# 5G Connected Edge Cloud for Industry 4.0 Transformation
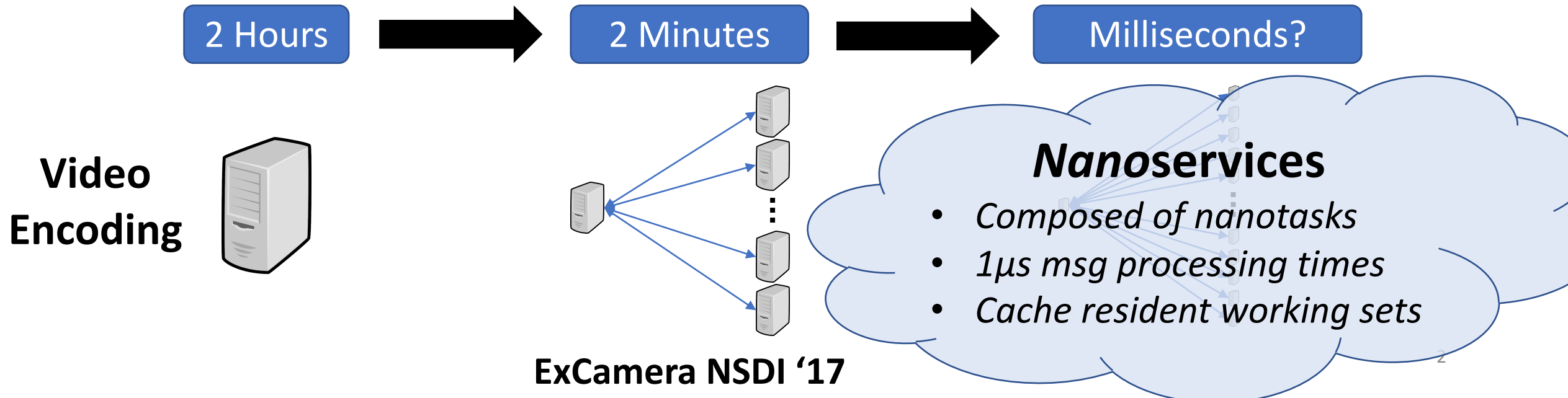
## ONF SPOTLIGHT

# The **nanoPU**:
# Redesigning the CPU-Network Interface to Minimize RPC Tail Latency

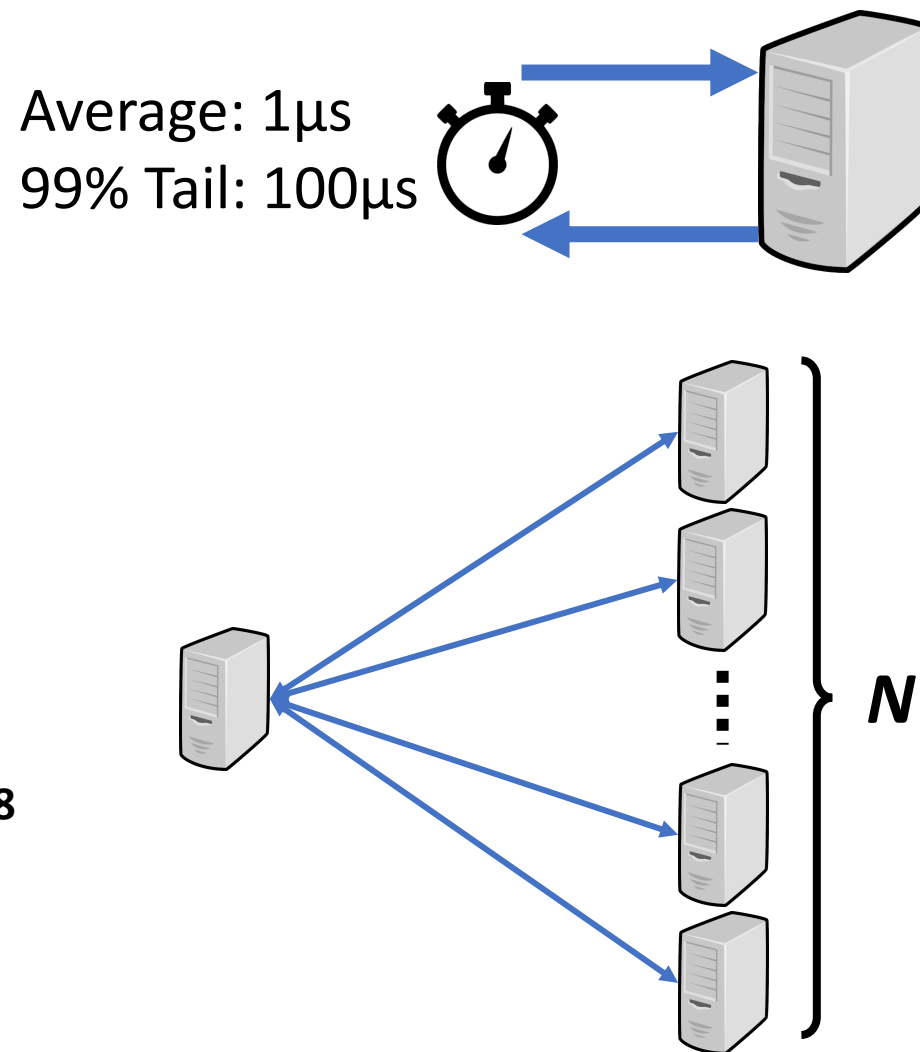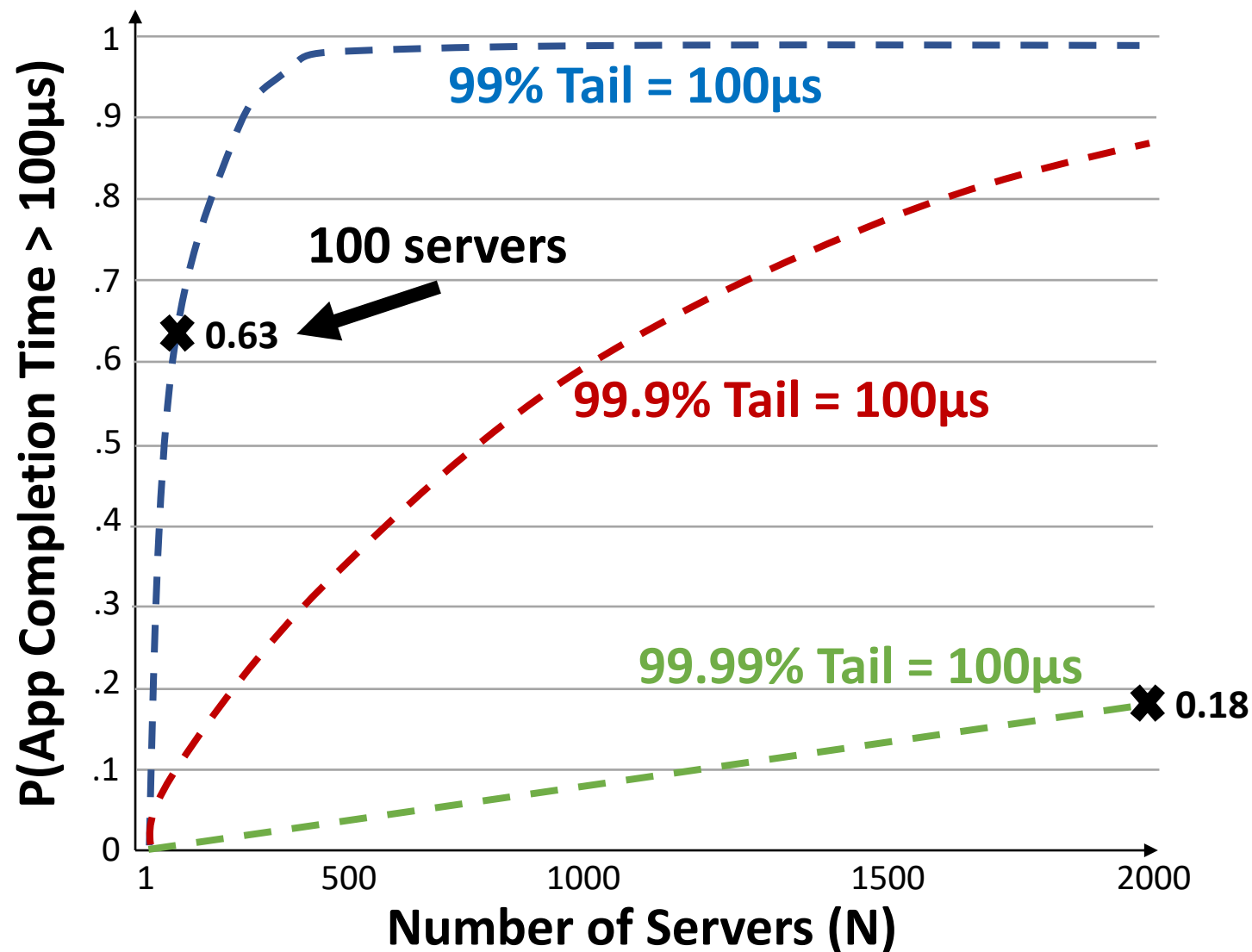**Stephen Ibanez**, Alex Mallery, Serhat Arslan, Theo Jepsen, Muhammad Shahbaz, Changhoon Kim, Nick McKeown

*Stanford University*

2020

# Towards Fine-Grained Computing

- Serverless computing (deployed as microservices) enables "fine-grained" computing on thousands of cores, reducing completion times:
  - Video encoding (ExCamera NSDI'17)
  - Object classification (Sprocket SoCC'18)
  - Software compilation (gg ATC'19)
  - MapReduce-style analytics (Pyren SoCC'17 , Flint CLOUD'18, Locus NSDI'19)

**2 Hours** → **2 Minutes** → **Milliseconds?**

**Video Encoding**

**ExCamera NSDI '17**

***Nanoservices***
- *Composed of nanotasks*
- *1μs msg processing times*
- *Cache resident working sets*

2

# Tail Latency Matters!



**99% Tail = 100μs**

**100 servers**

**✖ 0.63**

**99.9% Tail = 100μs**

**99.99% Tail = 100μs**

**✖ 0.18**

**P(App Completion Time > 100μs)**

1
.9
.8
.7
.6
.5
.4
.3
.2
.1
0

1    500    1000    1500    2000

**Number of Servers (N)**

Average: 1μs
99% Tail: 100μs

*N*

*Credit: L. Barroso, et al. "The Datacenter As a Computer"*, Ch. 2

# What causes high RPC tail latency?

**Primary Cause**

Poor job scheduling access to critical shared resources:
- *Network fabric resources*
- *CPU cores*
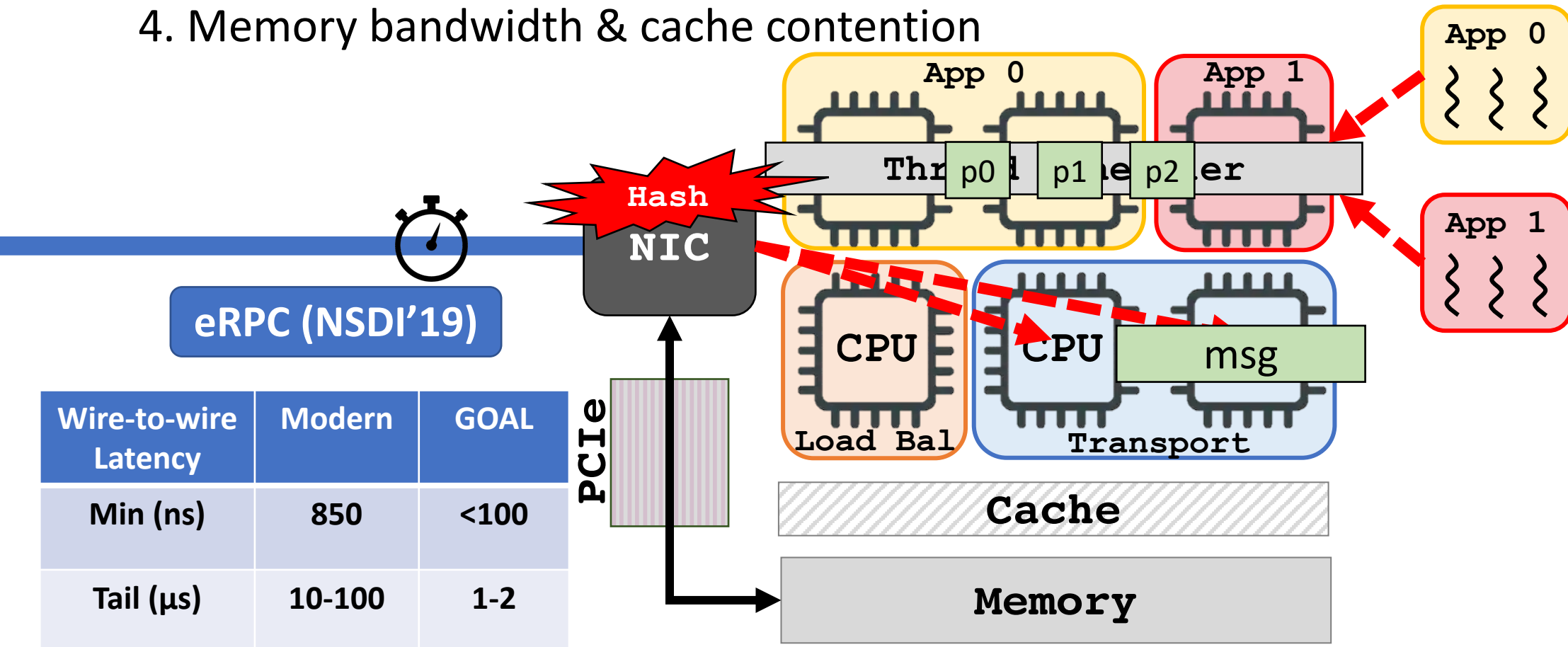- *Host memory bandwidth and cache space*

# Problems with Modern CPU/NIC/OS Designs

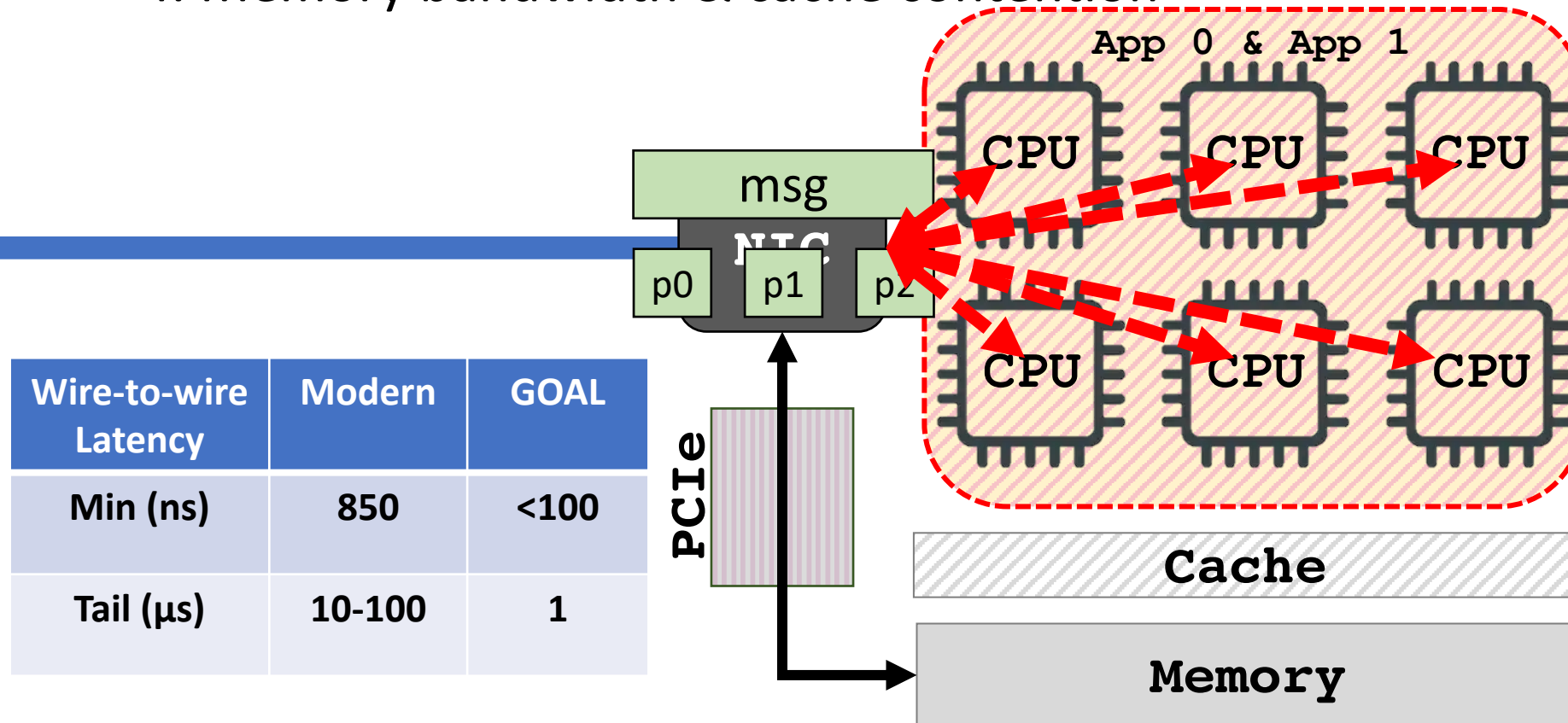1. Suboptimal congestion control for scheduling of the network fabric
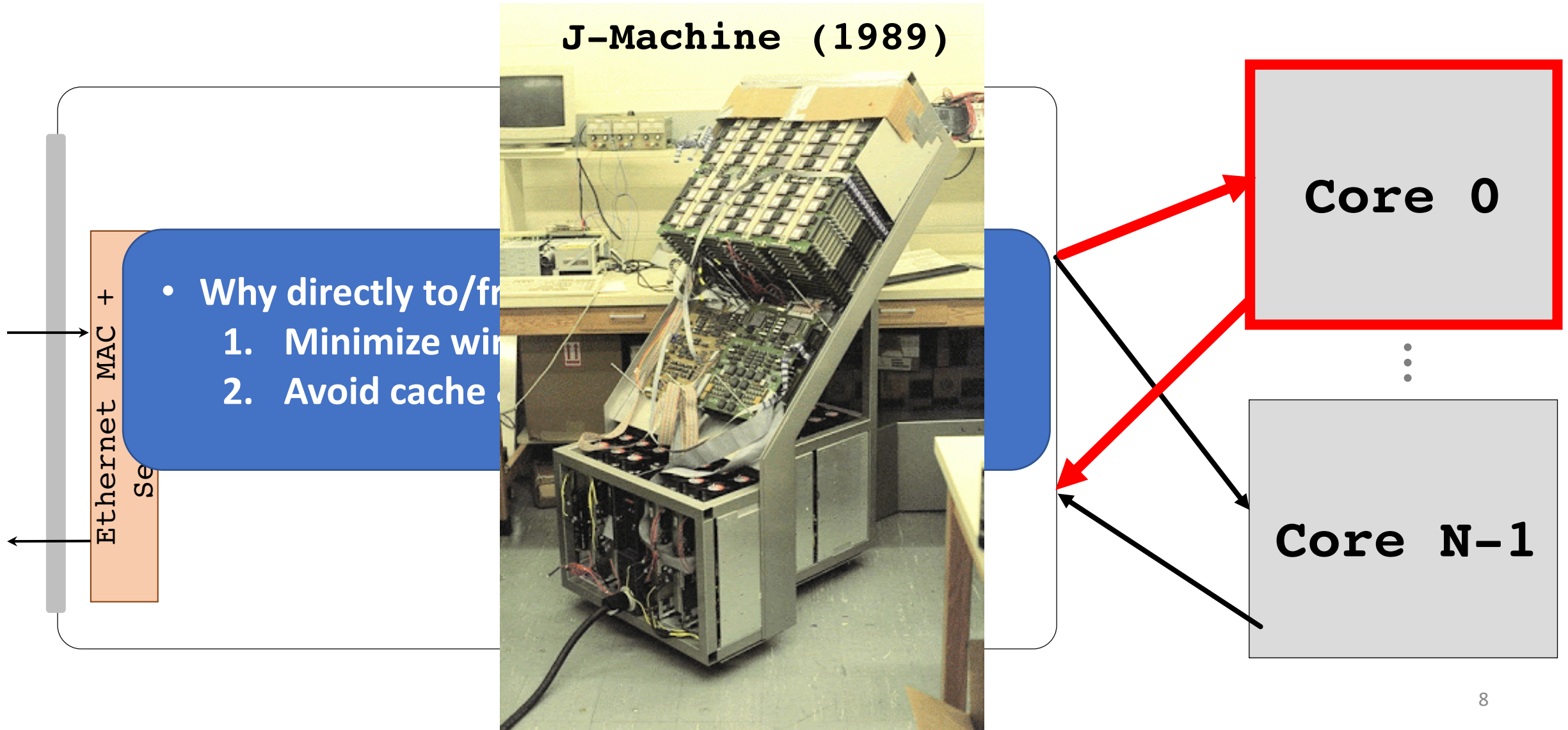
# Problems with Modern CPU/NIC/OS Designs

1. Suboptimal congestion control for scheduling of the network fabric
2. Inefficient load balancing across cores
3. Inefficient thread scheduling on each core
4. Memory bandwidth & cache contention



| Wire-to-wire Latency | Modern | GOAL |
|---|---|---|
| Min (ns) | 850 | <100 |
| Tail (µs) | 10-100 | 1-2 |

# Problems with Modern CPU/NIC/OS Designs

1. Suboptimal congestion control for scheduling of the network fabric
2. Inefficient load balancing across cores
3. Inefficient thread scheduling on each core
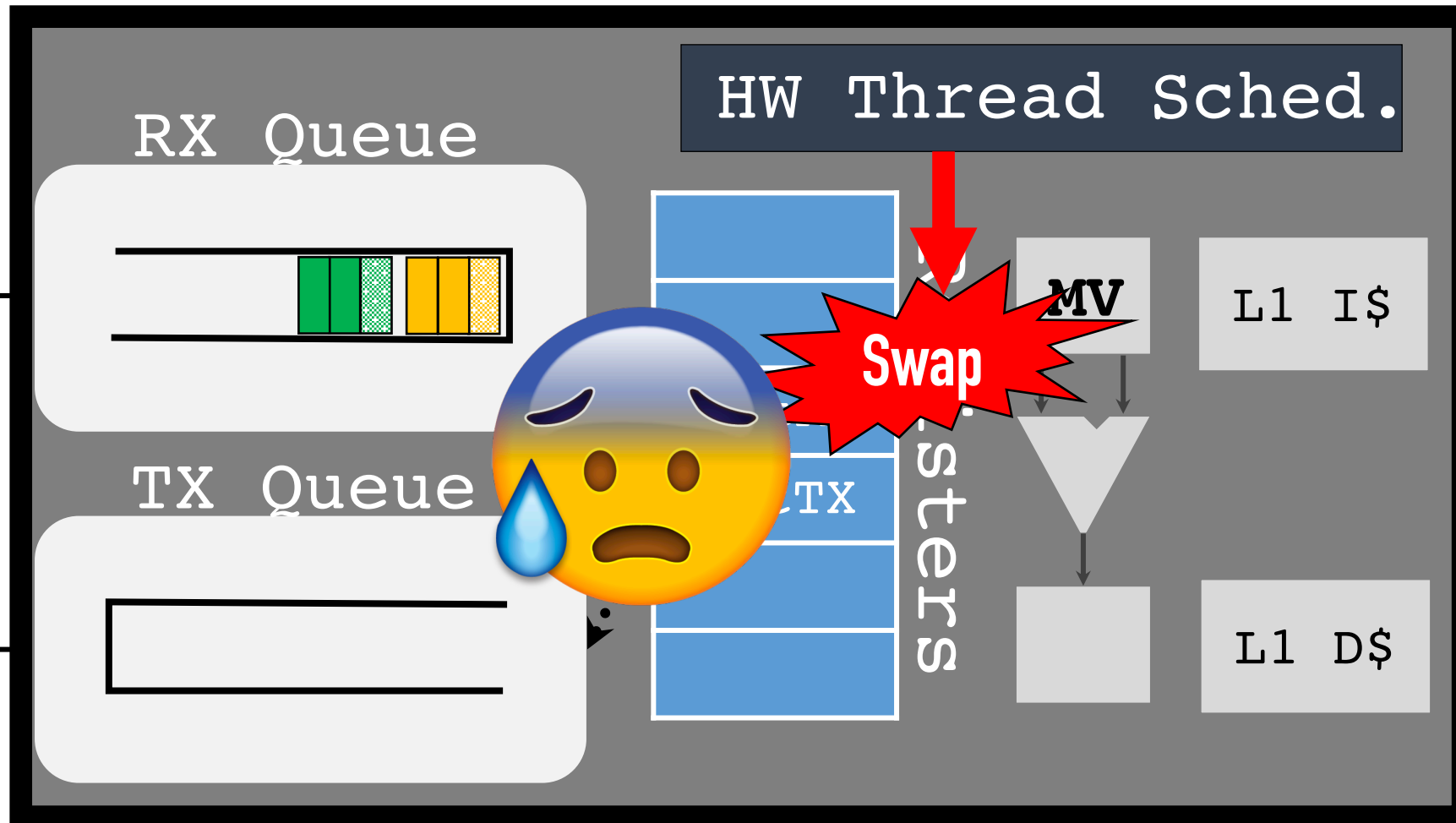4. Memory bandwidth & cache contention



App 0 & App 1

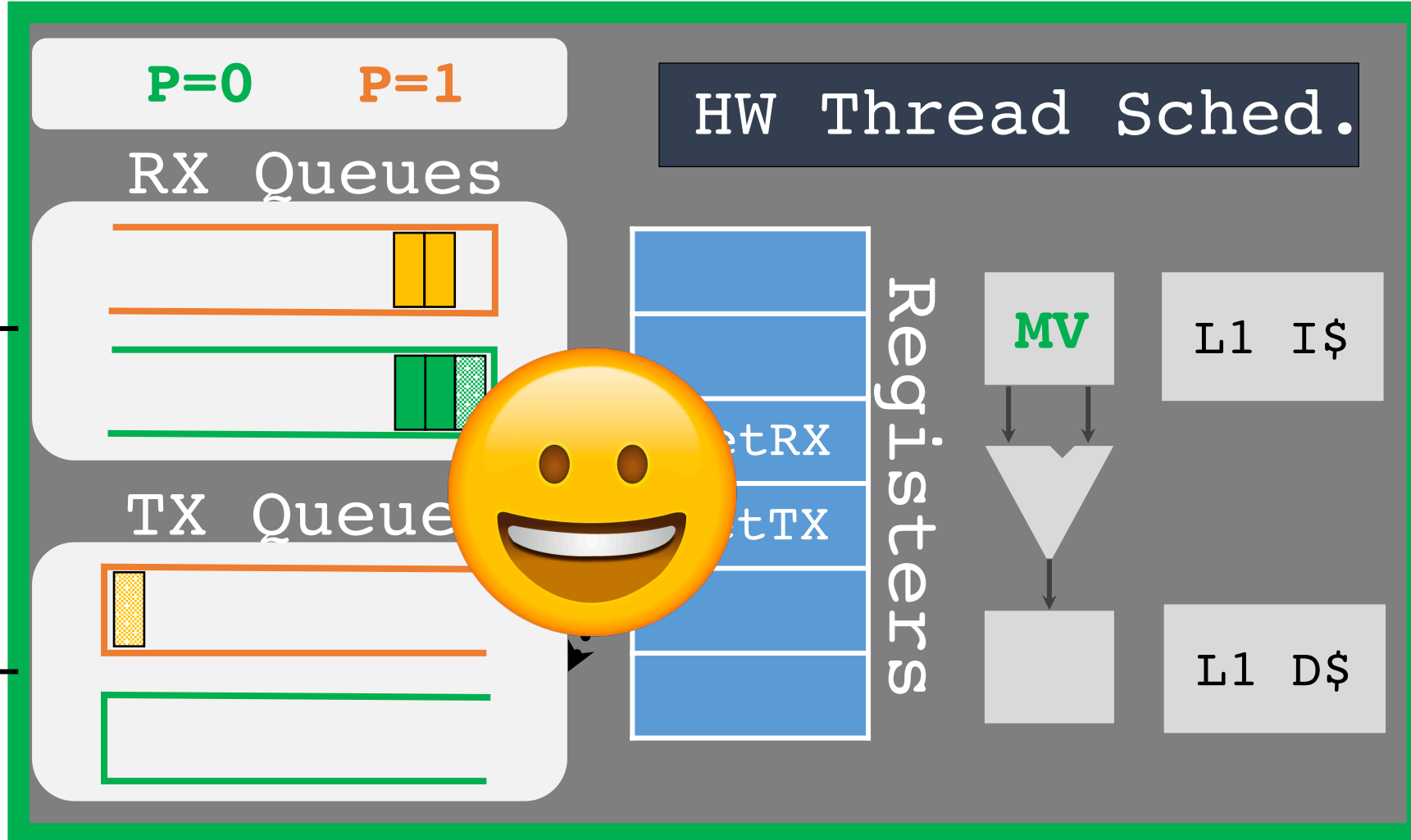CPU CPU CPU

CPU CPU CPU

msg

NIC

p0 p1 p2

PCIe

Cache

Memory

| Wire-to-wire Latency | Modern | GOAL |
|---|---|---|
| Min (ns) | 850 | <100 |
| Tail (µs) | 10-100 | 1 |

# The nanoPU Design



J-Machine (1989)

**Ethernet MAC +**

- Why directly to/fr
  1. Minimize wir
  2. Avoid cache

Core 0

⋮

Core N-1

# nanoPU's Register File Network Interface

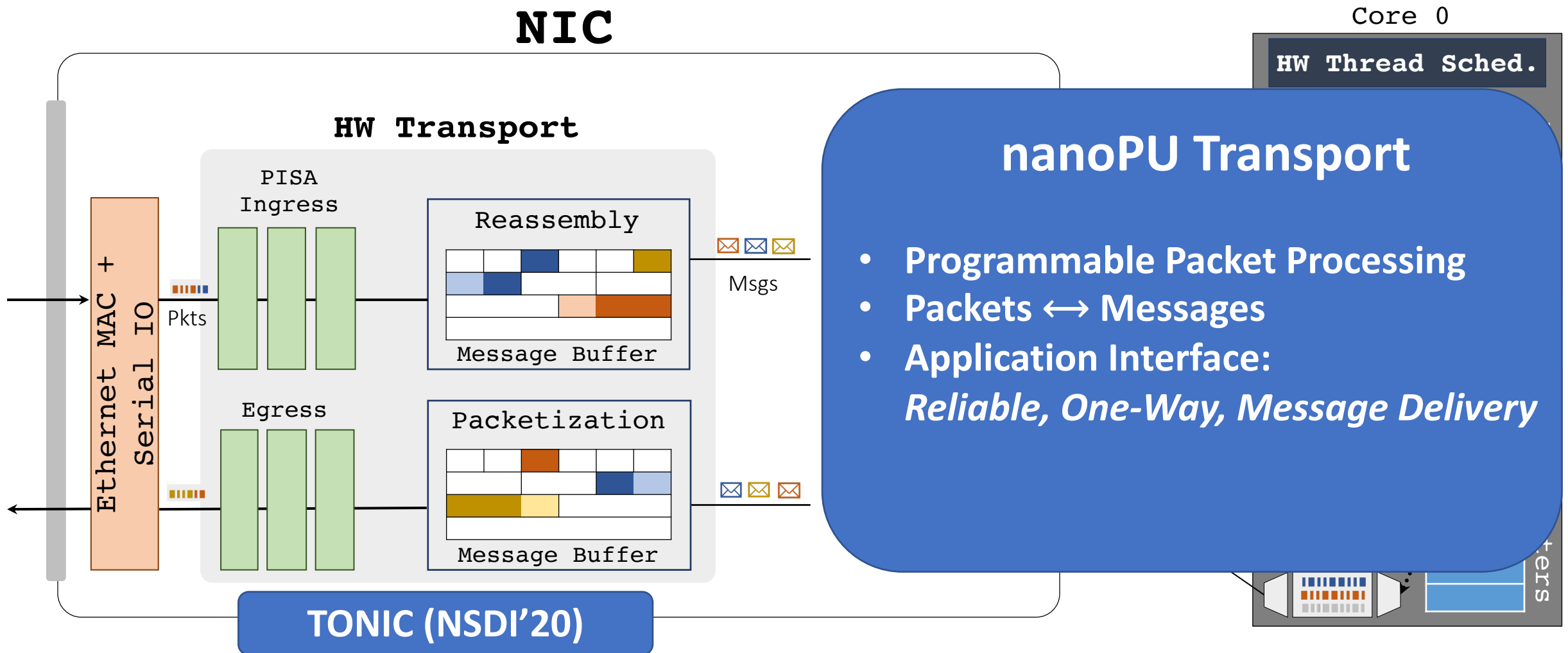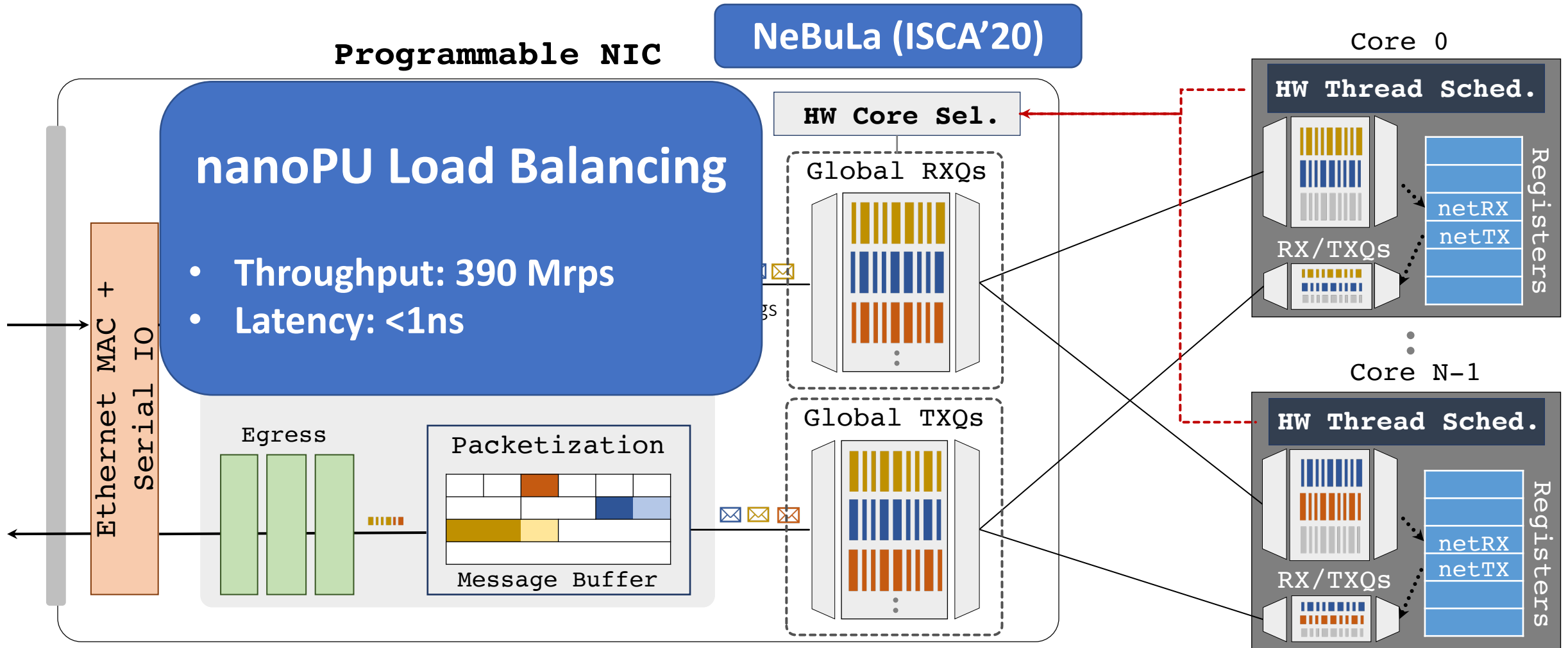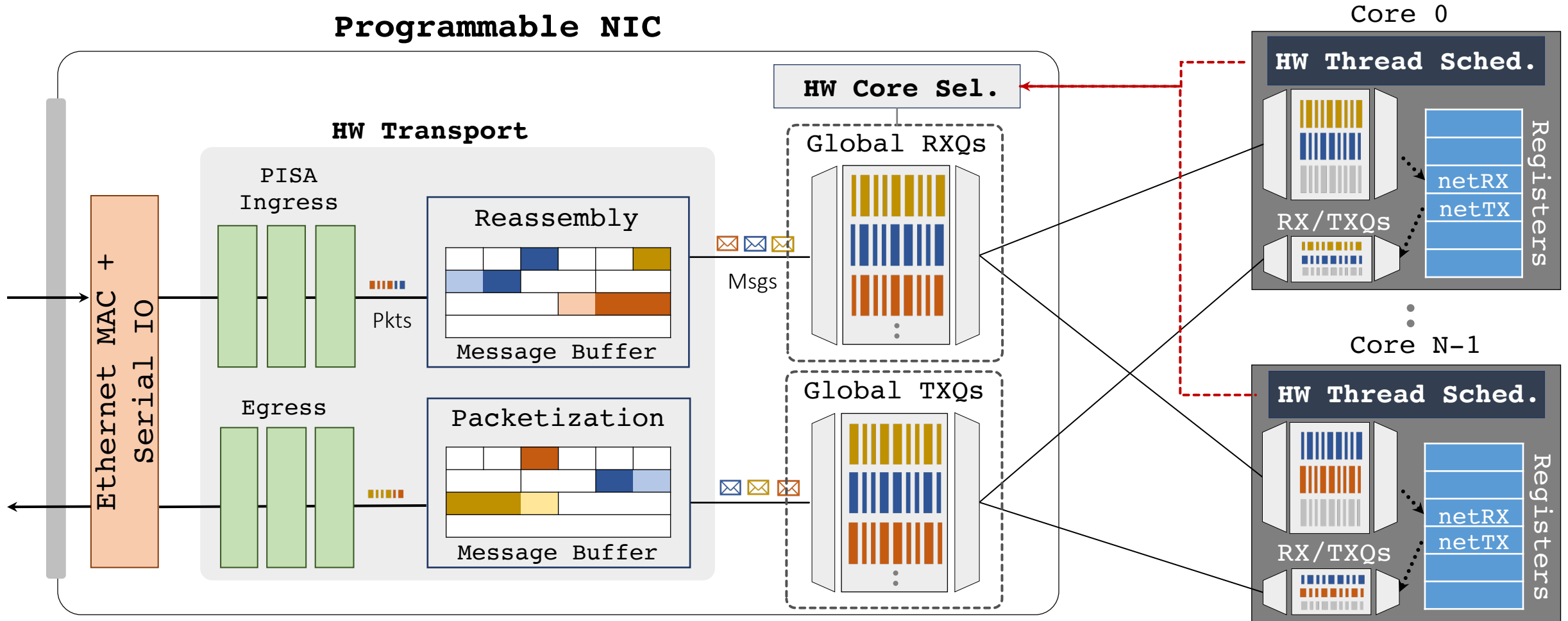# nanoPU's Register File Network Interface

# nanoPU Transport

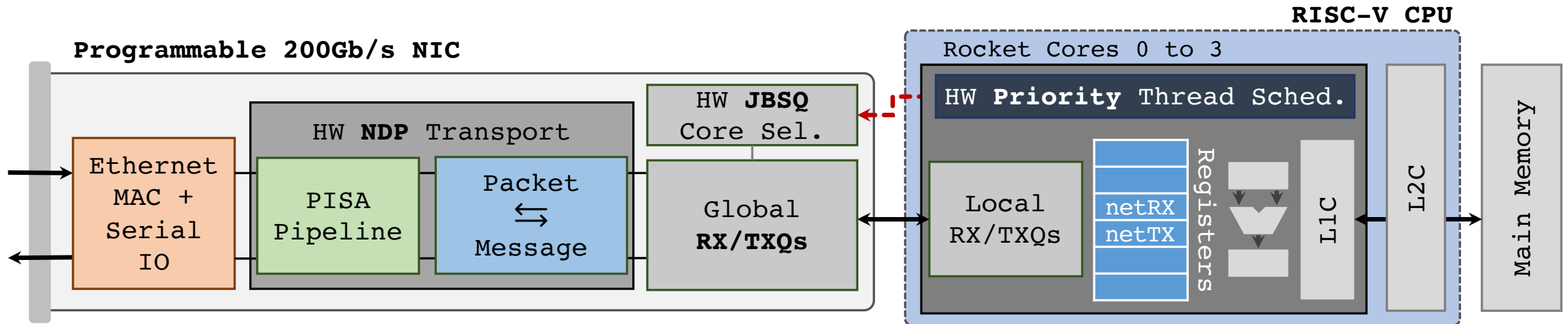# nanoPU Core Load Balancing

# The nanoPU Design

# nanoPU Prototype

- Quad-core nanoPU based on the open source RISC-V Rocket core
- 4,300 lines of Chisel code
- 1,200 lines of C and RISC-V assembly for custom *nano*kernel
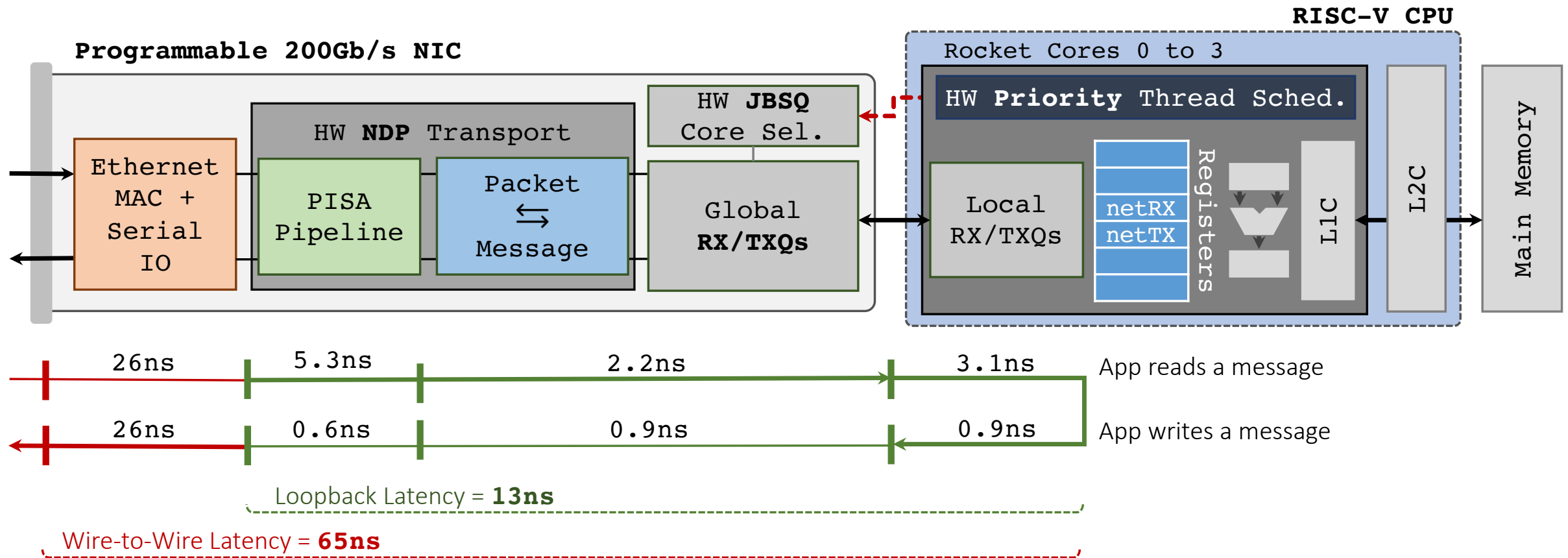- Implements NDP transport

# Evaluation Methodology

- nanoPU prototype running on AWS FPGAs
- Large-scale (hundreds of cores), cycle-accurate simulations with Firesim
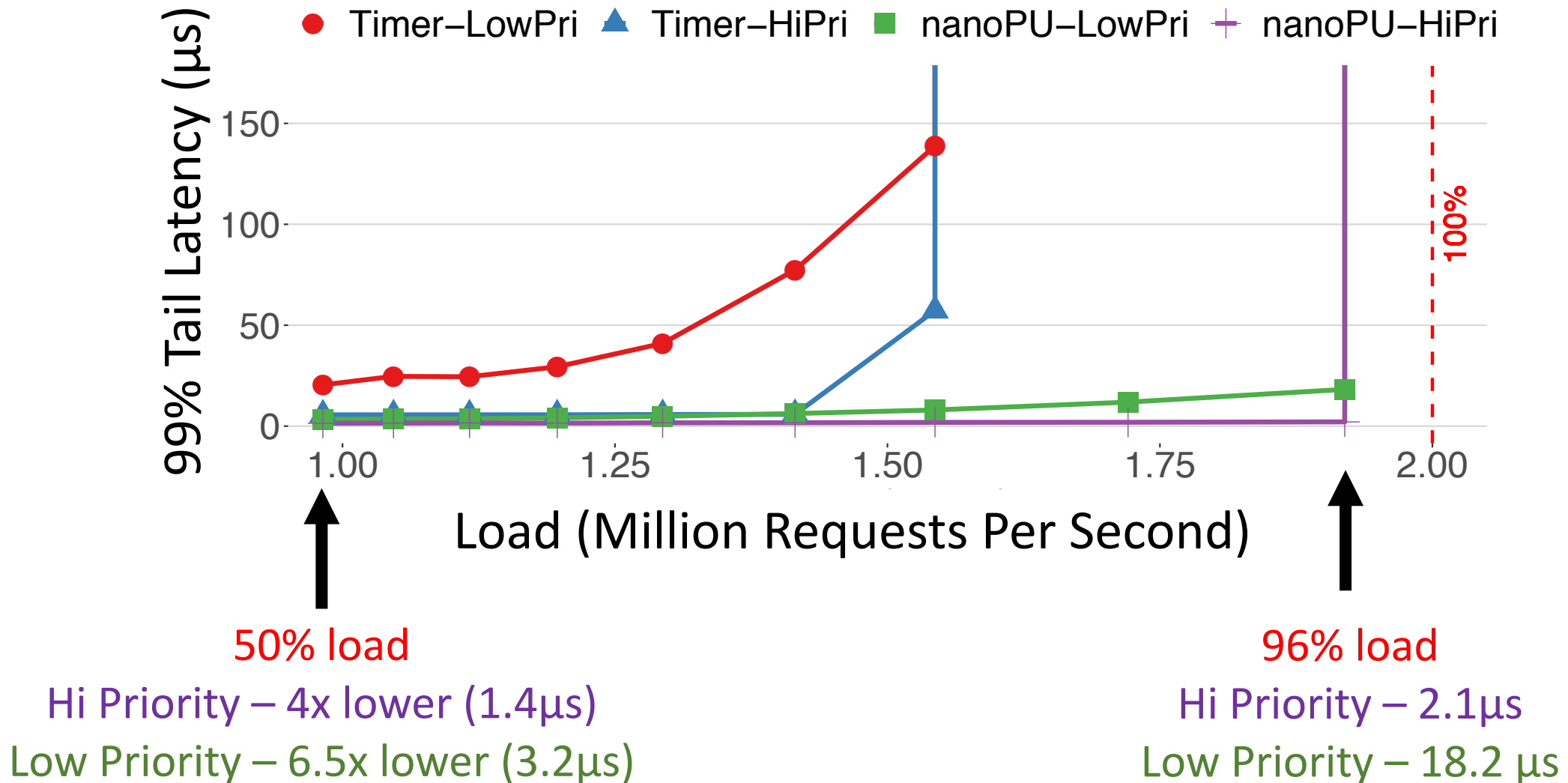- FPGA clock rate is 90MHz, simulated target clock rate is 3.2GHz



ISCA'18

# Microbenchmarks

# Thread Scheduling Evaluation

- Single core running two threads: high priority and low priority
- 500ns request processing time
- 10K requests per thread
- nanoPU HW thread scheduler vs. 5us timer-driven thread scheduler

# Thread Scheduling Evaluation



18

# Additional Evaluations



**Core Load Balancing**



**NDP Transport**
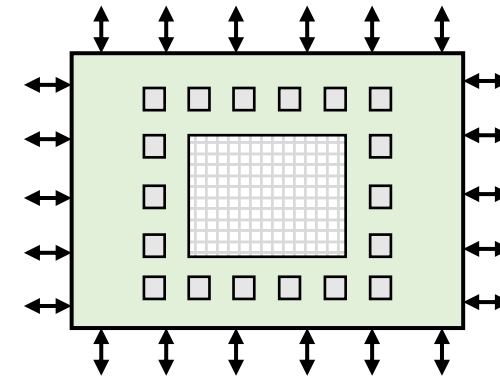
# Real Applications Running on the nanoPU

- MICA Key Value Store
- Chain Replication
- Raft Consensus
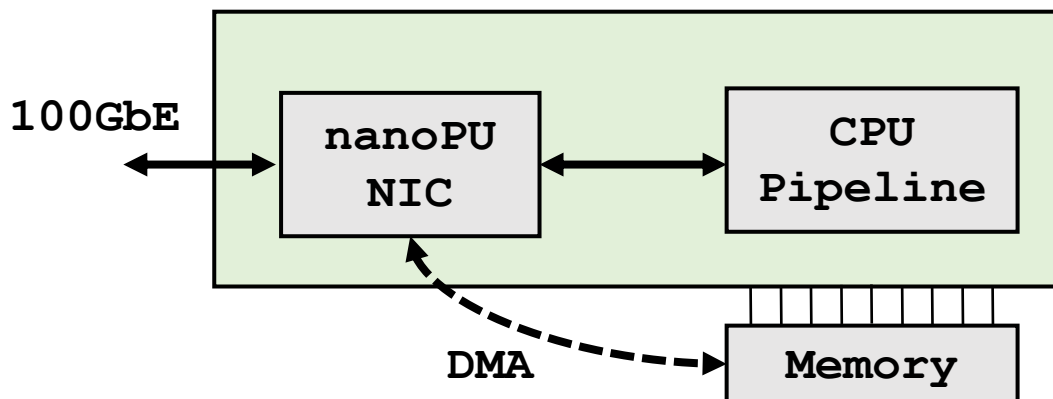
# nanoPU Deployment Possibilities

① nanoPU Cluster
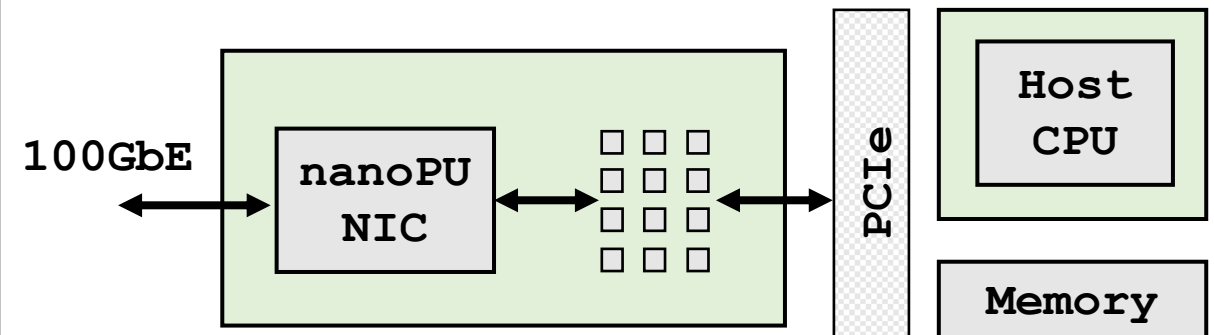
② nanoPU High IO Capacity Package

- 512 cores
- 256 x 100Gb/s

③ Modified Conventional CPU

100GbE ⟷ **nanoPU NIC** ⟷ **CPU Pipeline**

DMA ⟶ **Memory**

④ nanoPU SmartNIC

100GbE ⟷ **nanoPU NIC** ⟷ **PCIe** ⟷ **Host CPU**

**Memory**

# nanoPU Conclusions

## Key Takeaway:

To truly minimize average and tail RPC latency:

1. Fast path directly between network and CPU register file
2. Move key resource scheduling decisions to HW: **transport**, **load balancing, thread scheduling**

## Challenges:

- Need to rewrite applications
- Figure out how to use more sophisticated processors

# 5G Connected Edge Cloud for Industry 4.0 Transformation

**ONF SPOTLIGHT**

# Thank You

sibanez@stanford.edu