# *Towards "Intelligence, Storage, Network": Characterizing, Optimizing, and Outlooking*

## — A Systems Researcher's Perspective

Rong Chen

Institute of Parallel and Distributed Systems, SJTU

Huawei STW, 2023

*Joint work with Xingda, Xiating, Rongxin, Yuhan, Haibo, Binyu, and members of IPADS*

# Who AM I

Rong Chen (陈榕) / IPADS, SJTU

https://ipads.se.sjtu.edu.cn/rong_chen

► Research Interest: Building efficient, scalable, and reliable distributed systems

► Publications and awards in systems conferences (OSDI, SOSP, EuroSys)

► Huawei OlympusMons Pioneer Award, 2020

"Efficient Data Processing System based on New Heterogeneous Hardware"

# Disclaimers:
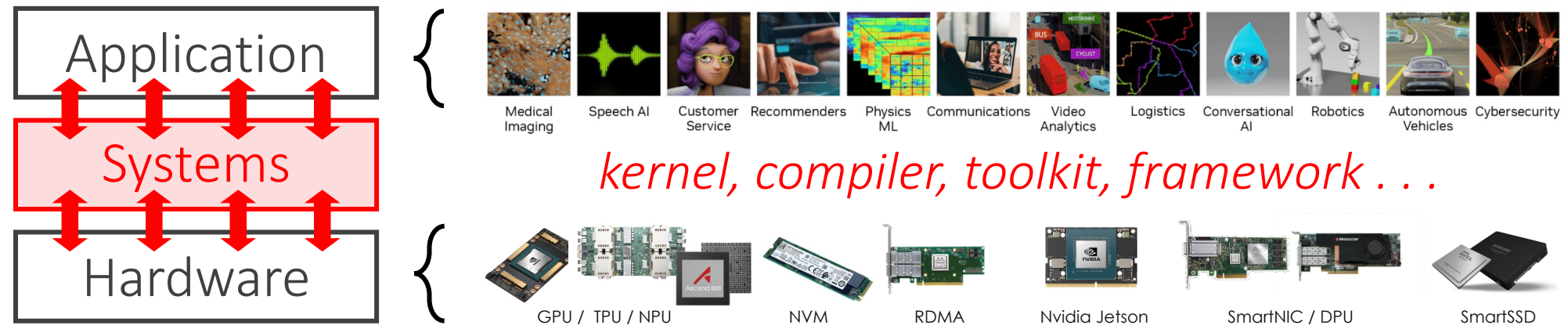
I am a Systems person, not a Network/Storage/AI expert ☺

# My view:

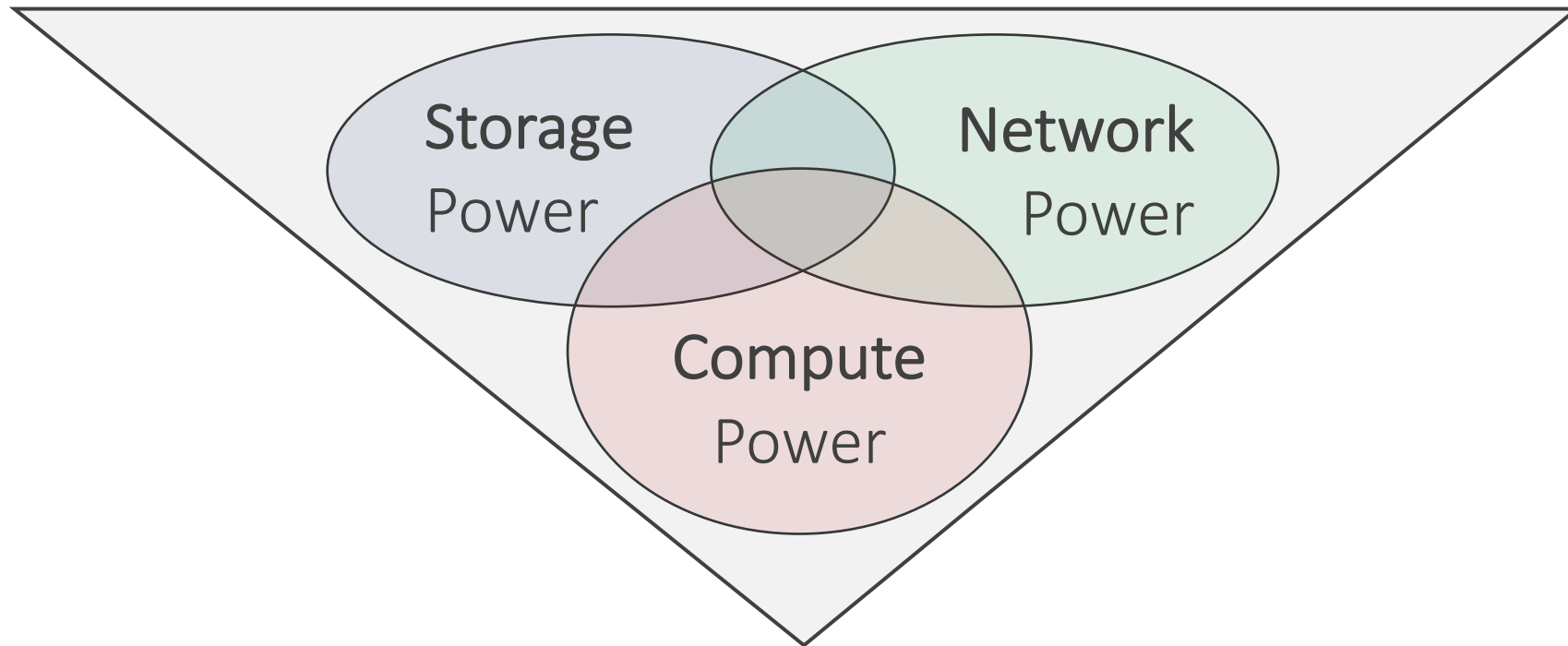How to (re)build high-performance system software stack
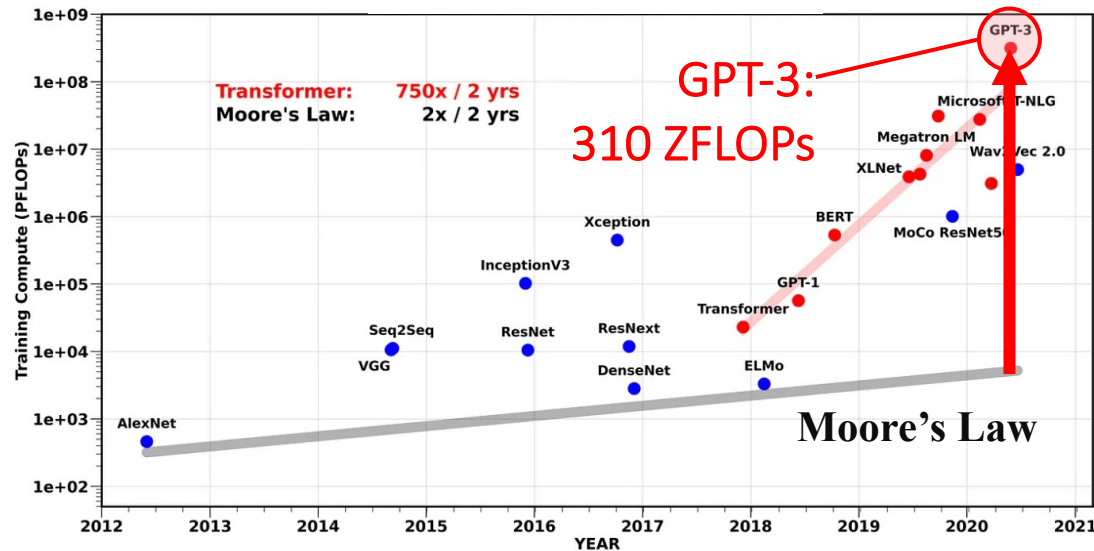by exploiting new hardware of "Intelligence, Storage, Network"
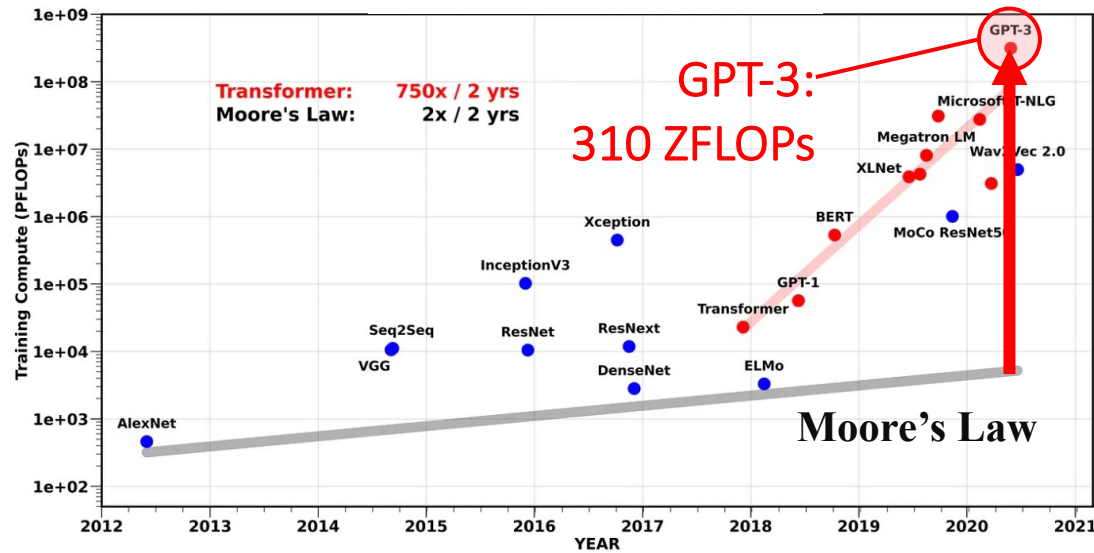


| Application |
| Systems |
| Hardware |

*kernel, compiler, toolkit, framework . . .*

Medical Imaging · Speech AI · Customer Service · Recommenders · Physics ML · Communications · Video Analytics · Logistics · Conversational AI · Robotics · Autonomous Vehicles · Cybersecurity

GPU / TPU / NPU · NVM · RDMA · Nvidia Jetson · SmartNIC / DPU · SmartSSD

# Application Demands

# Application Demands

## Compute Power



Source: "AI and Memory Wall", 2021. https://medium.com/riselab/ai-and-memory-wall-2cb4265cb0b8

# Application Demands

## Compute Power



Transformer: 750x / 2 yrs
Moore's Law: 2x / 2 yrs

GPT-3:
310 ZFLOPs

Moore's Law

## Storage Power



Transformer Size: 240x / 2 yrs
AI HW Memory: 2x / 2 yrs

Accelerator Memory

RecSys:
10 T Params

Source: "AI and Memory Wall", 2021. https://medium.com/riselab/ai-and-memory-wall-2cb4265cb0b8

# Application Demands



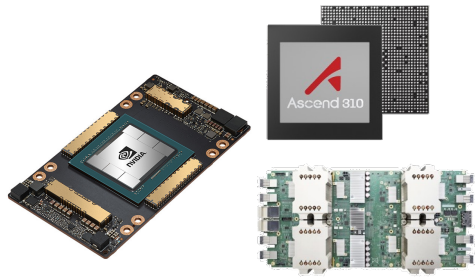Source: "AI and Memory Wall", 2021. https://medium.com/riselab/ai-and-memory-wall-2cb4265cb0b8

# My view:
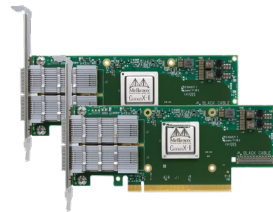
How to (re)build high-performance system software stack

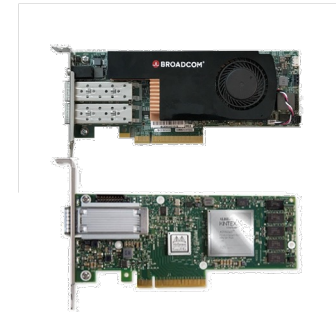by exploiting new hardware of "Intelligence, Storage, Network"
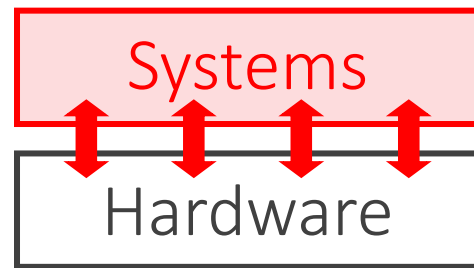
GPU / TPU / NPU     NVM     RDMA     DPU / SmartNIC     SmartSSD     Nvidia Jetson
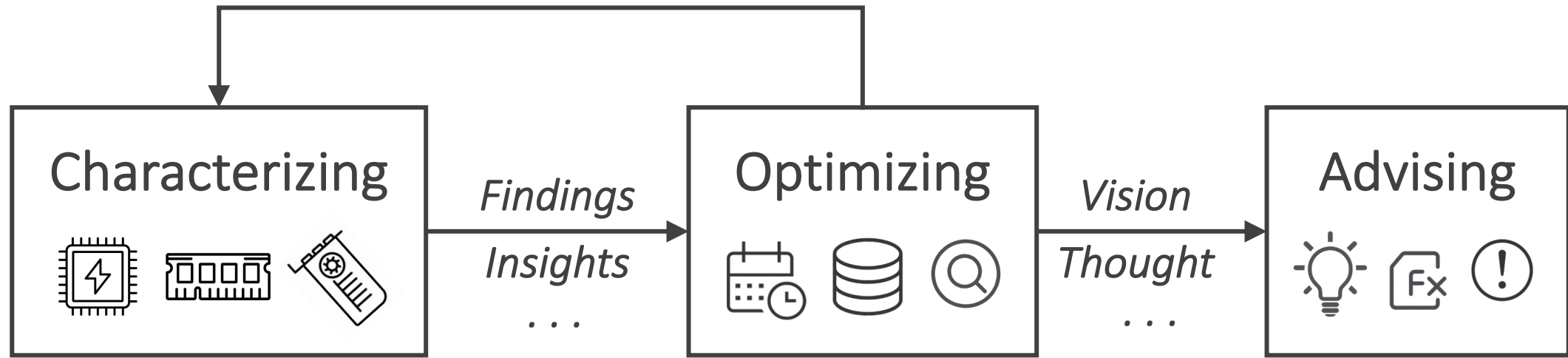
# Our Approach

Characterizing

*Findings*

*Insights*

. . .

Optimizing

*Vision*

*Thought*

. . .

Advising

Systems

Hardware

*kernel, compiler, toolkit, framework . . .*

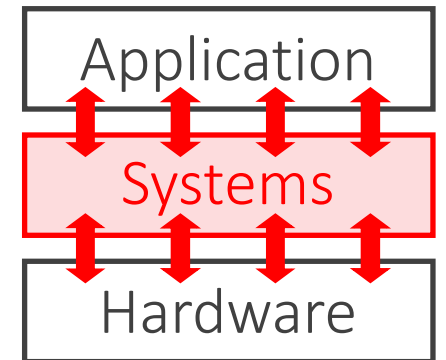GPU / TPU / NPU      NVM      RDMA      Nvidia Jetson      SmartNIC / DPU      SmartSSD

# Outline

Case #1: Collaborative offloading

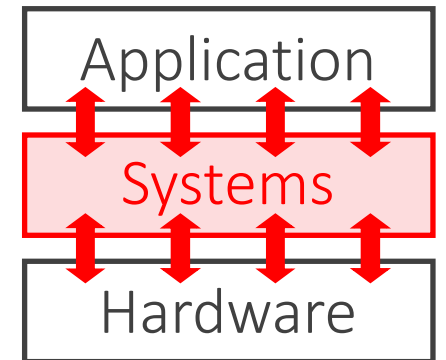Case #2: Cooperative offloading

Outlooking systems research for DPU

# Outline

**Case #1: Collaborative offloading**

Case #2: Cooperative offloading
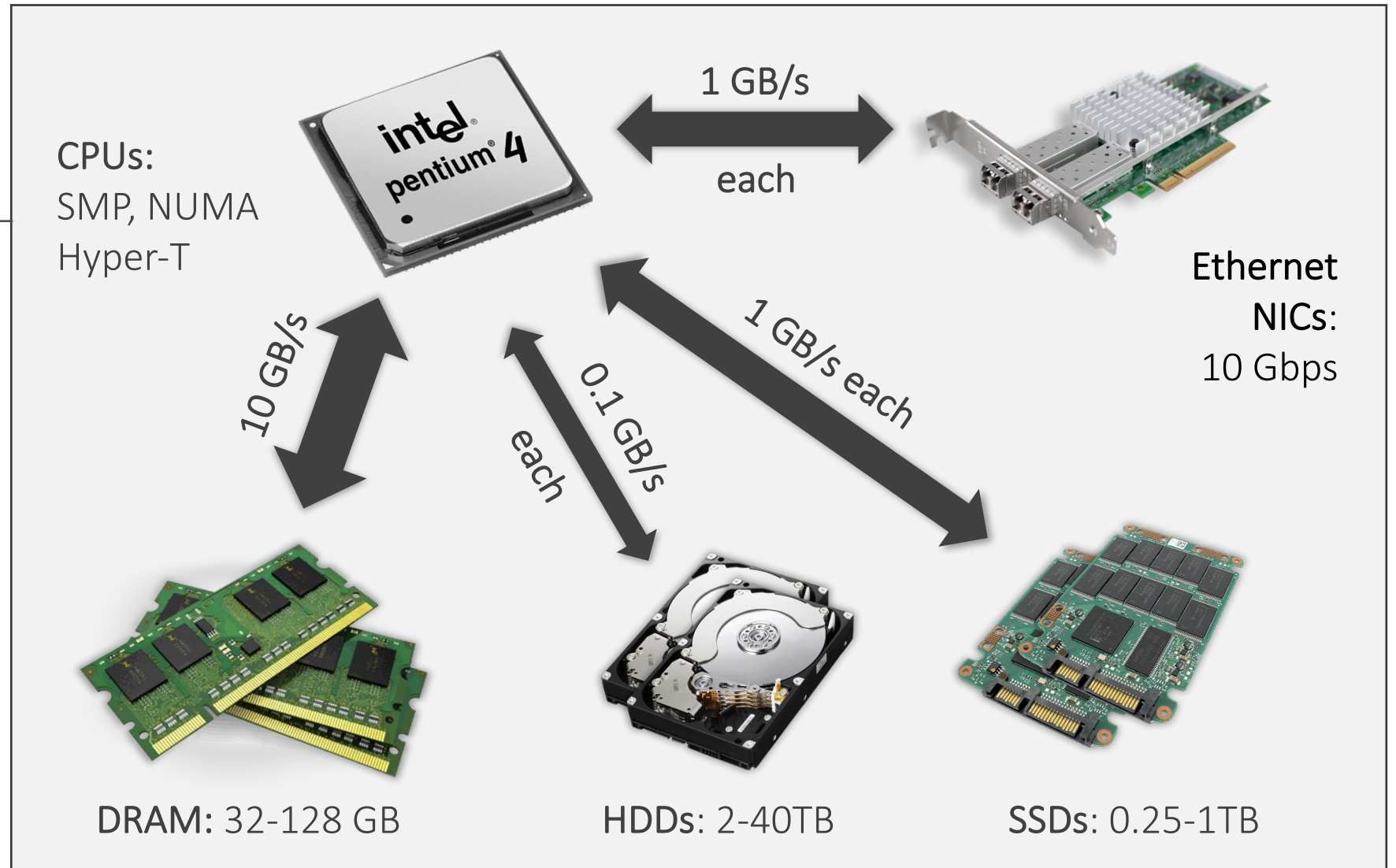
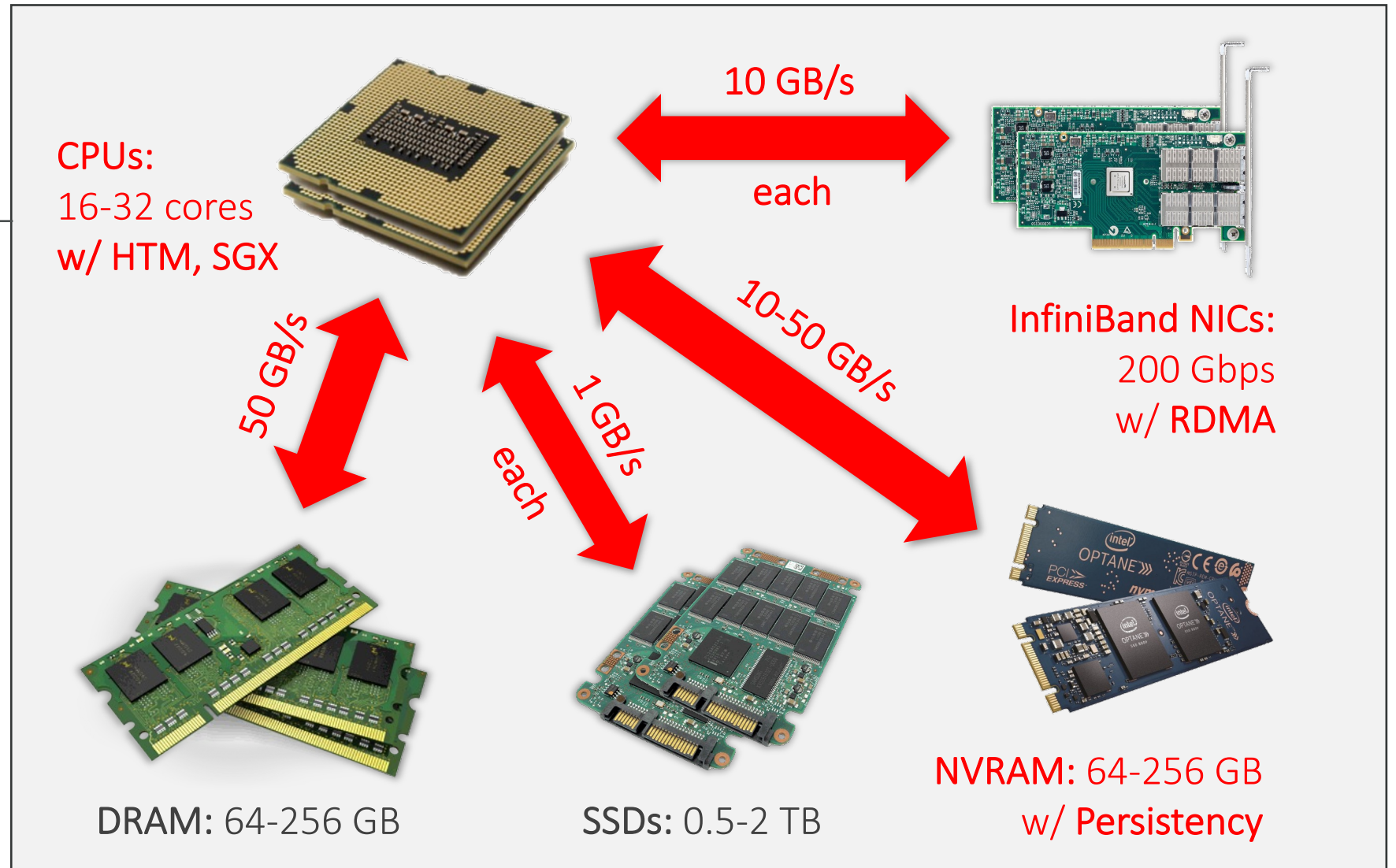Outlooking systems research for DPU

# Hardware in DC

Datacenter Server

CPUs:
SMP, NUMA
Hyper-T

1 GB/s each

10 GB/s

0.1 GB/s each

1 GB/s each

Ethernet NICs: 10 Gbps

DRAM: 32-128 GB

HDDs: 2-40TB

SSDs: 0.25-1TB

# Hardware in DC

**Datacenter Server**

**CPUs:**
16-32 cores
w/ HTM, SGX

10 GB/s each

**InfiniBand NICs:**
200 Gbps
w/ RDMA

50 GB/s

1 GB/s each

10-50 GB/s

**DRAM:** 64-256 GB

**SSDs:** 0.5-2 TB

**NVRAM:** 64-256 GB
w/ **Persistency**

# Common Practice: Offloading

Multicore
*Kernel/OS*

SGX
*Security*

Corey
OSDI'08

Multikernel
SOSP'09

Haven
OSDI'14

SCONE
OSDI'16

— 2008 — 2009 — 2010 — 2011 — 2012 — 2013 — 2014 — 2015 — 2016 — 2017 — 2018 — 2019 →

RDMA
*Distributed TX*

DrTM
SOSP'15

FaSST
OSDI'16

NOVA-Fortis
SOSP'17

ZoFS
SOSP'19

FaRM
SOSP'15

NVM
*File System*

SplitFS
SOSP'19

Systems
Community

# Opportunity: Collaboration
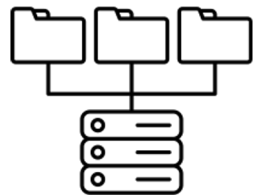
Different hardware devices can work together

▶ Case: RDMA NIC (RNIC) can directly access NVM

→ "Remote Persistent Memory"

▶ Scenarios: distributed logging in FS, TX, ..
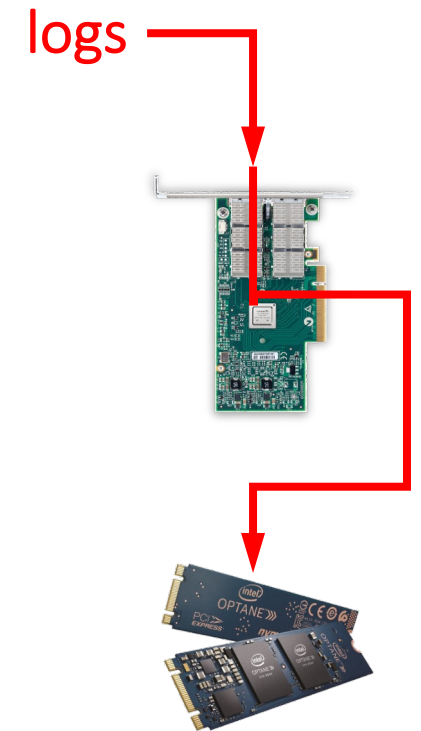
logs

Distributed Filesystems

Octopus [ATC'17]
Orion [FAST'19]

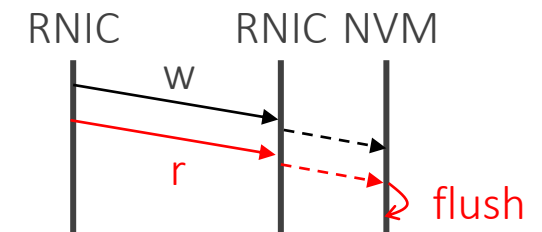Distributed Transactions

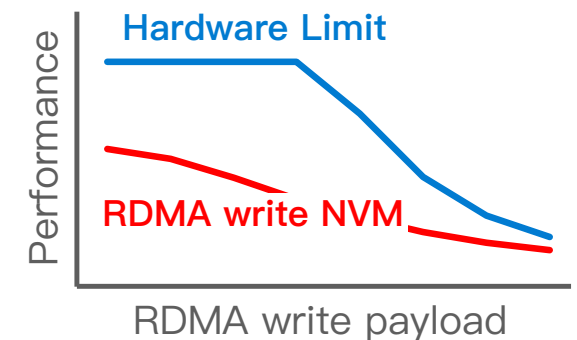DrTM+R [EuroSys'16]
FaRMv2 [SIGMOD'19]

# Challenge: Compatibility

*Functional* flaw: remote write is NOT persistent

▶ Solution[1]: **+ remote read** (two network roundtrips)

*Performance* pitfall: remote write is inefficient
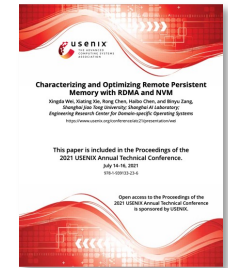
▶ **< 29%** of NVM thpt limit (15M vs. 52M reqs/s)

New hardware features are unaware of each other

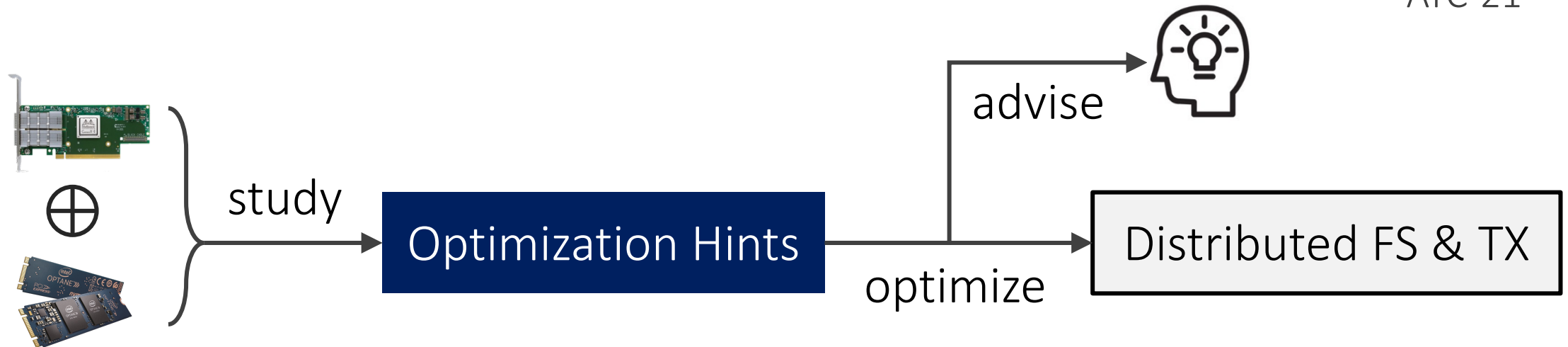[1] Intel. The librpmem library. https://pmem.io/pmdk/librpmem/

# Collaborative offloading for the concurrent use of RDMA & NVM

► Characterizing RDMA+NVM for optimization hints

► Case studies: distributed TX (DrTM+H) and FS (Octopus)
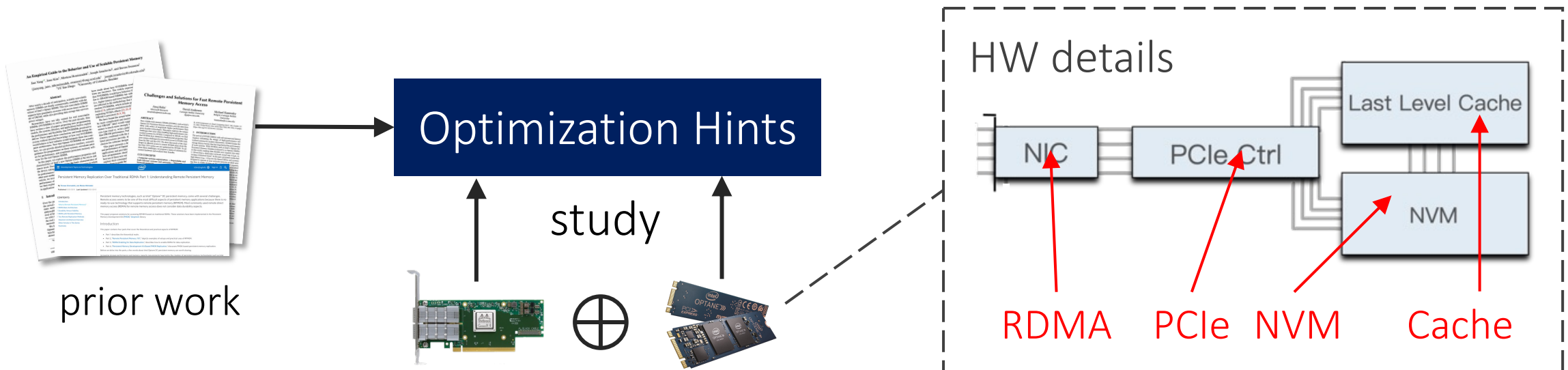
► Suggestions to RDMA/NVM hardware designers

USENIX
ATC'21

## Characterizing remote persistent memory w/ RDMA and NVM

► A systematic study of the collaboration btw. RDMA and NVM

► Tools: https://github.com/SJTU-IPADS/librdpma



prior work

Optimization Hints

study

HW details

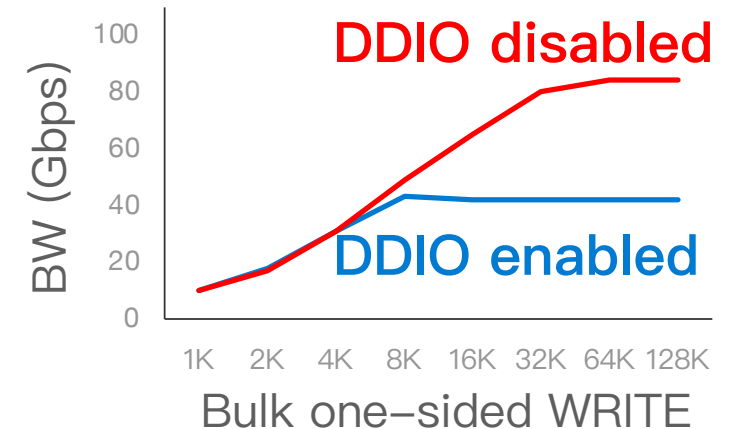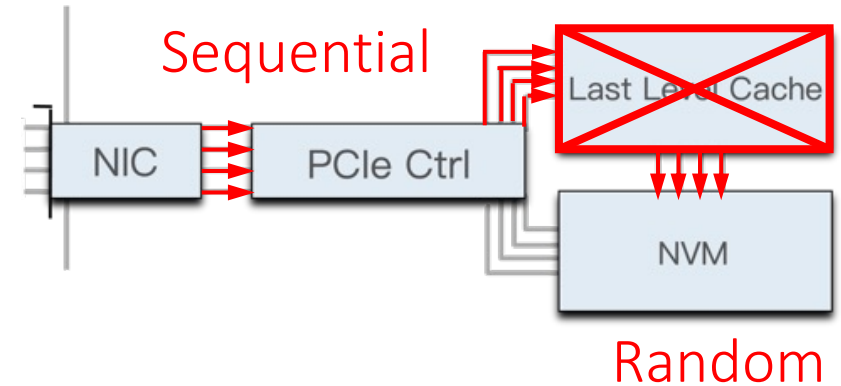RDMA   PCIe   NVM   Cache

# Example 1

## Optimization Hint

▶ Disable DDIO to skip LLC for large writes

## NVM feature

▶ Random I/O causes **write amplification**

## Performance pitfall

▶ RNIC **sequentially writes** the data to LLC

▶ Then, LLC **randomly evicts** the data to NVM

# Example 2

## Optimization Hint

► Use 64B granularity for small writes

Read-modify-write

NIC — PCIe Ctrl — Last Level Cache — NVM
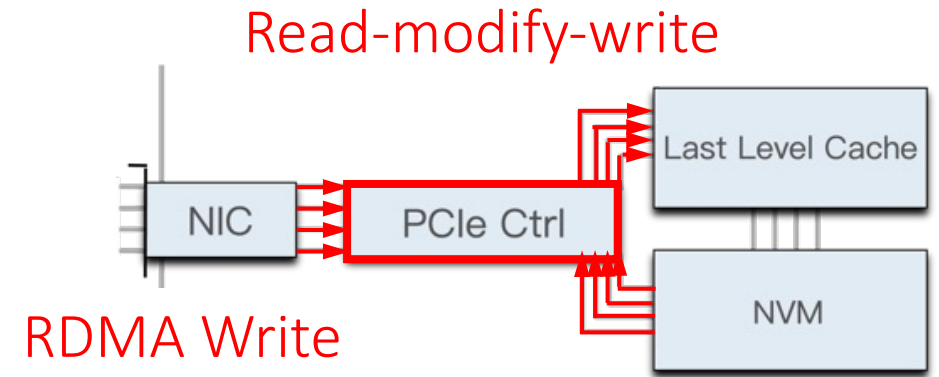
RDMA Write

## NVM feature

► Read-modify-write pattern (*PCIe partial-write*)

## Performance pitfall

► An **extra read** to NVM if write does not

fit **PCIe data word granularity** (e.g., 64B)

Throughput

**Hardware limit**

H1  H3  H5  H6

16B one-sided write

## Characterizing remote persistent memory w/ RDMA and NVM

► A design guideline: **9 optimization hints** in 3 aspects

► Achieve **87%** of NVM thpt limit (from 15M to 45M reqs/s)



prior work

Optimization Hints

study



Hardware Limit

+Opts

Performance

RDMA write NVM

RDMA write payload

# Optimizing

Applying our performance hints to existing RDMA-NVM systems

► **DrTM+H** (distributed TX) by **1.44×/2.09× for TPC-C/SmallBank**

► **Octopus** (distributed FS) by **2.40× for Data I/O**



Octopus
ATC'17

optimize

Optimization Hints

study

prior work

DrTM+H
OSDI'18

# Case Study: Distributed Transaction

## Applying our performance hints cumulatively on DrTM+H

| Hints | Optimizations |
|-------|---------------|
| H1 | Separate memory pool from different sockets to avoid cross-socket NVM access |
| H3 | Configure database with DDIO disabled |
| H4 | Use ntstore to optimize the commit phase |
| H5 | Align and pad logs/records larger than 256 B to XPLine granularity |
| H6+H7 | Align and pad logs/records smaller than 256 B to 64 B granularity |
| H8 | Implement a DRAM-based lock service for the validation phase |
| H9 | Implement remote persistent log with H9 in one roundtrip |

**Improve perf. & enable persist**

+NVM  H4 H1 H5 H3 H6 H7 H8 +Persist H9

Throughput

**2.09×**

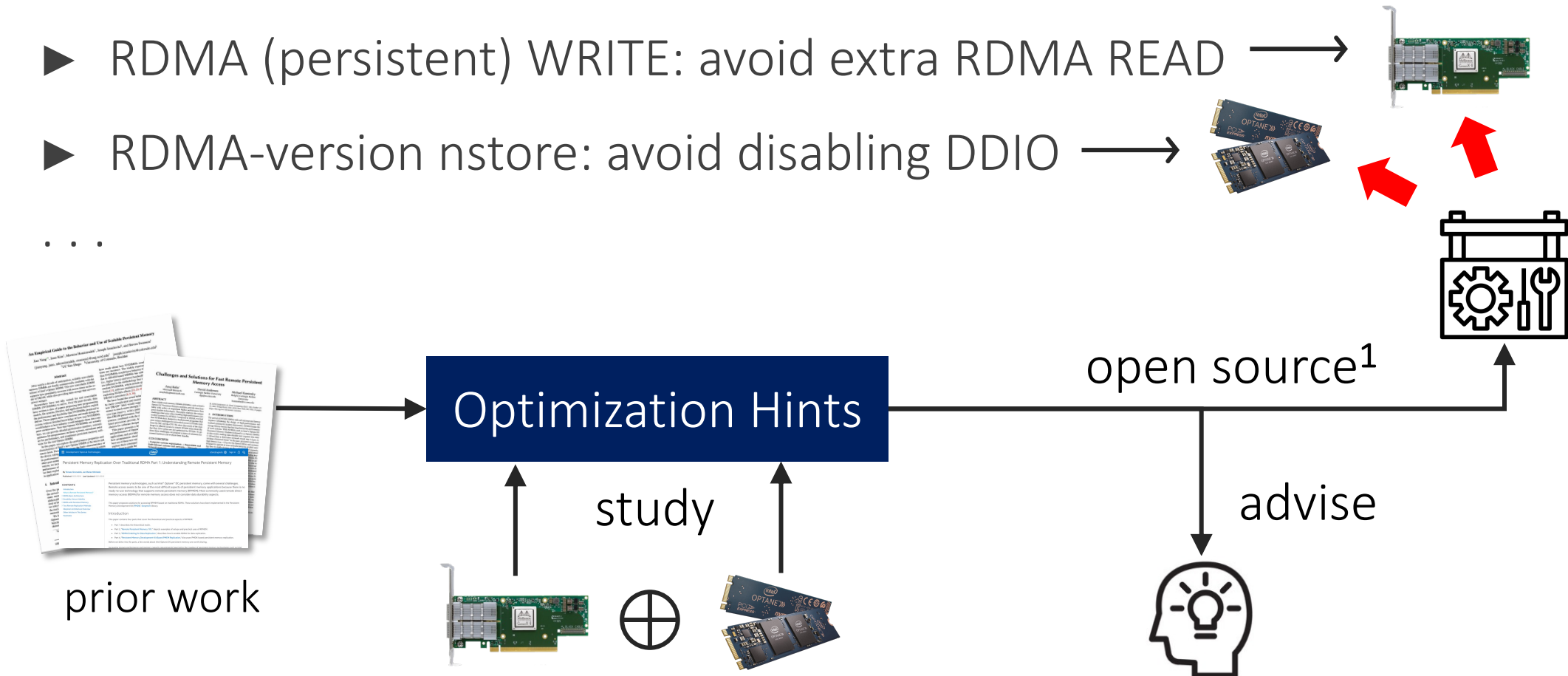Factor Analysis of SmallBank

**Optimization Hints** → optimize → DrTM+H OSDI'18

## Suggestions to hardware designers

► RDMA (persistent) WRITE: avoid extra RDMA READ  →

► RDMA-version nstore: avoid disabling DDIO  →

. . . .

**Optimization Hints**

open source[1]

study

advise

prior work

# Outline

**Intelligent** Hardware ➡ Network **+** Computation

Storage

CPU, FPGA, ASIC

. . .

SmartNIC

SmartSSD

Smart + X

Integrating multiple capabilities into a single device

▶ Typical case: DPU/SmartNIC (e.g., Nvidia BlueField)

## Integrating multiple capabilities into a single device

► **Typical case**: DPU/SmartNIC (e.g., Nvidia BlueField)

Nvidia
BlueField-2



- ConnectX-6 (2x 100Gbps)
- 16 GB of on-board DRAM
- ARM Cortex-A72 (8 cores)



**Intelligence**
ARM Cortex-A72
Accelerators

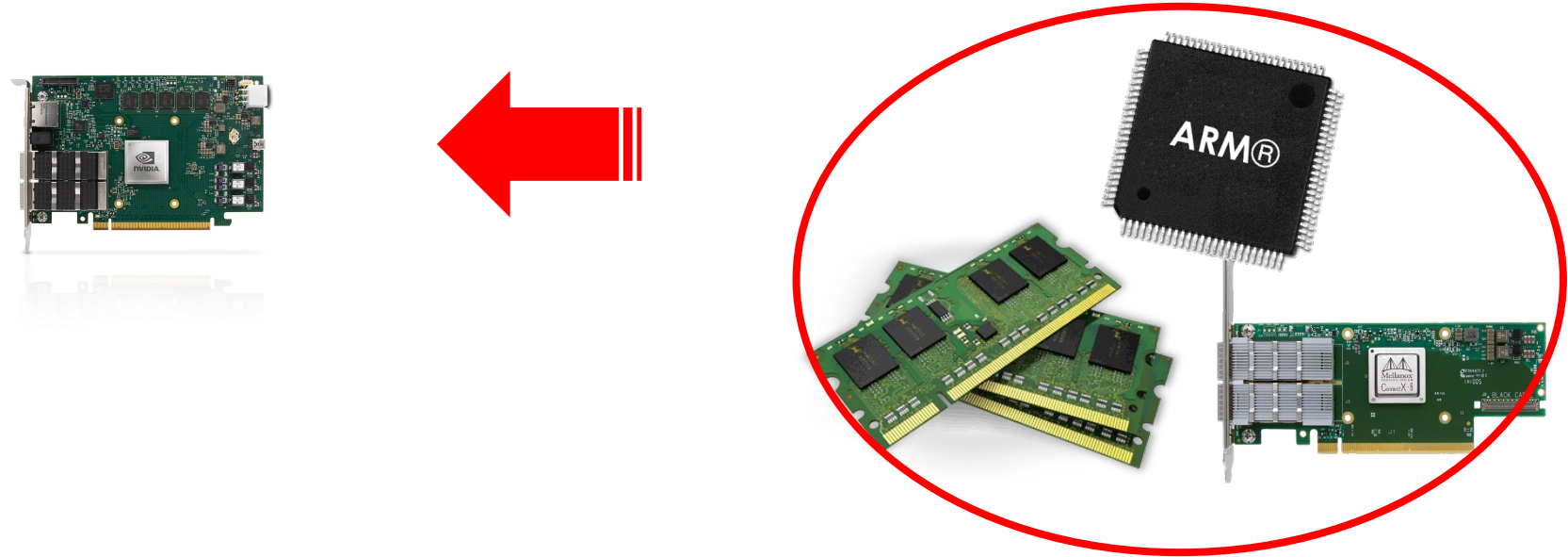**Network**
ConnectX-6
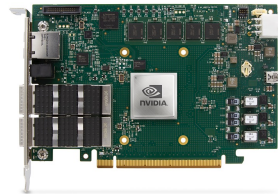
**Storage:**
16GB DDR4

BlueField 2

# New Trend : Capability Integration

## Integrating multiple capabilities into a single device

► **Typical case**: DPU/SmartNIC (e.g., Nvidia BlueField)

► **Good**: Innately immune to compatibility issues

► **Bad**: (much) higher cost, compared to RNIC



|  |  | BlueField-2[1] | ConnectX-6[2] |
|---|---|---|---|
| Price | 1.5× | $ 3615 | $ 2,440 |
| Space | 2.0× | 6.6 in. x 4.53 in. | 6.6 in. x 2.71 in. |
| Power | 3.2× | 75W | 23.6W |

[1] NVIDIA MBF2H516A-EEEOT BlueField-2
[2] NVIDIA MCX653106A-HDAT ConnectX-6

# Challenge: Underutilization

DPU is inferior in *every* single capability

- ▶ Wimpy cores (e.g., 8-core ARM) and small memory (e.g., 16GB)

- ▶ Net. perf. degradation (BF-2 vs. CX-6): latency (+6~30%), thpt (-15~36%)

## Case study: Get (k) in distributed key/value store (KVS)

### RNIC-based KVS

2x RDMA READs (1 for index, 1 for value)



### DPU-based KVS

1x SEND/RECV, offload indexing to DPU



YCSB C  THPT

14%

# Challenge: Underutilization

DPU is inferior in *every* single capability

▶ Wimpy cores (e.g., 8-core ARM) and small memory (e.g., 16GB)

▶ Net. perf. degradation (BF-2 vs. CX-6): latency (+6~30%), thpt (-15~36%)
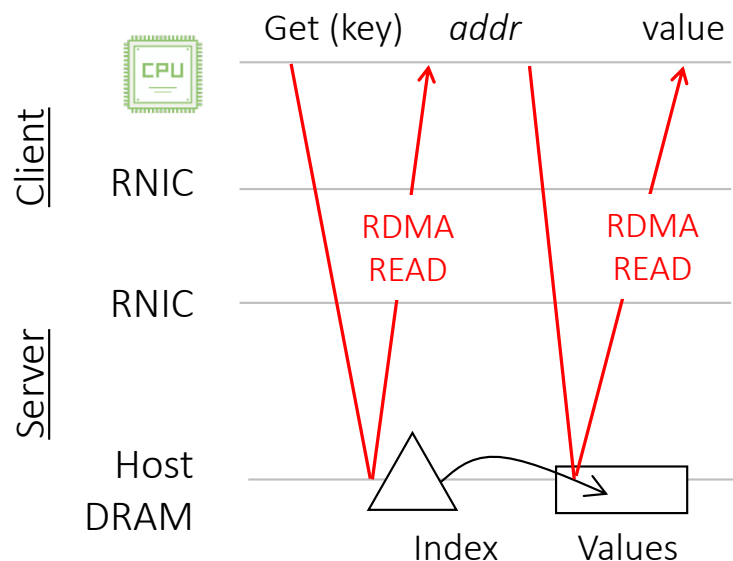
Existing systems only utilize a portion of DPU device

▶ Only NIC-Host path, treated as RNIC

▶ Only computing resource (SoC), treated as accelerator

A systematic way to fully utilize integrated capabilities

# Our work

## Cooperative offloading for fully utilizing DPU

► Characterizing: study offloading paths, rather than HW components

► A step-by-step optimization guideline for DS designer

► Case studies: DPU-accelerated distributed FS and KV

► Open-source toolkit: https://github.com/smartnickit-project

USENIX
OSDI'23

# Characterizing DPU (i.e., BlueField 2) in path level

► Study offloading paths, rather than HW components

► Four paths: NIC-Host (①), NIC-SoC (②), SoC-Host (③), SoC-only (④)

► Performance implications: bottlenecks, anomalies, and takeaways

# Example 1

## Findings

► NIC-Host is slower than RNIC

► Overhead: PCIe latency (300ns x4)

► Non-trivial for small request (1-2μs)

## Takeaway

► If only NIC-Host is used, select RNIC

as it is faster, cheaper, and saves power

RNIC          SoC

① 

RX/RX  PCIe1  PCIe Switch    DRAM

PCIe0

Host   Host DRAM   CPU CPU CPU CPU

RNIC ①   Latency (μs)

READ      WRITE     SND/RCV

(Payload = 64B)

# Example 2

## Findings

► NIC-SoC is faster than NIC-host (no PCIe0),

but still slower than RNIC (PCIe switch)

► SEND/RECV is much slow (wimpy SoC cores)



(Payload = 64B)

# Example 3

## Findings

▶ RDMA READ performance of NIC-SoC

collapses w/ large request (>=9MB)

## Advice: avoid large READ requests

▶ PCIe MTU: Host (512B) vs. SoC (128B)

▶ NIC-SoC READ: 4× PCIe packets

for large requests → HoL blocking



② WRITE

② READ

① WRITE
/READ

# Example 4

## Characterizing concurrent paths in DPU

▶ DPU is always underutilized when only using a single path

▶ Study the concurrent use of multiple offloading paths (e.g., ①+②)

## Takeaway

▶ Concurrent offloading can better utilize DPU, esp. when used in opposition directions (R+W)

▶ But, carefully avoid interference btw. paths, e.g., NIC cores (①+②) and PCIe switch (②+③)

## Our path-level DPU study

► A comprehensive perf. study on *offloading paths* (6) × *primitives* (3)

► 11 *findings/advice* for either individually using a single path or concurrently using multiple paths

study → path-level study

Advice/Findings

**Table 3:** The findings and advice from our study. Claims supported by sufficient evidence are denoted by **E**, while those supported by hypotheses are denoted by **H**.

| SNIC Paths | Findings/Advice | E/H |
|---|---|---|
| ① (§3.1) | Throughput of RDMA is lower than RNIC | H |
|  | Latency of RDMA is higher than RNIC | E |
| ② (§3.2) | One-sided RDMA performance is better | H |
|  | Avoid memory accesses to close addresses | E |
|  | Avoid large READ requests | H |
| ③/③* (§3.3) | RDMA overuses the PCIe bandwidth | E |
|  | Avoid large READ/WRTIE requests | H |
|  | Enable doorbell batching carefully for RDMA | E |
|  | Use DMA (③*) to improve PCIe utilization | E |
| ①+② (§4.1) | Improve throughput by using paths ① and ② concurrently (esp. in opposite directions) | H |
| ①/②+③ (§4.1) | Selectively offload traffic to ③ | E |

USENIX
OSDI'23

# Optimizing

## A step-by-step optimization guideline for system designers

1. Devise potential alternatives for DPU to support the given functionality, and optimize them based on our study

2. Evaluate and rank alternatives based on system-specific criteria

3. Select and combine alternatives in turn until DPU is saturated

study → path-level study

Advice/Findings → Guideline

# Case Study: Get(k) in Key/Value Store

## 1. Devise alternatives (A1-A5) and optimize them



A1     A2     A3     A4     A5

① READx2   ② S/R ③ READx2   ② S/R ④ ③ READ   ② READ ① READ   ② S/R ④ + ② READx2

Guideline — rebuild → DPU-accelerated KVS

## 2. Evaluate and rank alternatives based on high performance



Rank: A5 > A4 > A1 > A3 > A2

Guideline →(rebuild)→ DPU-accelerated KVS

## 3. Select and combine alternatives in turn until DPU is saturated



Rank: A5 > A4 > A1 > A3 > A2

A5+A4: use A5 first until SoC is saturated, and then select A4

Guideline —rebuild→ DPU-accelerated KVS

## Suggestions to hardware designers

► Support CXL to relieve the pressure on SoC cores

► Support ARM CCI (similar to DDIO on host CPU)

► Align PCIe MTU of SoC and Hots CPU

· · ·

Encourage hardware vendors to disclose more details of DPU

advise

study

path-level study

Advice/Findings

rebuild

Guideline

Distributed Systems

# Outline

Case #1: Collaborative offloading

Case #2: Cooperative offloading

Outlooking systems research for DPU

# Outlooking

## Which type of processor should be selected, SoC, FPGA, or ASIC?

► An inherent trade-off in programmability and performance

# Outlooking

## Which type of processor should be selected, SoC, FPGA, or ASIC?

► "Don't want to CHOOSE, want BOTH": SoC + FPGA/ASIC/...

| Marvell OCTEON | Intel IPU | Broadcom Stingray | NVIDIA BlueField |
|---|---|---|---|
| SoC (ARM) | Soc (Xeon) | SoC (ARM) | SoC (ARM) |
| VPP Accelerators | Agilex FPGA | ASIC Accelerators | ASIC Accelerators |
| | | | DPA (RISC-V) |

## How to measure integrated hardware components in DPU?

► New *metrics, benchmarks,* and *toolkits*?



Power the Next Wave of Applications with NVIDIA BlueField-3 DPUs

OVS Acceleration

UDF Acceleration

IPsec Acceleration

Evaluating performance/power-saving of accelerators for network functions

# Outlooking

## Which domain-specific accelerators deserve to be integrated?

### The killer applications of DPUs
datacenter networking, storage, security, and virtualization workloads

**Power the Next Wave of Applications with NVIDIA BlueField-3 DPUs**

- CPUs that are used for serial processing and running hyperthreaded applications.
- GPUs that excel at parallel processing and are optimized for accelerating modern workloads.
- DPUs that are ideal for infrastructure computing tasks; used to offload, accelerate, and isolate data center networking, storage, security, and manageability workloads.

HPC/AI

Cloud Computing

Storage

**Marvell's OCTEON 10 Challenges All Comers For DPU Supremacy**

- **Cloud and data center** servers to offload virtual overlay and cryptographic processing for multi-tenant VM, container, and storage services.
- **LTE and 5G vRAN** implementations when paired with Marvell's Fusion-O baseband processor providing a 5G and LTE-A PHY with the OCTEON used for CU or vRAN offload processing.
- **Enterprise router-firewall and SD-WAN** appliances using NFV service chaining to deliver L2/L3 forwarding, VPN termination, SPI, and new AI-based applications and security services.

Source: https://developer.nvidia.com/blog/power-the-next-wave-of-applications-with-nvidia-bluefield-3-dpus/
https://packetpushers.net/marvells-octeon-10-challenges-all-comers-for-dpu-supremacy/

# Outlooking

## Which domain-specific accelerators deserve to be integrated?

▶ Compression, encryption, virtualization, packet processing, . . .

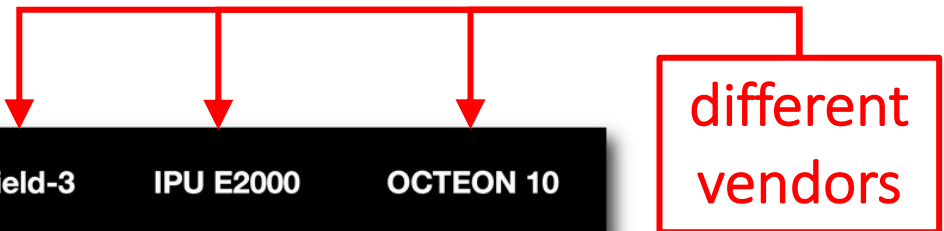| Acceleration | BlueField | BlueField-2 | BlueField-3 | IPU E2000 | OCTEON 10 |
|---|---|---|---|---|---|
| DMA | ✅ | ✅ | ✅ | | ✅ |
| Compress | | ✅ | ✅* | ✅ | |
| Erasure coding | | | ✅ | | |
| Regex | | ✅ | ✅ | | |
| Off-path encryption | ✅ | ✅ | | ✅ | |
| On-path encryption | ✅ | ✅ | ✅ | ✅ | ✅ |
| Packet processing | ✅ | ✅ | ✅ | ✅ | ✅ |
| Year | 2016 | 2020 | 2023 | 2023 | 2021 |

# Outlooking

## Which domain-specific accelerators deserve to be integrated?

► Compression, encryption, virtualization, packet processing, . . .

| Acceleration | BlueField | BlueField-2 | BlueField-3 | IPU E2000 | OCTEON 10 |
|---|---|---|---|---|---|
| DMA | ✅ | ✅ | ✅ | DIFF | ✅ |
| Compress | | ✅ | ✅* | ✅ | DIFF |
| Erasure coding | | | ✅ | DIFF | DIFF |
| Regex | | ✅ | ✅ | DIFF | DIFF |
| Off-path encryption | ✅ | ✅ | DIFF | ✅ | DIFF |
| On-path encryption | ✅ | ✅ | ✅ | ✅ | ✅ |
| Packet processing | ✅ | ✅ | ✅ | ✅ | ✅ |
| Year | 2016 | 2020 | 2023 | 2023 | 2021 |

different vendors

## Which domain-specific accelerators deserve to be integrated?

▶ Compression, encryption, virtualization, packet processing, . . .
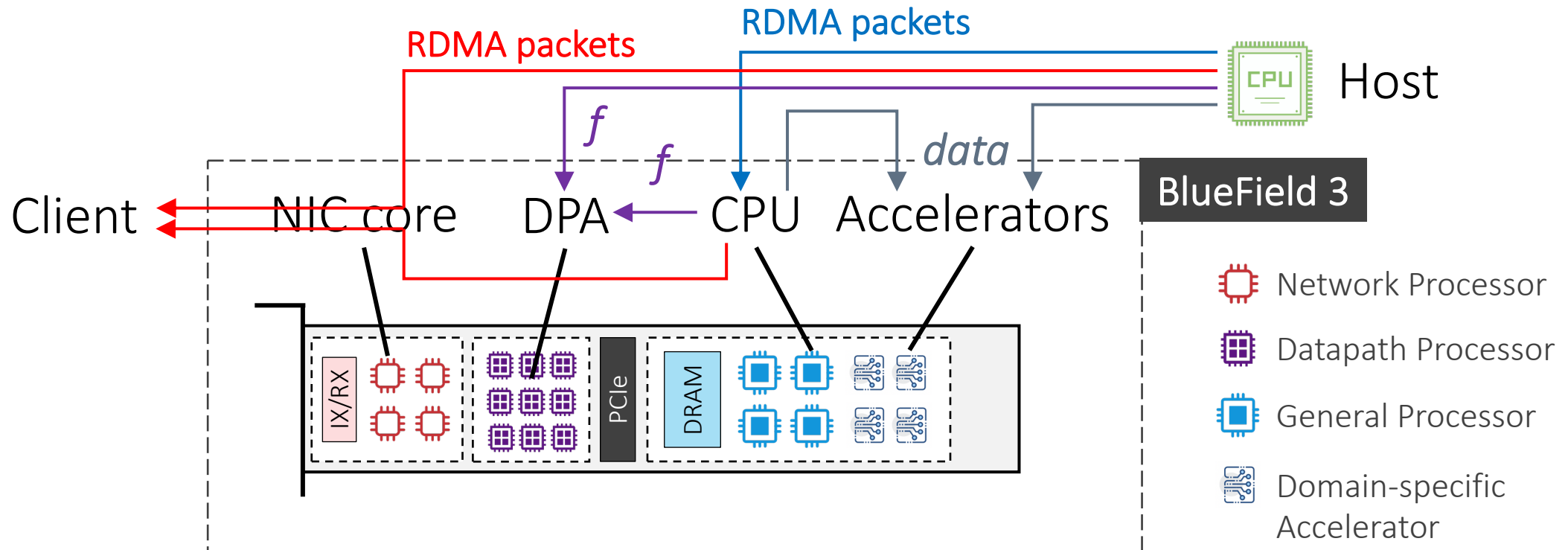
different versions

| Acceleration | BlueField | BlueField-2 | BlueField-3 | IPU E2000 | OCTEON 10 |
|---|---|---|---|---|---|
| DMA | ✅ | ✅ | ✅ | | ✅ |
| Compress | ADD ---→ ✅ | | ✅* | ✅ | |
| Erasure coding | | ADD ---→ ✅ | | | |
| Regex | ADD ---→ ✅ | | ✅ | | |
| Off-path encryption | ✅ | ✅ --✗--→ DEL | | ✅ | |
| On-path encryption | ✅ | ✅ | ✅ | ✅ | ✅ |
| Packet processing | ✅ | ✅ | ✅ | ✅ | ✅ |

| Year | 2016 | 2020 | 2023 | 2023 | 2021 |

# Outlooking

## How to unify system abstraction & programming interface?

► e.g., BlueField: PCIe accelerator vs. a standalone server

## How to unify system abstraction & programming interface?

► e.g., BlueField: PCIe accelerator vs. a standalone server

Hardware evolution:

*single capability breakthrough* & *multiple capability integration*

Our approach: characterizing, optimizing, and advising

► Collaborative offloading for multiple devices (e.g., RDMA & NVM)

► Cooperative offloading for intelligent devices (e.g., DPU)

Our outlook on systems research for DPU

See more at https://ipads.se.sjtu.edu.cn/rong_chen

THANK YOU