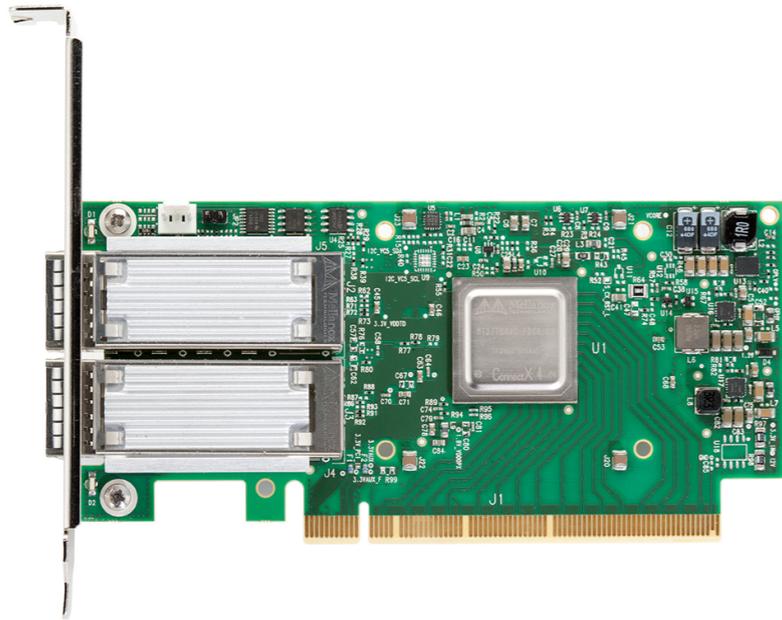# Design Guidelines
*for*
# High Performance RDMA Systems

**Anuj Kalia (CMU)**
Michael Kaminsky (Intel Labs)
David Andersen (CMU)

# RDMA is cheap (and fast!)
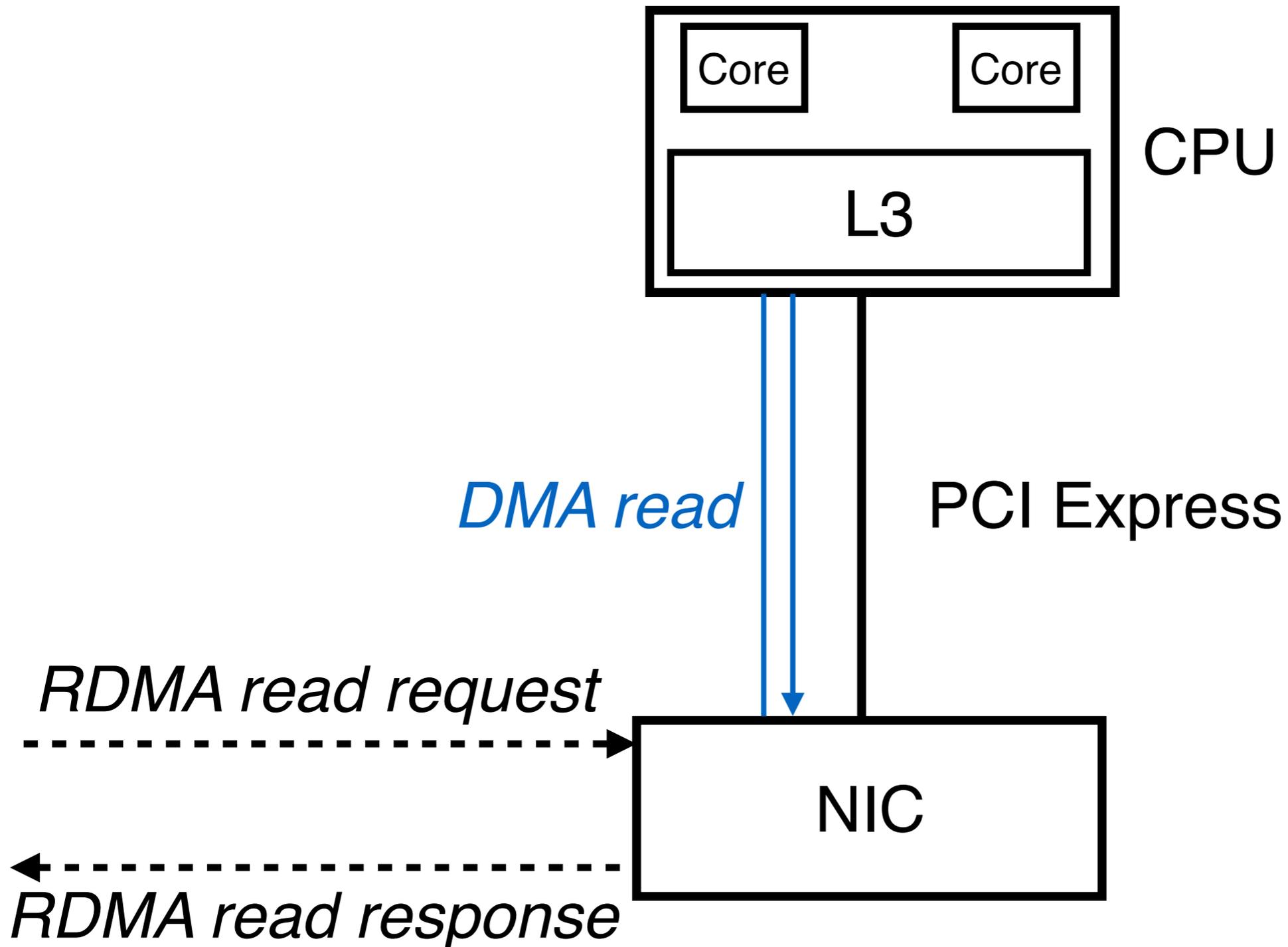
Mellanox Connect-IB
- 2x 56 Gbps InfiniBand
- ~2 $\mu$s RTT
- RDMA
- $1300

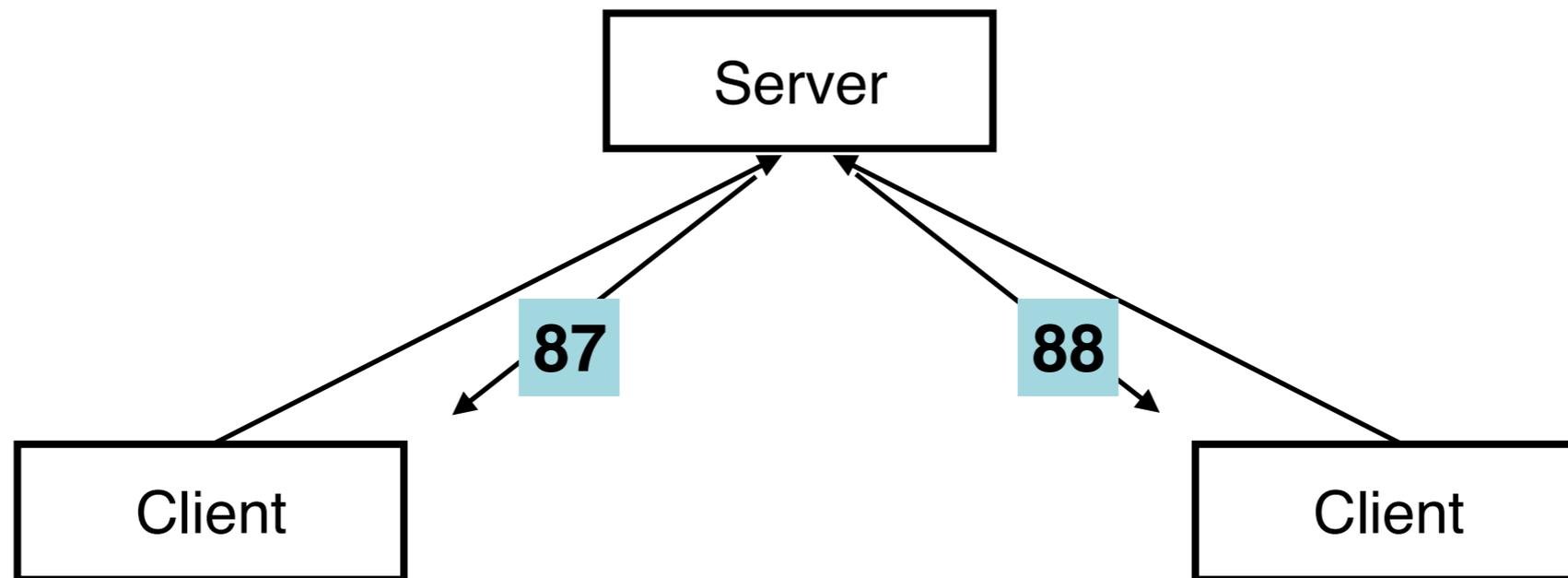**Problem**
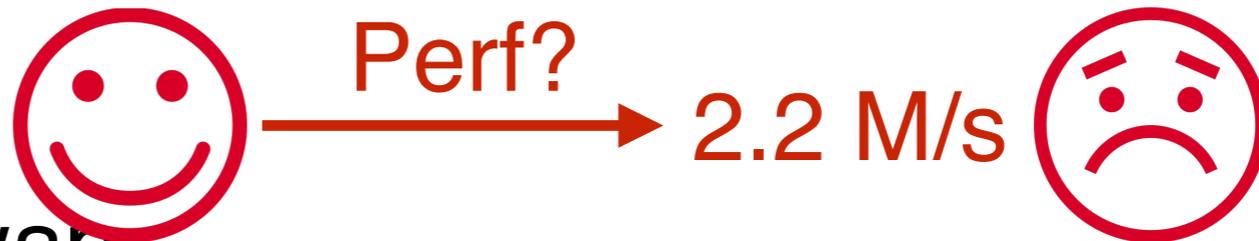Performance depends on complex low-level factors

# Background: RDMA read



*DMA read*

PCI Express

*RDMA read request*

NIC

*RDMA read response*

# How to design a sequencer?

```
                    ┌──────────┐
                    │  Server  │
                    └──────────┘
                   ↗↙          ↘↖
              87              88
          ↙                        ↘
   ┌──────────┐              ┌──────────┐
   │  Client  │              │  Client  │
   └──────────┘              └──────────┘
```

# Which RDMA ops to use?

Remote CPU bypass (one-sided)

- Read
- Write
- Fetch-and-add
- Compare-and-swap

Perf? → 2.2 M/s

Remote CPU involved (messaging, two-sided)

- Send
- Recv

# How we sped up the sequencer by 50X

# Large RDMA design space

**Operations**   | READ | WRITE | ATOMIC |      | SEND, RECV |

*Remote bypass (one-sided)*                     *Two-sided*

**Transports**   | Reliable |   | Unreliable |   | Connected |   | Datagram |

**Optimizations**   | Inlined |   | Unsignaled |   | Doorbell batching |

| WQE shrinking |   | 0B-RECVs |

# Guidelines

NICs have multiple processing units (PUs)

Avoid contention
Exploit parallelism

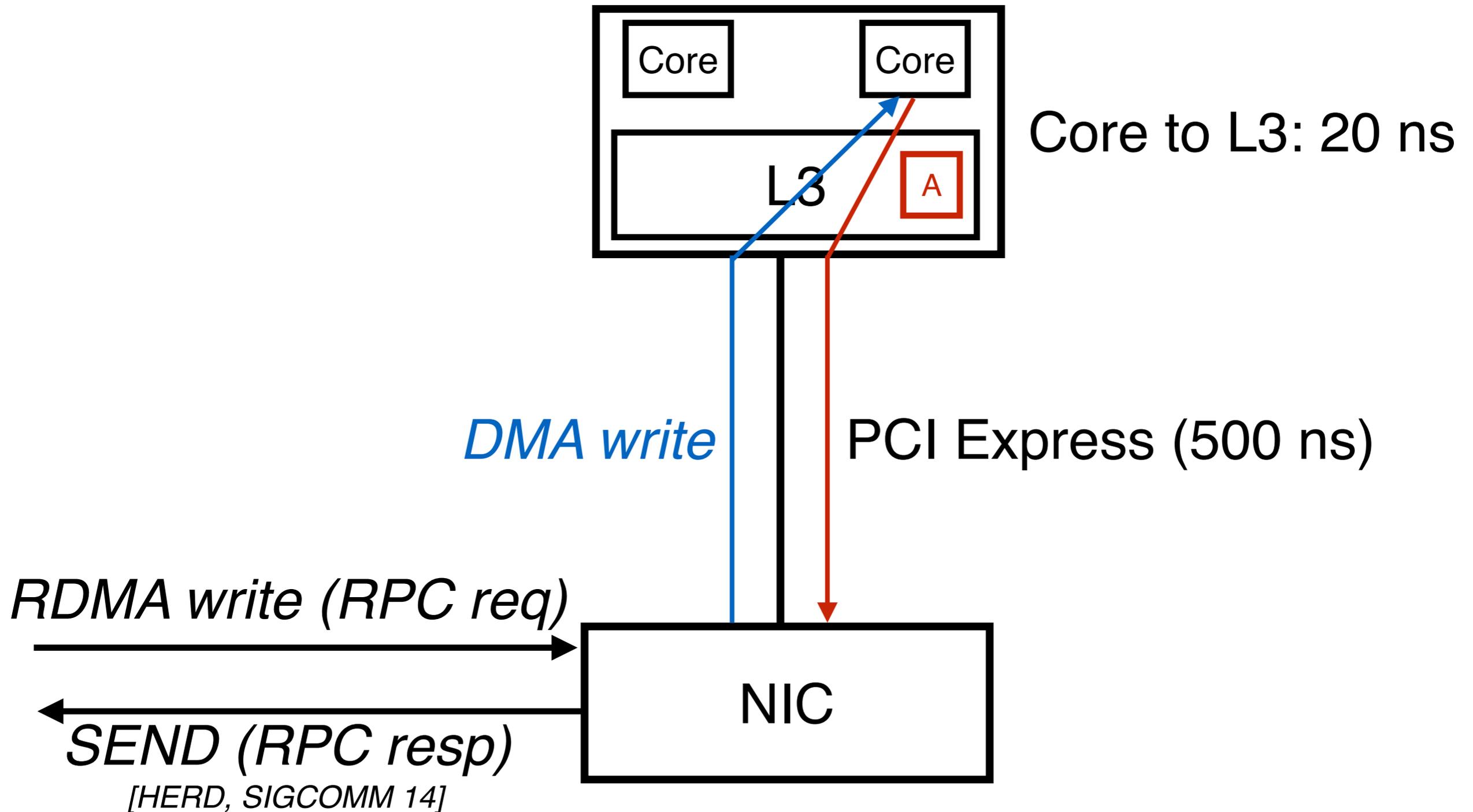PCI Express messages are expensive

Reduce CPU-to-NIC messages (MMIOs)
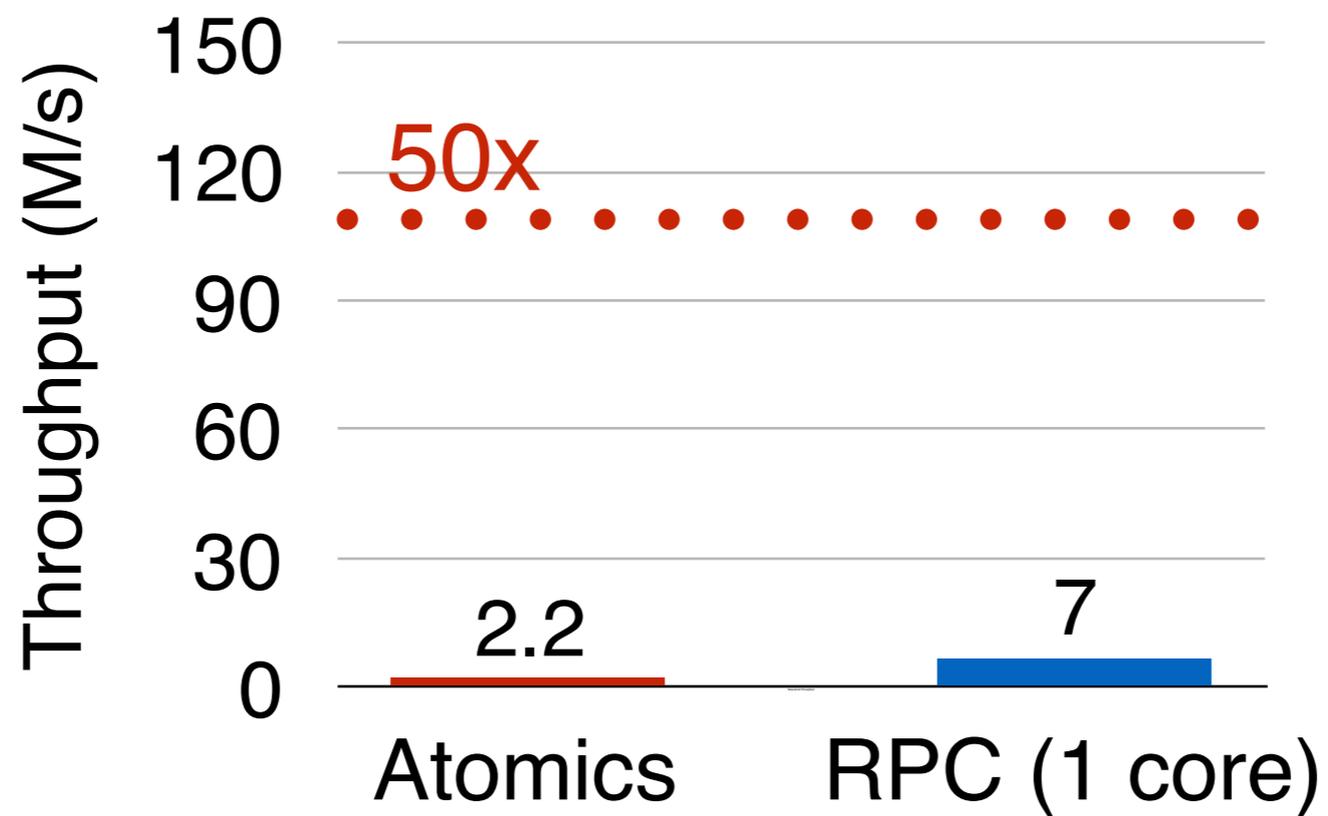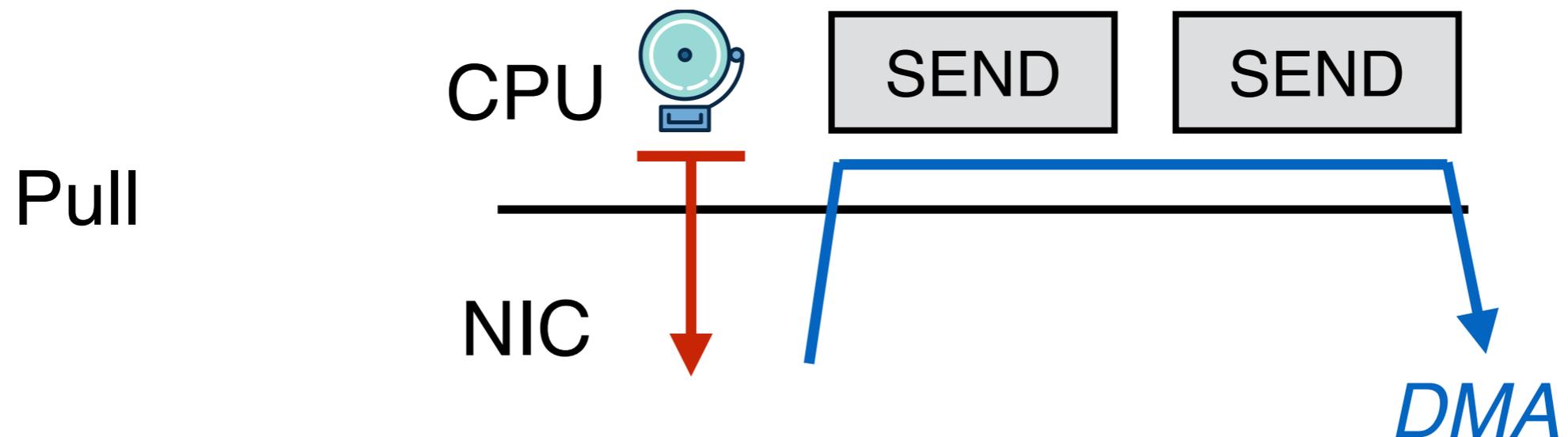Reduce NIC-to-CPU messages (DMAs)

# High contention w/ atomics

# Reduce contention: use CPU cores



Core to L3: 20 ns

DMA write

PCI Express (500 ns)

RDMA write (RPC req)

SEND (RPC resp)

[HERD, SIGCOMM 14]

# Sequencer throughput

# Reduce MMIOs w/ Doorbell batching

Push

CPU

SEND   SEND

NIC

MMIOs ⇒ lots of CPU cycles
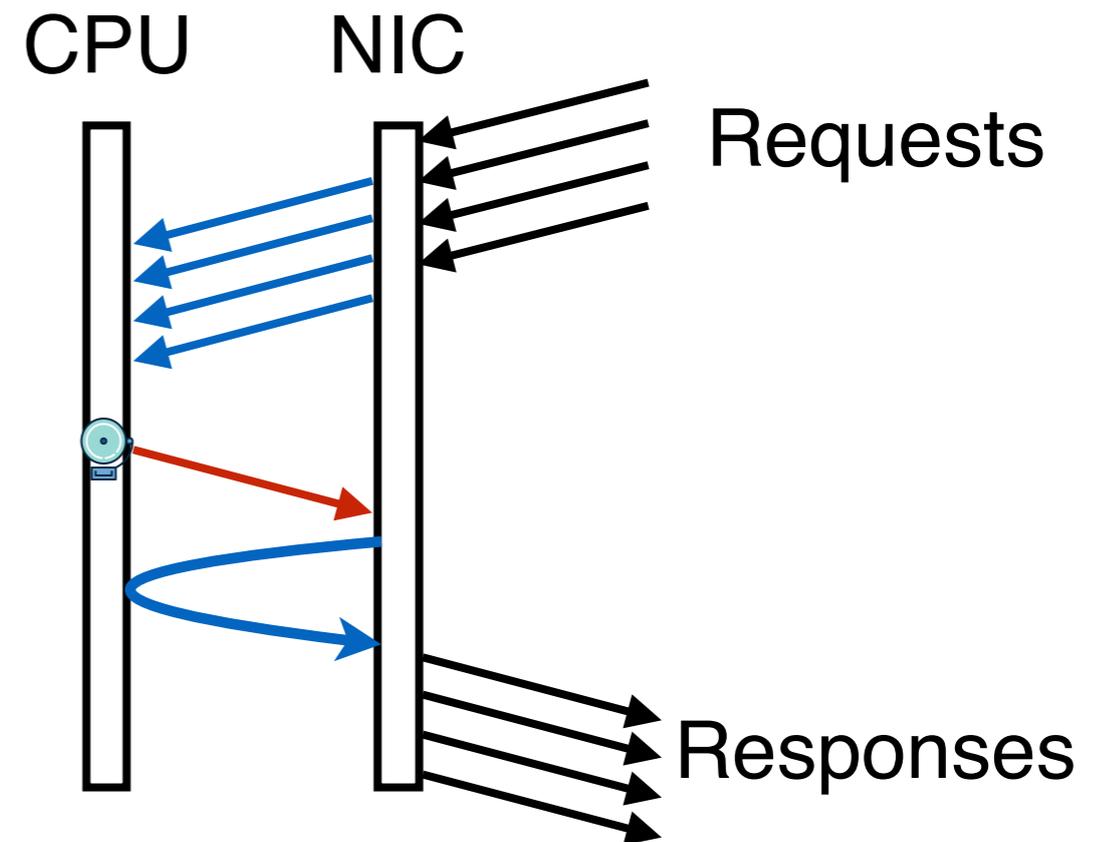
Pull

CPU

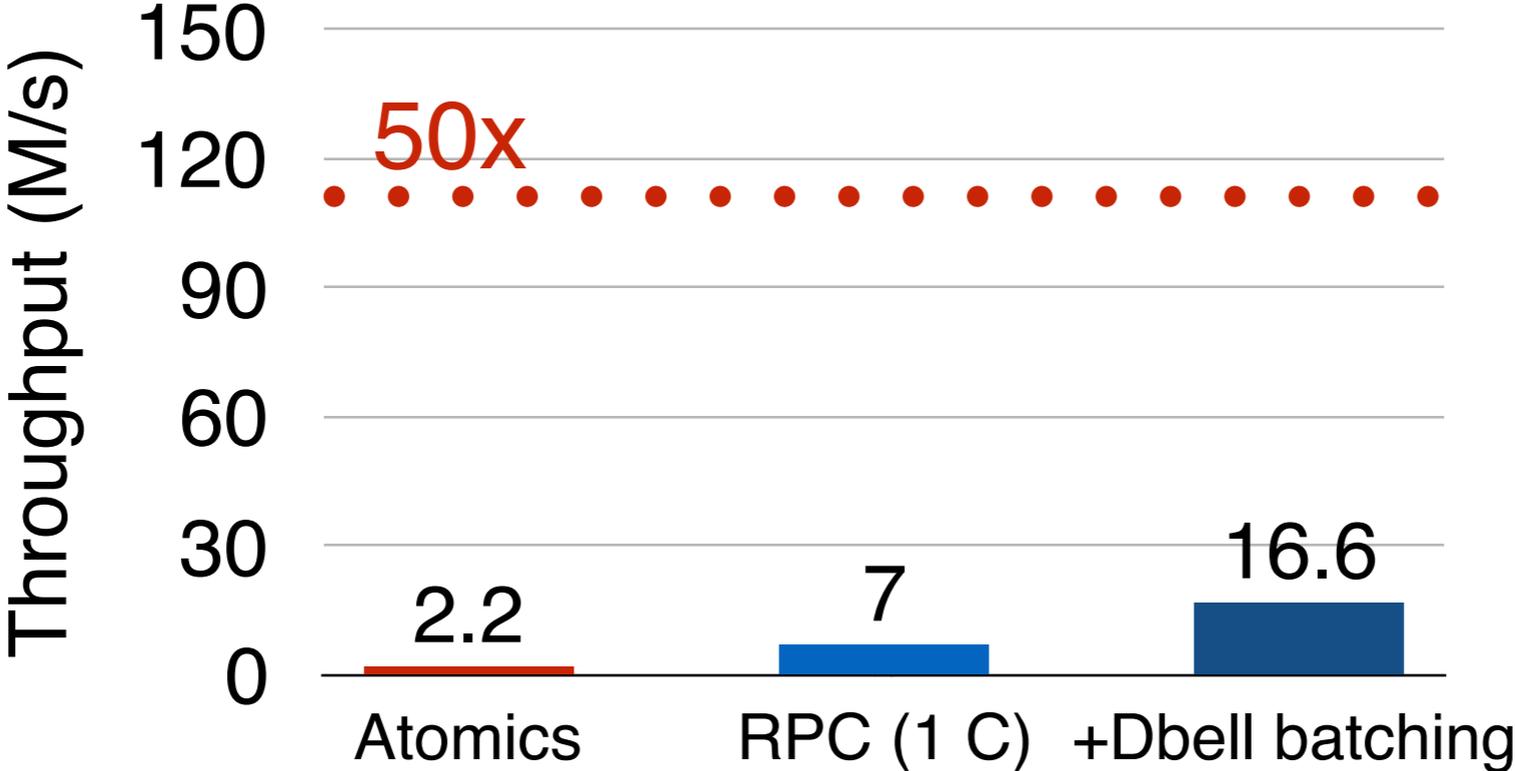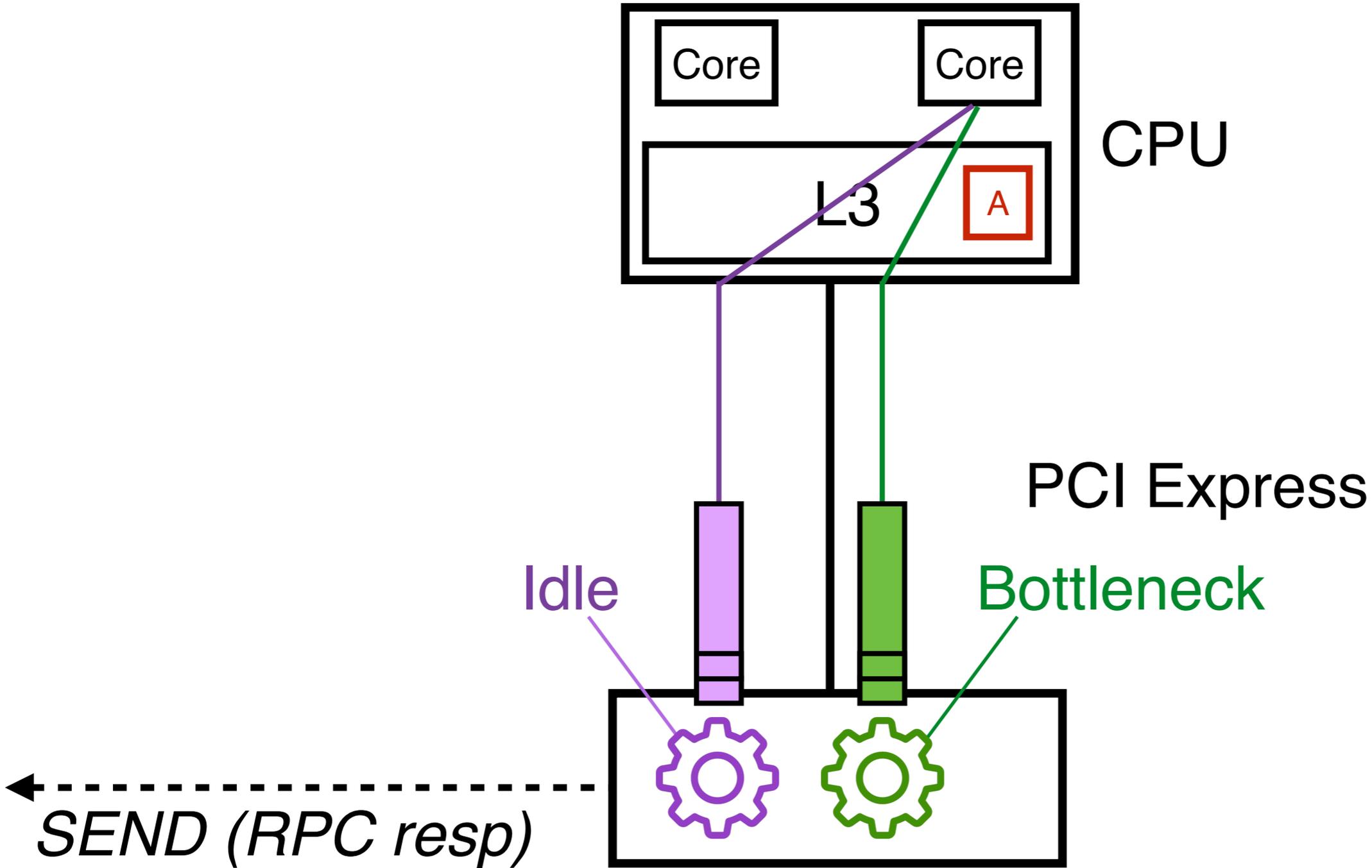SEND   SEND

NIC

*DMA*

# RPCs w/ Doorbell batching

**Push**

**Pull (Doorbell batching)**

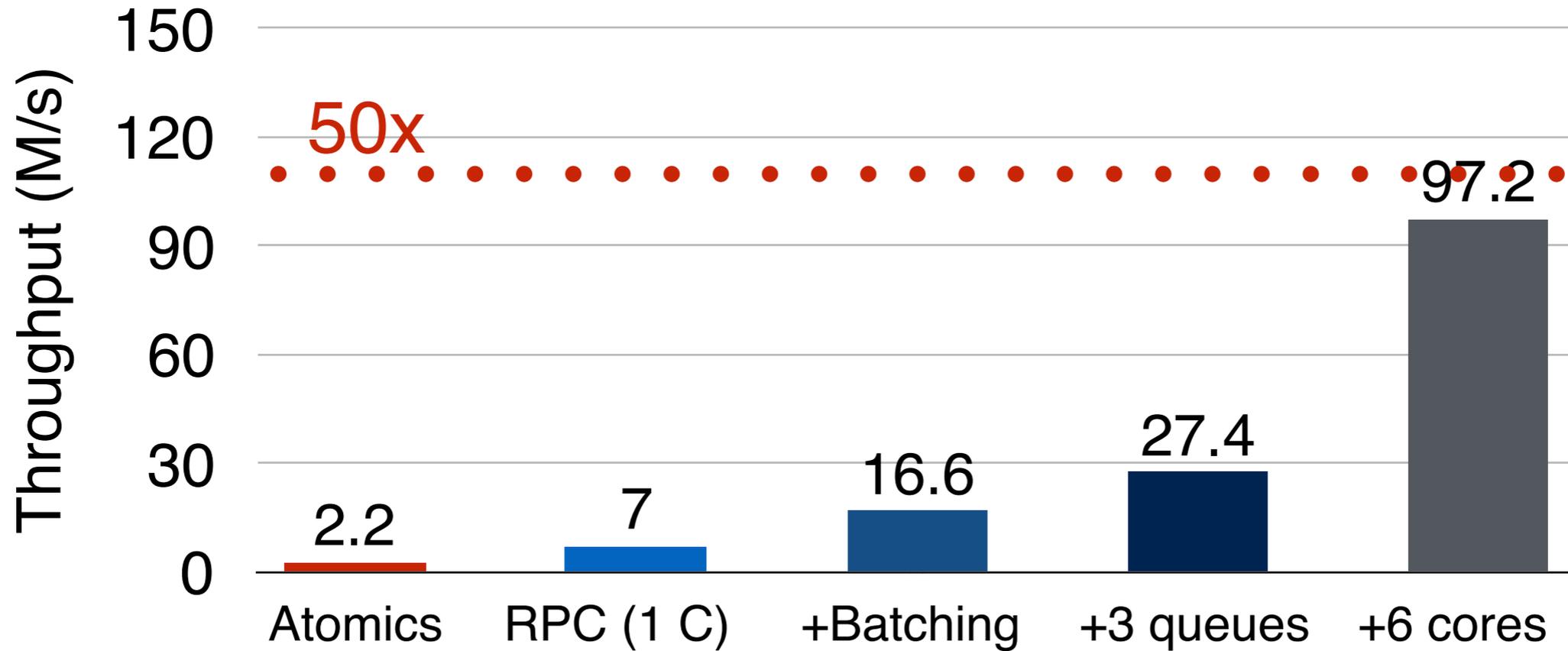Sequencer throughput

# Exploit NIC parallelism w/ multiQ



Core

Core

CPU

L3

A

PCI Express

Idle

Bottleneck

*SEND (RPC resp)*

# Sequencer throughput



Throughput (M/s)

150
120
90
60
30
0

50x

2.2
Atomics

7
RPC (1 C)

16.6
+Dbell batching

27.4
+3 queues

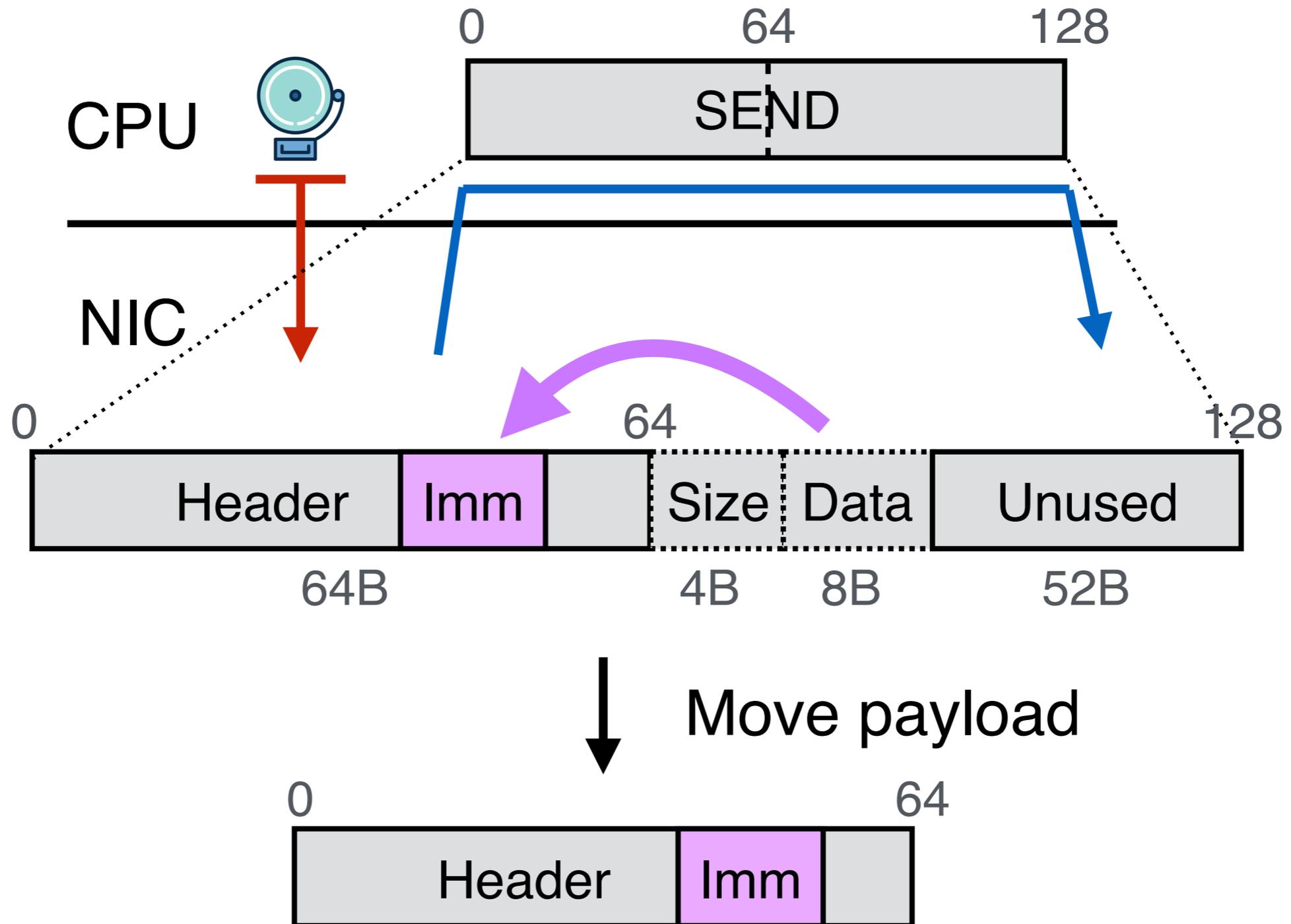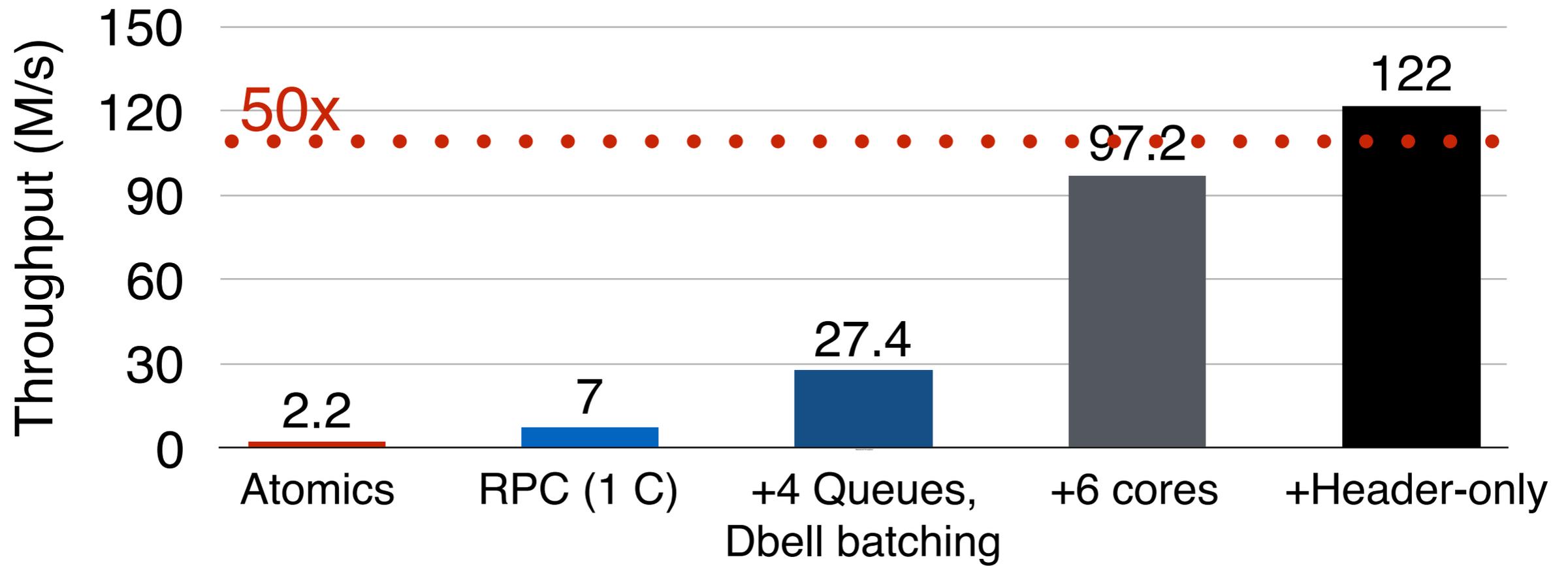# Sequencer throughput



**Bottleneck = PCIe DMA bandwidth (paper)**
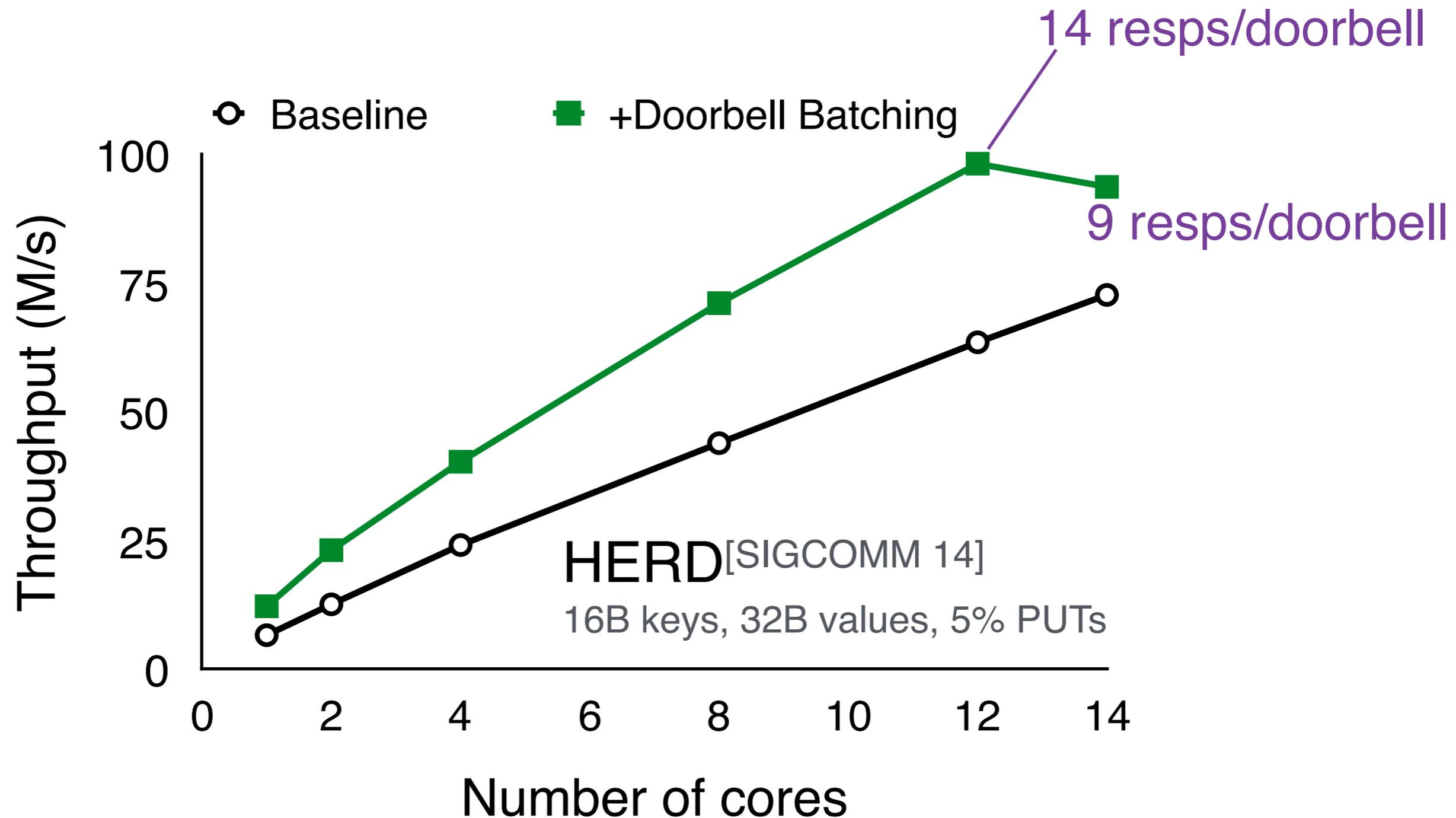
# Reduce DMA size: Header-only

# Sequencer throughput
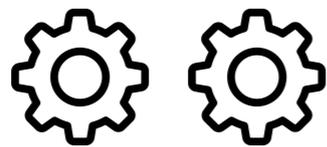
# Evaluation

- Evaluation of optimizations on 3 RDMA generations

- PCIe models, bottlenecks

- More atomics experiments

  - Example: atomic operations on multiple addresses

# RPC-based key-value store



**Figure:** Throughput (M/s) vs. Number of cores comparing Baseline (○) and +Doorbell Batching (■, green). Annotations: "14 resps/doorbell" at 12 cores (~99 M/s peak) and "9 resps/doorbell" at 14 cores. Caption in plot: HERD[SIGCOMM 14], 16B keys, 32B values, 5% PUTs.

# Conclusion

NICs have multiple processing units (PUs)

⚙ ⚙

Avoid contention
Exploit parallelism

PCI Express messages are expensive

Reduce CPU-to-NIC messages (MMIOs)
Reduce NIC-to-CPU messages (DMAs)

Code: https://github.com/anujkaliaiitd/rdma_bench