# *Everything Matters* in Programmable Packet Scheduling

**Albert Gran Alcoz**    Balázs Vass

Pooria Namyar    Behnaz Arzani

Gábor Rétvári    Laurent Vanbever

NSDI '25

**ETH** *zürich*    MŰEGYETEM 1782    USC    Microsoft

Packet scheduling defines which packet
to send next and when

# Researchers have proposed dozens of scheduling algorithms

Minimize flow completion times

Prioritize packets from short flows                    SRPT, PIAS

Enforce fairness

Send one packet from each class at a time              RR, WFQ

Minimize tail latency

Prioritize packets with high slack time                FIFO+, LSTF

# How can we deploy all scheduling algorithms?

Implement each of them on hardware ✗

ASICs lack sufficient resources

# How can we deploy all scheduling algorithms?

**Implement each of them on hardware**

ASICs lack sufficient resources                                    ✘

**Invent a universal packet scheduler**

No silver bullet in packet scheduling                              ✘

# How can we deploy all scheduling algorithms?

Implement each of them on hardware

ASICs lack sufficient resources

✗

Invent a universal packet scheduler

No silver bullet in packet scheduling

✗

Design an abstraction to represent all schedulers

# How can we deploy all scheduling algorithms?

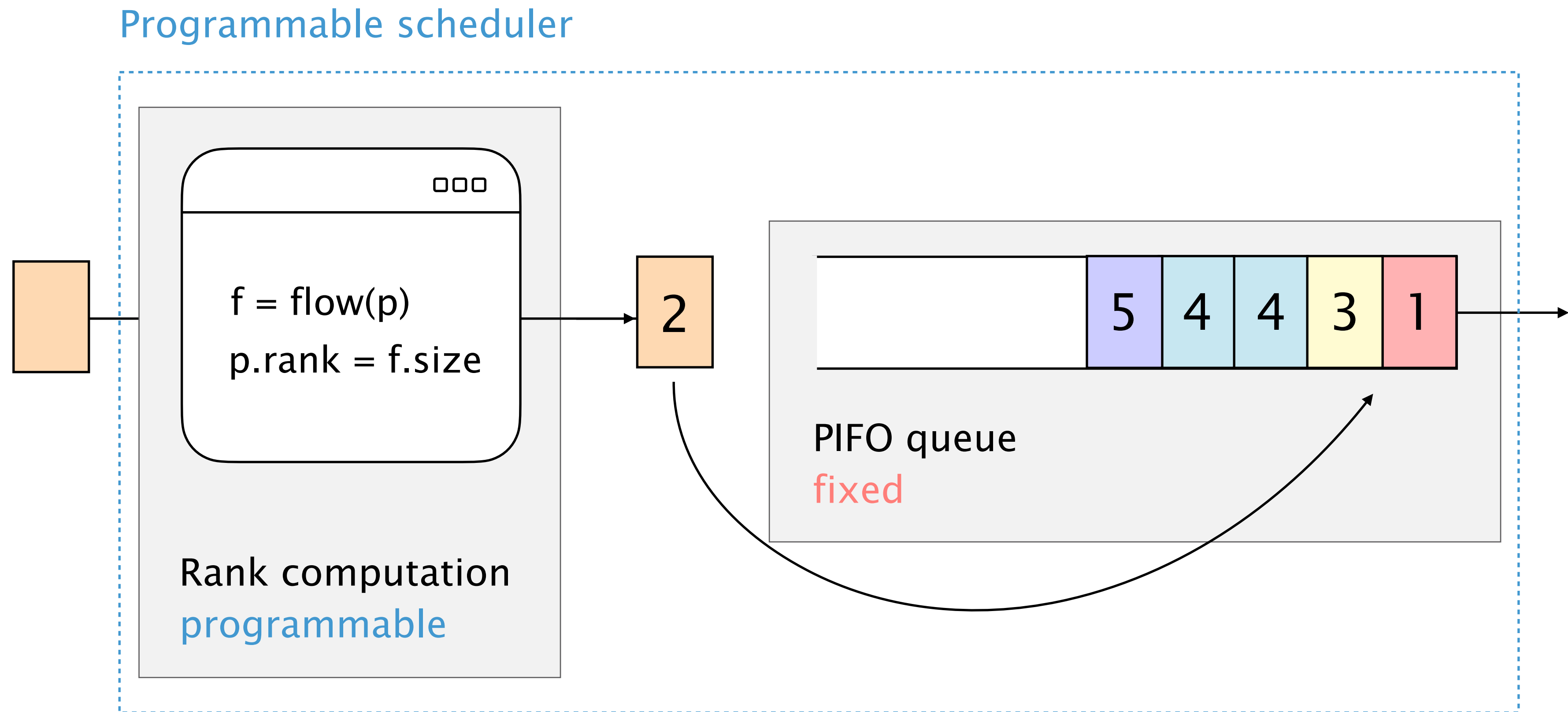Implement each of them on hardware

ASICs lack sufficient resources

✘

Invent a universal packet scheduler

No silver bullet in packet scheduling

✘

Programmable scheduling

✔

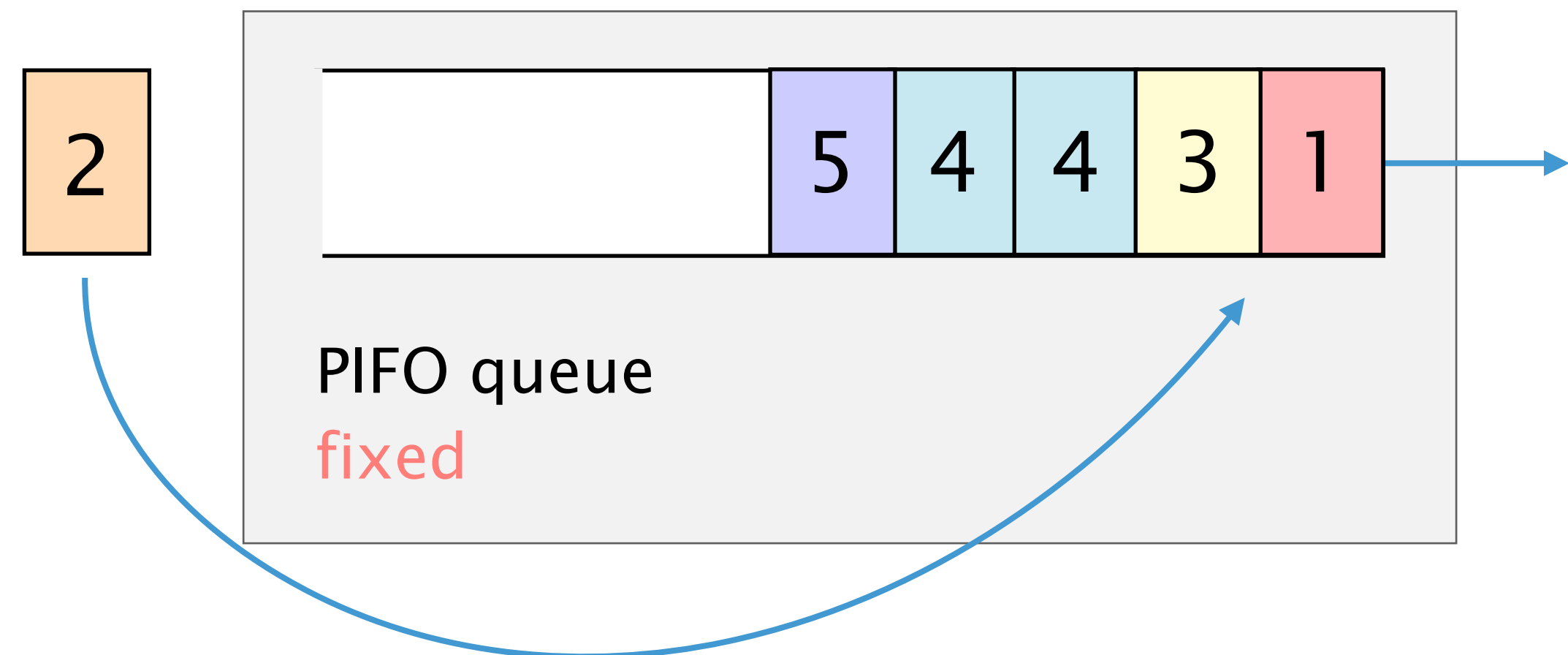# Push-In First-Out (PIFO) queues enable programmable scheduling



Programmable scheduler

f = flow(p)
p.rank = f.size

2

5 4 4 3 1

Rank computation
programmable

PIFO queue
fixed

# PIFO queues are characterized by
## two key behaviors

**Admission**

Enqueue packets with lowest ranks

**Scheduling**

Forward packets in rank order

# How to implement PIFO queues on hardware?

# How to implement PIFO queues on hardware?

New ASIC

High accuracy ✔

~200M $ ✘

Multiple years ✘

# How to implement PIFO queues on hardware?

|  | New ASIC |  | Programmable switches |  |
|---|---|---|---|---|
| High accuracy | ✔ |  |  |  |
| ~200M $ | ✘ | ~10K $ | ✔ |  |
| Multiple years | ✘ | Available today | ✔ |  |

# How to implement PIFO queues on hardware?

| New ASIC | | Programmable switches | |
|---|---|---|---|
| High accuracy | ✔ | Enough accuracy | ? |
| ~200M $ | ✘ | ~10K $ | ✔ |
| Multiple years | ✘ | Available today | ✔ |

# SP-PIFO approximates PIFO's scheduling using

## strict-priority queues

# SP-PIFO: Approximating Push-In First-Out Behaviors using Strict-Priority Queues

Albert Gran Alcoz
*ETH Zürich*

Alexander Dietmüller
*ETH Zürich*

Laurent Vanbever
*ETH Zürich*

## Abstract

Push-In First-Out (PIFO) queues are hardware primitives which enable programmable packet scheduling by providing the abstraction of a priority queue at line rate. However, implementing them at scale is not easy: just hardware designs (not implementations) exist, which support only about 1k flows.

In this paper, we introduce SP-PIFO, a programmable packet scheduler which closely approximates the behavior of PIFO queues using strict-priority queues—*at line rate, at scale, and on existing devices*. The key insight behind SP-PIFO is to dynamically adapt the mapping between packet ranks and available strict-priority queues to minimize the scheduling errors with respect to an ideal PIFO. We present a mathematical formulation of the problem and derive an adaptation technique which closely approximates the optimal queue mapping without any traffic knowledge.

We fully implement SP-PIFO in P4 and evaluate it on real

Figure 1: SP-PIFO approximates the behavior of PIFO queues

# SP-PIFO approximates PIFO's scheduling using

## strict-priority queues

One rank per queue



PIFO queue

Strict-priority queues

$\approx$

# SP-PIFO approximates PIFO's scheduling using strict-priority queues
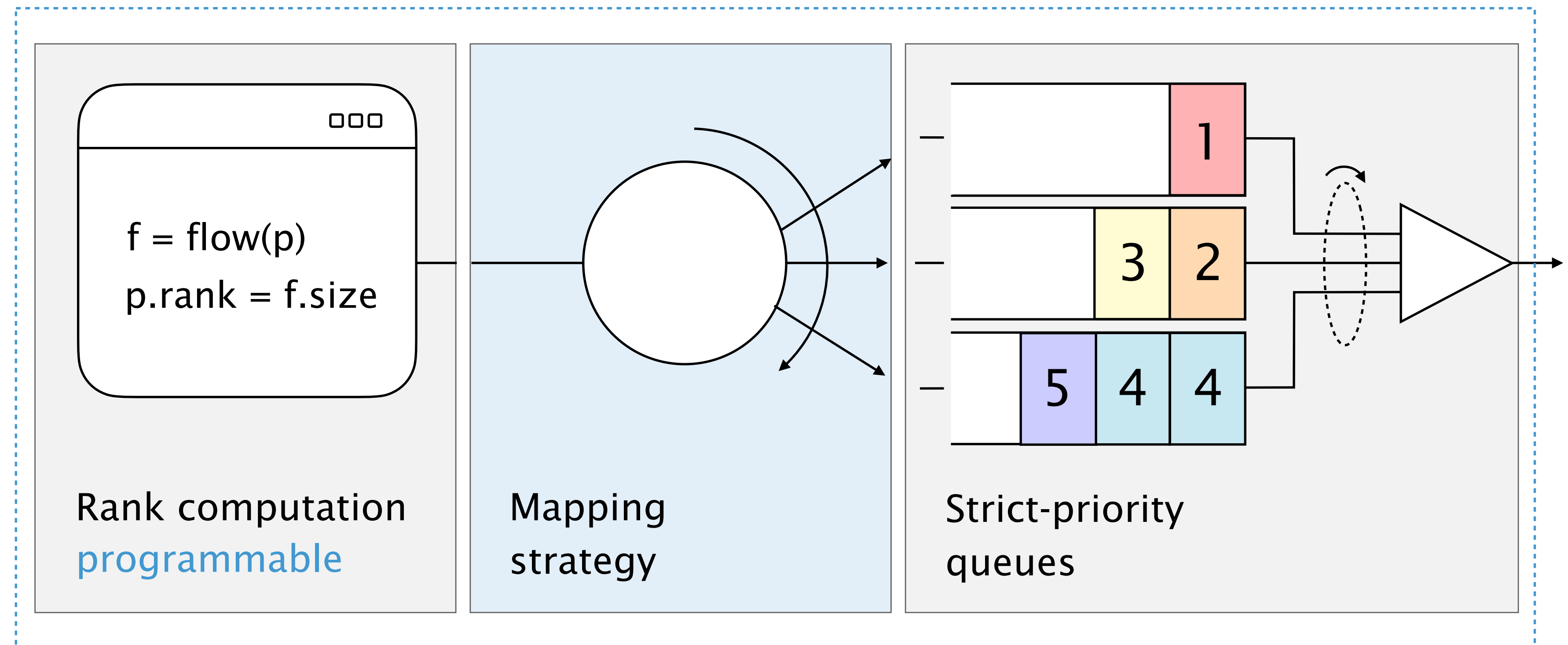


PIFO queue

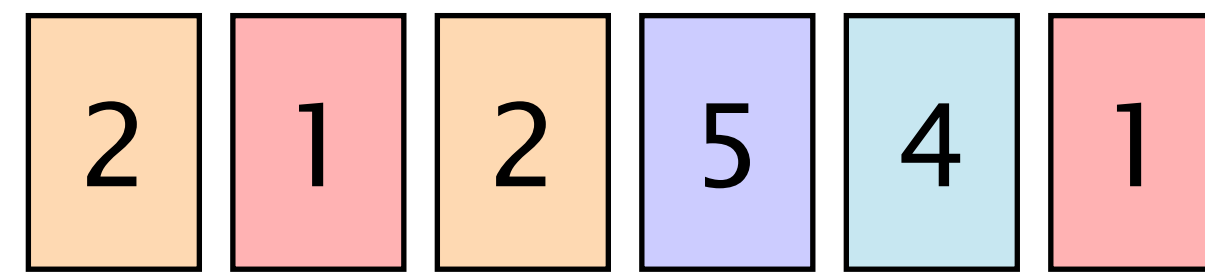In practice

Multiple ranks per queue

Inversion

Strict-priority queues

# SP-PIFO approximates PIFO's scheduling using strict-priority queues and a dynamic mapping strategy
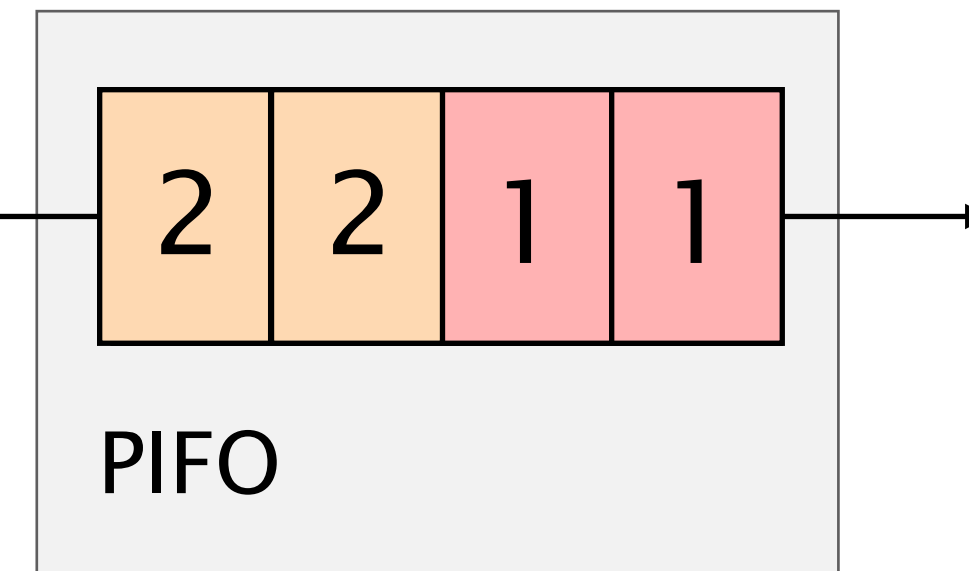


Programmable scheduler

$f = flow(p)$

$p.rank = f.size$

Rank computation
programmable

Mapping
strategy

Strict-priority
queues

# SP-PIFO approximates PIFO's scheduling, but not admission



Input sequence

2 1 2 5 4 1

Low-priority
packets dropped

PIFO

2 2 1 1

Input sequence

2 1 2 5 4 1

High-priority
packets dropped

SP-PIFO

2
5

2 1
5 4

Low-priority
packets enqueued

# AIFO approximates PIFO's admission on a single FIFO queue

SIGCOMM'21

# AIFO approximates PIFO's admission on a single FIFO queue



Programmable scheduler

f = flow(p)
p.rank = f.size

4  5

3  2  1

Rank computation
programmable

Admission
strategy

FIFO
Queue

# AIFO approximates PIFO's admission, but not scheduling

# Existing works only approximate one PIFO behavior

## Admission

Enqueue packets with lowest ranks

AIFO (SIGCOMM '21)

HCSFQ (NSDI '21)

AQ (SIGCOMM '23)

## Scheduling

Forward packets in rank order

SP-PIFO (NSDI '20)

PCQ (NSDI '20)

GearBox (NSDI '22)

QCluster (WWW '22)

Spring (INFOCOM '22)

# Existing works only approximate one PIFO behavior

**Admission**

Enqueue packets with lowest ranks

**Scheduling**

Forward packets in rank order

*"Everything matters"*

Can we approximate both PIFO behaviors

on existing programmable switches?

# Introducing…
# **PACKS**

A programmable scheduler
approximating both PIFO behaviors

# PACKS combines an admission- and a queue mapping-strategy

# PACKS combines an admission- and a queue mapping-strategy



Input sequence

2  1  2  5  4  1

$r < ?$

?

?

Admission Control

drop like PIFO

Queue Mapping

sort like PIFO

# PACKS combines an admission- and a queue mapping-strategy

Input sequence

| 2 | 1 | 2 | 5 | 4 | 1 |

r < ?

? ?

Admission Control

drop like PIFO

# PACKS combines an admission- and a queue mapping-strategy

Input sequence

| 2 | 1 | 2 | 5 | 4 | 1 |

6 packets

r < ?

?

?

Buffer availability

4 packets

# PACKS combines an admission- and a queue mapping-strategy

Input sequence

| 2 | 1 | 2 | 5 | 4 | 1 |

r < ?

? 

?

Buffer availability

B = 4 packets

Rank distribution (W)

| 1 | 2 |
| 1 | 2 |

| 4 | 5 |

B          Dropped

# PACKS combines an admission- and a queue mapping-strategy

Input sequence

2 1 2 5 4 1

r < 3

? ?

Buffer availability

B = 4 packets

$r_{drop} = 3$

Rank distribution (W)

1 2
1 2

4 5

B    Dropped

# PACKS combines an admission- and a queue mapping-strategy



Input sequence

| 2 | 1 | 2 | 5 | 4 | 1 |

r < 3

? ?

Buffer availability

B = 4 packets

Rank distribution (W)

| 1 | 2 |
| 1 | 2 |

4  5

B    Dropped

$r_{drop} = 3$

$\text{maximize} \quad r_{drop}$

$$s.t., W.quantile(r_{drop} - 1) \leq \frac{B}{|W|}$$

# PACKS combines an admission- and a queue mapping-strategy

Input sequence

2 1 2 5 4 1

r < 3

?

?

Queue Mapping

sort like PIFO

# PACKS combines an admission- and a queue mapping-strategy

Input sequence

| 2 | 1 | 2 | 5 | 4 | 1 |

r < 3

? ?

Rank distribution (W)

| 1 | 2 |
| 1 | 2 |

4 5

B          Dropped

Queue Availability

$B_1$ = 2 packets

$B_2$ = 2 packets

# PACKS combines an admission- and a queue mapping-strategy



Input sequence

2 | 1 | 2 | 5 | 4 | 1

r < 3

?

?

Rank distribution (W)

1 | 2
1 | 2
4 | 5

B1 | B2 | Dropped

Queue Availability

B1 = 2 packets

B2 = 2 packets

# PACKS combines an admission- and a queue mapping-strategy

Input sequence

| 2 | 1 | 2 | 5 | 4 | 1 |

$r < 3$

| 1 |
| 2 |

| 1 | 1 |
| 2 | 2 |

$q_1 = 1, \quad q_2 = 2$

Rank distribution (W)

| 1 | 2 |
| 1 | 2 |

| 4 | 5 |

B1   B2   Dropped

maximize $q_i$

$s.t., W.quantile(q_i) \leq \dfrac{\sum_{j=1}^{i} B_j}{|W|}$

# How to translate the algorithm to the online case?

How to monitor the rank distribution?

How to adapt to buffer dynamism?

How to account for workload shifts?

# How to translate the algorithm to the online case?

How to monitor the rank distribution?

Use a sliding window of latest ranks

How to adapt to buffer dynamism?

How to account for workload shifts?

# How to translate the algorithm to the online case?

How to monitor the rank distribution?

Use a sliding window of latest ranks

How to adapt to buffer dynamism?

Measure per-packet queue occupancy
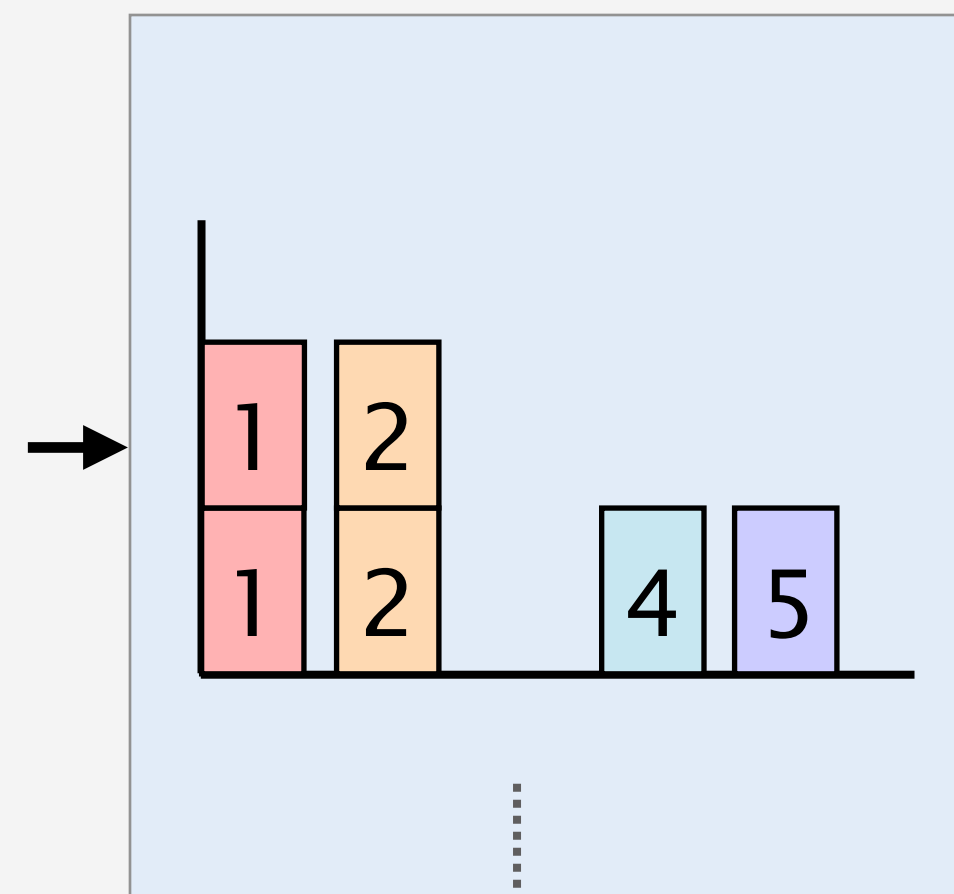
How to account for workload shifts?

$$W.quantile(r) \leq \frac{\sum_{j=1}^{i}(B_j - bj)}{B}$$

Per-packet

queue occupancy

# How to translate the algorithm to the online case?

**How to monitor the rank distribution?**

Use a sliding window of latest ranks

**How to adapt to buffer dynamism?**

Measure per-packet queue occupancy

**How to account for workload shifts?**

Allow a certain amount of bursts

$$W . quantile(r) \ \leq \ \alpha \cdot \frac{\sum_{j=1}^{i} (B_j - bj)}{B}$$
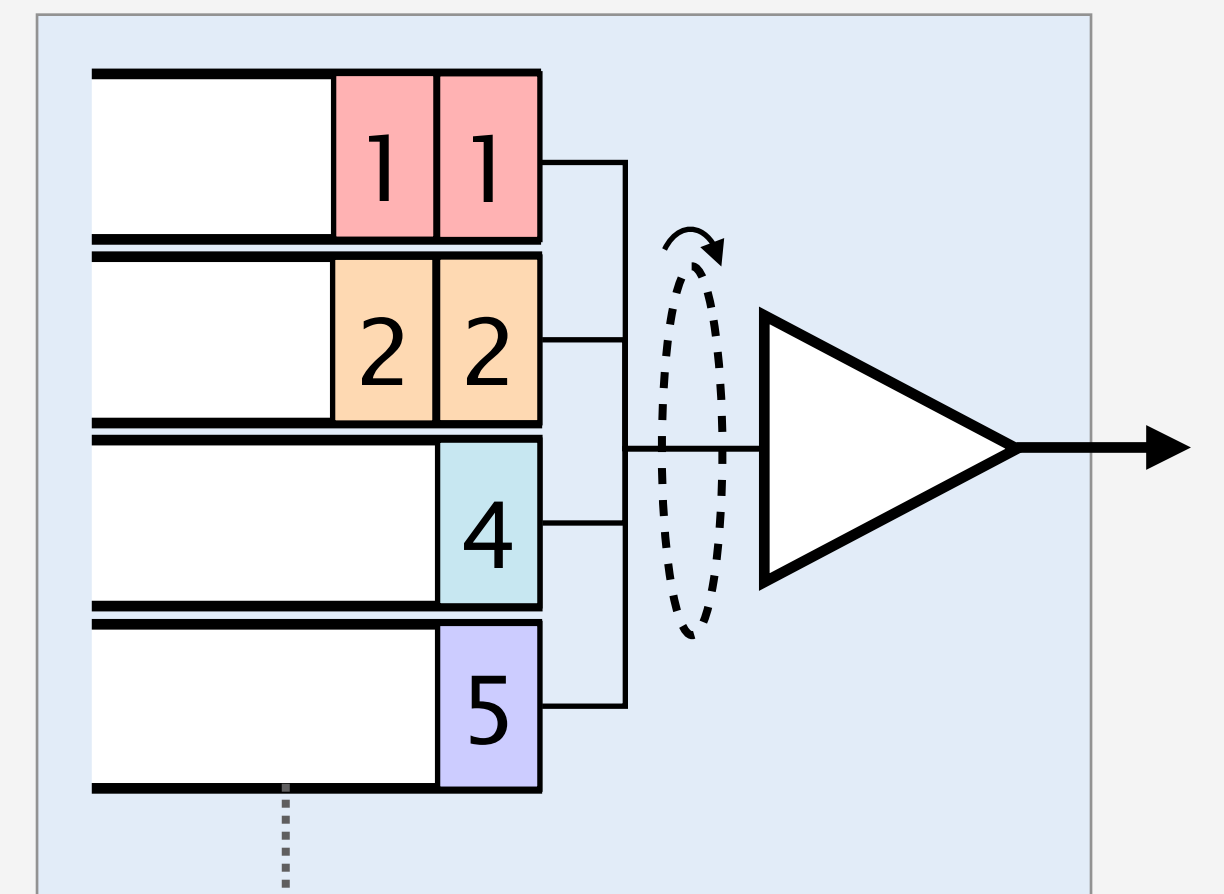
Burst

allowance

# PACKS

Sliding window
tracking

Buffer occupancy
 monitoring



$$W \cdot quantile(r) \ \leq \ \alpha \cdot \frac{\sum_{j=1}^{i} (B_j - bj)}{B}$$

Quantile

Queue occupancy

# PACKS

**Sliding window tracking**

**Admission and queue mapping**

**Buffer occupancy monitoring**



Scan top-down

Enqueue if:

$$W.quantile(r) \;\leq\; \alpha \cdot \frac{\sum_{j=1}^{i}(B_j - bj)}{B}$$

Quantile

Queue occupancy

# We evaluated PACKS on hardware and simulations

Packet-level simulation (NetBench)

Performance in approximating PIFO

Sensitivity to configuration parameters
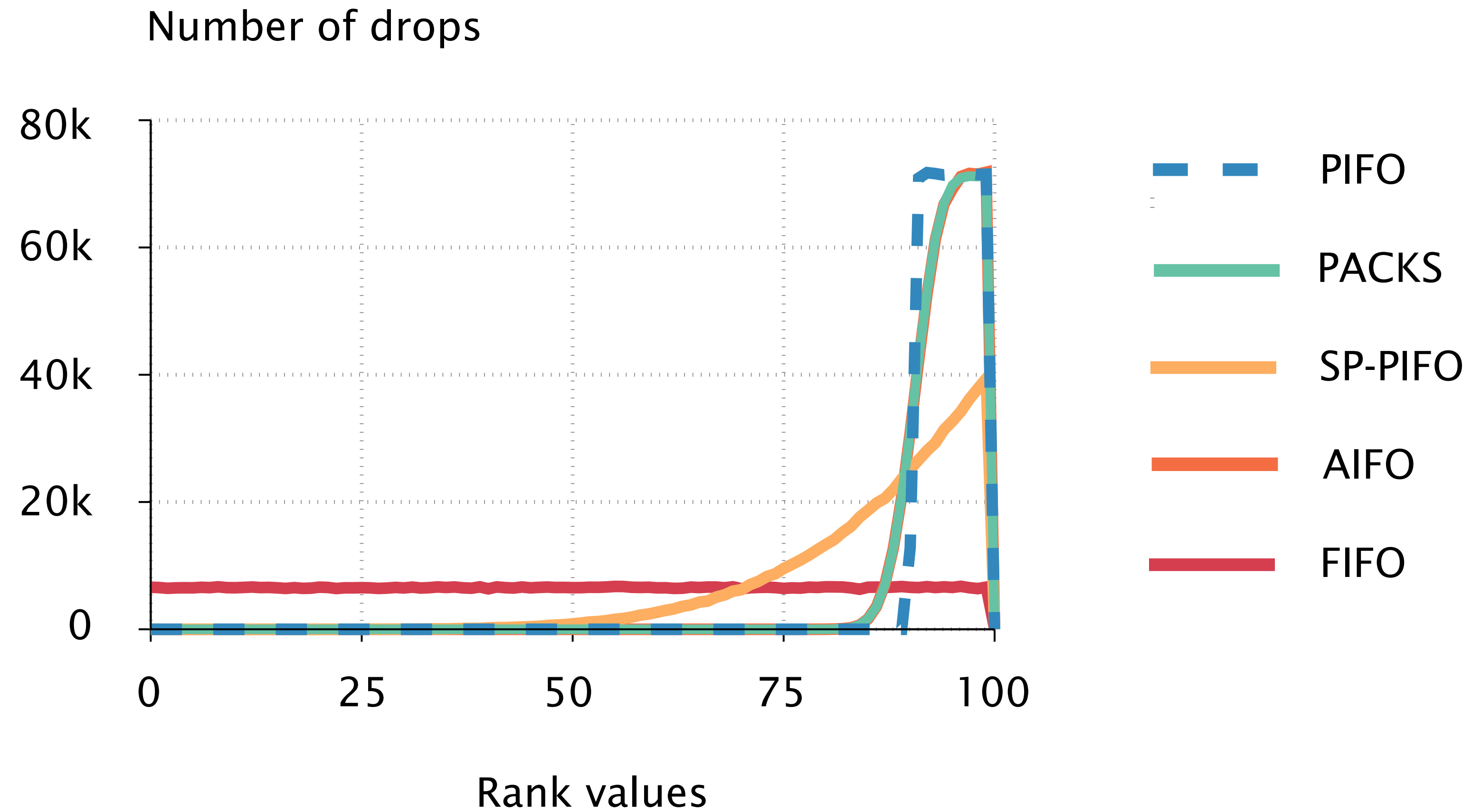
Practicality under pFabric and FQ scenarios

Hardware evaluation (Intel Tofino2)

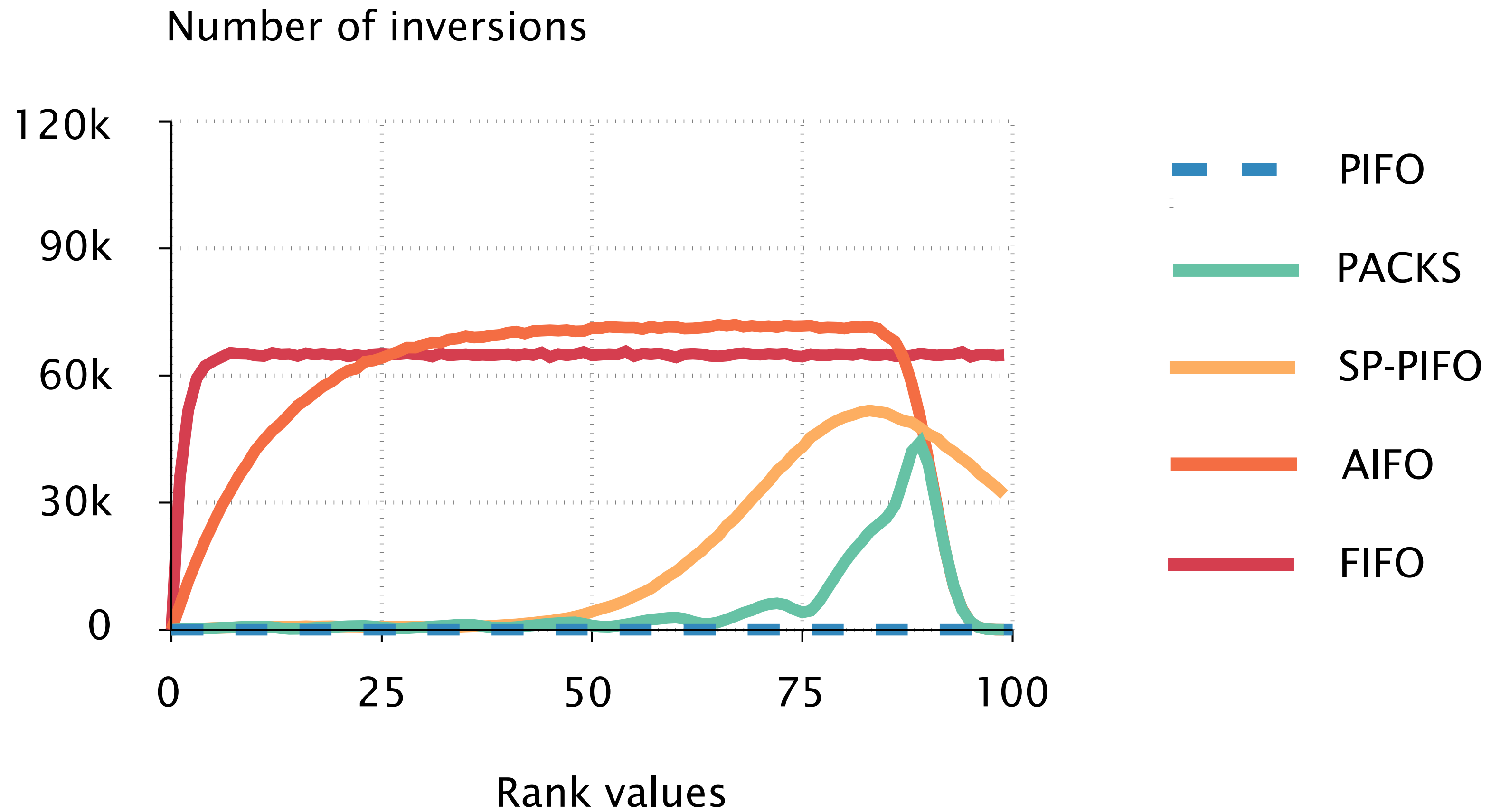Bandwidth allocation across priorities

Heuristic analysis (MetaOpt)
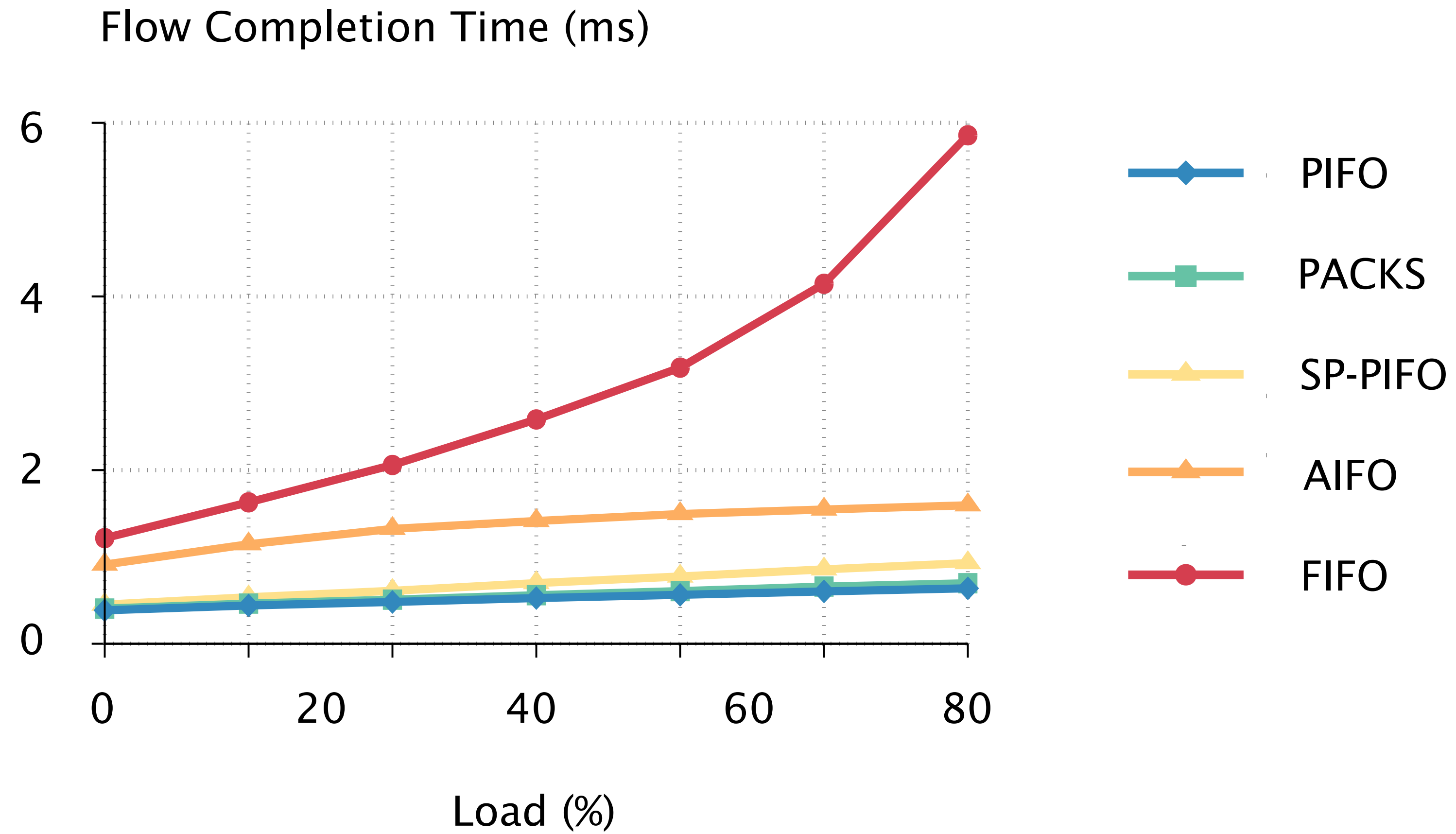
Adversarial workload analysis

# PACKS reduces packet drops by up to 60% compared to SP-PIFO

# PACKS reduces inversions by up to 7x and 15x compared to SP-PIFO and AIFO



Number of inversions

Rank values

PIFO
PACKS
SP-PIFO
AIFO
FIFO

# PACKS reduces mean FCTs by up to 33% and 2.6x compared to SP-PIFO and AIFO



Flow Completion Time (ms)

Load (%)

PIFO
PACKS
SP-PIFO
AIFO
FIFO

# *Everything Matters* in Programmable Packet Scheduling

**PACKS approximates PIFO's admission and scheduling behaviors**

at line rate, on existing programmable switches

**PACKS adapts to traffic workloads in real time**

using a sliding window and queue-aware policies

**PACKS outperforms existing approaches**

reducing drops, inversions, and flow completion times

github.com/nsg-ethz/packs