# Storage Systems (StoSys) XM_0092

## Lecture 11: CXL and io_uring

Animesh Trivedi

Autumn 2023, Period 1

VU
VRIJE
UNIVERSITEIT
AMSTERDAM

# Syllabus outline

1. ~~Welcome and introduction to NVM~~
2. ~~Host interfacing and software implications~~
3. ~~Flash Translation Layer (FTL) and Garbage Collection (GC)~~
4. ~~NVM Block Storage File systems~~
5. ~~NVM Block Storage Key-Value Stores~~
6. ~~Emerging Byte-addressable Storage~~
7. ~~Networked NVM Storage~~
8. ~~Trends: Specialization and Programmability~~
9. ~~Distributed Storage / Systems - I~~
10. ~~Distributed Storage / Systems - II~~
11. Emerging Topics ⬅

# Today is the last course lecture

We survived, it has been quite fun to teach this course

Hope you also had fun and learn a lot of advancements happening in the area of storage research

In coming days and weeks

- **Next Tuesday:** Milestone 5 interview - **sign up!**
- **Next Wednesday:** Guest Lecture from Nikolas
- **Afterwards:** Prepare for the exam - Good luck !
- **In the End:** We will ask for some feedback on the course
  - Me as a teacher
  - Broadly about the course - *you can be frank!*
  - ***Want to be the TA next year?***

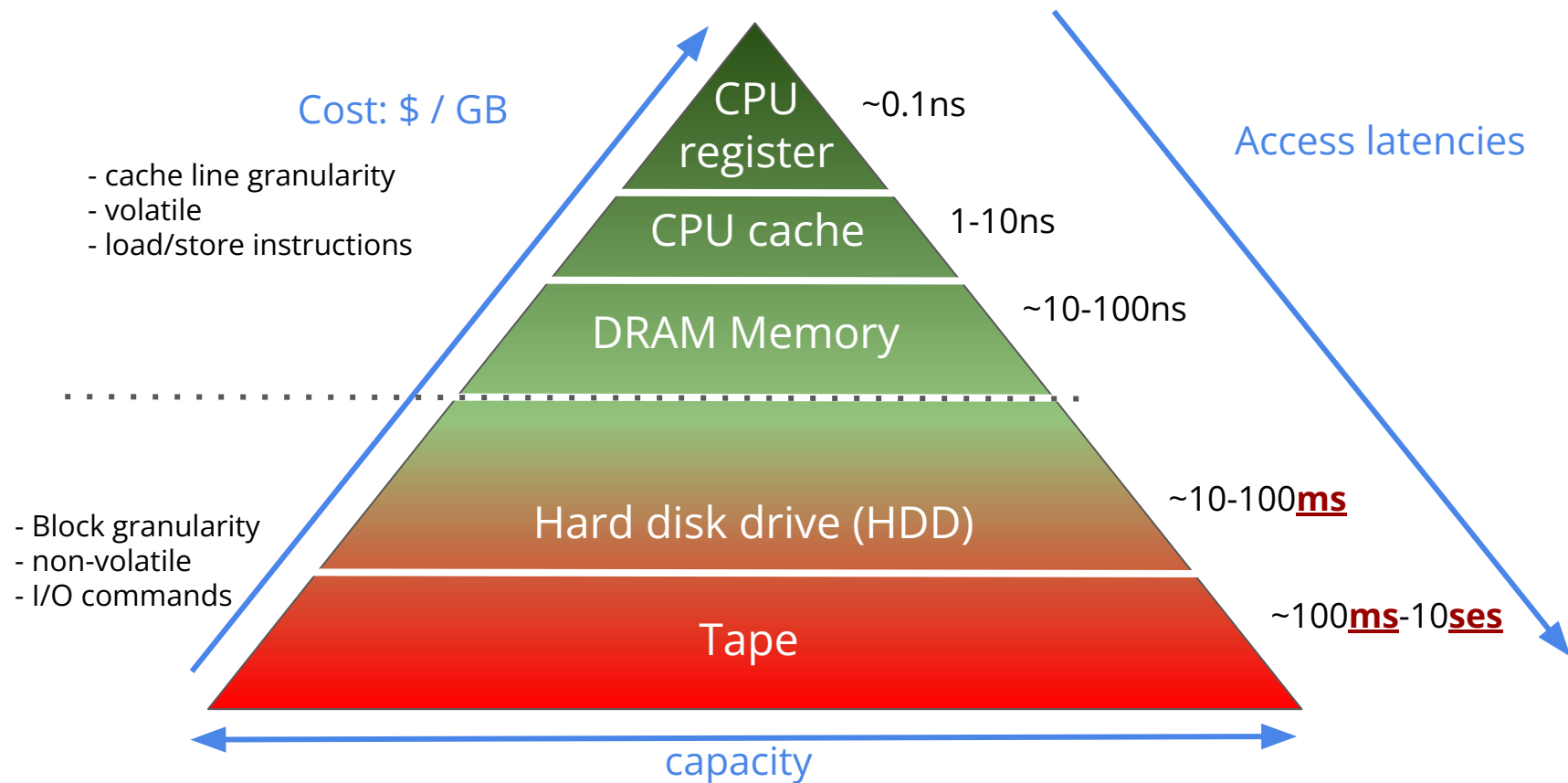# If you are interested in such research ...

Individual research projects (XM_405088)
- 6 or 12 ECTS credits
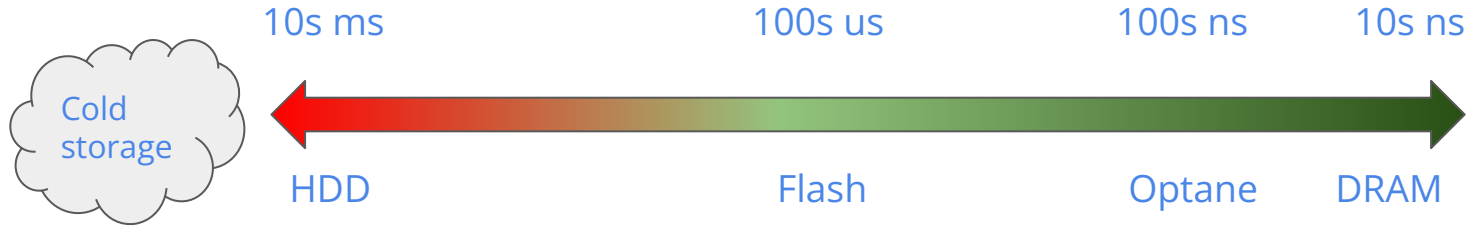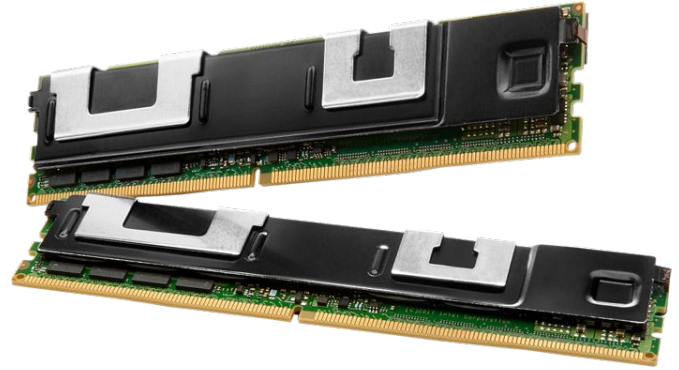
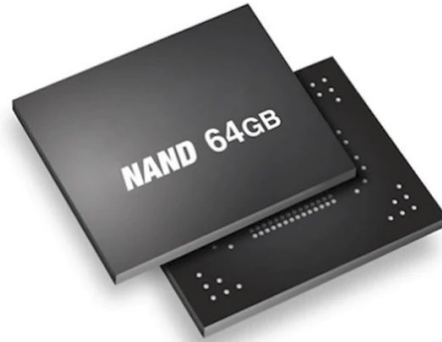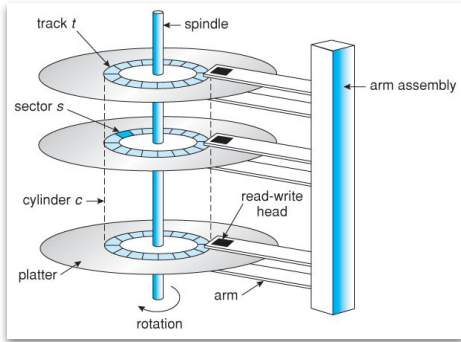Master projects / literature study

- Benchmarking the storage benchmarks
- io_uring/CXL research *(today's lecture)*
- Integrating NVM(e)/NVMoF storage in ML runtime to train large models
  (Swapping Tensors)
- Building computation storage device prototype in QEMU
- Virtualizing ZNS/NVMe devices
- Scheduling I/O operations for workload-specific optimizations
- *Your favorite idea ... I am broadly open to ideas from your side, pick a paper and lets discuss*
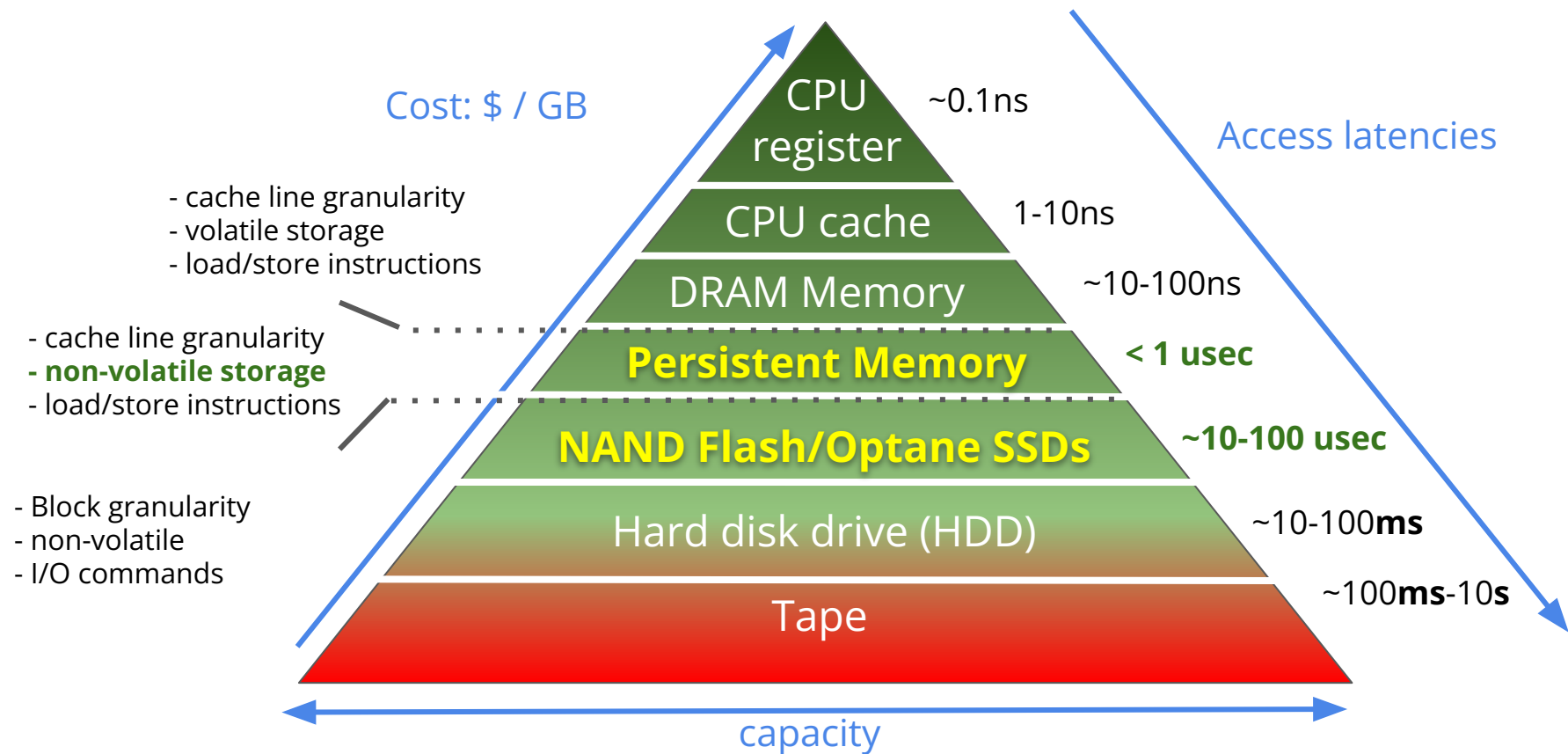
# The triangle of storage hierarchy

Cost: $ / GB

- cache line granularity
- volatile
- load/store instructions

Access latencies

CPU register  ~0.1ns

CPU cache  1-10ns

DRAM Memory  ~10-100ns

Hard disk drive (HDD)  ~10-100**ms**

Tape  ~100**ms**-10**ses**

- Block granularity
- non-volatile
- I/O commands

capacity

# Recap: From HDDs to Persistent Memories (PMem)



| 10s ms | 100s us | 100s ns | 10s ns |
|--------|---------|---------|--------|

Cold storage

| HDD | Flash | Optane | DRAM |
|-----|-------|--------|------|

# The (new) triangle of storage hierarchy

Cost: $ / GB

Access latencies

- cache line granularity
- volatile storage
- load/store instructions

- cache line granularity
- **non-volatile storage**
- load/store instructions

- Block granularity
- non-volatile
- I/O commands

**CPU register** — ~0.1ns

**CPU cache** — 1-10ns

**DRAM Memory** — ~10-100ns

**Persistent Memory** — < 1 usec

**NAND Flash/Optane SSDs** — ~10-100 usec

**Hard disk drive (HDD)** — ~10-100**ms**

**Tape** — ~100**ms**-10**s**

capacity

# **Multiple Emerging Topics (non-exhaustive)**

Domain-specific/specialized storage solutions

Storage virtualization, Disaggregation (end-to-end software-defined-*)

Quality-of-service in Storage Ecosystems (scheduling, multi-tenancy)

Energy Considerations

**CPU-free Computing** (re-thinking the computing architecture)
- [CPU-free Computing: A Vision with a Blueprint | Proceedings of the 19th Workshop on Hot Topics in Operating Systems](#)

**Hardware changes: Computer Express Link (CXL)**
- *Brief motivation and capabilities (without getting into too much hw/PCIe details)*

**New software APIs: `io_uring` (Linux, also being ported to other OSes)**
- *How is it different than other APIs and what options does it provide, performance implications*

# The Key Problems 1 / 2

**The CPU is the center of computing**

- direct memory access

- center of coherency

- controller of the devices

and the final coordinator and arbiter

**The CPU performance was fast!**



A more modern (simplified) setup



**Figure 2-1.** The organization of a simple computer with one CPU and two I/O devices.



Figure 1.4
**Hardware organization of a typical system.** CPU: Central Processing Unit, ALU: Arithmetic/Logic Unit, PC: Program counter, USB: Universal Serial Bus.

# The Key Problems 1 / 2

Compute offload
(ASIC, DSPs, FP)

DRAM memory

Ethernet, WiFi

Disk storage

CPU

# The Key Problems 1 / 2



**CPU cache management is non-trivial and complex
(even with same/similar homogeneous CPU architectures)**

# The Key Problems 1 / 2

Compute offload
(ASIC, DSPs, FP)

SW

SW

SW

DRAM memory

Ethernet, WiFi

CPU

Disk storage

# The Key Problems 1 / 2



**DRAM memory**

**SW**

**SW**

**SW**

CPU

**Ethernet, WiFi**

**Disk storage**

# The Key Problems 1 / 2



Elba (7 nm) Block Diagram

SW

SW

SW

CPU

DRAM memory

Disk storage

14

# The Key Problems 1 / 2



Elba (7 nm) Block Diagram

Cost-effective, Energy-efficient, and **Scalable Storage Computing** for Large-scale AI Applications. ACM Trans. Storage 16, 4, Article 21 (November 2020), 37 pages.
https://doi.org/10.1145/3415580

# The Key Problems 1 / 2



Elba (7 nm) Block Diagram

**SW**

**SW**

**SW**

CPU

*How are these two caches synchronized?*

# The Key Problems 1 / 2

**These accelerators can have :**

- Compute elements (specialized - FPGA, or general - ARM)
- Memory elements
- Storage chips
- Multi-level caches
- Outside connectivity

*Who manages "coherency", "data flow", "configuration", "management" of memories/caches/devices here? Software, hardware? Performance? Cost of development of new APIs, protocols?*

Elba

ARM®
ortex™-M7
back-end)

Flash memory
(channel #1)

Flash memory
(channel #2)

Flash memory
(channel #16)

# The Key Problems 2 / 2 : CPU - DRAM Coupling

malloc (3GB)

**SW**

malloc (1GB)

**SW**

malloc (2GB)

**SW**

malloc (0.5GB)

**SW**

CPU

**DRAM memory 8 GB / DIMM**

- **What happens to the remaining 1.5 GB DRAM?**
- **Do applications use all the DRAM what they ask for?**

# The Key Problems 2 / 2 : CPU - DRAM Coupling

1. Can not mix and match different DRAM technologies and generations
2. More performance means more capacity (need to buy more DIMMs)
3. Limit to how much DRAM can be packed in a single machine



**Very close coupling of CPU-DRAM (1) DRAM technology; (2) Density, capacity; and (3) Performance**

# The Key Problems 2 / 2 : CPU - DRAM Coupling





Figure 2: Memory stranding (§3.1). Stranding increases significantly as more CPU cores are scheduled. Error bars indicate the $5^{th}$ and $95^{th}$ percentiles (outliers in dots).

DRAM is a big power and cost factor in data center (up to ~40%)
A big part can remain underutilized
Azure with VMs : on average ~10% (but as high as ~30%)

Pond: CXL-Based Memory Pooling Systems for Cloud Platforms, ASPLOS 2023, https://doi.org/10.1145/3575693.3578835
TPP: Transparent Page Placement for CXL-Enabled Tiered-Memory. ASPLOS 2023, https://doi.org/10.1145/3582016.3582063

# The Key Problems 2 / 2 : CPU - DRAM Coupling



Figure 7: Application memory usage over last N mins.



Figure 11: Fraction of pages re-accessed at different intervals.

Not all pages allocation are used **<u>uniformly</u>:**

(1)  Only a small fraction of memory is <u>accessed</u> in 1-2 minutes window
(2)  For Web, almost 80% of the pages are <u>re-accessed</u> within a ten-minute interval but for warehouse it is 20%.
     *(do they all have to be in DRAM?)*

# Summary Problem

There has to be a better way to

- Manage non-CPU memories and caches (accelerators)
- Manage CPU-attached memories (allocation, disaggregate from the CPU)
- Expand beyond the CPU-attached memories

**+  Think of non-volatile memories …**
  - Persistent memories
  - Fast storage

Solution : **Compute Express Link (CXL)** *(the last protocol we will ever need)*

# Computer Express Link (CXL)



**A cache coherent Interconnect** between
- The CPU
- Accelerators
- Memory expansion cards

**Asymmetric protocol**

A set of standardized protocols defined on the top of PCIe 5.0 (PHY)
- Runs in the standard PCIe slots
- 32 GT/s, or 4 GB/lane ⇒ x32 card = **128 GB/sec**
- *Latencies approaching the NUMA CPU (with v6.0)*

| PCIe Specification | Data Rate per Lane (GT/s) | Encoding | x16 Unidirectional Bandwidth (GB/s) | Specification Ratification Year |
|---|---|---|---|---|
| 1.x | 2.5 | 8b/10b | 4 | 2003 |
| 2.x | 5 | 8b/10b | 8 | 2007 |
| 3.x | 8 | 128b/130b | 15.75 | 2010 |
| 4.0 | 16 | 128b/130b | 31.5 | 2017 |
| 5.0 | 32 | 128b/130b | 63 | 2019 |
| 6.0 | 64 | PAM4/FLIT | 128 | 2022 |

https://www.electronicdesign.com/technologies/embedded/article/21162617/cxl-coherency-memory-and-io-semantics-on-pcie-infrastructure
https://www.xda-developers.com/pcie-5/
https://www.rambus.com/blogs/pcie-6/

# Three CXL Protocols

**CXL.io**
- Mandatory for all hosts, and CXL supported devices
- Discovery, enumerations, capabilities (DMA, interrupts, IOV), and host physical address configuration
- Same in spirit to what any basic PCIe device would support

**CXL.mem**
- Enables (only) CPU to access device/accelerator memory in a cacheable manner
- Useful in DRAM expansion
- Device is not initiating any communication

**CXL.cache**
- The same as CXL.mem, but now devices can also access the CPU memory/caches
- Additional commands/requests for maintaining coherence among <u>all</u> copies

# **Three Classes of Devices**

# Three Generations of CXL Protocols

| Features | CXL 1.0 / 1.1 | CXL 2.0 | CXL 3.0 |
|---|---|---|---|
| Release date | 2019 | 2020 | 1H 2022 |
| Max link rate | 32GTs | 32GTs | 64GTs |
| Flit 68 byte (up to 32 GTs) | ✓ | ✓ | ✓ |
| Flit 256 byte (up to 64 GTs) | | | ✓ |
| Type 1, Type 2 and Type 3 Devices | ✓ | ✓ | ✓ |
| Memory Pooling w/ MLDs | | ✓ | ✓ |
| Global Persistent Flush | | ✓ | ✓ |
| CXL IDE | | ✓ | ✓ |
| Switching (Single-level) | | ✓ | ✓ |
| Switching (Multi-level) | | | ✓ |
| Direct memory access for peer-to-peer | | | ✓ |
| Enhanced coherency (256 byte flit) | | | ✓ |
| Memory sharing (256 byte flit) | | | ✓ |
| Multiple Type 1/Type 2 devices per root port | | | ✓ |
| Fabric capabilities (256 byte flit) | | | ✓ |

- CXL 3.0: Enabling composable systems with expanded fabric capabilities, October 6, 2022,
  https://www.computeexpresslink.org/_files/ugd/0c1418_998df4f459734f319e7a12cc2163b943.pdf
- Good overview, https://community.cadence.com/cadence_blogs_8/b/breakfast-bytes/posts/hot-chips-cxl-tutorial

Evolving Use Cases

**What can we do?** Expansion of DRAM, CPU-Memory Decoupling (multiple generation of devices), Memory Pooling and sharing, Single Logical Device (SLD → Exclusive to one CXL root) to Multiple Logical Device (MLD, connected to multiple CXL roots), Memory hot swapping ...

A look into the CXL device ecosystem and the evolution of CXL use cases,
https://0c141887-fbe4-4ec3-be17-adc8d70d3922.usrfiles.com/ugd/0c1418_037d4ba31f4b44cf9fcb37f5b36ae4d6.pdf

# Design a Distributed Cluster Running CXL



CXL 3.0 Fabric Architecture
- Interconnected Spine Switch System
- Leaf Switch NIC Enclosure
- Leaf Switch CPU Enclosure
- Leaf Switch Accelerator Enclosure
- Leaf Switch Memory Enclosure

*Multiple type of devices, Global Fabric Attached Memory (GFAM)*

CXL 3.0: Enabling composable systems with expanded fabric capabilities October 6, 2022
https://www.computeexpresslink.org/_files/ugd/0c1418_998df4f459734f319e7a12cc2163b943.pdf

28

# CXL.mem

Hello bytes, bye blocks: PCIe storage meets compute express link for memory expansion (CXL-SSD). https://doi.org/10.1145/3538643.3539745      29

# CXL.mem Expansion Device Example



1. PCIe enumeration and BAR mapping with, Host-Managed Device Memory (HDM) areas
2. Setup MMU and allocate the DRAM physical address from this area (software support)
3. Access happens, and the request is routed to the PCIe/CXL root

# CXL.mem Expansion Device Example



**Device Physical Address (DPA)**

*DRAM Translation Layer DTL ;)*
*See the ISCA'23 reference at the end of the slides*

**Multiple configurations** (1) striping across multiple devices, ports, roots; (2) allocation units…

# Transparent Page Placement (TPP)



(a) Without CXL    (b) With CXL

PCIe 6.0 latencies and bandwidth are approaching access to a remote NUMA CPU socket

**Challenge**: How to profile pages (at low-overheads) and put them in the right storage level in the CXL-enabled memory hierarchy

Pond: CXL-Based Memory Pooling Systems for Cloud Platforms, ASPLOS 2023, https://doi.org/10.1145/3575693.3578835
TPP: Transparent Page Placement for CXL-Enabled Tiered-Memory. ASPLOS 2023, https://doi.org/10.1145/3582016.3582063

# POND (ASPLOS'23): How to Disaggregate VM Memory

# Where does **Storage** Come into the Play?

*Any device can implement the CXL protocol*

- Use SSD as large capacity RAM
- **Byte*-addressable**
- Persistent

*64B addressable



Industry 1st CXL-based Storage Optimized for AI/ML



**Dual Mode Support**

Application · AI/ML embedding table lookups · Caching · Graph analytics

Linux

LBA | Load/Store

CXL.io — File system-based access supports legacy NVMe

CXL.mem — Load/store access for memory-mapped files

CXL Interface

Cache CTRL → DRAM

NAND

**Persistent Memory Mode**

Application · Metadata, · low-latency IO · High Priority · Indexing/Logging

Linux | RoCE

Load/Store

CXL.mem — Load/store access to DRAM Persistency cia flushes to NAND

CXL Interface ← External Battery

DRAM

NAND

# Emerging work: Quantifying and Hiding Flash Latencies

## Paper 1

**Hello Bytes, Bye Blocks: PCIe Storage Meets Compute Express Link for Memory Expansion (CXL-SSD)**

Myoungsoo Jung

*Computer Architecture and Memory Systems Laboratory,*
Korea Advanced Institute of Science and Technology (KAIST)
http://camelab.org

**ABSTRACT**

Compute express link (CXL) is the first open multi-protocol method to support cache coherent interconnect for different processors, accelerators, and memory device types. Even though CXL manages data coherency mainly between CPU memory spaces and memory on attached devices, we argue that it can also be useful to reform existing block storage as cost-efficient, large-scale working memory. Specifically, this paper examines three different sub-protocols of CXL from a memory expander viewpoint. It then suggests which device type can be the best option for PCIe storage to bridge its block semantics to memory-compatible, byte semantics. We then discuss how to integrate a storage-integrated memory expander into an existing system and speculate how much effect it does have on the system performance. Lastly, we visit various CXL network topologies and explore a new opportunity to efficiently manage the storage-integrated, CXL-based memory expansion.

## Paper 2

**Cache in Hand: Expander-Driven CXL Prefetcher for Next Generation CXL-SSDs**

Miryeong Kwon[*†], Sangwon Lee[*†], Myoungsoo Jung[*†]
[*]*Computer Architecture and Memory Systems Laboratory*, KAIST
[†]Panmnesia, inc.

**ABSTRACT**

Integrating compute express link (CXL) with SSDs allows scalable access to large memory but has slower speeds than DRAMs. We present ExPAND, an expander-driven CXL prefetcher that offloads last-level cache (LLC) prefetching from host CPU to CXL-SSDs. ExPAND uses a heterogeneous prediction algorithm for prefetching and ensures data consistency with CXL.mem's back-invalidation. We examine prefetch timeliness for accurate latency estimation. ExPAND, being aware of CXL multi-tiered switching, provides end-to-end latency for each CXL-SSD and precise prefetch timeliness estimations. Our method reduces CXL-SSD reliance and enables direct host cache access for most data. ExPAND enhances graph application performance by 3.5×, surpassing CXL-SSD pools with diverse prefetching strategies.

## Paper 3

**Overcoming the Memory Wall with CXL-Enabled SSDs**

Shao-Peng Yang, *Syracuse University*; Minjae Kim, *DGIST*; Sanghyun Nam, *Soongsil University*; Juhyung Park, *DGIST*; Jin-yong Choi, *FADU Inc.*; Eyee Hyun Nam, *FADU Inc.*; Eunji Lee, *Soongsil University*; Sungjin Lee, *DGIST*; Bryan S. Kim, *Syracuse University*

**Abstract**

This paper investigates the feasibility of using inexpensive flash memory on new interconnect technologies such as CXL (Compute Express Link) to overcome the memory wall. We explore the design space of a CXL-enabled flash device and show that techniques such as caching and prefetching can help mitigate the concerns regarding flash memory's performance and lifetime. We demonstrate using real-world application traces that these techniques enable the CXL device to have an estimated lifetime of at least 3.1 years and serve 68–91% of the memory requests under a microsecond. We analyze the limitations of existing techniques and suggest system-level changes to achieve a DRAM-level performance using flash.
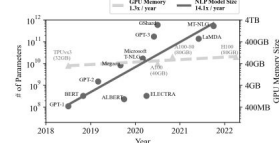
Figure 1: The trend in memory requirements for NLP applications [11, 30, 34, 43]. The number of parameters increases by a factor of 14.1× per year, while the memory capacity in GPUs only grows by a factor of 1.3× every year.

# Putting SSDs with CXL Memory Expander

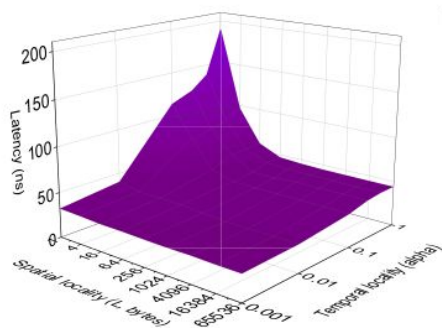Which type of device to use? Type-1, Type-2, or Type-3 when using SSD as memory expander?

**Type-3:**
- *(in CXL 1.0, 2.0)*: Only one Type-1 or Type-2 device allowed per CXL root, hence Type-3 are more scalable.
- Type-1/2 can be more complex, caches, all load/store requests require checking the cache states of PCIe storage computing complex
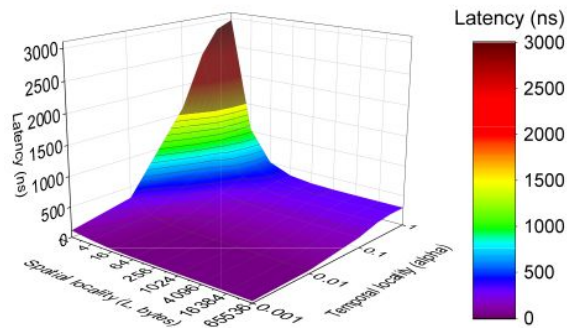
**Hence, a Type-3 device type is the ideal CXL device for a "memory expander"**

Hello bytes, bye blocks: PCIe storage meets compute express link for memory expansion (CXL-SSD). https://doi.org/10.1145/3538643.3539745
Cache in Hand: Expander-Driven CXL Prefetcher for Next Generation CXL-SSD. https://doi.org/10.1145/3599691.3603406

# CXL + Flash SSDs: Can Flash do it?



(a) LocalDRAM.    (b) CXL-SSD.

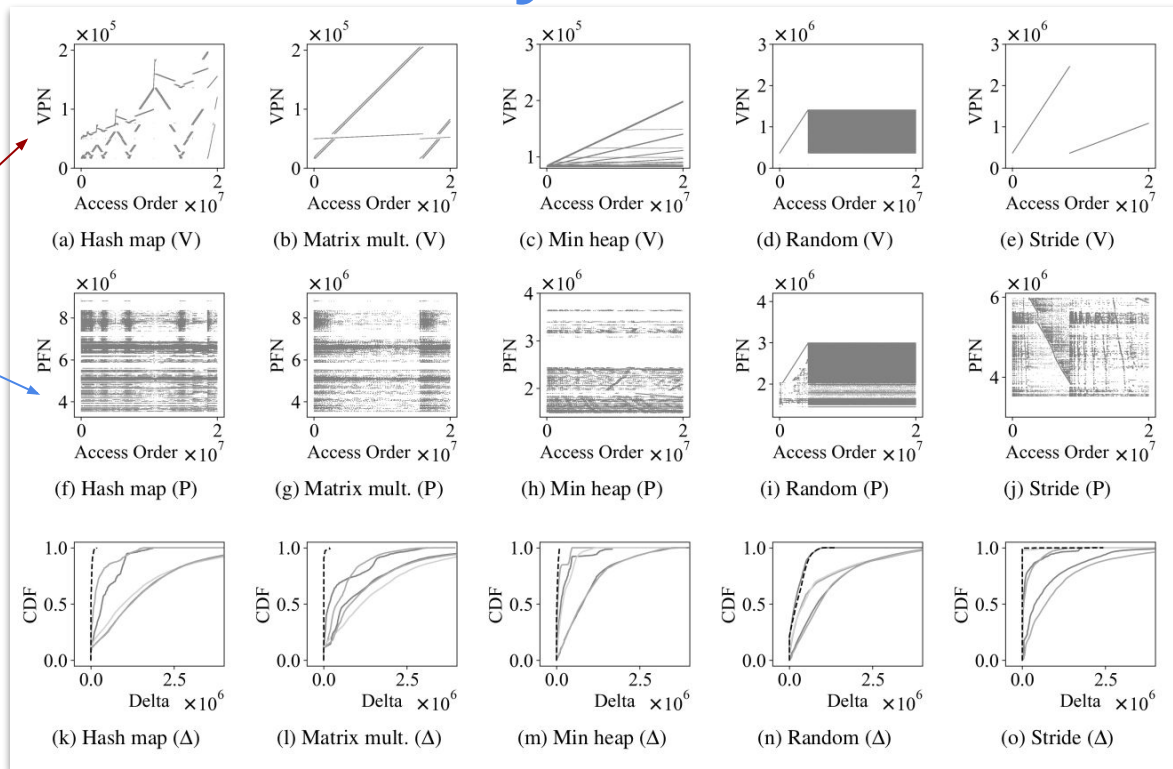**Can we use NAND flash SSDs as memory expander?**

- What latencies one get with the granularity mismatch?
  - Cache line : 64B, flash pages : 8-16 KiB
  - DRAM: 100s of nanoseconds, vs. flash in 10-100 microseconds
- What is the access pattern for common workloads?
- Can we optimize latencies in any manner? Prefetching, buffering, caching?
- How about flash P/E limitations? Can it endure small 64B writes?

# CXL-Enabled SSDs - Virtual vs. Physical Addresses

⇒ Shows that the access pattern at
the virtual address level
do not correspond to
the physical address level

**Why?**
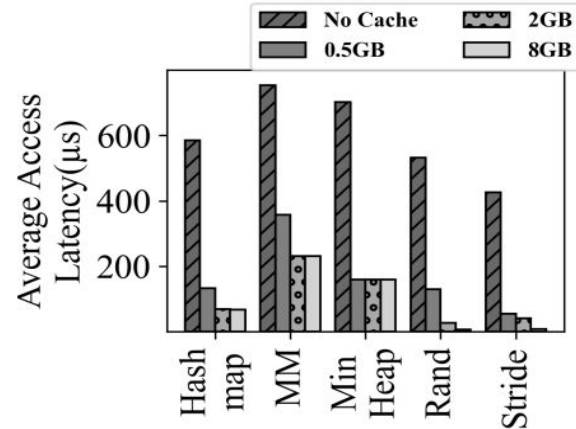
Just basic prefetching is not effective to hide latencies



(a) Hash map (V)  (b) Matrix mult. (V)  (c) Min heap (V)  (d) Random (V)  (e) Stride (V)

(f) Hash map (P)  (g) Matrix mult. (P)  (h) Min heap (P)  (i) Random (P)  (j) Stride (P)

(k) Hash map (Δ)  (l) Matrix mult. (Δ)  (m) Min heap (Δ)  (n) Random (Δ)  (o) Stride (Δ)

# Impact of Caching

Underline{Inter-arrival time} of 64B requests
has a huge impact

- Queuing delays w/o cache
- Small amount of cache helps (0.5GB)



(a) Average access latency



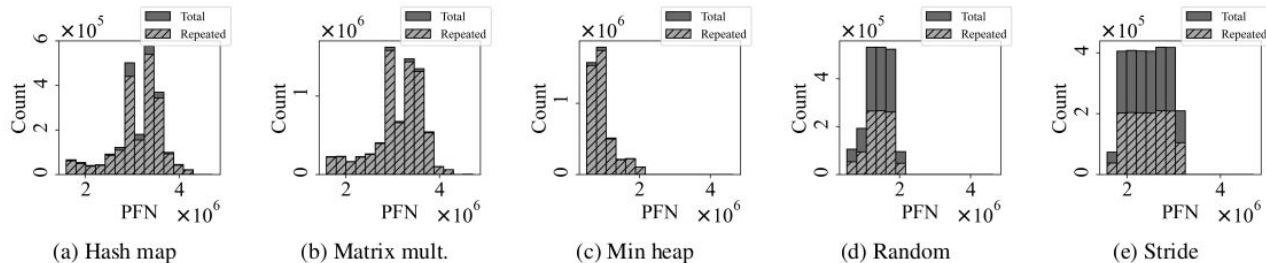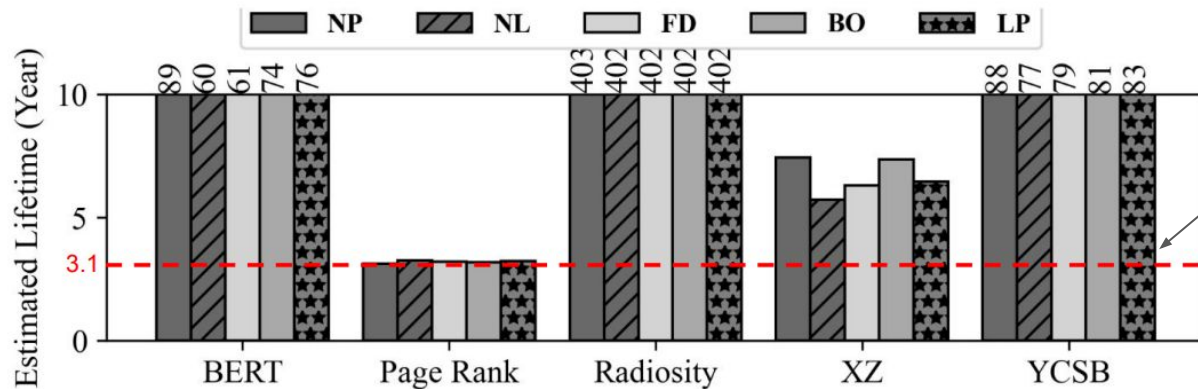(a) Hash map   (b) Matrix mult.   (c) Min heap   (d) Random   (e) Stride

Figure 6: Flash memory read count for physical memory frames. The solid bar represents the total number of reads, while the shaded bar, the number of repeated reads. A repeated read is a read request to an outstanding read request.
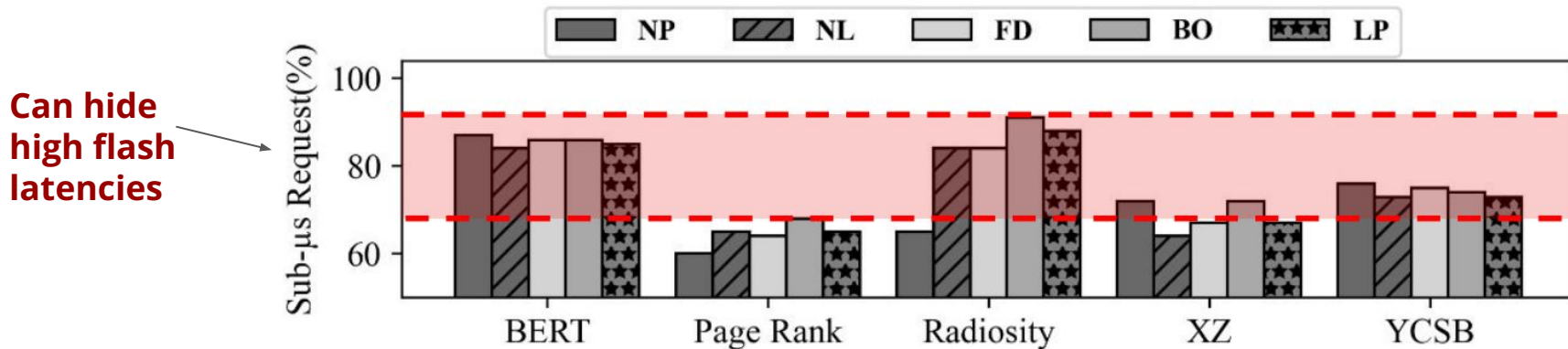
Lots of repeated accesses for the same page!

Multiple 64B requests go into the same flash page **(Keep track of it)**
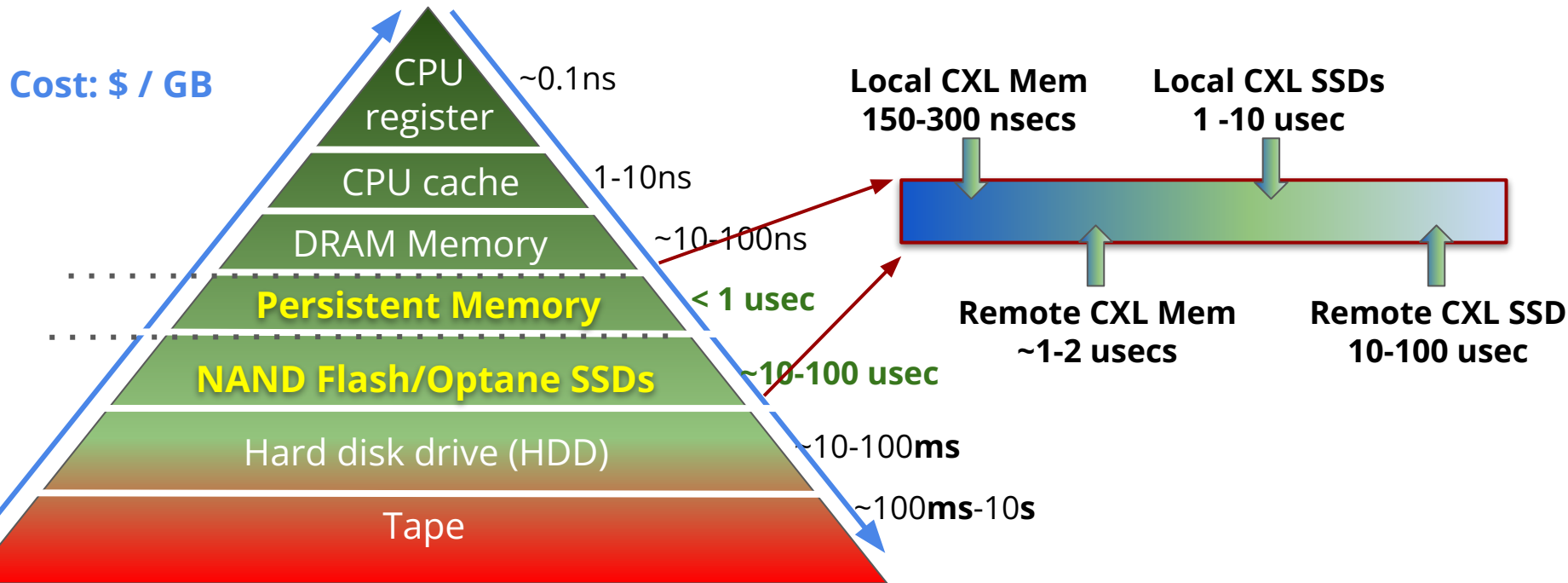
# Workload-level Performance



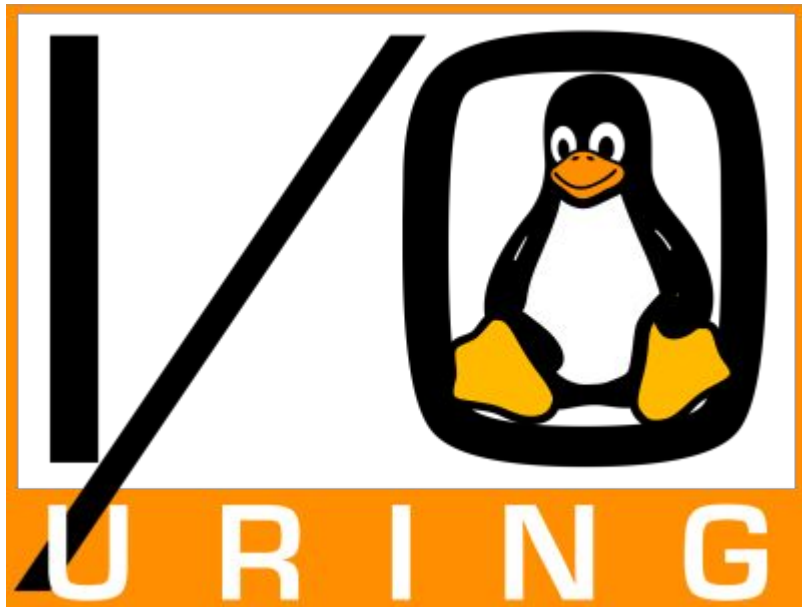**3.1 years (min) with the workloads**
SLC flash, 100K P/E)

Can hide high flash latencies

# The New(er) Triangle of Storage-Memory Continuum

Cost: $ / GB

CPU register — ~0.1ns

CPU cache — 1-10ns

DRAM Memory — ~10-100ns

Persistent Memory — < 1 usec

NAND Flash/Optane SSDs — ~10-100 usec

Hard disk drive (HDD) — ~10-100**ms**

Tape — ~100**ms**-10**s**

**Local CXL Mem 150-300 nsecs**

**Local CXL SSDs 1 -10 usec**

**Remote CXL Mem ~1-2 usecs**

**Remote CXL SSD 10-100 usec**

**Instead of discrete steps, it is a continuous spectrum now: Continuum**

# io_uring : What is it and why you should care?



What is io_uring?, https://unixism.net/loti/what_is_io_uring.html

# The Long Debate: How to get Concurrency?

**Threads** versus **Events** (Asynchronous)

Blocking I/O

Asynchronous I/O

**Non-Blocking I/O** and *Asynchronous* **I/O** are two different things!

`while(not_done!)`

# Linux I/O Options



Standard POSIX I/O **blocking** read/write calls:
- https://man7.org/linux/man-pages/man2/read.2.html
- https://man7.org/linux/man-pages/man2/write.2.html

Make I/O calls **non-blocking** : set O_NONBLOCK flag on the file descriptor
- https://man7.org/linux/man-pages/man2/fcntl.2.html (`O_NONBLOCK)`

**Asynchronous I/O** on Linux : libaio and POSIX AIO
- https://github.com/littledan/linux-aio
- Example of how to use libaio: https://github.com/axboe/fio/blob/master/engines/libaio.c

Asynchronous programming. Blocking I/O and non-blocking I/O, https://luminousmen.com/post/asynchronous-programming-blocking-and-non-blocking

# AIO Issues

SIGNAL based delivery of completion
- Preemption and context switch
- Needs care for signal-safe function execution

Linux' AIO works truly "asynchronously" under very restricted conditions:

- works only with O_DIRECT modes (alignment, and size restrictions)
- works only when the file's metadata is available
        (otherwise blocks until the metadata is fetched)
- can block based on device's queue capacity
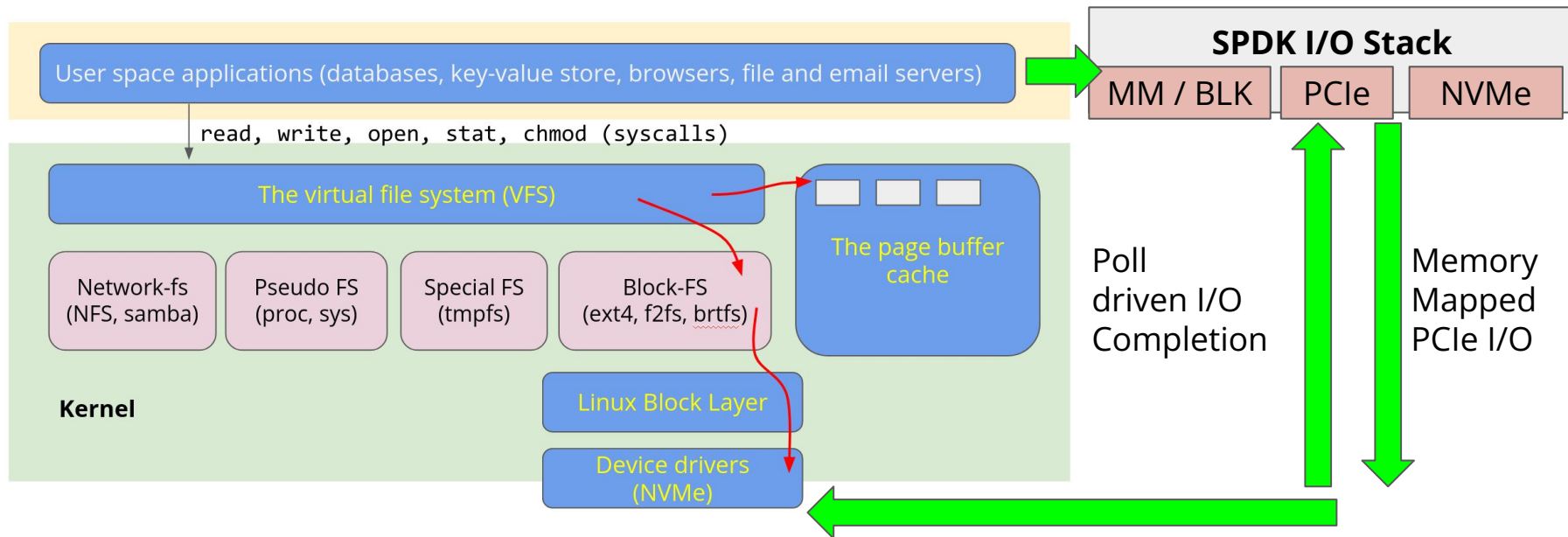- needs to memcpy of I/O metadata (~100 bytes)

Good introduction: https://unixism.net/loti/async_intro.html and https://kernel.dk/io_uring.pdf

Signal Handling in a Multi-Threaded Application in Linux, https://www.baeldung.com/linux/signals-multi-threaded-app
Re: [PATCH 09/13] aio: add support for async openat(), https://lwn.net/Articles/671657/

Archive-link: Article, Thread

On Mon, Jan 11, 2016 at 2:07 PM, Benjamin LaHaise <bcrl@kvack.org> wrote:
> Another blocking operation used by applications that want aio
> functionality is that of opening files that are not resident in memory.
> Using the thread based aio helper, add support for IOCB_CMD_OPENAT.

So I think this is ridiculously ugly.

AIO is a horrible ad-hoc design, with the main excuse being "other,
less gifted people, made that design, and we are implementing it for
compatibility because database people - who seldom have any shred of
taste - actually use it".

But AIO was always really really ugly.

# Cost of these Interfaces

## TABLE I: Categories of system-call techniques

| Kind | Mechanism | Examples | per sys request | | cost[ns] |
| --- | --- | --- | --- | --- | --- |
| | | | traps | csw | |
| Sync | Blocking | read(), write() | 1 | 2 | $955 \pm 1069$ |
| Sync | Non-Blocking | SOCK_NONBLOCK & epoll() | $[1, 3]$ | $[2, 6]$ | $1656 \pm 1318$ |
| Async | Callback | POSIX AIO [13] | 1 | 2, 3 | $6224 \pm 12\,232$ |
| Async | Queue-based | Linux AIO | $]0, 2]$ | $]1, 4]$ | $1922 \pm 1467$ |

# Skip the OS Complexity: The SPDK Stack

**SPDK I/O Stack**

| MM / BLK | PCIe | NVMe |

User space applications (databases, key-value store, browsers, file and email servers)

`read, write, open, stat, chmod (syscalls)`

The virtual file system (VFS)

The page buffer cache

| Network-fs (NFS, samba) | Pseudo FS (proc, sys) | Special FS (tmpfs) | Block-FS (ext4, f2fs, brtfs) |

**Kernel**

Linux Block Layer

Device drivers (NVMe)

Poll driven I/O Completion

Memory Mapped PCIe I/O

- A user-space I/O framework for NVMe devices (only)
- Block-level abstraction (no file system, but there are research prototypes)
- Has user-space mapped drivers (https://spdk.io/doc/userspace.html)
- Designed for light-weight I/O, best performance (eschews many core OS features)

# SPDK can have the Highest Performance

**6 Million Ops/core**

**~30 Million Ops/server**



SPDK NVMe BDEV IOPS
4KiB Random Read, QD=192, 1 CPU Core



4KiB Bdevperf Core Scaling Performance

2 CPU sockets, Intel(R) Xeon(R) Gold 6348 CPU @ 2.60GHz
22x Kioxia® KCM61VUL3T20 3.2TBs (FW: 0105) (10 on CPU NUMA Node 0, 12 on CPU NUMA Node 1)

SPDK NVMe BDEV Performance Report Release 23.05, June 2023,
https://ci.spdk.io/download/performance-reports/SPDK_nvme_bdev_perf_report_2305.pdf

# Intricately Linked Issues

What is the system call interface

What is the kernel threading model

Signal vs queuing

What is the cost of scheduling, context switching

Management of concurrency

Programming languages (error handling)

49

# Background Reading on this Topic



50

# Storage APIs: Recap

| Applications |
|:---:|

syscalls

| Linux Kernel |
|:---:|

| NVMe Device |
|:---:|

| Applications |
|:---:|

SPDK

| NVMe Device |
|:---:|

| Applications |
|:---:|

SQ   CQ

| Linux Kernel |
|:---:|

| NVMe Device |
|:---:|

**Libaio:**
- + Async I/O
- + Any files/FSes
- + Any device: HDD, NVMe

- − Async only with direct I/O
- − Performance
- − Metadata management

**SPDK:**
- + Performance
- + Close application integration
- + No syscall or interrupts

- − Only NVMe
- − No kernel assistance
- − Scalability and brittle

`io_uring`

**Best of both worlds?**

# io_uring: A Structured Approach to Asynchronous I/O



**Producer-consumer pattern**
- SQ: producer = application (tail), consumer = kernel (head)
- CQ: producer = kernel (tail), consumer = application (head)

Head and tail pointers manipulation with exclusive write ownership

What's new with io_uring, 2022, https://kernel.dk/axboe-kr2022.pdf    52

# io_uring: A Structured Approach to Asynchronous I/O



Application

1

2
**syscall**
io_uring_enter

2

SQ

CQ

7

6 **Kernel**

3

io_uring

5

4

*interrupt*

```
Request:
    ● File descriptor
    ● Offset
    ● Size
    ● (also vector)
    ● ...
```

```
Response:
    ● I/O status
    ● Context
      (user-defined)
    ● Size of the I/O
```

Applications can
- **Async I/O**
- I/O on any fd type (+net)
- Queue requests (batch)
- Vector I/O
- Optimize (fixed FD, pin)

53

# The three new Syscalls

1.  **io_uring_setup:** This call is for creating the ring structure (queue-depth, I/O completion and notification modes)
    a.  <u>Completion</u> polling by the kernel on the device (IORING_SETUP_**IOPOLL)**
    b.  Kernel polling for <u>submission </u>(IORING_SETUP_**SQPOLL**, zero system call)

2.  **io_uring_enter:** This call enters the kernel and tells it to process I/O requests (any type and extensible, not just storage I/O)
    a.  Networking, ZNS, Programmable storage and more
    b.  Replacement for the ioctl() call: a private interface between a device driver and application

3.  **io_uring_register**: This call is for registering specific fd, buffers, file ranges that are being used frequently to put them on an optimized fast path

See: `/usr/include/linux/io_uring.h` file for the full structure and call definitions

# Three Modes of Operations



**(a)** io_uring (default)

**(b)** with completion polling

**(c)** with submission polling

55

Diego Didona, Jonas Pfefferle, Nikolas Ioannou, Bernard Metzler, and Animesh Trivedi. 2022. Understanding modern storage APIs: a systematic study of libaio, SPDK, and io_uring. In Proceedings of the 15th ACM International Conference on Systems and Storage (SYSTOR '22). https://doi.org/10.1145/3534056.3534945

Zebin Ren and Animesh Trivedi. 2023. Performance Characterization of Modern Storage Stacks: POSIX I/O, libaio, SPDK, and io_uring. In Proceedings of the 3rd Workshop on Challenges and Opportunities of Efficient and Performant Storage Systems (CHEOPS '23). https://doi.org/10.1145/3578353.3589545

# Benchmarking Setup

**Setup 1 [Systor'22]:**

- 2x Intel® Xeon® E5-2630 (Sandy Bridge), 10 cores/socket ⇒ 20 CPU cores

- 20 Intel® DC P3600 400GB NVMe <u>Flash SSDs</u> ⇒ ~6 Million IOPS

**Setup 2 [CHEOPS'23]:**

- 2x Intel® Xeon® Silver 4210R (Cascade Lake), 10 cores/socket ⇒ 20 CPU cores

- 7× Intel Corporation 900P NVMe <u>Optane SSD</u> ⇒ 4.2 Million IOPS

# Number of System Calls



**Doing I/O with zero system calls!**

# Results: Efficiency (<u>single</u> CPU core)



*io_uring sits between libaio and SPDK*

*Performance collapses with the kernel polling*

Systor'22

# Analysis

*[Interesting]* 8 milliseconds constant latency for all queue depths!

Poor scheduling, and CPU sharing - **Careful!**

SPDK is still 5x more efficient

# Result: Efficiency with <u>TWO</u> CPU cores

[ aio  <  iou  <  iou with polling  <  iou with kernel poll  <  SPDK ]

Normal service order can be resumed (**but** at the cost of 2x CPU cores)!

# Results: Scalability

CPU Cores used

**io_uring kernel polling:** Performance collapses when the number of poller CPU threads increases beyond the cores

**CPU efficiency is still bad:** 10x more CPU cores needed

# io_uring : Programming Ecosystem

- liburing : https://github.com/axboe/liburing
  - 3x syscall based programming can be tricky, hence, a high(er)-level library

**List of manual pages**

| | | | | | | |
|---|---|---|---|---|---|---|
| [en] IO_URING_CHECK_VERSION(3) | [en] io_uring_get_events(3) | [en] io_uring_prep_linkat(3) | [en] io_uring_prep_readv(3) | [en] io_uring_prep_symlink(3) | [en] io_uring_register(2) | [en] io_uring_sq_space_left(3) |
| [en] IO_URING_VERSION_MAJOR(3) | [en] io_uring_get_probe(3) | [en] io_uring_prep_madvise(3) | [en] io_uring_prep_readv2(3) | [en] io_uring_prep_symlinkat(3) | [en] io_uring_register_buf_ring(3) | [en] io_uring_sqe_set_data(3) |
| [en] IO_URING_VERSION_MINOR(3) | [en] io_uring_get_sqe(3) | [en] io_uring_prep_mkdir(3) | [en] io_uring_prep_recv(3) | [en] io_uring_prep_sync_file_range(3) | [en] io_uring_register_buffers(3) | [en] io_uring_sqe_set_data64(3) |
| [en] __io_uring_buf_ring_cq_advance(3) | [en] io_uring_major_version(3) | [en] io_uring_prep_mkdirat(3) | [en] io_uring_prep_recv_multishot(3) | [en] io_uring_prep_tee(3) | [en] io_uring_register_buffers_sparse(3) | [en] io_uring_sqe_set_flags(3) |
| [en] io_uring(7) | [en] io_uring_minor_version(3) | [en] io_uring_prep_msg_ring(3) | [en] io_uring_prep_recvmsg(3) | [en] io_uring_prep_timeout(3) | [en] io_uring_register_buffers_tags(3) | [en] io_uring_sqring_wait(3) |
| [en] io_uring_buf_ring_add(3) | [en] io_uring_opcode_supported(3) | [en] io_uring_prep_msg_ring_cqe_flags(3) | [en] io_uring_prep_recvmsg_multishot(3) | [en] io_uring_prep_timeout_remove(3) | [en] io_uring_register_buffers_update_tag(3) | [en] io_uring_submit(3) |
| [en] io_uring_buf_ring_advance(3) | [en] io_uring_peek_cqe(3) | [en] io_uring_prep_msg_ring_fd(3) | [en] io_uring_prep_remove_buffers(3) | [en] io_uring_prep_timeout_update(3) | [en] io_uring_register_eventfd(3) | [en] io_uring_submit_and_get_events(3) |
| [en] io_uring_buf_ring_cq_advance(3) | [en] io_uring_prep_accept(3) | [en] io_uring_prep_msg_ring_fd_alloc(3) | [en] io_uring_prep_rename(3) | [en] io_uring_prep_unlink(3) | [en] io_uring_register_eventfd_async(3) | [en] io_uring_submit_and_wait(3) |
| [en] io_uring_buf_ring_init(3) | [en] io_uring_prep_accept_direct(3) | [en] io_uring_prep_multishot_accept(3) | [en] io_uring_prep_renameat(3) | [en] io_uring_prep_unlinkat(3) | [en] io_uring_register_file_alloc_range(3) | [en] io_uring_submit_and_wait_timeout( |
| [en] io_uring_buf_ring_mask(3) | [en] io_uring_prep_cancel(3) | [en] io_uring_prep_multishot_accept_direct(3) | [en] io_uring_prep_send(3) | [en] io_uring_prep_write(3) | [en] io_uring_register_files(3) | [en] io_uring_unregister_buf_ring(3) |
| [en] io_uring_check_version(3) | [en] io_uring_prep_cancel64(3) | [en] io_uring_prep_nop(3) | [en] io_uring_prep_send_set_addr(3) | [en] io_uring_prep_write_fixed(3) | [en] io_uring_register_files_sparse(3) | [en] io_uring_unregister_buffers(3) |
| [en] io_uring_close_ring_fd(3) | [en] io_uring_prep_close(3) | [en] io_uring_prep_openat(3) | [en] io_uring_prep_send_zc(3) | [en] io_uring_prep_writev(3) | [en] io_uring_register_files_tags(3) | [en] io_uring_unregister_eventfd(3) |
| [en] io_uring_cq_advance(3) | [en] io_uring_prep_close_direct(3) | [en] io_uring_prep_openat2(3) | [en] io_uring_prep_send_zc_fixed(3) | [en] io_uring_prep_writev2(3) | [en] io_uring_register_files_update(3) | [en] io_uring_unregister_iowq_aff(3) |
| [en] io_uring_cq_has_overflow(3) | [en] io_uring_prep_connect(3) | [en] io_uring_prep_openat_direct(3) | [en] io_uring_prep_sendmsg(3) | [en] io_uring_queue_exit(3) | [en] io_uring_register_files_update_tag(3) | [en] io_uring_unregister_ring_fd(3) |
| [en] io_uring_cq_ready(3) | [en] io_uring_prep_fadvise(3) | [en] io_uring_prep_openat2_direct(3) | [en] io_uring_prep_sendmsg_zc(3) | [en] io_uring_queue_init(3) | [en] io_uring_register_iowq_aff(3) | [en] io_uring_wait_cqe(3) |
| [en] io_uring_cqe_get_data(3) | [en] io_uring_prep_fallocate(3) | [en] io_uring_prep_poll_add(3) | [en] io_uring_prep_sendto(3) | [en] io_uring_queue_init_params(3) | [en] io_uring_register_iowq_max_workers(3) | [en] io_uring_wait_cqe_nr(3) |
| [en] io_uring_cqe_get_data64(3) | [en] io_uring_prep_fgetxattr(3) | [en] io_uring_prep_poll_multishot(3) | [en] io_uring_prep_shutdown(3) | [en] io_uring_recvmsg_cmsg_firsthdr(3) | [en] io_uring_register_ring_fd(3) | [en] io_uring_wait_cqe_timeout(3) |
| [en] io_uring_cqe_seen(3) | [en] io_uring_prep_files_update(3) | [en] io_uring_prep_poll_remove(3) | [en] io_uring_prep_socket(3) | [en] io_uring_recvmsg_cmsg_nexthdr(3) | [en] io_uring_register_sync_cancel(3) | [en] io_uring_wait_cqes(3) |
| [en] io_uring_enter(2) | [en] io_uring_prep_fsetxattr(3) | [en] io_uring_prep_poll_update(3) | [en] io_uring_prep_socket_direct(3) | [en] io_uring_recvmsg_name(3) | [en] io_uring_setup(2) | |
| [en] io_uring_enter2(2) | [en] io_uring_prep_fsync(3) | [en] io_uring_prep_provide_buffers(3) | [en] io_uring_prep_socket_direct_alloc(3) | [en] io_uring_recvmsg_out(3) | [en] io_uring_setup_buf_ring(3) | |
| [en] io_uring_for_each_cqe(3) | [en] io_uring_prep_getxattr(3) | [en] io_uring_prep_read(3) | [en] io_uring_prep_splice(3) | [en] io_uring_recvmsg_payload(3) | [en] io_uring_sq_ready(3) | |
| [en] io_uring_free_buf_ring(3) | [en] io_uring_prep_link(3) | [en] io_uring_prep_read_fixed(3) | [en] io_uring_prep_statx(3) | [en] io_uring_recvmsg_payload_length(3) | | |
| [en] io_uring_free_probe(3) | [en] io_uring_prep_link_timeout(3) | | | [en] io_uring_recvmsg_validate(3) | | |

- Active research in leveraging io_uring in DBs, key-value store, etc.
- Applicability beyond storage as the "core" kernel-application interfacing API

# What you should know from this lecture

What is CXL and what key problems does it solve

What is different types of CXL protocols, device types, and generational features

What does flash + CXL allow us to do

What is asynchronous and non-block I/O, and what different APIs support them

What is io_uring? What are the different operation completion modes it support

What are the performance implications of these modes

**The New(er) Triangle of Storage-Memory Continuum**

# To Conclude

**Storage Research is fundamentally changing and reshaping what kind of systems we can build tomorrow**

- Performance
- Abstractions
- Efficiency
- Programmability
- Cost
- Scalability

This course came out of this report ;)

## Data Storage Research Vision 2025

Report on NSF Visioning Workshop held May 30–June 1, 2018

George Amvrosiadis[†], Ali R. Butt[¶], Vasily Tarasov[‡], Erez Zadok[*], Ming Zhao[§]

Irfan Ahmad, Remzi H. Arpaci-Dusseau, Feng Chen, Yiran Chen, Yong Chen, Yue Cheng, Vijay Chidambaram, Dilma Da Silva, Angela Demke-Brown, Peter Desnoyers, Jason Flinn, Xubin He, Song Jiang, Geoff Kuenning, Min Li, Carlos Maltzahn, Ethan L. Miller, Kathryn Mohror, Raju Rangaswami, Narasimha Reddy, David Rosenthal, Ali Saman Tosun, Nisha Talagala, Peter Varman, Sudharshan Vazhkudai, Avani Waldani, Xiaodong Zhang, Yiying Zhang, and Mai Zheng.

[†]Carnegie Mellon University, [¶]Virginia Tech, [‡]IBM Research, [*]Stony Brook University, [§]Arizona State University

February 2019

**Executive Summary**

With the emergence of new computing paradigms (e.g., cloud and edge computing, big data, Internet of Things (IoT), deep learning, etc.) and new storage hardware (e.g., non-volatile memory (NVM), shingled-magnetic recording (SMR) disks, and kinetic drives, etc.), a number of open challenges and research issues need to be addressed to ensure sustained storage systems efficacy and performance. The wide variety of applications demand that the fundamental design of storage systems should be revisited to support application-specific and application-defined semantics. Existing standards and abstractions need to be reevaluated; new sustainable data representations need to be designed to support emerging applications. To take advantage of hardware advancements, new storage software designs are also necessary in order to maximize overall system efficiency and performance.

Therefore, there is a urgent need for a consolidated effort to identify and establish a vision for storage systems research and comprehensive techniques that provide practical solutions to the storage issues facing the information technology community. To address this need, the National Science Foundation's (NSF) "Visioning Workshop on Data Storage Research 2025" brought together a number of storage researchers from academia, industry, national laboratories, and federal agencies to develop a collective vision for future storage research, as well as to prioritize

# The New(er) Triangle of Storage-Memory Continuum

**Cost: $ / GB**

CPU register — ~0.1ns

CPU cache — 1-10ns

DRAM Memory — ~10-100ns

**Persistent Memory** — **< 1 usec**

**NAND Flash/Optane SSDs** — **~10-100 usec**

Hard disk drive (HDD) — ~10-100**ms**

Tape — ~100**ms**-10**s**

**Local CXL Mem 150-300 nsecs**

**Local CXL SSDs 1 -10 usec**

**Remote CXL Mem ~1-2 usecs**

**Remote CXL SSD 10-100 usec**

**Instead of discrete steps, it is a continuous spectrum now: Continuum**

# Further Reading - CXL (1 or 2)

- CXL Consortium, https://www.computeexpresslink.org/
- CXL resources, https://www.computeexpresslink.org/resource-library
- Linux CXL driver code: https://elixir.bootlin.com/linux/latest/source/drivers/cxl
- Debendra Das Sharma, and others, An Introduction to the Compute Express Link (CXL) Interconnect, **2023**, https://arxiv.org/abs/2306.11227
- Hasan Al Maruf, and others. TPP: Transparent Page Placement for CXL-Enabled Tiered-Memory. In Proceedings of the 28th ACM ASPLOS **2023**. https://doi.org/10.1145/3582016.3582063
- Myoungsoo Jung. **2022**. Hello bytes, bye blocks: PCIe storage meets compute express link for memory expansion (CXL-SSD). In Proceedings of the 14th ACM HotStorage '22, https://doi.org/10.1145/3538643.3539745
- Miryeong Kwon, Sangwon Lee, and Myoungsoo Jung. 2023. Cache in Hand: Expander-Driven CXL Prefetcher for Next Generation CXL-SSD. In Proceedings of the 15th ACM HotStorage '23, https://doi.org/10.1145/3599691.3603406
- Huaicheng Li, and others. Pond: CXL-Based Memory Pooling Systems for Cloud Platforms. In Proceedings of the 28th ACM ASPLOS 2023, https://doi.org/10.1145/3575693.3578835
- Shao-Peng Yang and others. Overcoming the Memory Wall with CXL-Enabled SSDs, USENIX ATC **2023**, https://www.usenix.org/conference/atc23/presentation/yang-shao-peng
- Donghyun Gouk and others, Direct Access, High-Performance Memory Disaggregation with DirectCXL, USENIX ATC **2022**, https://www.usenix.org/conference/atc22/presentation/gouk

# Further Reading - CXL (2 of 2)

- CXL-ANNS: Software-Hardware Collaborative Memory Disaggregation and Computation for Billion-Scale Approximate Nearest Neighbor Search, USENIX ATC 2023, https://www.usenix.org/conference/atc23/presentation/jang
- Marcos K. Aguilera, and others. 2023. Memory disaggregation: why now and what are the challenges. SIGOPS Oper. Syst. Rev. 57, 1 (June **2023**), 38–46. https://doi.org/10.1145/3606557.3606563
- Hasan Al Maruf and Mosharaf Chowdhury. 2023. Memory Disaggregation: Advances and Open Challenges. SIGOPS Oper. Syst. Rev. 57, 1 (June **2023**), 29–37. https://doi.org/10.1145/3606557.3606562
- Jianguo Wang and Qizhen Zhang. **2023**. Disaggregated Database Systems. In Companion of the **2023** International Conference on Management of Data (SIGMOD '23). https://doi.org/10.1145/3555041.3589403
- Wenjing Jin, and others. DRAM Translation Layer: Software-Transparent DRAM Power Savings for Disaggregated Memory. In Proceedings of the 50th Annual International Symposium on Computer Architecture (**ISCA '23**). https://doi.org/10.1145/3579371.3589051
- What's the Difference Between CXL 1.1 and CXL 2.0? https://www.electronicdesign.com/technologies/embedded/article/21249351/cxl-consortium-whats-the-difference-between-cxl-11-and-cxl-20
- QEMU CXL setup, https://www.qemu.org/docs/master/system/devices/cxl.html
- How To Map a CXL Endpoint to a CPU Socket in Linux, https://stevescargall.com/blog/2022/12/27/how-to-map-a-cxl-endpoint-to-a-cpu-socket-in-linux/

# Further Reading - io_uring (1 of 2)

- Efficient IO with io_uring, https://kernel.dk/io_uring.pdf
- What's new with io_uring, https://kernel.dk/axboe-kr2022.pdf
- An Introduction to the io_uring Asynchronous I/O Framework, https://blogs.oracle.com/linux/post/an-introduction-to-the-io-uring-asynchronous-io-framework
- Zebin Ren and Animesh Trivedi. 2023. Performance Characterization of Modern Storage Stacks: POSIX I/O, libaio, SPDK, and io_uring. In Proceedings of the 3rd Workshop on Challenges and Opportunities of Efficient and Performant Storage Systems (CHEOPS '23). Association for Computing Machinery, New York, NY, USA, 35–45. https://doi.org/10.1145/3578353.3589545
- Diego Didona, Jonas Pfefferle, Nikolas Ioannou, Bernard Metzler, and Animesh Trivedi. 2022. Understanding modern storage APIs: a systematic study of libaio, SPDK, and io_uring. In Proceedings of the 15th ACM International Conference on Systems and Storage (SYSTOR '22). Association for Computing Machinery, New York, NY, USA, 120–127. https://doi.org/10.1145/3534056.3534945
- Simon A. F. Lund, Philippe Bonnet, Klaus B. A. Jensen, and Javier Gonzalez. 2022. I/O interface independence with xNVMe. In Proceedings of the 15th ACM International Conference on Systems and Storage (SYSTOR '22). Association for Computing Machinery, New York, NY, USA, 108–119. https://doi.org/10.1145/3534056.3534936
- Sidharth Sundar, William Simpson, Jacob Higdon, Caeden Whitaker, Bryan Harris, and Nihat Altiparmak. 2023. Energy Implications of IO Interface Design Choices. In Proceedings of the 15th ACM Workshop on Hot Topics in Storage and File Systems (HotStorage '23). Association for Computing Machinery, New York, NY, USA, 58–64. https://doi.org/10.1145/3599691.3603411

# Further Reading - io_uring (2 of 2)

- Ringing in a new asynchronous I/O API, https://lwn.net/Articles/776703/
- [PATCHSET v5] io_uring IO interface, https://lore.kernel.org/linux-block/20190116175003.17880-1-axboe@kernel.dk/
- Gabriel Haas and Viktor Leis. 2023. What Modern NVMe Storage Can Do, and How to Exploit it: High-Performance I/O for High-Performance Storage Engines. Proc. VLDB Endow. 16, 9 (May 2023), 2090–2102. https://doi.org/10.14778/3598581.3598584
- Hugh C. Lauer and Roger M. Needham. 1979. On the duality of operating system structures. SIGOPS Oper. Syst. Rev. 13, 2 (April 1979), 3–19. https://doi.org/10.1145/850657.850658
- John Ousterhout, Why Threads Are A Bad Idea (for most purposes), https://web.stanford.edu/~ouster/cgi-bin/papers/threads.pdf
- Rob von Behren, Jeremy Condit, and Eric Brewer. 2003. Why events are a bad idea (for high-concurrency servers). In Proceedings of the 9th conference on Hot Topics in Operating Systems - Volume 9 (HOTOS'03). USENIX Association, USA, 4. https://dl.acm.org/doi/10.5555/1251054.1251058
- Philipp Haller, Martin Odersky, Scala Actors: Unifying thread-based and event-based programming, 2008, https://doi.org/10.1016/j.tcs.2008.09.019.
- A 5 part series on the asynchronous nature of I/O, OS, and concurrency: https://blog.acolyer.org/2014/12/08/on-the-duality-of-operating-system-structures/
- μTune: Auto-Tuned Threading for OLDI Microservices, https://www.usenix.org/conference/osdi18/presentation/sriraman
- Linux Asynchronous I/O, https://oxnz.github.io/2016/10/13/linux-aio/
- Linux-aio, https://github.com/littledan/linux-aio