# *Seer: Enabling Future-Aware Online Caching in Networked Systems*

**Jason Lei,** Vishal Shrivastav
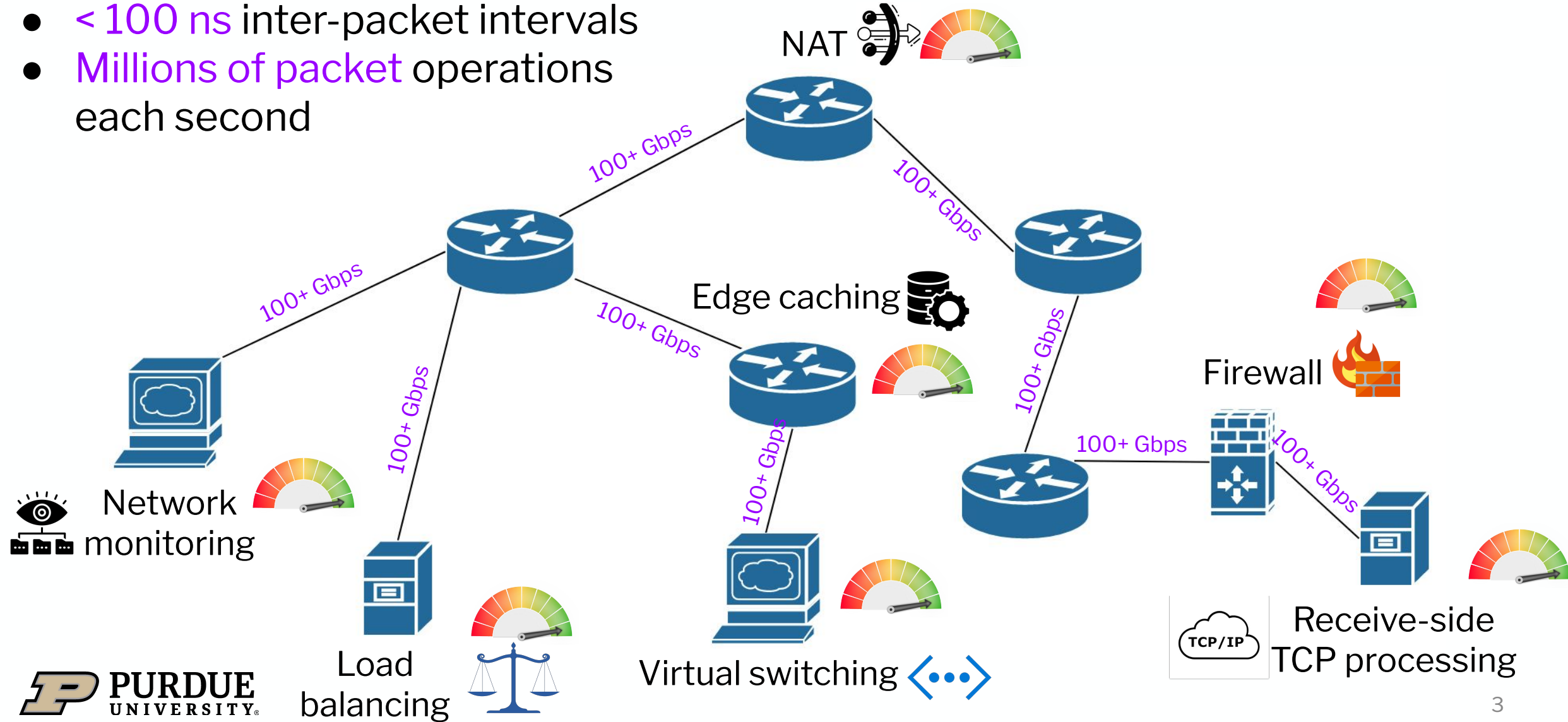
PURDUE UNIVERSITY®

nsdi '24

# Motivation: *State-Intensive High-Speed Network Applications*



NAT

- 300-600 MB tables
  to serve common network
  applications [TEA, SIGCOMM '20]

Edge caching

Firewall

Network
monitoring

Load
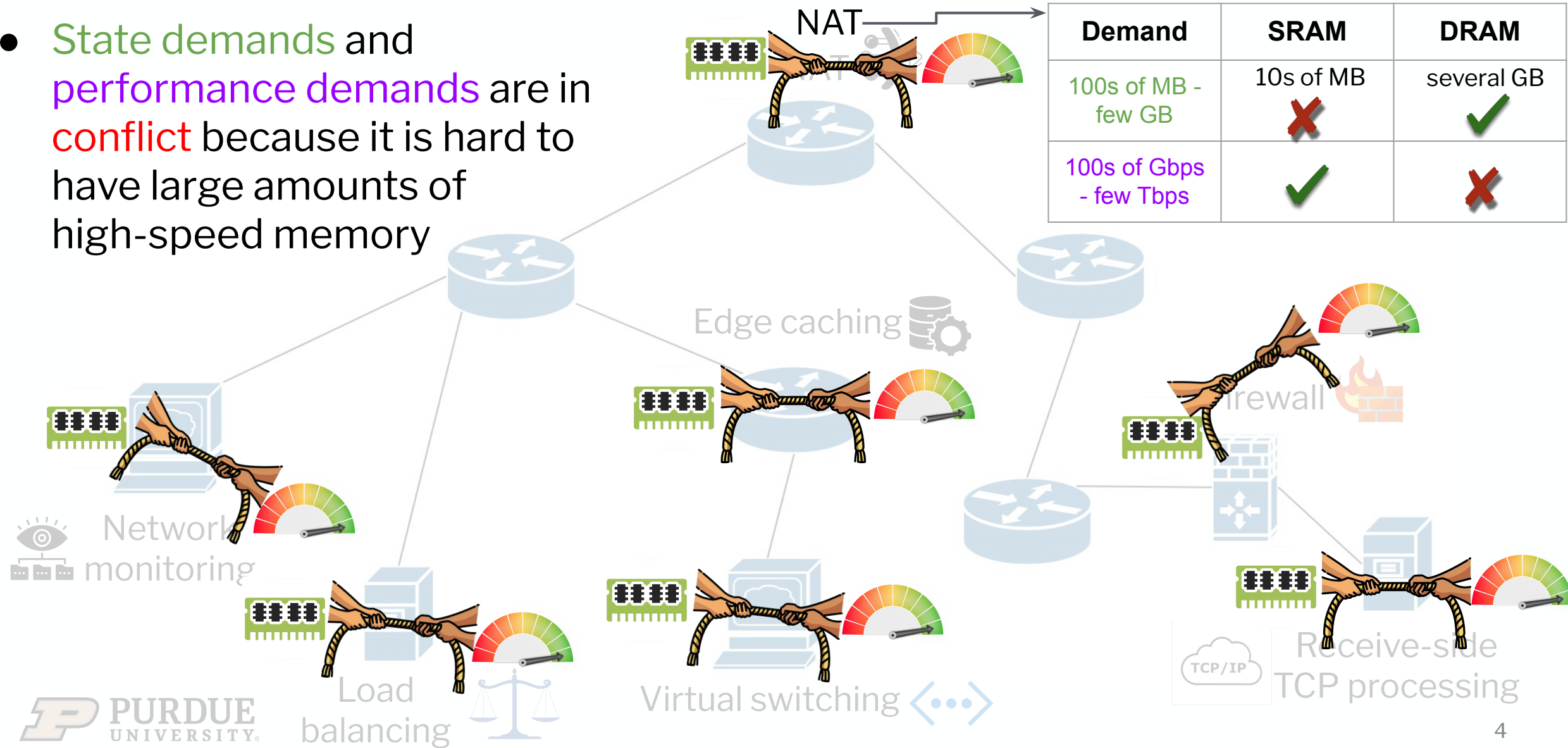balancing

Virtual switching

Receive-side
TCP processing

# Motivation: *State-Intensive* *High-Speed* *Network Applications*

- < 100 ns inter-packet intervals
- Millions of packet operations each second

NAT

100+ Gbps

100+ Gbps

Edge caching

100+ Gbps

100+ Gbps

100+ Gbps

100+ Gbps

Firewall

100+ Gbps

100+ Gbps

Network monitoring

Load balancing

Virtual switching

TCP/IP

Receive-side TCP processing

PURDUE UNIVERSITY®

3

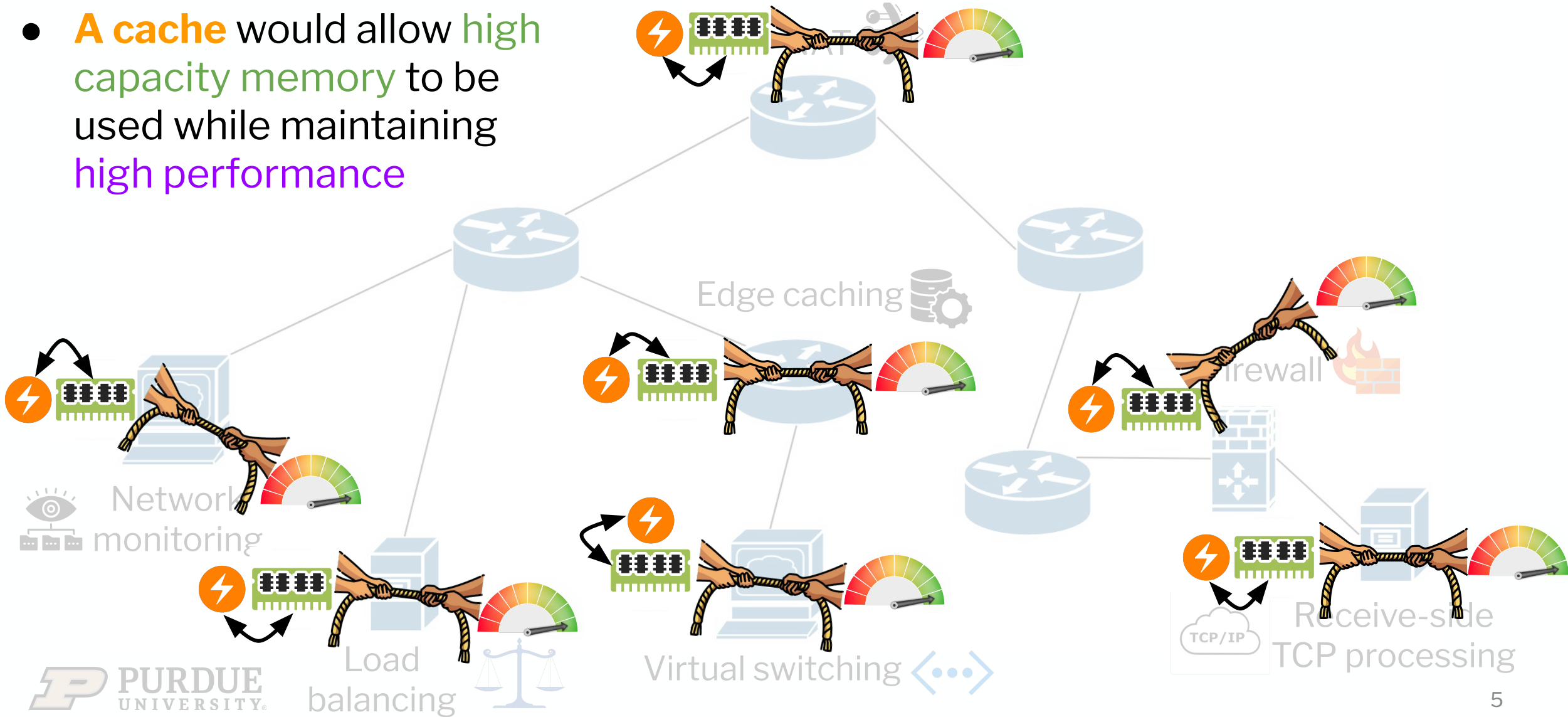# *Motivation: State-Intensive High-Speed Network Applications*

- State demands and performance demands are in conflict because it is hard to have large amounts of high-speed memory

| Demand | SRAM | DRAM |
|---|---|---|
| 100s of MB - few GB | 10s of MB ✗ | several GB ✓ |
| 100s of Gbps - few Tbps | ✓ | ✗ |

NAT

Edge caching

Firewall

Network monitoring

Load balancing

Virtual switching

Receive-side TCP processing

PURDUE UNIVERSITY

4

# Motivation: *State-Intensive* *High-Speed* *Network Applications*

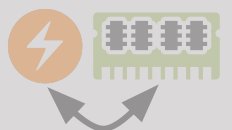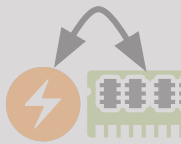- **A cache** would allow high capacity memory to be used while maintaining high performance

Edge caching

Firewall

Network monitoring

Load balancing

Virtual switching

Receive-side TCP processing

*Motivation: State-Intensive High-Speed Network Applications*

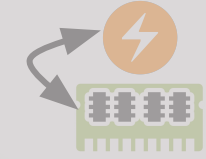- **A cache** would allow high capacity memory to be used while maintaining high
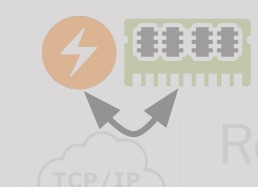
**Takeaway:**

**High-speed state-intensive** network applications require **efficient caching**

NAT

Network monitoring

Load balancing

Virtual switching

Receive-side TCP processing

PURDUE UNIVERSITY
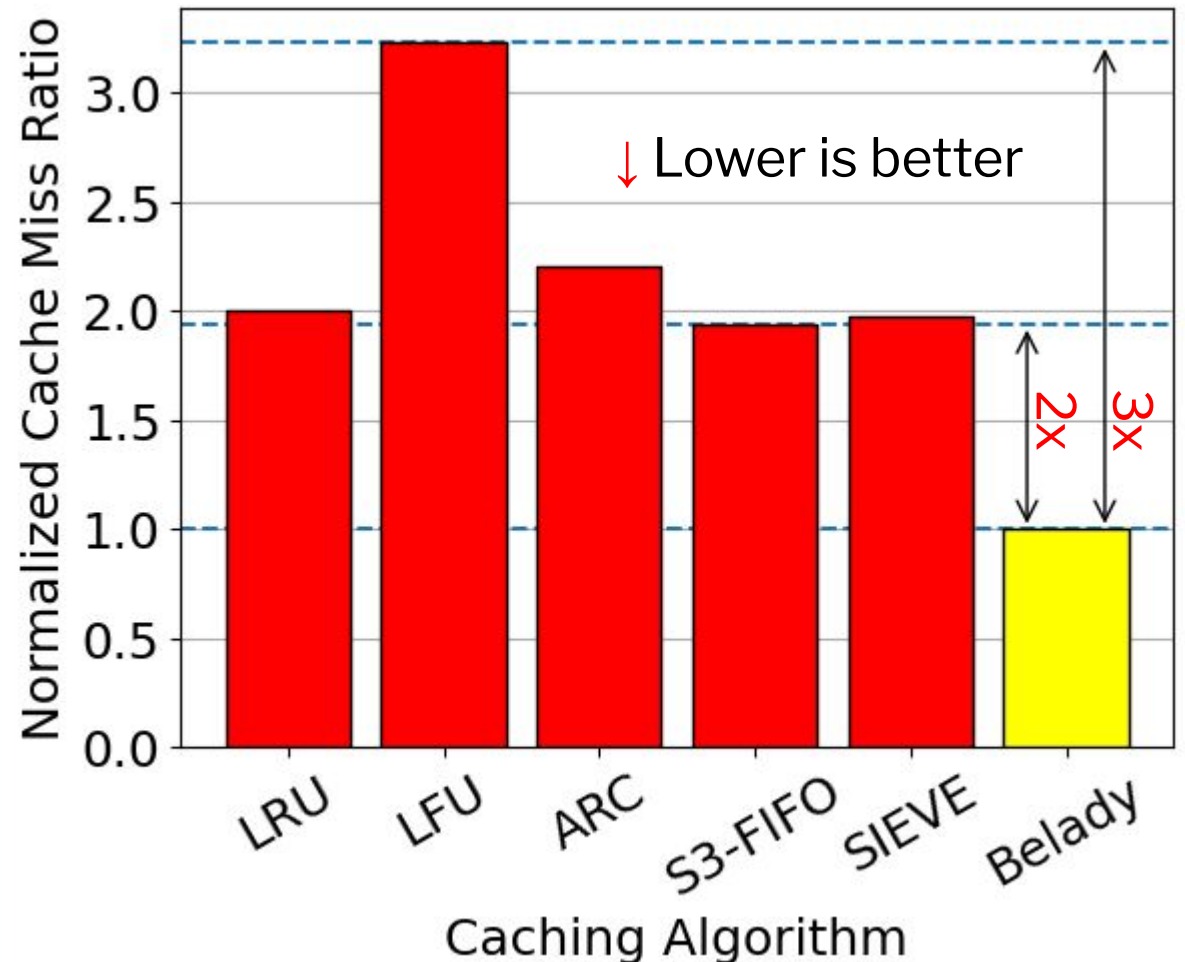
# *State of Practice Falls Short of Ideal Case*

- No shortage of online caching algorithms - LRU, LFU, ARC, CLOCK, S3-FIFO, SIEVE, etc.

- All fall short of optimal offline caching algorithm (Belady) by a significant margin

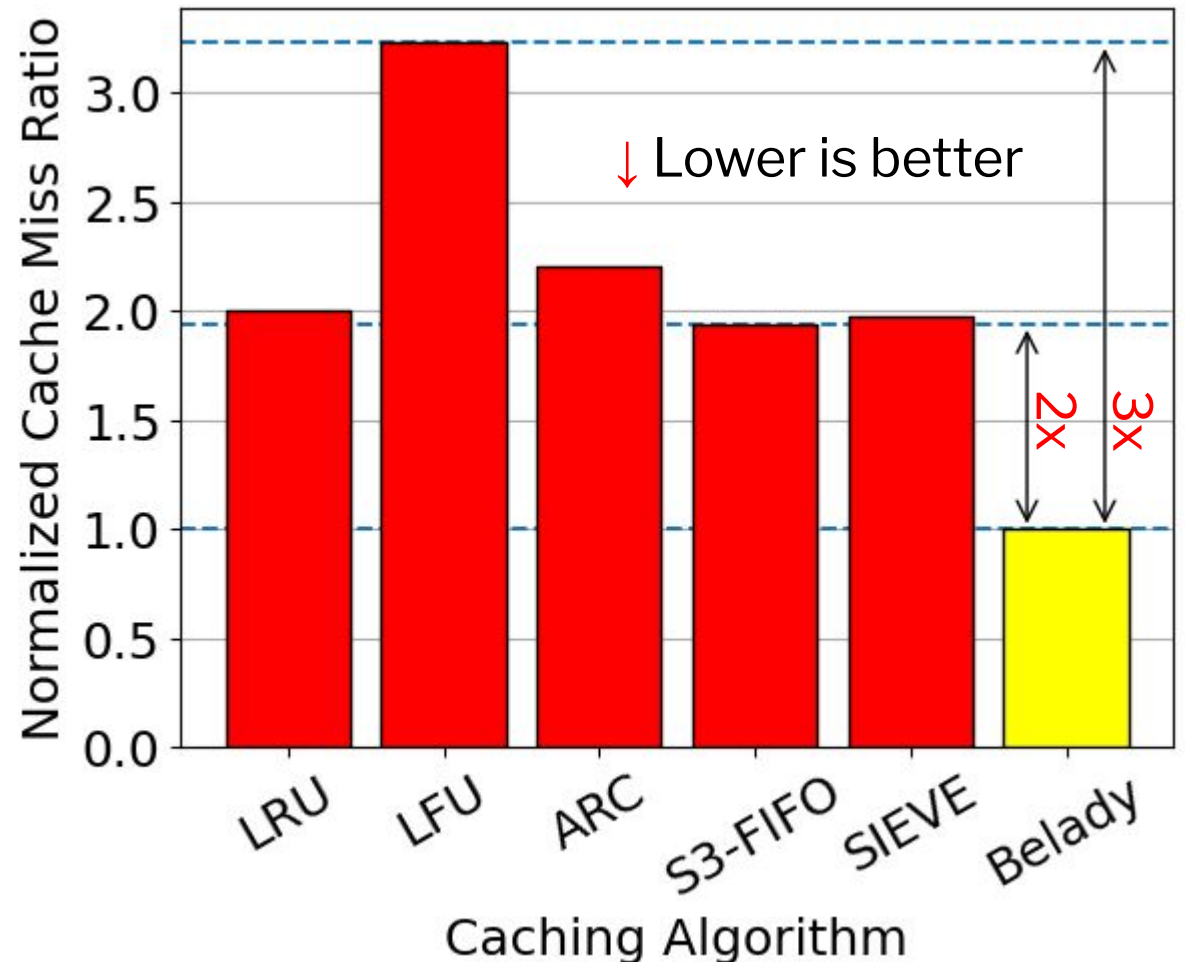  - Ranging from 2-3x higher cache miss ratio



**Setup:**
2-tier fat-tree network with 144 nodes running in-network load balancing application with websearch workload

# State of Practice Falls Short of Ideal Case

- No shortage of online caching algorithms - LRU, LFU, ARC, CLOCK, S3-FIFO, SIEVE, etc.

- All fall short of optimal offline caching algorithm (Belady) by a significant margin

  - Ranging from 2-3x higher cache miss ratio

- ML-based solutions are prone to mispredictions



↓ Lower is better

**Setup:**
2-tier fat-tree network with 144 nodes running in-network load balancing application with websearch workload

**Fundamental Cause of Performance Gap:**

Offline algorithm (Belady) uses knowledge of future state accesses to make optimal caching decisions, but ...

**Traditional online caching algorithms** lack **accurate awareness of future state accesses**

**Setup:**
2-tier fat-tree network with 144 nodes running in-network load balancing application with websearch workload

PURDUE UNIVERSITY®
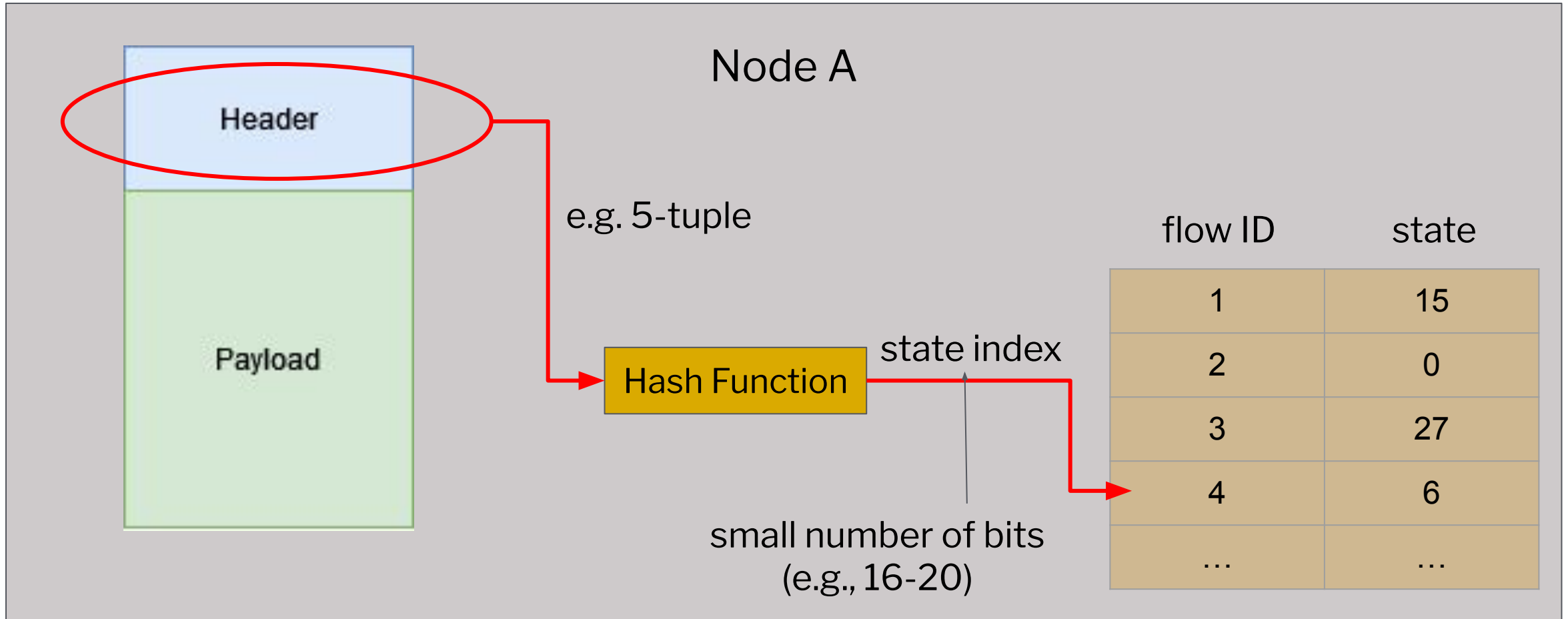
# *Key Research Question*

How can **future-aware caching** be realized accurately **in practice (online setting)?**

# *Key Insight*

*Traditional online caching assumes future-awareness is challenging.* ***However…***

## Networked setting presents unique opportunities to provide very accurate visibility into future state accesses!
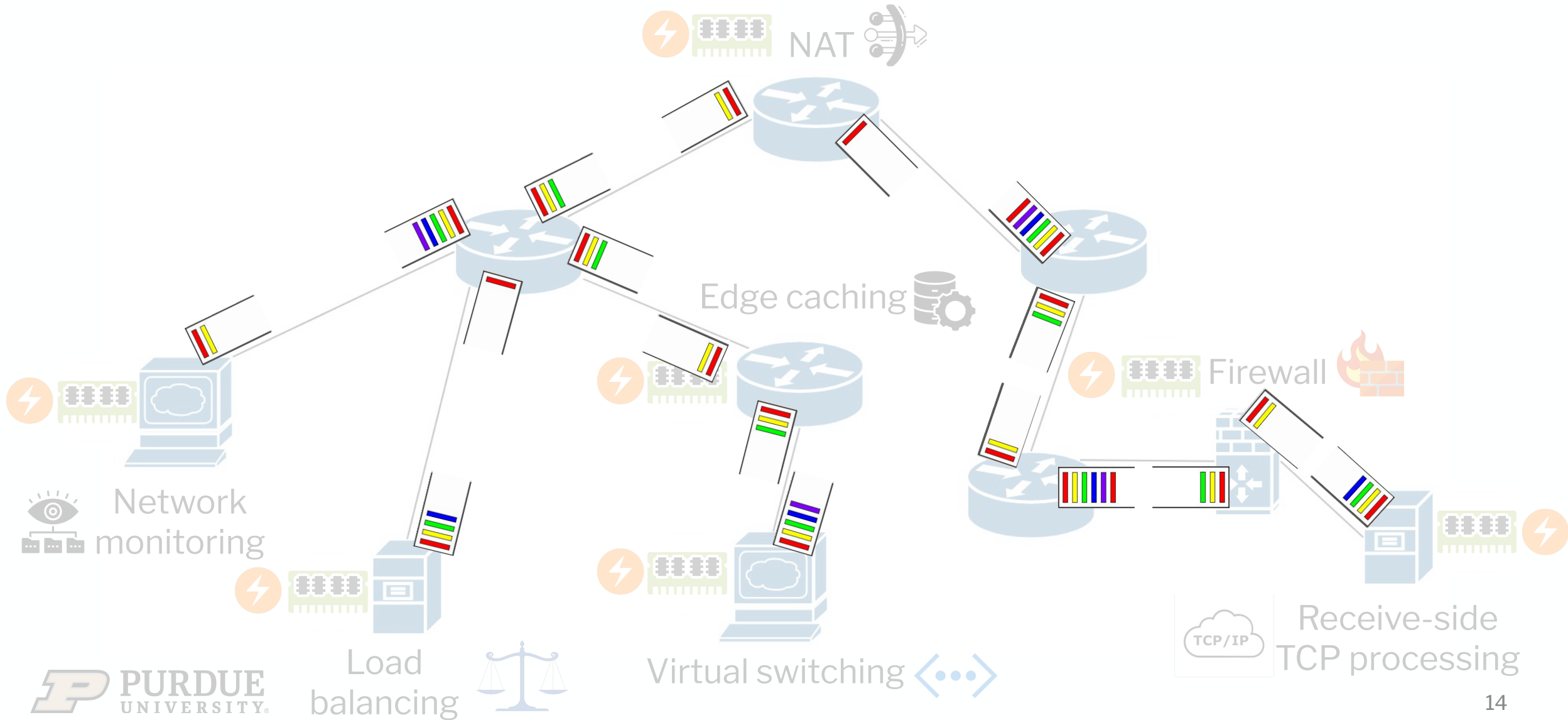
# Insight 1: Header-Based State Indexing

Node A



Header

Payload

e.g. 5-tuple

Hash Function

state index

small number of bits
(e.g., 16-20)

| flow ID | state |
|---------|-------|
| 1 | 15 |
| 2 | 0 |
| 3 | 27 |
| 4 | 6 |
| … | … |

Node A

**Takeaway 1:**

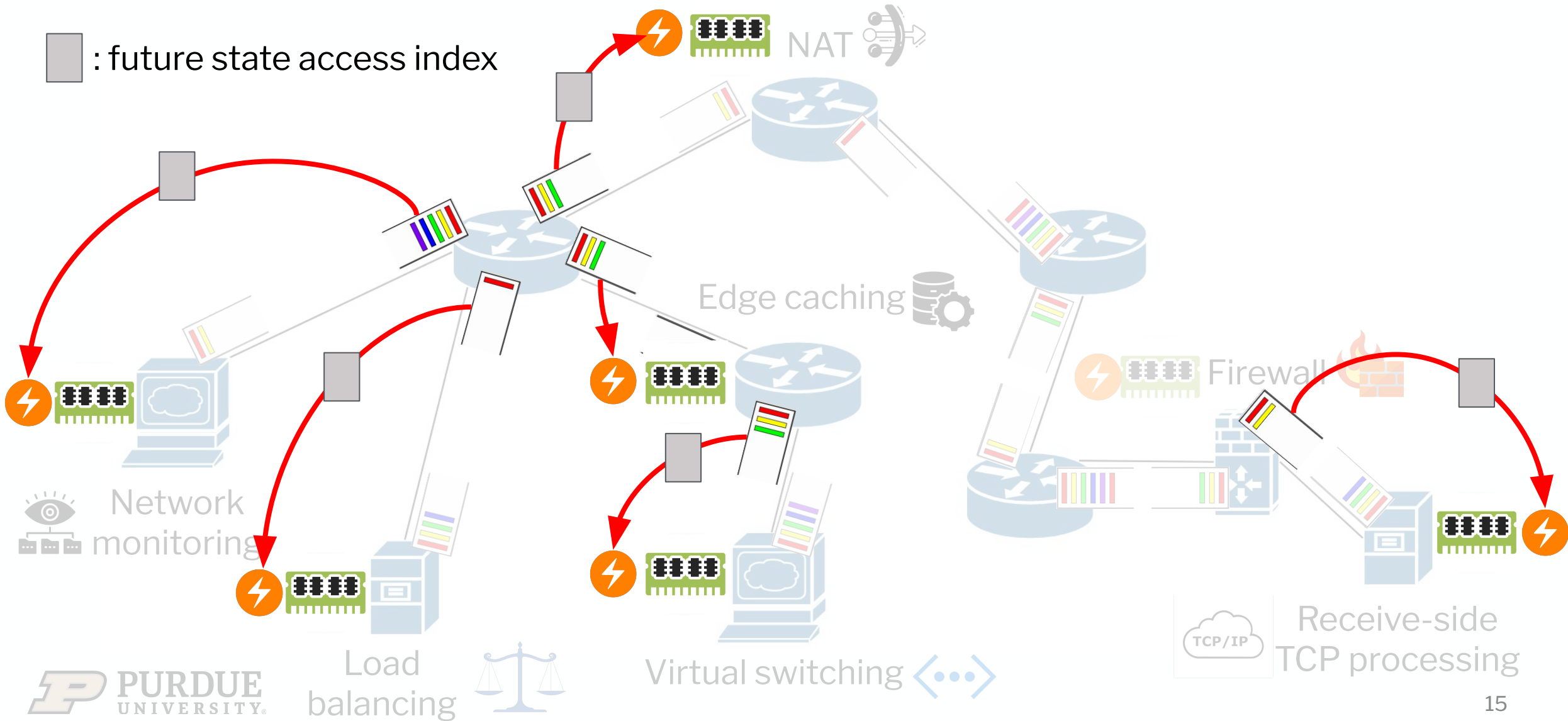**State access indices** are carried in incoming packet headers, and can be encoded using **a small number of bits**

(e.g., 16-20)

…        …

PURDUE
UNIVERSITY.

# Insight 2: Network Delays Create Opportunities



NAT

Edge caching

Firewall

Network monitoring

Load balancing

Virtual switching

Receive-side TCP processing

PURDUE UNIVERSITY

14

# Insight 2: Network Delays Create Opportunities



: future state access index

NAT

Edge caching

Firewall

Network monitoring

Load balancing

Virtual switching

Receive-side TCP processing

⬜ : future state access index

**Takeaway 2:**

**Delays in the network** can be leveraged to forward **state index information in advance!**

NAT

monitoring

Load balancing

Virtual switching

TCP/IP

Receive-side TCP processing

PURDUE UNIVERSITY.
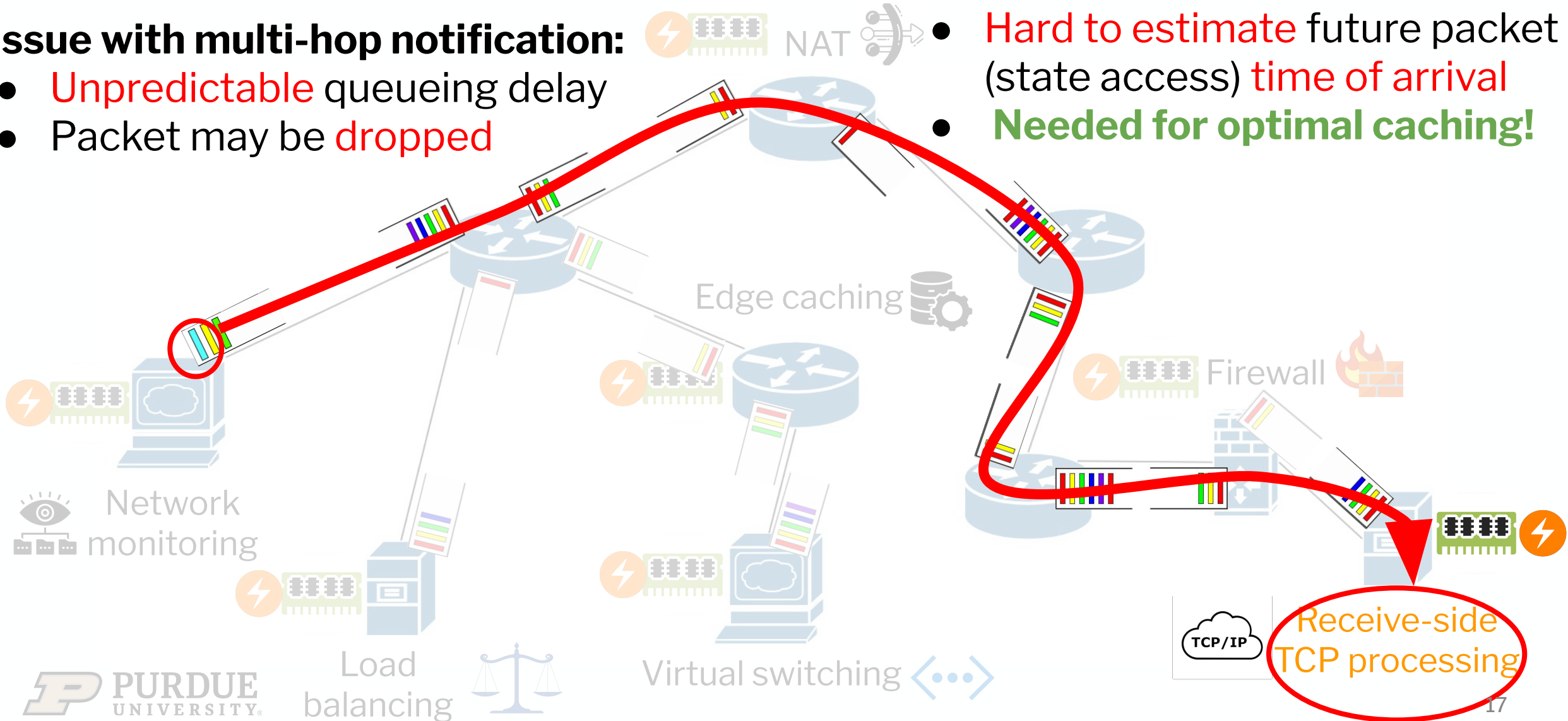
16

# Insight 3: Neighbors Are Most Accurate Notifiers

**Issue with multi-hop notification:**
- Unpredictable queueing delay
- Packet may be dropped

- Hard to estimate future packet (state access) time of arrival
- **Needed for optimal caching!**

NAT

Edge caching

Firewall

Network monitoring

Load balancing

Virtual switching

TCP/IP

Receive-side TCP processing

PURDUE UNIVERSITY.

# Insight 3: Neighbors Are Most Accurate Notifiers

**Issue with multi-hop notification:**
- Unpredictable queueing delay
- Pa...

- Hard to estimate future packet (state access) time of arrival
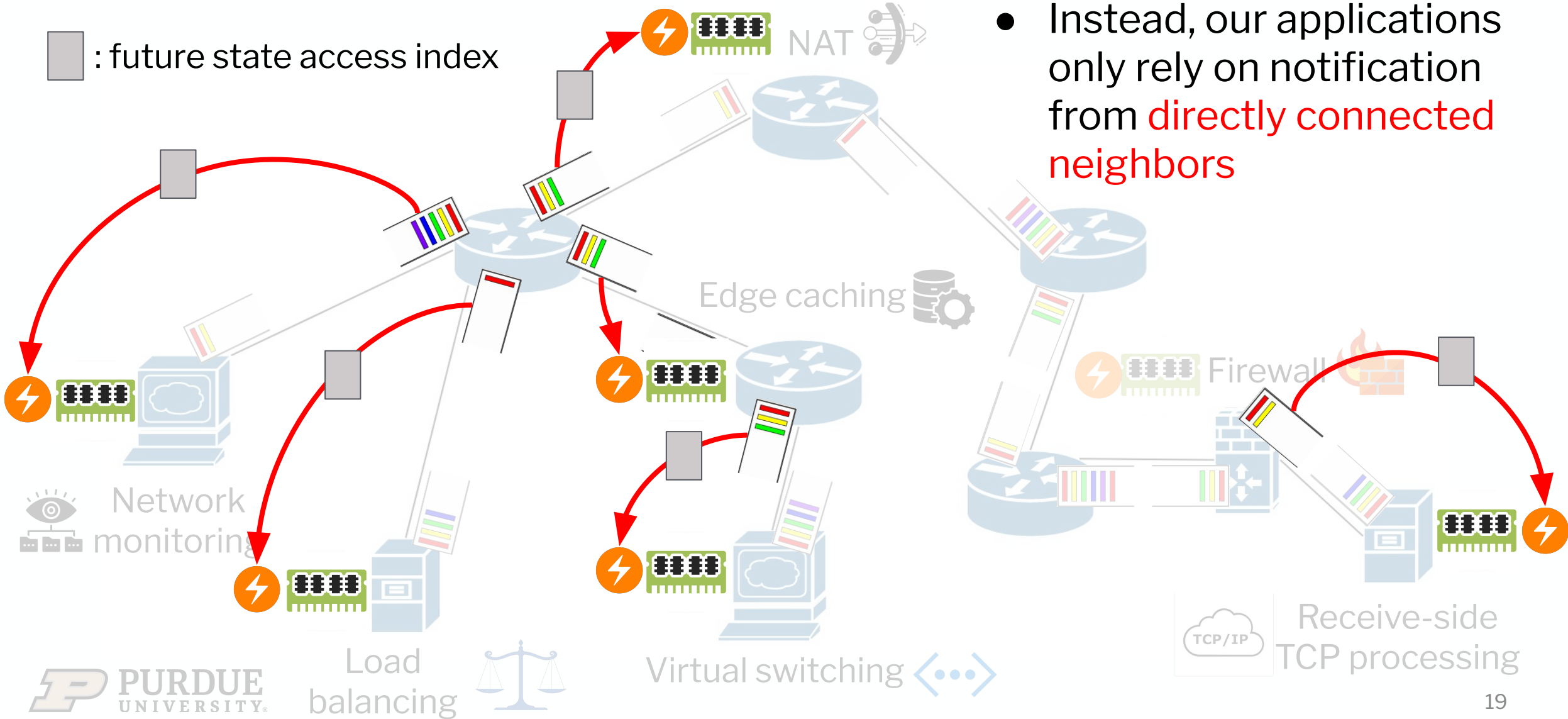- Needed for optimal caching!

**Takeaway 3:**

**Multi-hop notification** provides **inaccurate estimation** of future state **access time**, but **optimal algorithm (Belady)** heavily relies on it

NAT

Load balancing

Virtual switching

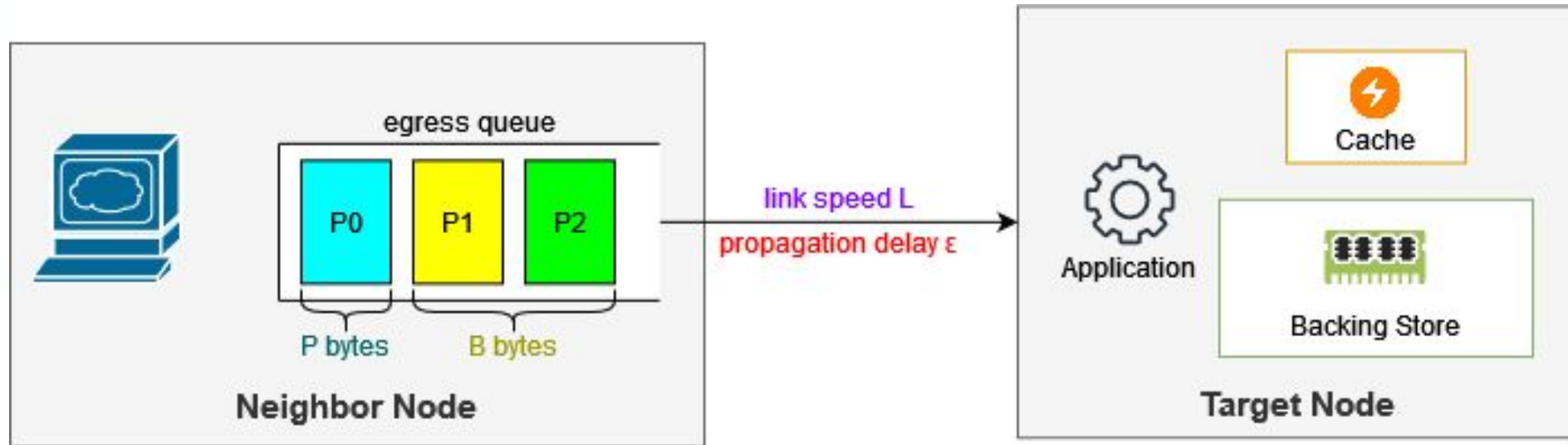TCP/IP

Receive-side TCP processing

PURDUE UNIVERSITY.

# Insight 3: Neighbors Are Most Accurate Notifiers



⬜ : future state access index

- Instead, our applications only rely on notification from directly connected neighbors

NAT

Edge caching

Firewall

Network monitoring

Load balancing

Virtual switching

Receive-side TCP processing

PURDUE UNIVERSITY

19

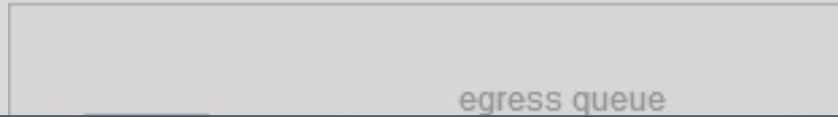# Insight 3: Neighbors Are Most Accurate Notifiers

# *Insight 3:  Neighbors Are Most Accurate Notifiers*



$$T = t_0 + (P + B)/L + \varepsilon$$

- **T**: time when packet will reach target node
- **$t_0$**: current time at target node
- **P**: size of packet in question P0
- **B**: bytes of queued data in front of packet P0
- **L**: link speed
- **ε**: link propagation delay

egress queue

**Takeaway 4:**

**Directly connected neighbors** provide a perfectly accurate **estimation of future state access times!**

- $t_0$: current time at target node
- **P**: size of packet in question P0
- **B**: bytes of queued data in front of packet P0
- **L**: link speed
- **ε**: link propagation delay

PURDUE
UNIVERSITY®

22

# *Putting It All Together*

1. **State access indices** are carried in incoming packet headers, and can be encoded using **a small number of bits.**

2. **Delays in the network** can be leveraged to forward state index information **in advance.**

3. **Directly connected neighbors** provide a **perfectly accurate estimate** of future state access times.

# *Our Contributions*

## *Seer: A Future-Aware Online Caching System*

1. Low-Overhead **Future State Access Notification**

2. Future-Aware **Cache Manager**

3. Fast **Hardware Implementation**

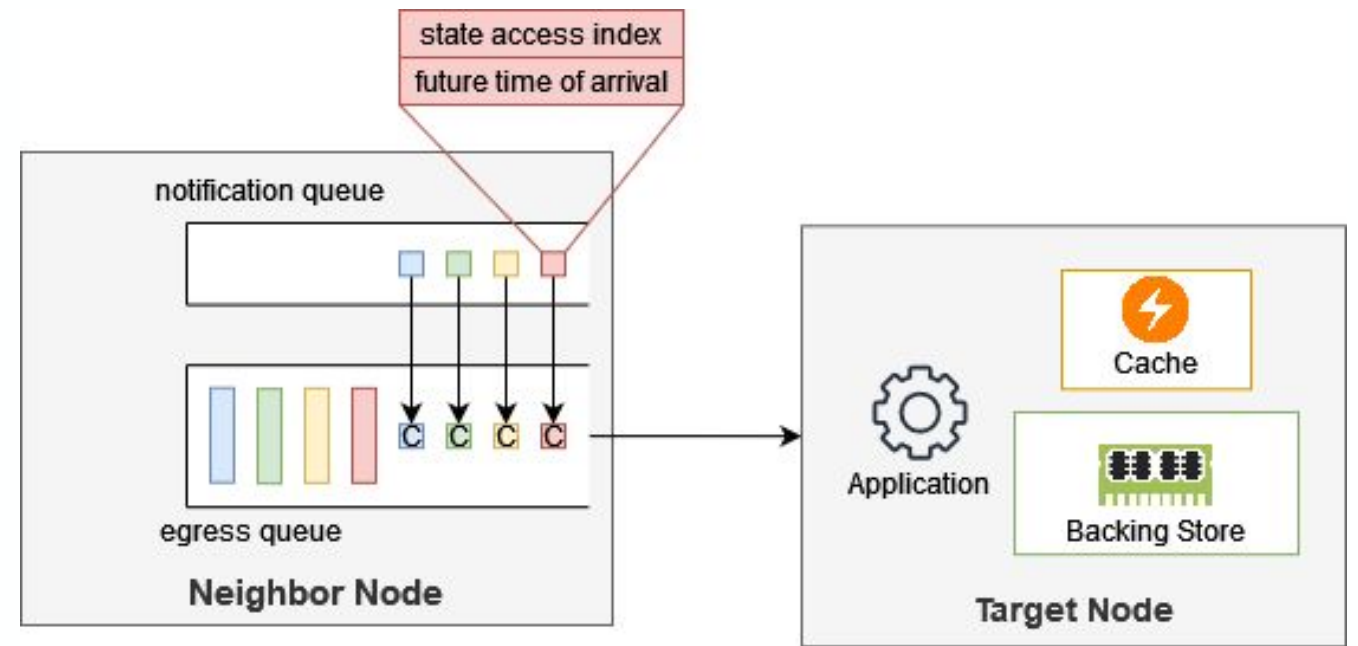# (1) Low-Overhead *Future State Access Notification*

**Future State Access Notification**

for each packet contains:

- State access index
- Future time of arrival of corresponding packet

**Naive solution for notification:**

**control packets**

- High bandwidth overhead – one control packet per data packet
  - If all pkts are minimum-sized, can consume half of total bandwidth!



State access notifications carried in control packets

# (1) Low-Overhead *Future State Access Notification*
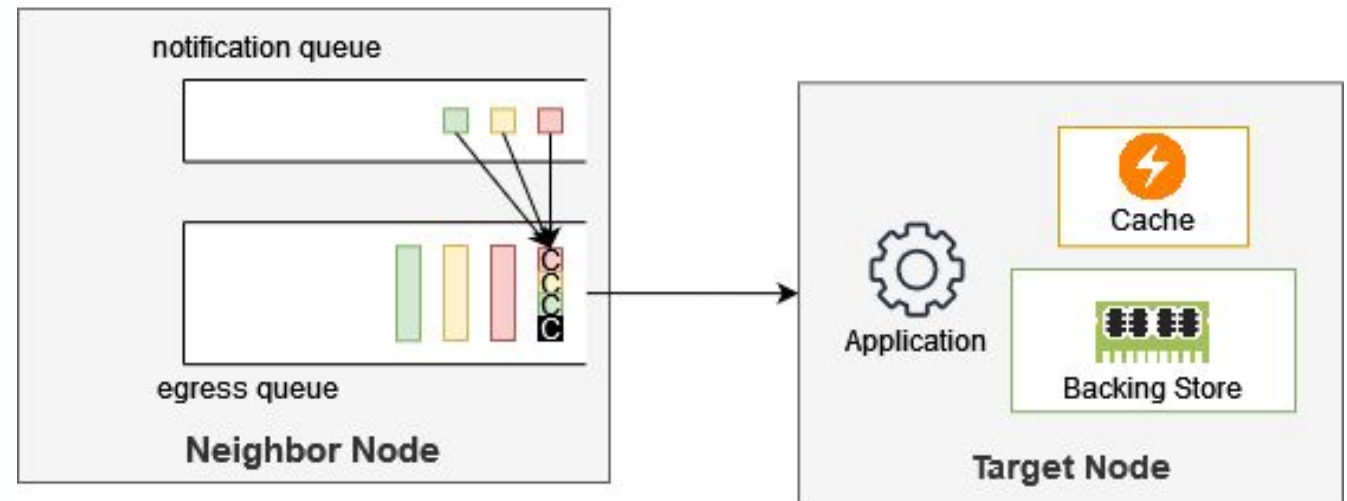
**Future State Access Notification**

for each packet contains:

- State access index
- Future time of arrival of corresponding packet

**Naive solution for notification: control packets**

- High bandwidth overhead – one control packet per data packet
  - If all pkts are minimum-sized, can consume half of total bandwidth!

Batching reduces number of control packets but can delay notification while it waits for batch

# (1) Low-Overhead *Future State Access Notification*
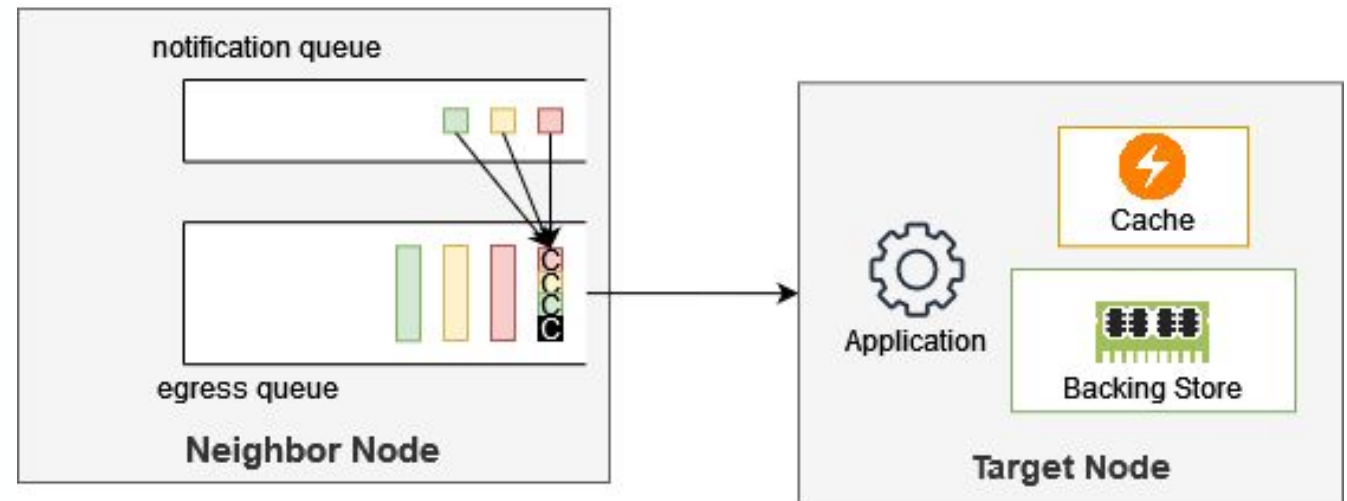
**Future State Access Notification**

for each packet contains:

- State access index
- Future time of arrival of corresponding packet

**Naive solution for notification:**

**control packets**

- High bandwidth overhead – one control packet per data packet
  - If all pkts are minimum-sized, can consume half of total bandwidth!

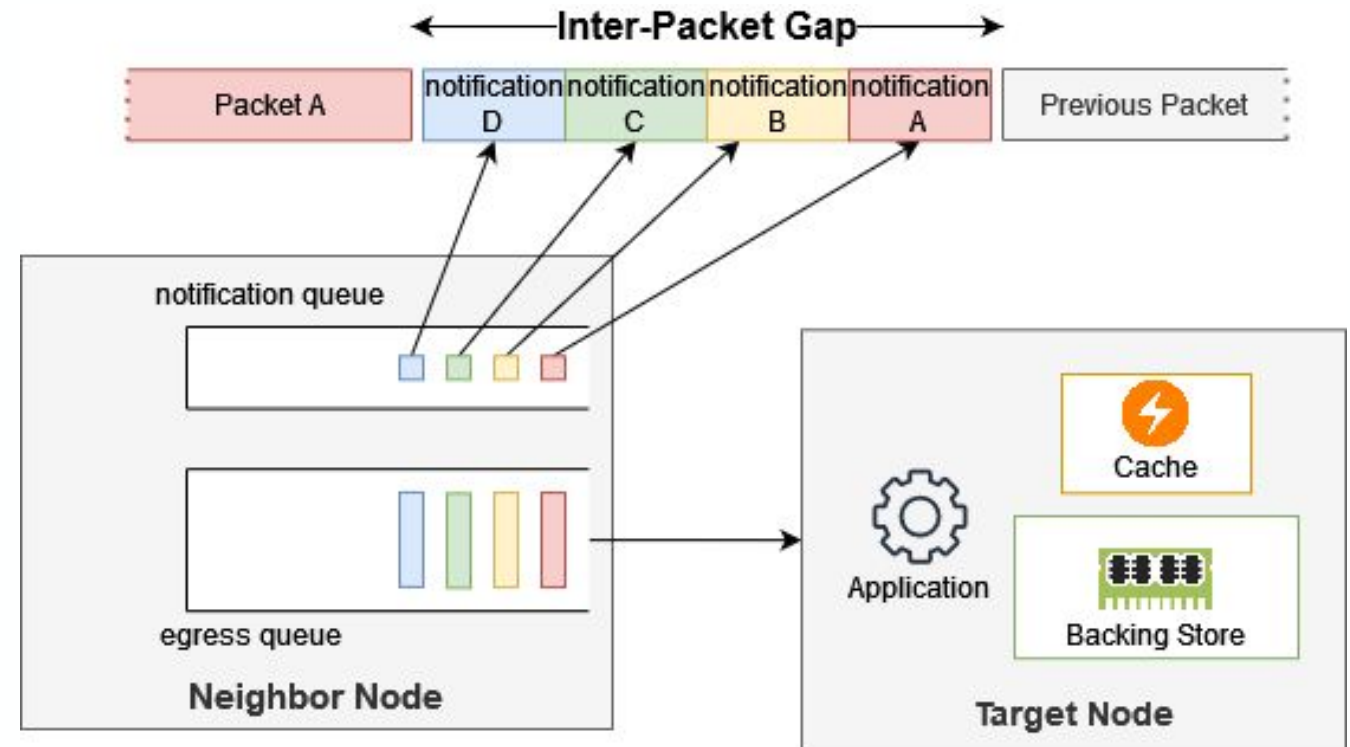Batching reduces number of control packets but can delay notification while it waits for batch



**How to send future state notifications in a timely manner and with low overhead?**

# (1) Low-Overhead *Future State Access Notification*
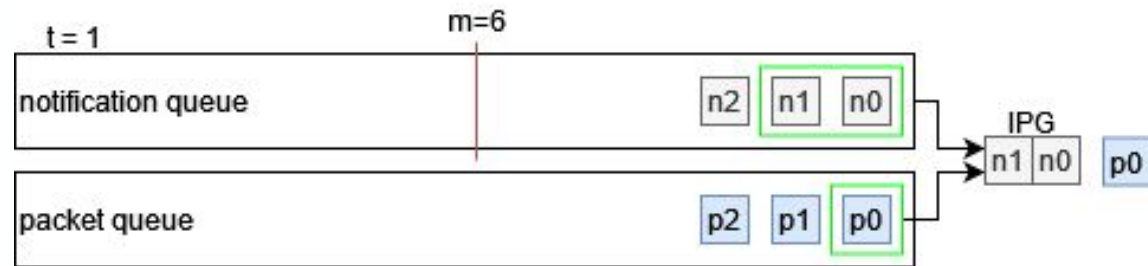
- **Send notifications in IPG**
  - Ethernet PHY enforces a minimum of 96 bit **inter-packet gap (IPG)** between packets
  - Can carry multiple packets' state access notification within a single IPG
  - **Zero bandwidth overhead**
  - **Limitation**: Limited # of bits for communication
    - Limits rate at which packet notifications can be sent

Send notifications using IPG between packets

# (1) Optimization: *Opportunistic* Batching

- Send notifications over IPG under normal scenarios
- Send a control packet when notification queue exceeds <span style="color:darkred">configurable parameter *m*</span>

# *(1)  Optimization: Opportunistic Batching*

- Send notifications over IPG under normal scenarios
- Send a control packet when notification queue exceeds configurable parameter **m**
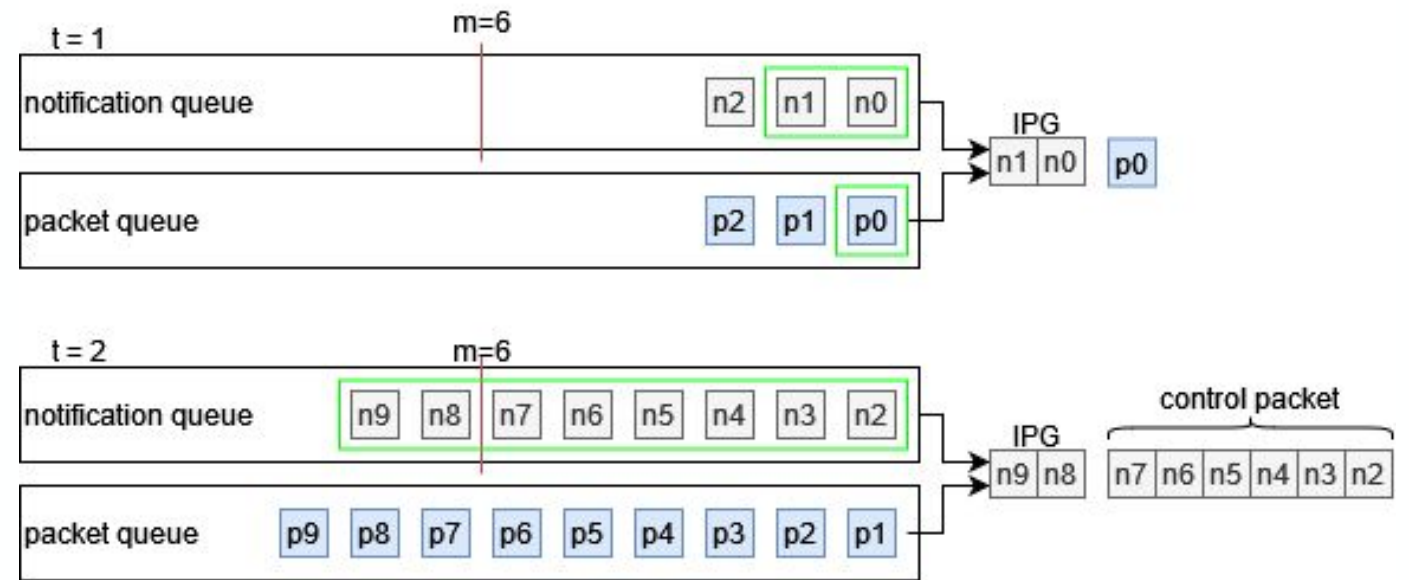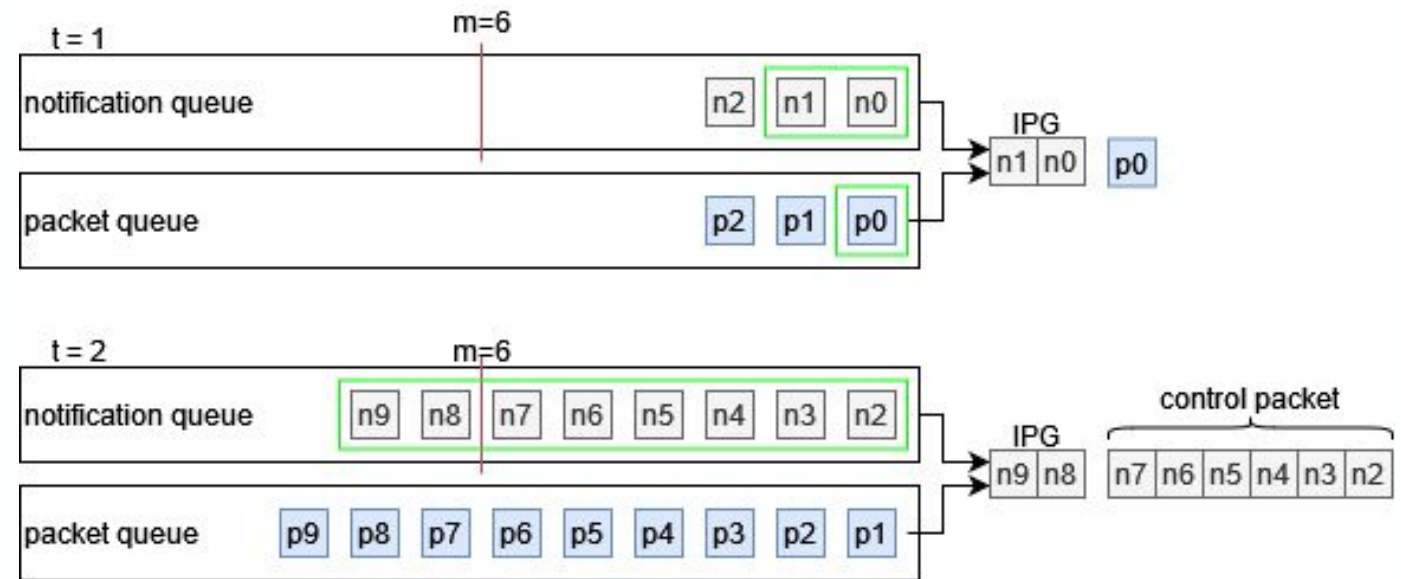  - By configuring *m*, we control bandwidth overhead of control pkts

# (1) Optimization: *Opportunistic* Batching

- Send notifications over IPG under normal scenarios
- Send a control packet when notification queue exceeds configurable parameter *m*
  - By configuring *m*, we control bandwidth overhead of control pkts
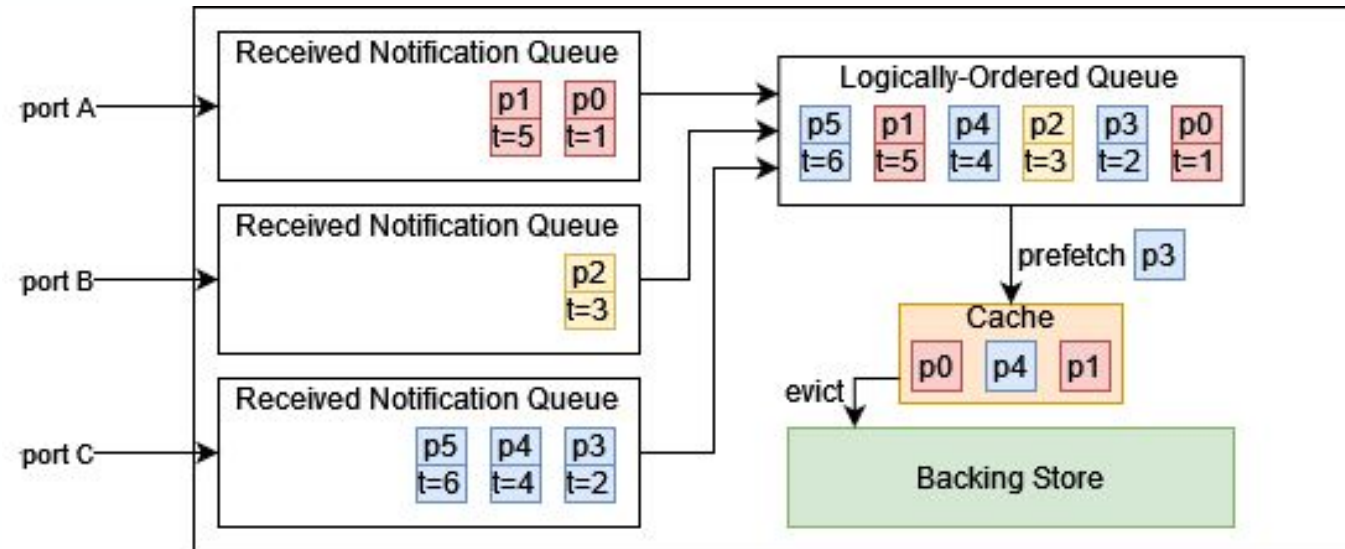


**Best of both worlds:**
Timely notification with low bandwidth overhead and no batching delay

# *(2) Future-Aware Cache Manager*

- Cache manager uses received future state access notifications to make smarter **prefetching** and cache **eviction** decisions

- Cache manager consists of two components:
  - Future-Aware **Prefetching**
  - Future-Aware **Cache Eviction**

# (2) Future-Aware *Prefetching*

- **Goal:** Fetch state in order of **predicted time of access**
  - One received notification queue per input port
  - Combine into **one logically sorted queue** based on future access time
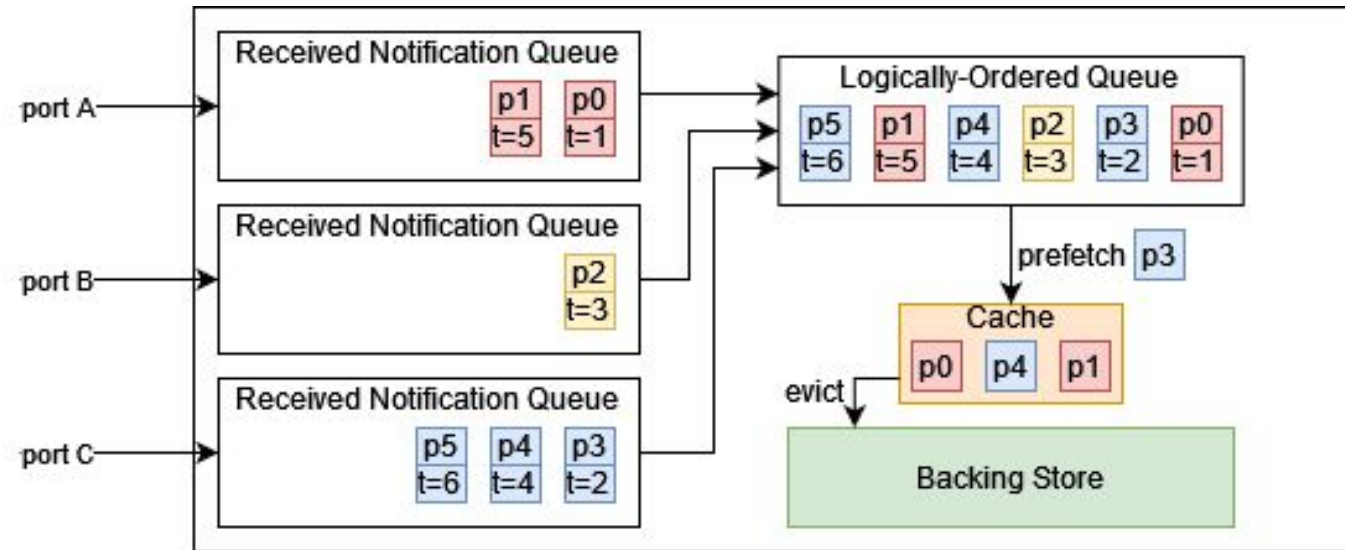  - Fetch soonest state not in cache

# (2) *Future-Aware Prefetching*

- **Goal:** Fetch state in order of **predicted time of access**
  - One received notification queue per input port
  - Combine into **one logically sorted queue** based on future access time
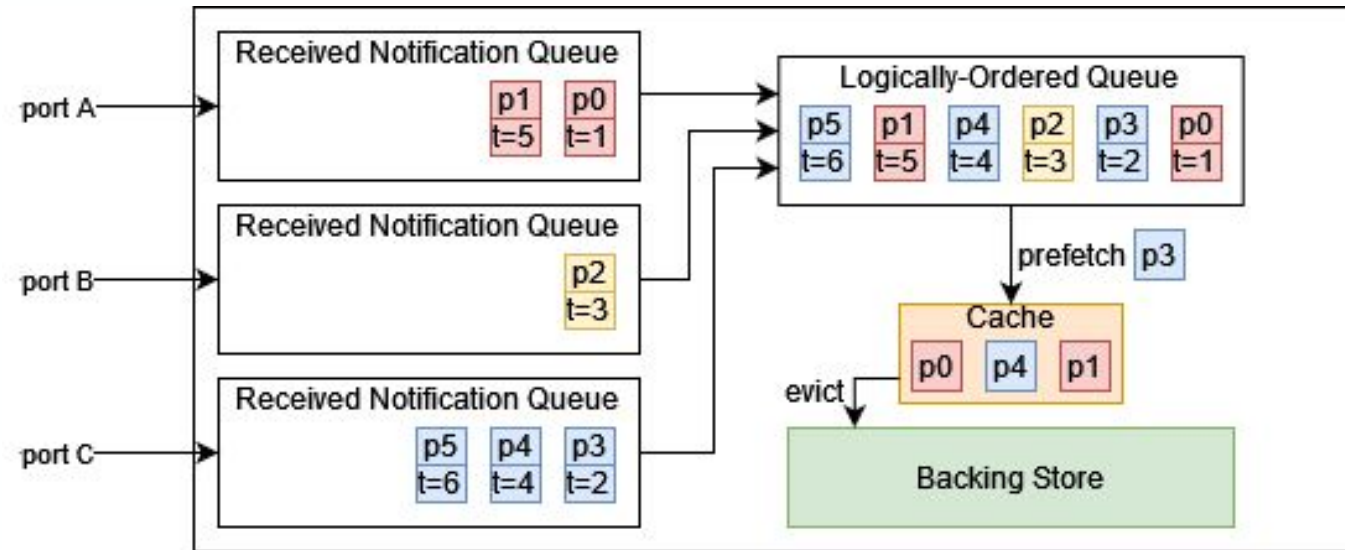  - Fetch soonest state not in cache



(3) **Fast Hardware Implementation**

$$t_{prefetch} = 1 + \log(P) + k \text{ clock cycles}$$

P: number of ports

k: cache set size

# (2) *Future-Aware Prefetching*

- **Goal:** Fetch state in order of **predicted time of access**
  - One received notification queue per input port
  - Combine into **one logically sorted queue** based on future access time
  - Fetch soonest state not in cache
  - If cache is full → **eviction algorithm**



(3) **Fast** *Hardware Implementation*

$$t_{prefetch} = 1 + \log(P) + k \text{ clock cycles}$$

P: number of ports

k: cache set size

# (2) Future-Aware *Cache Eviction*

- **Goal:** Emulate **Belady's algorithm** as closely as possible:
  - Evict an entry that will be accessed furthest in the future

- **Challenge:** Knowledge of only a **partial set** of future state accesses
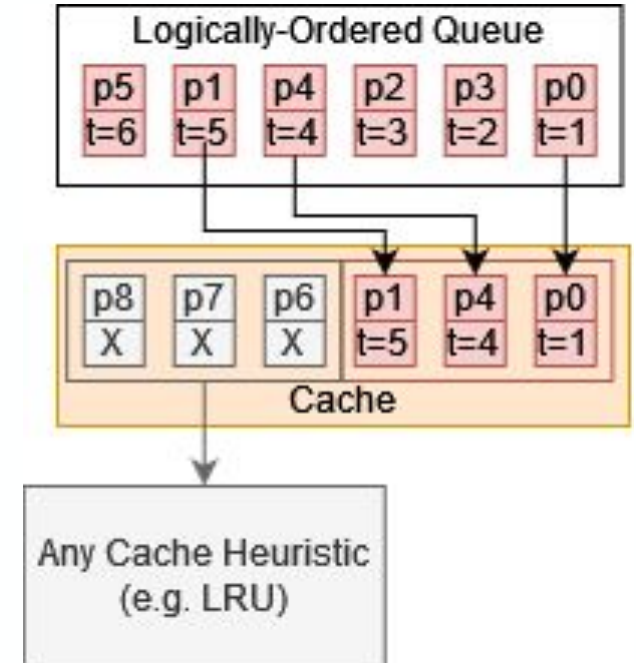
# *(2) Future-Aware Cache Eviction*

- **Our solution:**
  - Split cache into two sets: objects with known access time vs unknown access time
  - Prioritize evicting objects with unknown access time using **any** cache heuristic



Logically-Ordered Queue

| p5 t=6 | p1 t=5 | p4 t=4 | p2 t=3 | p3 t=2 | p0 t=1 |

| p8 X | p7 X | p6 X | p1 t=5 | p4 t=4 | p0 t=1 |

Cache

Any Cache Heuristic (e.g. LRU)

- **Bounded performance**
  - Worst case: **Caching heuristic**

# (2) Future-Aware *Cache Eviction*

- **Our solution:**
  - Split cache into two sets: objects with known access time vs unknown access time
  - Prioritize evicting objects with unknown access time using **any** cache heuristic
  - When cache solely contains objects with known access time, evict according to Belady's algorithm

- **Bounded performance**
  - Worst case: **Caching heuristic**
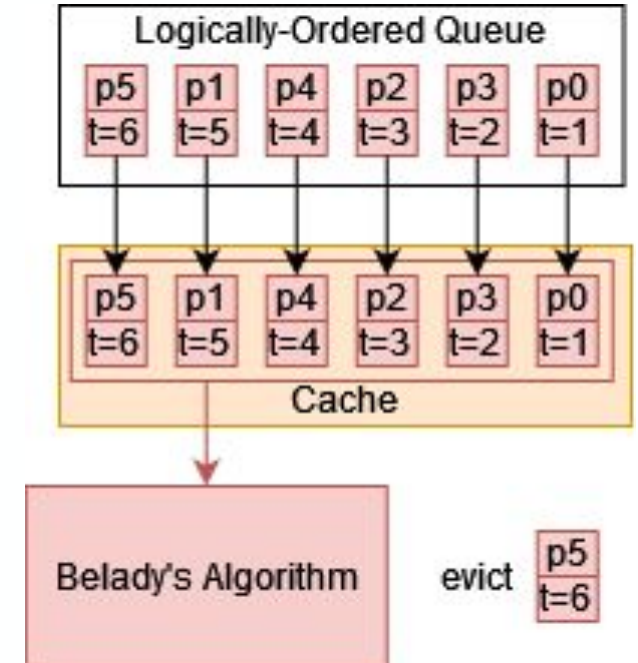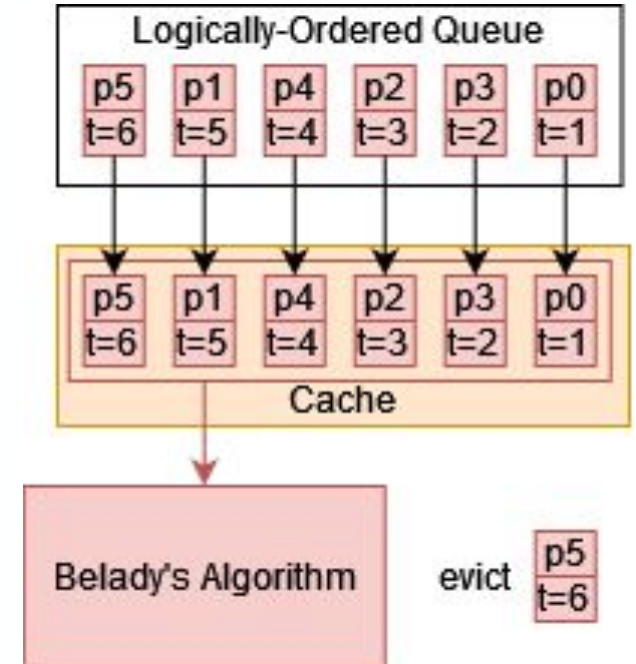  - Best case: **Belady's algorithm**

# *(2) Future-Aware Cache Eviction*

- **Our solution:**
  - Split cache into two sets: objects with known access time vs unknown access time
  - Prioritize evicting objects with unknown access time using ***any*** cache heuristic
  - When cache solely contains objects with known access time, evict according to Belady's algorithm

- **Bounded performance**
  - Worst case: **Caching heuristic**
  - Best case: **Belady's algorithm**



*(3)* ***Fast Hardware Implementation***

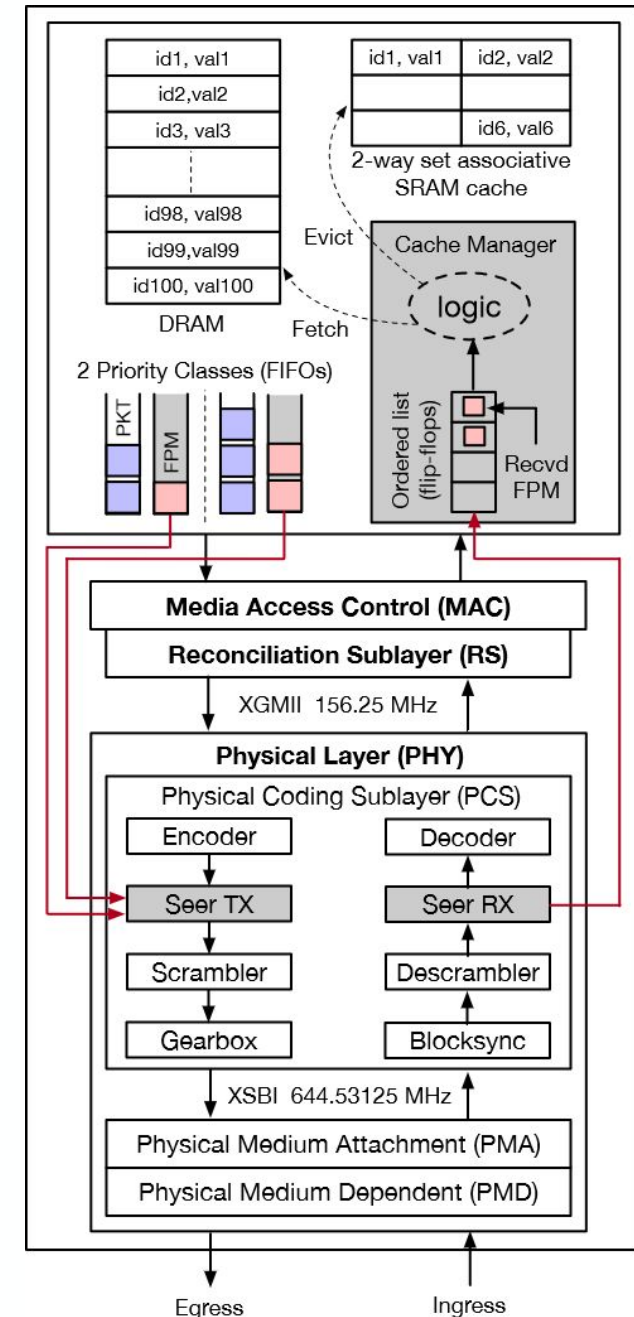$$t_{evict} = k \text{ clock cycles}$$
k: cache set size
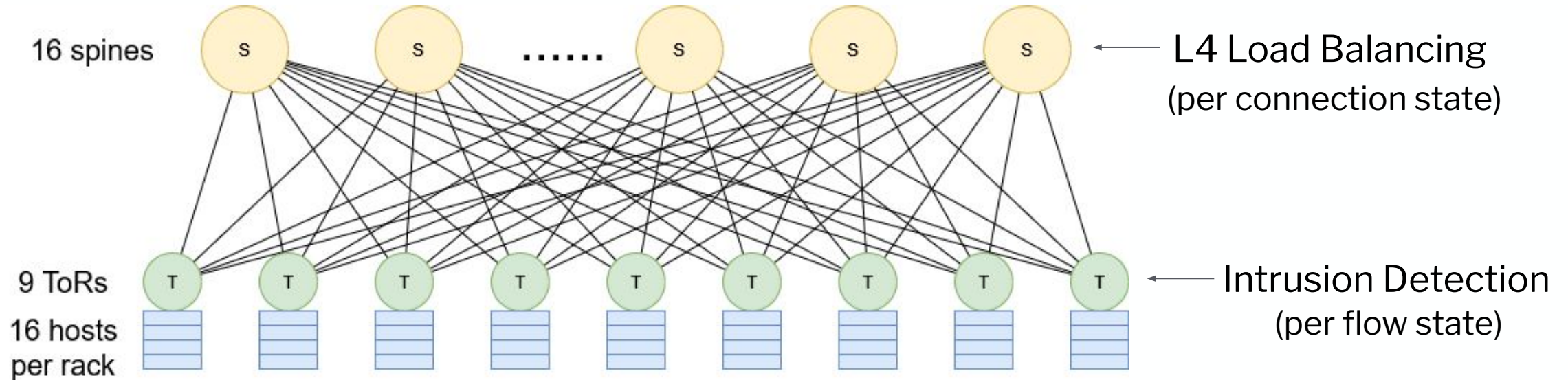
# *Prototype*

- **FPGA prototype**
  - Altera Stratix V FPGA: 234 K adaptive logic modules, 52 Mbits SRAM, four 10 Gbps network ports

- **Seer modifies Ethernet physical layer (PHY) to access IPG**
  - Replaces default idle 0 values in IPG with state access notification

# *Evaluation Setup*

**Applications:**



16 spines ...... S S S S S → L4 Load Balancing (per connection state)

9 ToRs — 16 hosts per rack → Intrusion Detection (per flow state)

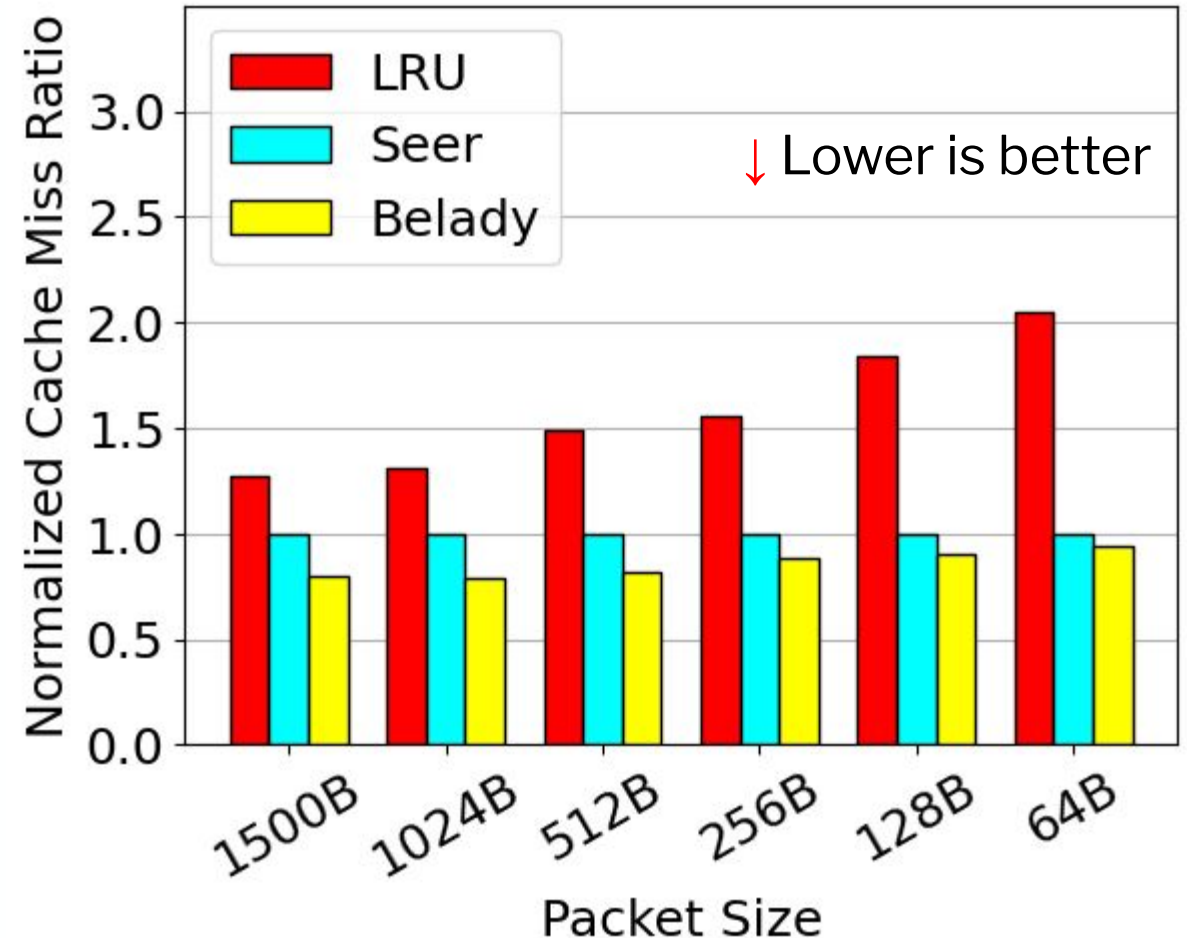| Packet-level simulator in C | | |
|---|---|---|
| • Two-tier Fattree topology <br>     ○ 16 spine switches <br>     ○ 9 racks <br>     ○ 16 hosts / rack (total 144) <br>     ○ Full bisection bandwidth | • 100 Gbps links <br> • 100 ns per-hop propagation delay <br> • 100 ns backing memory access latency <br> • 96 bit inter-packet gap | • DCTCP congestion control <br> • ECMP load balancing <br> • Switches support ECN <br><br> • **Evaluation metric:** <br>     ○ **Cache miss ratio** |

**PURDUE** UNIVERSITY®
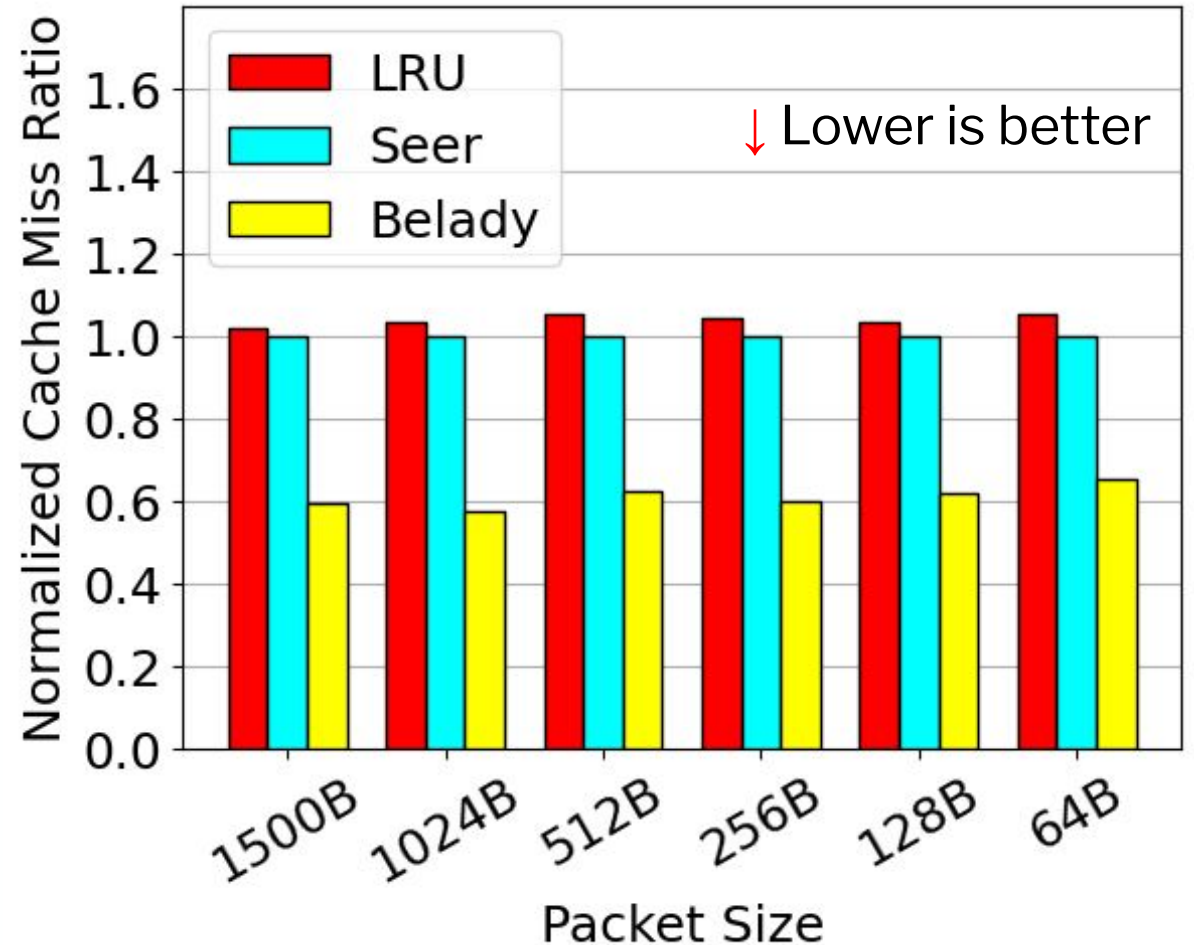
# *Evaluation: Good Case for Seer*

- **Incast** Traffic Pattern:
  - Incast traffic results in most queueing at neighbor node
  - Provides furthest visibility into future state accesses
- Seer remains within 7-20% of Belady
- Seer performs 20-100% better than LRU



Performance for each packet size normalized w.r.t. corresponding Seer performance
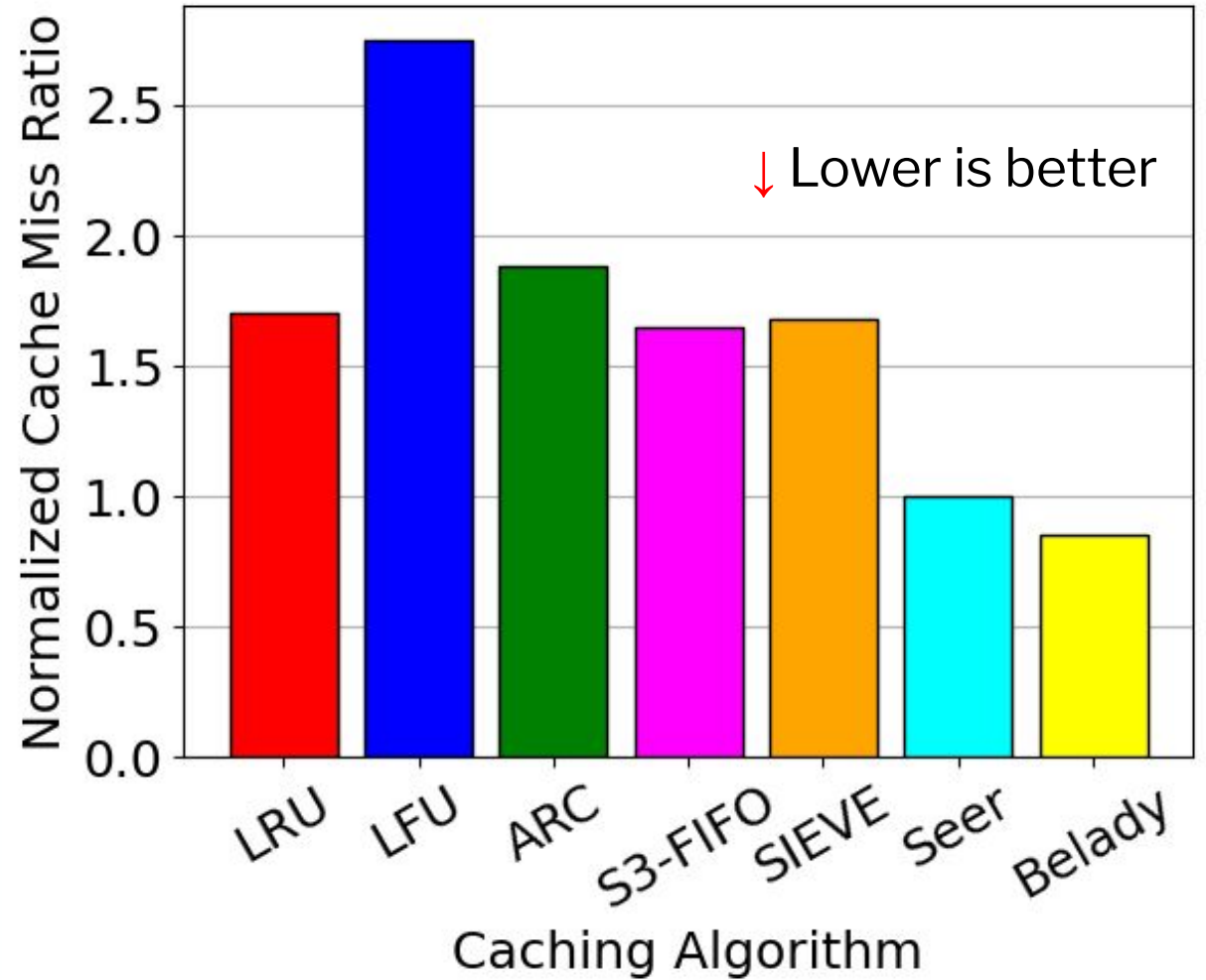
# *Evaluation: Bad Case for Seer*

- **Permutation** Traffic Pattern:
  - Permutation traffic over full bisection bandwidth fattree network results in least queueing at neighbor node
  - Provides least visibility into future state accesses
- Seer remains within 35-40% of Belady
- Seer performs 2-5% better than LRU

Performance for each packet size normalized w.r.t. corresponding Seer performance

# *Evaluation: Realistic Workload*

- **Websearch** workload
  - Representative of datacenter workload
  - Heavy-tailed flow size distribution

- 60-180% lower cache miss ratio for Seer compared to state-of-art

- **Flow completion time (FCT) show similar trend:**
  - Seer reduces FCT by 25-75% compared to LRU



↓ Lower is better

Normalized w.r.t. Seer

# *Conclusion*

- Seer enables future-aware online caching in a networked system

- Seer makes three key technical contributions:
  - Low-Overhead Protocol for Future State Access Notification
  - Design of Future-Aware Cache Manager
  - Fast Hardware Implementation

- Seer performs close to optimal offline caching in practice, with worst case performance bounded by state-of-the-art caching heuristic

PURDUE
UNIVERSITY®

# Thank You

Any questions?

PURDUE UNIVERSITY

nsdi '24