# A DBMS-centric Evaluation of BlueField DPUs on Fast Networks
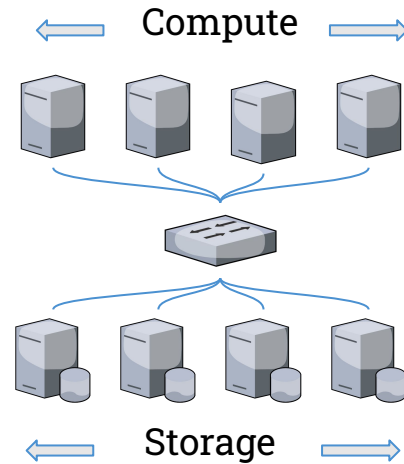
Technical University of Darmstadt

**Lasse Thostrup,** Daniel Failing, Tobias Ziegler & Carsten Binnig

DM
DATA MANAGEMENT
TU DARMSTADT

1

# State of DBMSs

- Target **disaggregated** setups

- Scale compute & memory independently

- Networks are increasingly on the **hot-path**

- Fast networks & RDMA provides **state-of-the-art** scale-out performance

Compute

Storage

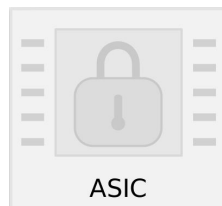# Programmable Networks

**Devices**

SmartNICs

Data Processing Units (DPU)
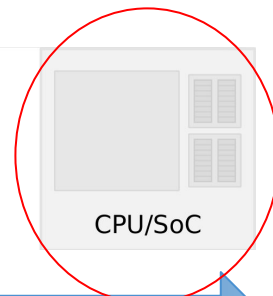
Programmable Switches

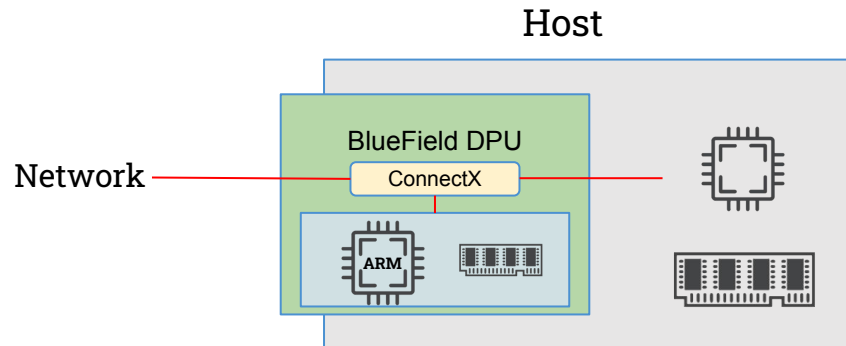P4

**Programmability & Processing**

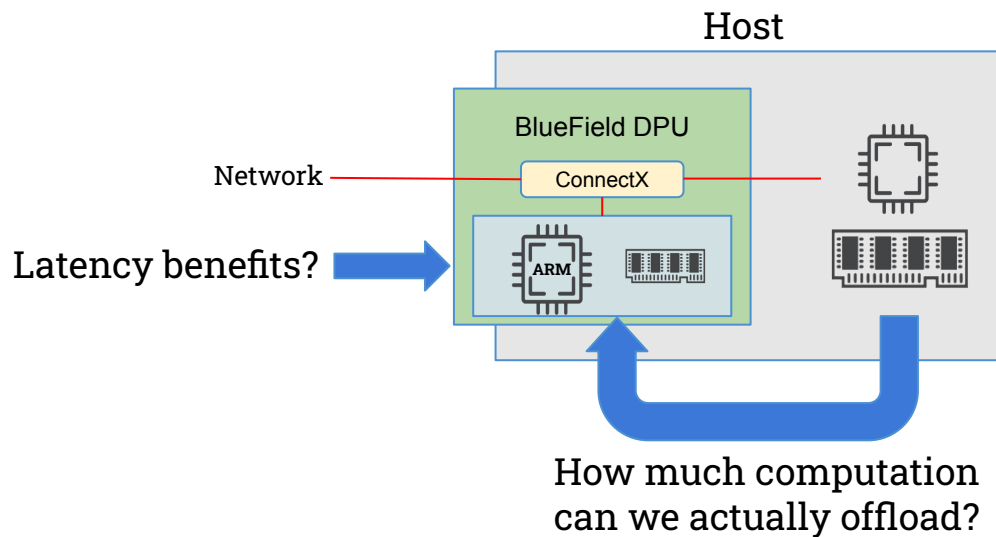ASIC

FPGA

CPU/SoC

Efficiency ← → Flexibility

# Nvidia BlueField DPU

- **Typical applications**: networking, storage, security

- Equipped with various **hardware accelerators**

- Embedded **ARM** subsystem:
  - 8 cores ARM A-72 CPU
  - 16 GB DRAM

# Motivation

Host

BlueField DPU

Network

ConnectX

ARM

Latency benefits?

How much computation
can we actually offload?
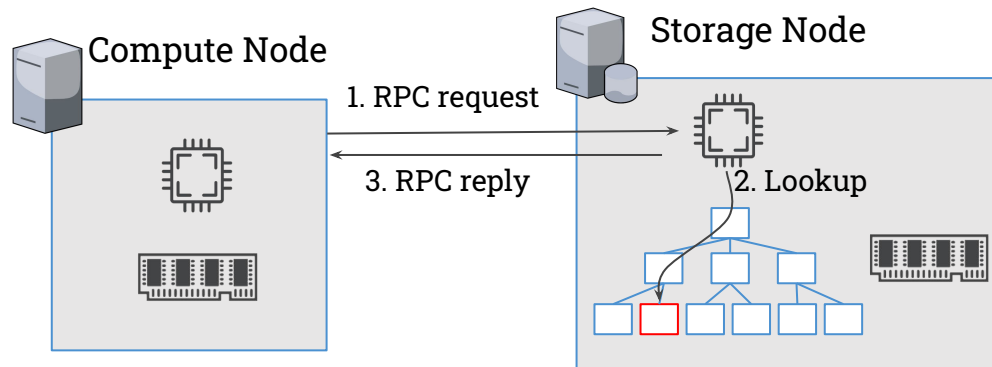
**We consider two use-cases:**
- Remote B-tree
- Remote Sequencer

**Expectations**:

- Added computational power and memory yields faster processing
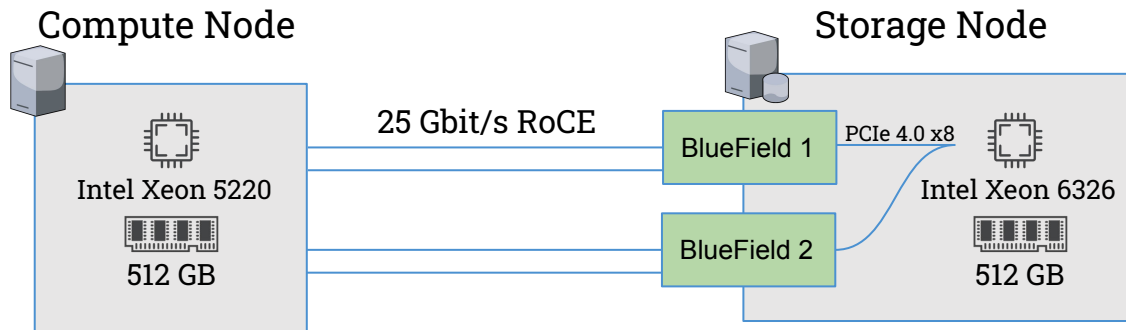- Closer proximity to the network yields lower latency

# Remote B-tree

- Main data-structure in OLTP databases

- Access B-tree through Remote Procedure Calls (RPC)

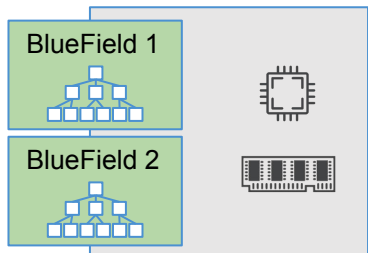- Efficient RPC-framework with state-of-the-art optimizations

# Experimental Setup

- Mirrors disaggregated memory setup

- 50 Gbit/s between Compute Node & BlueFields
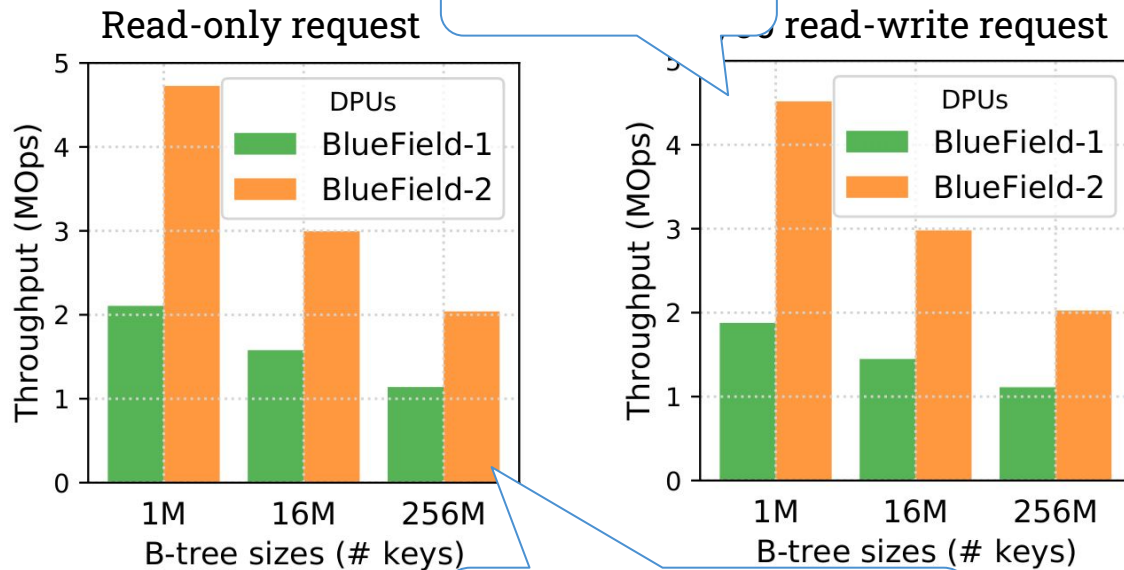
- Use 8 threads on Storage Node & BlueFields

# Remote B-tree w/ RPC - Throughput

# Remote B-tree w/ RPC - Throughput

# Remote B-tree w/ RPC - Latency

# RDMA Latency

## RDMA latency (*ib_send_lat*)



- ~0.4 µs lower latency for the BlueField without inlining
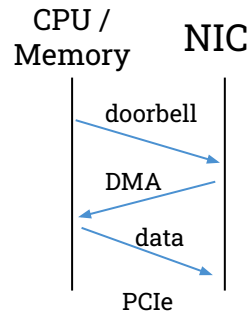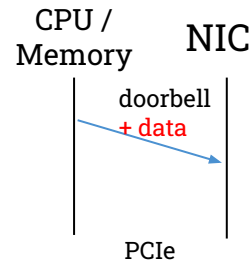
- Inlining reduces the latency gap (enabled by default)

- With 64 B messages, practically no latency difference
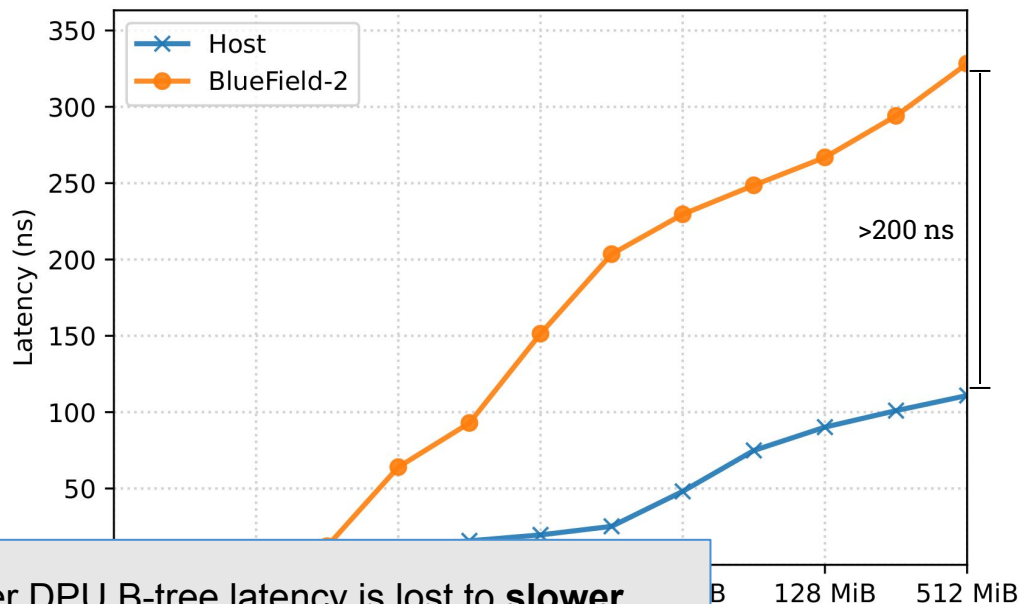
**Without inlining**          **With inlining**

# Memory Latency

- Random reads in increasing memory block sizes

- Different cache sizes:
  - BlueField LLC: **6 MB**
  - Host LLC: **24 MB**

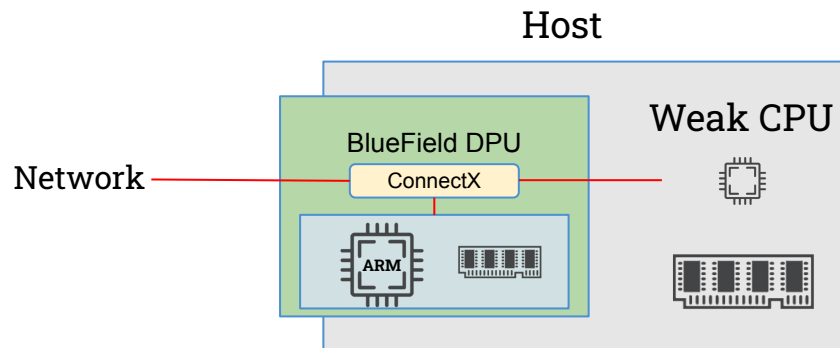- BlueField-2 inhibits much higher latency when spilling out of the cache

### Memory latency (*tinymembench*)



Legend:
- Host
- BlueField-2

y-axis: Latency (ns) — 50, 100, 150, 200, 250, 300, 350

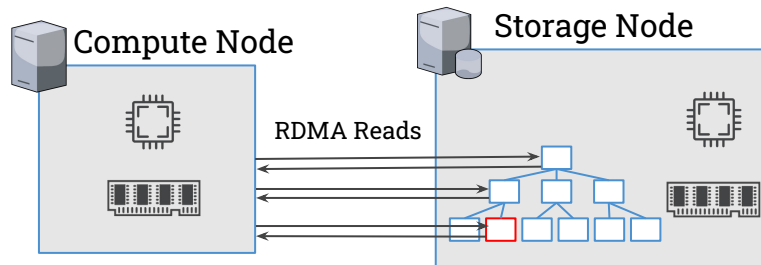x-axis: B, 128 MiB, 512 MiB

>200 ns

The **expected** lower DPU B-tree latency is lost to **slower memory access** in B-tree lookups and a similar network latency.

# Near-zero Computation

- Storage servers typically employed with very **little compute power**

- RPC heavily involves remote CPU → Instead: **One-sided RDMA**

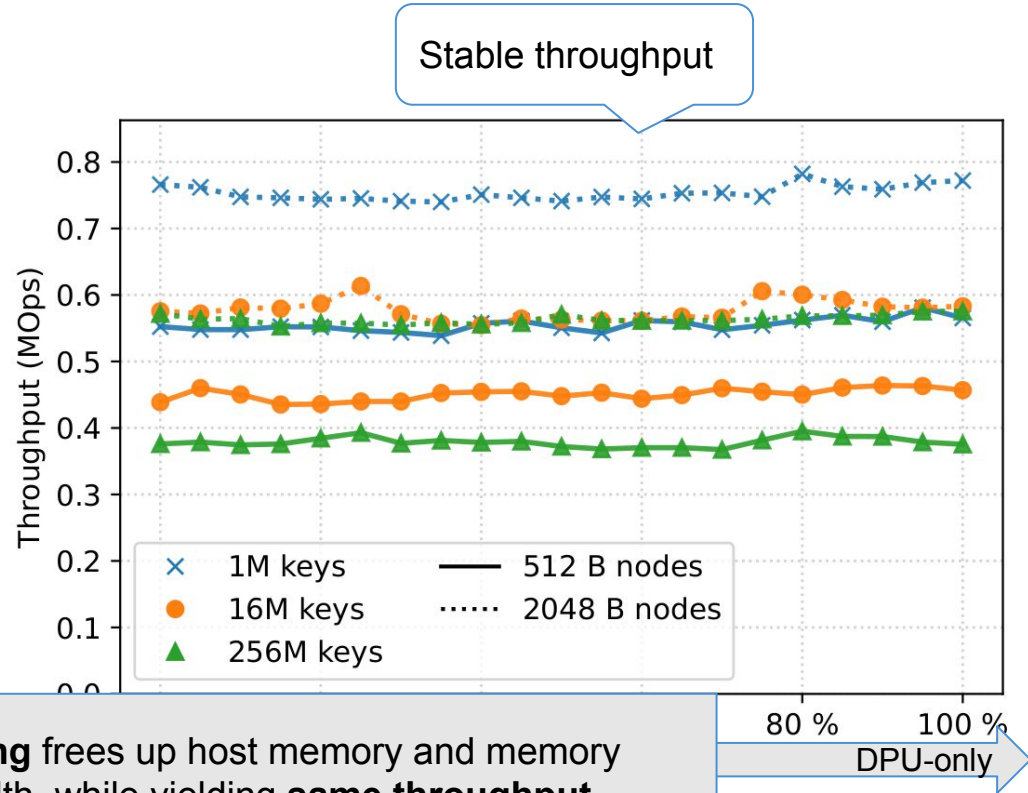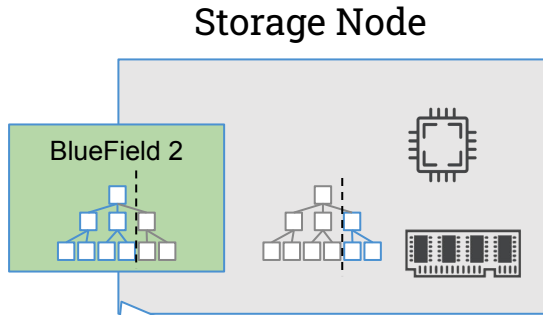- One-sided RDMA accesses have **no remote CPU load**
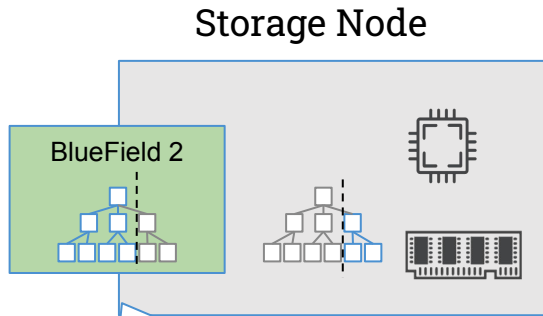
# Remote B-tree w/ One-sided RDMA



- Performs the binary search for child on Compute Node

- A network round-trip for each level

- **No CPU load** on Storage Node or BlueField
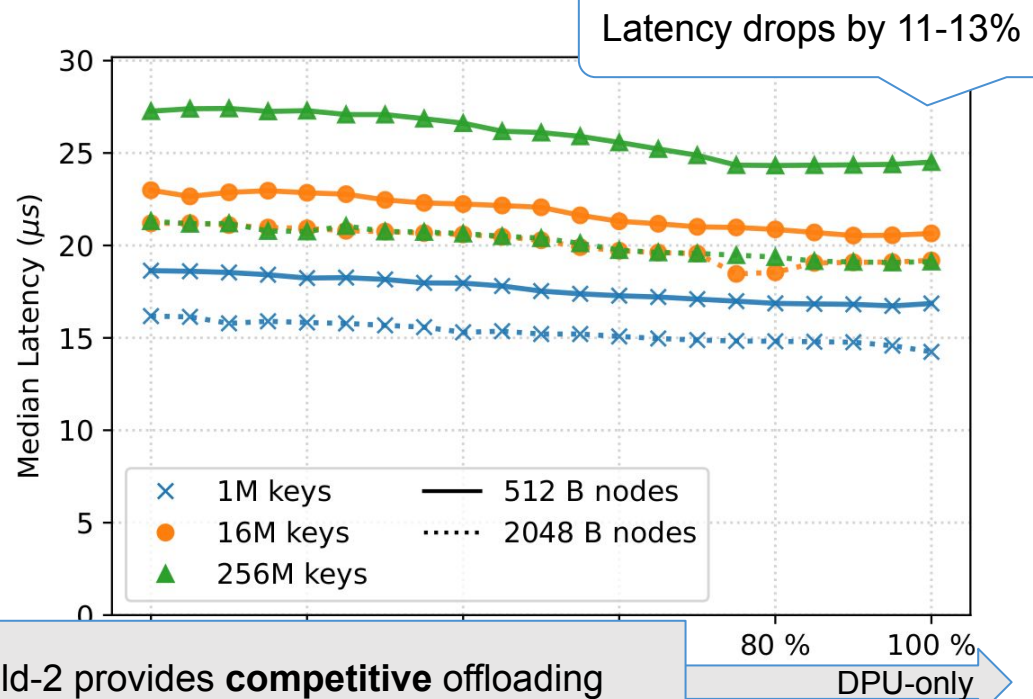
# Use case 1: Remote B-tree w/ One-sided RDMA

# Use case 1: Remote B-tree w/ One-sided RDMA

# Use case 2: Remote Sequencer w/ One-sided RDMA
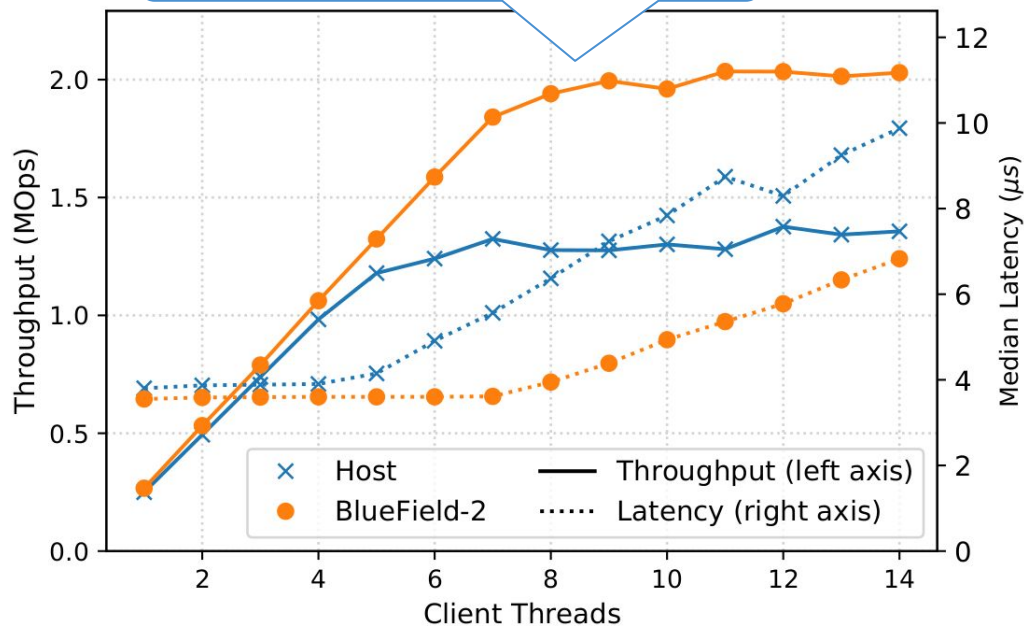


Storage Node

BlueField 2

What is the **throughput** and **latency** when incrementing a remote counter with RDMA fetch & add?

Almost 50% throughput increase

- × Host
- ● BlueField-2
- —— Throughput (left axis)
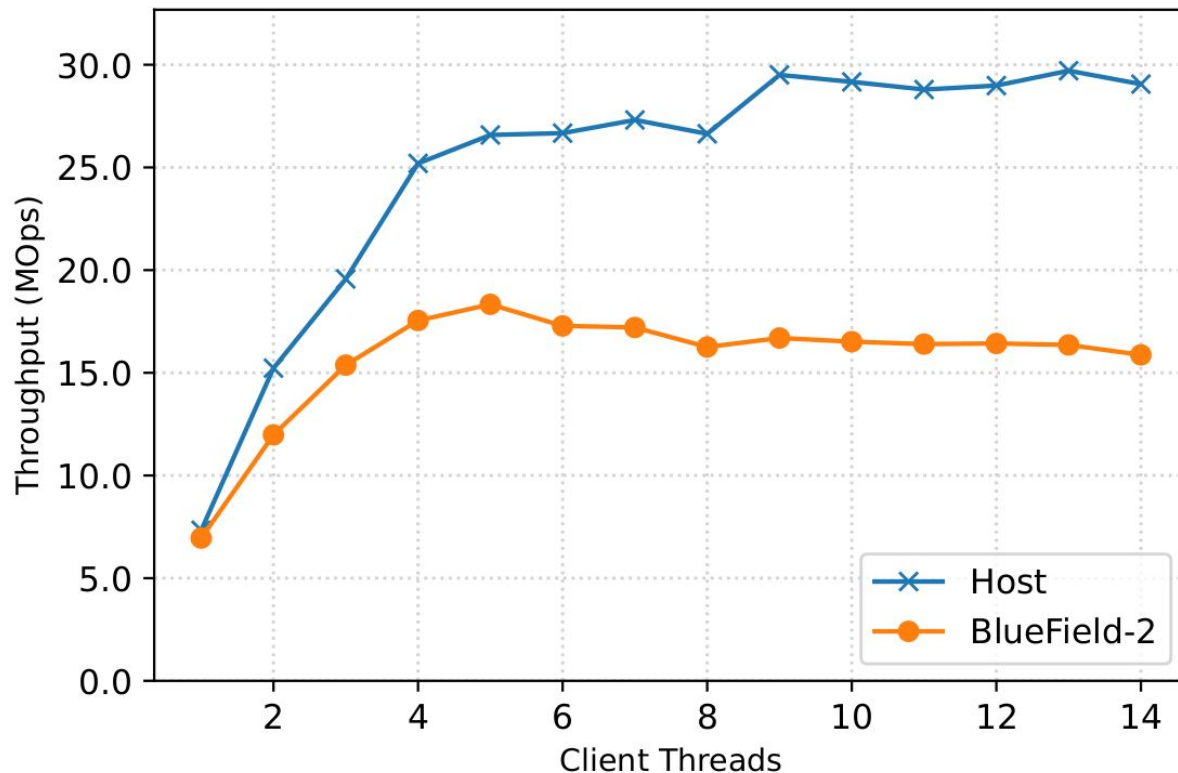- ⋯⋯ Latency (right axis)

# Conclusion

The BlueField DPU can **accelerate** typical DBMS operations

- Main benefits shown for one-sided accesses

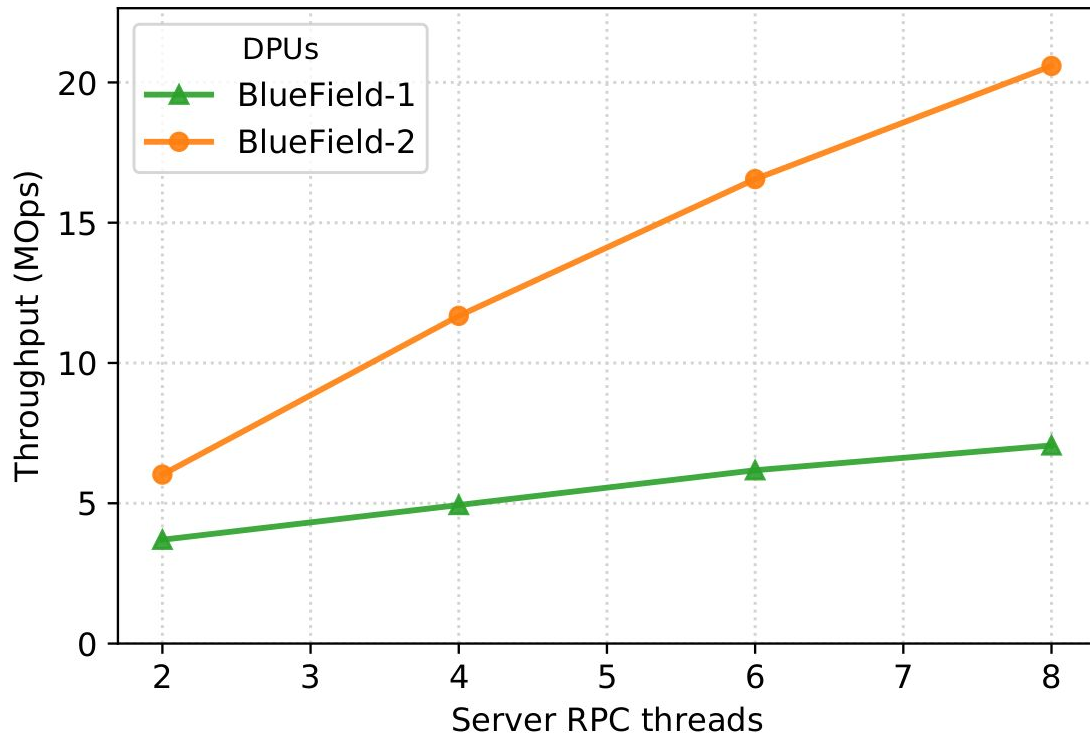- Two-sided operations can provide speed-ups when used in union with host

**Future work**

- Explore other BlueField hardware accelerators (compression, encryption & NVMe)
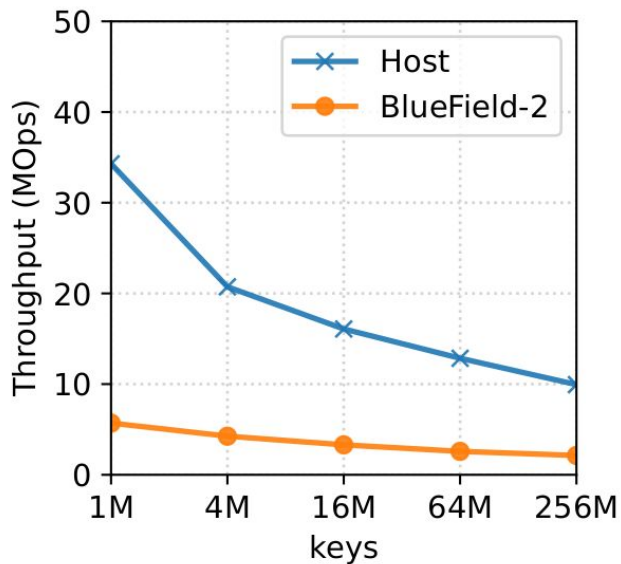
- BlueField 3 → faster CPU & memory
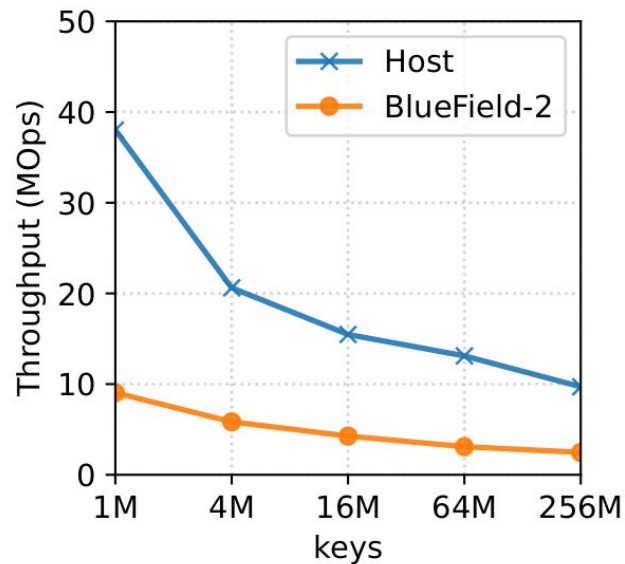
# Backup: Throughput of RPC sequencer

# Backup: No-Op RPC BlueField-1 vs BlueField-2

# Backup: Throughput of local B-tree



(a) Updates

(b) Lookups