# Operating Systems and C
# Fall 2022
# **11. File**

# Memory Abstraction

recall fundamental abstractions
- interpreter     during lecture 4
- **memory**     **today (again)!**
- communication    at end of course

still a memory abstraction

logical name

physical address

WRITE(name, value)

WRITE(address, value)

**Associativity Layer**

**Location-addressed Memory**

READ(name)

READ(address)

Associative Memory

mapping

memory
API

Source: Saltzer and Kaashoek

# Computer Hardware

recall hardware:
- ~~CPU~~
- **memory (~~main memory~~, disk)**
- **communication**

CPU store/load only via. Main memory.
How CPU accesses data on disk?

CPU

Registers

PC

ALU

Bus interface

System bus

Memory bus

I/O bridge

Main memory

I/O bus

USB controller

Graphics adapter

Disk controller

Mouse  Keyboard

Display

Disk

*hello executable stored on disk*

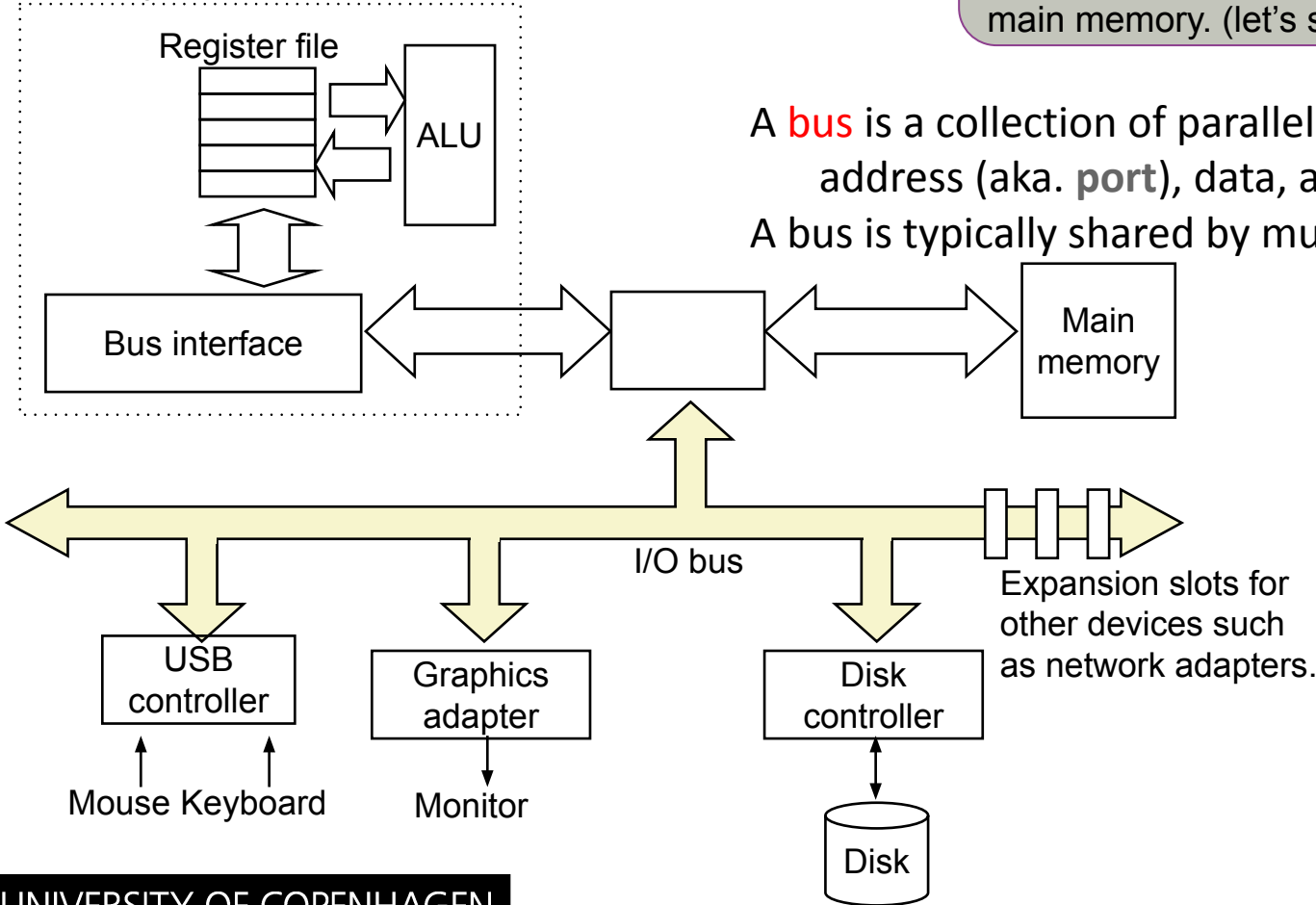Expansion slots for other devices such as network adapters

focus today

# I/O Bus

**answer:**
CPU sends instruction to disk controller.
(which implements **memory abstraction**).
disk controller transfers to (read) / from (write)
main memory. (let's see this in action)
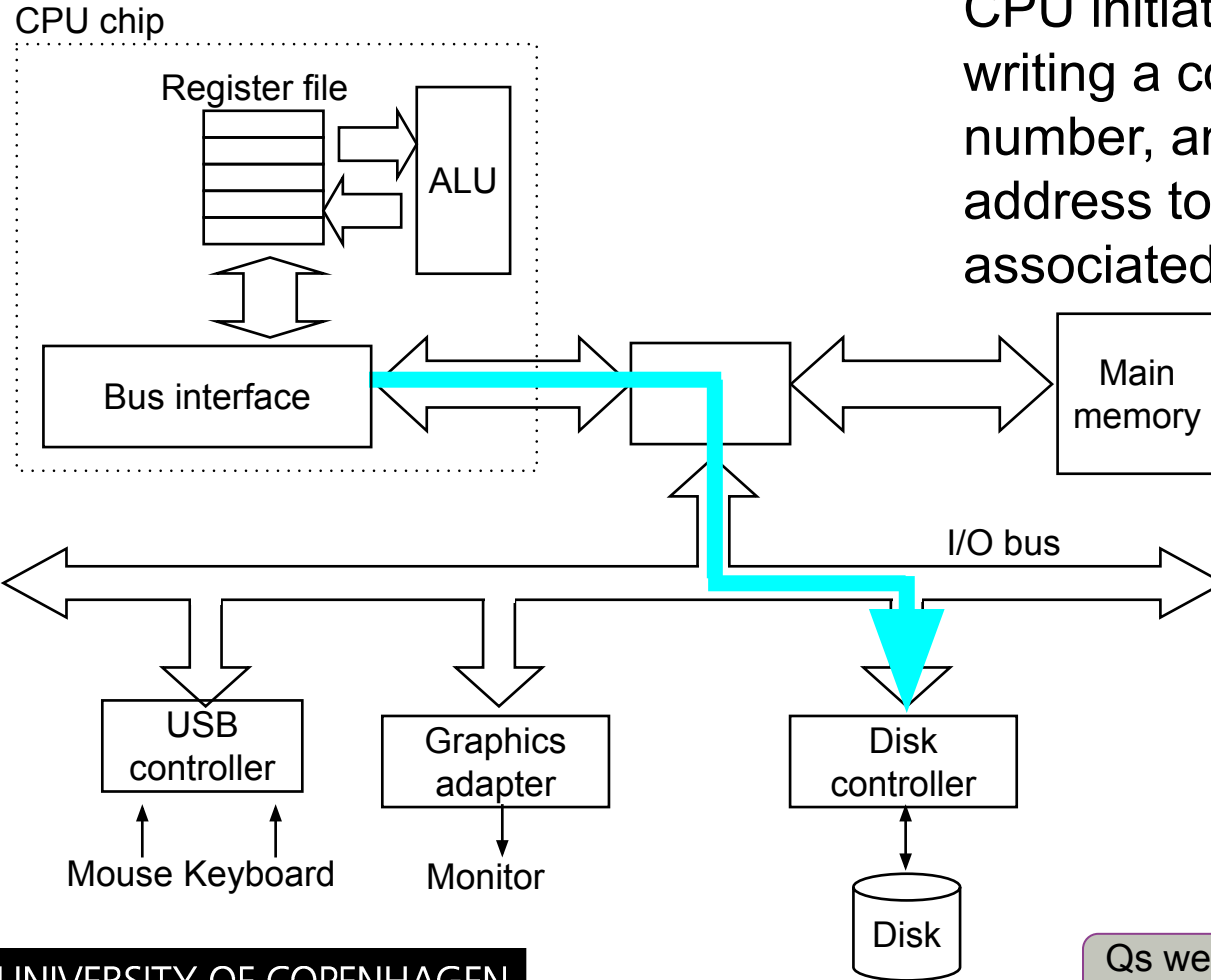
A bus is a collection of parallel wires that carry address (aka. **port**), data, and control signals.
A bus is typically shared by multiple devices.

CPU chip

Register file

ALU

Bus interface

Main memory

I/O bus

Expansion slots for other devices such as network adapters.

USB controller

Graphics adapter

Disk controller

Mouse Keyboard

Monitor

Disk

IT UNIVERSITY OF COPENHAGEN

CPU chip

Register file

ALU

Bus interface

Main memory

I/O bus

USB controller

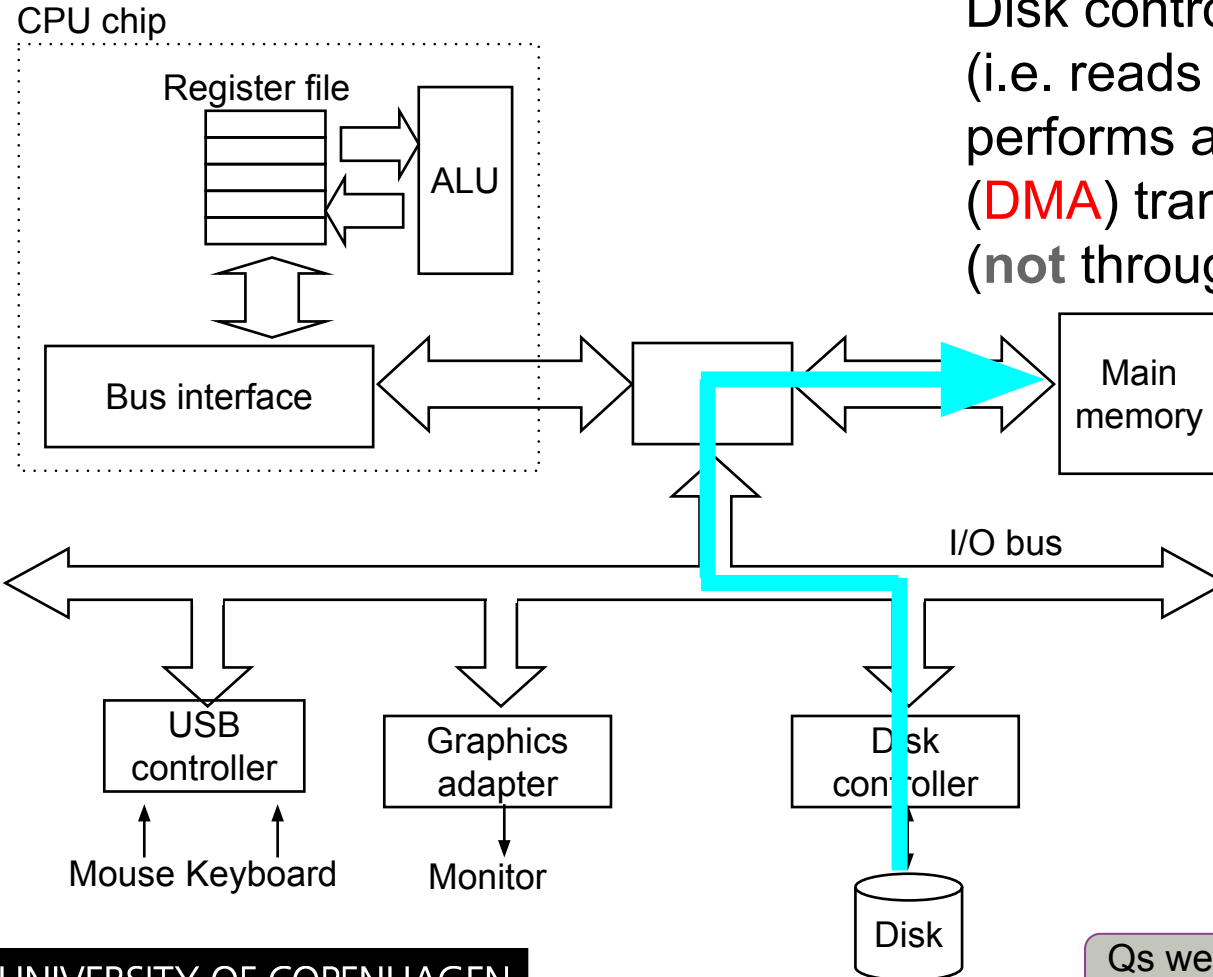Graphics adapter

Disk controller

Mouse Keyboard

Monitor

Disk

CPU initiates a disk read by writing a command, logical block number, and destination memory address to a port (address) associated with disk controller.

What does a disk controller do?
How do hosts interact with disk controllers?
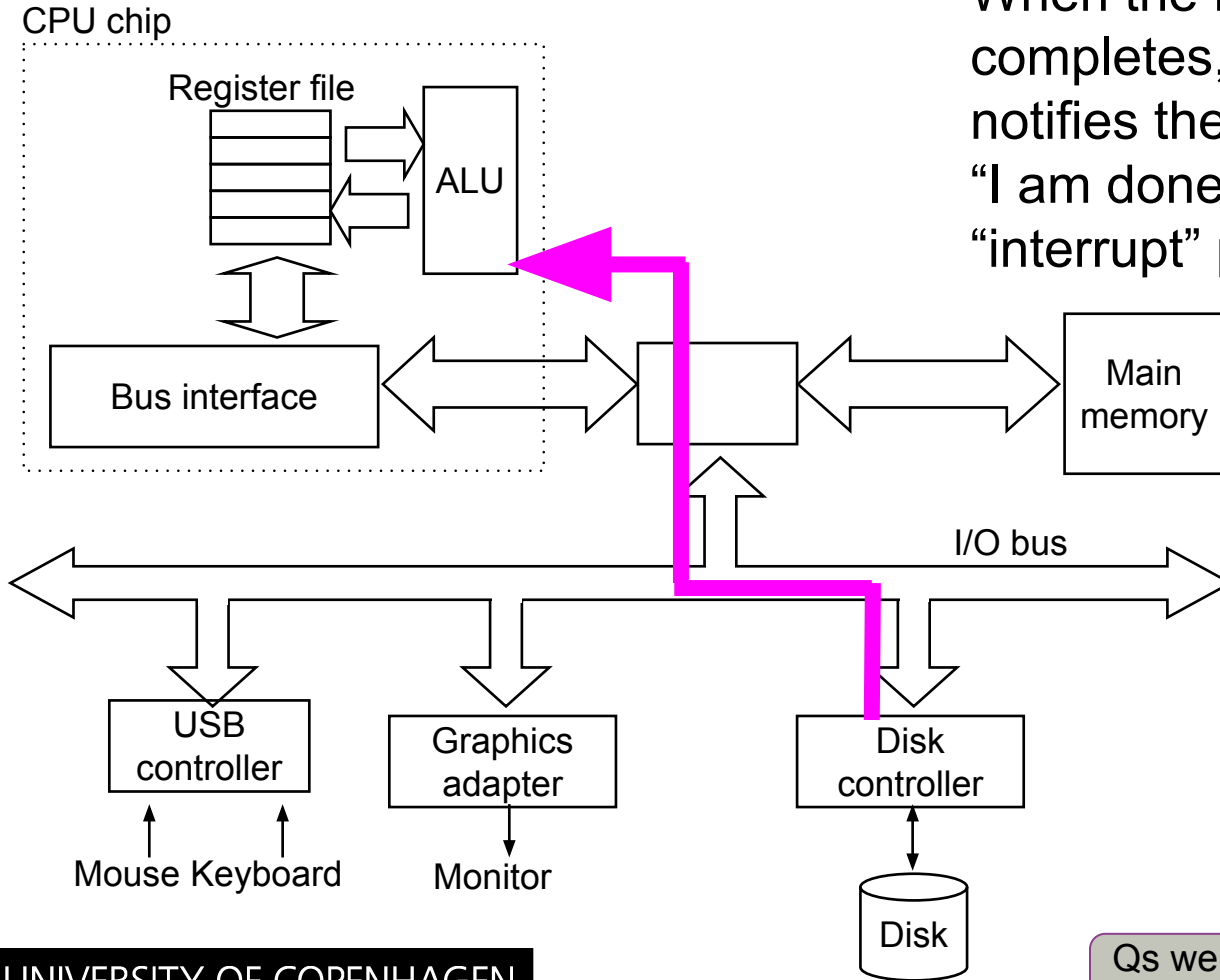
Qs we will address

Disk controller does its work (i.e. reads the sector) and performs a direct memory access (DMA) transfer into main memory. (**not** through the CPU)

How does a disk controller perform a DMA?

CPU chip

Register file

ALU

Bus interface

Main memory

I/O bus

USB controller

Graphics adapter

Disk controller

Mouse Keyboard

Monitor

Disk

Qs we will address

# Reading a Disk Sector (3)

When the DMA transfer completes, the disk controller notifies the CPU with an *interrupt* "I am done" (i.e., asserts a special "interrupt" pin on the CPU)

CPU chip

Register file

ALU

Bus interface

Main memory

I/O bus

USB controller

Mouse Keyboard

Graphics adapter

Monitor

Disk controller

Disk

How are I/Os handled on the host?
How are I/Os exposed to programmers?

Qs we will address

IT UNIVERSITY OF COPENHAGEN

1. **File System**

   file is the main abstraction for data that is stored

   - How are I/Os exposed to programmers?

   - How are I/Os handled on the host?

2. Storage devices

   because they are cool

   - What does a disk controller do?

   - How do hosts interact with disk controllers?

   - How does a disk controller perform a DMA?
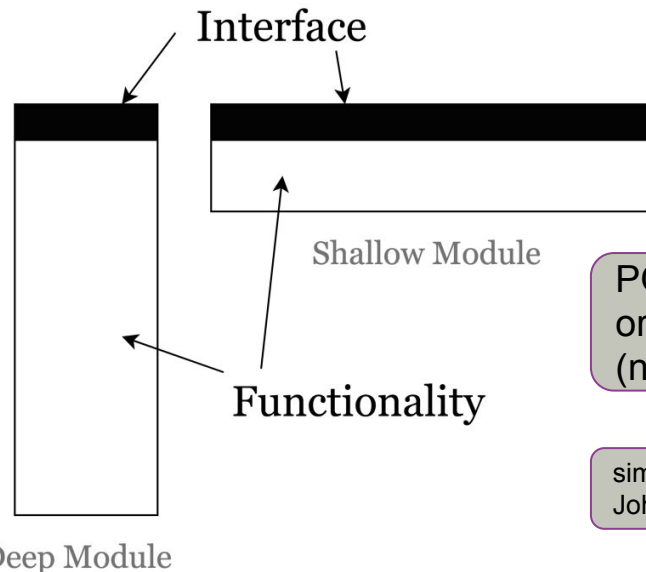
3. Computational Storage

like memory!

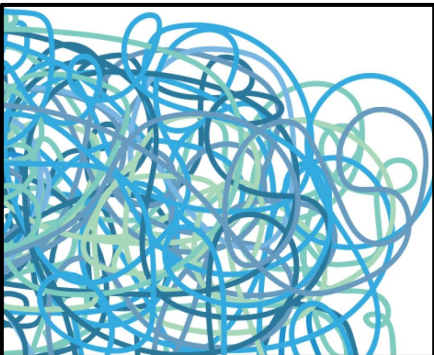# A file is an array of bytes

File interface: create/delete, open/close, read/write



Interface

Shallow Module

Functionality

Deep Module

POSIX: unified unix specification.
original POSIX interface is beautiful.
(not part of C language, but std lib)

simple interface, hiding a lot of complexity.
John Ousterhout, professor at Stanford.

A PHILOSOPHY OF SOFTWARE DESIGN JOHN OUSTERHOUT

HAGEN

https://www.youtube.com/watch?v=bmSAYlu0NcY&feature=youtu.be

# Outline

## File System Layers

What happens when you open a file or when you read from a file?

## Linux File System Components

How is the Linux file system organized?

# Layering and Naming

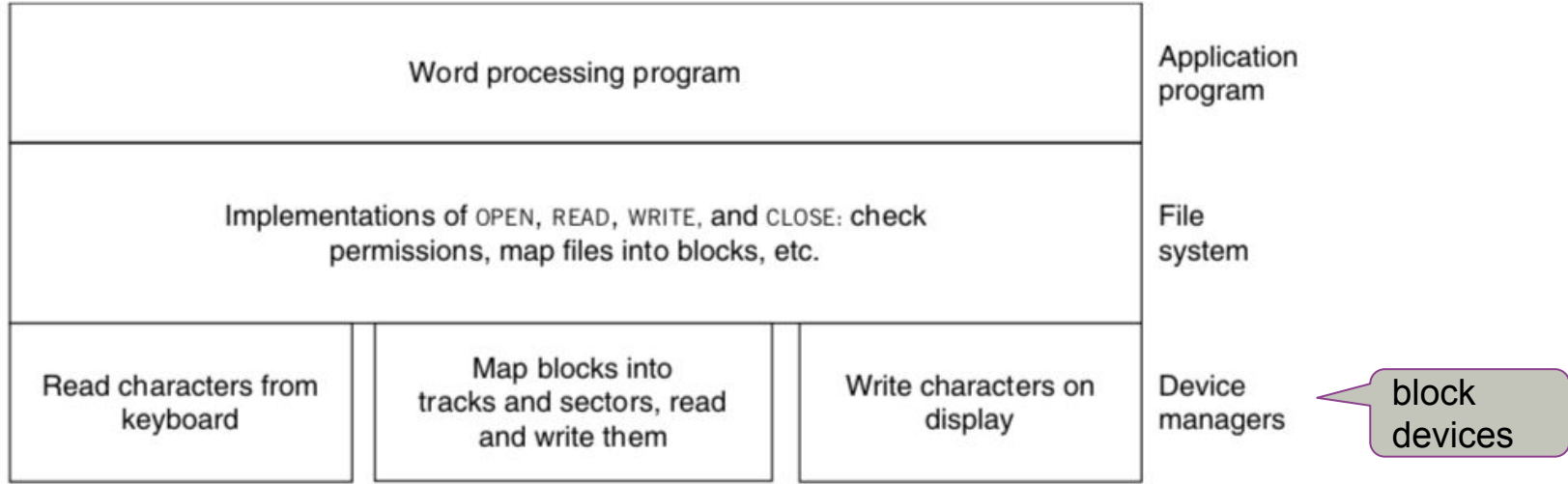| Word processing program | | | Application program |
| Implementations of OPEN, READ, WRITE, and CLOSE: check permissions, map files into blocks, etc. | | | File system |
| Read characters from keyboard | Map blocks into tracks and sectors, read and write them | Write characters on display | Device managers |

block devices

Figure from Principles of Computer System Design, Saltzer /Kaashoek

# Block Layer

A **block device** is an **array of blocks**.
To each block is associated a number,
a Logical Block Address (LBA)

```
procedure BLOCK_NUMBER_TO_BLOCK (integer b) returns block
    return device[b]
```

sound familiar?
virtual memory! pages!
(quantized data)
block is a unit of transfer.
(associativity layer maps block
number to actual block)

hard disk drives today
(incl. SSDs) are block devices.

IT UNIVERSITY OF COPENHAGEN

# File Layer

How to represent files?
Each **file** is a **collection of disk blocks**
(more abstractly (haha), an **array of bytes**)

```
structure inode {
    integer block_numbers[N];    // the numbers of the blocks that constitute the file
    integer size;                // the size of the file in bytes
}
```

*inode* ("index node") is a collection of block numbers (associated to the file), and their collective size.
(we need this level of indirection)

```
procedure INDEX_TO_BLOCK_NUMBER (instance of inode i, integer index) returns integer{
    return i.block_numbers[index];
}
```

```
procedure INODE_TO_BLOCK ( integer offset, instance of inode i) returns instance of block
    o ← offset / BLOCKSIZE;
    b ← INDEX_TO_BLOCK_NUMBER(inode, o);
    return BLOCK_NUMBER_TO_BLOCK(b);
}
```

# Inode Name Layer

How to avoid carrying inodes around?

level of indirection

number the inodes!
inode_number to inode table (map),
carry this table around.

```
procedure INODE_NUMBER_TO_INODE(integer inode_number) returns instance of inode{
      return inode_table[inode_number];
}


procedure INODE_NUMBER_TO_BLOCK (integer offset, integer inode_number)
                                              returns instance of block {
      structure inode i ← INODE_NUMBER_TO_INODE (inode_number);
      o ← offset / BLOCKSIZE;
      b ← INDEX_TO_BLOCK_NUMBER (i, o);
      return BLOCK_NUMBER_TO_BLOCK (b);
}
```

`return INODE_TO_BLOCK(offset,i);`

# File Name Layer

File system state:
inode_table

**Representing directories**

directory is also an inode.
now we have 2 types of inodes.

```
structure inode{
    integer block_numbers[N];   // the numbers of the blocks that constitute the file
    integer size;               // the size of the file in bytes
    integer type;               // type of file: regular file, directory,...
}
```

each block stores
many inode nums

**User-friendly names**

| File name | Inode number |
|-----------|--------------|
| program   | 10           |
| Paper     | 12           |

when you work with files, you
don't work with inode numbers.
you work with filenames.
need mapping from filename to
inode number.

in dir, we store, **alongside an
inode number**, the *filename* of
that inode.

# File Name Layer

Directory lookup

```
procedure NAME_TO_INODE_NUMBER (character string filename, integer dir) returns integer {
        return LOOKUP (filename, dir);
}


procedure LOOKUP (character string filename, integer dir) returns integer {
        instance of block b;
        instance of inode i ← INODE_NUMBER_TO_INODE (dir);
        if i.type ≠ DIRECTORY then return FAILURE;
        for offset from 0 to i.size − 1 do {
                b ← INODE_NUMBER_TO_BLOCK (offset, dir);
                if STRING_MATCH (filename, b) then {
                        return INODE_NUMBER (filename, b);// return inode number for filename
                }
                offset ← offset + BLOCKSIZE;             // increase offset by block size
        }
        return FAILURE;
}
```

(inode number)

if filename occurs in b,

then return the inode num that's written next to the filename

STRING_MATCH, INODE_NUMBER implementation not shown

# Path Name Layer

Hierarchy of Directories

```
procedure PATH_TO_INODE_NUMBER (character string path, integer dir) returns integer{
    if (PLAIN_NAME (path)) return NAME_TO_INODE_NUMBER (path, dir);
    else {
            dir ← LOOKUP (FIRST (path), dir);
            path ← REST (path);
            return PATH_TO_INODE_NUMBER (path, dir);
    }
}
```

# Absolute Path Name Layer

File system state:
inode_table
Process state:
wd

Change working directory

**procedure** CHDIR (*path* **character string**) { *wd* ← PATH_TO_INODE_NUMBER (*path*, *wd*); }

How to name a file regardless of the current working directory?

**procedure** GENERALPATH_TO_INODE_NUMBER (**character string** *path*) **returns integer**{
    **if** (*path*[0] = "/") **return** PATH_TO_INODE_NUMBER(*path*, 1);
    **else return** PATH_TO_INODE_NUMBER(*path*, *wd*);
}

(root inode number)

# Unix File System Naming Scheme
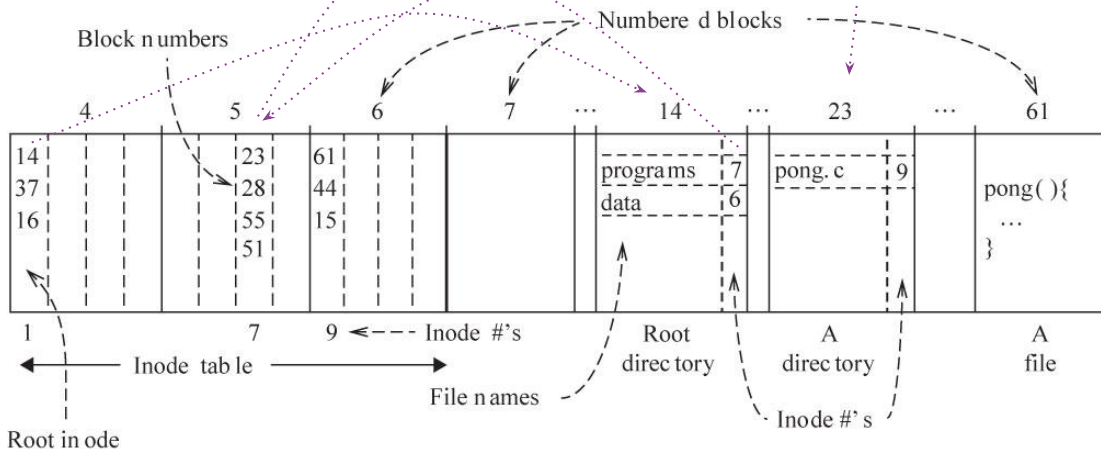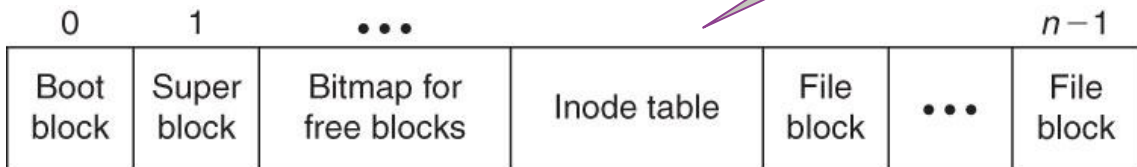
File system state:
    inode_table
Process state:
    wd

**Disk** Layout for a file system

must be persistent, else won't know which files are there

# Symbolic Link Layer

How about flexible management of files?

```
LINK (from_name, to_name);
UNLINK (from_name);
```

```
structure inode{
    integer block_numbers[N];
    integer size;
    integer type;
    integer refcnt;
}
```

once no path refers to inode, it can be garbage collected.

# Symbolic Link Layer

How to create links across file systems
(where the inode numbers are not unique)?

```
procedure PATHNAME_TO_INODE (character string filename) returns instance of inode{
        instance of inode i;
        inode_number ← GENERALPATH_TO_INODE_NUMBER (filename);
        i ← INODE_NUMBER_TO_INODE (inode_number);
        if i.type = SYMBOLIC then
                i = GENERALPATH_TO_INODE_NUMBER (COERCE_TO_STRING (i.block_numbers))
        return i;
}
```

# Symbolic Link Layer

10s

How to attach new disks to a file system?

MOUNT ("/dev/fd1", "/floppy")

disk     mount point

(1) represents a file system
(2) device and root inode for the given file system

Inode pinned in memory for <u>usb</u>

(1) name of parent inode, i.e., usb

Inode pinned in memory for <u>/dev/usb1</u>

10s

| Layer | Names | Values | Context | Name-mapping algorithm | |
|---|---|---|---|---|---|
| Symbolic link | Path names | Path names | The directory hierarchy | PATHNAME_TO_GENERAL_PATH | user-oriented names |
| Absolute path name | Absolute path names | Inode numbers | The root directory | GENERALPATH_TO_INODE_NUMBER | |
| Path name | Relative path names | Inode numbers | The working directory | PATH_TO_INODE_NUMBER | |
| File name | File names | Inode numbers | A directory | NAME_TO_INODE_NUMBER | machine-user interface |
| Inode number | Inode numbers | Inodes | The inode table | INODE_NUMBER_TO_INODE | machine-oriented names |
| File | Index numbers | Block numbers | An inode | INDEX_TO_BLOCK_NUMBER | |
| Block | Block numbers | Blocks | The disk drive | BLOCK_NUMBER_TO_BLOCK | |

Source: Saltzer and Kaashoek

File system state:
  inode_table
  file_table
Process state:
  fd_table
  wd

### Which files are in use?
### `file_table`

| File name | Inode number | cursor |
|-----------|--------------|--------|
| program   | 10           | 64     |
| Paper     | 12           | 0      |

Cursor is the first byte that will be accessed by the next read or write operation.

### Which files is each process using?
### `fd_table`

Mapping from file descriptors into the **file_table**.
(file descriptors are per-process. natural numbers; 0 is stdin, 1 is stdout, 2 is stderr, …)

Multiple processes can have a file open with different cursors, and
Multiple processes can have a file open sharing a cursor (fork; **fd_table** shared)

# API: inode

```
structure inode {
        integer block_numbers[N];   // the number of blocks that constitute the file
        integer size;               // the size of the file in bytes
        integer type;               // type of file: regular file, directory, symbolic link
        integer refcnt;             // count of the number of names for this inode
        integer userid;             // the user ID that owns this inode
        integer groupid;            // the group ID that owns this inode
        integer mode;               // inode's permissions
        integer atime;              // time of last access (READ, WRITE,...)
        integer mtime;              // time of last modification
        integer ctime;              // time of last change of inode
```

we did not talk about e.g. access control

# API Calls: Open

Wd
file_table
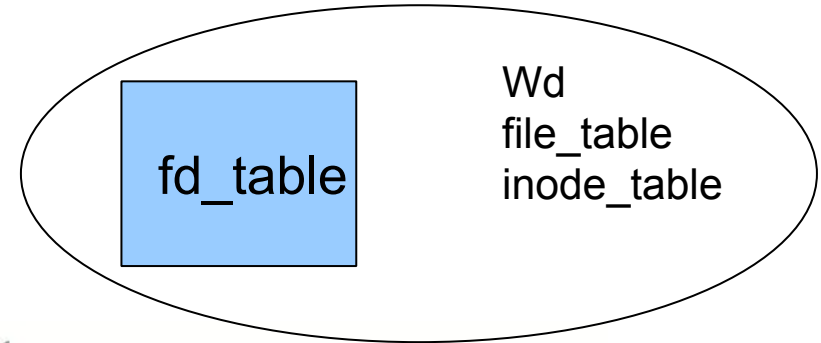inode_table

fd_table

```
procedure OPEN (character string filename, flags, mode) {
      inode_number ← PATH_TO_INODE_NUMBER (filename, wd);
      if inode_number = FAILURE and flags = O_CREATE then { // Create the file?
             inode_number ← CREATE (filename, mode);        // Yes, create it.
      } else return FAILURE;
      inode ← INODE_NUMBER_TO_INODE (inode_number);
      if PERMITTED (inode, flags) then {   // Does this user have the required permissions?
             file_index ← INSERT (file_table, inode_number);
             fd ← FIND_UNUSED_ENTRY (fd_table);// Yes, find entry in file descriptor table
             fd_table[fd] ← file_index;         // Record file index for the file descriptor
             return fd;                   // Return fd
      } else return FAILURE;             // No, return a failure
}
```

# API Calls: Read

## Process state

| File name | Inode number | cursor |
|-----------|--------------|--------|
| program | 10 | 64 |
| Paper | 12 | 0 |

fd_table

Wd
file_table
inode_table

```
procedure READ (fd, reference buf, n) {
        file_index ← fd_table[fd];
        cursor ← file_table[file_index].cursor;
        inode ← INODE_NUMBER_TO_INODE (file_table[file_index].inode_number);
        m = MINIMUM (inode.size – cursor, n);
        atime of inode ← NOW ();
        if m = 0 then return END_OF_FILE;
        for i from 0 to m – 1 do {
                b ← INODE_NUMBER_TO_BLOCK (i, inode_number);
                COPY (b, buf, MINIMUM (m – i, BLOCKSIZE));
                i ← i + MINIMUM (m – i, BLOCKSIZE);
        }
        file_table[file_index].cursor ← cursor + m;
        return m;
}
```

a word about the Linux file system.

## File System Layers

What happens when you open a file or when you read from a file?

## Linux File System Components

How is the Linux file system organized?

# Linux File System

**User Space**

Process   Process

**Kernel Space**

IO System Call

Virtual File System

Block Layer

Device Driver

Submission   Completion

struct bio

struct request

uniform way for different file systems to hook up to Linux. (POSIX)

**Block Device HW**

Block Device Controller

# Linux I/O System Calls

- creat, open, read, write, close, lseek
- fsync
- link, unlink
- stat, lstat, fstat
- access, umask, chmod, chown, utime
- ioctl

libraries for this

there are also *async I/O* system calls. (e.g. aio_read) and ways to batch system calls (io_submit, …)

# Linux Virtual File System

The virtual file system defines the generic file system interface and data structures:
file, dentry, inode, vfsmount, super_block.

Each specific file system provides a specific implementation:
block-based FS (ext4, btrfs), network FS (NFS, ceph),
stackable FS, pseudo FS (sysfs),
special purpose FS (tmpfs)

# Linux VFS

they all respect the VFS setup.



file table,
cache,
inode table, …

don't be afraid of VFS; it's just inodes, inode table, etc.

https://www.starlab.io/blog/introduction-to-the-linux-virtual-filesystem-vfs-part-i-a-high-level-tour

# Linux Legacy Block Layer

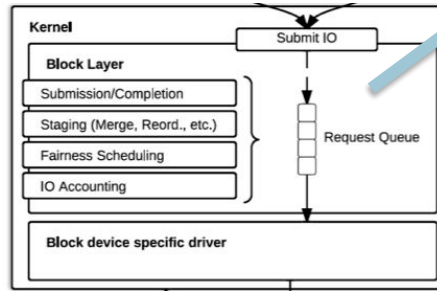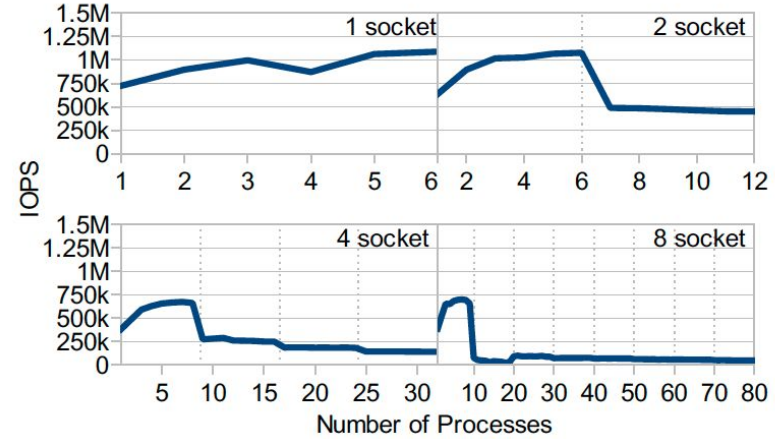block layer is taking block I/O requests, and issuing those to HW. how hard can that be?



**Userspace**

Process    Process

**Kernel**    Submit IO

**Block Layer**

Submission/Completion

Staging (Merge, Reord., etc.)

Fairness Scheduling    Request Queue

IO Accounting

**Block device specific driver**

Status / Completion Interrupt

Single queue (e.g. SATA) capable hardware device

in 2013, they started making block layer smart: disk receives many random I/Os ⇒ disk wants to reorder them to be sequential (lump them together to be smart about how you access disk).

devices are async, yet I/O from FS is sync.
mapping from synchronous to asynchronous was done by the block layer.

https://kernel.dk/blk-mq.pdf
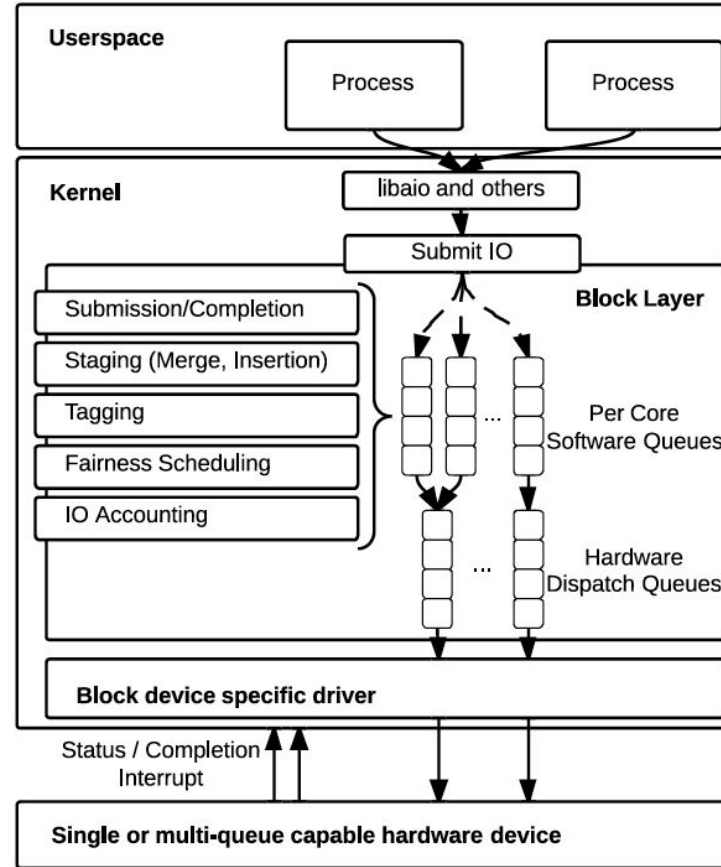
# Scalability Problem

| Platform (Intel) | Sandy Bridge-E | Westmere-EP | Nehalem-EX | Westmere-EX |
|---|---|---|---|---|
| Processor | i7-3930K | X5690 | X7560 | E7-2870 |
| Num. of Cores | 6 | 12 | 32 | 80 |
| Speed (Ghz) | 3.2 | 3.46 | 2.66 | 2.4 |
| L3 Cache (MB) | 12 | 12 | 24 | 30 |
| NUMA nodes | 1 | 2 | 4 | 8 |

https://kernel.dk/blk-mq.pdf

# Linux mlqblk Block Layer

multi-queue block layer.
i.e. per-core queue in software.
sound good?

block layer maintainer

Jens Axboe's design

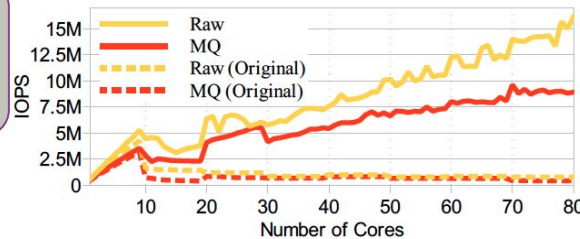https://kernel.dk/blk-mq.pdf

hmm, that yielded no speedup.



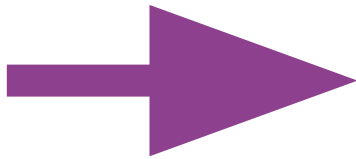last 10 years of work in Linux IO stack: remove latency to increase throughput to IO devices

IT UNIVERSITY OF COPENHAGEN

**User Space**

Process

Process

**Kernel Space**

IO System Call

Virtual File System

struct bio

Block Layer

struct request

Device Driver

Submission    Completion

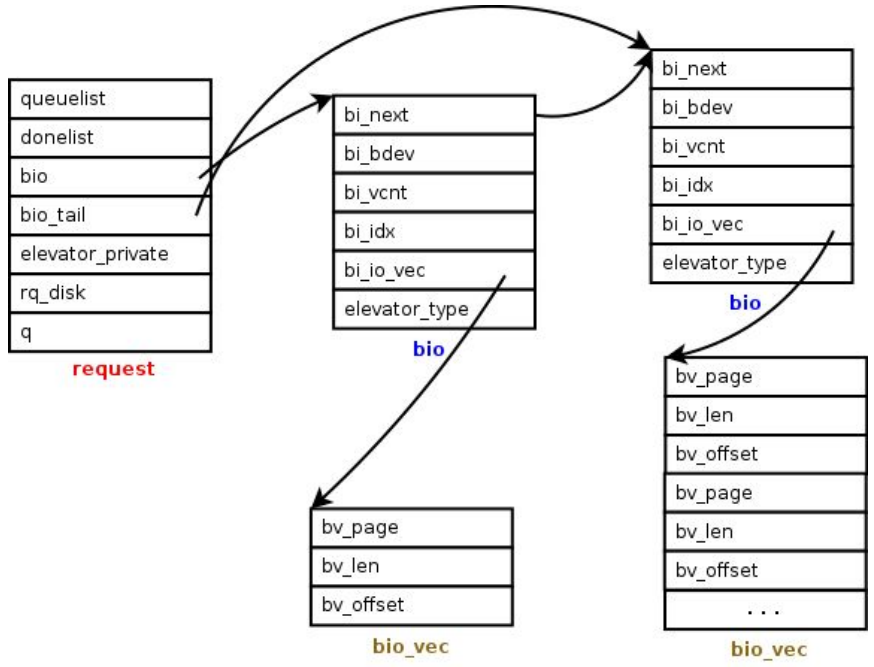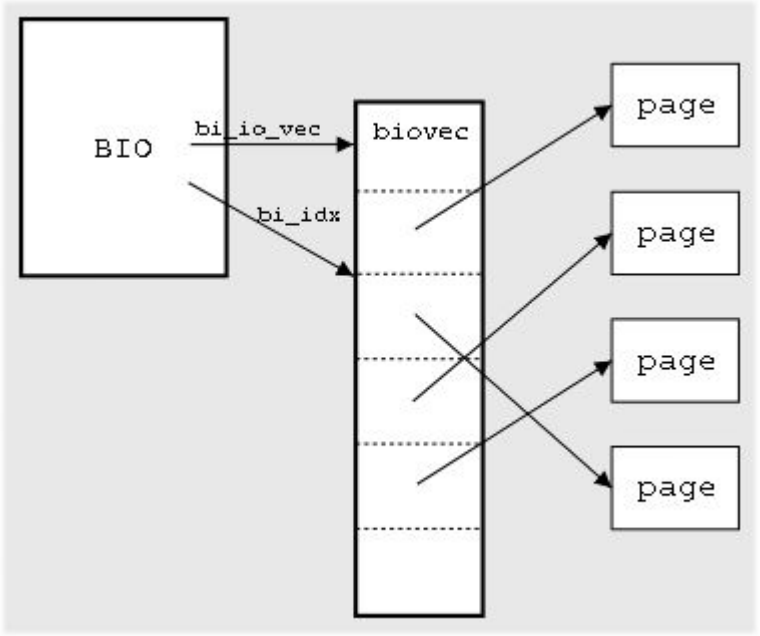**Block Device HW**

Block Device Controller

10s

- How are I/Os represented?
  - Data structures: bios and requests
- How are I/Os submitted?
  How are I/O completions handled?
  - What is the storage interface?
  - Put differently: What is the abstraction of the underlying storage devices?

(just wanted to mention this, not dwell on it)

10s



http://elixir.free-electrons.com/linux/latest/source/Documentation/block/biodoc.txt

https://www.kernel.org/doc/Documentation/block/request.txt

# Outline

1.  **File System**
    - How are I/Os exposed to programmers?
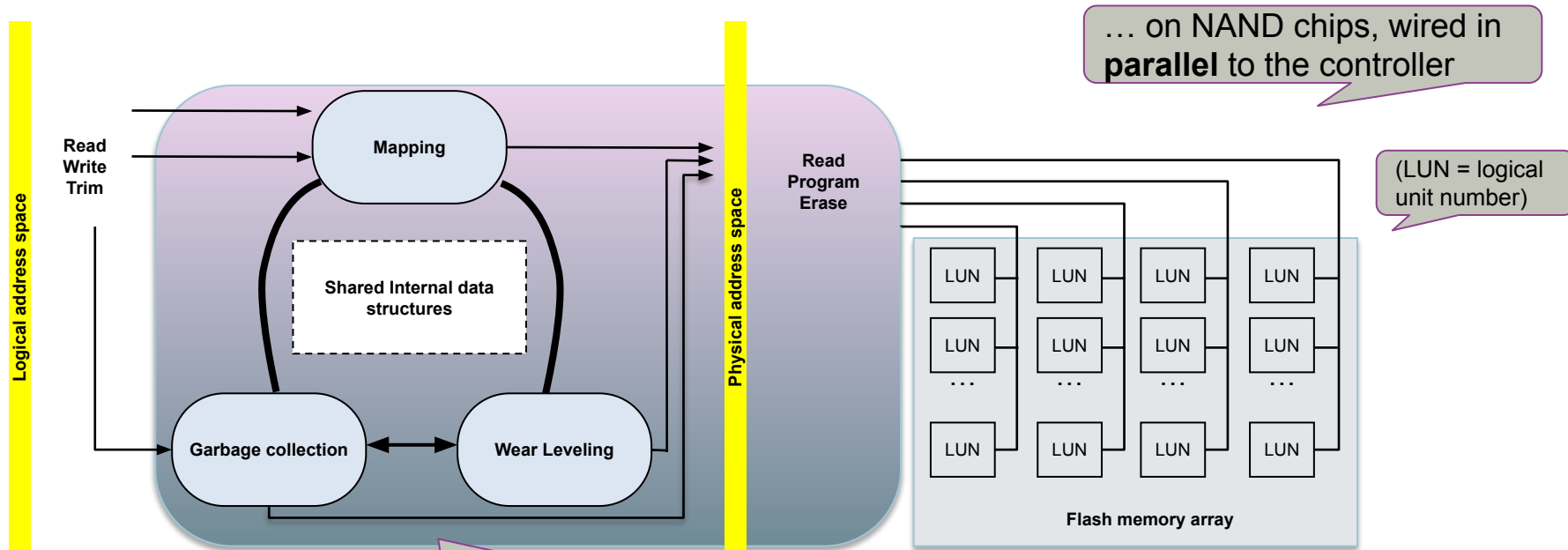    - How are I/Os handled on the host?
2.  **Storage devices**
    - What does a disk controller do?
    - How do hosts interact with disk controllers?
    - How does a disk controller perform a DMA?
3.  Computational Storage

# What does a SSD controller do?

1. Handles interactions with host
2. Maps logical ops onto physical reads, writes, erase



… on NAND chips, wired in **parallel** to the controller

(LUN = logical unit number)

Flash Translation Layer (FTL)

(added sophistication)

(wear leveling = equi-distribute the wear of disk)

(Copenhagen!)

much work on clever mapping (grouping) that lowers need for garbage collection

(in SSD, random access is just as fast as sequential access, assuming well-managed parallelism)

10s

Host

Interconnect

Storage
Device

- Physical Interconnect

  SATA / AHCI
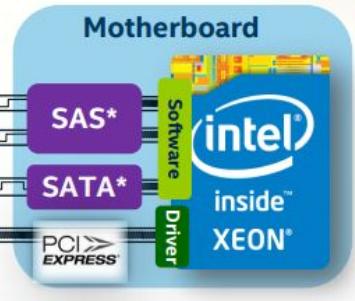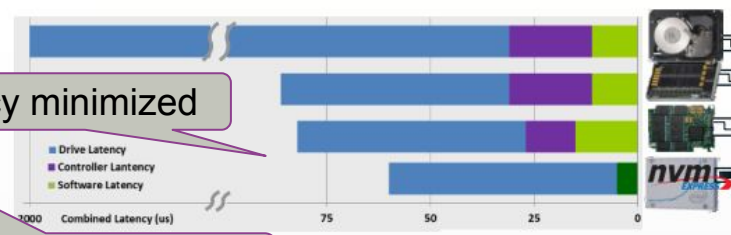
  PCIe

  Ethernet

  …

- Protocol

  SATA

  **NVMe**

  NVMf

# NVMe

## Why NVM Express?*
Standardized interface for non-volatile memory

**Motherboard**

SAS*
SATA*
PCI EXPRESS

Software
Driver

intel inside XEON

- Drive Latency
- Controller Lantency
- Software Latency

Combined Latency (us)   75   50   25   0

latency minimized

(why: with NVMe, memory is shared on the PCIe bus)

**Host**

| Controller Management | Core 0 | Core 1 | Core n |
|---|---|---|---|

Controller Management: Admin Submission Queue, Admin Completion Queue

Core 0: I/O Submission Queue, I/O Completion Queue

Core 1: I/O Submission Queue, I/O Submission Queue, I/O Completion Queue

Core n: I/O Submission Queue, I/O Completion Queue

MSI-X    MSI-X    MSI-X    MSI-X

**NVMe Controller**

(driver on host)

queues (submission, completion) shared between host and device
(access memory on host from device)

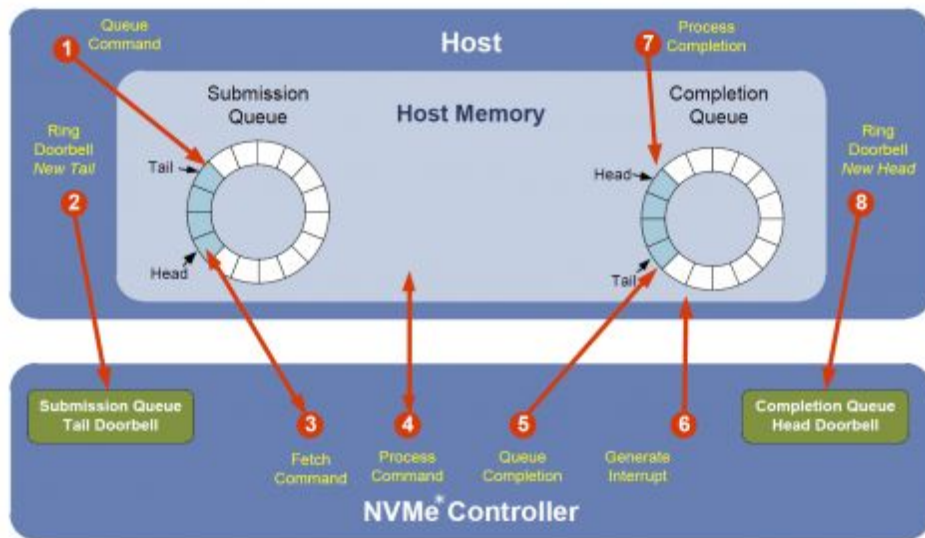(on the device)

### Simple Command Set – Optimized for NVM

| Admin Commands |
|---|
| Create I/O Submission Queue |
| Delete I/O Submission Queue |
| Create I/O Completion Queue |
| Delete I/O Completion Queue |
| Get Log Page |
| Identify |
| Abort |
| Set Features |
| Get Features |
| Asynchronous Event Request |
| Firmware Activate (optional) |
| Firmware Image Download (opt) |
| Format NVM (optional) |
| Security Send (optional) |
| Security Receive (optional) |

| NVM I/O Commands |
|---|
| Read |
| Write |
| Flush |
| Write Uncorrectable (optional) |
| Compare (optional) |
| Dataset Management (optional) |
| Write Zeros (optional) |
| Reservation Register (optional) |
| Reservation Report (optional) |
| Reservation Acquire (optional) |
| Reservation Release (optional) |

**Only 10 Admin and 3 I/O commands required**

IDF14

# Interconnect – NVMe

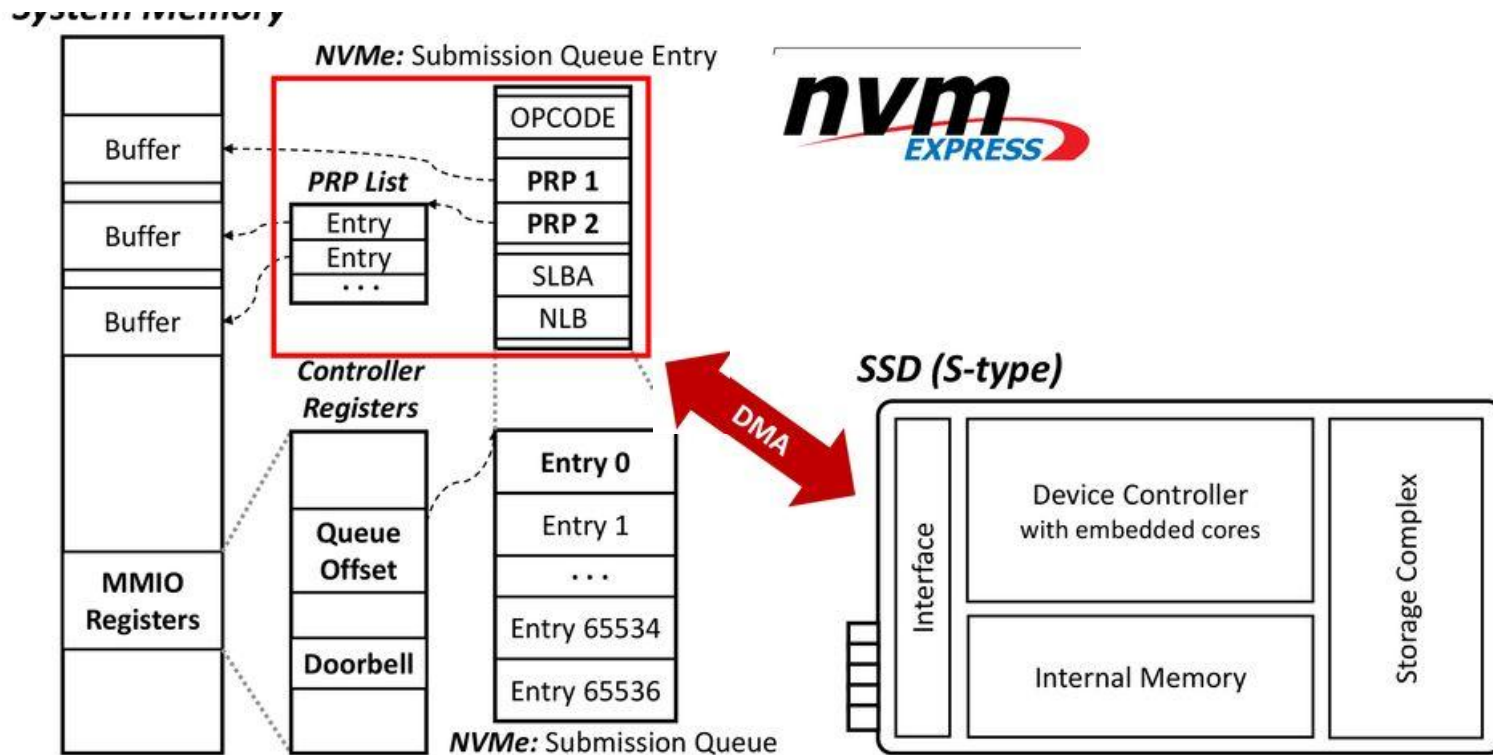| | AHCI | NVMe |
|---|---|---|
| **Maximum Queue Depth** | 1 command queue 32 commands per Q | 64K queues 64K Commands per Q |
| **Un-cacheable register accesses (2K cycles each)** | 6 per non-queued command 9 per queued command | 2 per command |
| **MXI-X and Interrupt Steering** | Single interrupt; no steering | 2K MSI-X interrupts |
| **Parallelism & Multiple Threads** | Requires synchronization lock to issue command | No locking |
| **Efficiency for 4KB Commands** | Command parameters require two serialized host DRAM fetches | Command parameters in one 64B fetch |
| **Driver Support** | Typically in-box | Installed with device |

if you want to be fast with how you do IO, then you have to manipulate, from your program, submissions and completions.

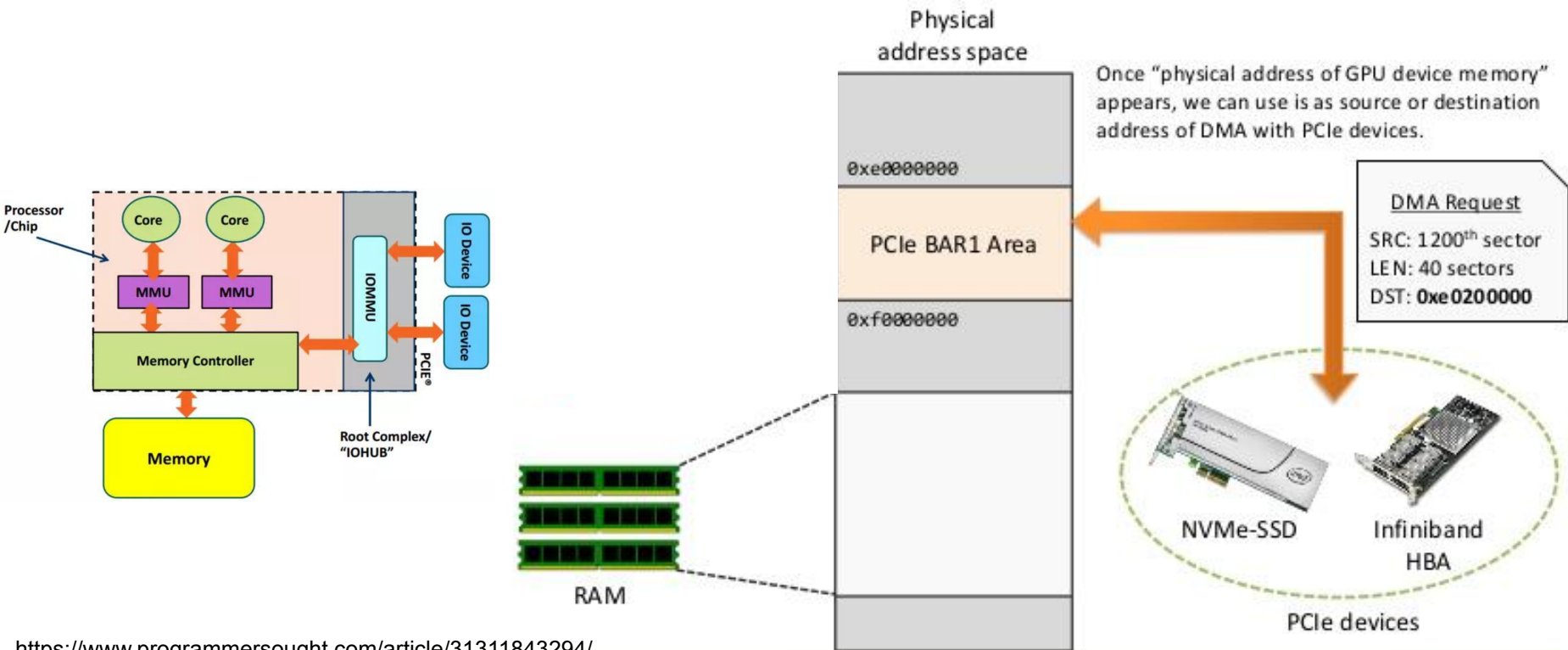https://www.flashmemorysummit.com/English/Collaterals/Proceedings/2013/20130812_PreConfD_Marks.pdf

# How to deal with data transfers?

where things are stored
(recall virtual memory)

30s

https://slideplayer.com/slide/14772624/   11.11.2020 · 45

# How to deal with data transfers?



Processor/Chip
Core  Core
MMU  MMU
Memory Controller
Memory
IOMMU
Root Complex/ "IOHUB"
PCIE®
IO Device
IO Device

Physical address space

Once "physical address of GPU device memory" appears, we can use is as source or destination address of DMA with PCIe devices.

0xe0000000

PCIe BAR1 Area

0xf0000000

RAM

DMA Request
SRC: 1200th sector
LEN: 40 sectors
DST: 0xe0200000

NVMe-SSD  Infiniband HBA

PCIe devices

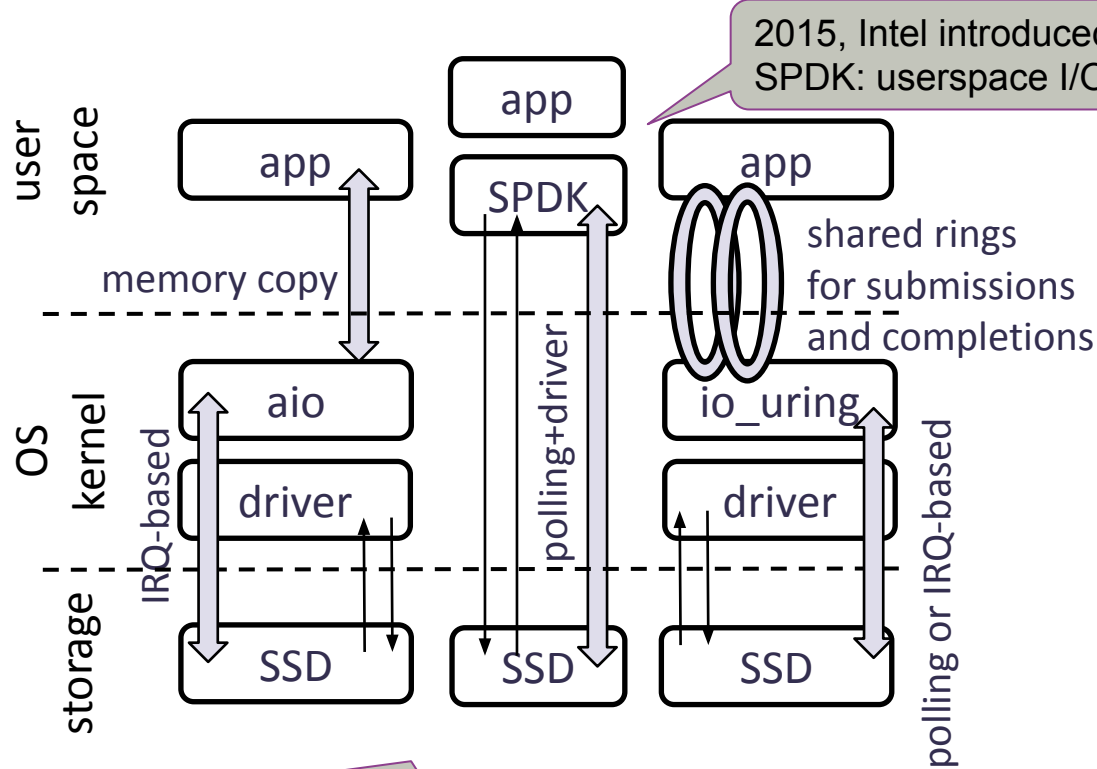18    PGconf.ASIA 2019 - Full-throttle running on Terabytes log-data    ᴴHeteroDB

# NVMe Interfaces

usually, data must be transferred from userspace to kernelspace before it can be transferred to a device. (copy, ctx switch, …)
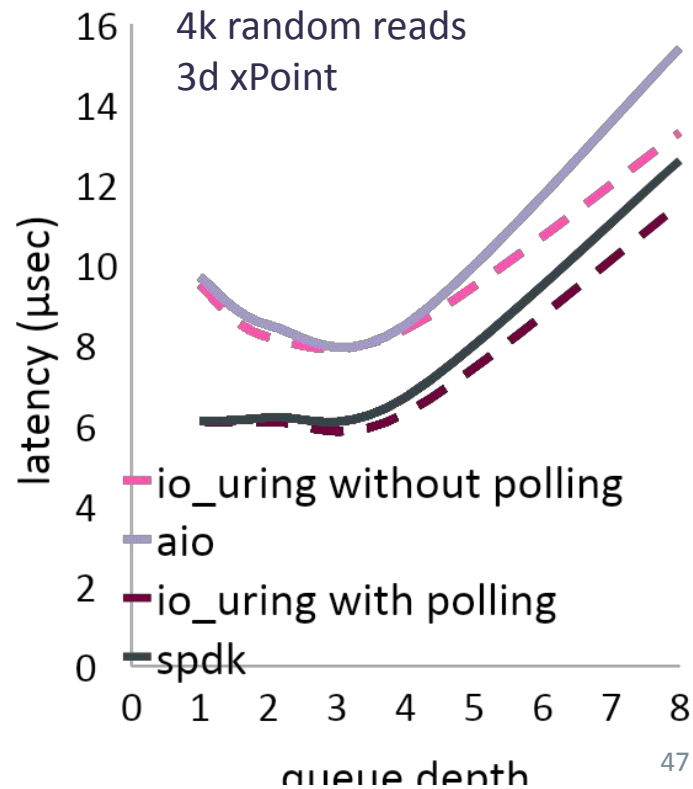
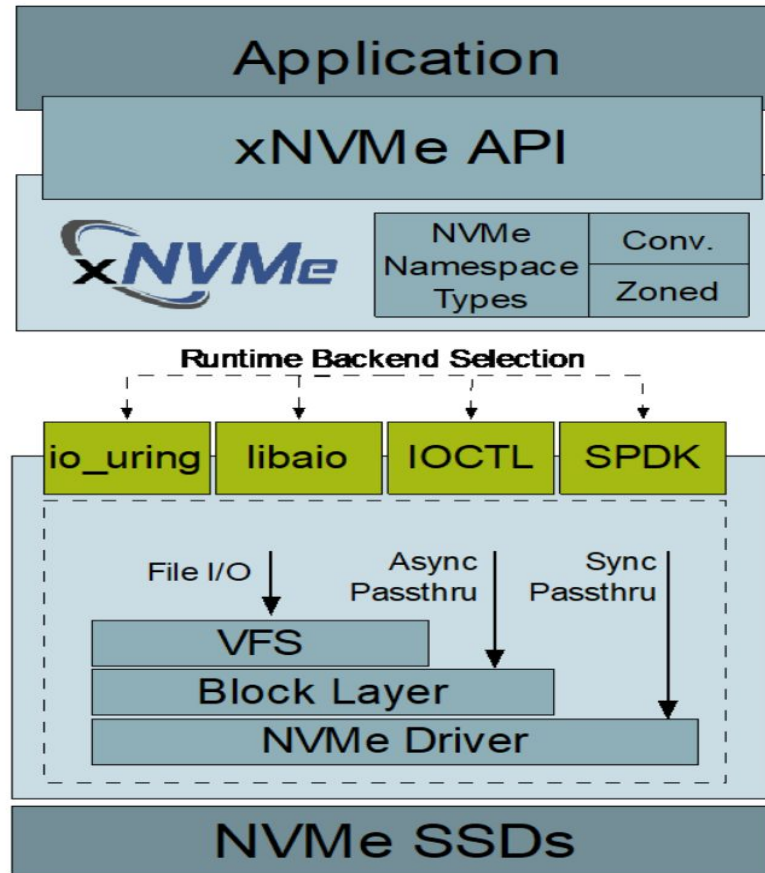sources: Faster IO through io_uring & Efficient I/O with io_uring & J.Axboe

2015, Intel introduced SPDK: userspace I/O.

SPDK matches io_uring & aio in performance.



user space

app

memory copy

app

SPDK

app

shared rings for submissions and completions

OS kernel

aio

driver

IRQ-based

polling+driver

io_uring

driver

polling or IRQ-based

storage

SSD

SSD

SSD

direct mapping from userspace to datastructures on device.

4k random reads
3d xPoint

latency (μsec)

16
14
12
10
8
6
4
2
0

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

0  1  2  3  4  5  6  7  8

queue depth

Samsung (Copenhagen)

abstracts over these; uniform API

# Logical Address Space

Recall, the file systems slides a few slides ago …

A block device is an array of blocks.
To each block is associated a number,
a Logical Block Address (LBA)
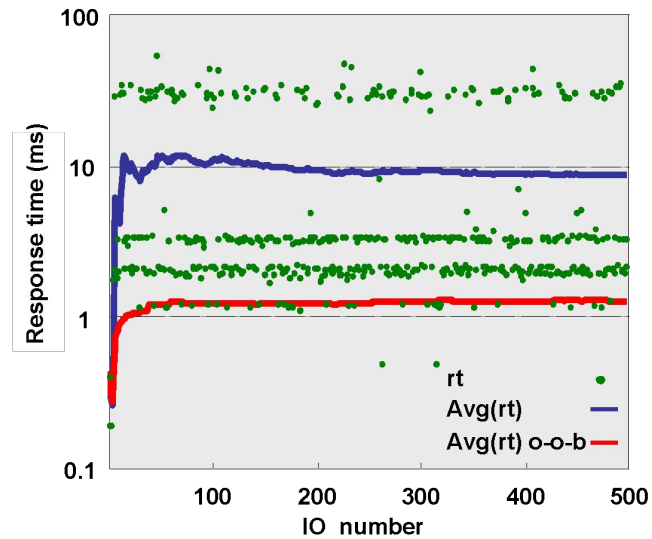
# Block device

**CIDR 2009**

Measuring Samsung SSD RW performance
- Out-of-the-box … and after filling the device!!! (similar behavior on Intel SSD)



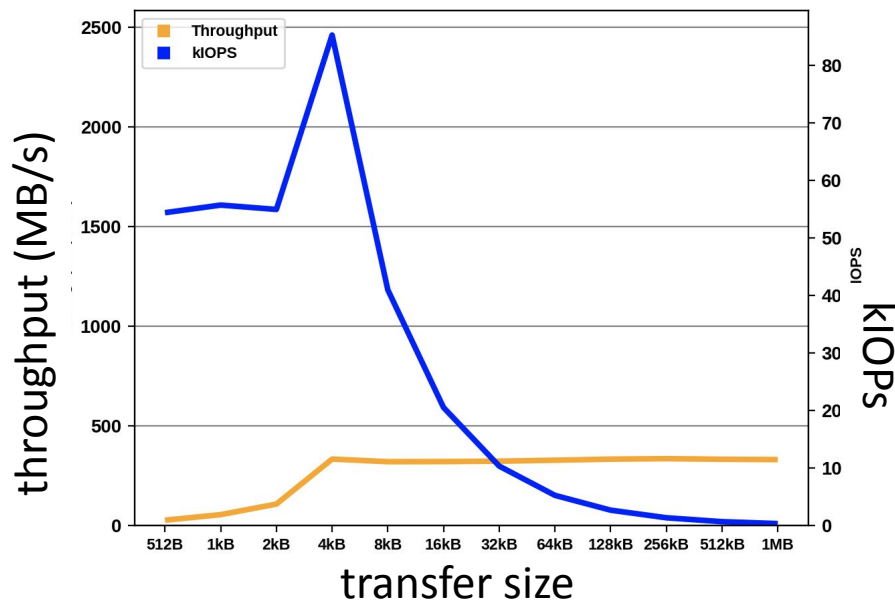*Random Writes – Samsung SSD*
*Out of the box*

*Random Writes – Samsung SSD*
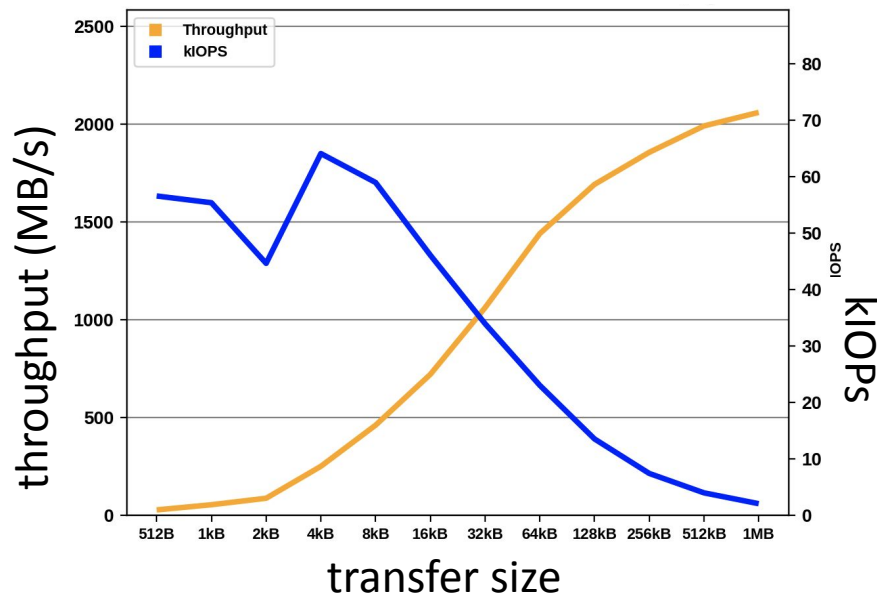*After filling the device*

# Performance contract?

random writes- source: AnandTech 2019



Samsung SSD with Z-NAND



Intel Optane

## No intrinsic performance characteristics for SSDs (equipped with a generic FTL)

(flash translation layer)
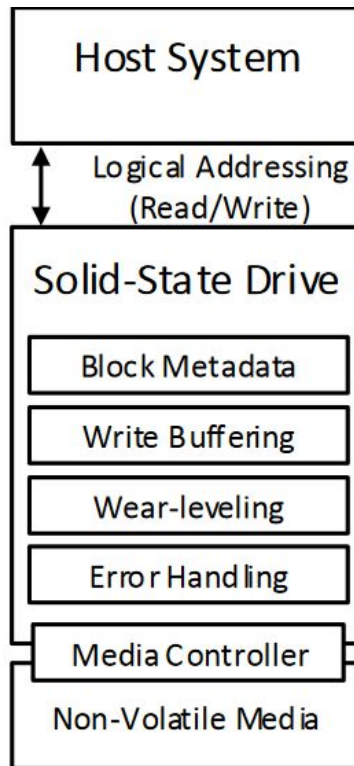
# Open-Channel

**Block**    **Open-Channel**

Physical address space exposed
- host can make decisions about
  ***data placement*** & ***I/O scheduling***

SSD management split between
- back-end (embedded on SSD)
  ***block metadata*** & ***wear levelling***
  (for warrantee)
- front-end (host-based) FTL
  ***mapping of logical to physical***
  ***address spaces, overprovisioning,***
  & ***garbage collection***

**Block**

Host System

Logical Addressing
(Read/Write)

Solid-State Drive

Block Metadata

Write Buffering

Wear-leveling

Error Handling

Media Controller

Non-Volatile Media

**Open-Channel**

Host System

Write Buffering

Physical Addressing
(Read/Write/Erase)

Open-Channel
Solid-State Drive

Block Metadata

Wear-leveling

Error Handling

Media Controller

Non-Volatile Media

# Open-Channel

alibaba, Western Digital, Microsoft, … but didn't become part of NVMe standard.



LightNVM: The Linux Open-Channel SSD Subsystem

Matias Bjørling, *CNEX Labs, Inc. and IT University of Copenhagen; Javier Gonzalez, CNEX Labs, Inc.;* Philippe Bonnet, *IT University of Copenhagen*

https://www.usenix.org/conference/fast17/technical-sessions/presentation/bjorling

**Taking control of SSDs with LightNVM**

Matias Bjørling · 1st
Director, Emerging System Architectures at Western Digital
Copenhagen, Capital Region, Denmark · 500+ connections ·

Principal Software Engineer | SSDR R&D Center Lead
Samsung Electronics
Jan 2019 – Present · 10 mos
Copenhagen Area, Capital Region, Denmark
I established and lead Samsung Semiconductor Denmark Research (SSDR) – Samsung's Memory Solutions first R&D center in Europe and fifth worldwide.

Javier González

**CNEX LABS AND BROADCOM WIN MOST INNOVATIVE FLASH MEMORY TECHNOLOGY AWARD AT FLASH MEMORY SUMMIT 2017**

*CNEX Labs Open-Channel SSD technology and Broadcom's NetXtreme S-Series SOC provide unprecedented IO isolation and scalability for hyperscale and cloud service providers*

OPEN CHANNEL SSDs:
EXTENDING SPDK'S REACH

**Alibaba Open Channel Ecosystem**

As Alibaba's **strategic partner** on Open Channel SSDs, Intel has worked with Alibaba extensively since 2017 to co-develop and co-validate this innovative solution. Alibaba's strength as a leading cloud service provider combined with Intel's strength as the leading memory and storage innovator puts us in a position to deliver the industry's 1st Open Channel SSD product.
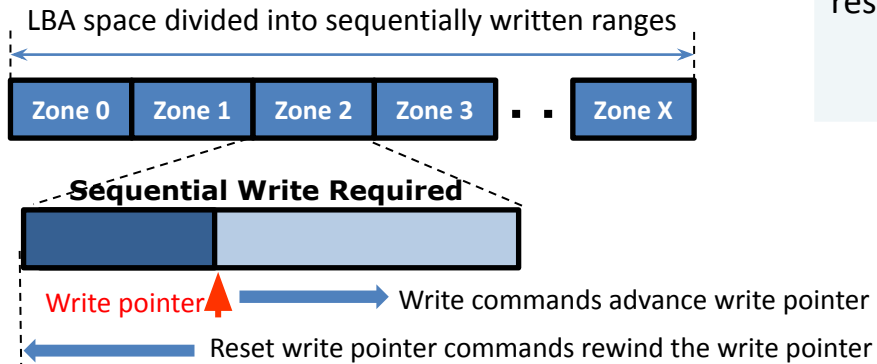
- Alibaba is collaborating with major vendors in industry to build an ecosystem for Open Channel SSD
- Share development & debug resources
- Reduce time & complexity for SSD qualification in Alibaba
- Massive deployment in 2019

© Alibaba Group 2018

# Zoned Namespaces (ZNS)

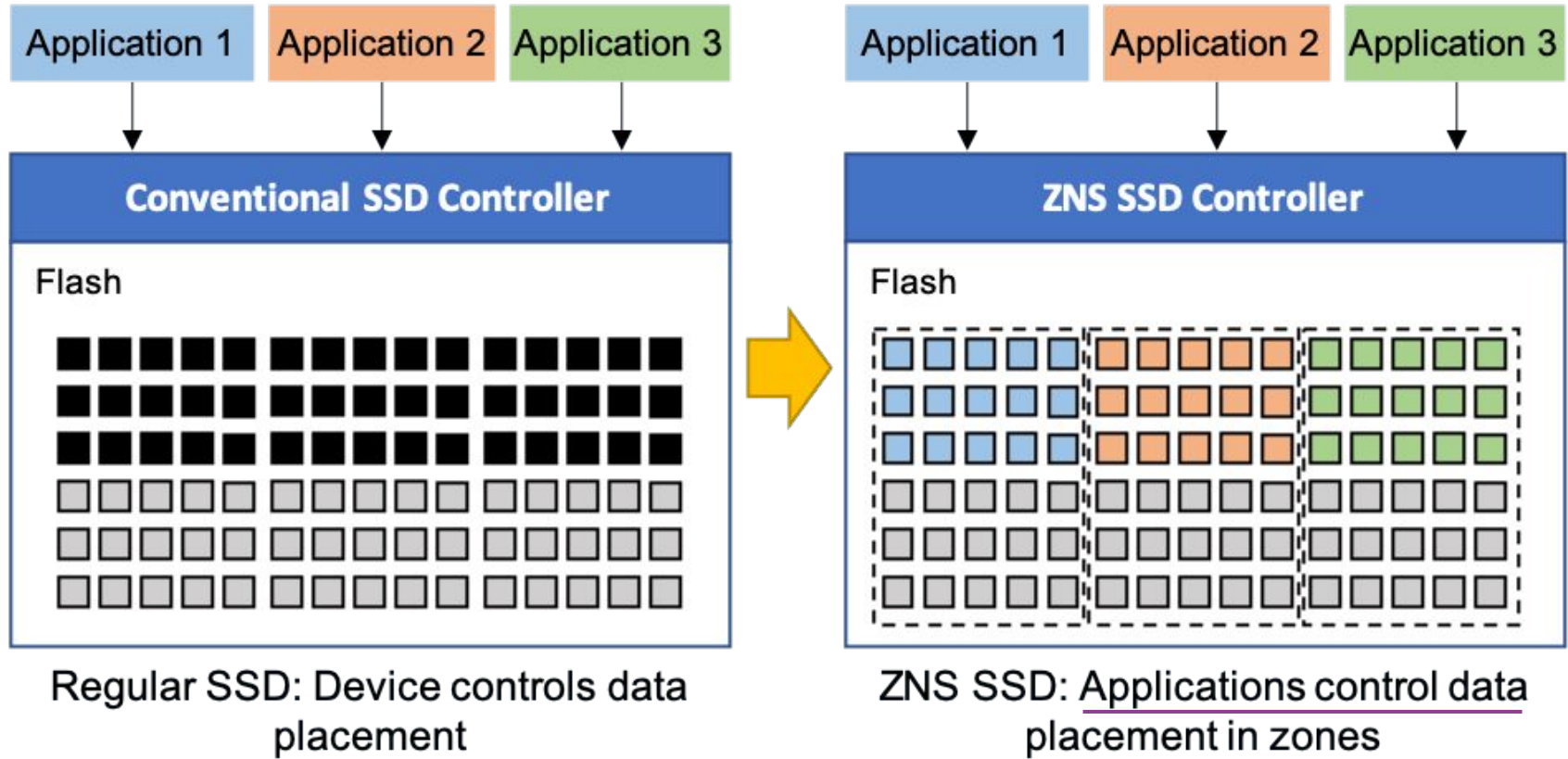| Mode | Host managed | Host aware |
|---|---|---|
| Host responsibility | Explicit zone transition Write command Write pointer per zone | Implicit zone transition Append command at large queue depth |
| SSD responsibility | Zone placement I/O scheduling | Zone placement I/O scheduling sequential writes within a zone |

you: how to map workload onto zones, disk: manage parallelism best possible

LBA space divided into sequentially written ranges

| Zone 0 | Zone 1 | Zone 2 | Zone 3 | . . | Zone X |

**Sequential Write Required**

Write pointer ▲ ➔ Write commands advance write pointer

⬅ Reset write pointer commands rewind the write pointer

one level of control on device.
(not as good as open-channel)

30s

https://nvmexpress.org/new-nvmetm-specification-defines-zoned-namespaces-zns-as-go-to-industry-technology/



IT UNIVERSITY OF COPENHAGEN

# Outline

1. **File System**
   - How are I/Os exposed to programmers?
   - How are I/Os handled on the host?
2. **Storage devices**
   - What does a disk controller do?
   - How do hosts interact with disk controllers?
   - How does a disk controller perform a DMA?
3. **Computational Storage**

> conventional SSD: a lot of complexity on device; gets in way of host.
> open-channel: host did too much work. ZNS is a compromise, still not perfect.

> idea: you *should* be able to program your SSDs.

# Computational Storage

## 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
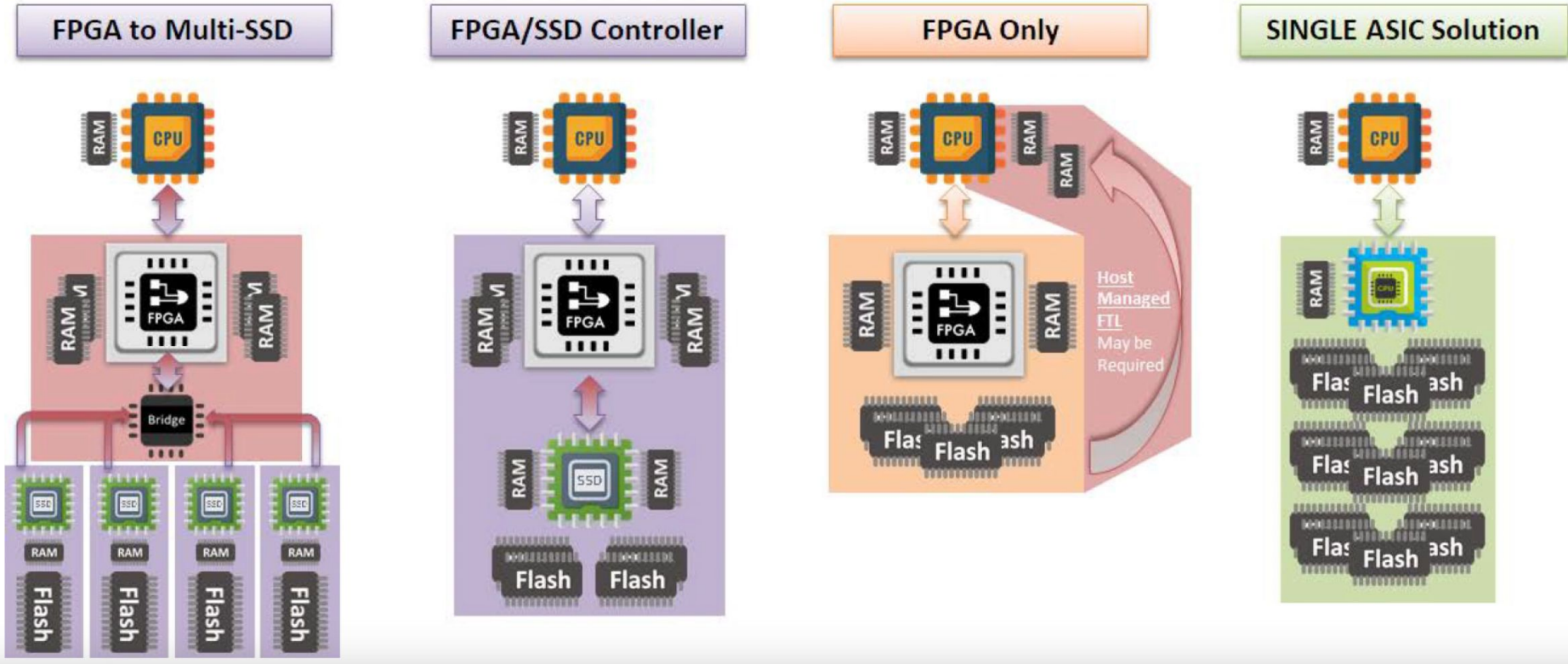  - » need programming environment in controller

Jim Gray, NASD Talk, 6/8/98
http://jimgray.azurewebsites.net/jimgraytalks.htm

## Computational storage = Computation on the IO Path

# Computational storage

storage industry has defined an architecture for computational storage. NVMe will be standard.



**Specialized storage interface + functionality offload**

# Communication Abstraction

SEND(link_name, outgoing_message_buffer)

RECEIVE(link_name, incoming_message_buffer)

## Communication Link

Source: Saltzer and Kaashoek

# Take-Aways

File abstraction is one of Unix enduring contribution. Beautiful example of a deep module. You should be able to describe the data structures and name mapping steps involved in file system operations.

NVMe as host/storage interface. NVMe manages completion/submission queues. NVMe namespaces include block device and zones.

With computational storage, storage devices are moving from a memory to a communication abstraction.