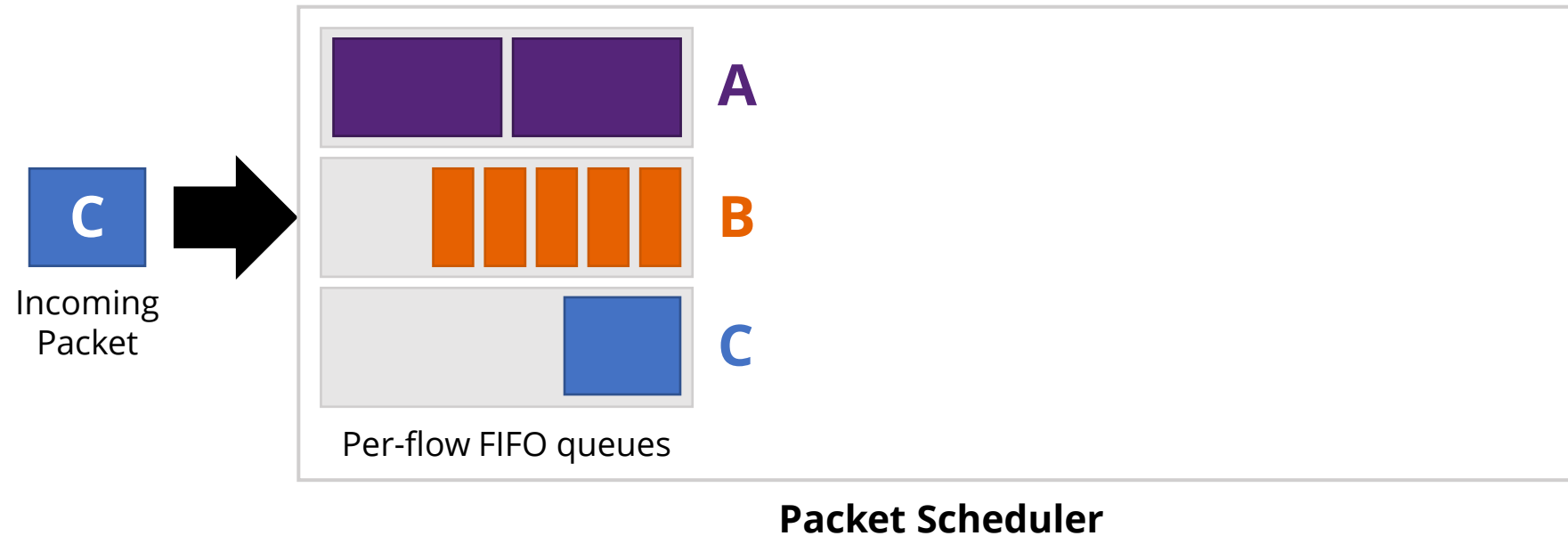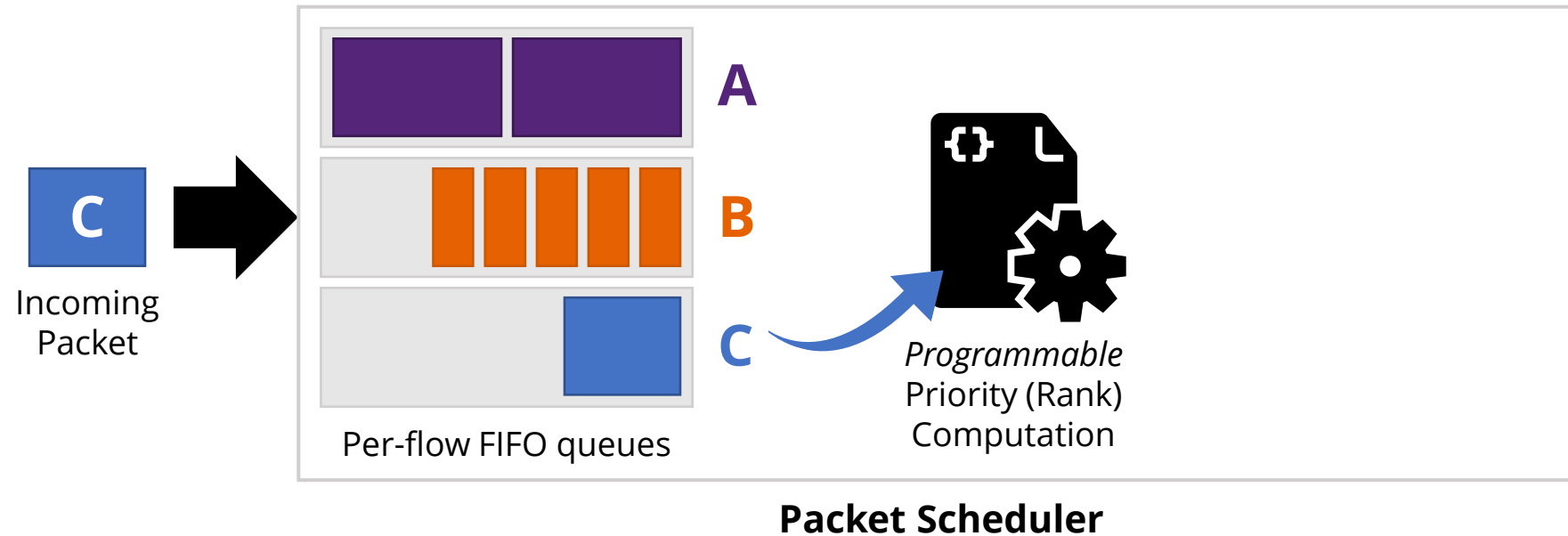
# *Programmable* Packet Scheduling

Programmable Packet Scheduling at Line Rate [SIGCOMM '16]
→ **Push-In First-Out (PIFO)**



Incoming Packet

A

B

C

Per-flow FIFO queues

*Programmable* Priority (Rank) Computation

**Packet Scheduler**
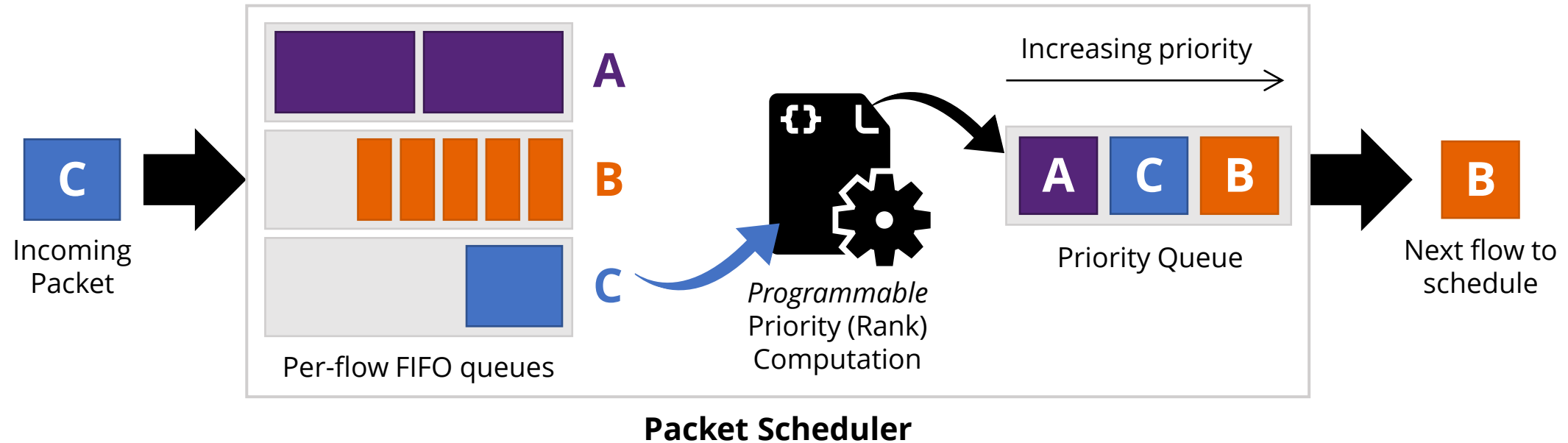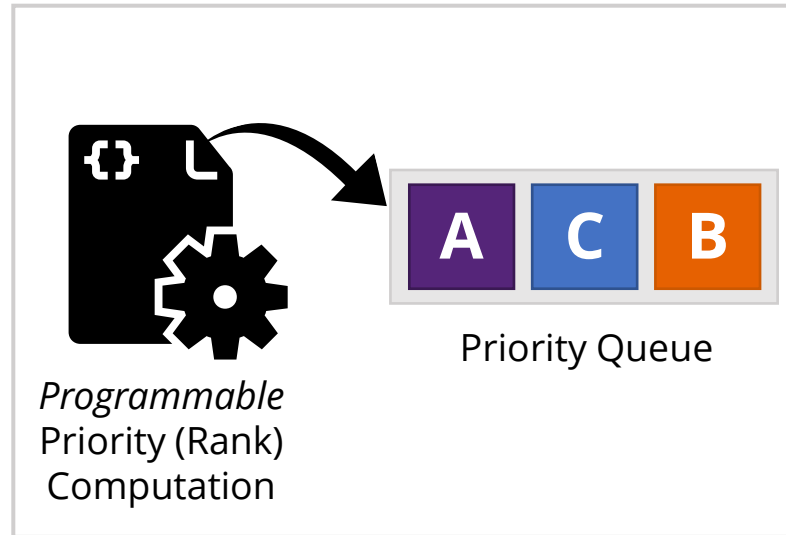
# *Programmable* Packet Scheduling

Programmable Packet Scheduling at Line Rate [SIGCOMM '16]
→ **Push-In First-Out (PIFO)**



Incoming Packet

Per-flow FIFO queues

A

B

C

*Programmable* Priority (Rank) Computation

Increasing priority

Priority Queue
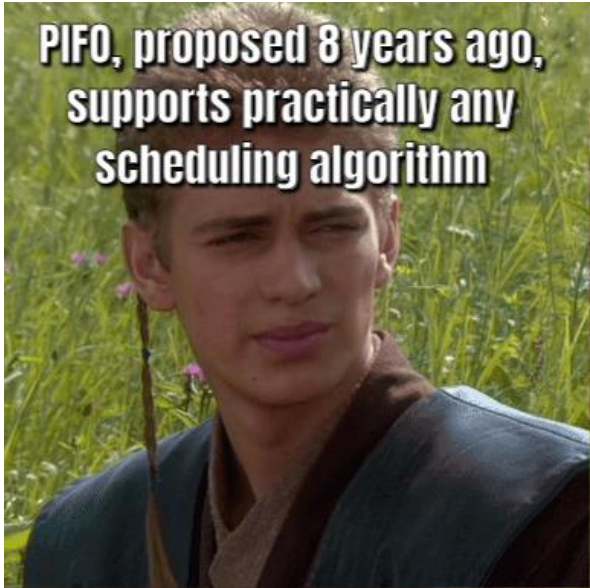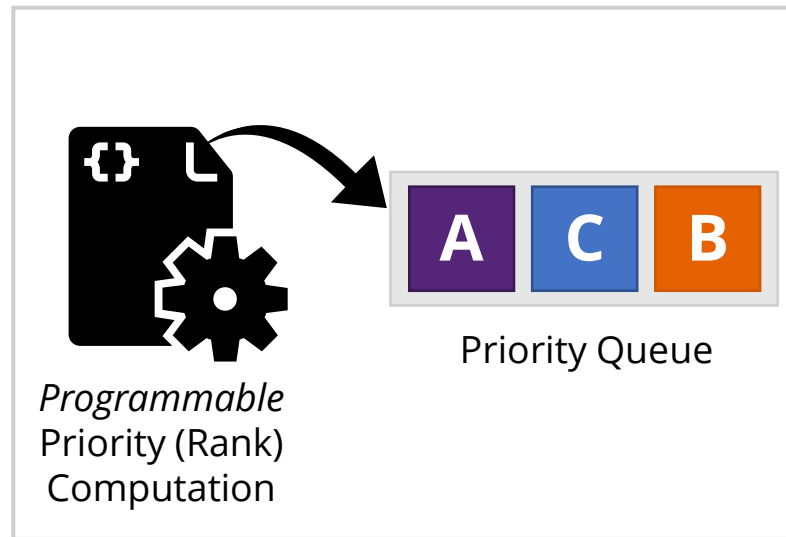
A C B

Next flow to schedule

**Packet Scheduler**

# *Programmable* Packet Scheduling

Programmable Packet Scheduling at Line Rate [SIGCOMM '16]
→ **Push-In First-Out (PIFO)**



**PIFO Abstraction**

# *Programmable* Packet Scheduling

Programmable Packet Scheduling at Line Rate [SIGCOMM '16]
→ **Push-In First-Out (PIFO)**



*Programmable*
Priority (Rank)
Computation

Priority Queue

**PIFO Abstraction**

At the heart of PIFO is a
**hardware priority queue** that
provides, at minimum, enqueue
and dequeue-max functionality

PIFO's vision is hampered by throughput, scalability, and resource
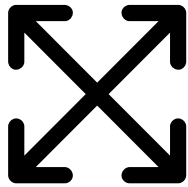overhead issues associated with existing priority queue designs

# This Talk

- Minimum requirements for scheduling in switches and SmartNICs
- State-of-the-art priority queue designs are infeasible
- How do we get there?
- Evaluation

# This Talk

- **Minimum requirements for scheduling in switches and SmartNICs**
- State-of-the-art priority queue designs are infeasible
- How do we get there?
- Evaluation

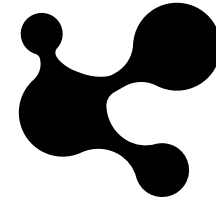# Minimum requirements for scheduling in switches and SmartNICs

**Flow Count Scalability**

Support flow counts representative of modern networks: **O(100K)**

**Single-Instance Performance**

Sustain packet rates corresponding to today's line rates: **100Gbps+ (148.8 Mpps)**
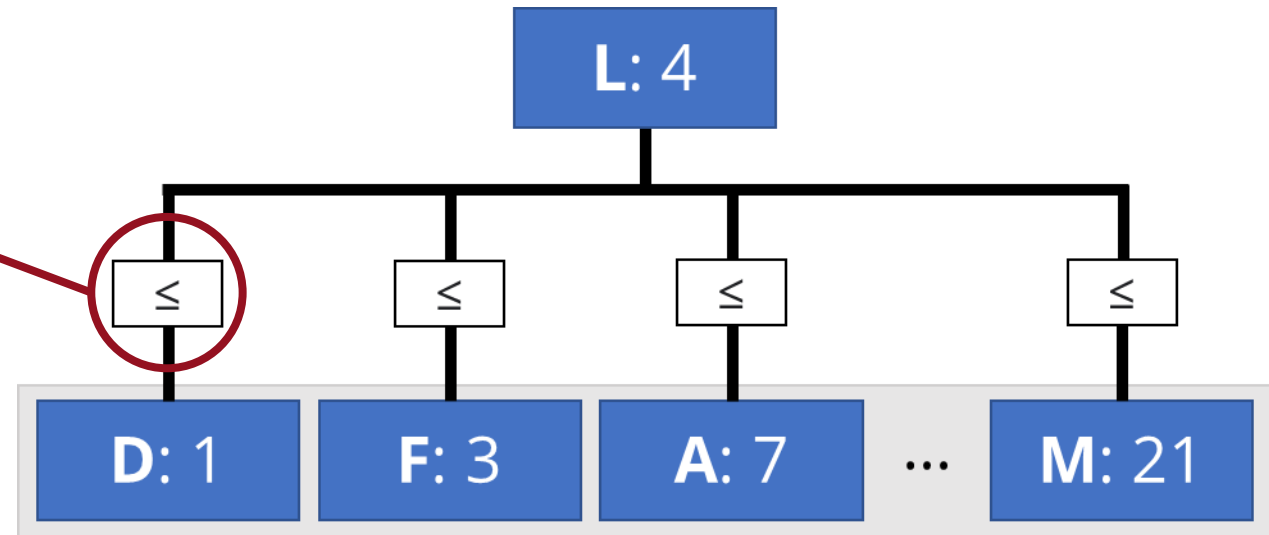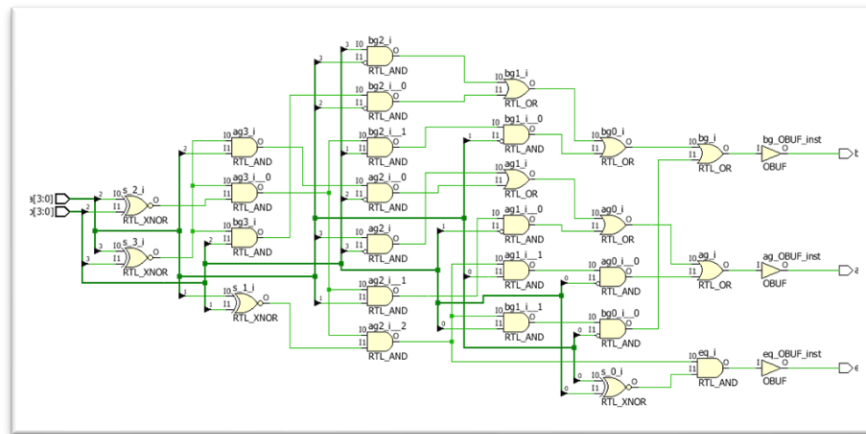
**Logical Partitioning**

Statistically multiplex a single, **physical** priority queue between many independent **logical** priority queues

# This Talk

- Minimum requirements for scheduling in switches and SmartNICs
- **State-of-the-art priority queue designs are infeasible**
- How do we get there?
- Evaluation

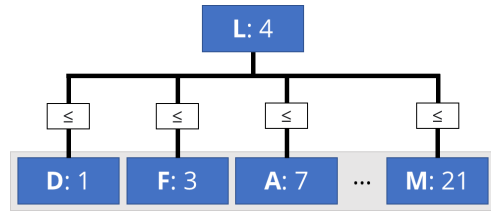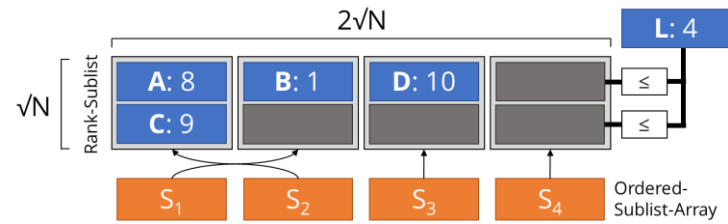# Existing designs are infeasible

**PIFO** [SIGCOMM '16]



N comparators followed by priority decoding to decide where to insert the next entry → supports at most 4K flows
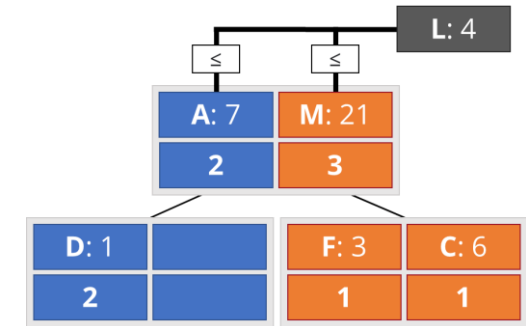
Poor Scalability

# Existing designs are infeasible



**PIFO** [SIGCOMM '16]

**PIEO** [SIGCOMM '19]
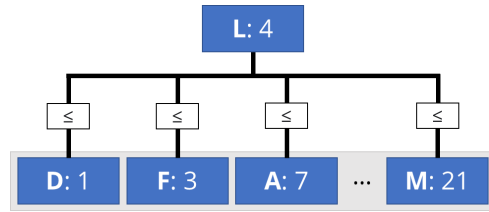
**BMW-Tree** [SIGCOMM '23]

Poor Scalability

**4K** flows

**64K** flows

**100K+** flows

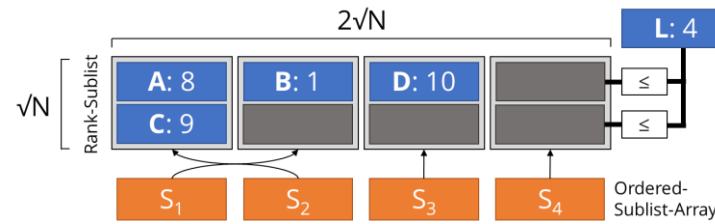**Flow Count Scalability**

# Existing designs are infeasible



**PIFO** [SIGCOMM '16]

Poor Scalability

**4K** flows



**PIEO** [SIGCOMM '19]

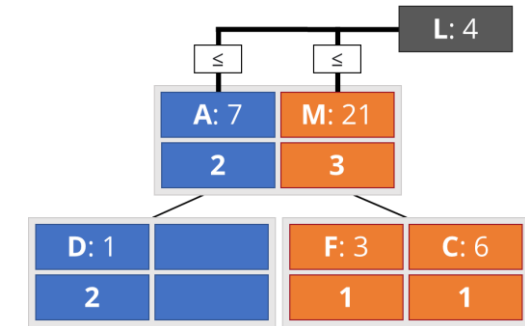Poor Performance

**64K** flows

15 Mpps (10% of line rate at 100 Gbps) on an **FPGA SmartNIC**



**BMW-Tree** [SIGCOMM '23]

No Logical Partitioning

**100K+** flows

1.5 – 6X chip area to implement on a 32-port **Switch ASIC**

# This Talk

- Minimum requirements for scheduling in switches and SmartNICs
- State-of-the-art priority queue designs are infeasible
- **How do we get there?**
- Evaluation

# Key Idea

If the priority span is bounded, we can achieve all 3 properties (**scalability**, **performance**, and **logical partitioning**) using non-comparison-based sorting.

# Integer Priority Queueing (IPQ)



*Priority buckets*

How to perform **dequeue-max**?
- Iteratively checking each bucket is slow!

# Integer Priority Queueing (IPQ)



How to perform **dequeue-max**?
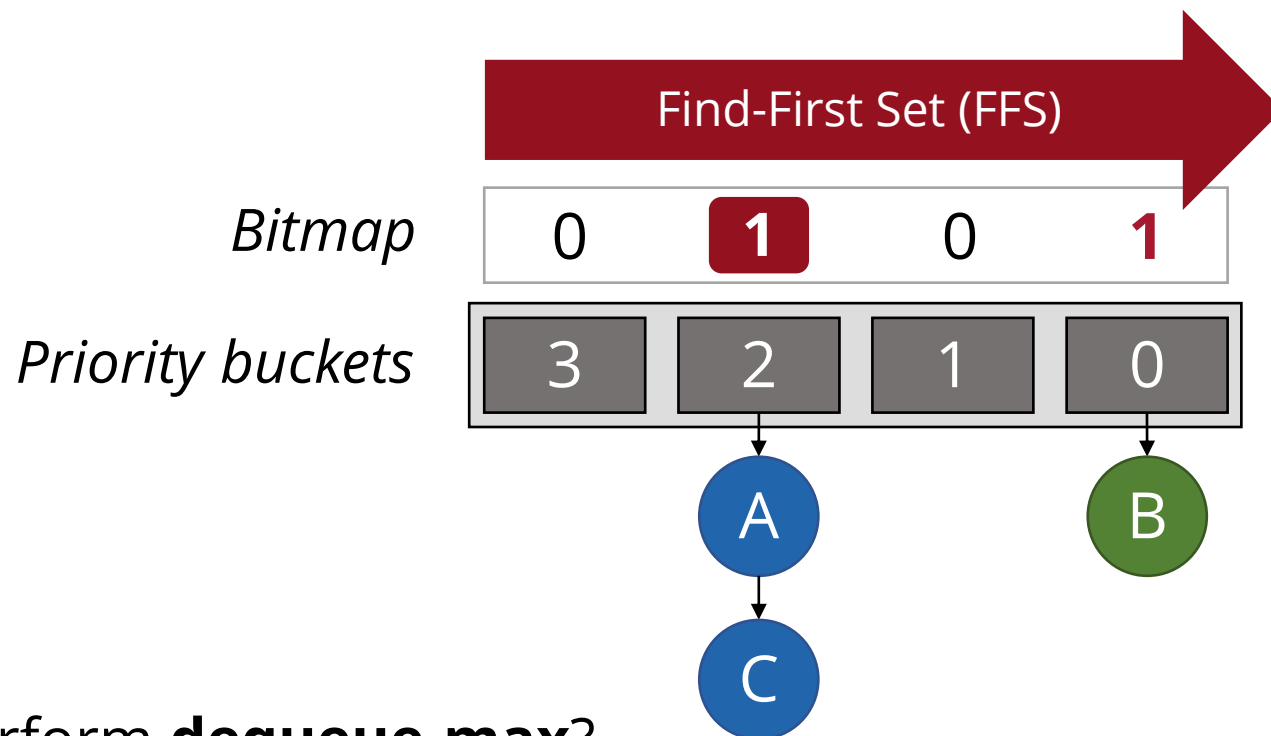- Augment with a *bitmap* encoding bucket occupancy, then use *Find-First Set*

# Integer Priority Queueing (IPQ)



How to perform **dequeue-max**?
- Augment with a *bitmap* encoding bucket occupancy, then use *Find-First Set*

# Integer Priority Queueing (IPQ)



Bitmap | 0 ... | | 1

Priority buckets | 32K ... | 3 | 2 | 1 | 0

C          B

What if we need to support a **huge number of priorities** (*e.g.*, 32K)?
- Can't do FFS on a 32K-bit bitmap

# Integer Priority Queueing (IPQ)

# Integer Priority Queueing (IPQ)

# Integer Priority Queueing (IPQ)



Bit is set if **any priority bucket in this subtree** is not empty

Bitmap Tree

Bit is set if **priority bucket** is not empty

Priority buckets

# Integer Priority Queueing (IPQ)

# Integer Priority Queueing (IPQ)



Data-structure is called a **Hierarchical Find-First Set (HFFS) Queue**

# Bitmapped Bucket Queue (BBQ)

Data-structure is called a **Hierarchical Find-First Set (HFFS) Queue**

Many software systems use FFS-based priority queueing (*e.g.*, Linux scheduler for process scheduling, and Eiffel [NSDI '19] for packet scheduling).

Our insight is that this data-structure is amenable to a *high-performance, **fully-pipelined** hardware implementation*.

# ⤢ Scalability comes for free

BBQ uses an IPQ-based design, breaking the dependence between queue size and run-time complexity of operations.

*Scalability "falls out" of this high-level design choice.*

# ⏱ BBQ is *highly optimized* for performance

(1) Need a high **operating frequency** ($f_{max}$)

(2) Need to maximize **operations-per-cycle** (OPC)

Performance (ops/sec) = $f_{max}$ (cycles/sec) × OPC (ops/cycle)

# 🎯 BBQ is *highly optimized* for performance

(1) Need a high **operating frequency** ($f_{max}$)

> BBQ achieves this by using a <u>deep pipeline</u> where individual stages are designed to do both *little* and *roughly equal* amounts of work.

(2) Need to maximize **operations-per-cycle** (OPC)

$$\text{Performance (ops/sec)} = f_{max} \text{ (cycles/sec)} \times \text{OPC (ops/cycle)}$$
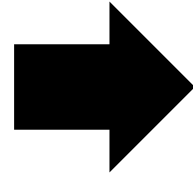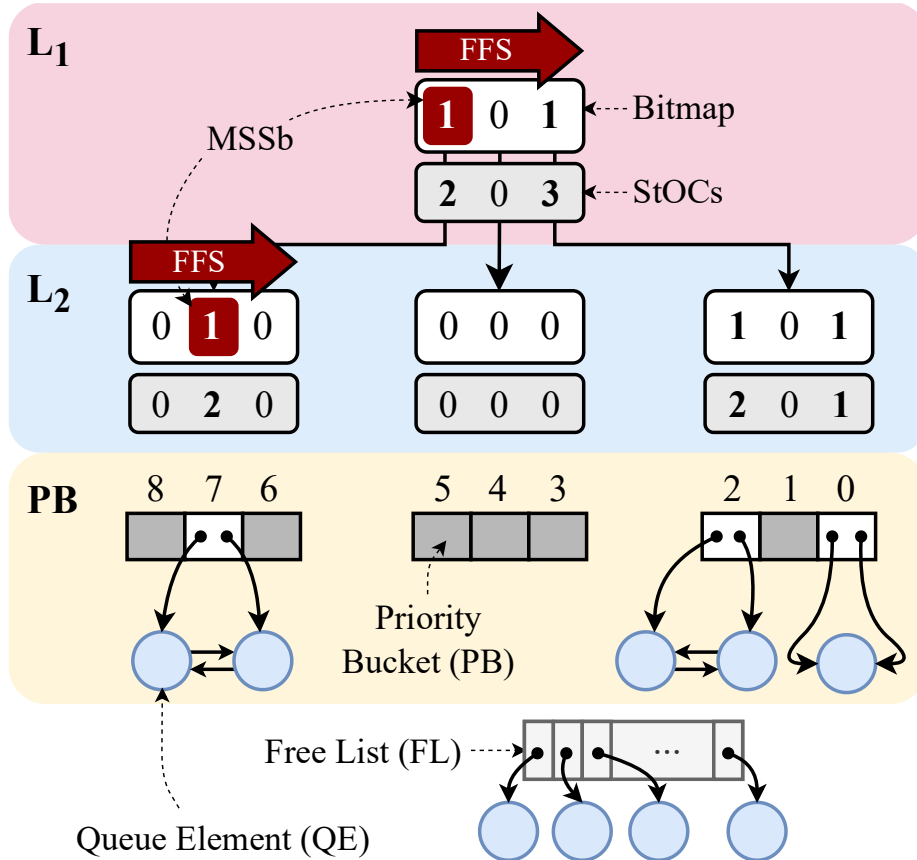
# BBQ is *highly optimized* for performance



| Cycle | Description | PHR |
|---|---|---|
| 0 | Register inputs<br><br>If ENQUEUE:<br>$F \leftarrow$ FreeList.POP() // Pop free list | |
| 1 | Compute $L_1$ bitmap index<br>↳ Read the corresponding $L_1$ StOC | $L_1$ |
| 2 | Compute, Write: $L_1$ StOC ⤳ $L_1$ bitmap<br>Read $L_2$ bitmap | |
| 3 | // Read delay for $L_2$ bitmap | |
| 4 | Compute $L_2$ bitmap index<br>↳ Read the corresponding $L_2$ StOC | $L_2$ |
| 5 | // Read delay for $L_2$ StOC | |
| 6 | Compute, Write: $L_2$ StOC ⤳ $L_2$ bitmap<br>Read the corresponding PB | |
| 7 | // Read delay for PB | |
| 8 | If DEQUEUE:<br>(a) Read $X \leftarrow$ QE_DATA[PB.TAIL$^{new}$]<br>(b) Read $Y \leftarrow$ QE_PREV[PB.TAIL$^{new}$] | |
| 9 | // Read delay for QE_DATA and QE_PREV | |
| 10 | If ENQUEUE: // Enqueue at the HEAD<br>(a) QE_DATA[$F$] $\leftarrow$ Data to enqueue<br>(b) Write QE_NEXT[$F$] $\leftarrow$ PB.HEAD<br>(c) Write QE_PREV[PB.HEAD] $\leftarrow$ F<br>(d) Write PB.HEAD $^{new}$ $\leftarrow$ F<br><br>If DEQUEUE: // Dequeue from TAIL<br>(a) FreeList.PUSH(PB.TAIL)<br>(b) Write PB.TAIL$^{new}$ $\leftarrow$ Y<br>(c) Output $X$ | PB |

Deep (11-stage) hardware pipeline for a 2-level BBQ

# 🏎 BBQ is *highly optimized* for performance

(1) Need a high **operating frequency** ($f_{max}$)

> BBQ achieves this by using a <u>deep pipeline</u> where individual stages are designed to do both *little* and *roughly equal* amounts of work
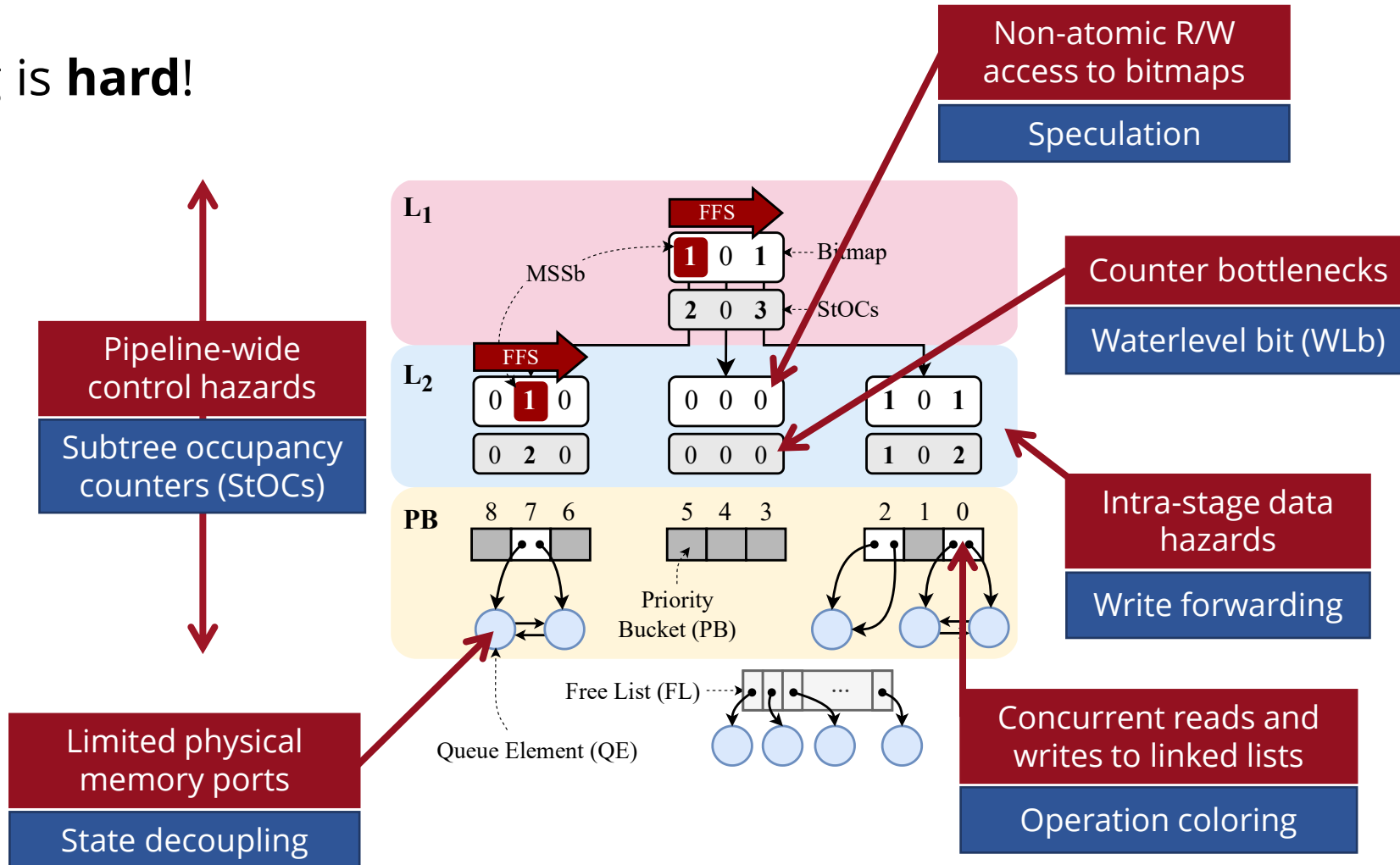
(2) Need to maximize **operations-per-cycle** (OPC)

> BBQ realizes a <u>fully-pipelined</u> architecture (OPC of 1) agnostic of workload

$$\text{Performance (ops/sec)} = f_{max} \text{ (cycles/sec)} \times \text{OPC (ops/cycle)}$$

# 🎛️ BBQ is *highly optimized* for performance

Pipelining is **hard**!



Non-atomic R/W access to bitmaps

Speculation

Counter bottlenecks

Waterlevel bit (WLb)

Pipeline-wide control hazards

Subtree occupancy counters (StOCs)

Intra-stage data hazards

Write forwarding

Limited physical memory ports

State decoupling

Concurrent reads and writes to linked lists

Operation coloring

# BBQ supports logical partitioning with zero fragmentation and performance overhead



Priority buckets

Logical BBQs

# BBQ supports logical partitioning with zero fragmentation and performance overhead



Steering

Priority buckets
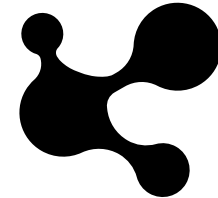
# How does BBQ meet our requirements?

### Flow Count Scalability

**IPQ-based design** allows scaling to O(100K) flows

### Single-Instance Performance

**Highly optimized, fully-pipelined design** allows BBQ to support 150 Mpps (100 Gbps) on FPGAs and 1.5 Bpps (1 Tbps) on switch ASICs
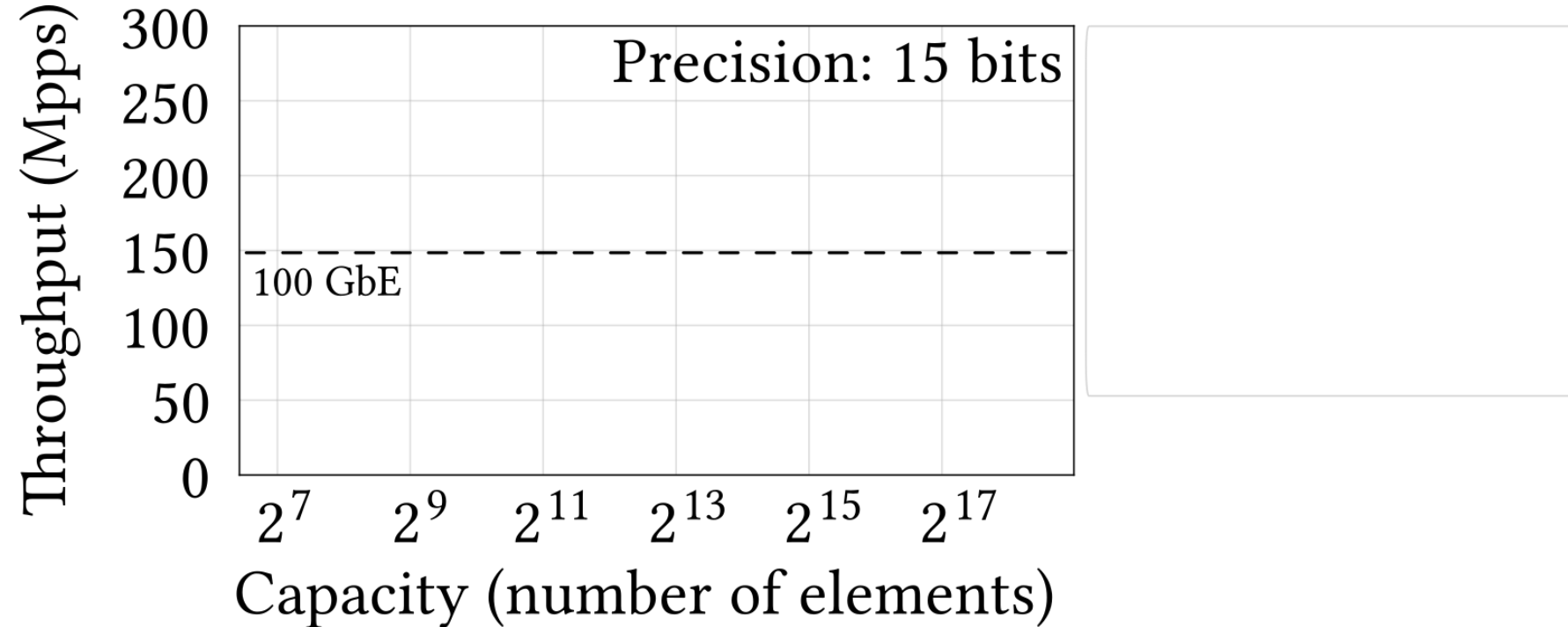
### Logical Partitioning

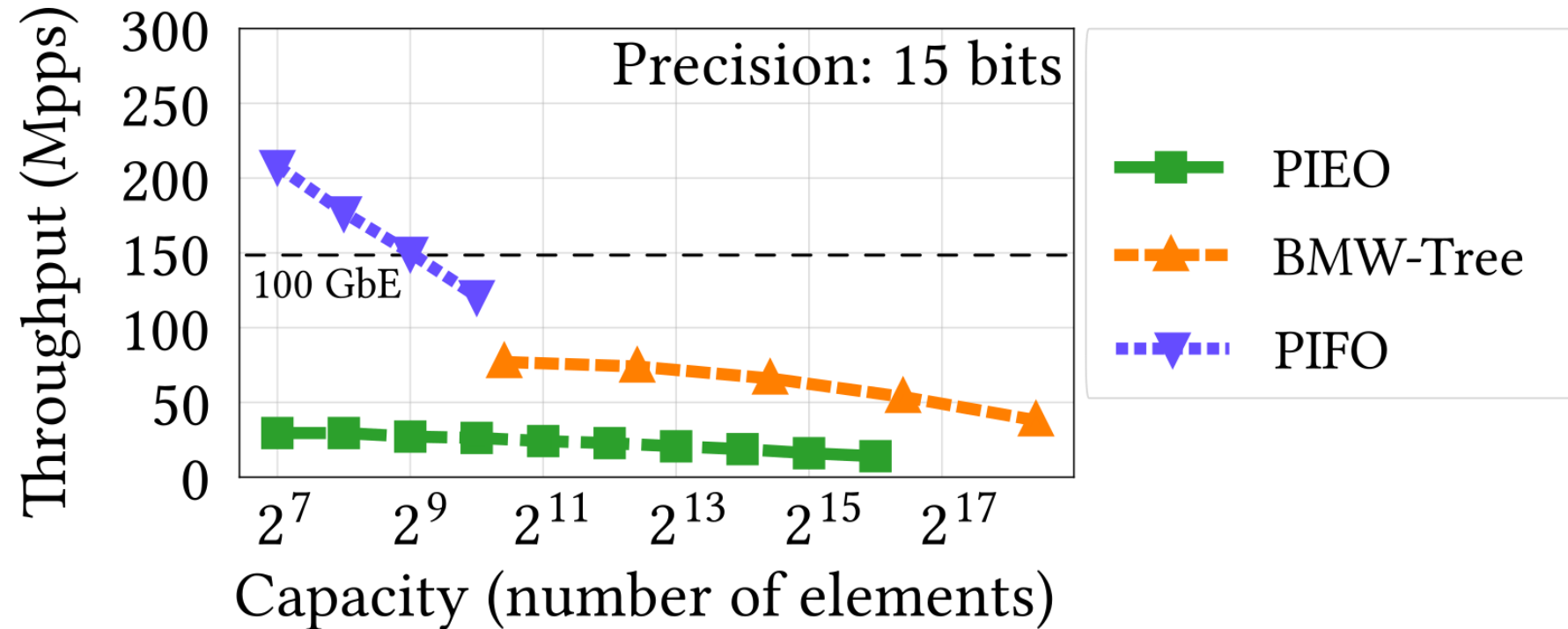BBQ's **unique priority index structure** allows logical partitioning with no performance overhead

# This Talk

- Minimum requirements for scheduling in switches and SmartNICs
- State-of-the-art priority queue designs are infeasible
- How do we get there?
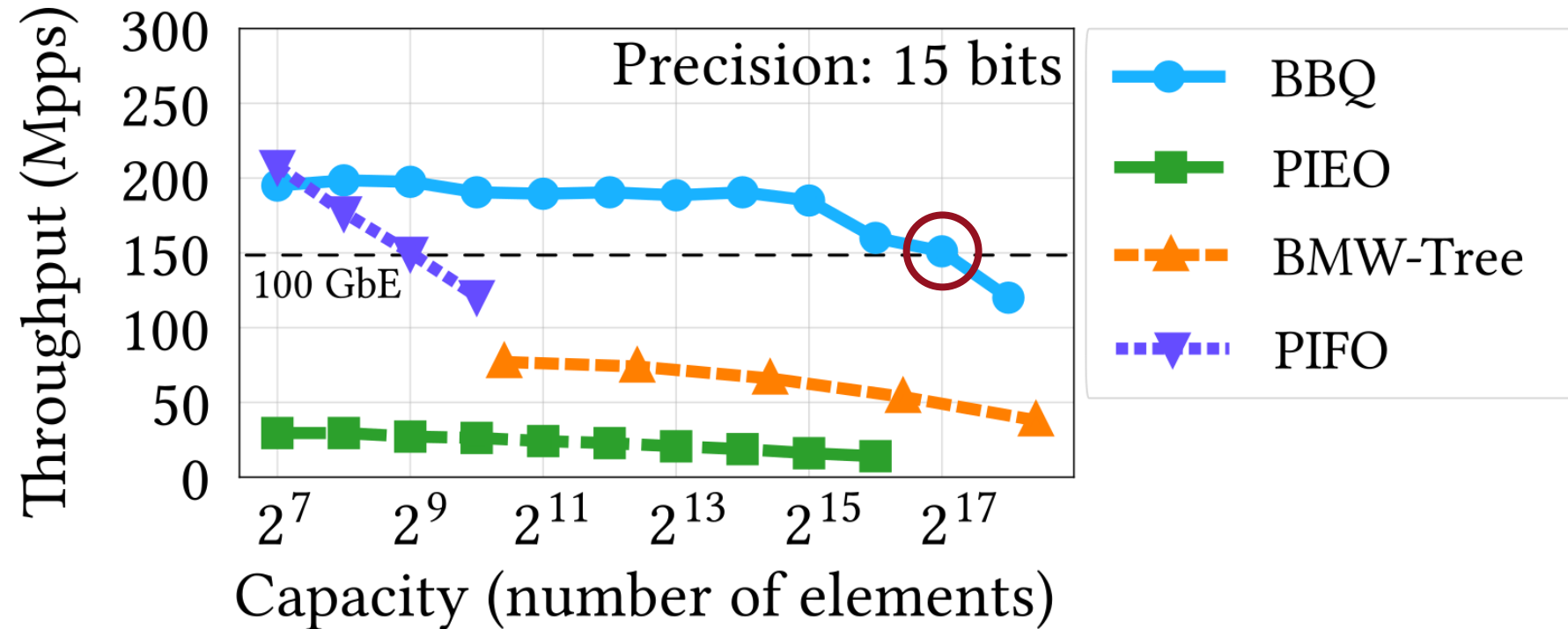- **Evaluation**

# **BBQ Evaluation (FPGA)**: Performance

# BBQ Evaluation (FPGA): Performance

# BBQ Evaluation (FPGA): Performance



On a Stratix 10 FPGA, BBQ sustains **100 Gbps line rate** (148.8 Mpps) with 100K+ elements and 32K priorities, 3X the packet rate of state-of-the-art designs.

# Conclusion

Existing hardware priority queues do not meet the stringent requirements imposed by modern schedulers. We design BBQ, an IPQ that – for the first time – makes it *feasible* to implement priority packet scheduling on line rate switches and SmartNICs.

BBQ supports 100K+ entries and 32K priorities at 100 Gbps line-rate (148.8 Mpps) on an FPGA, and 1 Tbps (1.5 Bpps) on an ASIC.
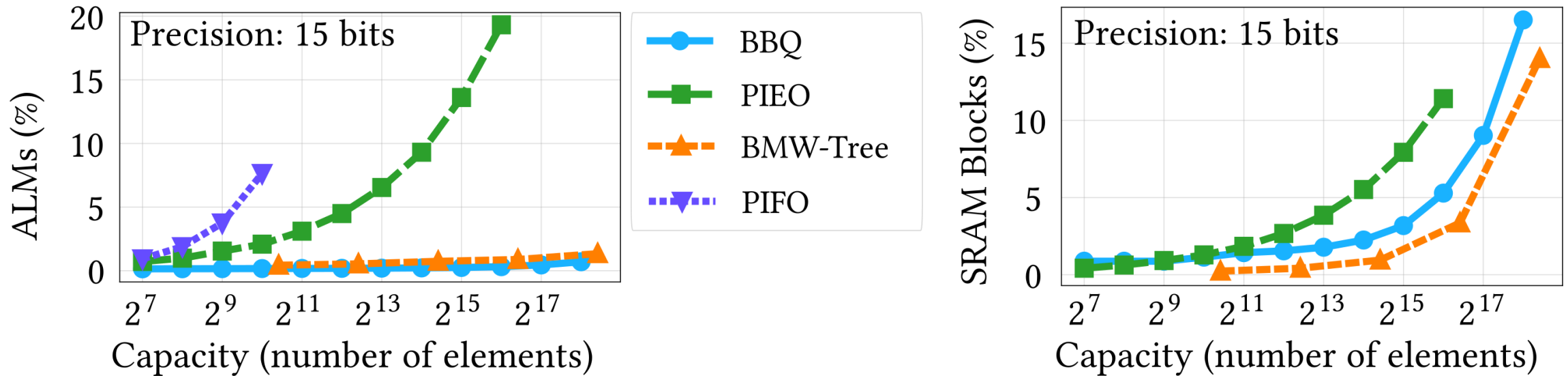
Open-source code:
https://github.com/cmu-snap/BBQ

# BBQ Evaluation (FPGA): Resources



BBQ requires very few ALMs. Its SRAM usage is between PIEO and BMW-Tree (but requires fewer copies to meet the same performance target)

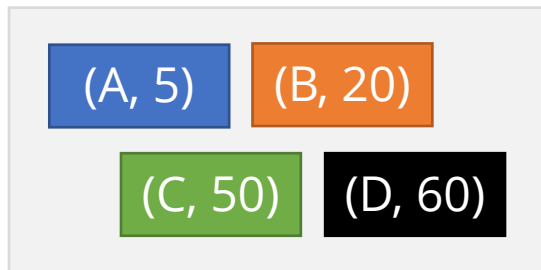# How does BBQ fit in the context of approximate priority queue designs?

1. BBQ is *complementary* to approaches that assume a small set of priority queues as a hardware primitive (SP-PIFO, PCQ, Sifter)
   - Accuracy improves with more strict-priority queues
   - BBQ's priority index structure (bitmap tree) is an efficient priority decoder, which is how we can scale to larger priority spans

2. BBQ's design shows that is possible to scale to a large number of queue elements without sacrificing accuracy

# Accuracy

Want the execution output of BBQ to be identical to an "ideal" priority queue... *unfortunately, pipelining breaks this property!*

# Accuracy

Want the execution output of BBQ to be identical to an "ideal" priority queue... *unfortunately, pipelining breaks this property!*



BBQ

Pipeline

# Accuracy

Want the execution output of BBQ to be identical to an "ideal" priority queue... *unfortunately, pipelining breaks this property!*

# Accuracy

Want the execution output of BBQ to be identical to an "ideal"
priority queue... *unfortunately, pipelining breaks this property!*

# Accuracy

Want the execution output of BBQ to be identical to an "ideal" priority queue... *unfortunately, pipelining breaks this property!*
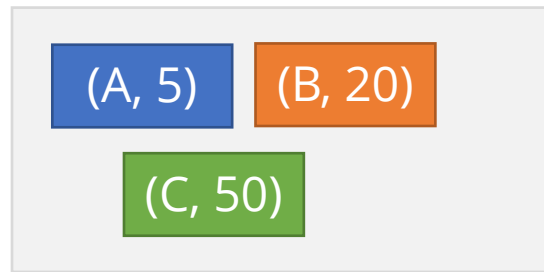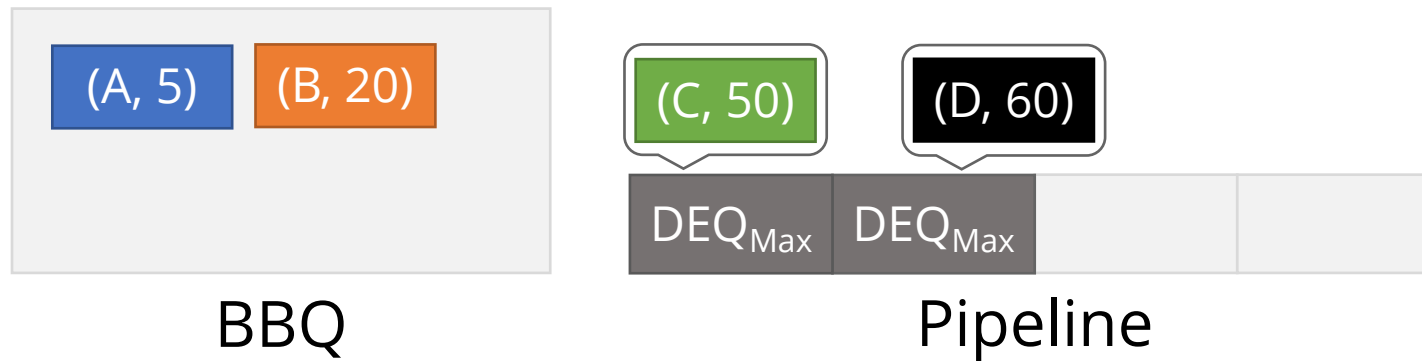
BBQ:

(A, 5)   (B, 20)

**BBQ**

Pipeline:

(C, 50)   (D, 60)

$DEQ_{Max}$   $DEQ_{Max}$

**Pipeline**

The highest-priority element in the system is not always in the BBQ, creating potential scheduling inaccuracies.

# Accuracy

We prove that combining BBQ with a tiny PIFO recovers accuracy. The composite design has all the nice properties of BBQ, but without the pipeline latency.

THEOREM 1 (**PRIORITY SET INVARIANT FOR BBQ$_\odot$**). *In a BBQ$_\odot$ instance composed of a BBQ with pipeline latency $p$ cycles and a PIFO of size $k > p$, the top $(k - p)$ highest-priority elements are always in the PIFO.*