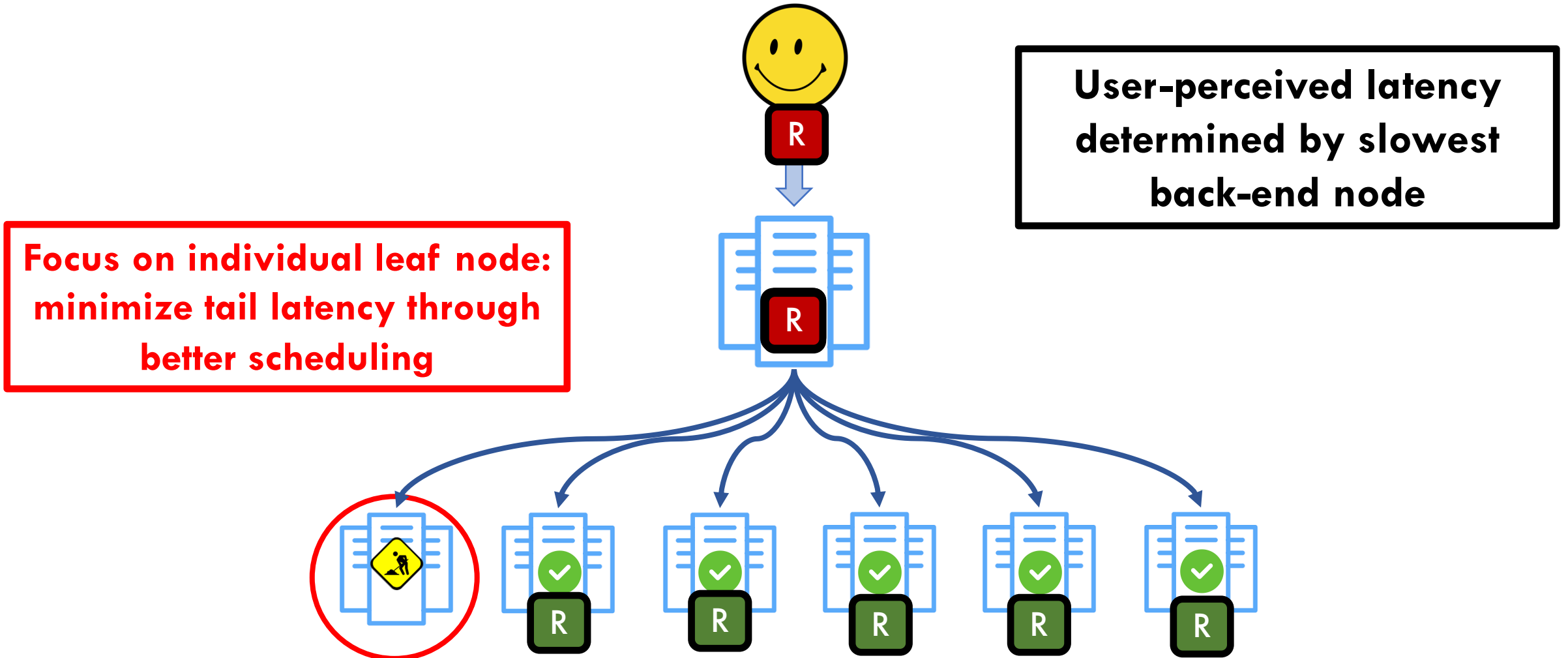# Shinjuku: Preemptive Scheduling for Microsecond-Scale Tail Latency

**Kostis Kaffes**, Timothy Chong, Jack Tigar Humphries,

Adam Belay, David Mazières, Christos Kozyrakis
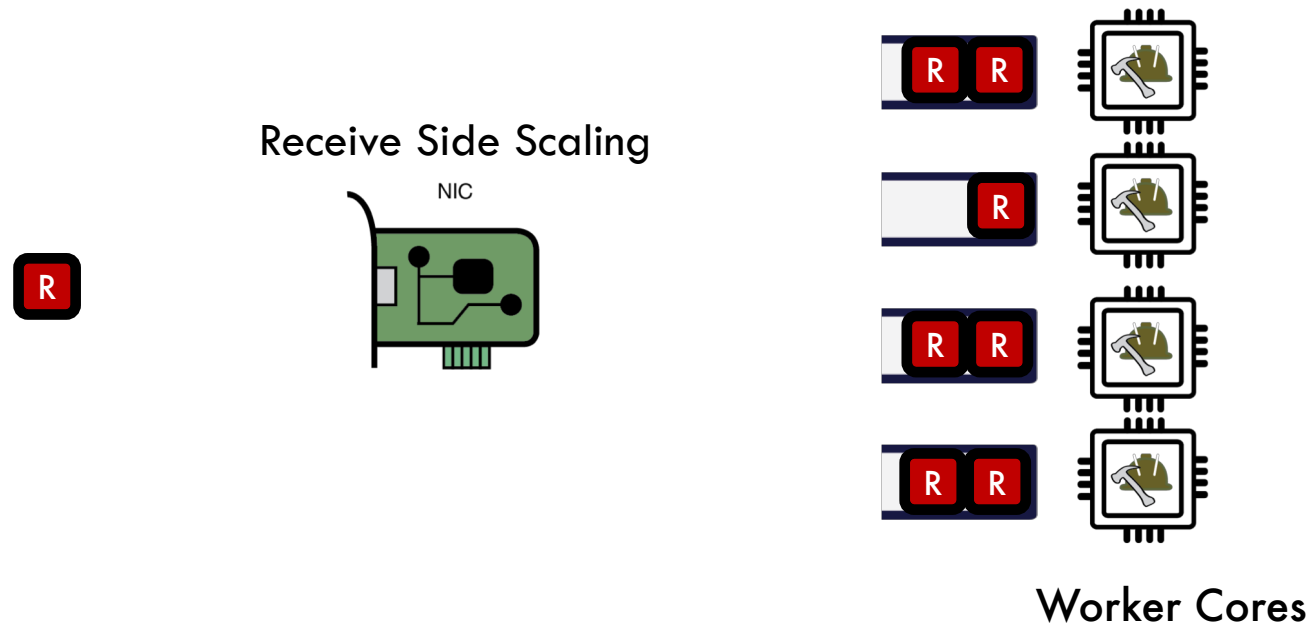
# Tail latency matters for datacenter workloads



User-perceived latency determined by slowest back-end node

Focus on individual leaf node: minimize tail latency through better scheduling

# Achieving low tail latency at microsecond scale is hard

**Problem**: High OS overheads

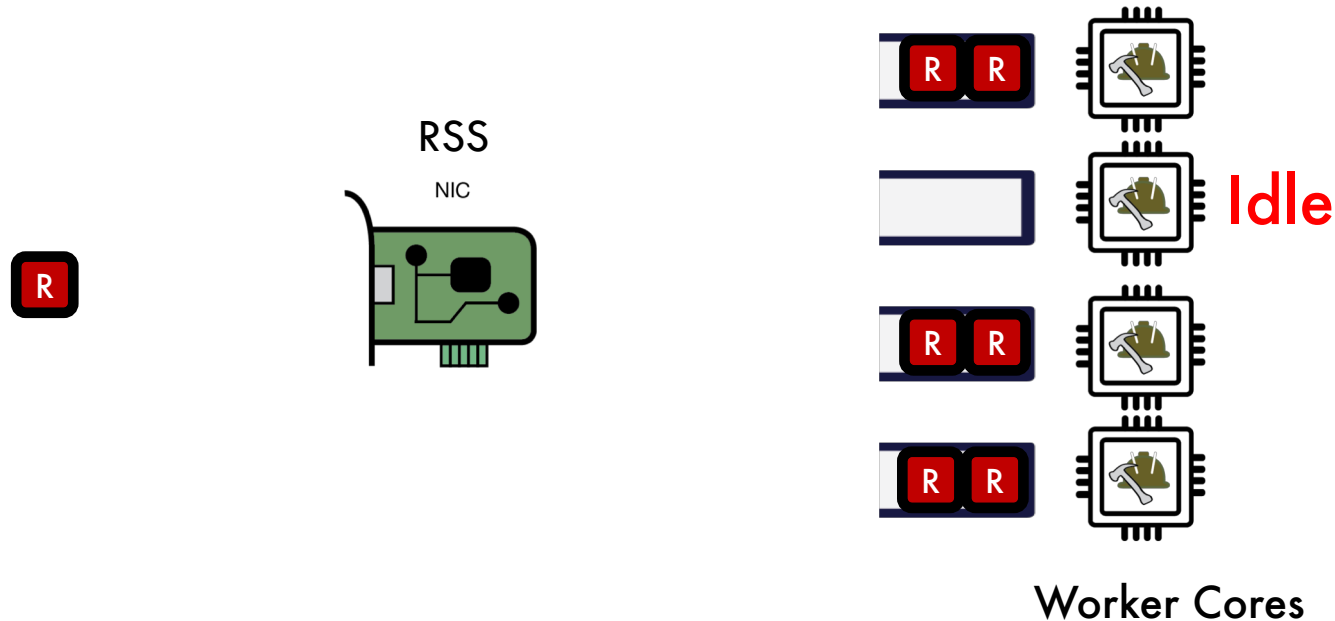**Solution**: OS Bypass, polling (no interrupts), run-to-completion (no scheduling)

**D**istributed Queues + **F**irst **C**ome **F**irst **S**erve scheduling

**d-FCFS** (DPDK, IX, Arrakis)

Receive Side Scaling

NIC

Worker Cores

# Achieving low tail latency at microsecond scale is hard

**Problem**: Queue imbalance because d-FCFS is not work conserving
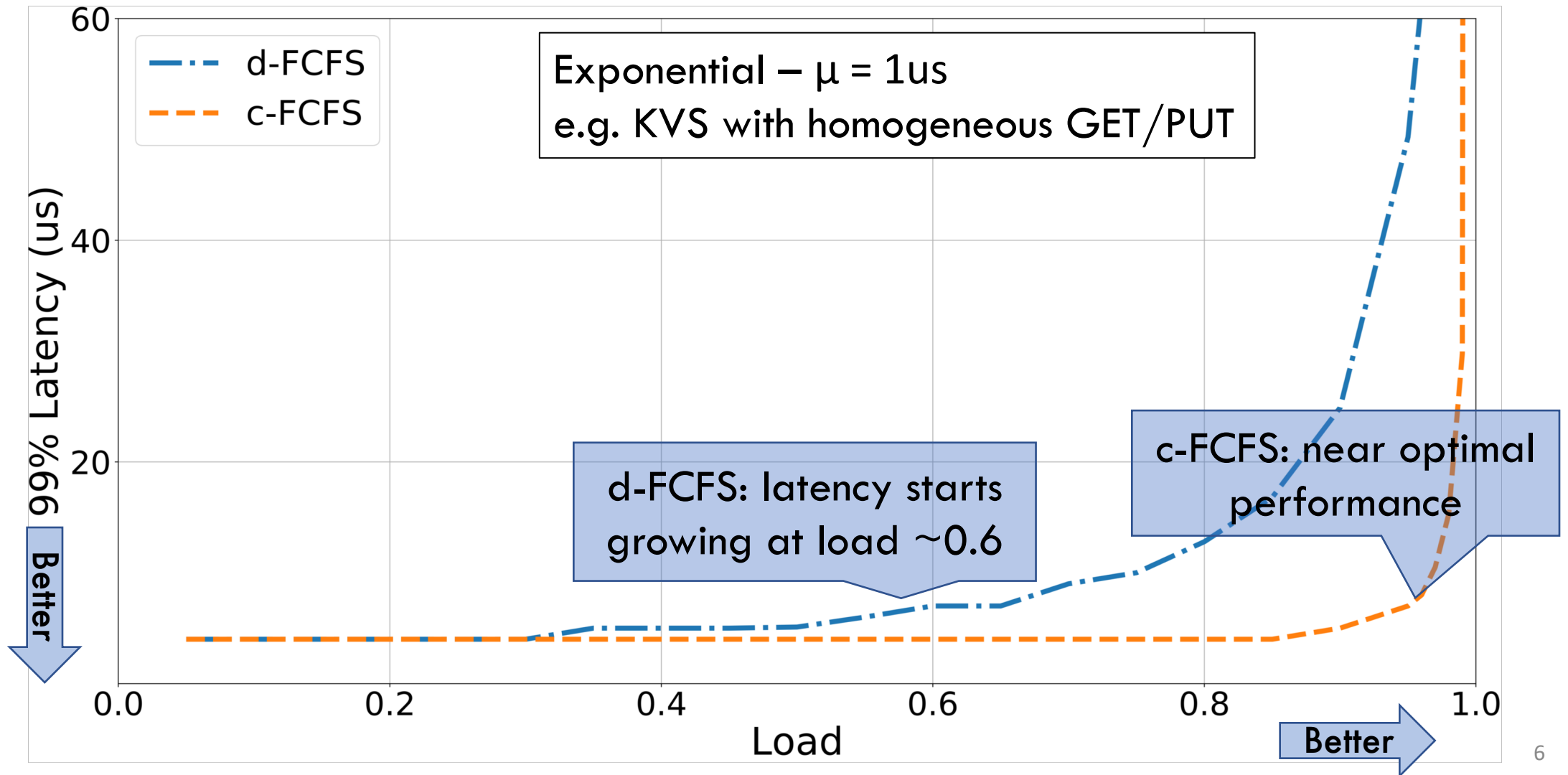
RSS

NIC

**Idle**

Worker Cores

# Achieving low tail latency at microsecond scale is hard

**Problem**: Queue imbalance because d-FCFS is not work conserving
**Solution**: Centralized queue - **c-FCFS**



NIC

**Approximation: d-FCFS + stealing e.g., ZygOS**

Worker Cores

# Ideal centralized queue is better in simulation



Exponential – μ = 1us
e.g. KVS with homogeneous GET/PUT

d-FCFS: latency starts growing at load ~0.6

c-FCFS: near optimal performance

Better

Better

# Is FCFS good enough when task duration varies?



Bimodal – 99.5% 0.5us – 0.5% 500us
e.g. KVS with some RANGE queries

c-FCFS: latency increases even for low load

# Problem: Short requests get stuck behind long ones



NIC

All cores are hogged by long requests

# What if we could use the same preemptive scheduling as Linux?



Bimodal – 99.5% 0.5us – 0.5% 500us
e.g. KVS with some RANGE queries

Legend:
- d-FCFS
- c-FCFS
- PS - 1ms

Y-axis: 99% Latency (us)
X-axis: Load

PS–1ms: latency increases even for low load (same as c-FCFS)

Better (down arrow)
Better (right arrow)

# Solution: What if we could use preemptive scheduling but at usec scale?



Bimodal – 99.5% 0.5us – 0.5% 500us
e.g. KVS with some RANGE queries

99% Latency (us)

Better

d-FCFS
c-FCFS
PS - 5us
PS - 1ms

PS-5us: near optimal performance with fast preemption

Load

Better

# Insights

Effective scheduling for tail latency requires:

- Centralized queue
- Preemption
- Scheduling policies tailored for each workload

Problem: Microsecond scale requires

- Millions of queue accesses per second
- Preemption as often as every 5us
- Light-weight scheduling policies

# Solution: Shinjuku

A single address-space operating system that achieves microsecond-scale tail latency for all types of workloads regardless of variability in task duration
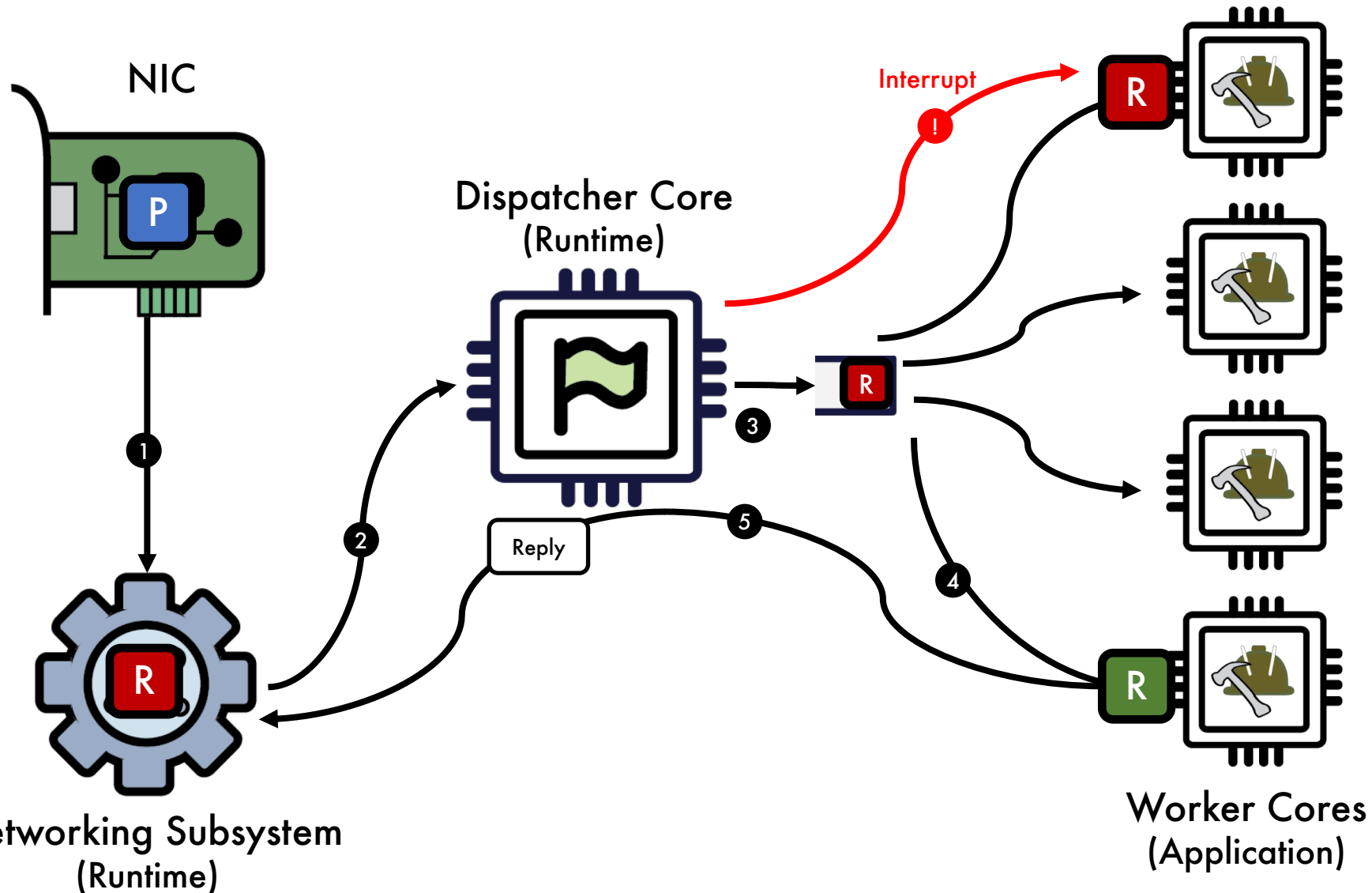
Key Features:

- Dedicated core for scheduling and queue management
- Leverage hardware support for virtualization for fast preemption
- Very fast context switching in user space
- Match scheduling policy to task distribution and target latency

# Outline

- <span style="color:red">Shinjuku Design</span>

- Preemption Mechanisms

- Scheduling Policies

- Evaluation

# Shinjuku Design



NIC

Dispatcher Core
(Runtime)

Interrupt

Reply

Networking Subsystem
(Runtime)

Worker Cores
(Application)

1 Process packets and generate application-level requests

2 Pass requests to centralized dispatcher using shared memory

3 Add requests to centralized queue

4 Schedule requests to worker cores using shared memory

5 Send replies back to clients through the networking subsystem

! Interrupt long running requests and schedule other requests from the queue

14

# Outline

- Shinjuku Design

- <span style="color:red">Preemption Mechanisms</span>

- Scheduling Policies

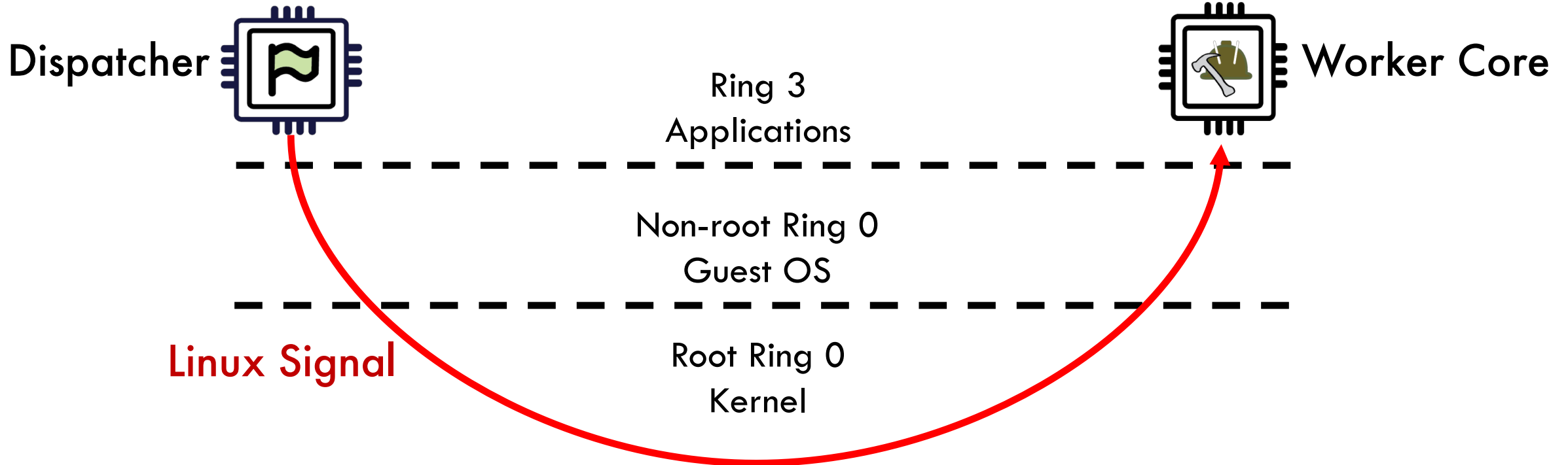- Evaluation

# Minimizing Preemption Overhead

**Sender Overhead**
2084 cycles

**Receiver Overhead**
2523 cycles

Linux Signal

Dispatcher

Worker Core

Ring 3
Applications

Non-root Ring 0
Guest OS

Linux Signal

Root Ring 0
Kernel

# Minimizing Preemption Overhead

|  | Sender Overhead | Receiver Overhead |
|---|---|---|
| Linux Signal | 2084 cycles | 2523 cycles |
| Hardware Interrupts | 2081 cycles | 2662 cycles |



Dispatcher

Worker Core

Ring 3
Applications

Non-root Ring 0
Guest OS

LOCAL APIC

LOCAL APIC

Root Ring 0
Kernel

VMExit

VMExit

# Minimizing Preemption Overhead

|  | Sender Overhead | Receiver Overhead |
|---|---|---|
| Linux Signal | 2084 cycles | 2523 cycles |
| Hardware Interrupts | 2081 cycles | 2662 cycles |
| no VMExits | 298 cycles | 1212 cycles |

-85%   -52%

Dispatcher                                          Worker Core

Ring 3
Applications

Non-root Ring 0
Guest OS

LOCAL APIC          LOCAL APIC

Root Ring 0
Kernel

VMExit                                          VMExit

Map APIC to dispatcher's
address space

Posted Interrupts

14

# Minimizing Preemption Overhead

| Sender Overhead | | Receiver Overhead | |
|---|---|---|---|
| 2084 cycles | | 2523 cycles | |
| 2081 cycles | -85% | 2662 cycles | -52% |
| 298 cycles | | 1212 cycles | |

Linux Signal
Hardware Interrupts
no VMExits

Dispatcher

Worker Core

5us Time Slice

LOCAL APIC

Root Ring 0
Kernel

LOCAL APIC

Map APIC to dispatcher's
address space

Posted Interrupts

# Outline

- Shinjuku Design

- Preemption Mechanisms
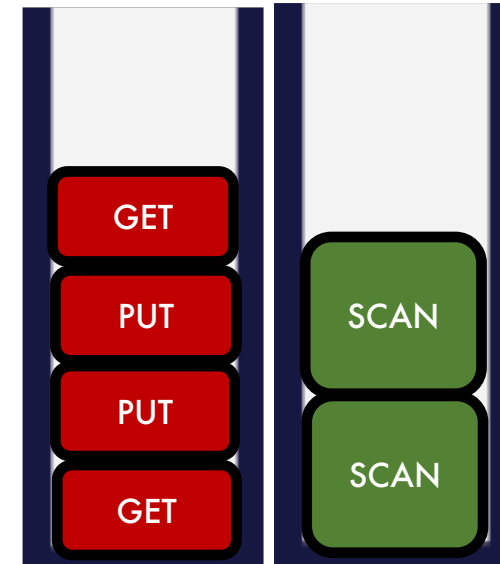
- Scheduling Policies

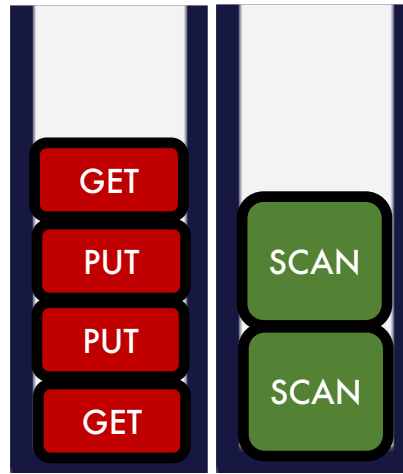- Evaluation

# Scheduling policy

Case 1

Case 2

1) Which queue select from?
2) Where to preempted requests?

GET
PUT
PUT
GET
SCAN

Single Queue (SQ)

GET
PUT
PUT
GET

SCAN
SCAN

Multiple Queues (MQ)

# Queue Selection Policy



Multiple Queues (MQ)

**Policy:** Select the queue with the highest ratio: $\frac{\textit{Waiting Time}}{\textit{Target Latency}}$

**Short requests:** Initially low Target Latency ➜ High Ratio

**Long requests:** Eventually high Waiting Time ➜ High Ratio

# Outline

- Shinjuku Design

- Preemption Mechanisms

- Scheduling Policies

- <span style="color:red">Evaluation</span>

# Evaluation

## Systems

**Shinjuku** – Centralized preemptive scheduling

    14 Logical Cores for workers

    1 Physical Core for both networker and dispatcher (1 Logical Core each)

**IX** – d-FCFS
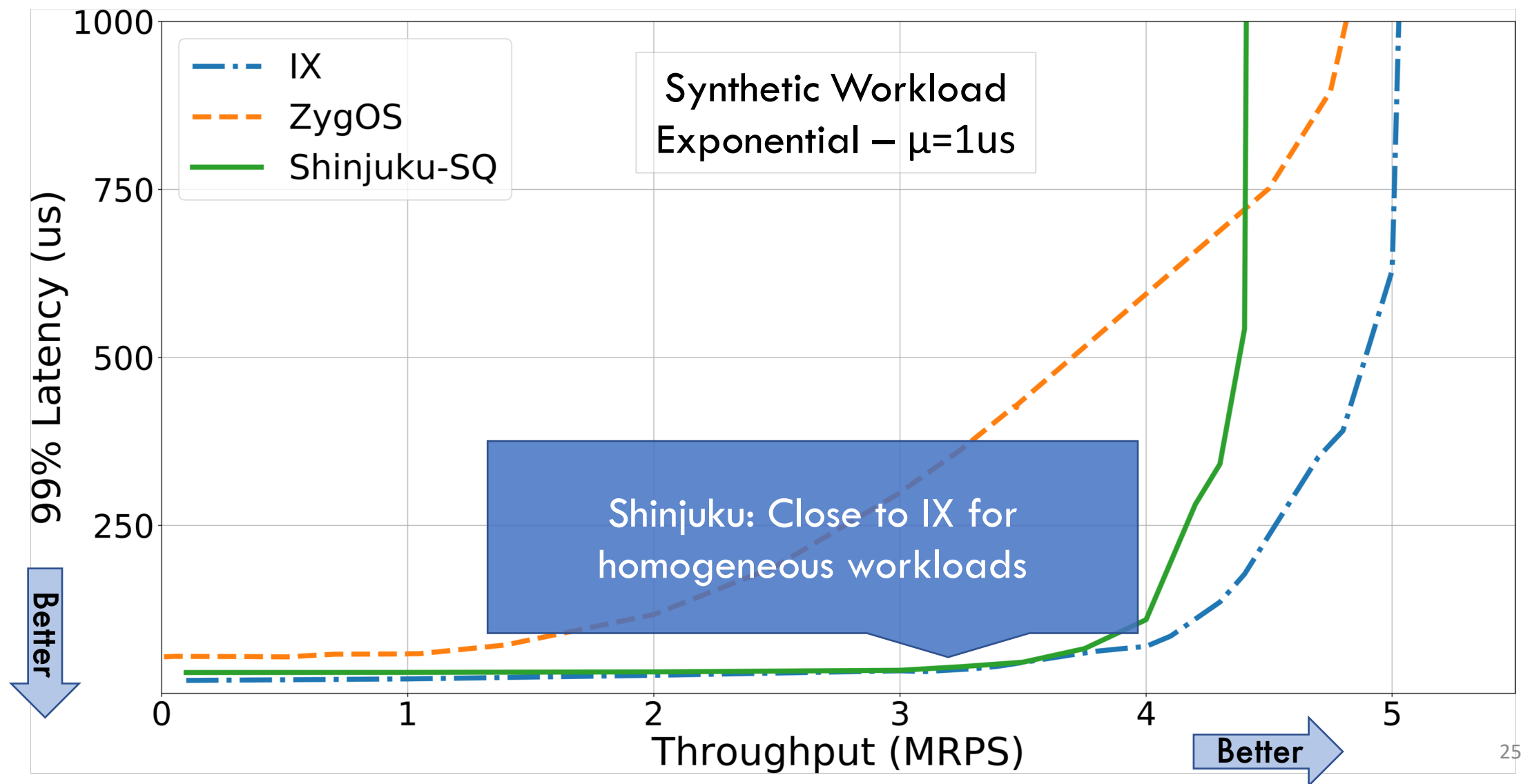
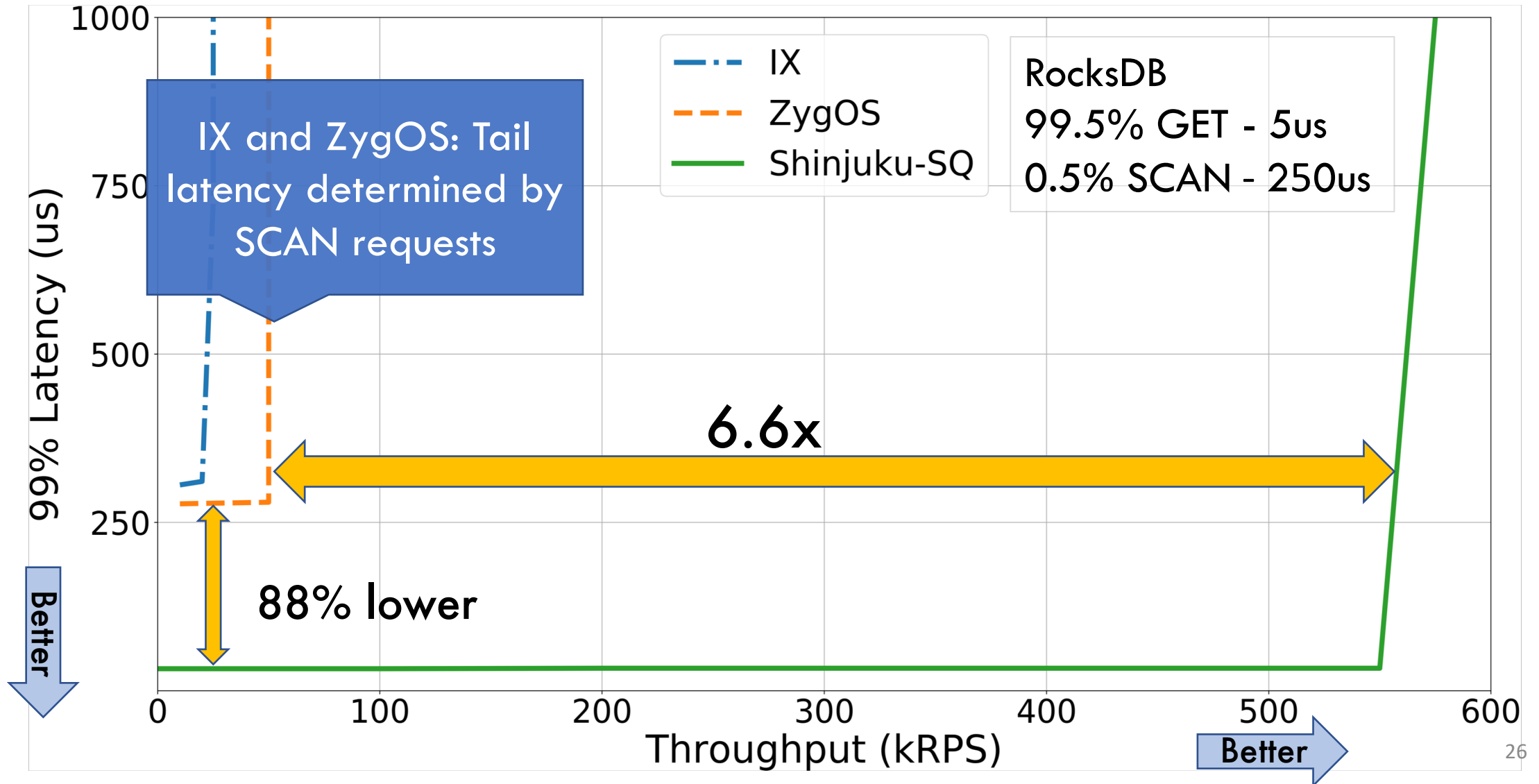**ZygOS** – d-FCFS + work stealing

    16 Logical Cores for workers

## Workloads

Synthetic benchmark with different service time distributions

RocksDB - in-memory database

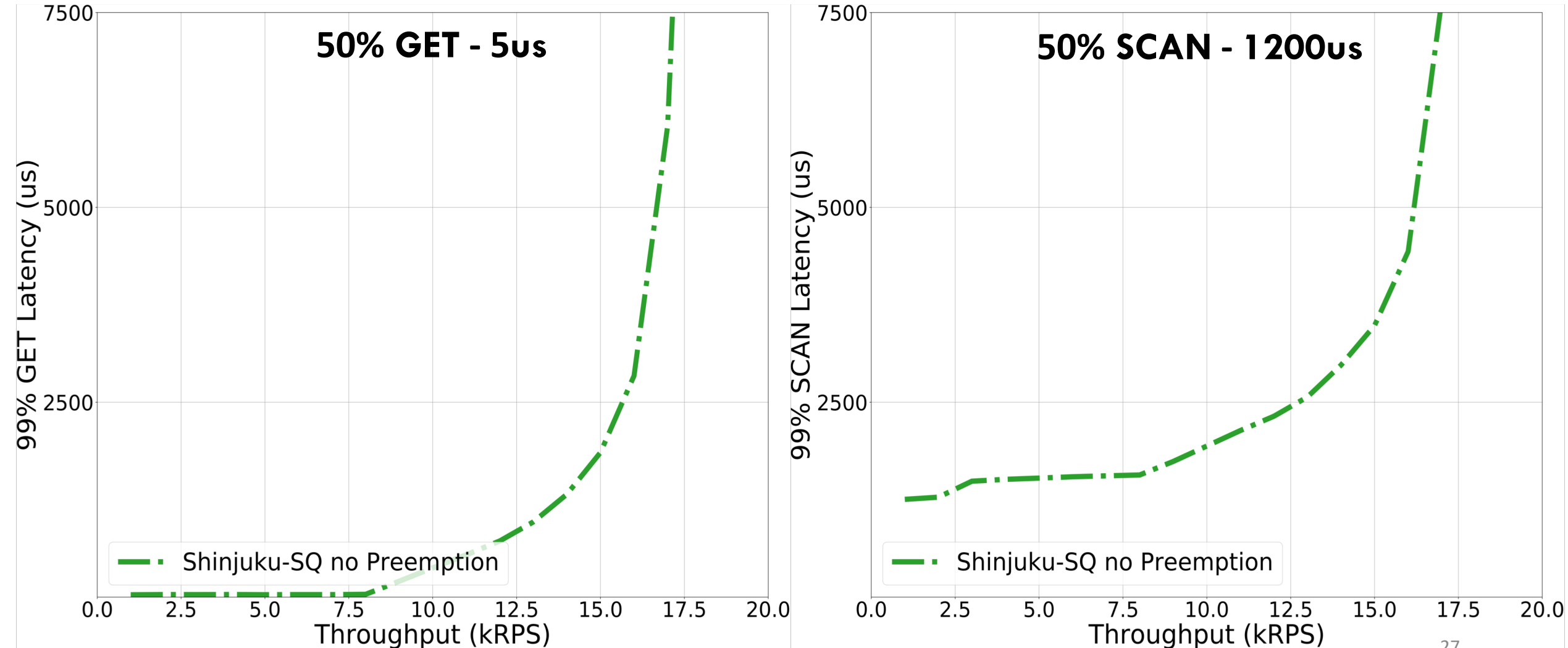# Shinjuku under low variability



99% Latency (us)

1000

750

500

250

- IX
- ZygOS
- Shinjuku-SQ

Synthetic Workload
Exponential – μ=1us

Shinjuku: Close to IX for
homogeneous workloads

Better

0   1   2   3   4   5

Throughput (MRPS)

Better

# Shinjuku under high variability

# How important is each optimization?

## Single Queue no Preemption



**50% GET - 5us**

**50% SCAN - 1200us**

Shinjuku-SQ no Preemption

# How important is each optimization?
## Single Queue with Preemption



**50% GET - 5us**

**50% SCAN - 1200us**

Preemption offers flatter latency for some loss of throughput

Left chart — y-axis: 99% GET Latency (us), x-axis: Throughput (kRPS)
- Shinjuku-SQ no Preemption
- Shinjuku-SQ with Preemption

Right chart — y-axis: 99% SCAN Latency (us), x-axis: Throughput (kRPS)
- Shinjuku-SQ no Preemption
- Shinjuku-SQ with Preemption

# How important is each optimization?
## Multiple Queues with Preemption



**50% GET - 5us**

**50% SCAN - 1200us**

Multi-queue policy recovers the lost throughput

Legend (left and right plots):
- Shinjuku-SQ no Preemption
- Shinjuku-SQ with Preemption
- Shinjuku-MQ

Left plot axes: 99% GET Latency (us) vs Throughput (kRPS)
Right plot axes: 99% SCAN Latency (us) vs Throughput (kRPS)
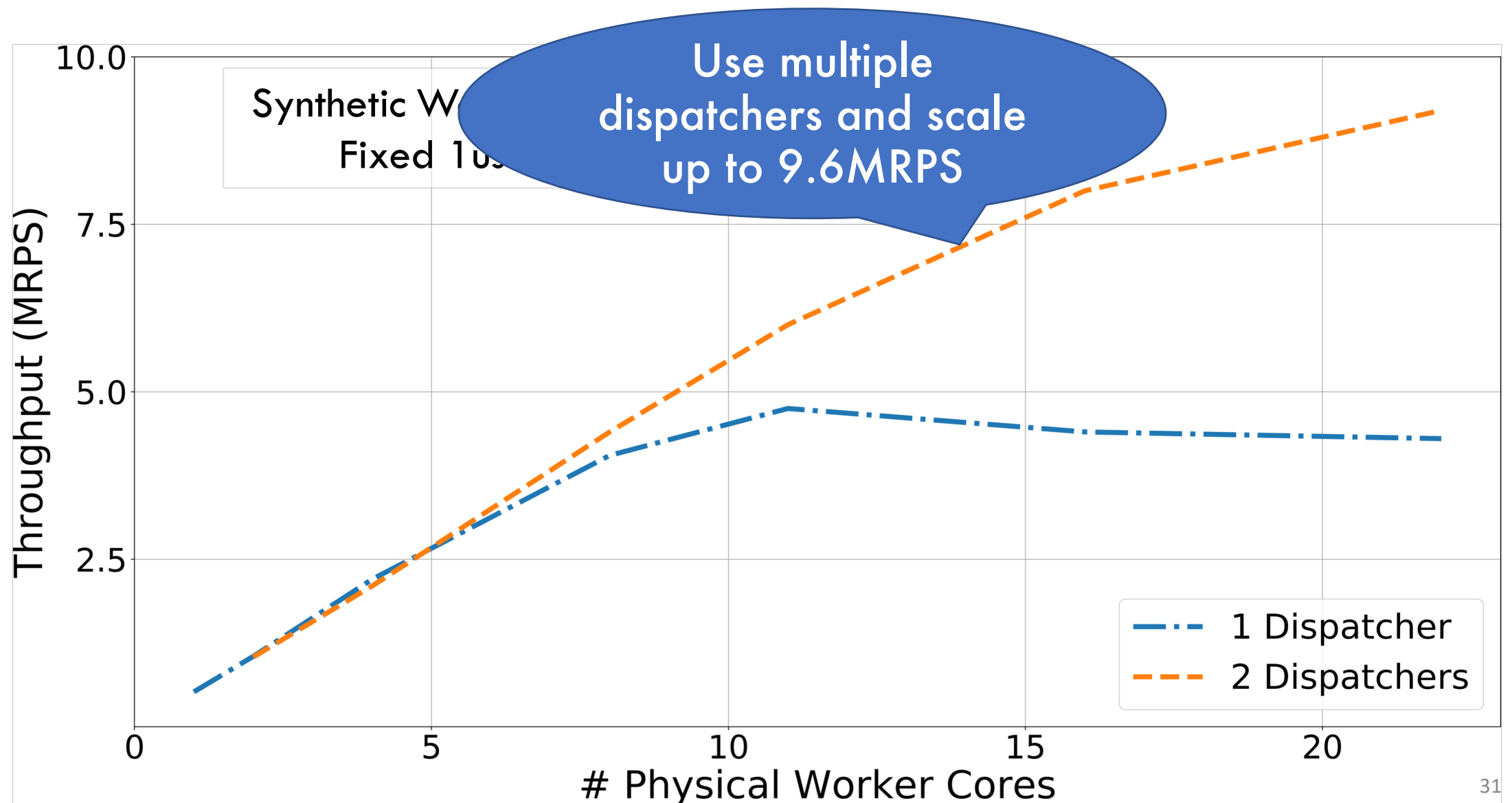
# Does Shinjuku scale?

# Does Shinjuku scale?

# More details in the paper

- Fast context switching

- How Shinjuku supports high line rates

- Placement policy of interrupted requests

- The problems of RSS-only scheduling of requests to cores

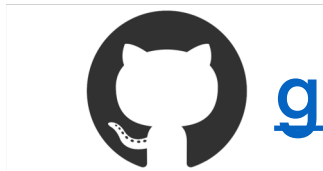- More performance analysis

# Conclusion

Low tail latency for general workloads requires:

- Preemptive Scheduling

- Centralized Queueing

- Flexible Scheduling Policies
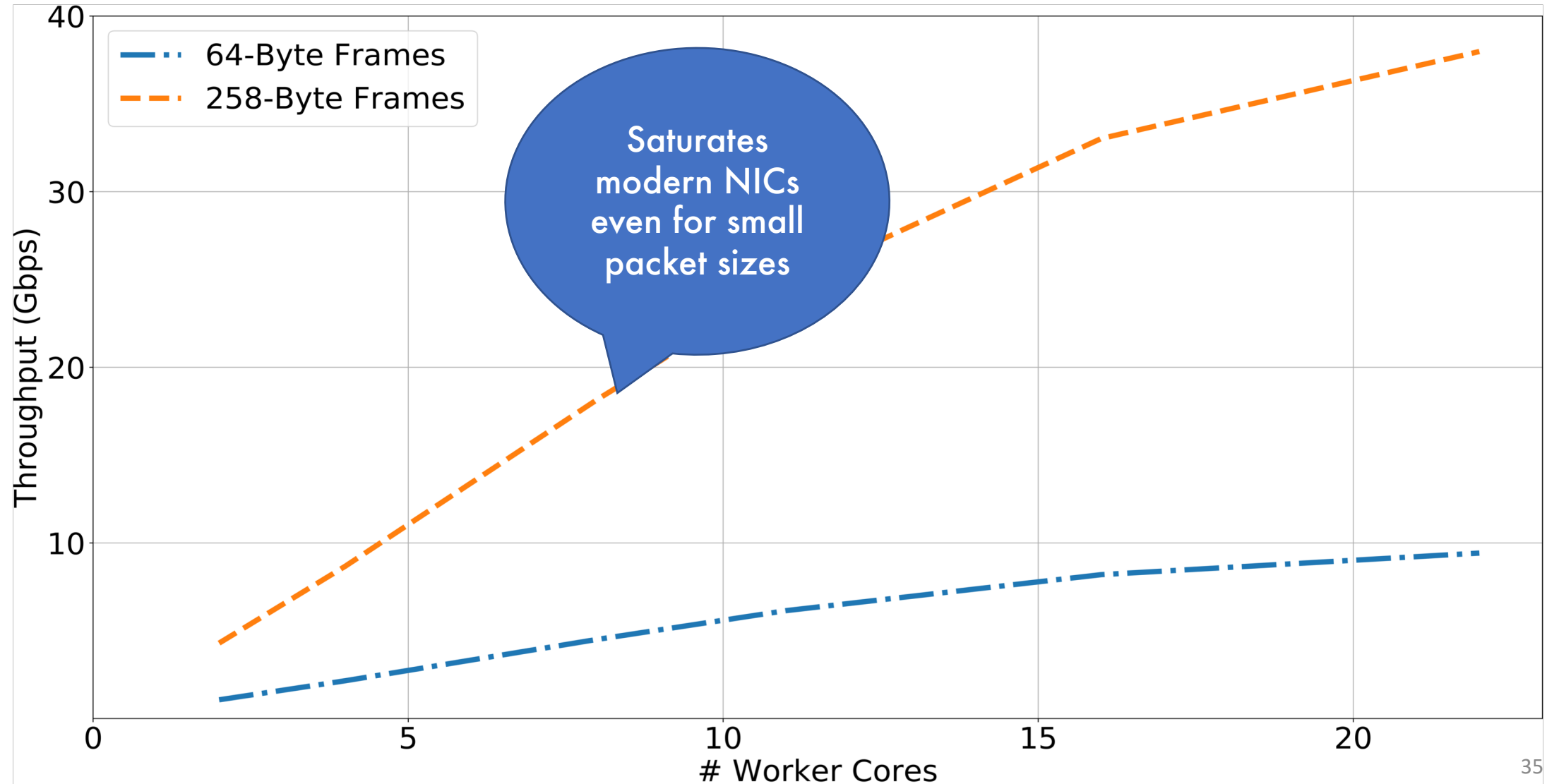
Shinjuku meets these demands at microsecond scale:

- Scalable centralized queue using dedicated core

- Preemption every 5us
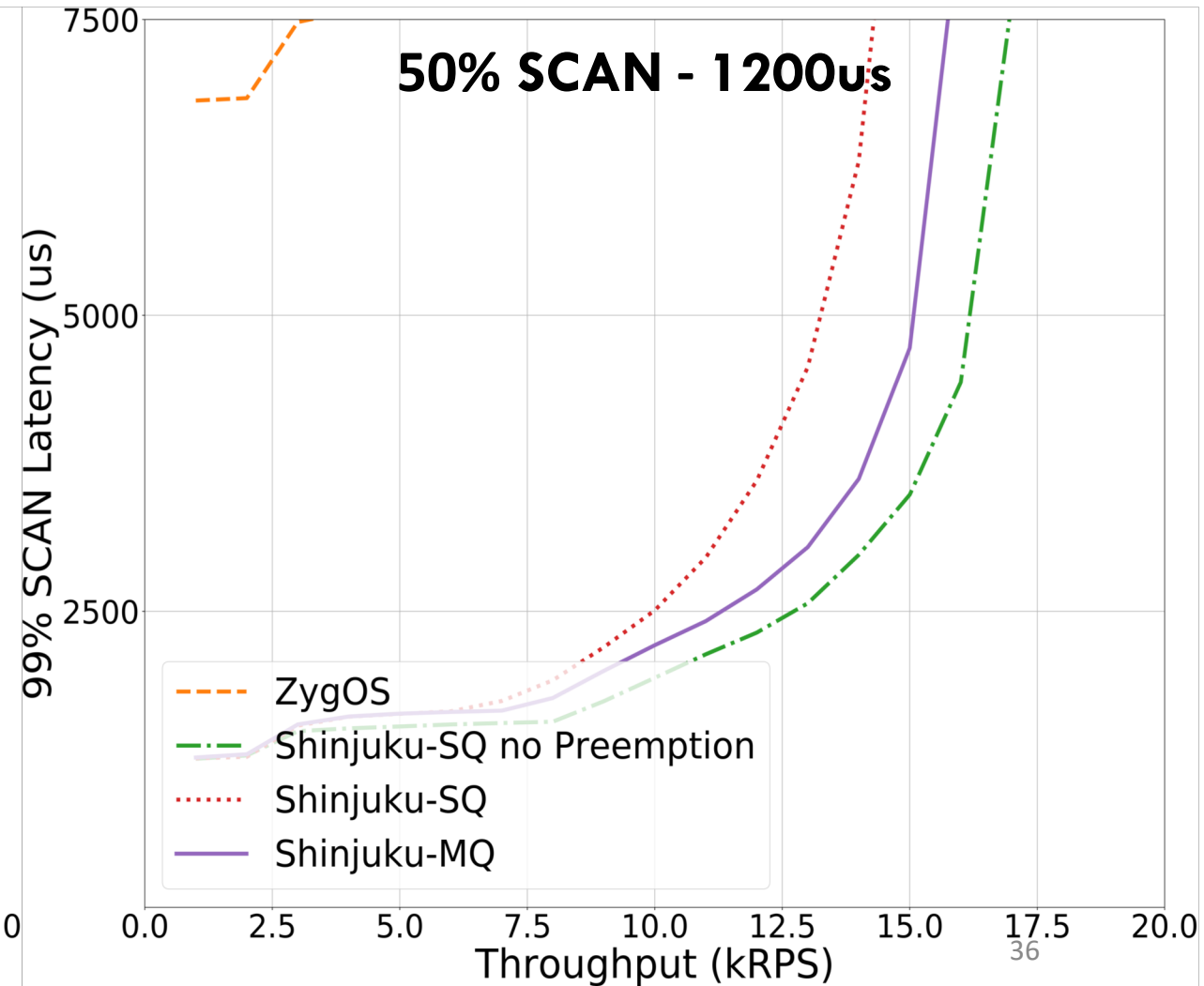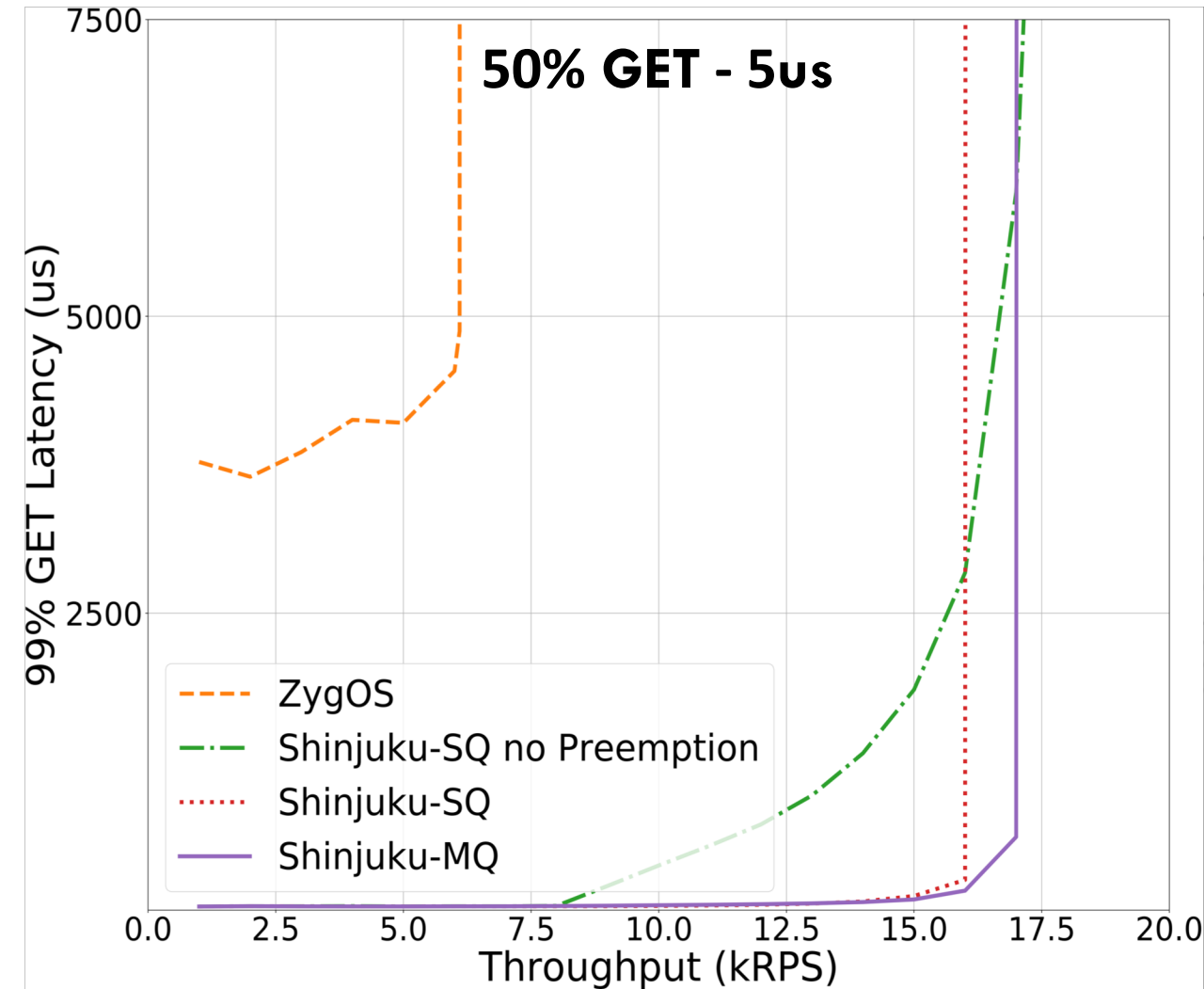
- Latency-driven scheduling policies

[github.com/stanford-mast/shinjuku](github.com/stanford-mast/shinjuku)

# Backup

# Shinjuku Network Scaling

# How important is each optimization?



50% GET - 5us

50% SCAN - 1200us

99% GET Latency (us) vs Throughput (kRPS)

99% SCAN Latency (us) vs Throughput (kRPS)

Legend:
- ZygOS
- Shinjuku-SQ no Preemption
- Shinjuku-SQ
- Shinjuku-MQ

Slowdown=
$$\frac{Total\ Latency}{Service\ Time}$$

# Time slice matters

Legend:
- 5us
- 20us
- 50us
- 100us

Synthetic Workload
Bimodal
50% 1us – 50% 100us

y-axis: 99% Slowdown (100, 75, 50, 25)

x-axis: Throughput (kRPS) (0, 50, 100, 150, 200, 250)

Better (downward arrow)
Better (rightward arrow)