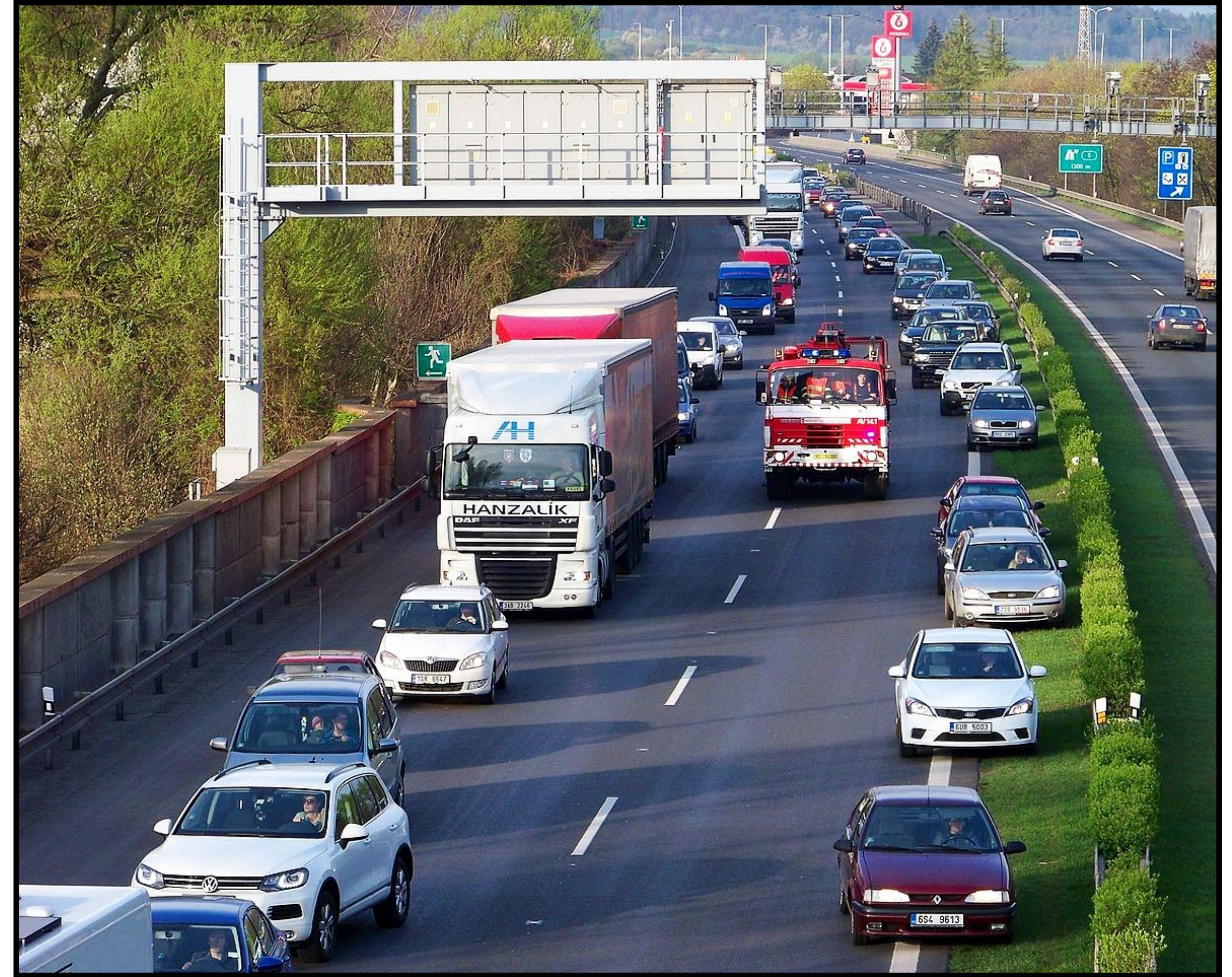**In-Network Congestion Management
for Security and Performance**
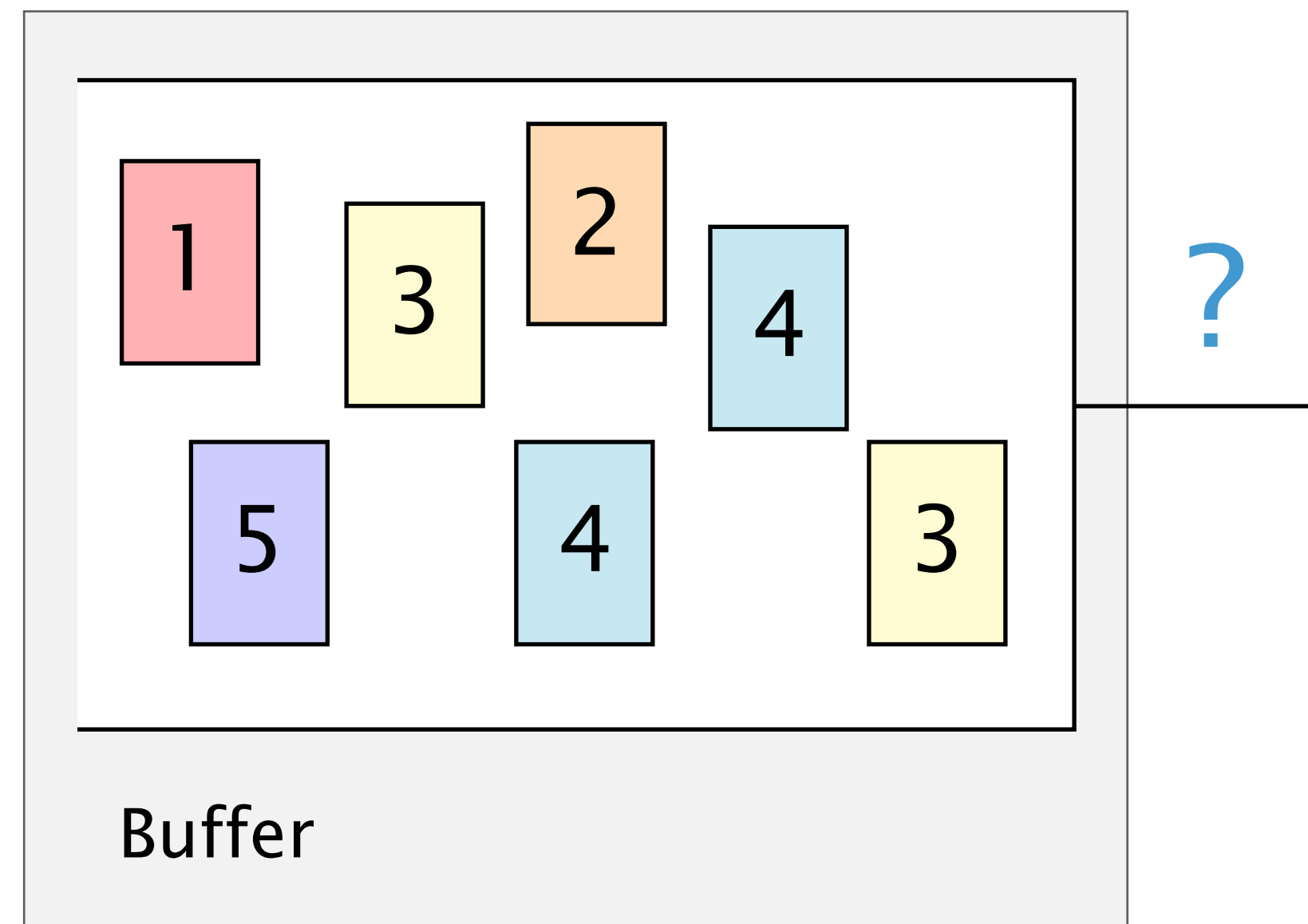
Albert Gran Alcoz
September 18 2024

# Packet scheduling

# Packet scheduling defines what packet should we send next and when



Buffer

# 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

# A universal scheduling algorithm does not exist

## Universal Packet Scheduling

Radhika Mittal[†]        Rachit Agarwal[†]        Sylvia Ratnasamy[†]        Scott Shenker[†‡]

[†]UC Berkeley        [‡]ICSI

### Abstract

In this paper we address a seemingly simple question: *Is there a universal packet scheduling algorithm?* More precisely, we analyze (both theoretically and empirically) whether there is a single packet scheduling algorithm that, at a network-wide level, can perfectly match the results of *any* given scheduling algorithm. We find that in general the answer is "no". However, we show theoretically that the classical Least Slack Time First (LSTF) scheduling algorithm comes closest to being universal and demonstrate empirically that LSTF can closely replay a wide range of scheduling algorithms in realistic network settings. We then evaluate whether LSTF can be used *in practice* to meet various network-wide objectives by looking at popular performance metrics (such as mean FCT, tail packet delays, and fairness); we find that LSTF performs comparable to the state-of-the-art for each of them. We also discuss how LSTF can be used in conjunction with active queue management schemes (such as CoDel) without changing the core of the network.

## 1 Introduction

There is a large and active research literature on novel packet scheduling algorithms, from simple schemes such as priority scheduling [31], to more complicated mech-

We can define a universal packet scheduling algorithm (hereafter UPS) in two ways, depending on our viewpoint on the problem. From a theoretical perspective, we call a packet scheduling algorithm *universal* if it can replay any *schedule* (the set of times at which packets arrive to and exit from the network) produced by any other scheduling algorithm. This is not of practical interest, since such schedules are not typically known in advance, but it offers a theoretically rigorous definition of universality that (as we shall see) helps illuminate its fundamental limits (i.e., which scheduling algorithms have the flexibility to serve as a UPS, and why).

From a more practical perspective, we say a packet scheduling algorithm is universal if it can achieve different desired performance objectives (such as fairness, reducing tail latency, minimizing flow completion times). In particular, we require that the UPS should match the performance of the best known scheduling algorithm for a given performance objective. [1]

The notion of universality for packet scheduling might seem esoteric, but we think it helps clarify some basic questions. If there exists no UPS then we should *expect* to design new scheduling algorithms as performance objectives evolve. Moreover, this would make a strong argument for switches being equipped with programmable

# How to deploy all scheduling algorithms?

✗ Generality

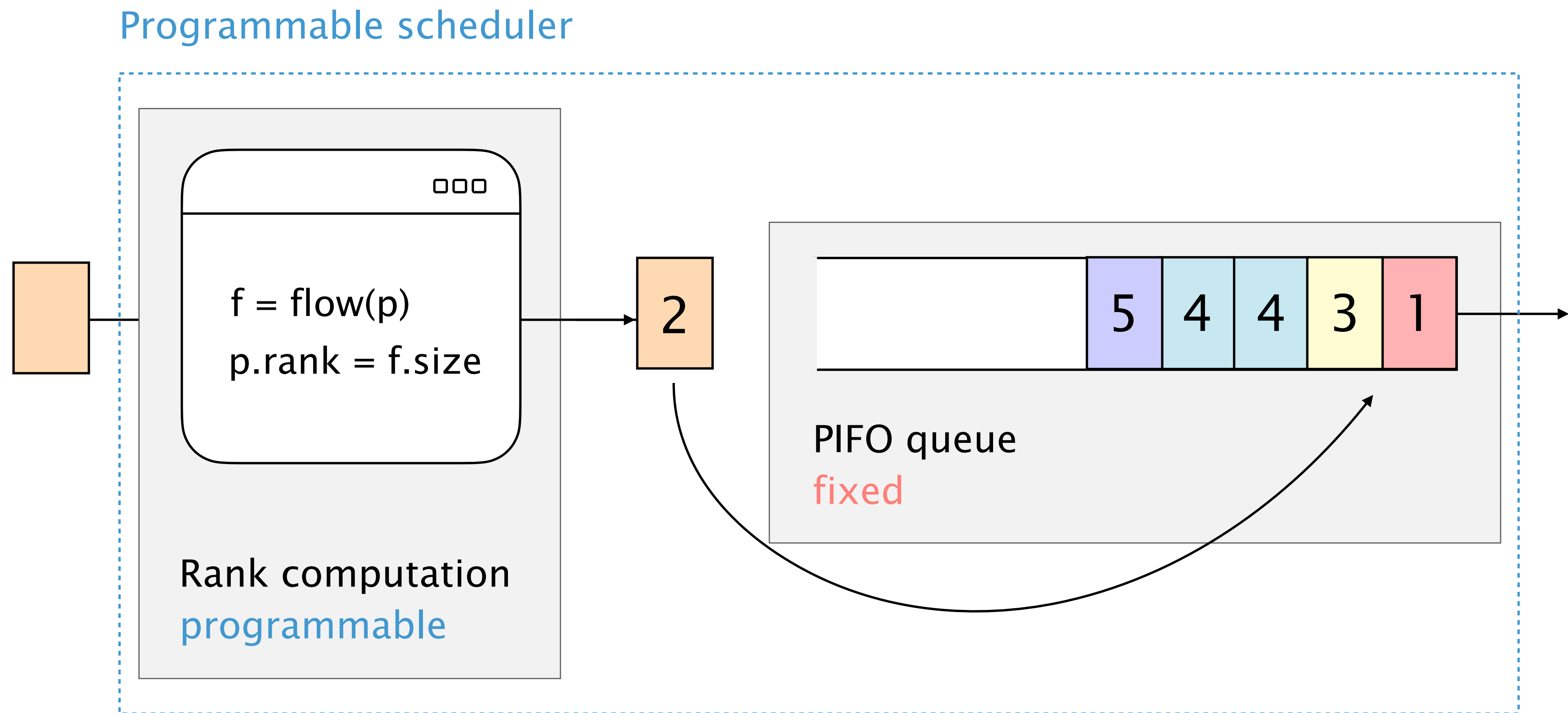Universal packet scheduler

Flexibility

Customized algorithms

# How to deploy all scheduling algorithms?

✗ **Generality**

Universal packet scheduler

Programmable
scheduling

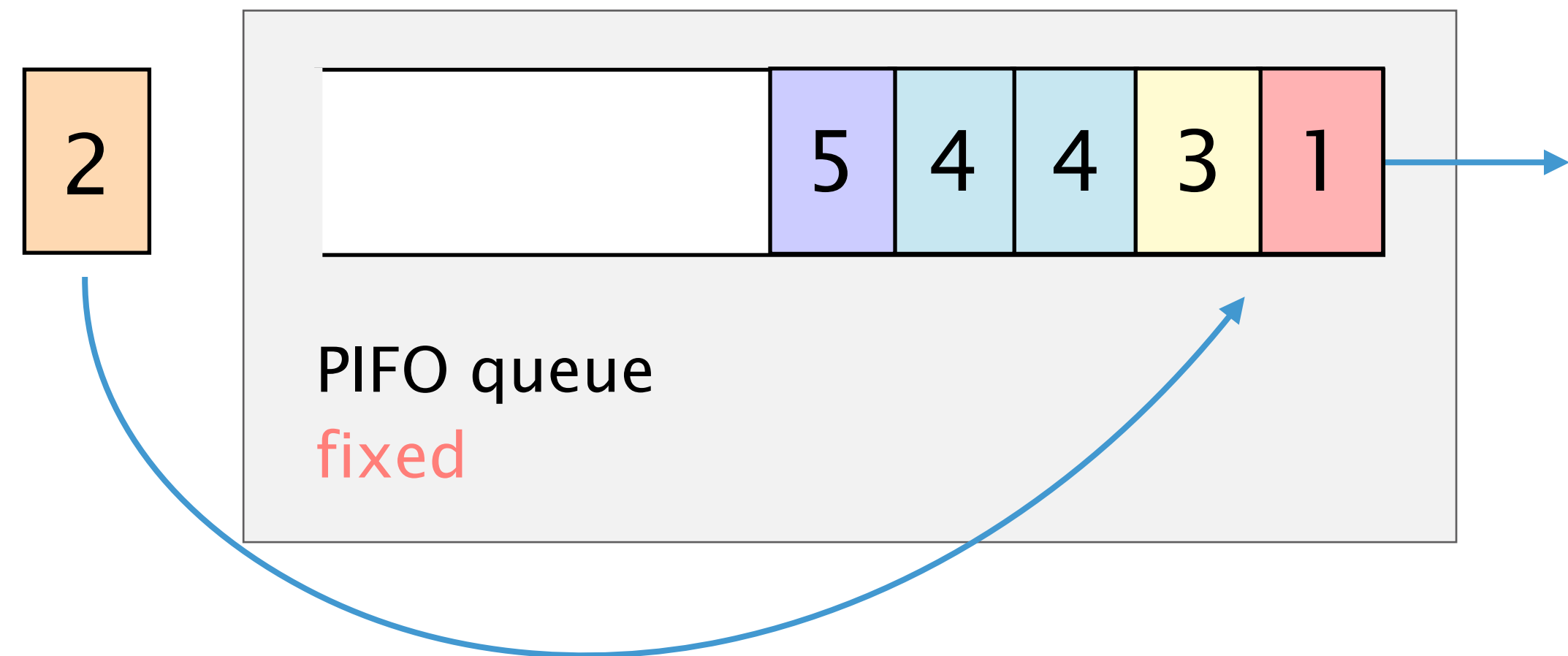# Push-In First-Out (PIFO) queues enable programmable scheduling

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

Enqueue packets with the lowest ranks

Forward packets in rank order



2

5 4 4 3 1

PIFO queue
fixed

# How to implement PIFO queues on hardware?

New ASIC

High performance  ✔

~200M $  ✘

Multiple years  ✘

# How to implement PIFO queues on hardware?

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

# How to implement PIFO queues on hardware?

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

## Objective

Enable programmable scheduling on existing devices

to improve the Internet's performance and security

How to enable programmable scheduling
on existing devices?

How to use it to improve
the Internet's security?

| SP-PIFO |
| :---: |
| [NSDI'20] |

| PACKS |
| :---: |
| [NSDI'25] |

| ACC-Turbo |
| :---: |
| [SIGCOMM'22] |

Approximating

PIFO's scheduling

Incorporating

PIFO's admission

Mitigating

DDoS attacks

How to enable programmable scheduling
on existing devices?

How to use it to improve
the Internet's security?

| SP-PIFO | PACKS | ACC-Turbo |
|---------|-------|-----------|
| [NSDI'20] | [NSDI'25] | [SIGCOMM'22] |

Approximating
PIFO's scheduling

Incorporating
PIFO's admission

Mitigating
DDoS attacks

# We can approximate PIFO queues using strict-priority queues

PIFO queue

Strict-priority queues

# We can approximate PIFO queues using strict-priority queues

Multiple ranks per queue



PIFO queue

$\approx$

Inversion

Strict-priority queues

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



Programmable scheduler

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

Rank computation
programmable

Adaptation
strategy

Strict-priority
queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues

Mapping
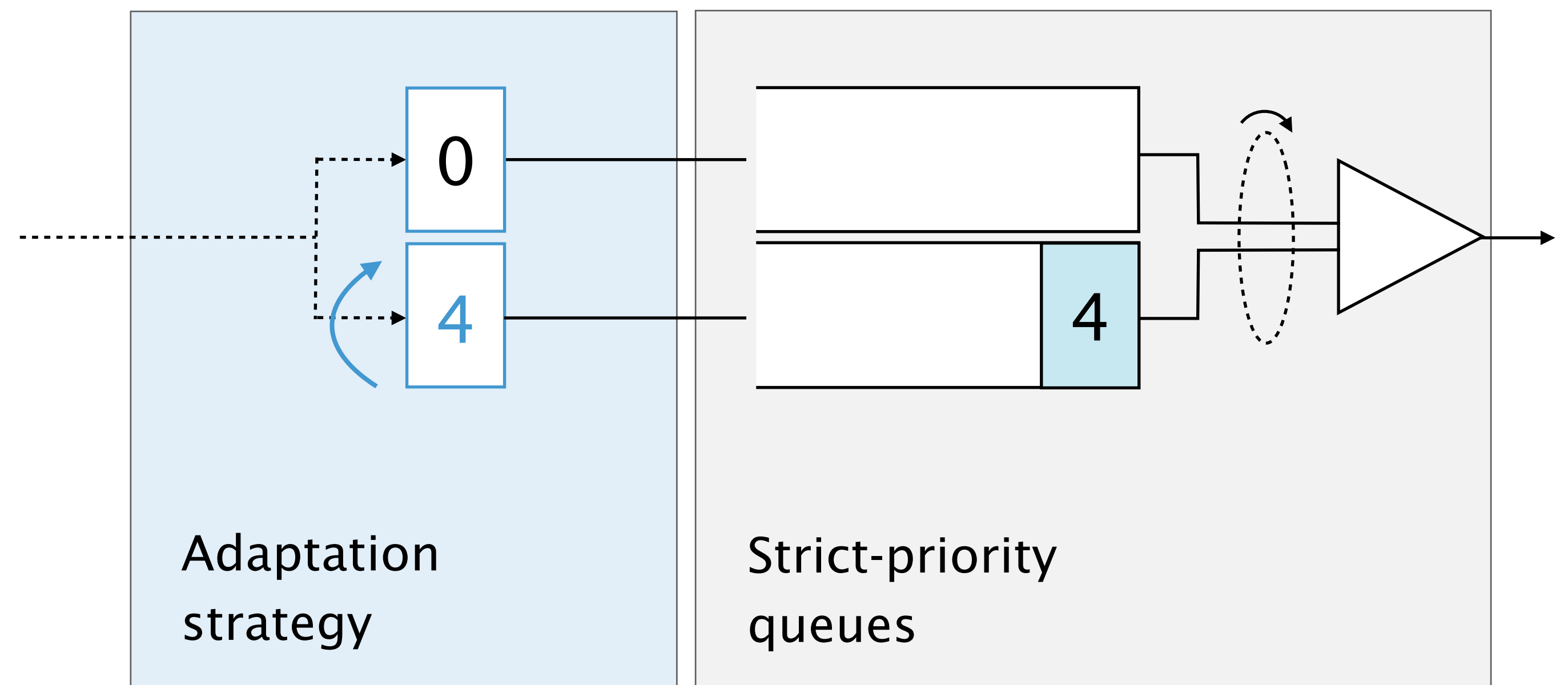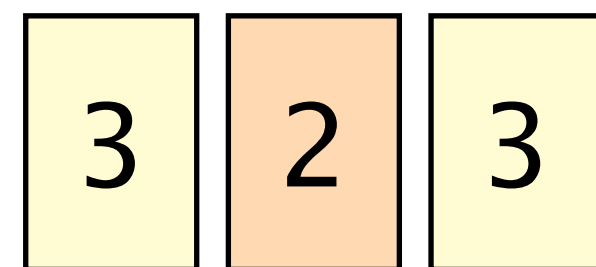
Scan bottom-up, enqueue if rank >= bound



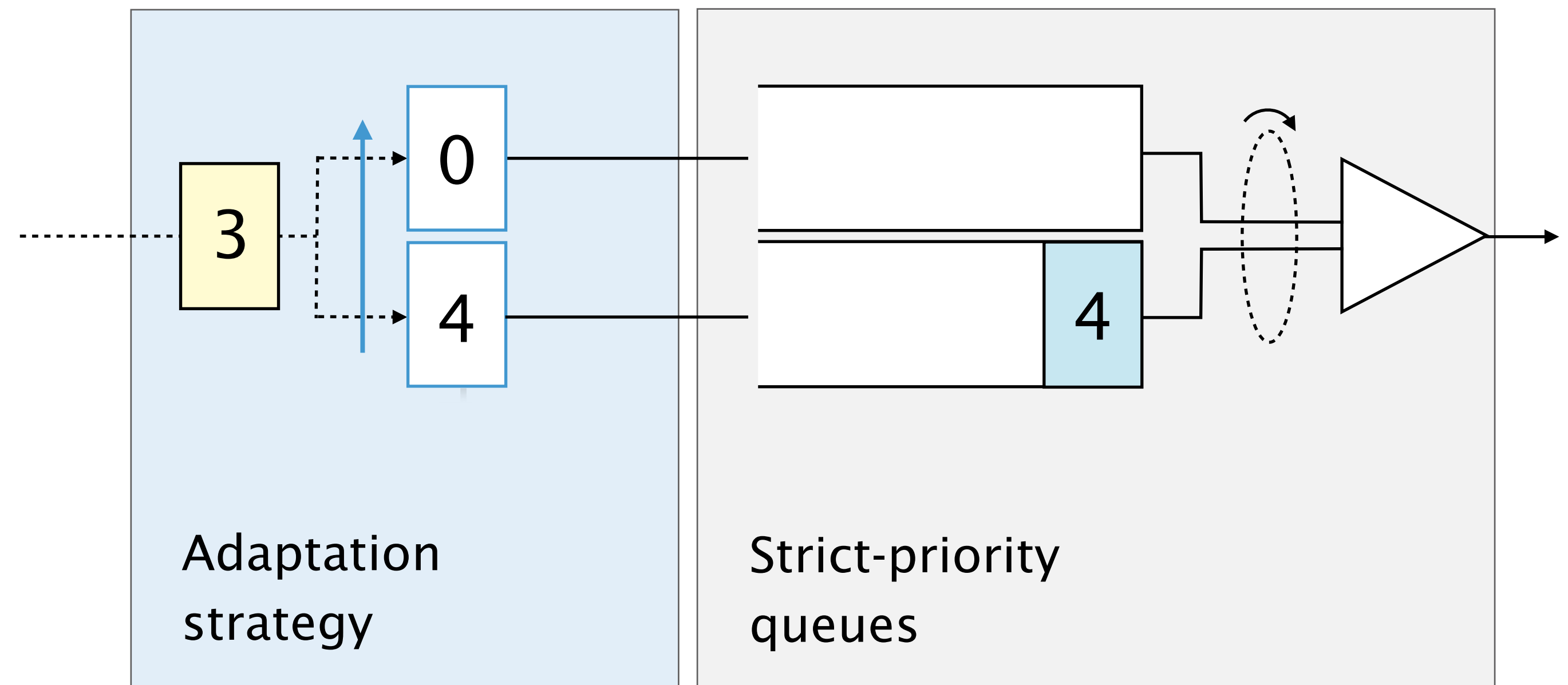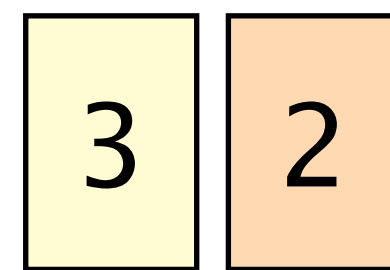Input sequence

3  2  3  4

Adaptation strategy

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues

Mapping
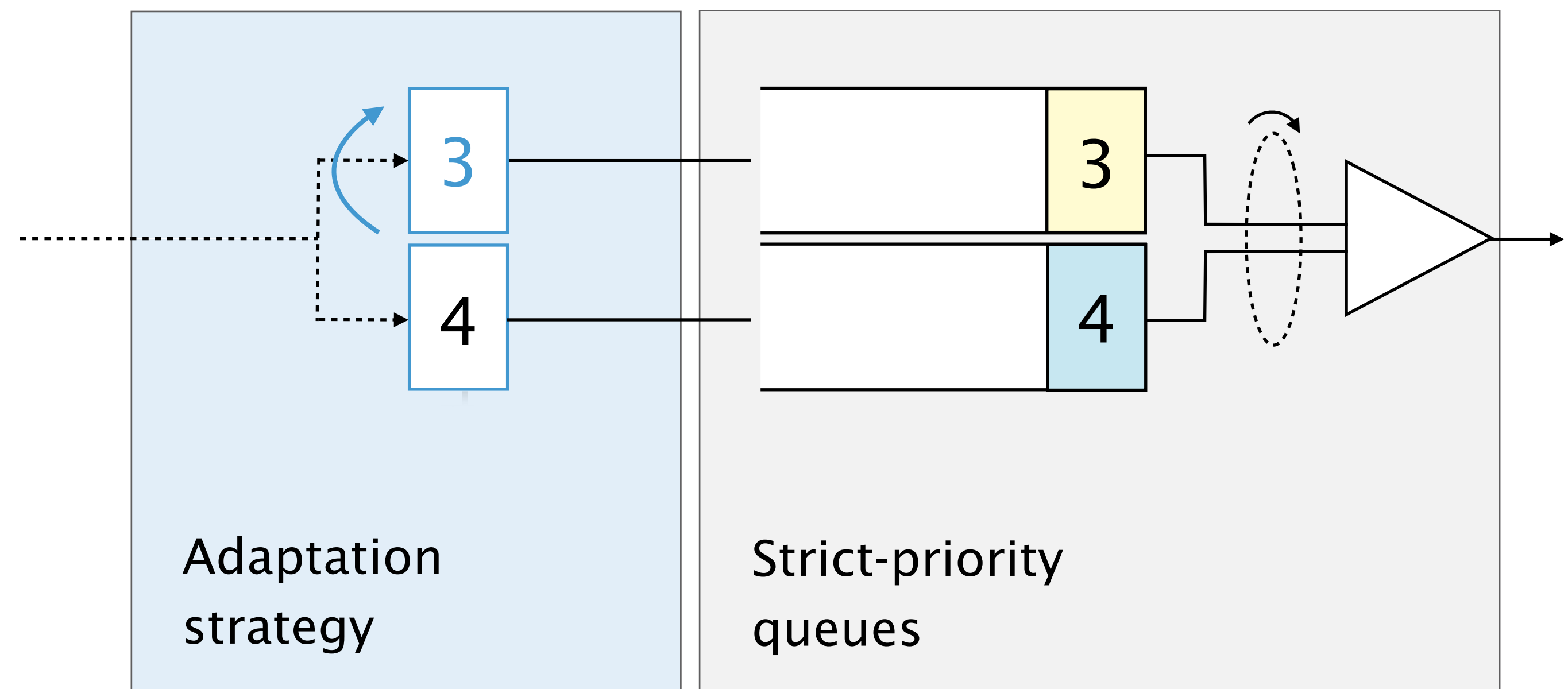
Scan bottom-up, enqueue if rank >= bound

Input sequence

3 2 3

4

0

0

Adaptation strategy

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues



Mapping

Scan bottom-up, enqueue if rank >= bound

Input sequence

3 2 3

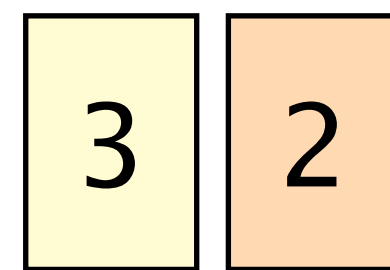Adaptation strategy

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues

Input sequence

3 2 3

Push-up adaptation

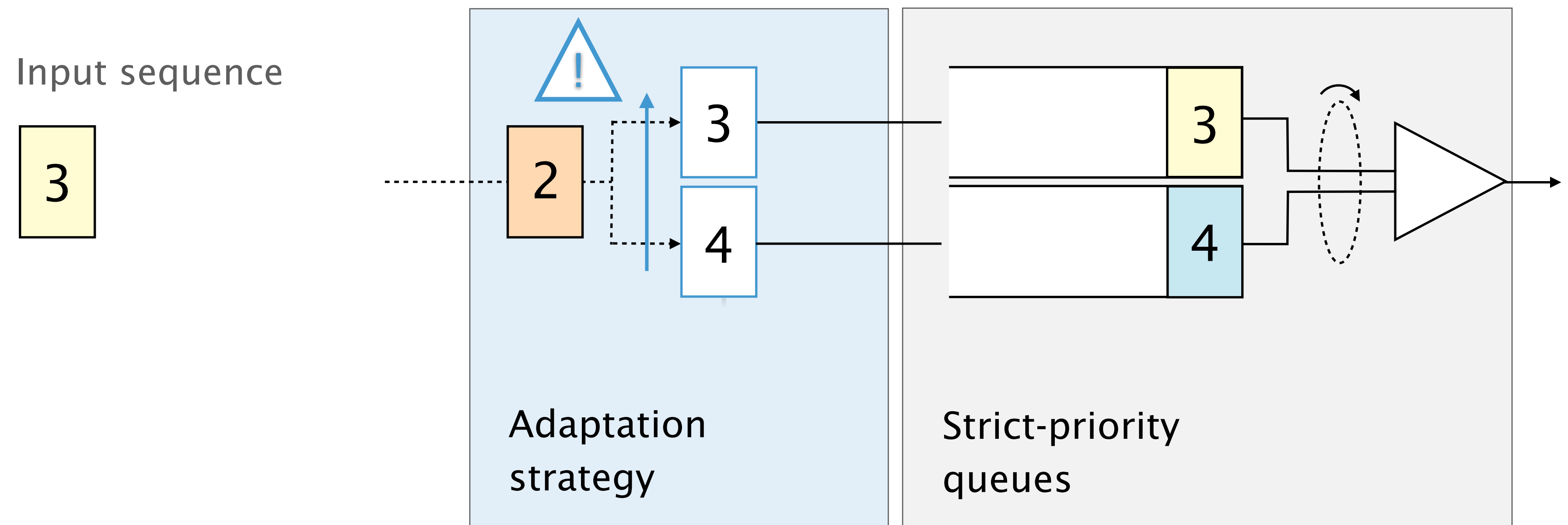Set bound to packet rank after enqueue

0

4

Adaptation strategy

4

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues

Mapping

Scan bottom-up, enqueue if rank >= bound



Input sequence

Adaptation strategy

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues



Push-up adaptation

Set bound to packet rank after enqueue

Input sequence

3 2

3

4

Adaptation strategy

3

4

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues

Push-down adaptation
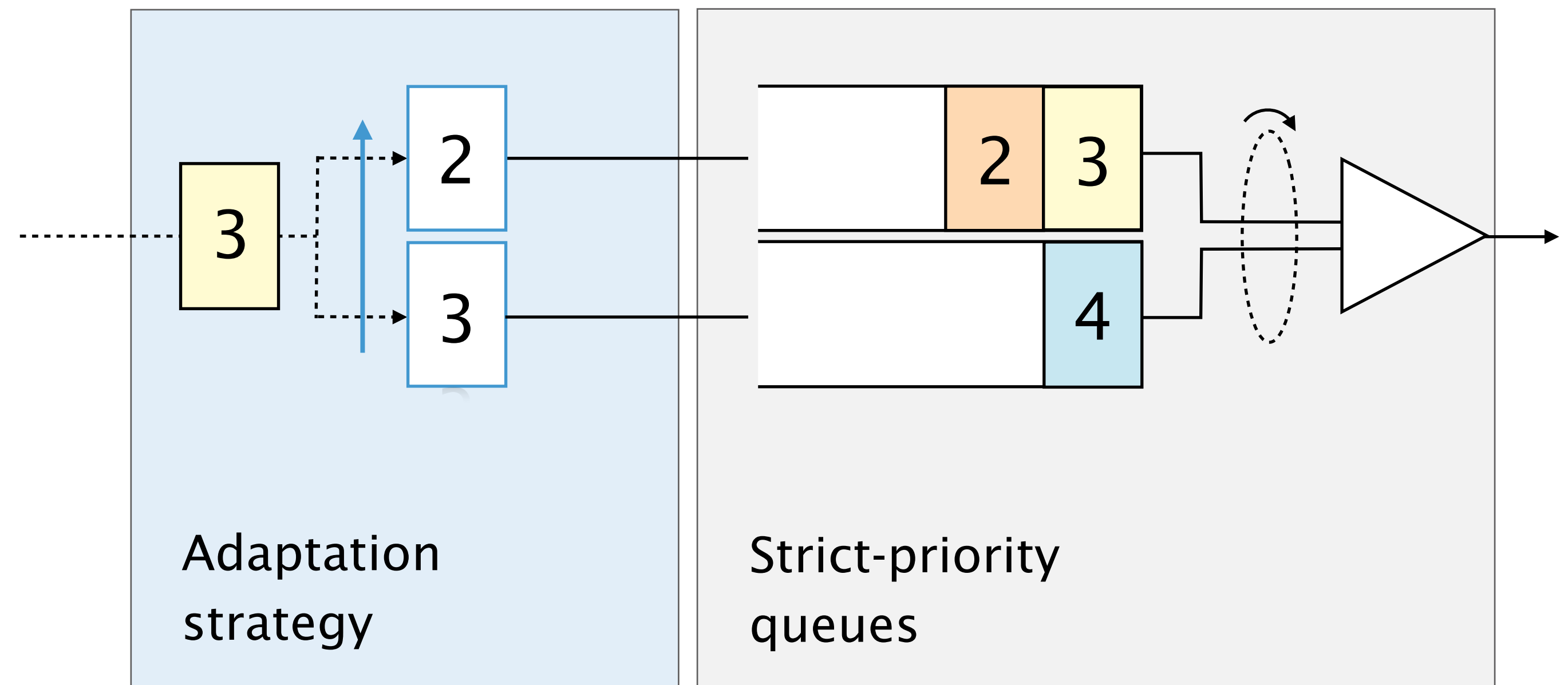
Decrease all bounds after inversion, by inversion cost

Input sequence

3

3

4

Adaptation strategy

2 3

4

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks to strict-priority queues



**Push-down adaptation**

Decrease all bounds after inversion, by inversion cost
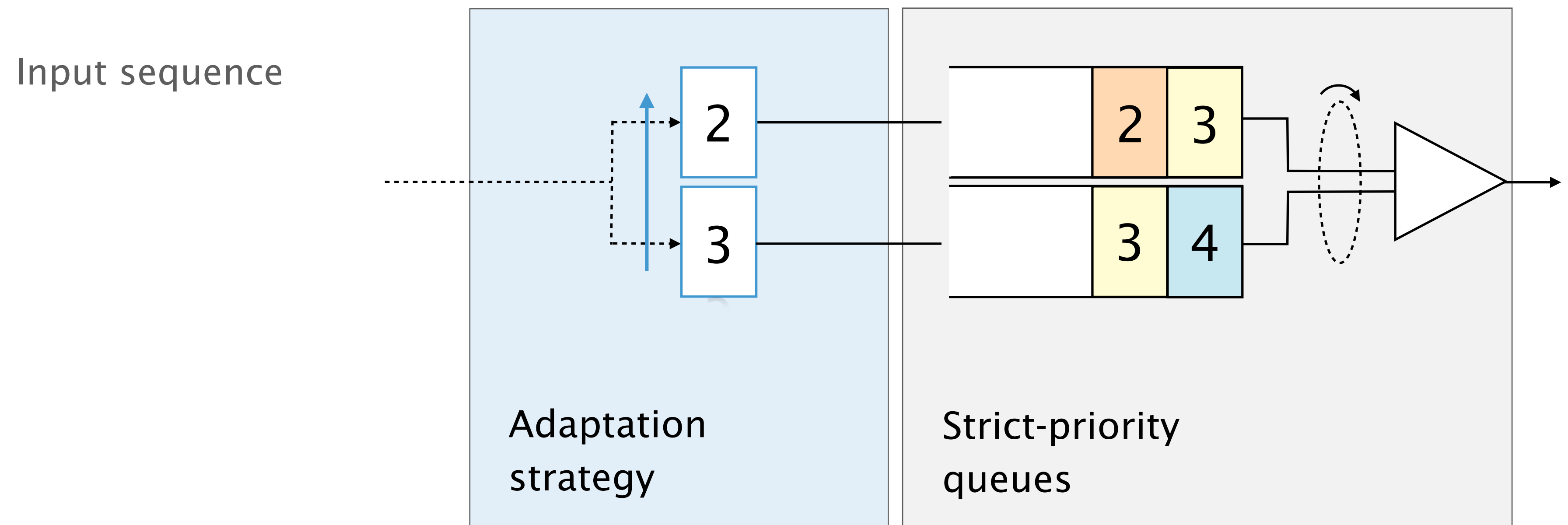
Input sequence

3

2

3

Adaptation strategy

2 3

4

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks

# to strict-priority queues



Mapping

Scan bottom-up, enqueue if rank >= bound

Input sequence

Adaptation strategy

Strict-priority queues

# SP-PIFO adapts the mapping of packet ranks
# to strict-priority queues

Mapping
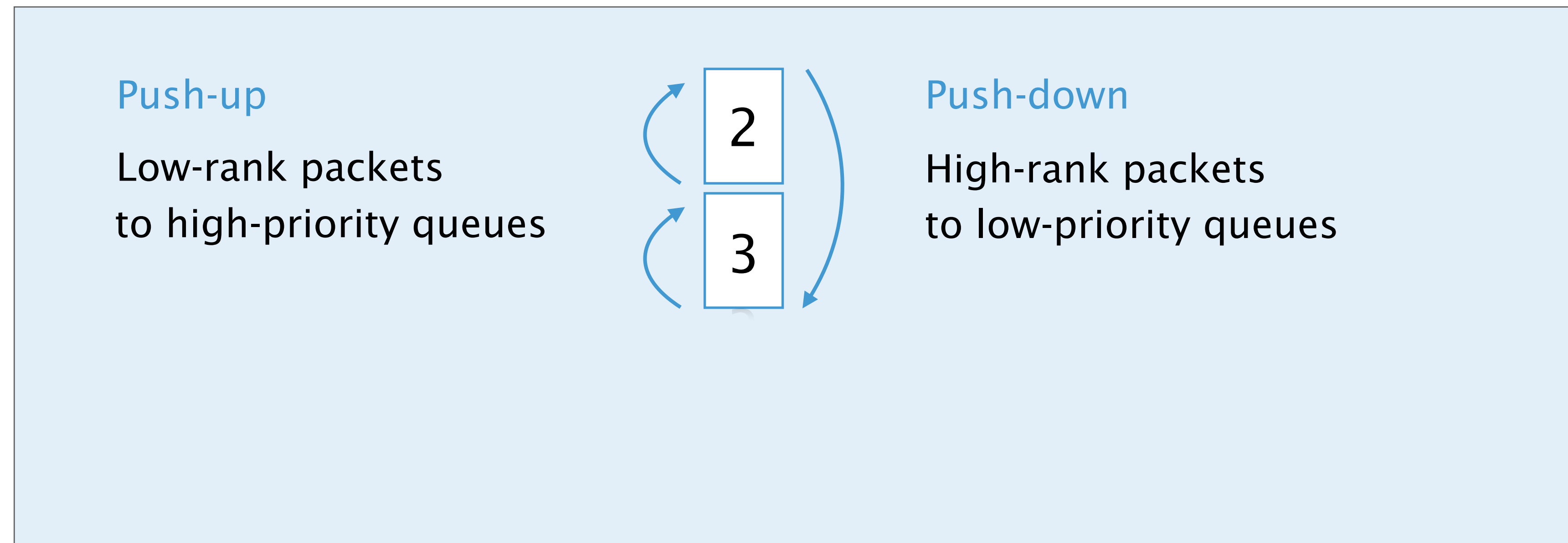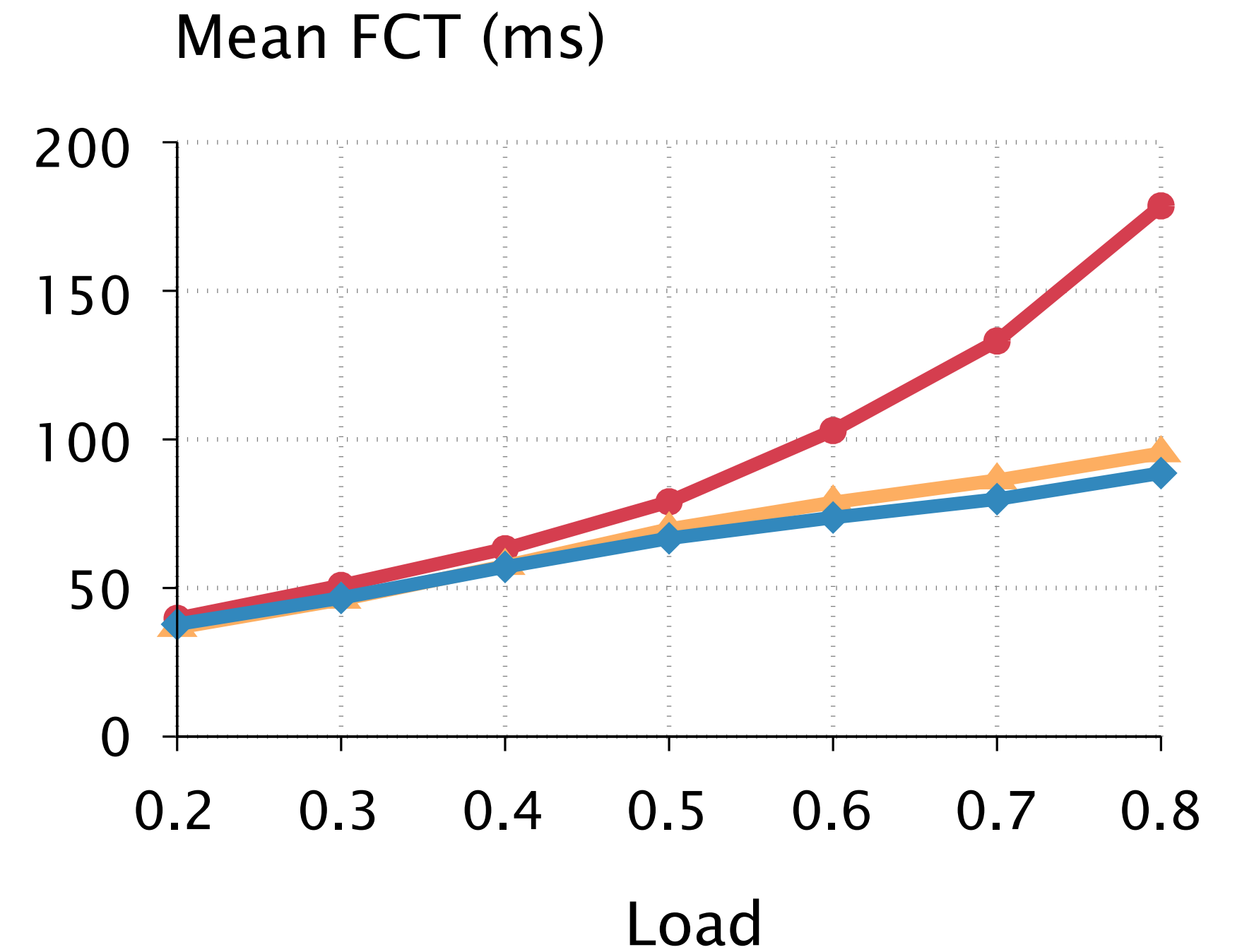
Scan bottom-up, enqueue if rank >= bound

Input sequence



Adaptation
strategy

Strict-priority
queues

# SP-PIFO adapts the mapping of packet ranks
# to strict-priority queues

Push-up

Low-rank packets
to high-priority queues

2

3

Push-down

High-rank packets
to low-priority queues

# SP-PIFO allows us to minimize flow completion times (FCTs)



Mean FCT (ms)

Small flows <100KB

Mean FCT (ms)

Big flows ≥1MB

How to enable programmable scheduling
on existing devices?

How to use it to improve
the Internet's security?

SP-PIFO

[NSDI'20]

Approximating
PIFO's scheduling

PACKS

[NSDI'25]

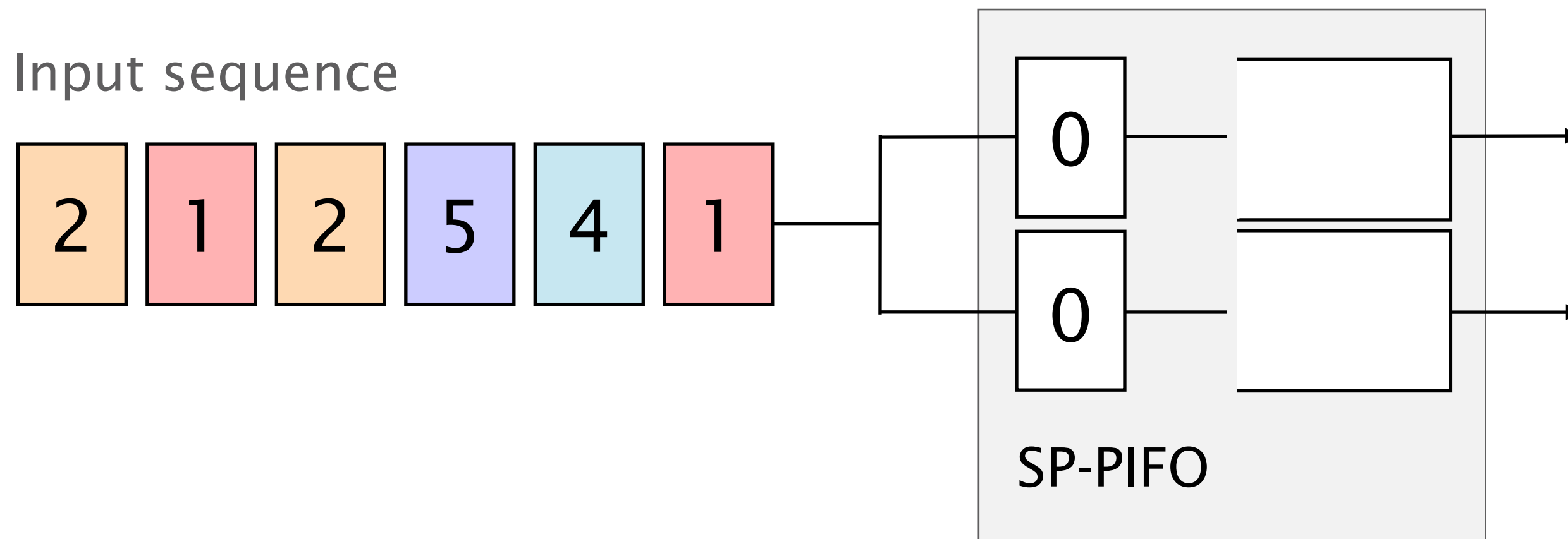Incorporating
PIFO's admission

ACC-Turbo

[SIGCOMM'22]

Mitigating
DDoS attacks

# PIFO's admission prevents the dropping of important packets
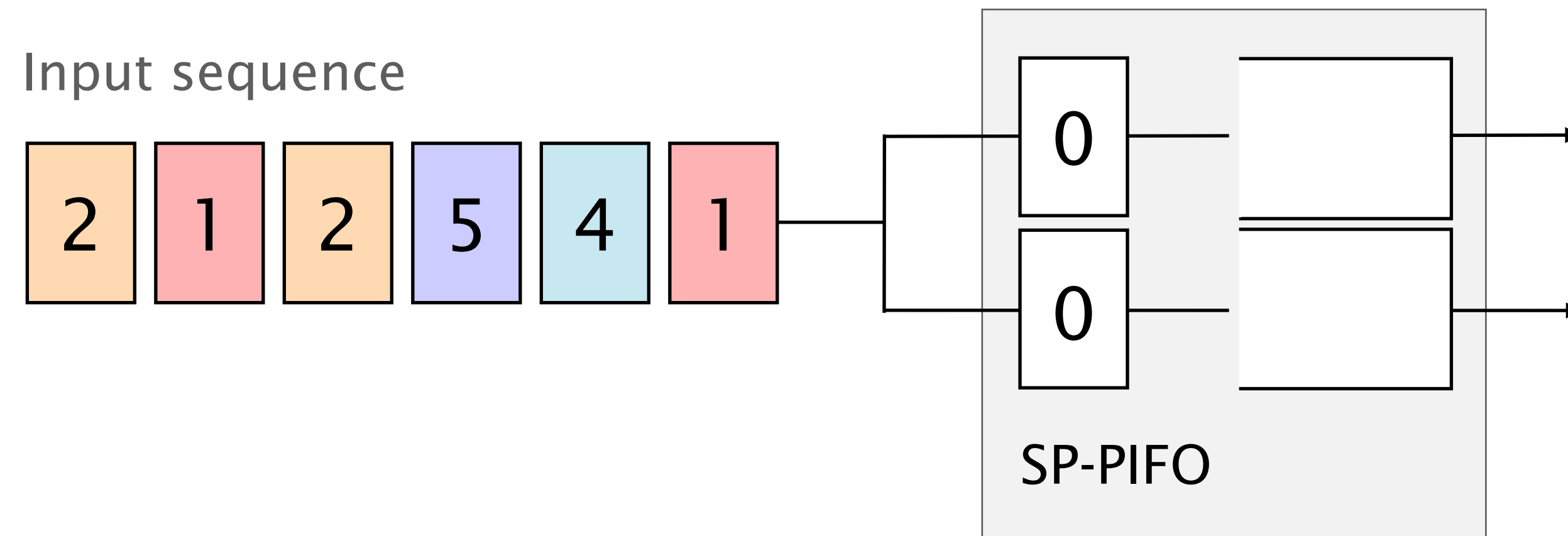
Input sequence



PIFO queue
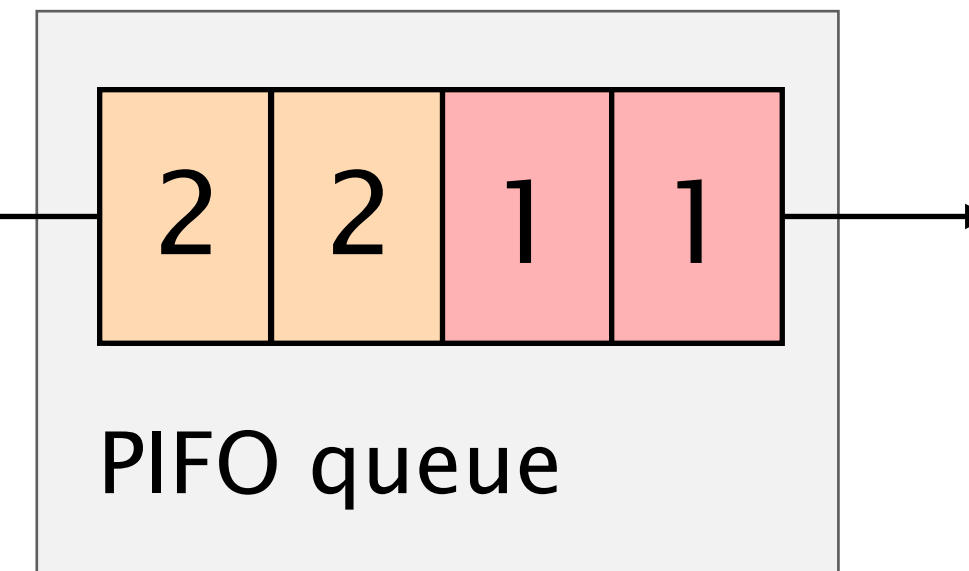
Input sequence

SP-PIFO

# PIFO's admission prevents the dropping of important packets

Input sequence

Dropped

PIFO queue

Input sequence

SP-PIFO

# PIFO's admission prevents the dropping of important packets



Input sequence

2 1 2 5 4 1

Dropped

2 2 1 1

PIFO queue

Input sequence

2 1 2 5 4 1

Important packets are dropped

SP-PIFO

2
5

2 1

5 4
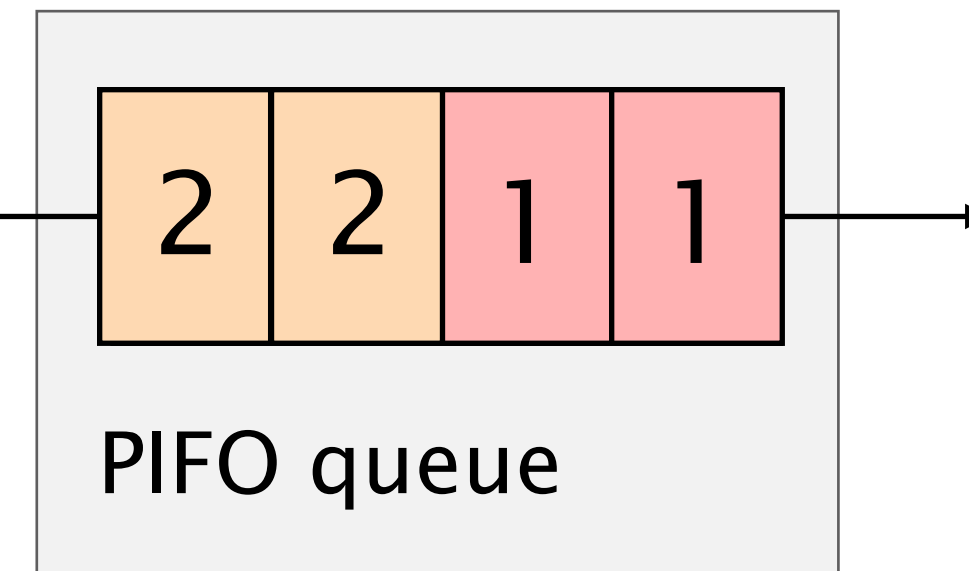
Non-important packets take buffer space

# We need to preemptively block non-important packets

# We need to preemptively block non-important packets

Input sequence

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

r < ?

?

?

PACKS

# PACKS monitors the rank distribution and the queue occupancy

# PACKS monitors the rank distribution and the queue occupancy

Input sequence

2 1 2 5 4 1

r < ?

?

?

PACKS

Rank distribution (W)

1 2

1 2

4 5

B

Dropped

Buffer availability

B = 4 packets

# PACKS monitors the rank distribution and the queue occupancy

Input sequence

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

r < 3

?

?

$r_{drop} = 3$

Buffer availability

B = 4 packets

Rank distribution (W)

| 1 | 2 |
| 1 | 2 |

| 4 | 5 |

B          Dropped

# PACKS monitors the rank distribution and the queue occupancy

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

$r_{drop} = 3$

maximize $r_{drop}$

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

# PACKS monitors the rank distribution and the queue occupancy

Input sequence

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

$r < 3$

? ?

Queues availability

$B_1$ = 2 packets

$B_2$ = 2 packets

Rank distribution (W)

| 1 | 2 |
| 1 | 2 |

4  5

Dropped

# PACKS monitors the rank distribution and the queue occupancy

Input sequence

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

$r < 3$

? ?

Rank distribution (W)

| 1 | 2 |
| 1 | 2 |
B1   B2   |   4   5   Dropped

Queues availability

B1 = 2 packets

B2 = 2 packets

# PACKS monitors the rank distribution and the queue occupancy

## SP–PIFO

Per-packet heuristic

No traffic knowledge

No queue information

## PACKS

Window-based

Rank-distribution aware

Queue-occupancy aware

## SP-PIFO

Scheduling ✔

## PACKS

Scheduling ✔

Admission ✔

# PACKS reduces inversions by up to 7x and drops by up to 60% with respect to SP-PIFO



Number of inversions

Number of drops

How to enable programmable scheduling
on existing devices?

How to use it to improve
the Internet's security?

| SP-PIFO | PACKS | ACC-Turbo |
|---|---|---|
| [NSDI'20] | [NSDI'25] | [SIGCOMM'22] |

Approximating
PIFO's scheduling

Incorporating
PIFO's admission

Mitigating
DDoS attacks

Pulse-wave DDoS attacks are

a new type of network-layer DDoS attack

# Pulse-wave DDoS attacks are

# a new type of network-layer DDoS attack

Target
a critical link



User

Attacker

Target Link

Critical
Services

# Pulse-wave DDoS attacks are

## a new type of network-layer DDoS attack

Target
a critical link

Volumetric
(Gbps)

Multiple
attack vectors



User

Attacker

Target Link

Critical
Services

# Pulse-wave DDoS attacks are
# a new type of network-layer DDoS attack

Target
a critical link

Volumetric

(Gbps)

Multiple
attack vectors

Short high–rate pulses

User

Attacker

Target Link

Critical
Services

# Pulse-wave DDoS attacks are an extreme case of congestion

High
Throughput

Different vectors (NTP, DNS …)

07:00          07:30          08:00          08:30     Time

Short duration: ~ 1 min

Damian Menscher, Google 2021

Pulse-wave DDoS attacks exploit
the limitations of existing defenses

# Pulse-wave DDoS attacks exploit
# the limitations of existing defenses

Narrow
attack coverage

*Signature-based*
*Access-control lists*

In–Network
Defense

User

Attacker

Target Link

Critical
Services

# Pulse-wave DDoS attacks exploit
# the limitations of existing defenses

Narrow
attack coverage

Filter-based
Rerouting-based

Drastic
mitigation



In–Network
Defense

User

Attacker

Target Link

Critical
Services

# Pulse-wave DDoS attacks exploit
# the limitations of existing defenses

# Pulse-wave DDoS attacks exploit
# the limitations of existing defenses

Narrow
attack coverage

*If count > threshold:*
*mitigate;*

Offline
Appliance

Drastic
mitigation

In–Network
Defense

Slow
reaction time

Target Link

User

Critical
Services

Risk of
misconfiguration

Attacker

# A pulse-wave DDoS defense needs to be …

**Fast**

reaction

**Generic**

detection

# A pulse-wave DDoS defense needs to be …

Fast

reaction

In-network, at line rate

with limited resources

Generic

detection

Unsupervised techniques

with uncertainty

# A pulse-wave DDoS defense needs to be …

Fast
reaction

In-network, at line rate
with limited resources

Generic
detection

Unsupervised techniques
with uncertainty

Risk of false positives

# A pulse-wave DDoS defense needs to be …

| | | |
|---|---|---|
| **Fast** reaction | In-network, at line rate with limited resources | |
| **Generic** detection | Unsupervised techniques with uncertainty | Risk of false positives |
| **Safe** mitigation | Limited impact under misclassification | |

# A pulse-wave DDoS defense needs to be …

| | | |
|---|---|---|
| **Fast** | In-network, at line rate | |
| reaction | with limited resources | |
| | | Risk of false positives |
| **Generic** | Unsupervised techniques | |
| detection | with uncertainty | |
| | | |
| **Safe** | Limited impact | Filtering ✘ |
| mitigation | under misclassification | Rerouting ✘ |

# A pulse-wave DDoS defense needs to be …

Fast
reaction

In-network, at line rate
with limited resources

Generic
detection

Unsupervised techniques
with uncertainty

Safe
mitigation

Limited impact
under misclassification

Risk of false positives

Programmable
scheduling

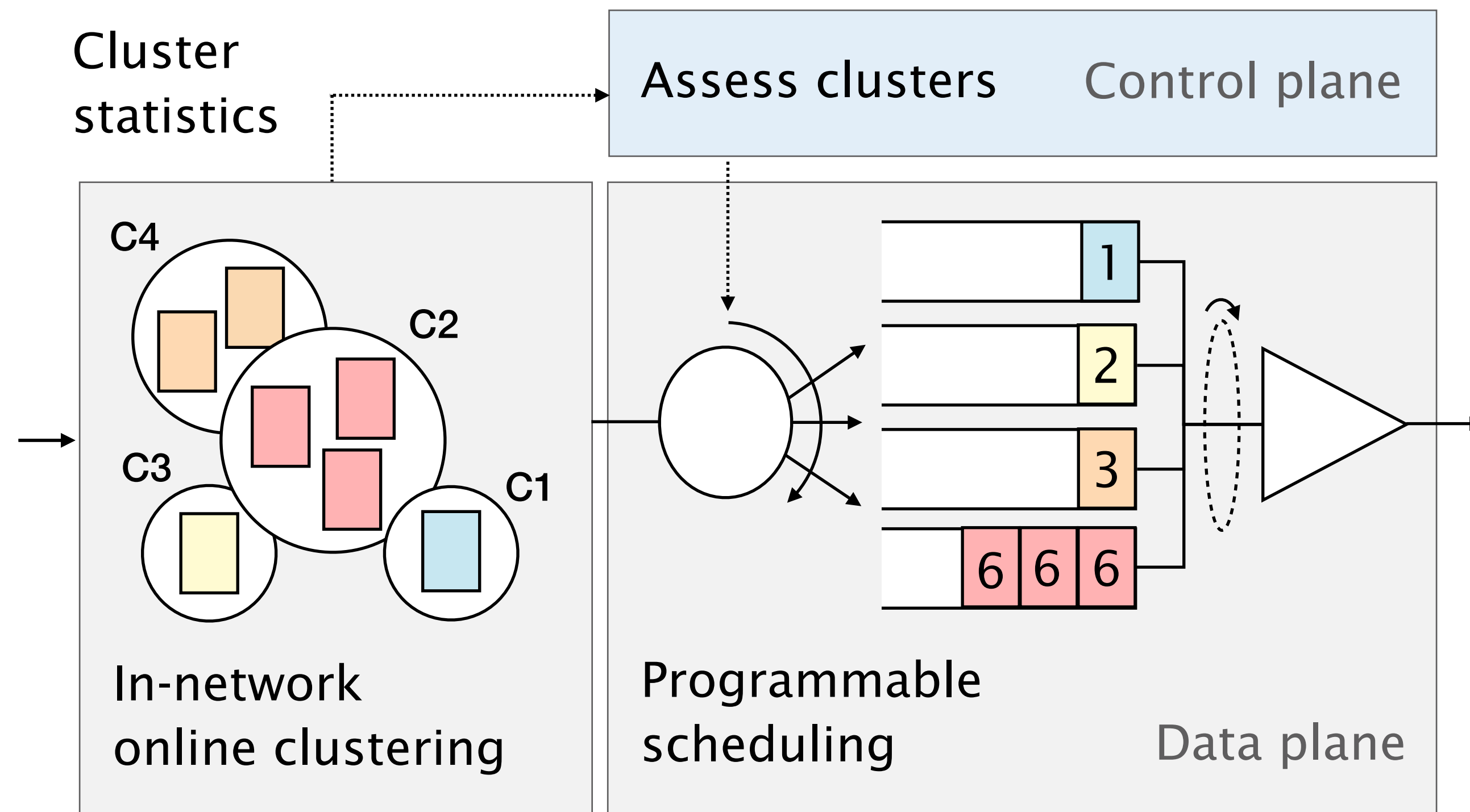# Programmable scheduling is a safe mitigation technique

Leverages the whole uncertainty spectrum

with fine-grained scheduling policies

Only drops under congestion

starting by most-malicious packets

Does not require activation

can be always-on

# ACC-Turbo utilizes online clustering and programmable scheduling

# ACC-Turbo outperforms existing defenses and mitigates pulse-wave DDoS attacks

No defense

# ACC-Turbo outperforms existing defenses and mitigates pulse-wave DDoS attacks

ACC-Turbo

# ACC-Turbo outperforms existing defenses and mitigates pulse-wave DDoS attacks

ACC-Turbo



Output Attack ▬▬▬ Output Benign ▬▬▬

Throughput

Time

Deprioritize
(Not filter)

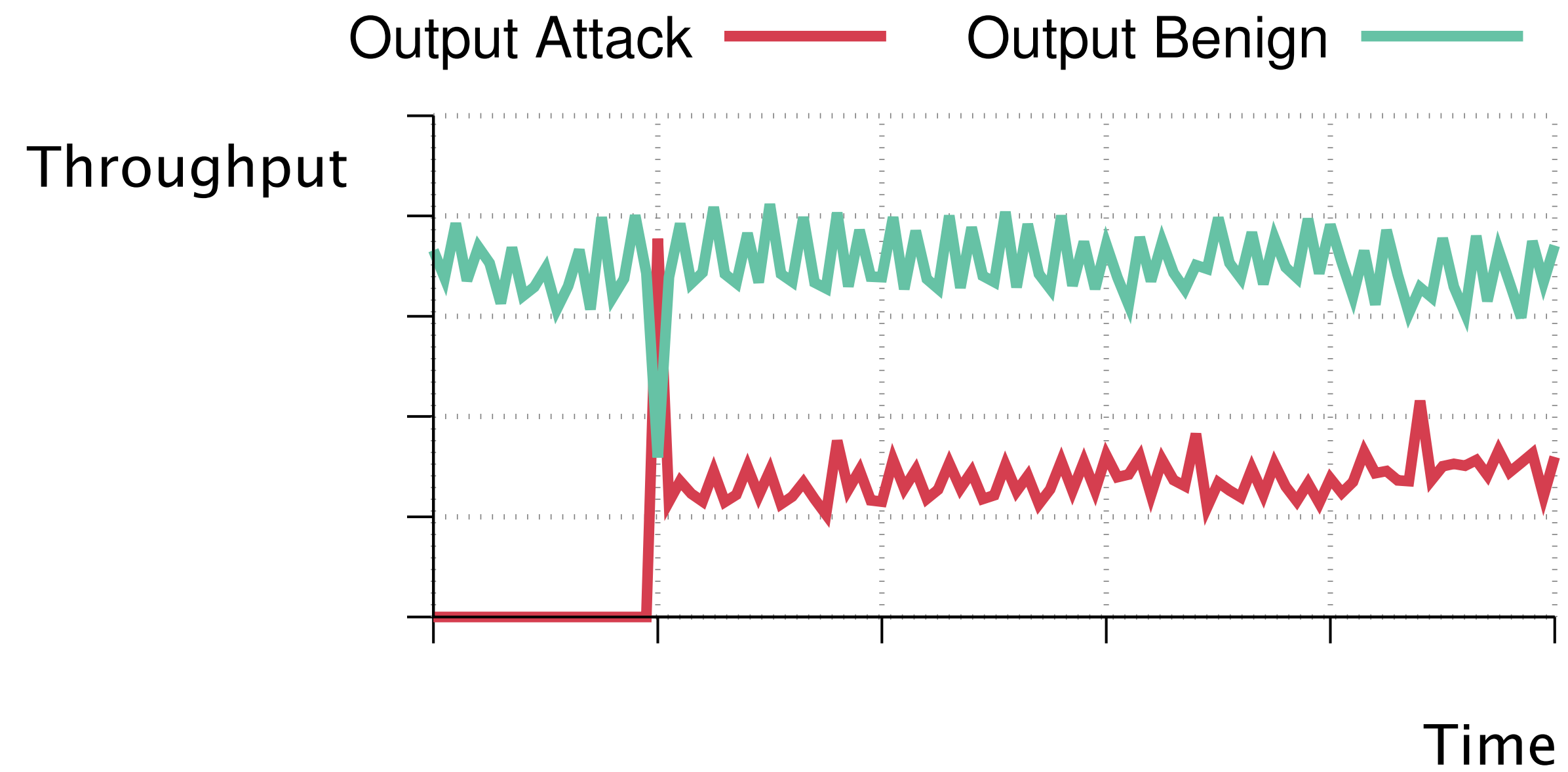# ACC-Turbo outperforms existing defenses and mitigates pulse-wave DDoS attacks

ACC-Turbo

How to enable programmable scheduling
on existing devices?

How to use it to improve
the Internet's security?

| SP-PIFO | PACKS | ACC-Turbo |
|---|---|---|
| [NSDI'20] | [NSDI'25] | [SIGCOMM'22] |

Approximating
PIFO's scheduling

Incorporating
PIFO's admission

Mitigating
DDoS attacks

# In-Network Congestion Management for Security and Performance

How to enable programmable scheduling
on existing devices?

How to use it to improve
the Internet's security?

| SP-PIFO | PACKS | ACC-Turbo |
|---------|-------|-----------|
| [NSDI'20] | [NSDI'25] | [SIGCOMM'22] |

Approximating

PIFO's scheduling

Incorporating

PIFO's admission

Mitigating

DDoS attacks

# Selected publications

**NSDI '20**        SP-PIFO: Approximating Push-In First-Out Behaviors using Strict-Priority Queues
A. Gran Alcoz, A. Dietmüller, L. Vanbever

**SIGCOMM '22**    Aggregate-Based Congestion Control for Pulse-Wave DDoS Defense
A. Gran Alcoz, M. Strohmeier, V. Lenders, L. Vanbever

**HotNets '23**    QVISOR: Virtualizing Packet Scheduling Policies
A. Gran Alcoz, L. Vanbever

**SIGCOMM '24**    Principles for Internet Congestion Management
L. Brown, A. Gran Alcoz, F. Cangialosi, A. Narayan, M. Alizadeh, H. Balakrishnan,
E. Friedman, E. Katz-Bassett, A. Krishnamurthy, M. Schapira, S. Shenker

**NSDI '25**        Everything Matters in Programmable Packet Scheduling
A. Gran Alcoz, B. Vass, P. Namyar, B. Arzani, G. Rétvári, L. Vanbever