DASYA

www.dasya.dk

DAPHNE

https://daphne-eu.github.io/

IT UNIVERSITY OF COPENHAGEN
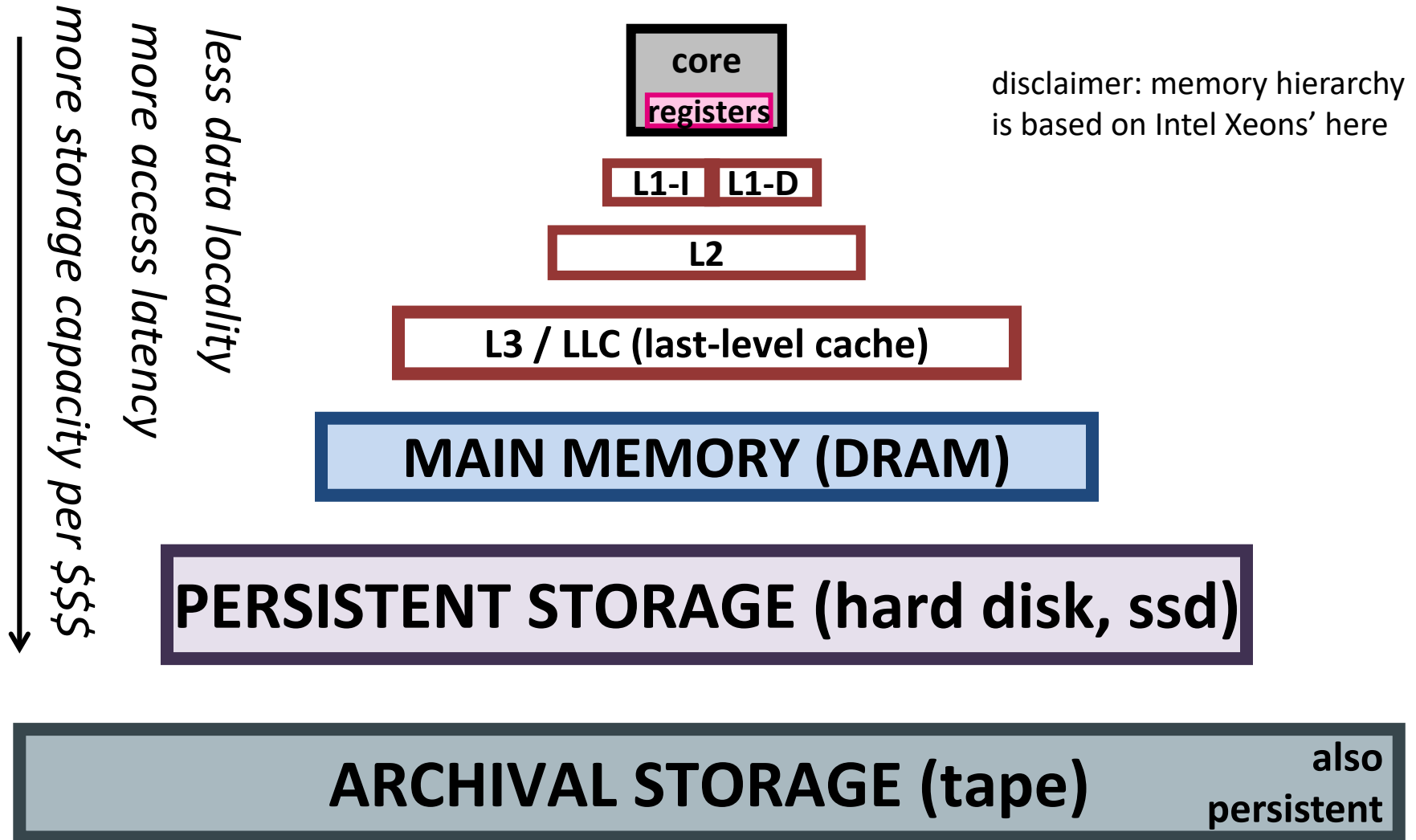
www.itu.dk

# data-intensive systems in the microsecond era

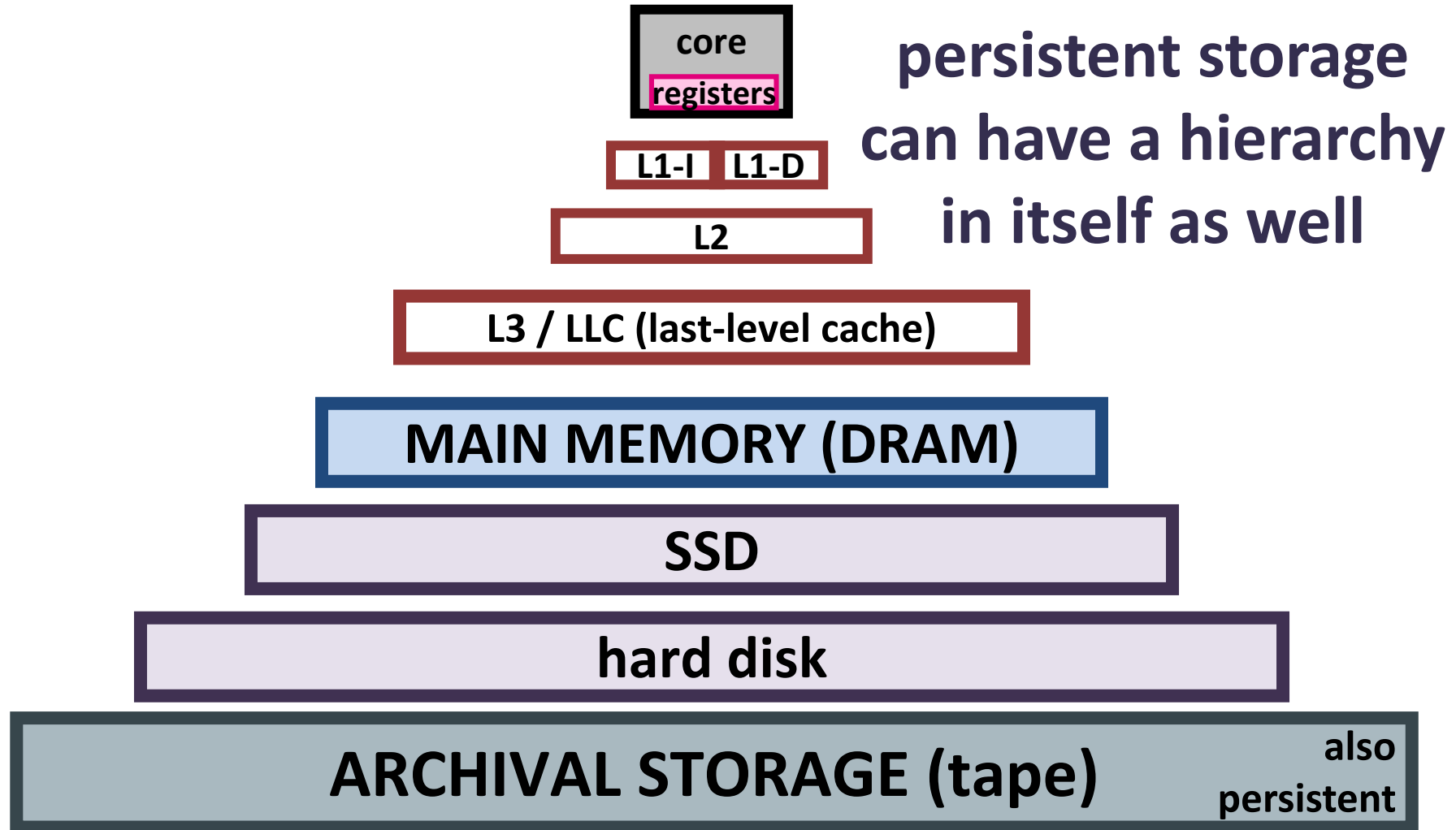## Pınar Tözün

pito@itu.dk, www.pinartozun.com

work done in collaboration
with Philippe Bonnet @ ITU

# (typical) storage hierarchy

*more storage capacity per $$$*
*more access latency*
*less data locality*

**core**
**registers**

disclaimer: memory hierarchy is based on Intel Xeons' here

**L1-I** **L1-D**

**L2**

**L3 / LLC (last-level cache)**

**MAIN MEMORY (DRAM)**

**PERSISTENT STORAGE (hard disk, ssd)**

**ARCHIVAL STORAGE (tape)**

**also persistent**

2

# (typical) storage hierarchy

**core**
**registers**

**L1-I**  **L1-D**

**L2**

**L3 / LLC (last-level cache)**

**MAIN MEMORY (DRAM)**

**SSD**

**hard disk**

**ARCHIVAL STORAGE (tape)**

**persistent storage can have a hierarchy in itself as well**

**also persistent**

3

# (typical) storage hierarchy

**distributed setting (e.g., cluster of machines)**

**core**

**registers**

**persistent storage can have a hierarchy in itself as well**

**L1-I** **L1-D**

**L2**

**L3 / LLC (last-level cache)**

**MAIN MEMORY (DRAM)**

**local disk**

**remote disk**

**ARCHIVAL STORAGE (tape)**

**also persistent**

# latency to fetch data

**core**
**registers**

L1-I  L1-D

L2

L3 / LLC (last-level cache)

**MAIN MEMORY (DRAM)**

**NVMe SSD**

**hard disk**

**ARCHIVAL STORAGE (tape)**

also persistent

1 cycle

~4 cycles

~10 cycles

~30-60 cycles

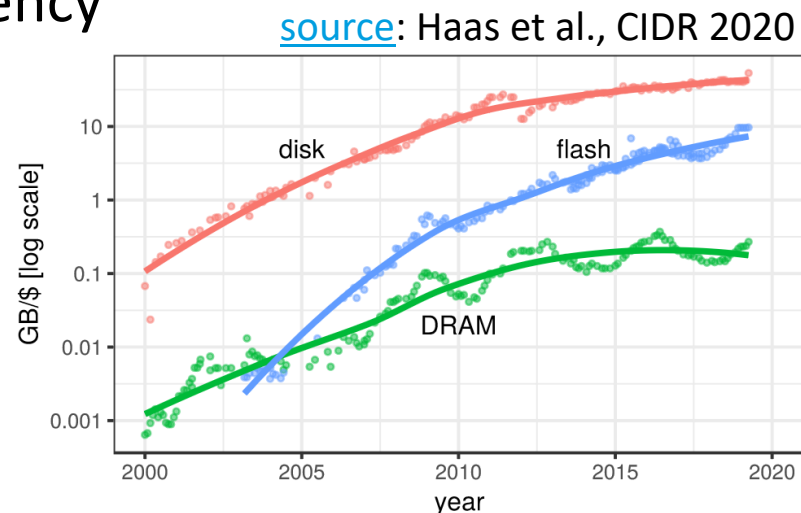~60ns or ~100-200 cycles

~10 μsec

~5ms

~100sec

5

# why focus on SSDs in this talk?

➜ except for SSDs, each layer stayed almost stable the last decade in terms of latency

- improvements on SSD internals
- from SAS/SATA to PCIe
- linux block IO improvements e.g., multiqueue

➜ improved price/capacity

source: Haas et al., CIDR 2020



➜led to several SSD-optimized data systems
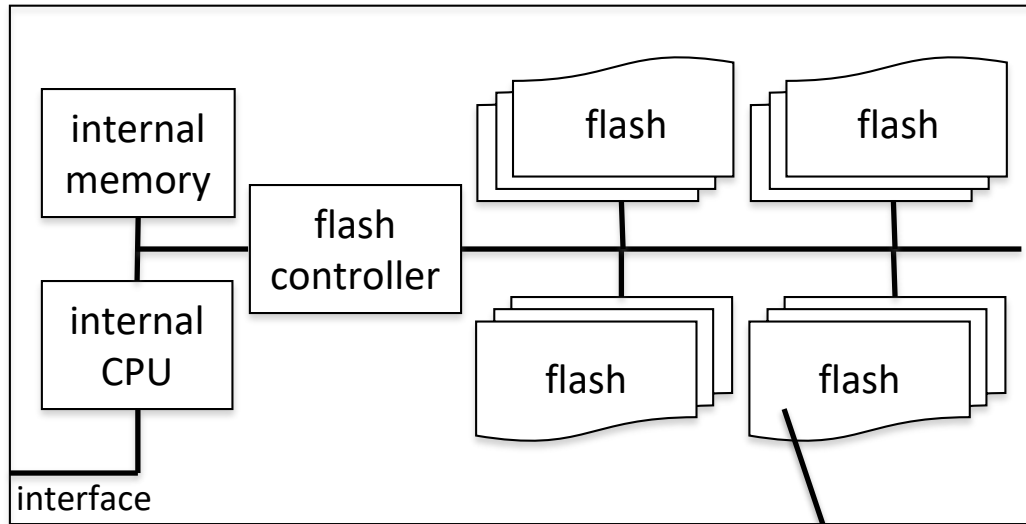
- RocksDB, BwTree, LeanStore, Umbra …

## increasing shift from pure in-memory-optimized to SSD-optimized data systems!

# agenda

- SSD internals & state of affairs today
- emerging SSD & computational storage landscape
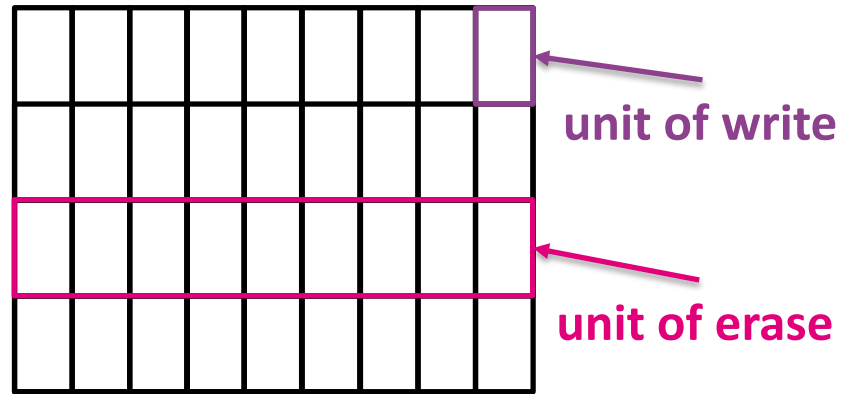
# solid-state disk (SSD)

interconnected flash chips

hard disk compatible API

compared to hard disks
- efficient random access
- internal parallelism

# flash chips

**unit of write**

**unit of erase**

**flash translation layer (FTL) hides the internal complexities of flash chips from end-users**

**but knowing them can lead to smarter software design**

**cannot override a unit before erasing it first**

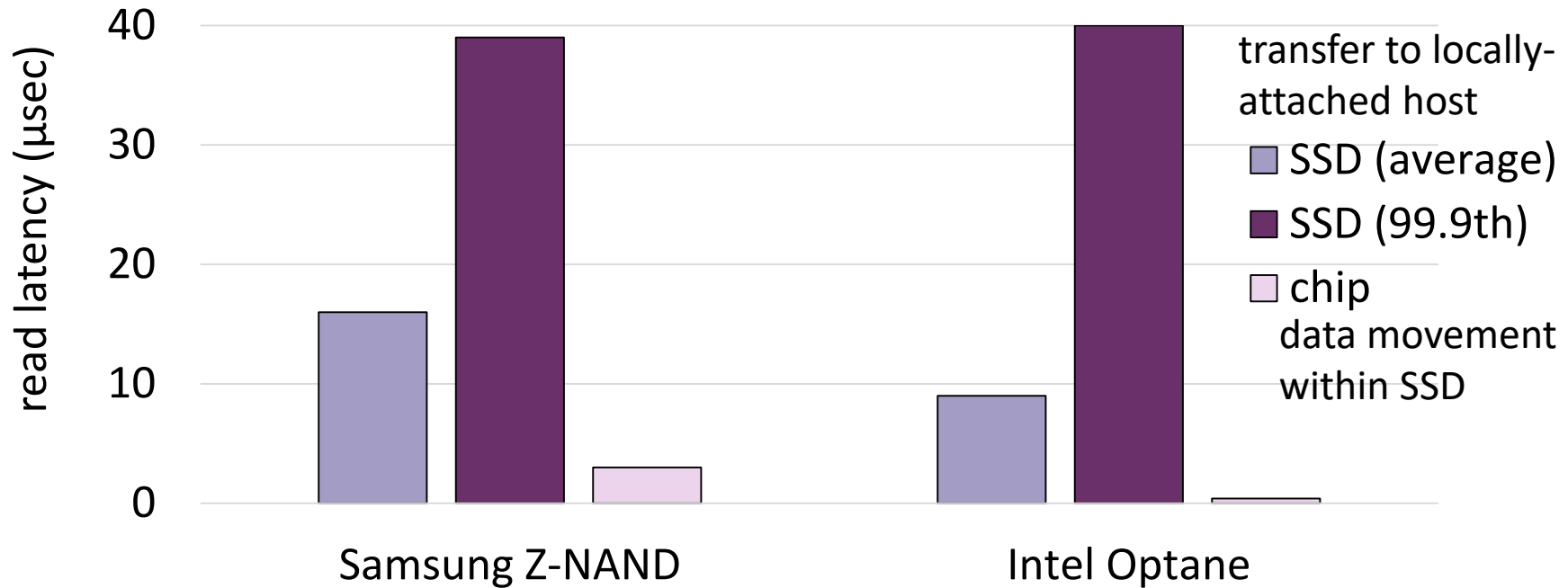**garbage collection –** for not used blocks so we can rewrite them

**write amplification** = data physically written / data logically written >= 1

writing data might cause rewrites & garbage collection

**wear leveling –** some cells/blocks die over time

**unpredictable read/write latencies**

# SSDs in the μsec era
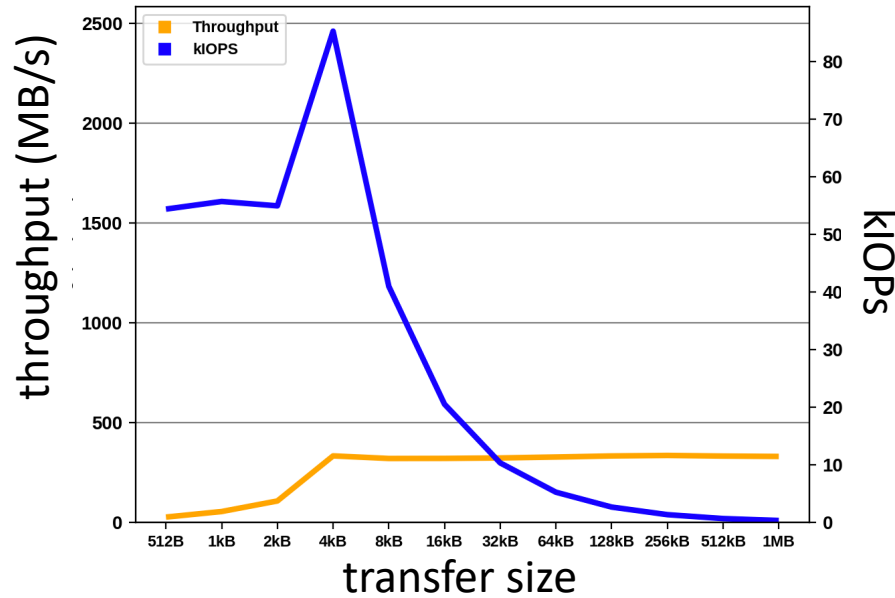
4K random read using fio - sources: [1, 2, 3]



**SSDs equipped with Z-NAND & Optane deliver at best 5x & 20x the read latency of the underlying storage chip, respectively.**
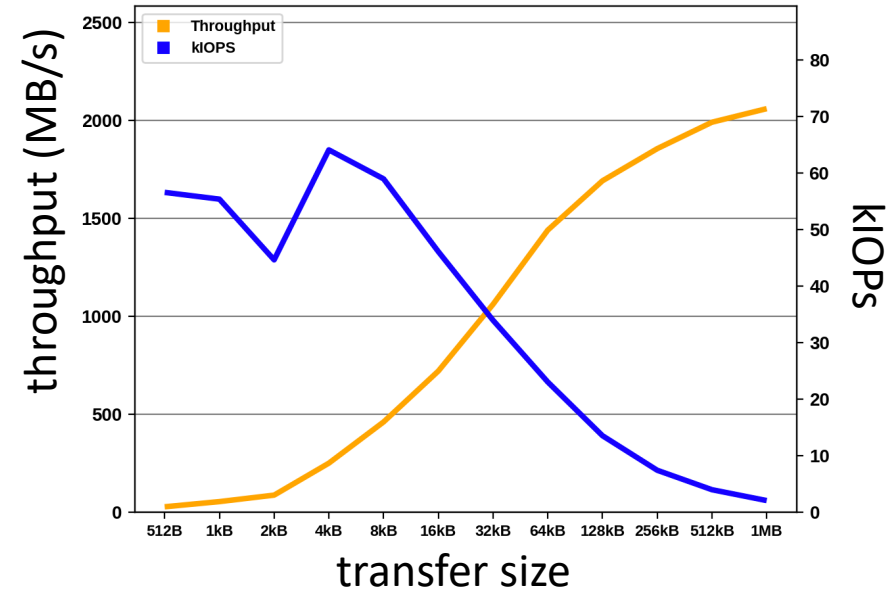
# FTLs in the μsec era ..

random writes- source: AnandTech



Samsung SSD with Z-NAND



Intel Optane

# FTLs in the μsec era ..

random writes- source: AnandTech



Samsung SSD with Z-NAND

Intel Optane

Samsung SSD with NAND

## … have drastic impact on throughput!

# linux IOs in the μsec era

sources: Faster IO through io_uring &
Efficient I/O with io_uring & J.Axboe



app

app

app

SPDK

app

shared rings
for submissions
and completions

user space

memory copy

OS kernel

aio

io_uring

driver

driver

IRQ-based

polling+driver

SSD

SSD

SSD

storage

polling or IRQ-based

4k random reads
3d xPoint

latency (μsec)

16

12

8

4

0

0    2    4    6    8
queue depth

– – io_uring without polling
—— aio
– – io_uring with polling
—— spdk

**separation of control & data plane in linux now
zero copy & minimized synchronization overhead**

**the benefits of fast storage wasted by**
**- data movement overheads**
**(from device to host & across network)**
**- black-box generic flash-translation layers**
**- multitude of software layers**

**how do we prevent these?**

# agenda

- SSD internals & state of affairs today
- emerging SSD & computational storage landscape

# computational storage

## back when I was a kid

**Put Everything
in Future (Disk) Controllers
(it's not "if", it's "when?")**

Jim Gray

http://www.research.Microsoft.com/~Gray

Acknowledgements:
**Dave Patterson** explained this to me a year ago
**Kim Keeton**
**Erik Riedel** } Helped me sharpen
these arguments
**Catharine Van Ingen**

1

**Basic Argument for x-Disks**
- Future disk controller is a super-computer.
  - 1 bips processor
  - 128 MB dram
  - 100 GB disk plus one arm
- Connects to SAN via high-level protocols
  - RPC, HTTP, DCOM, Kerberos, Directory Services,….
  - Commands are RPCs
  - management, security,….
  - Services file/web/db/… requests
  - Managed by general-purpose OS with good dev environment
- Move apps to disk to save data movement
  - need programming environment in controller

Jim Gray, NASD Talk, 6/8/98
http://jimgray.azurewebsites.net/jimgraytalks.htm
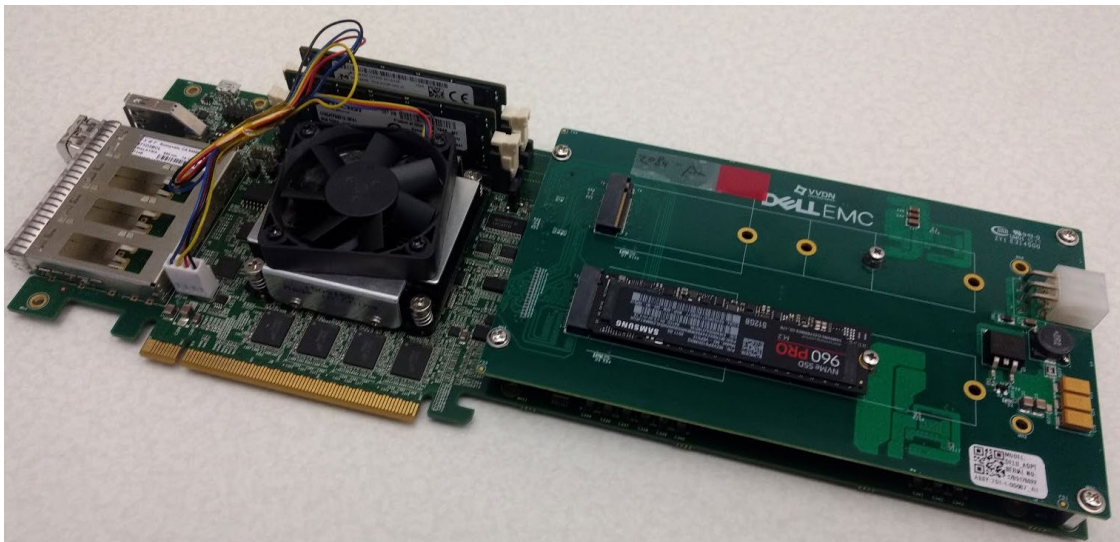
## = computation on the IO path

# computational storage

back when I joined ITU



**8-core ARMv8 processor**

**32GB DRAM**

**2TB+ of NVM via M.2 slots**

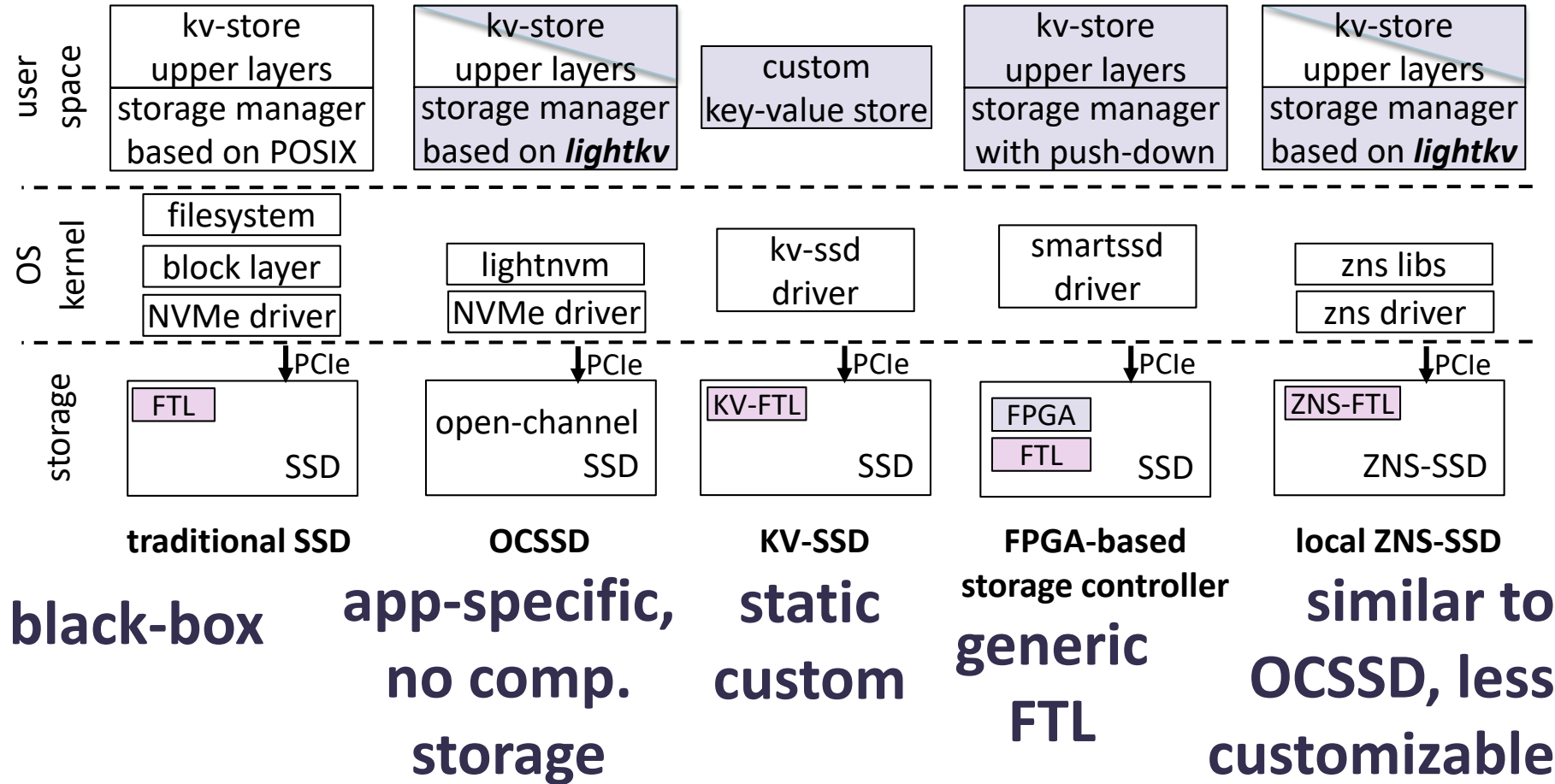**4x 10Gb Ethernet**

Dragon Fire Card (DFC)
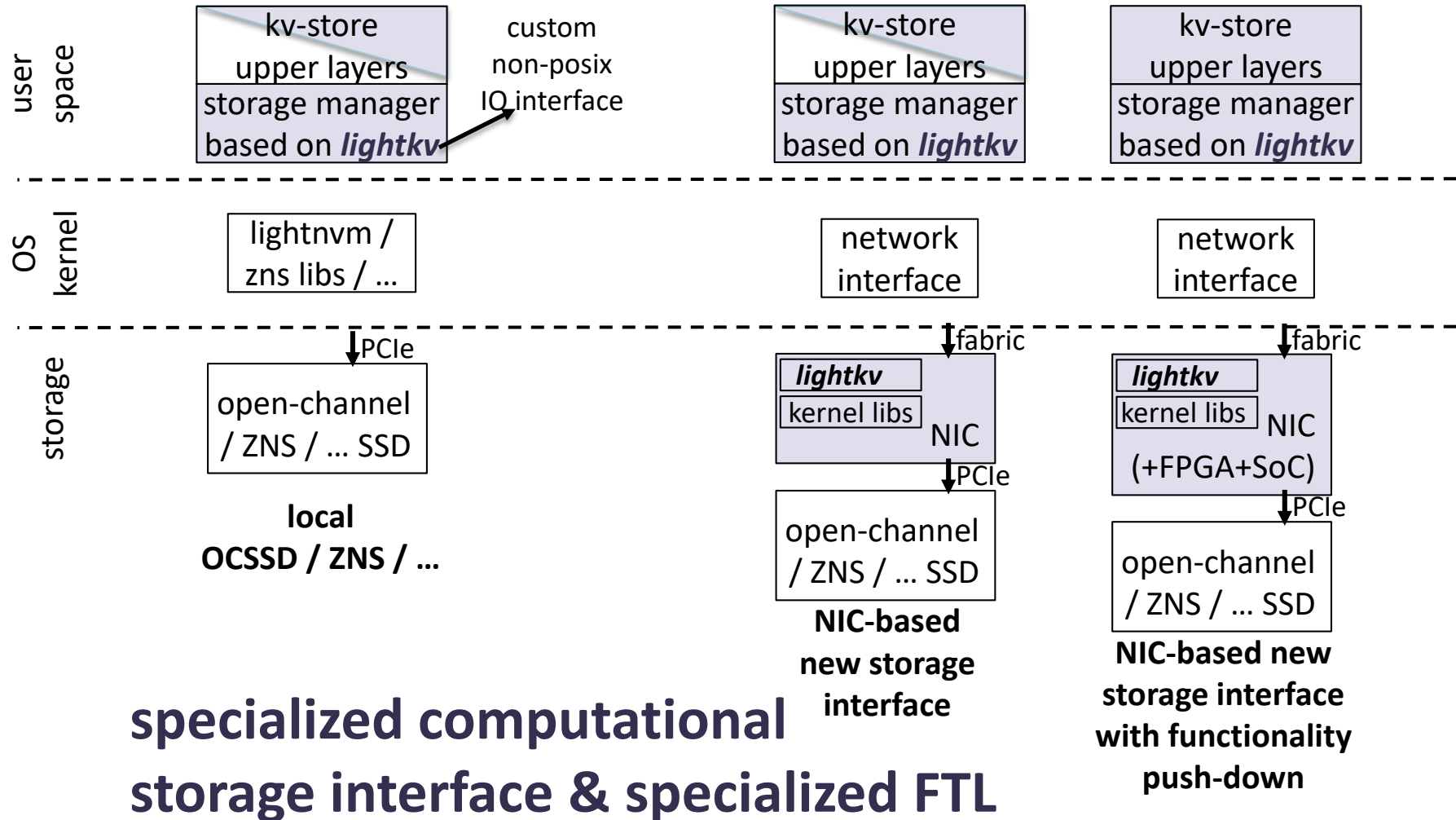https://github.com/DFC-OpenSource/

- Future disk controller is a super-computer.
  - » 1 bips processor
  - » 128 MB dram
  - » 100 GB disk plus one arm

# SSD landscape – local

**kv-store needs to change when you start app-specific storage management & pushing functionality down!**



|  | traditional SSD | OCSSD | KV-SSD | FPGA-based storage controller | local ZNS-SSD |
|---|---|---|---|---|---|
| **user space** | kv-store upper layers / storage manager based on POSIX | kv-store upper layers / storage manager based on *lightkv* | custom key-value store | kv-store upper layers / storage manager with push-down | kv-store upper layers / storage manager based on *lightkv* |
| **OS kernel** | filesystem / block layer / NVMe driver | lightnvm / NVMe driver | kv-ssd driver | smartssd driver | zns libs / zns driver |
| **storage** | PCIe — FTL / SSD | PCIe — open-channel SSD | PCIe — KV-FTL / SSD | PCIe — FPGA / FTL / SSD | PCIe — ZNS-FTL / ZNS-SSD |

**black-box**  **app-specific, no comp. storage**  **static custom**  **generic FTL**  **similar to OCSSD, less customizable**

# SSD landscape – disaggregated



**user space**

| kv-store | custom | kv-store | kv-store |
|---|---|---|---|
| upper layers | non-posix | upper layers | upper layers |
| storage manager | IO interface | storage manager | storage manager |
| based on *lightkv* | | based on *lightkv* | based on *lightkv* |

**OS kernel**

| lightnvm / zns libs / … | | network interface | network interface |

**storage**

↓PCIe     ↓fabric     ↓fabric

| open-channel / ZNS / … SSD | *lightkv* kernel libs NIC | *lightkv* kernel libs NIC (+FPGA+SoC) |

↓PCIe     ↓PCIe

**local OCSSD / ZNS / …**

| | open-channel / ZNS / … SSD | open-channel / ZNS / … SSD |

**NIC-based new storage interface**

**NIC-based new storage interface with functionality push-down**

**specialized computational storage interface & specialized FTL**

19

# programming SSDs

The VLDB Journal 2021

**Host-Based LS**

Host: LPID → SSDAdr
SSD: SSDAdr → FlashAdr

**SSD-Based LS**

LPID
MAP
LPID → FlashAdr
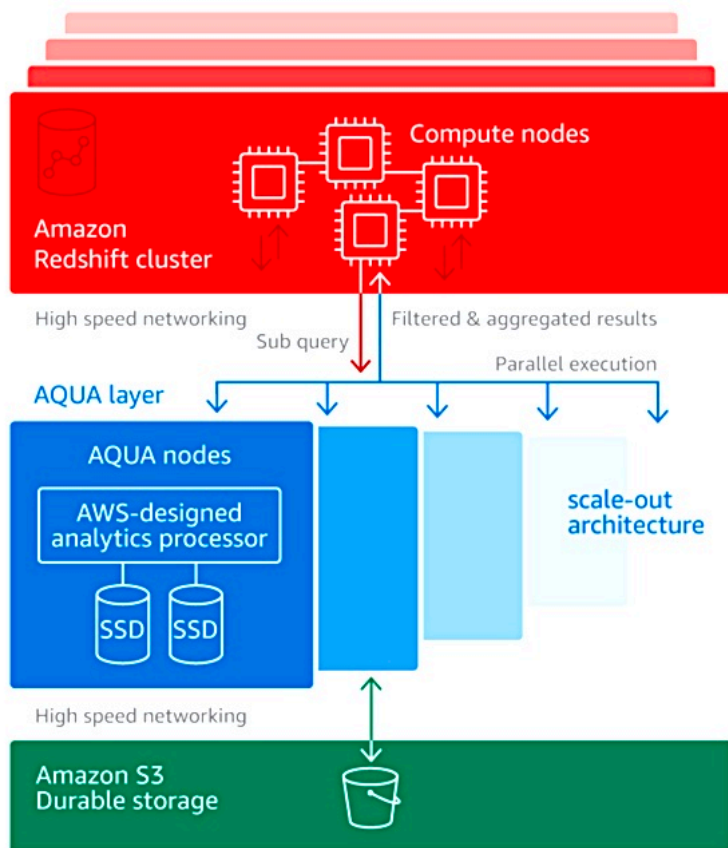
- Philippe Bonnet and Ivan Picoli in collaboration with MSR

- programming a storage controller using OX framework on an OCSSD



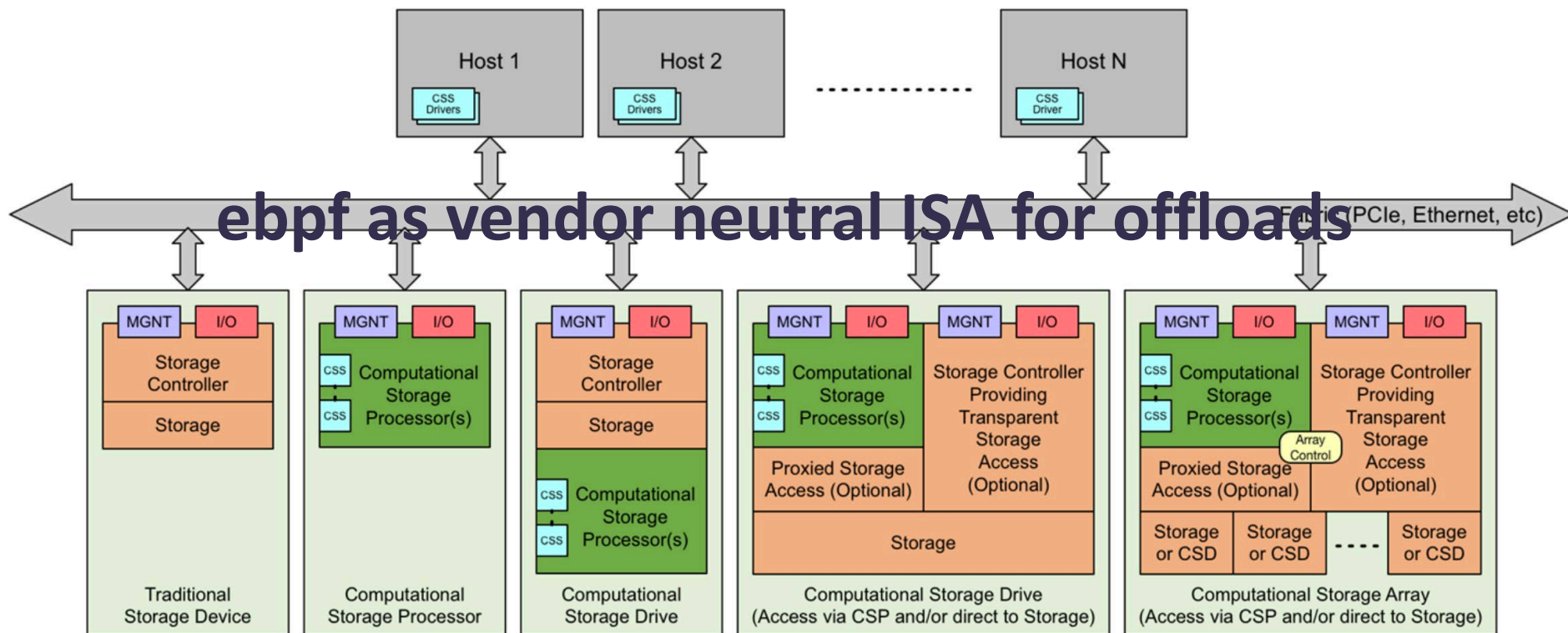SSD with transactional Batch I/O interface

**BwTree-specific FTL!**

# AWS AQUA

Advanced Query Accelerator
- near-data processing from AWS
(also called computational storage)
- announced in 2019
  (see video if interested)
- they are using SSDs and FPGAs
  at the AQUA layer
- goal: to reduce network traffic
  by reducing data movement

# envisioned architectures

SNIA. Computational Storage Architecture and Programming Model. V0.5, Rev 1. Aug 2020
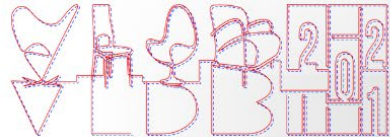


**ebpf as vendor neutral ISA for offloads**

**being standardized in NVMe (expected in 2022)**

# conclusion

- data management community increasingly shifts from pure in-memory optimized to SSD-optimized

- NVMe SSDs aren't a uniform class of devices

- expanding range of standardized storage interfaces (block, ZNS, KV, OCSSD)
  ➜ **the storage interface is a design choice**

- computational storage enables the definition of even more specialized storage interfaces

**need for co-design of storage engine – FTL – SSD**

VLDB 2021 looking for student volunteers, contact us if interested!

**47th International Conference on Very Large Data Bases**

Copenhagen, Denmark - August 16-20, 2021

COVID-19 - VLDB 2021 Hybrid is still on - <u>Read more</u>

➤ Get to attend the top international data management conference!
➤ Get insight into inner workings of a conference
➤ Contribute as virtual or on-site volunteer

You can help with
- Registration desk support
- Microphone duty for on-site discussions
- Registering participants in conference app
- Check program artefacts (videos, posters,…) –2-4 weeks prior to conference

Check out vldb.org/2021

Contact volunteer chair Ira Assent: ira@cs.au.dk