

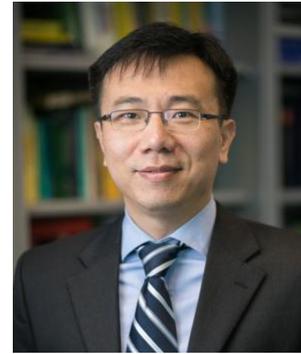
TCP \approx RDMA: CPU-efficient Remote Storage Access with i10



Jaehyun Hwang



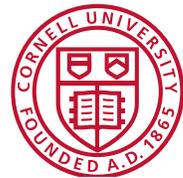
Qizhe Cai



Ao Tang

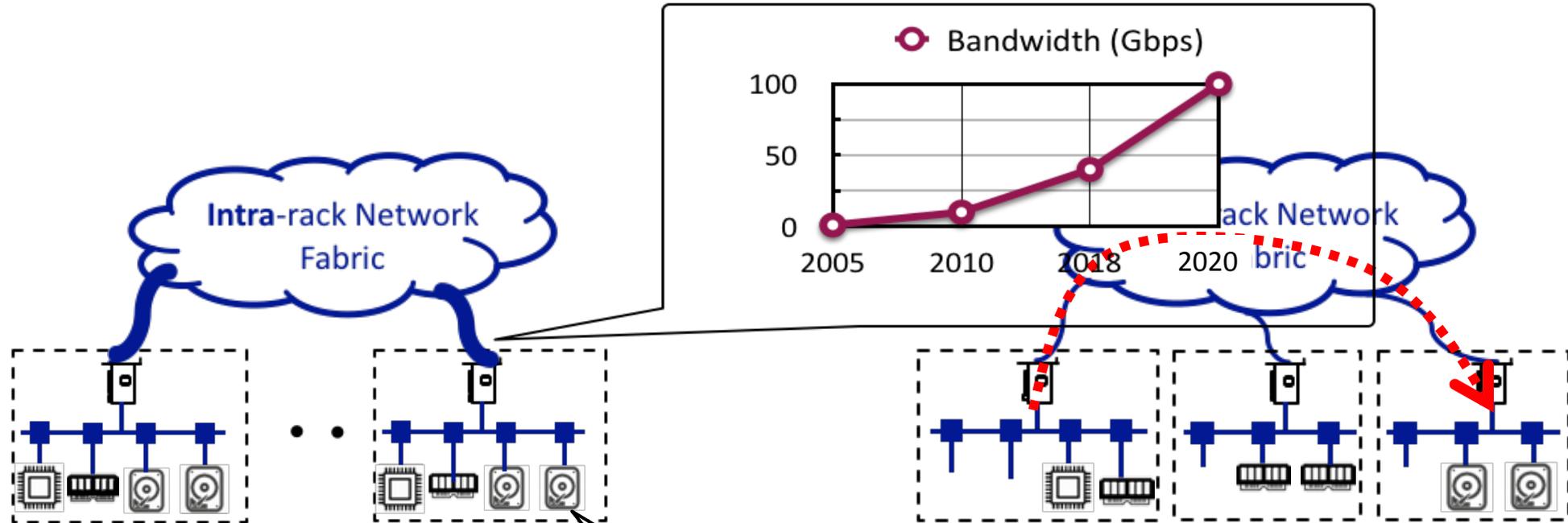


Rachit Agarwal



Cornell University

i10 Motivation: Two trends in Remote I/O



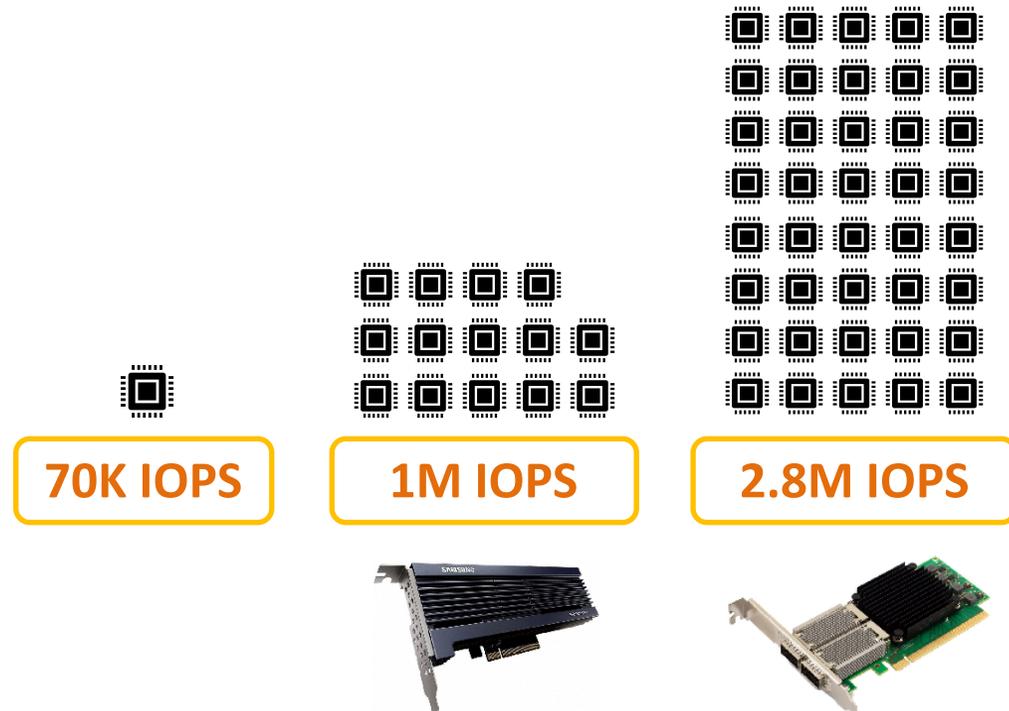
- Move from SCSI to NVMe
- Emergence of NVMe SSDs:
 - >1M IOPS (read), >400K IOPS (write)

1. High-throughput NVMe SSD, high-bandwidth access links

2. Apps use disaggregated storage via Remote I/O

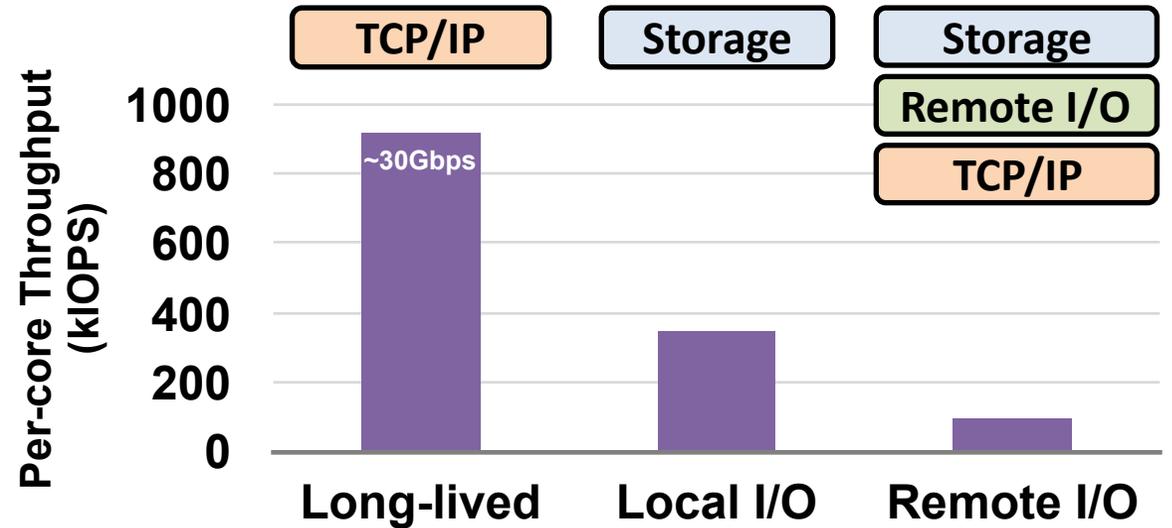
What do these trends mean for Remote I/O?

- Software stack (iSCSI) performance



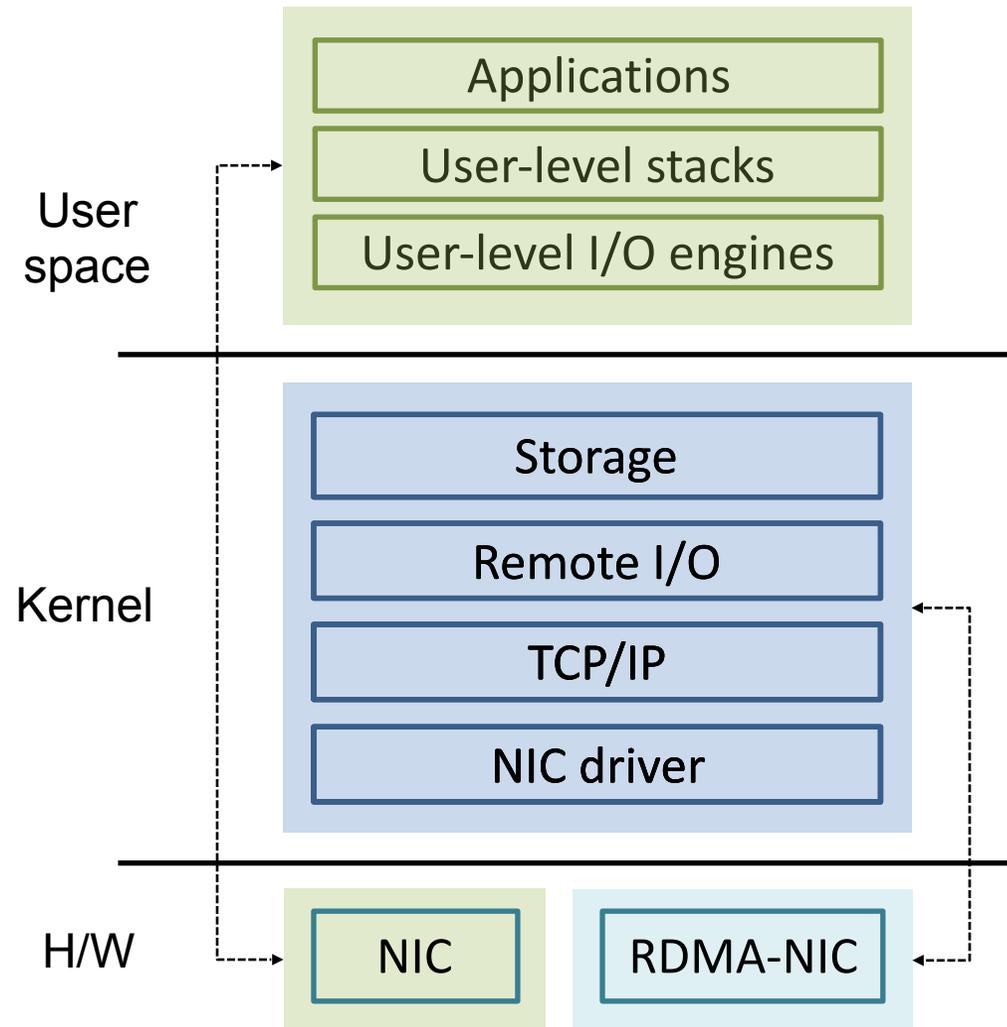
Bottlenecks pushed back to OS, software stack!

- Storage + network overlap



Bottlenecks at the boundary of storage and network stacks!

Previous Approaches



User-space storage and network stacks

- Storage + Remote I/O (user) + DPDK
- Performance: Good!

In-Kernel (NVMe-over-TCP)

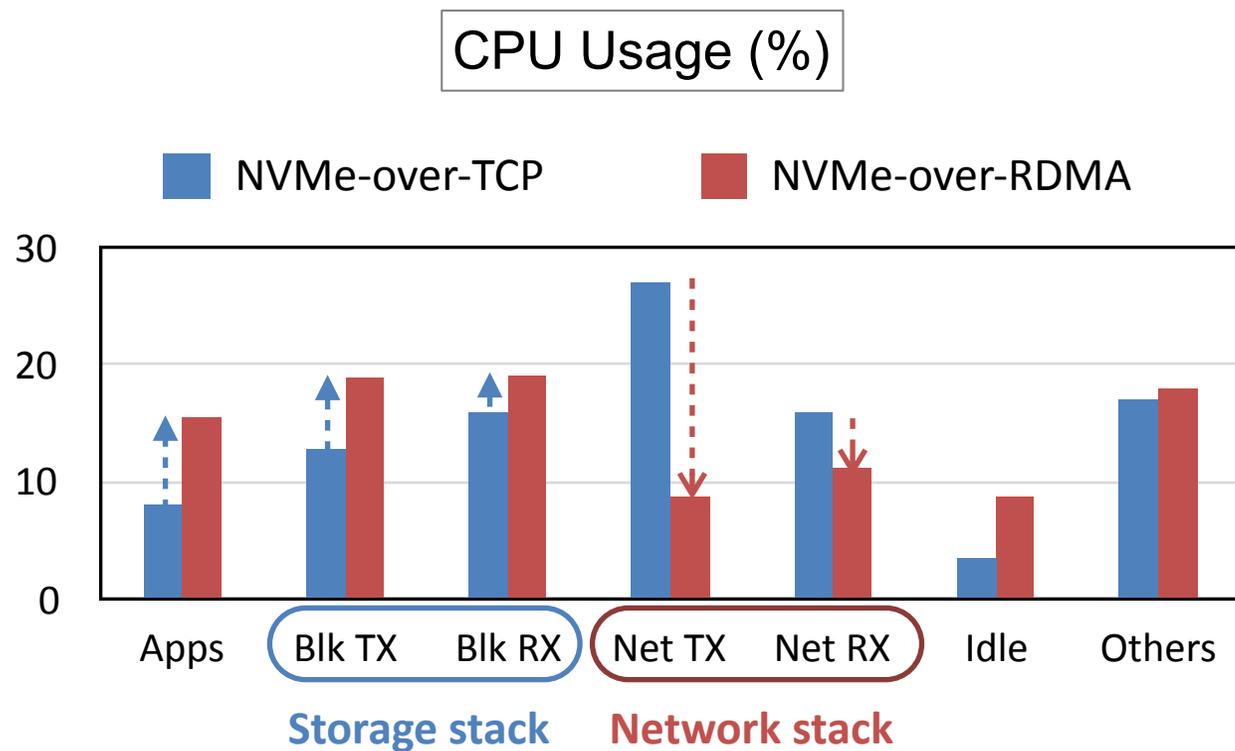
- Storage + Remote I/O + TCP (all in the kernel)
- Performance: Not-so-good!

Fundamental?

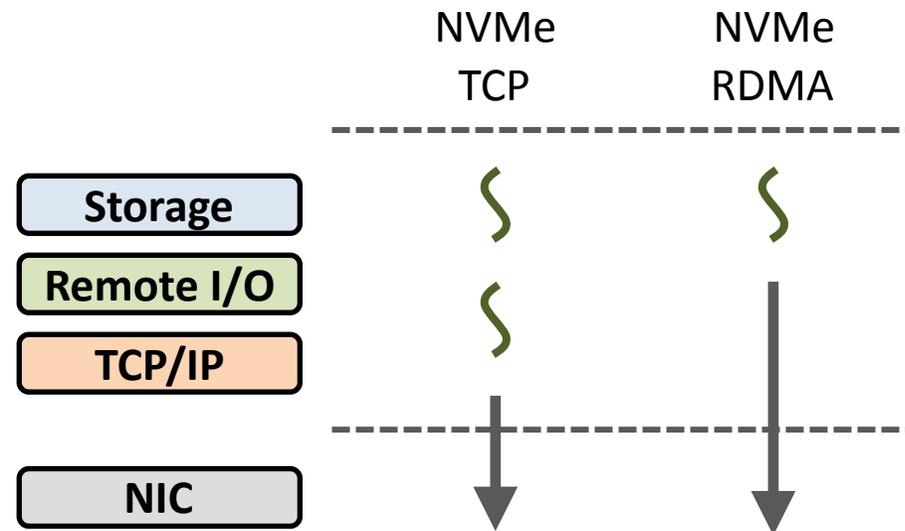
NVMe-over-RDMA

- Storage + Remote I/O (kernel) + RDMA
- Performance: Good!

Remote Storage Access Overheads: TCP vs. RDMA



Network processing overhead!



Context switching overhead!

i10 Summary

- A new **remote I/O stack** implemented **entirely in the kernel**
 - No changes in apps, no changes in TCP/IP stack, no changes in hardware
- **Throughput-per-core similar to NVMe-over-RDMA**
 - Latency within 1.7x of RDMA (for SSD accesses)

- **Three *simple* ideas**

- **i10-lane**: dedicated resources
- **i10-caravans**: request and data batching
- **Delayed doorbells**: Interrupt coalescing

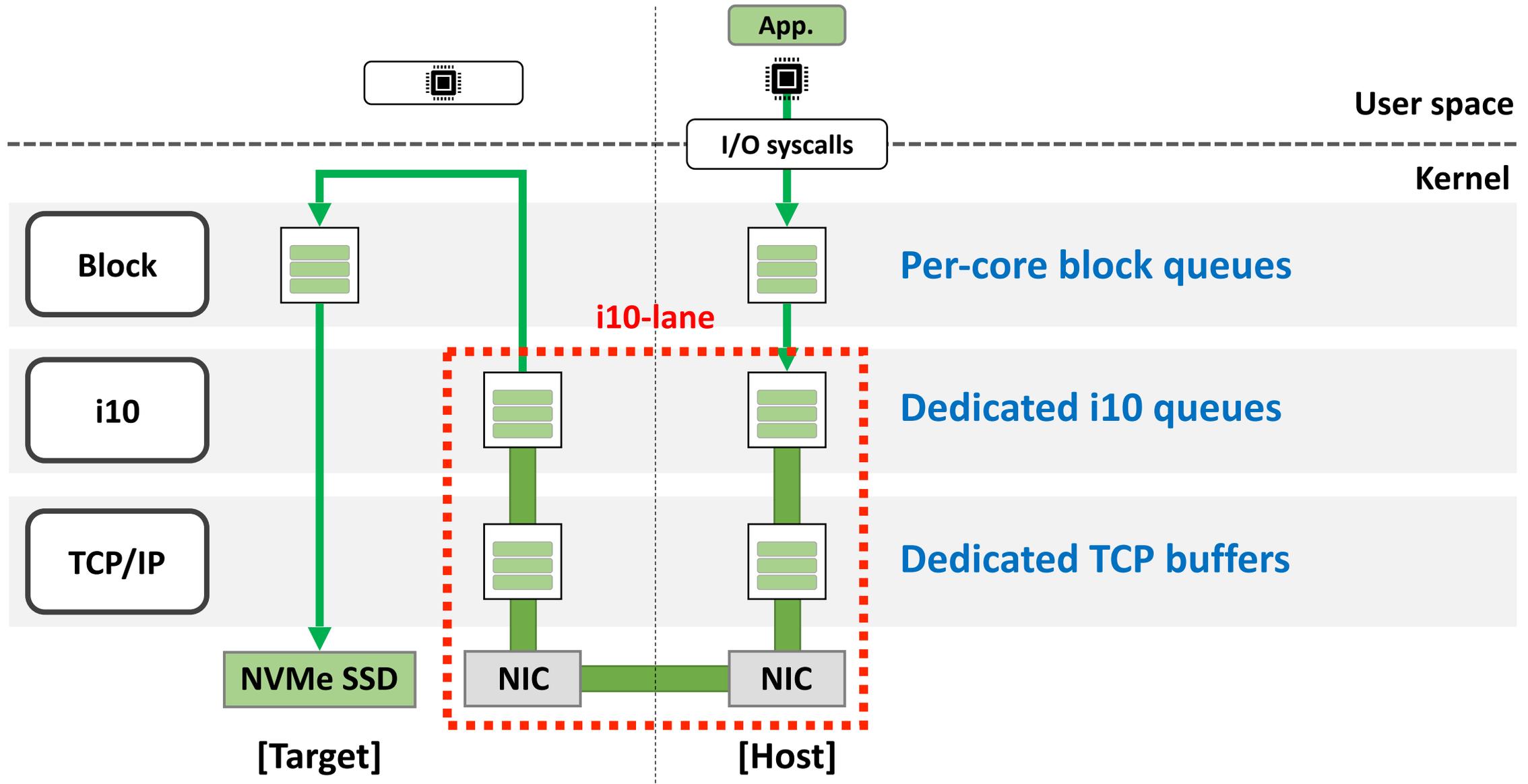
Minimize network processing



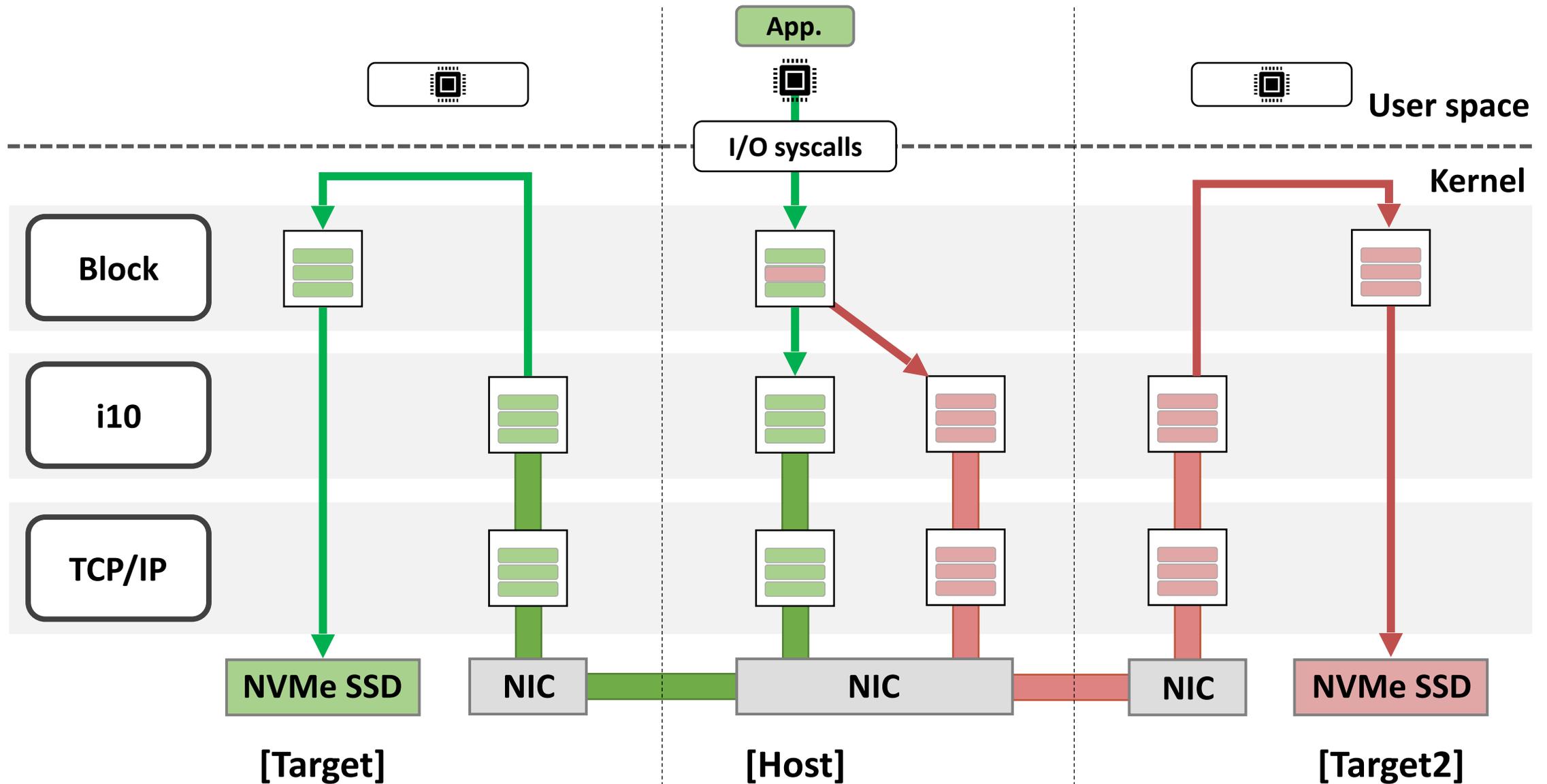
Minimize context switching



i10-lane: Dedicated per-(core, target) pair “pipe”

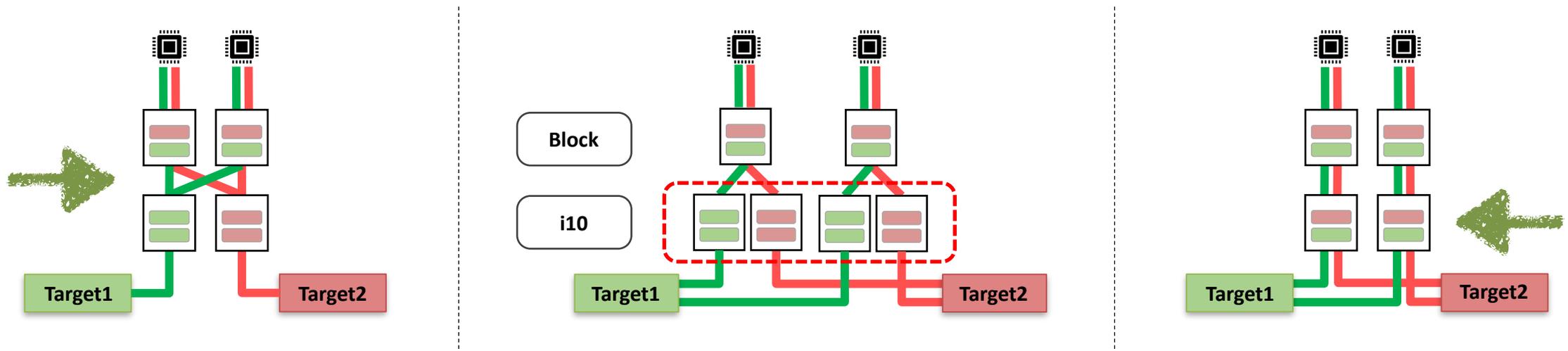


i10-lane: Dedicated per-(core, target) pair “pipe”



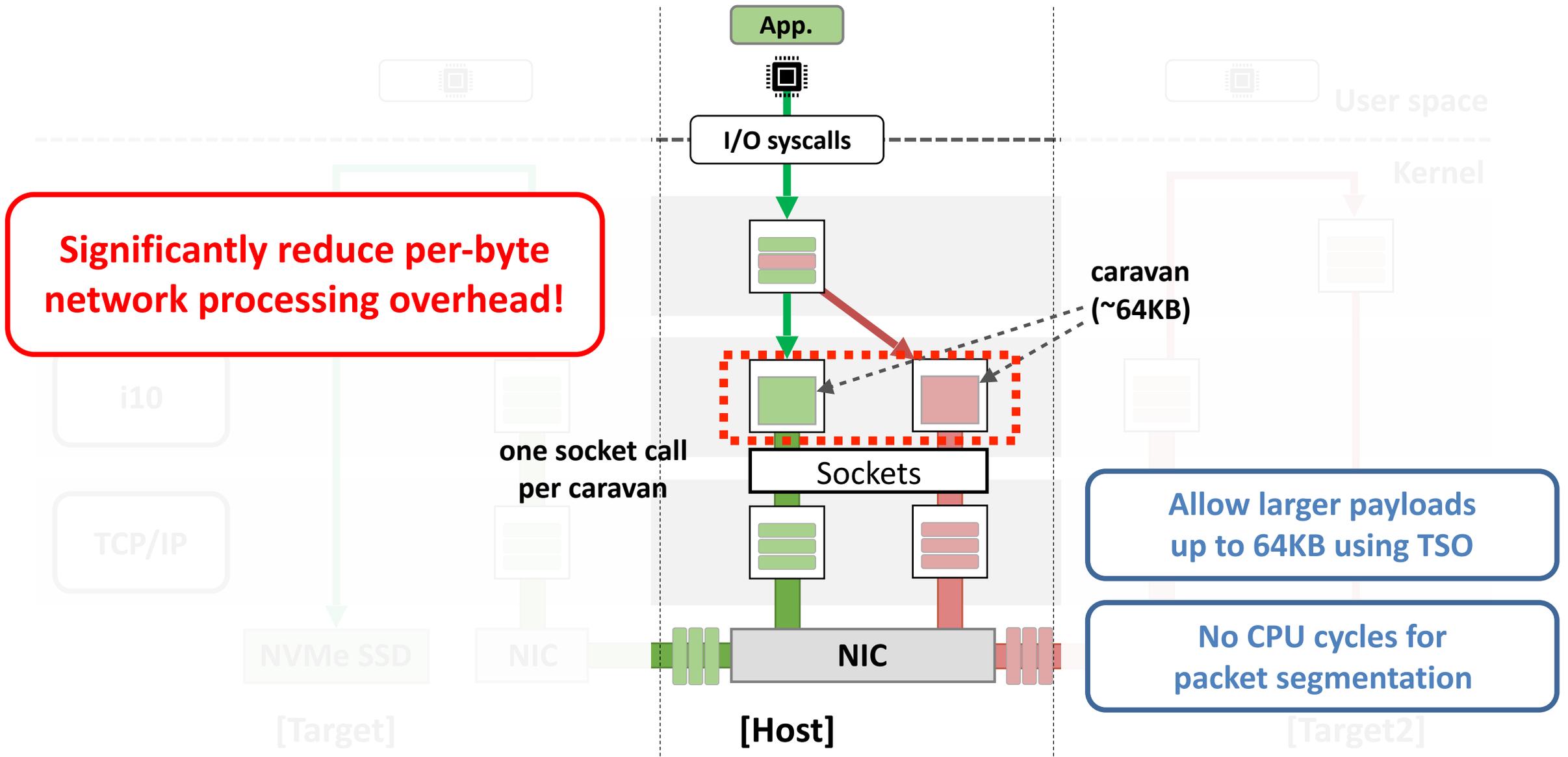
i10-lane: Why per-(core, target) pair lanes?

- **Per-target:** Too much contention
- **Per-core:** Fewer batching opportunities

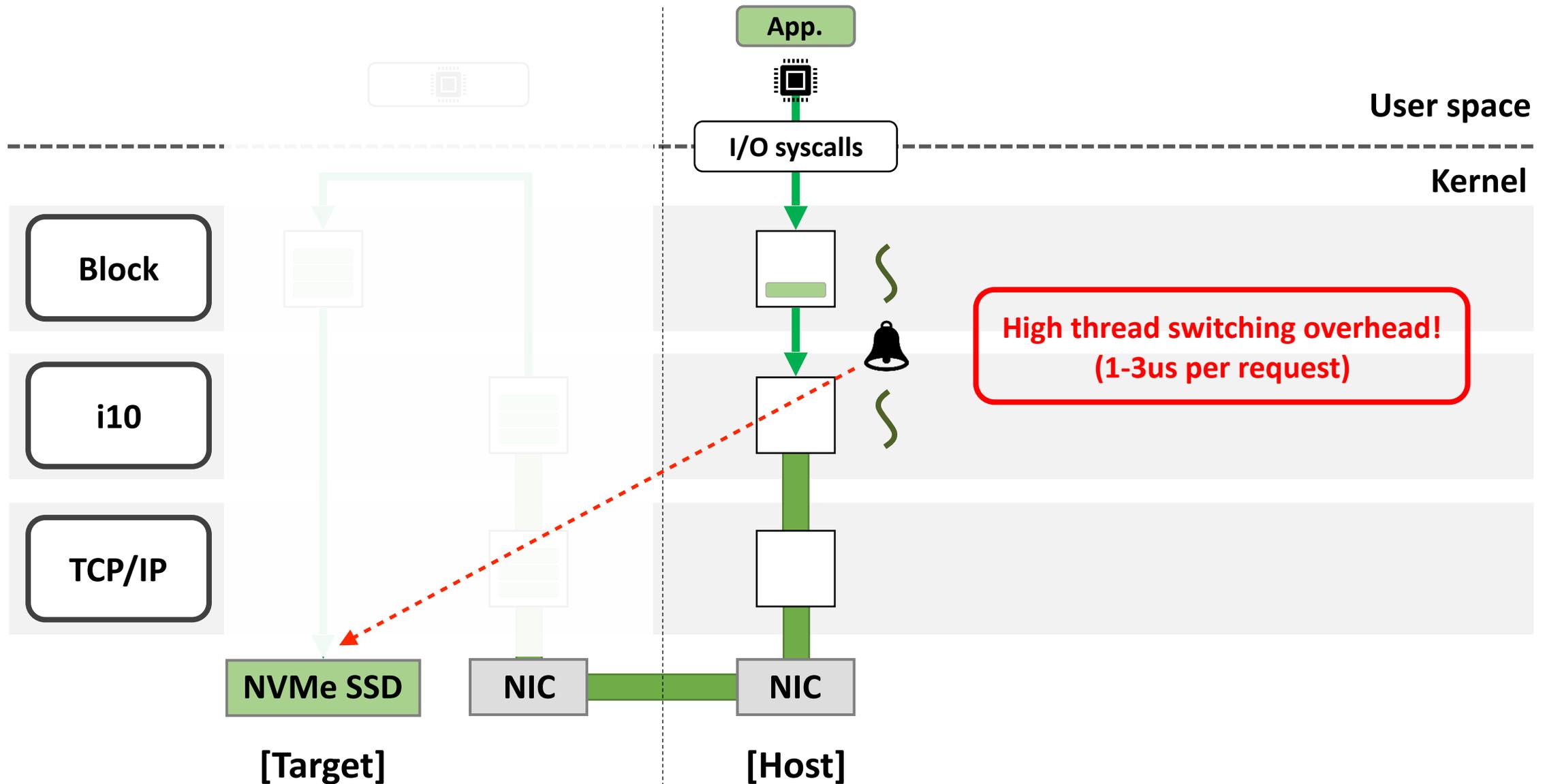


All requests in each i10 queue are destined to the same target over the same TCP connection

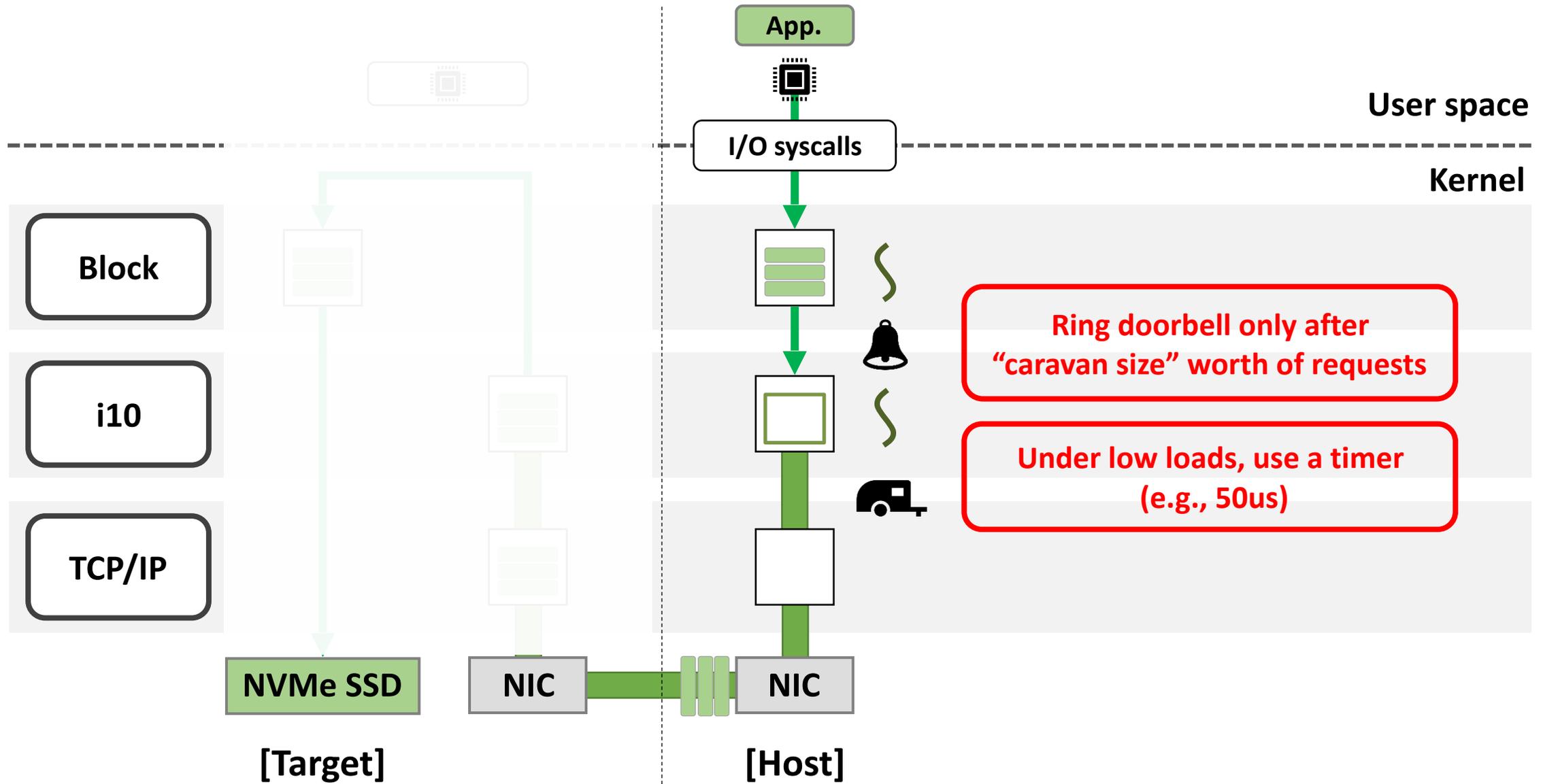
i10 Caravans: i10-lanes enable efficient batching



Context switching in Remote I/O (without i10)



Delayed doorbells: Minimizing context switching



i10 Evaluation Setup

- **Two 24-core servers connected directly**
 - 100Gbps Mellanox CX-5
 - No switches in middle — ensure bottlenecks in the kernel
- **NVMe-device at both servers**
 - ~700k IOPS (read), ~400k IOPS (write)
 - ~100us read latency
- **No specialized hardware functionalities used in i10 evaluation**
 - For hardware and software configuration, see the paper.

i10 Evaluation: how does i10 performance ...

- ... compare to NVMe-over-RDMA?

- Metrics of interest: throughput per core, average latency, tail latency

- ... compare to user-space stacks?

Please see our paper!

- ... vary with different workloads, hardware and applications?
 - **What do we expect from TCP \approx RDMA ?**

- read/write ratios

- Delay

Throughput

Latency

- Aggregation size

Answer: Comparable (or better)



Not terrible (<1.7X)



- I/O Request sizes

- Storage device access latency

- Real applications

- Number of target devices

- ... scale with number of cores?

- ... depend on various design aspects (lanes, caravans, delayed doorbell)?

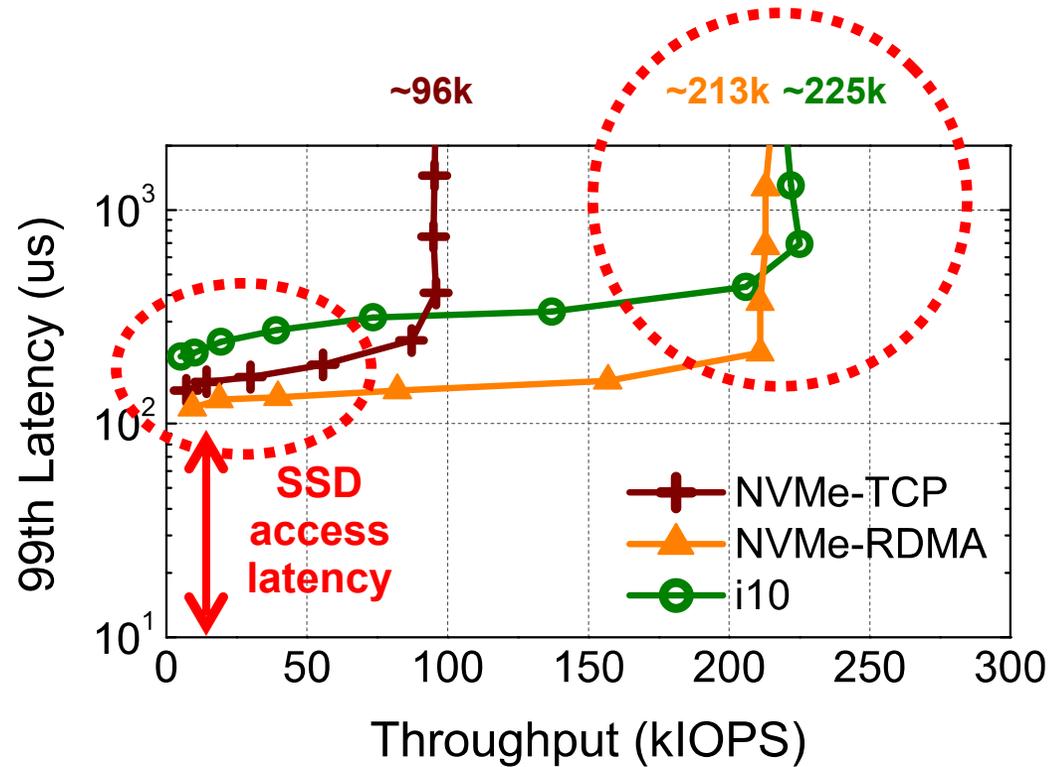
Teacher : why do you have the same answers ?

Student : because we have the same questions

Teacher :



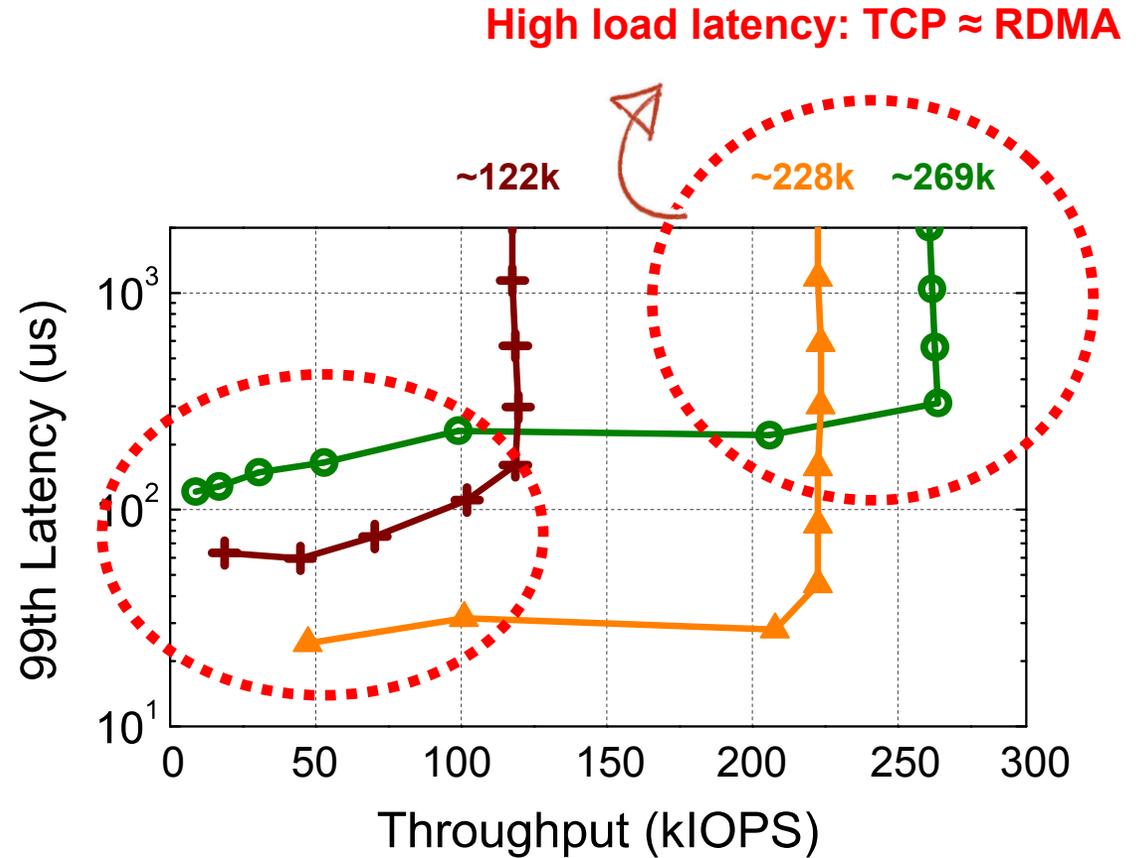
Single core performance



NVMe SSD

TCP ≈ RDMA:

- Throughput: **Comparable**
- Tail latency: **<1.7X**

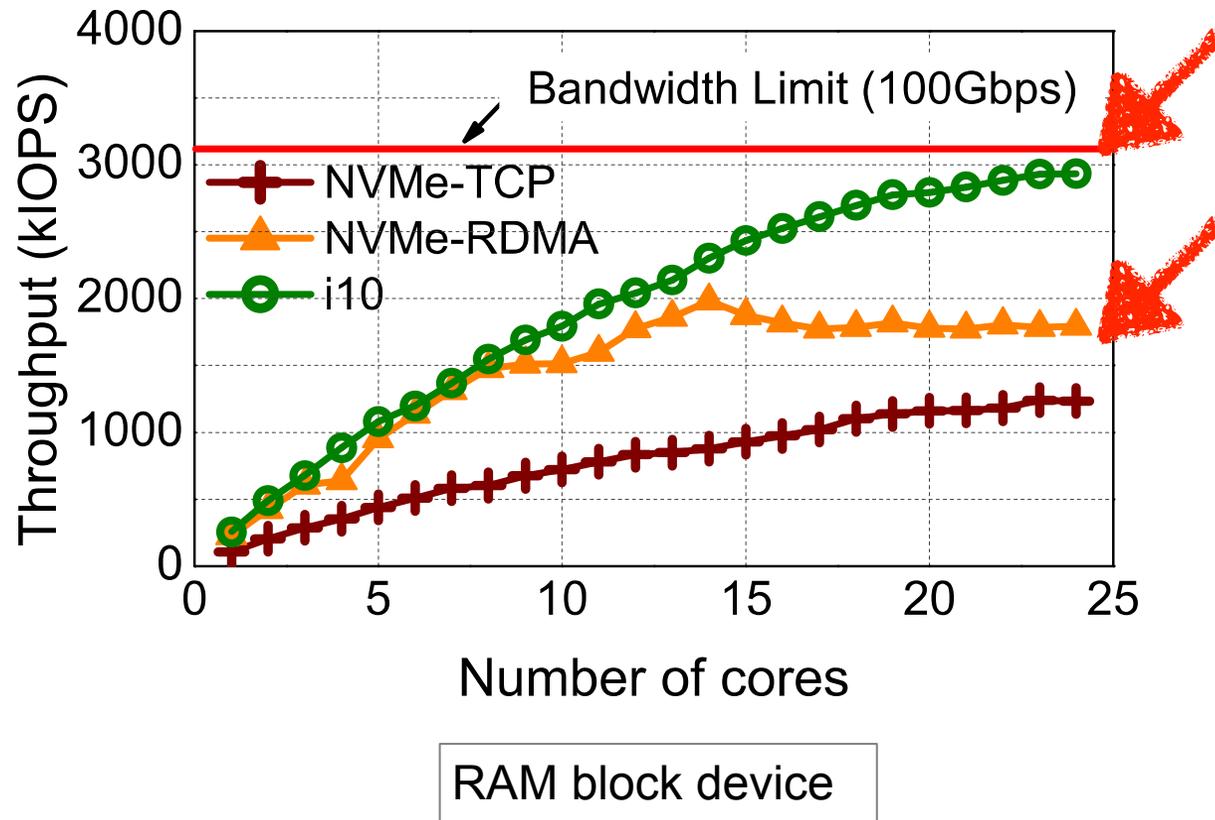


RAM block device

TCP ≈ RDMA:

- Throughput: **Comparable (or better)**
- Tail latency: **+97us**

Scalability with number of cores

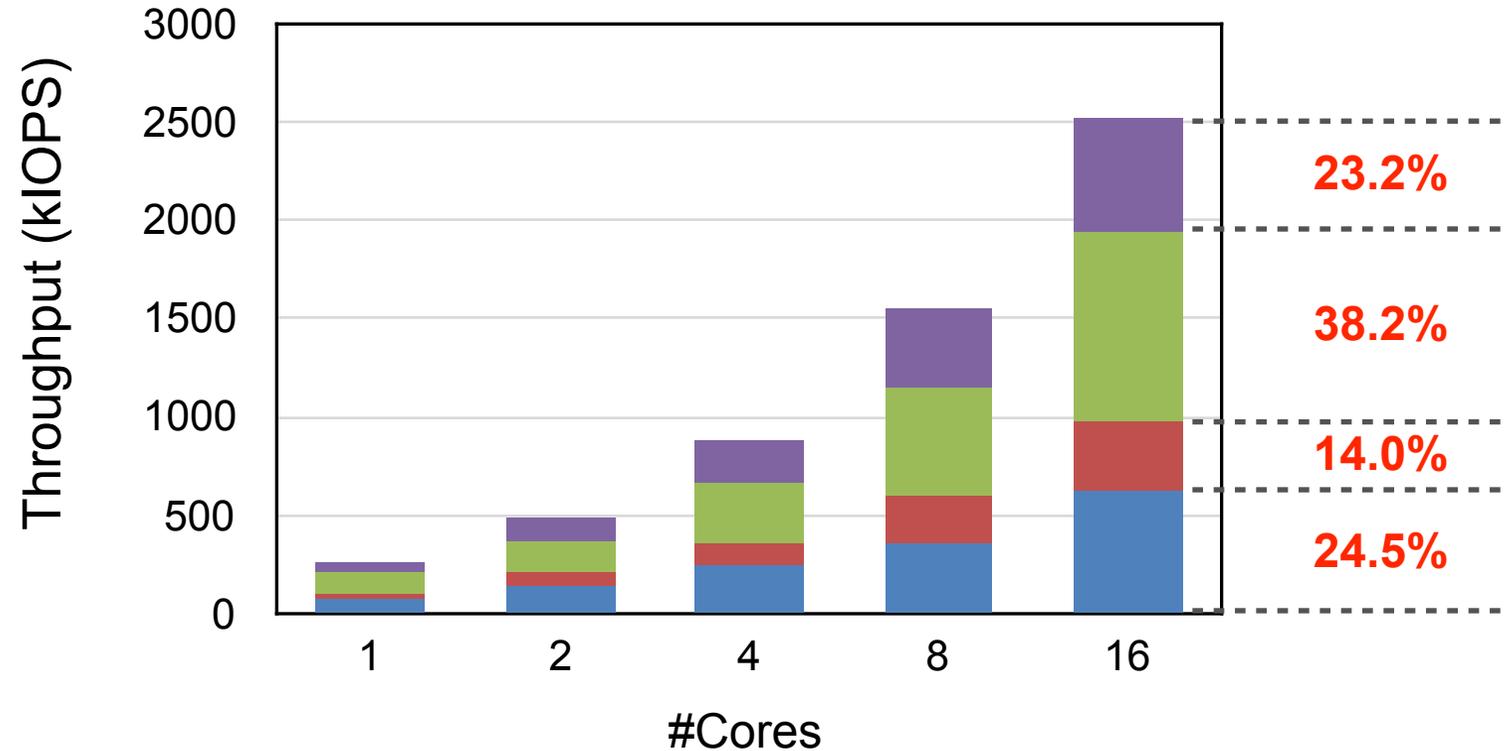


TCP ≈ RDMA:

- Throughput: **Scales similar (~14 cores) or better**
- Seems related to hardware scalability

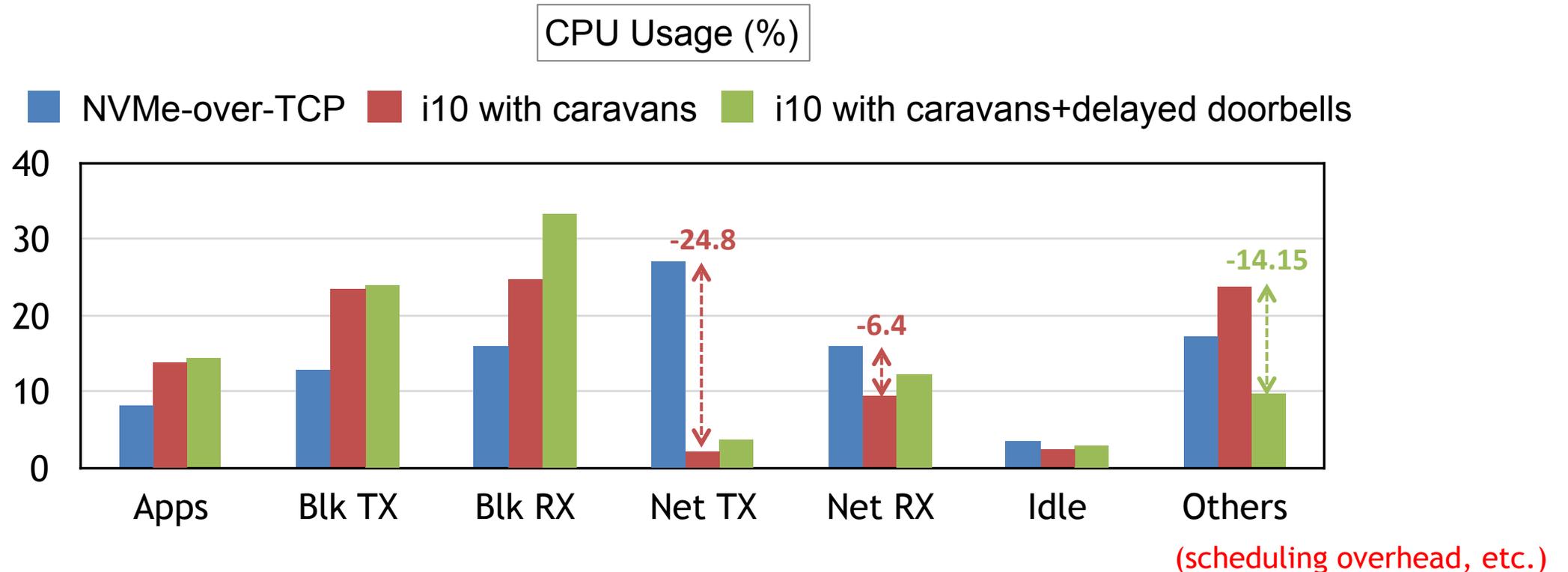
Benefits from individual design components

■ i10-lane ■ TSO/GRO + Jumbo ■ i10 Caravans ■ Delayed Doorbells



Each of the design component contributes to i10 performance

Understanding performance improvement



i10 improves over NVMe-over-TCP by using

Fewer cycles for network processing (Net Tx/Rx) and scheduling (Others)

More cycles for Applications, and block layer operations (Blk Tx/Rx)

i10 kernel implementation

Edit

Manage topics

6 commits

1 branch

Kernel implementation

releases

1 contributor

Branch: master

New pull request

Create new file

Upload files

Find file

Clone or download

Further evaluation

Test scripts ...

All are available at:

<https://github.com/i10-kernel/>

jaehyun-hwang Update README.md

Latest commit 0672ce3 6 days ago

drivers/nvme

7 days ago

include/linux

7 days ago

scripts

6 days ago

README.md

7 days ago

README.md

TCP \approx RDMA: CPU-efficient Remote Storage Access with i10

i10 is a new in-kernel remote storage I/O stack for high-performance network and storage hardware. Our i10 design offers a number of benefits: