

Data storage

- A challenge in **every age**: from the first stone carvings to the digital revolution.
- **Principles**: safety, validity, integrity.
- **Data**: piece of information.
- **Digital data**: bits, bytes, k, M, G, T, P, E ...

Virtualization

- A way of **decoupling data storage** from its physical appearance.
- At **any levels** of data storage hierarchy.
- Not a magic word, but a set of real technical solutions: e.g. volume management, device mapping, clustered file systems.
- Layers **hide details** beneath, **offer services** above.

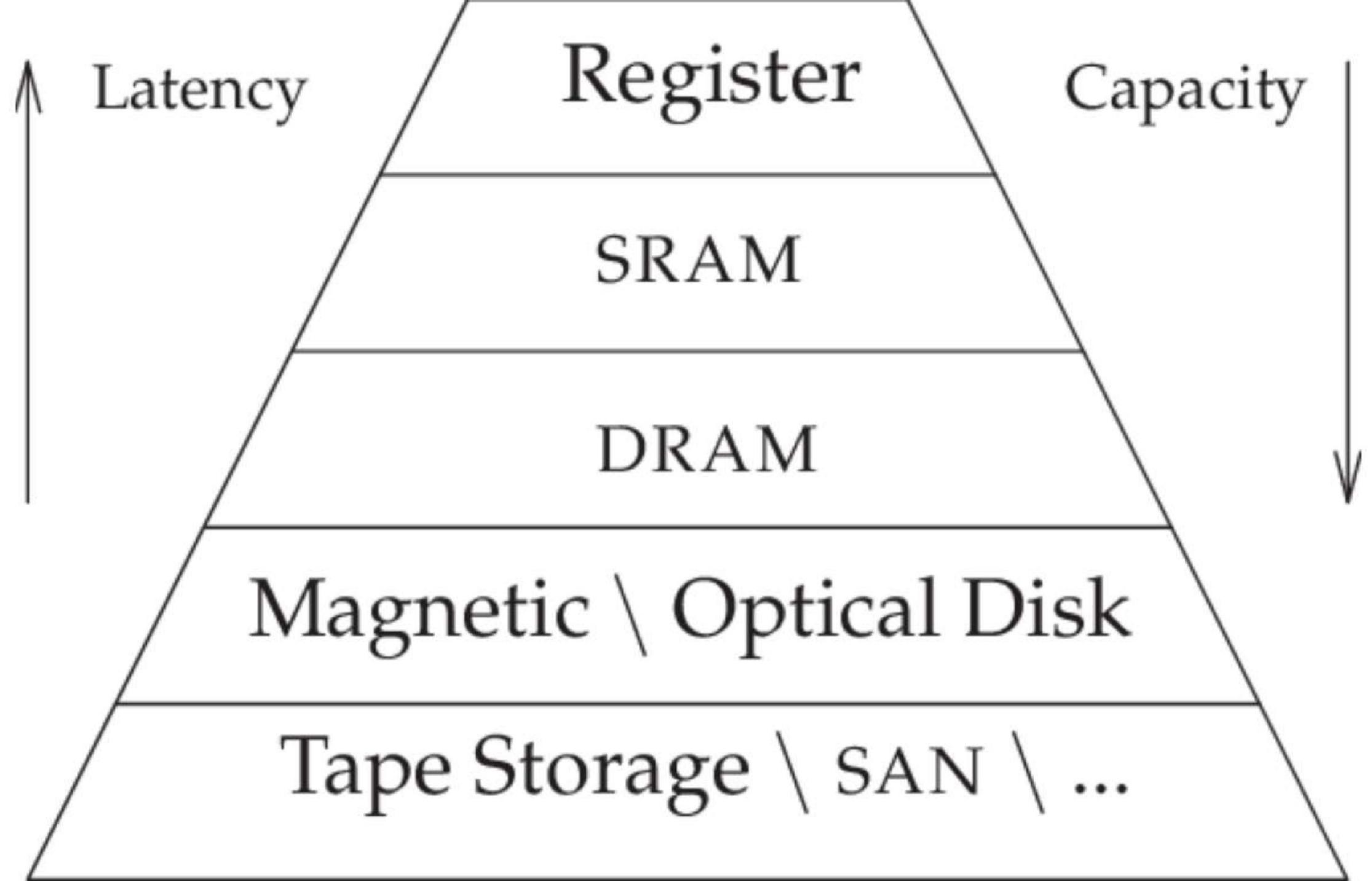
[1] https://en.wikipedia.org/wiki/Storage_virtualization

Virtualization methods

Block device based: it creates a mapping between the provider (storage) and user of data (host) on block level

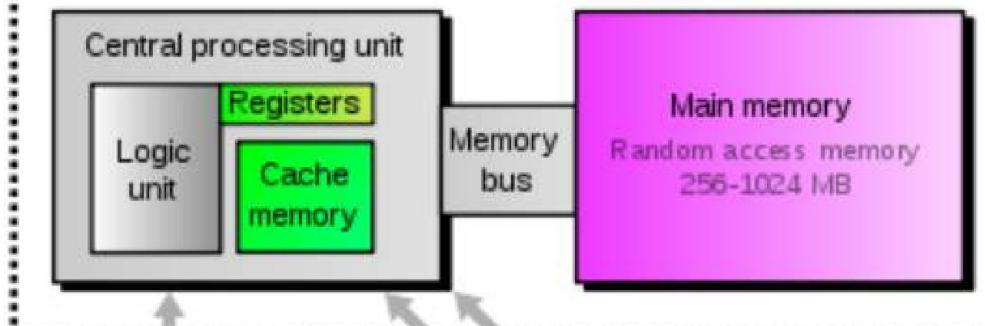
- Host based (LVM, LDM);
- Storage device based: RAID controllers;
- Network based: iSCSI, FC;

File system based: it offers generalized access to different file systems in a uniform way, e.g. VFS.

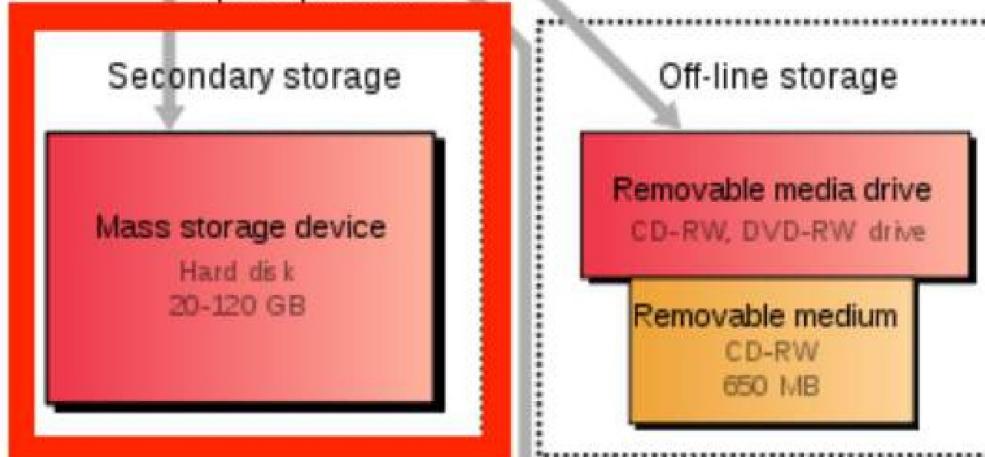


Capacity

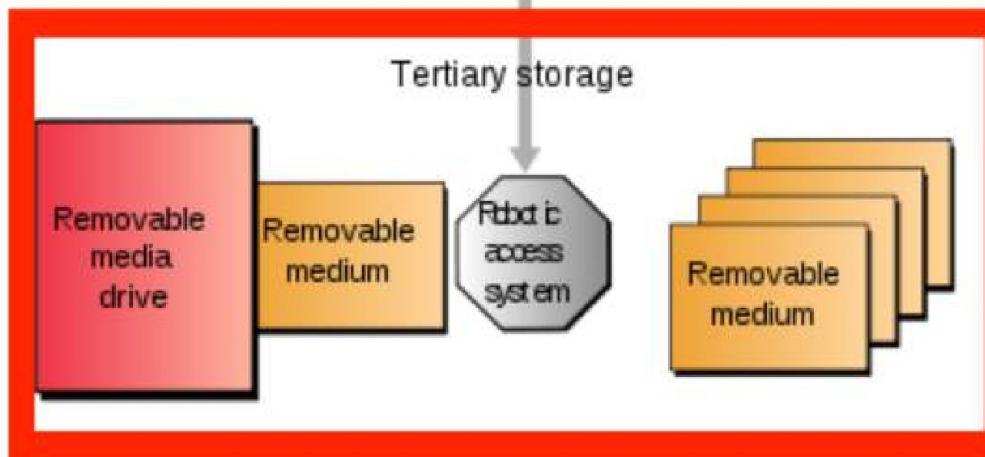
Primary storage



Input/output channels



Tertiary storage



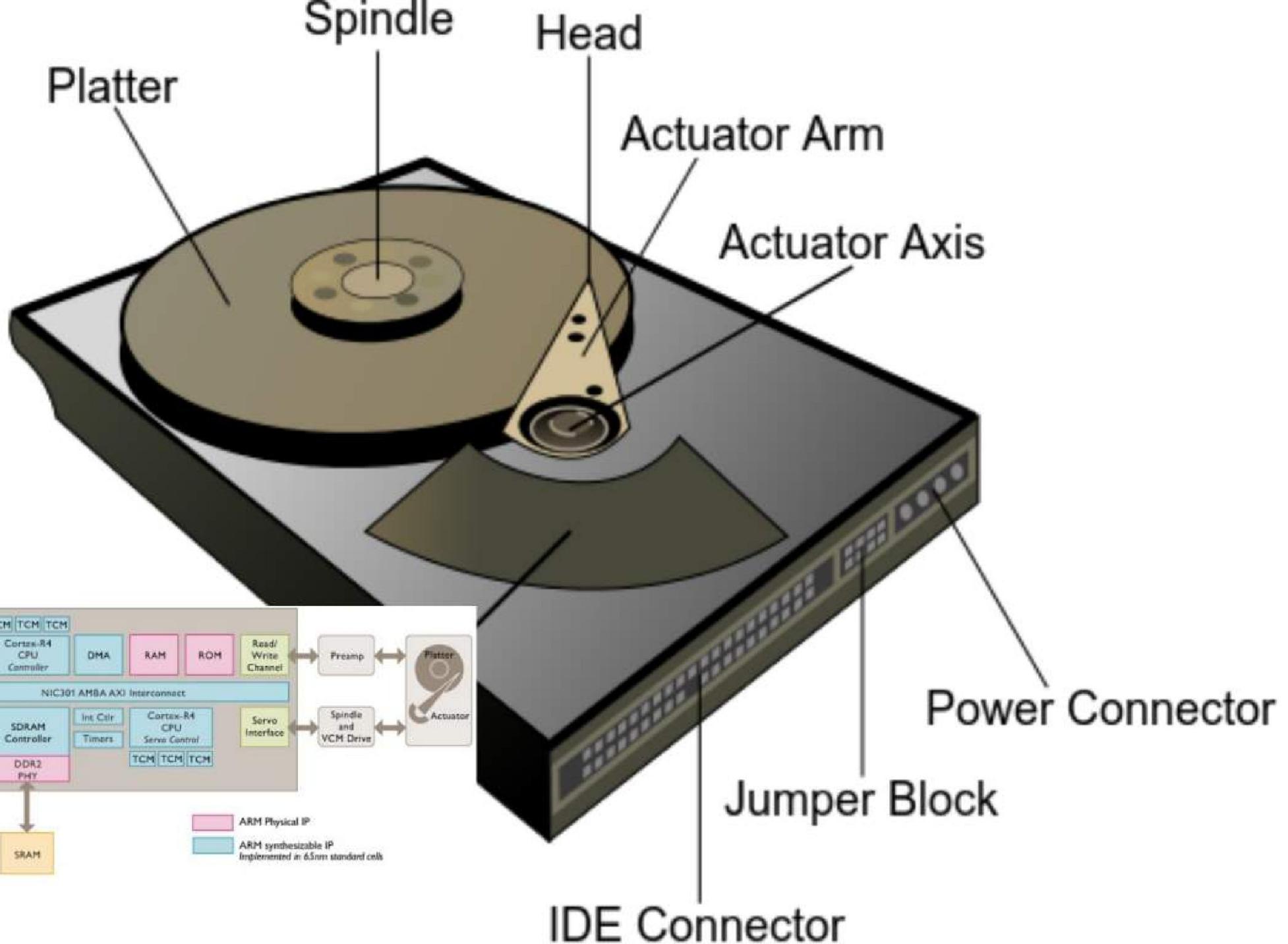
High

Volume

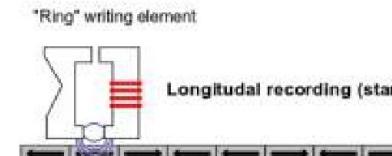
Storage

Storage

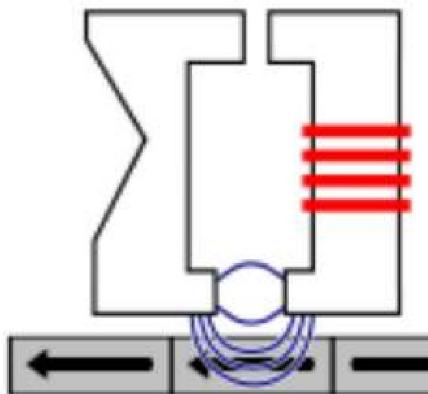
Physical
(dis)



- **Magnetic recording:**
 - Magnetic cover on the platters;
 - **Direction changes** in magnetization store bits;
- Disk controller is used to control the rotation as well as the head positioning;
- Data is stored in **blocks** from 512 Byte (classic) to 4096 Byte (modern);
- **Architecture:** media - buffer - interface;
- **Longitudinal** and **vertical** recording.



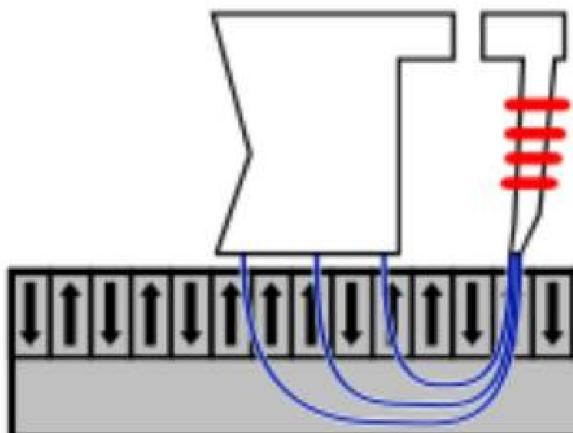
"Ring" writing element



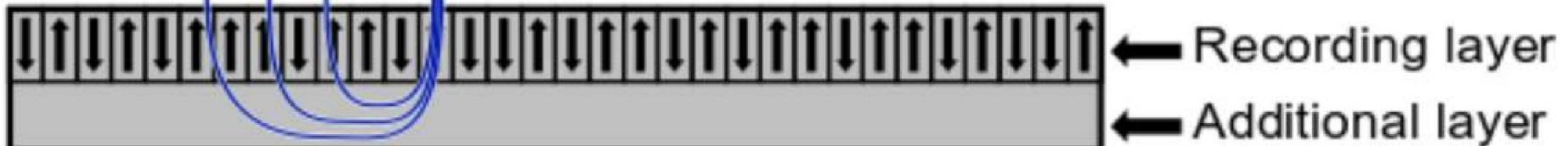
Longitudinal recording (standard)



"Monopole" writing element



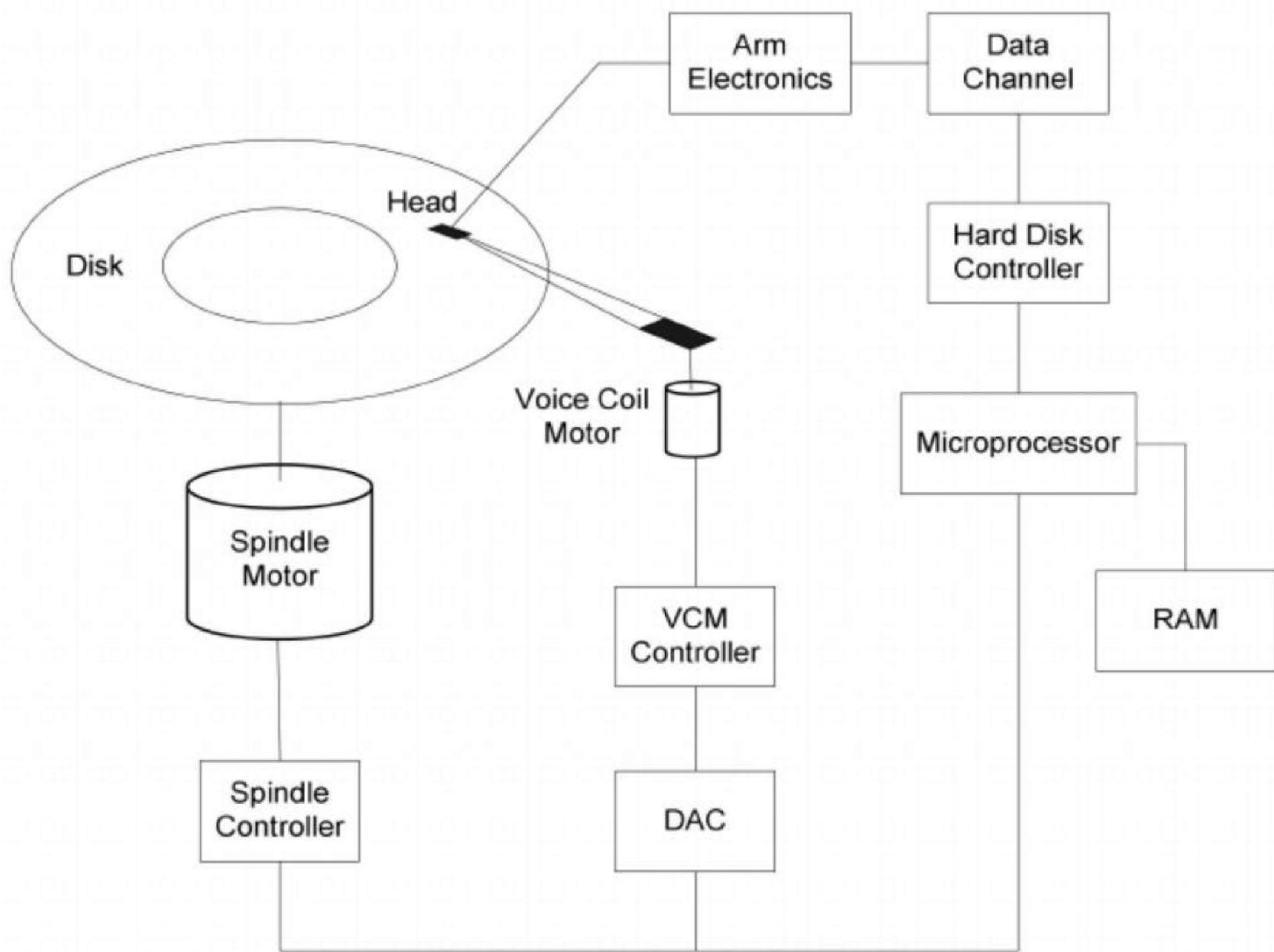
Perpendicular recording



Characterized by

- **capacity** (Bytes, 1-10 Tbytes)
 - OS reports less for many reasons;
 - error correction: 10% of space;
 - redundancy + FS structures;
 - decimal or binary prefix;
 - **speed** (rpm, 4200-15000 rpm)
 - **latency** (sec, 2-6 ms);
 - **throughput** (bit/s, 1-3 Gbit/s);
 - **form factor** (2.5" vs 3.5");
 - **energy consumption** (W);
- **Design** (IDE, SATA, SAS, 5400-15000 rpm, 2.5" or 3.5" form factor);
 - **Laptops** (IDE, 2.5" form factor, 4200-5400 rpm, 10-15 W power consumption);
 - **Enterprise** (SAS, 3.5" form factor, 10000-15000 rpm, 100-200 W power consumption);

- **Desktop:** 4-12 Tbytes, 3.5", 0.5 Gbit/s, 5400-10000 rpm;
- **Laptop:** 1-4 Tbytes, 2.5", 0.5 Gbit/s, 4200-7200 rpm;
- **Enterprise:** 1-12 Tbytes, any, 1.6 Gbit/s, 10000-15000 rpm;



Interfaces:

the physical + logical link

- IDE, EIDE;
- ATA, PATA, SATA;
- SCSI;
- SAS;
- FC.

Magnetic tapes

- Has a long history, since 1951, IBM + DEC were pioneers.
- Allows **sequential data storage**.
- Operations: record, play, fast forward, rewind.
- BOT ... EOT.
- Media write: **linear** vs. **helical** tracks.
- **Exposed tapes** ... **cartridges** (with flash memory).
- Access time: 3 magnitude longer than those of HDD's.
- **Compression**: redundant data compressed in a limited window buffer (2:1 typical).
- **Encryption**: for protecting data. **Key management!**



Linear Tape Open (LTO):

LTO1, 100 GB, 200 GB, 20 MB/s;
LTO2, 200 GB, 400 GB, 40 MB/s;

...
LTO7, **6 TB**, 15 TB, 750 MB/s;
LTO8, **12 TB**, 30 TB, 900 MB/s;

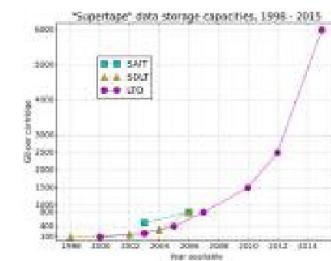
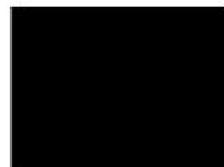
...
One step backward compatibility.

Error correction: CRC + repeated read.
Data blocks: Mbytes of block size.
Designed for data archiving, i.e. 15-30 years.
Data structures: data + metadata.
.tar format.
Labels, BAR-codes.
Cleaning tapes.



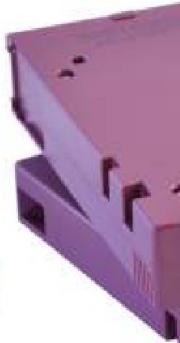
Tape drives + tape libraries:

- Take tapes in parallel.
- Collect a huge set of cartridges.
- Offer **standard interfaces** to hosts.



[4] https://en.wikipedia.org/wiki/Tape_drive

- Has a long history, since 1951, IBM + DEC were pioneers.
- Allows **sequential data** storage.
- Operations: record, play, fast forward, rewind.
- BOT ... EOT.
- Media write: **linear** vs. **helical** tracks.
- **Exposed tapes** ... **cartridges** (with flash memory).
- Access time: 3 magnitude longer than those of HDD's.
- **Compression**: redundant data compressed in a limited window buffer (2:1 typical).
- **Encryption**: for protecting data. **Key management!**



Linear Tape Open (LTO):

LTO1, 100 GB, 200 GB, 20 MB/s;

LTO2, 200 GB, 400 GB, 40 MB/s;

...

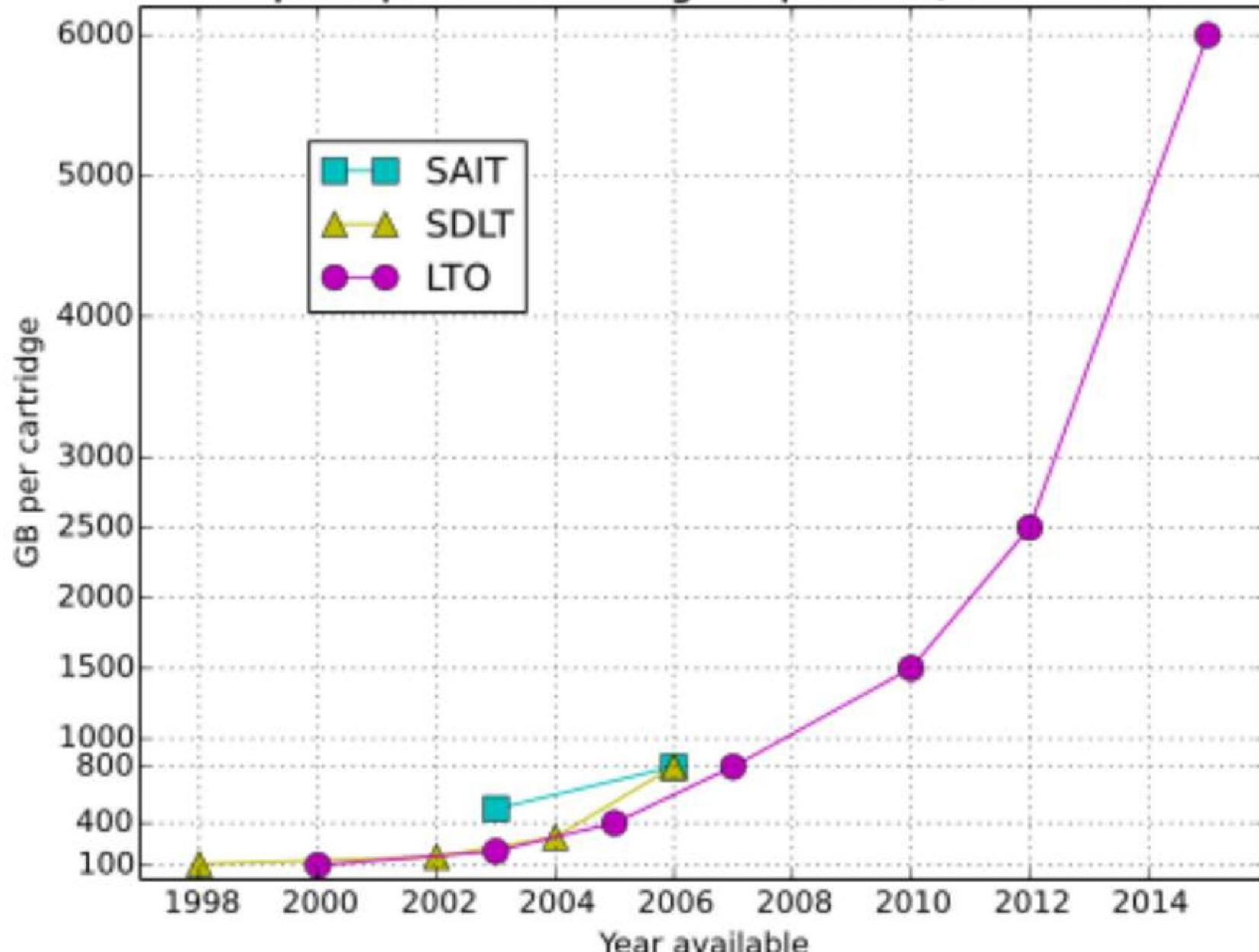
LTO7, **6 TB**, 15 TB, 750 MB/s;

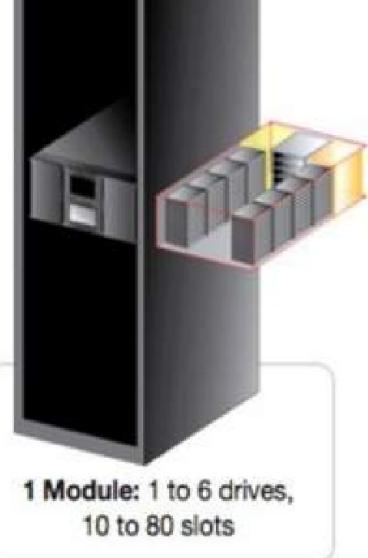
LTO8, **12 TB**, 30 TB, 900 MB/s;

....

One step backward compatibility.

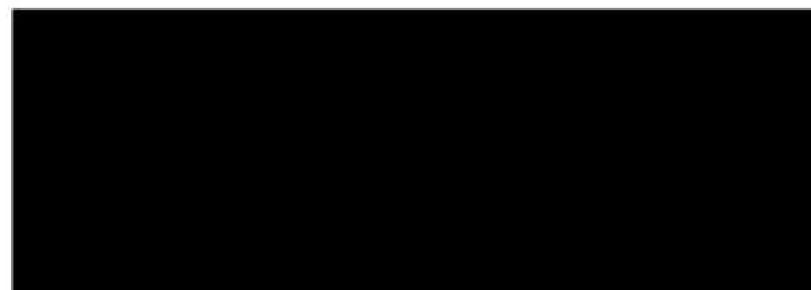
"Supertape" data storage capacities, 1998 - 2015

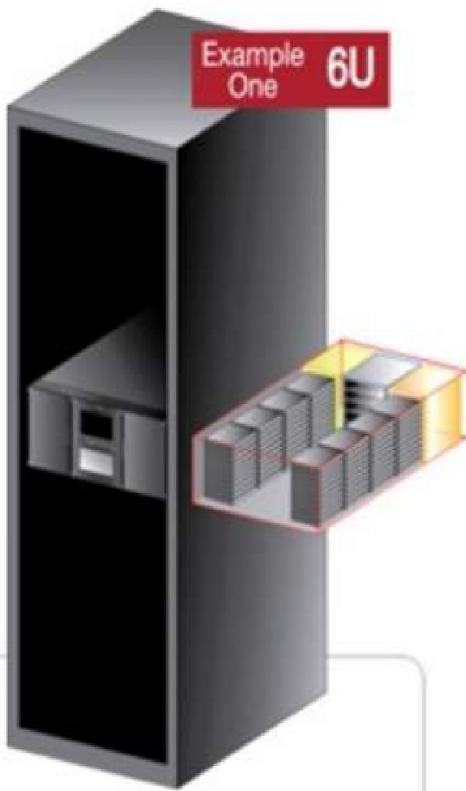




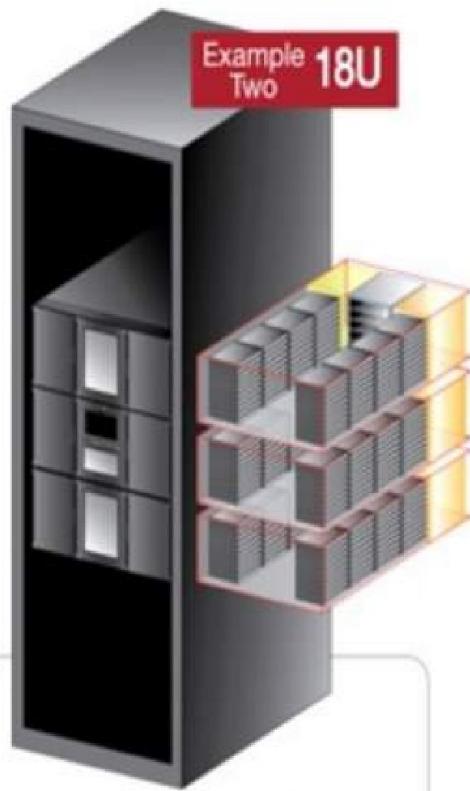
Tape drives + tape libraries:

- Take tapes in parallel.
- Collect a huge set of cartridges.
- Offer **standard interfaces** to hosts.

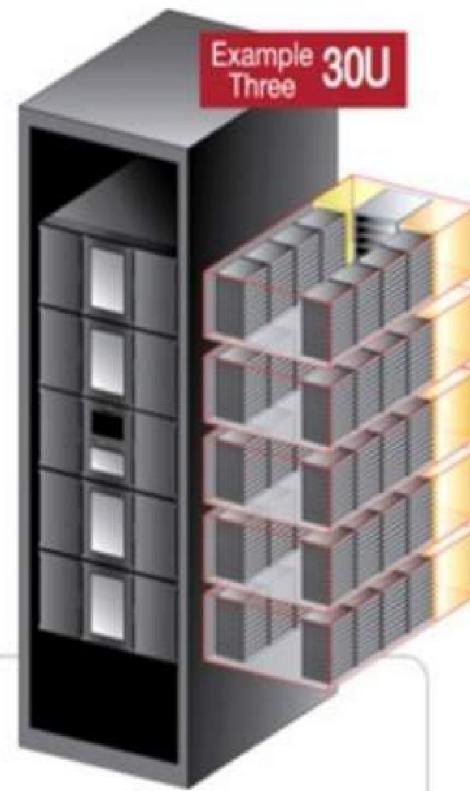




1 Module: 1 to 6 drives,
10 to 80 slots



3 Modules: 1 to 18 drives,
10 to 240 slots

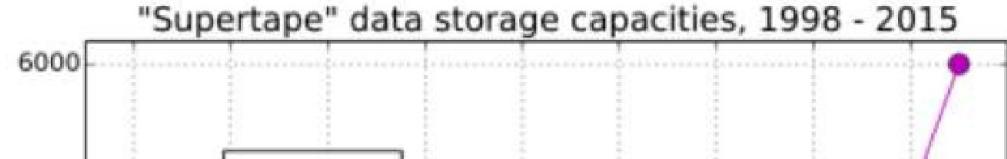


5 Modules: 1 to 30 drives,
10 to 400 slots



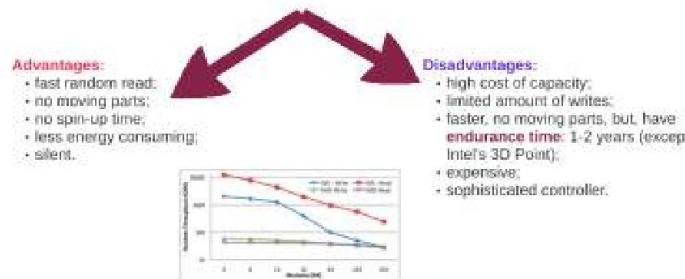
7 Modules: 1 to 42 drives,
10 to 560 slots

- Error correction: CRC + repeated read.
- Data blocks: Mbytes of block size.
- Designed for data archiving, i.e. 15-30 years.
- Data structures: data + metadata.
- .tar format.
- Labels, BAR-codes.
- Cleaning tapes.

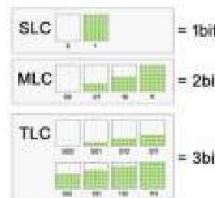


Solid State Drives

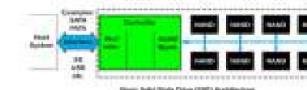
- A sort of **flash memory** (NAND), meant to replace (suplement) HDDs. Same size, similar form factors (exception M.2), same interfaces:
 - **SAS**, 12 Gbit/s;
 - **PCI-e 3.0**, 31 Gbit/s;
- Non-volatile, stores data when losing power.
- Stores data **as electrical charges**.



- Information is stored in **floating gate transistors**, i.e. read and written, holds charge for a long period.
- **Levels** of data storage:
 - single level cell: **1 bit/cell**, fast, endurable (50000-100000 erases), small;
 - 2-level cell: **2 bits/cell**, slower, less endurable (1000-10000 erases), medium;
 - 4-level cell: **4 bits/cell**, slower, less endurable, large.
- **Capacities**: 32 Gbyte - 4 Tbyte.
- Combined HDD + SSD devices.

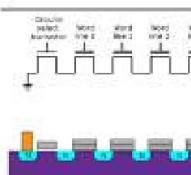


- Controller:**
- R/W cache + battery;
 - bad block mapping;
 - wear leveling;



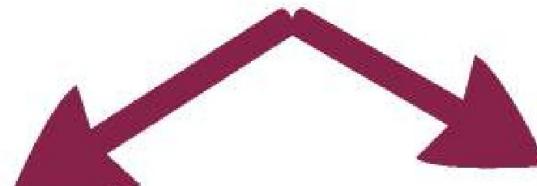
Solid State Drives

- A sort of **flash memory** (NAND), meant to replace (suplement) HDDs. Same size, similar **form factors** (exception M.2), **same interfaces**:
 - **SAS**, 12 Gbit/s;
 - **PCI-e 3.0**, 31 Gbit/s;
- Non-volatile, stores data when losing power.
- Stores data **as electrical charges**.



Advantages:

- fast random read;
- no moving parts;
- no spin-up time;
- less energy consuming;
- silent.



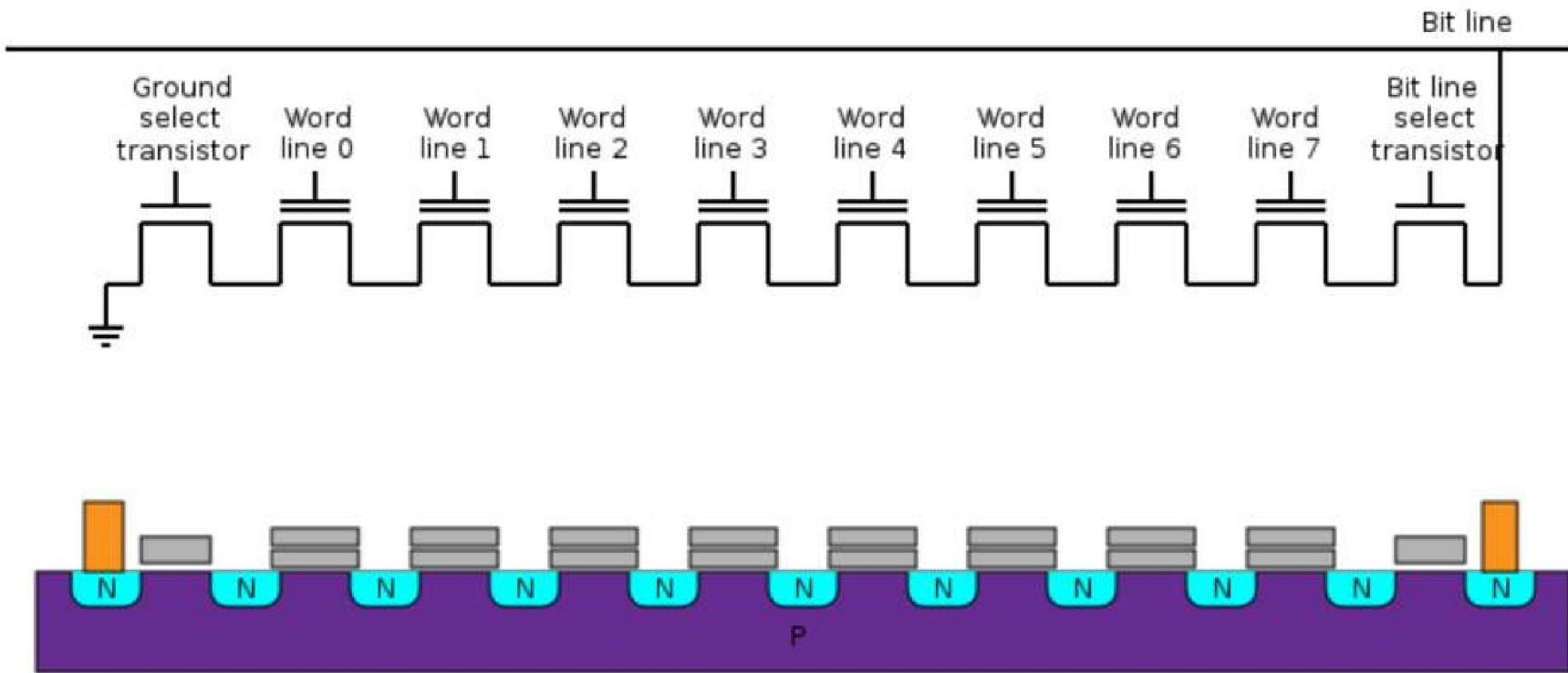
Disadvantages:

- high cost of capacity;
- limited amount of writes;
- faster, no moving parts, but, have **endurance time**: 1-2 years (except Intel's 3D Point);

Controller:

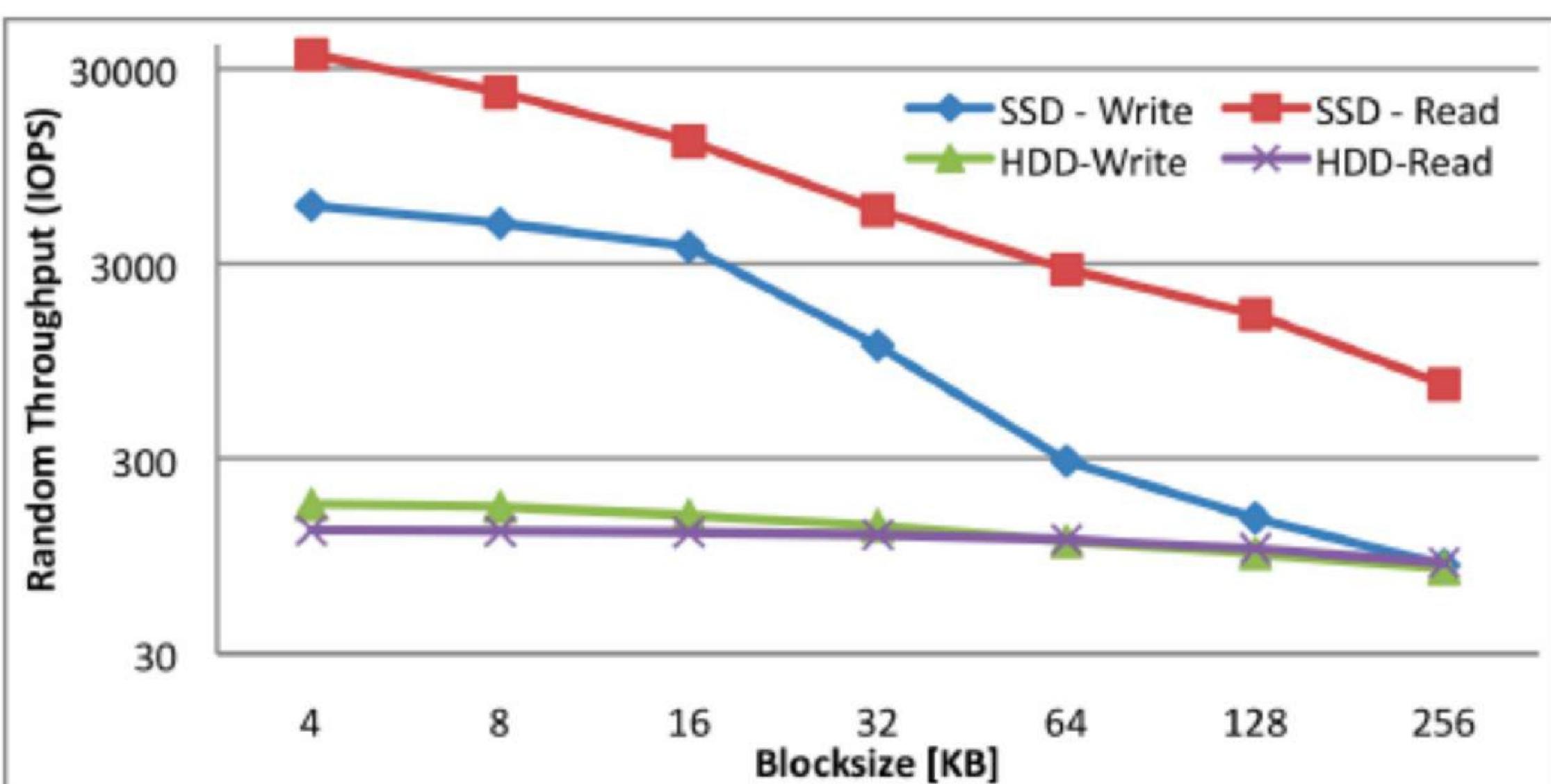
- R/W cache
- bad block management
- wear leveling

6



Advantages:

- fast random read;
- no moving parts;
- no spin-up time;
- less energy consuming;
- silent.

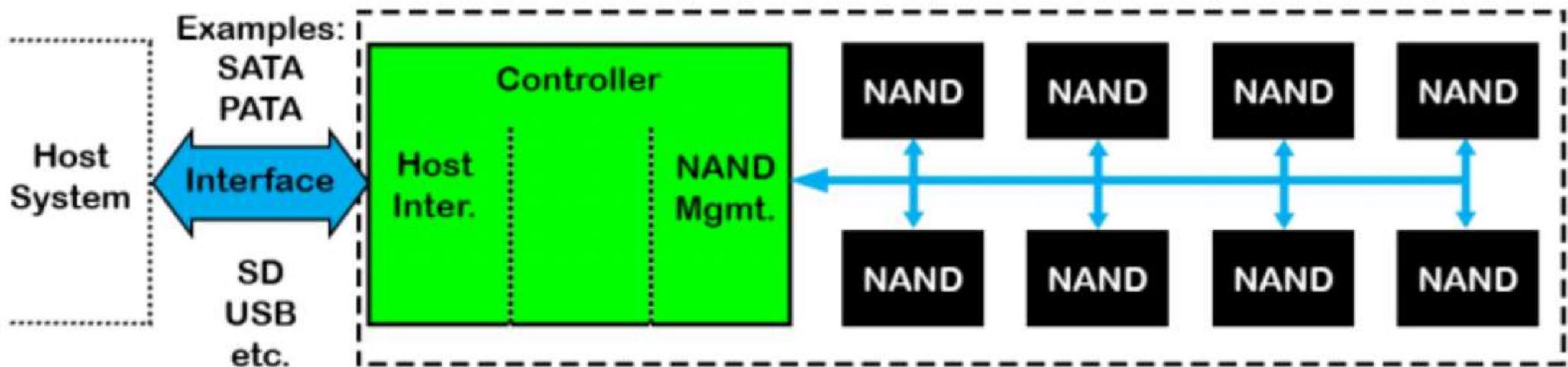


Disadvantages:

- high cost of capacity;
- limited amount of writes;
- faster, no moving parts, but, have **endurance time**: 1-2 years (except Intel's 3D Point);
- expensive;
- sophisticated controller.

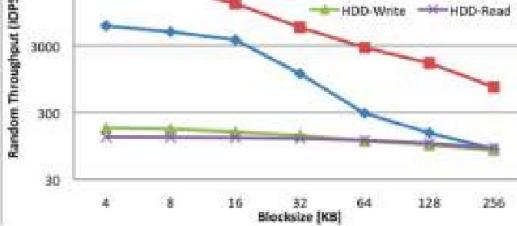
Controller:

- R/W cache + battery;
- bad block mapping;
- wear leveling;

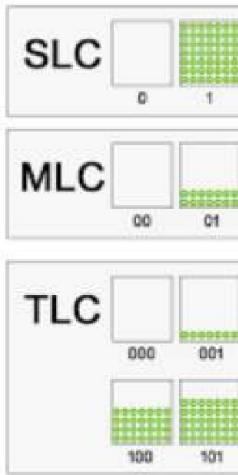


Basic Solid State Drive (SSD) Architecture

- sophisticated controller.

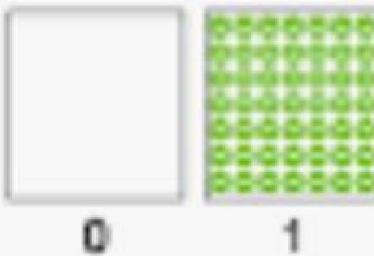


- Information is stored in **floating gate transistors**, i.e. read and written, holds charge for a long period.
- **Levels** of data storage:
 - single level cell: **1 bit/cell**, fast, endurable (50000-100000 erases), small;
 - 2-level cell: **2 bits/cell**, slower, less endurable (1000-10000 erases), medium;
 - 4-level cell: **4 bits/cell**, slower, less endurable, large.
- **Capacities:** 32 Gbyte - 4 Tbyte.
- Combined HDD + SSD devices.



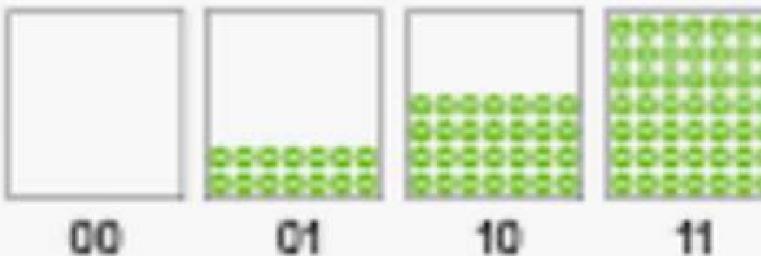
[4] <https://en.wikipedia.org>

SLC



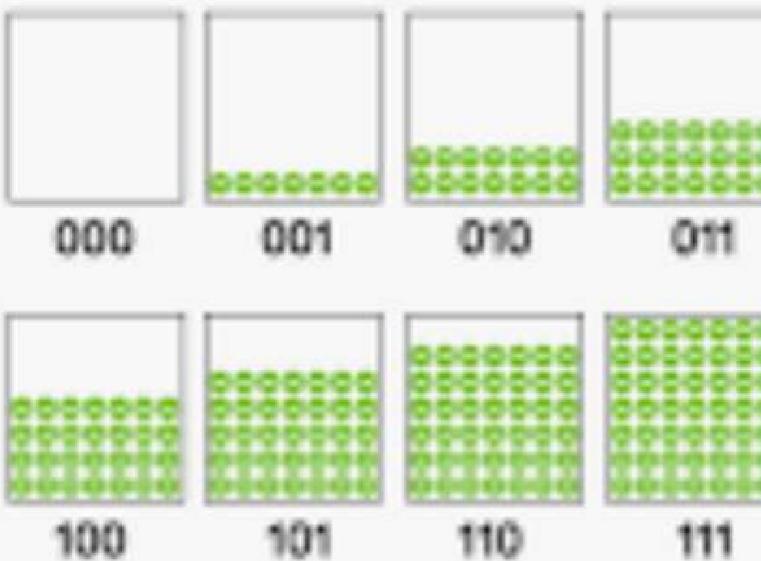
= 1bit

MLC



= 2bit

TLC



= 3bit

Storage structures

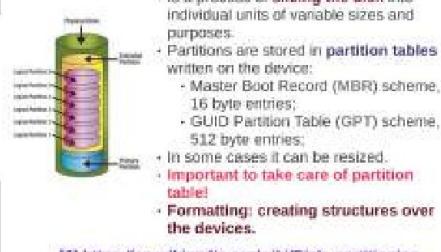
compound storage devices built upon simple storage media

Devices

- Storage components appear as **devices** on operating system level.
- Corresponding drivers (= interface between the kernel and the devices) take care of them.
- Linux kernel: **/dev** files
- major numbers: kernel, SCSI, sd, 8, 65, etc.
- minor numbers: used by the driver, for sd it indicates disk and partition (0 disk, 1, part, 2 part, 16 disk, etc);

[5] <https://www.cetly.com/library/viewtime-device-drivers/229000001/1-00102.html>

Disk partitioning



- Is a process of slicing the disk into individual units of variable sizes and purposes.
- Partitions are stored in **partition tables** written on the device:
 - Master Boot Record (MBR) scheme, 16 byte entries;
 - GUID Partition Table (GPT) scheme, 512 byte entries;
- In some cases it can be resized.
- **Important to take care of partition table!**
- **Formatting:** creating structures over the devices.

[6] https://en.wikipedia.org/wiki/Disk_partitioning

RAID



- Back-up & Array of independent physical Data
- What are expenses? Maintenance vs. PC costs.
- Separate access disk, Single logical disk.
- Free garment storage. Shared data storage, parity calculation.
- Different levels:
 - Redundancy
 - Error recovery.
 - Data loss tolerance.
 - RAID levels.
- Use the corresponding RAID to the purpose!
- RAID can be built with:
 - Many disks or fast.
 - Device subdvision, device hierarchy.



Practice #3: Creating RAID

http://www.rootz.it/content/Direct_Attached_Storage

```
sudo install mdadm  
mdadm --create /dev/md0 --level 1  
--raid-devices 2 /dev/sda /dev/sdb  
  
cat /proc/mdstat  
mdadm --detail /dev/md0  
ls -l /dev/md0  
cat /proc/mdstat  
mdadm --detail /dev/md0  
ls -l /dev/md0  
  
mdadm -f /dev/sda1 /dev/md0  
mdadm --remove /dev/sda1 /dev/md0  
mdadm -A /dev/sda1 /dev/md0  
mdadm -add /dev/sda1 /dev/md0  
mdadm -add /dev/sdb1 /dev/md0  
  
mdadm -stop /dev/md0  
mdadm --assemble -scan  
  
mdadm --build /dev/md0 --level=1  
--raid-devices 2 /dev/sda1 /dev/sdb1  
mdadm -stop /dev/sda1 /dev/sdb1  
mdadm -assemble -scan
```

```
mdadm -c /dev/sda1 /dev/sdb1  
mdadm --assemble -scan  
  
mdadm --create /dev/md0 --level 5  
--raid-devices 3 /dev/sda1 /dev/sdb1 /dev/sdc1  
mdadm --add /dev/sdc1 /dev/md0  
mdadm -grow /dev/sdc1 --raid-devices 4  
mdadm -detail /dev/md0  
  
mdadm -f /dev/sda1 /dev/md0  
mdadm -grow /dev/sdc1 --raid-devices 4  
mdadm -detail /dev/md0  
  
mdadm -f /dev/sda1 /dev/md0  
mdadm -c /dev/sda1 /dev/sdb1  
mdadm --create /dev/md0 --level 5  
--raid-devices 3 /dev/sda1 /dev/sdb1 /dev/sdc1  
mdadm --add /dev/sdc1 /dev/md0  
mdadm -grow /dev/sdc1 --level 6  
mdadm -detail /dev/md0  
  
mdadm /dev/sdc1 -f /dev/sdc1  
mdadm /dev/sdc1 -stop  
mdadm /dev/sdc1 -add /dev/sdc1
```

Just a Bunch of Disks (JBoD)

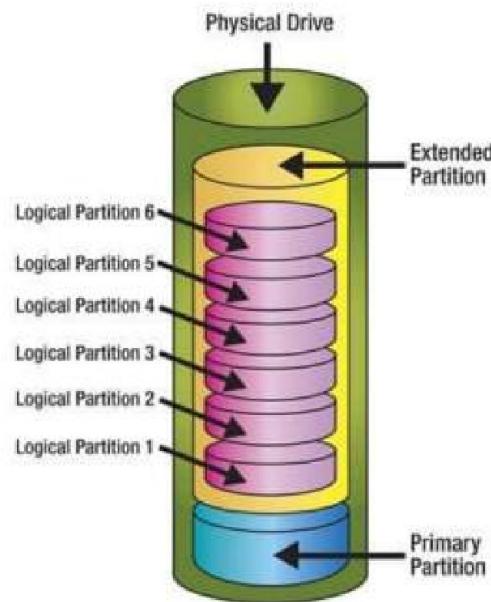
= disk expansion units

- It offers very simple disk organizing features;
- No significant added value, individual disks are connected through **intelligent expansion units**.
- These units are built up of as follows:
 - basic enclosure;
 - backplane with one or more expander chips;
 - redundant power supply;
 - proper connection, like SAS or SATA;
- Simple operations may apply:
 - concatenation = spanning;
 - mirroring;

Devices

- Storage components appear as **devices** on operating system level.
- Corresponding drivers (= interface between the kernel and the devices) take care of them.
- Linux kernel: **Idev** files
- major numbers: kernel, SCSI, **sd**, 8, 65, etc.
- minor numbers: used by the driver, for sd it indicates disk and partition (0 disk, 1, part, 2 part, 16 disk, etc);

Disk partitioning



- Is a process of **slicing the disk** into individual units of variable sizes and purposes.
- Partitions are stored in **partition tables** written on the device:
 - Master Boot Record (MBR) scheme, 16 byte entries;
 - GUID Partition Table (GPT) scheme, 512 byte entries;
- In some cases it can be resized.
- **Important to take care of partition table!**
- **Formatting: creating structures over the devices.**

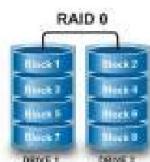
[6] https://en.wikipedia.org/wiki/Disk_partitioning

RAID

- Redundant Array of Inexpensive/Independent Disks.
- What was expensive? Mainframe vs. PC worlds.
- Spread data across disks. **Single logical disk**.
- Reengineered structure: **Striped data storage, parity calculation**.
- Designed for:
 - **speedup** devices;
 - **size** increase;
 - **error recovery**.
- Redundancy: **useful** vs. **useless**.
- RAID levels...
- **Use the corresponding RAID to the purpose!**
- SCSI RAIDs do better. Why?
 - multiple devices on bus;
 - device subdivision, device hierarchy.



RAID0:
 - Concatenation and striping disks;
Benefit: speedup.
Drawback: failure of any device causes the whole structure to fail;



RAID1:
 - Mirroring, no striping;
 - Identically written disks;
Benefits: security, data written 2 places;
Drawback: performance, slower R/W.



RAID2:
 - It has historical relevance;
 - Bit-level striping and parity;
 - Parity is written on dedicated drive;



RAID3:
 - Has no practical relevance.
 - Block-level striping and parity;
 - Byte-level striping and parity;
 - Dedicated parity disk;



RAID4:
 - Has no practical relevance.
 - Block-level striping;
 - Distributed parity;
Benefits: fast IO, secure for 1 device failure;
Drawback: slow rebuild time;



RAID5:
 - Most widely used;
 - Block-level striping;
 - Distributed parity;
Benefits: fast IO, secure for 1 device failure;
Drawback: slow rebuild time;



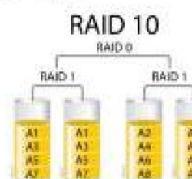
RAID6:
 - Widely used;
 - Block-level striping;
 - Distributed parities (2);
Benefits: fast IO, secure for 2 devices failure;
Drawbacks: even slower rebuild time, 2 extra devices;



Hot spare!
 ≠
RAID

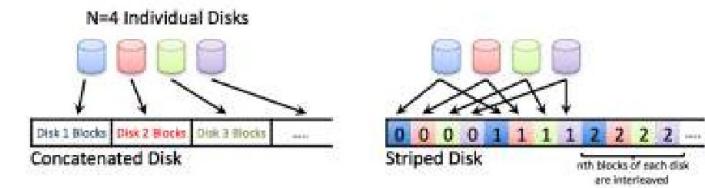


Special RAID levels:
 - RAID 1+0: mirrors first, then concatenate;
 - RAID 0+1: concatenates first, then mirrors;



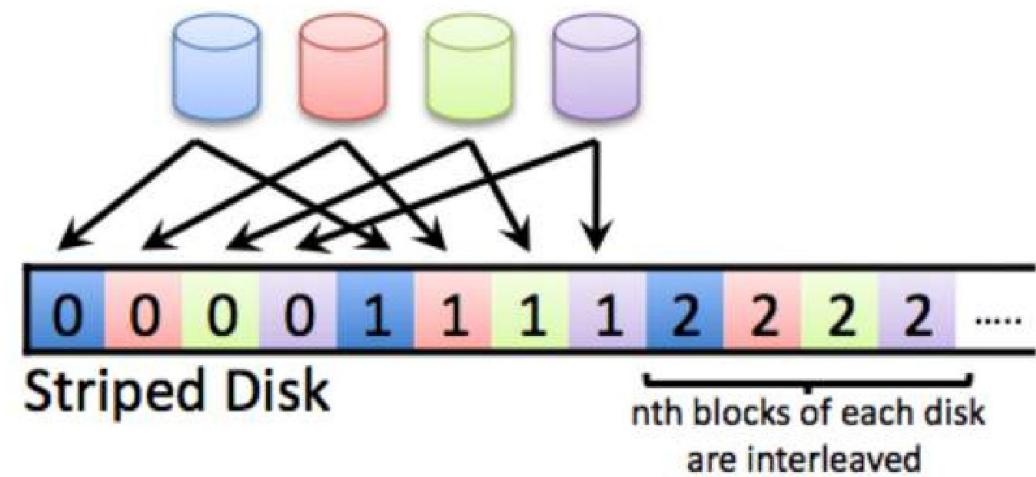
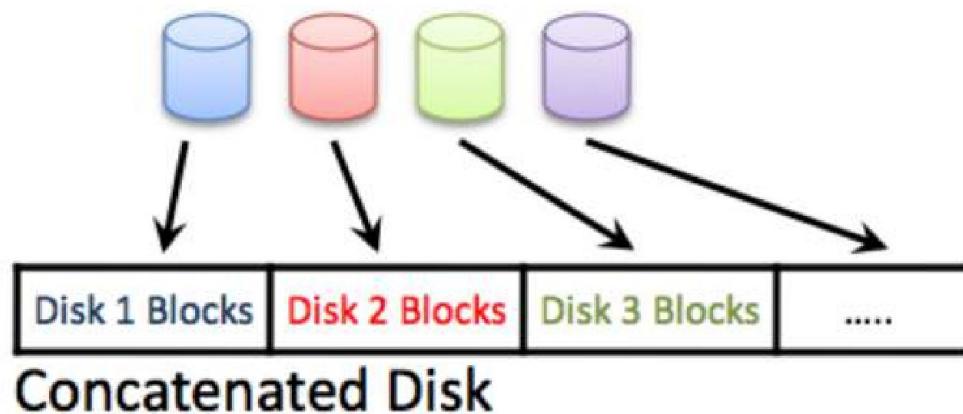
[7] https://docs.fedoraproject.org/en-US/Fedora/14/html/Storage_Administration_Guide/index.html
 [8] <https://media.techtarget.com/searchNetworking/Downloads/IncidentResponseChapter10.pdf>

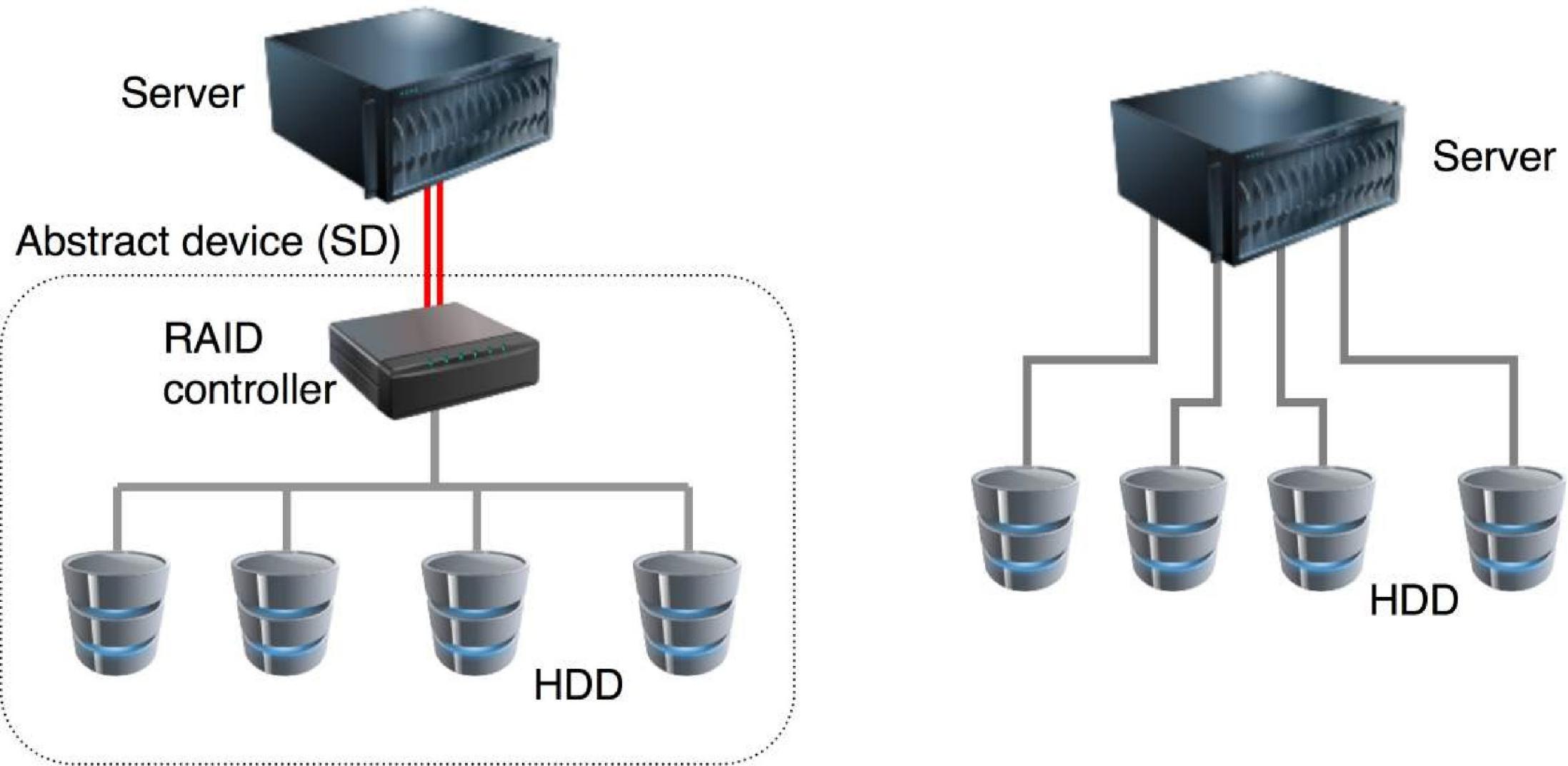
- Redundant Array of Inexpensive/Independent Disks.
- What was expensive? Mainframe vs. PC worlds.
- Spread data across disks. **Single logical disk.**
- Reengineered structure: **Striped data storage, parity calculation.**
- Designed for:
 - **speedup** devices;
 - **size** increase;
 - **error recovery**.
- Redundancy: **useful** vs. **useless**.
- RAID levels...
- **Use the corresponding RAID to the purpose!**
- SCSI RAIDs do better. Why?
 - multiple devices on bus;
 - device subdivision, device hierarchy.



Hot spare
≠
RAID

N=4 Individual Disks





RAID0:

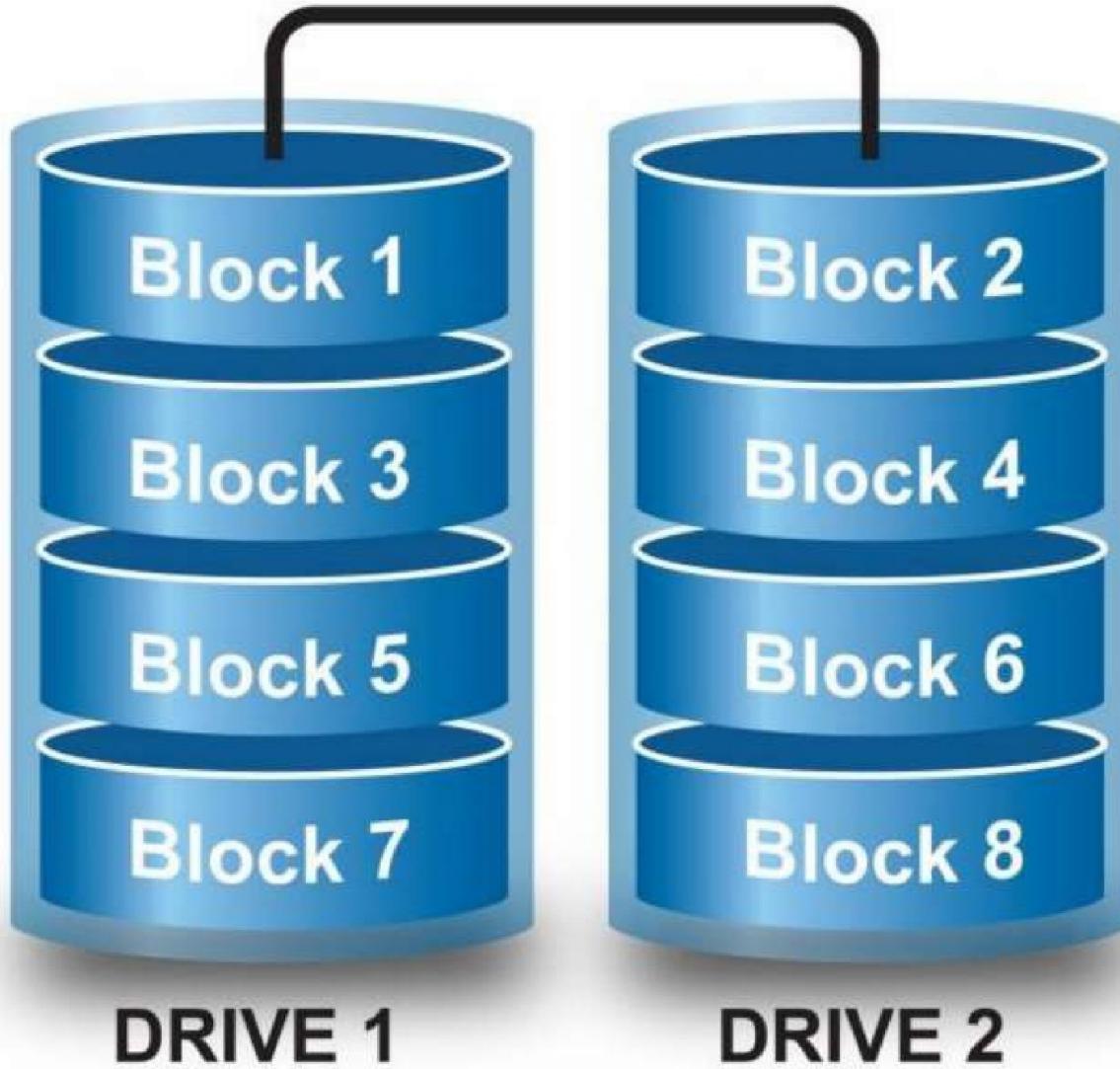
- Concatenation and striping disks;
- **Benefit:** speedup.
- **Drawback:** failure of any device causes the whole structure to fail;

RAID1:

- Mirr
- Iden
- **Ben**
- writt
- **Draw**
- slow

RAID 0

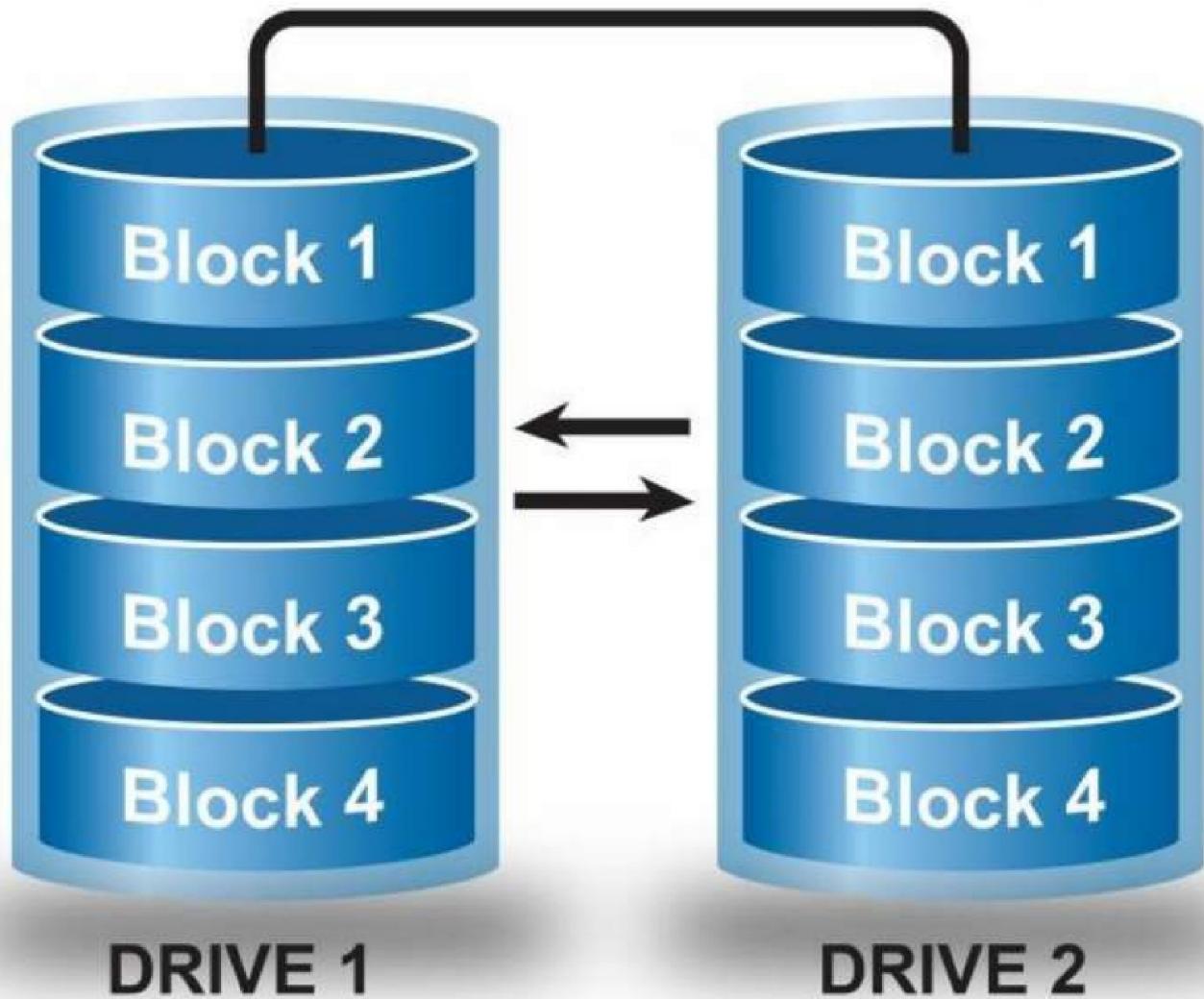
RAID 0



RAID1:

- Mirroring, no striping;
- Identically written disks;
- **Benefits:** security, data written 2 places;
- **Drawback:** performance, slower R/W;

RAID 1



Mirrored Data to both Drives

[7] <http://www.123RF.com>

striping;
en disks;
rity, data
;
formance,

RAID2:

- It has historical relevance;
- Bit-level striping and parity;
- Parity is written on dedicated drive;

RAID3:

- Has practical relevance;
- Byte striping with parity;
- Dedicated parity disk;

RAID3:

- Has no practical relevance.
- Byte-level striping and parity;
- Dedicated parity disk;

RAID4:

- Has no relevance
- Block-level striping and parity;
- Dedicated parity disk;
- Benefits
- Drawbacks

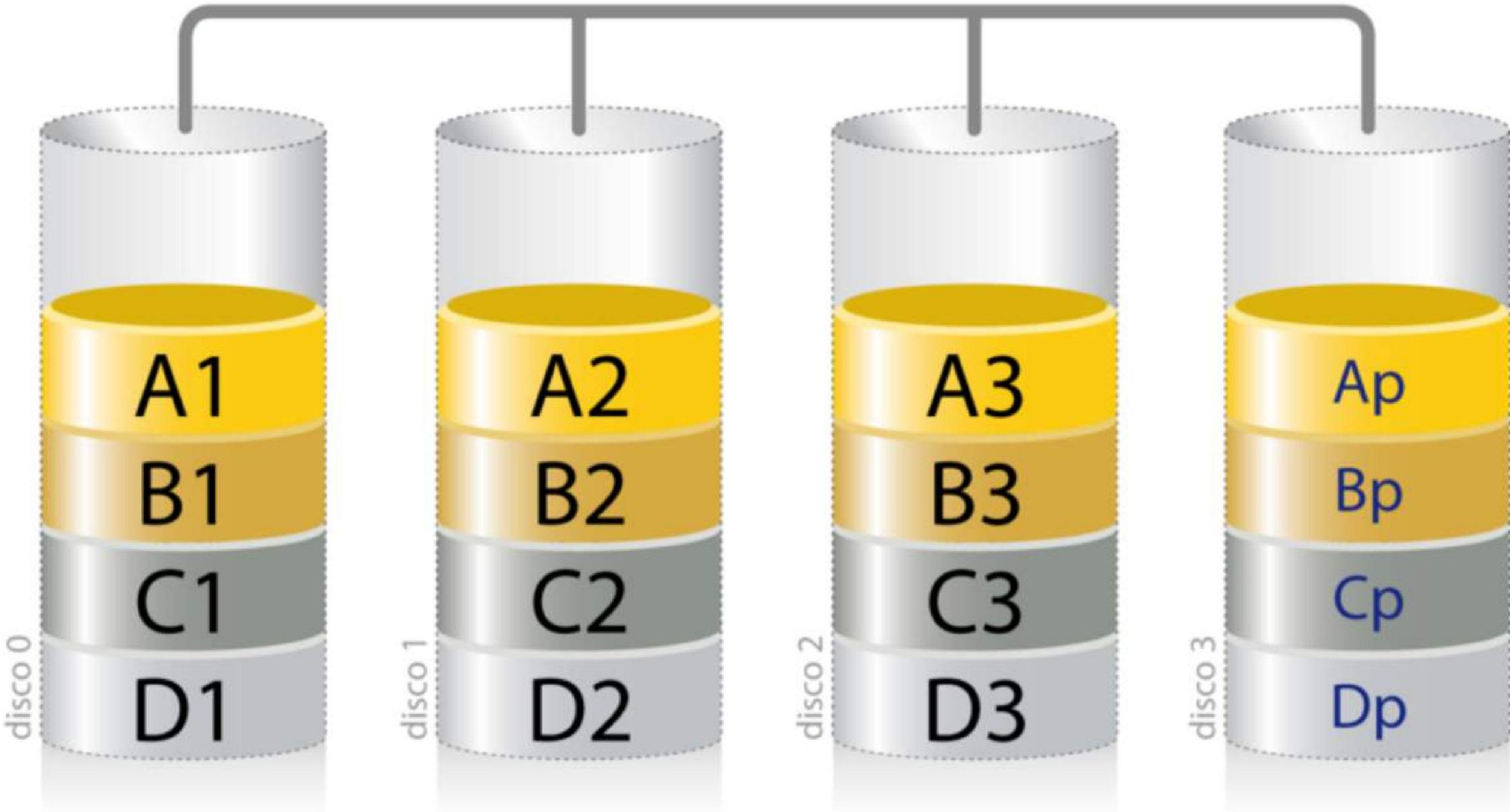
RAID4:

- Has no practical relevance.
- Block-level **striping** and **parity**;
- **Dedicated parity disk**;
- Benefits: fast IO;
- Drawback: extra parity disk;

RAID5:

- Most
- Block
- Distr
- **Bene**
secu
failure
- **Draw**
rebuil

RAID 4



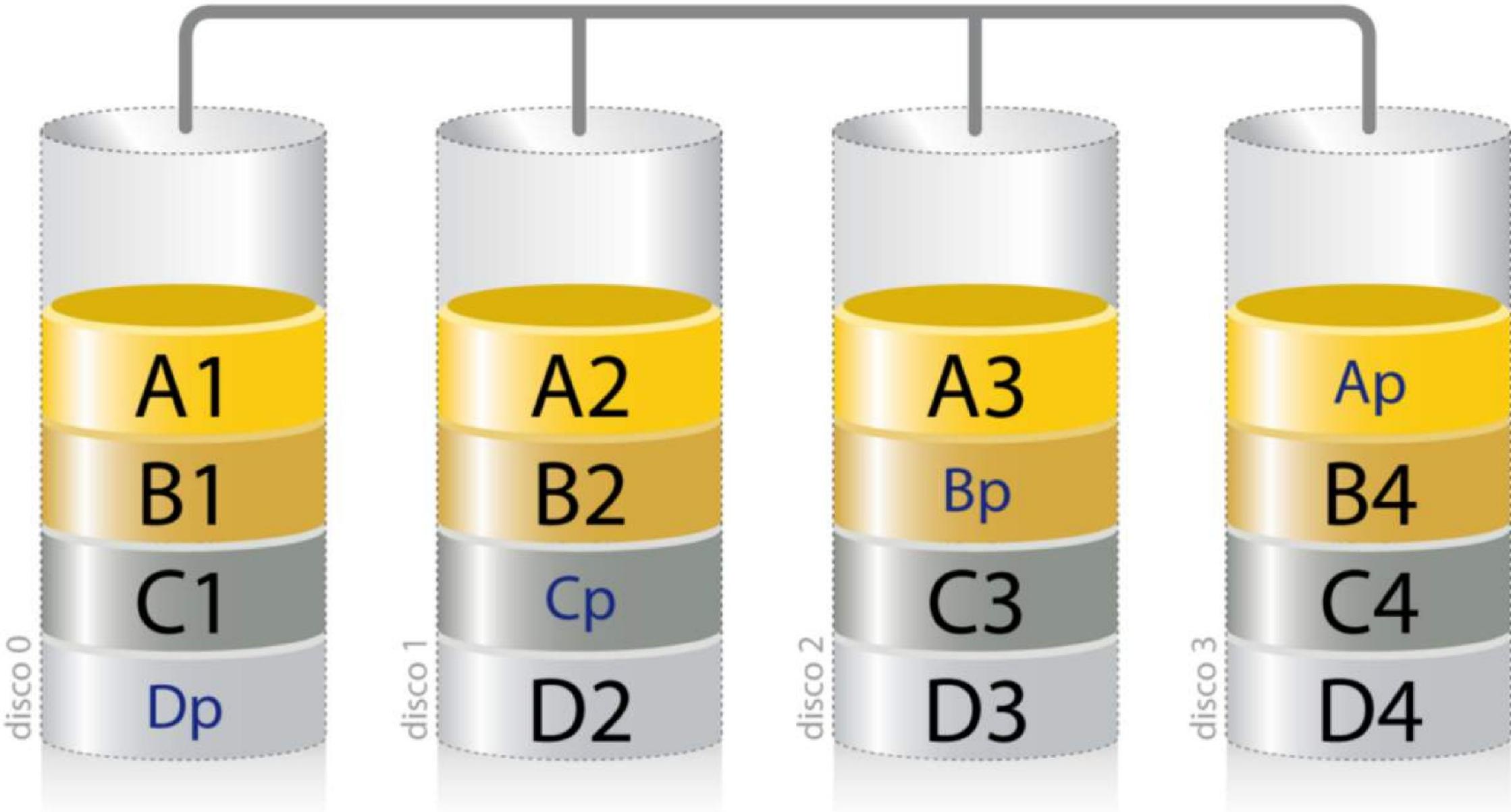
RAID5:

- Most widely used;
- Block-level striping;
- Distributed parity;
- **Benefits:** fast IO, secure for 1 device failure;
- **Drawback:** slow rebuild time;

RAID6:

- Wide
- Blo
- Dist
- **Ben**
secu
failu
• **Draw**
slow
extra

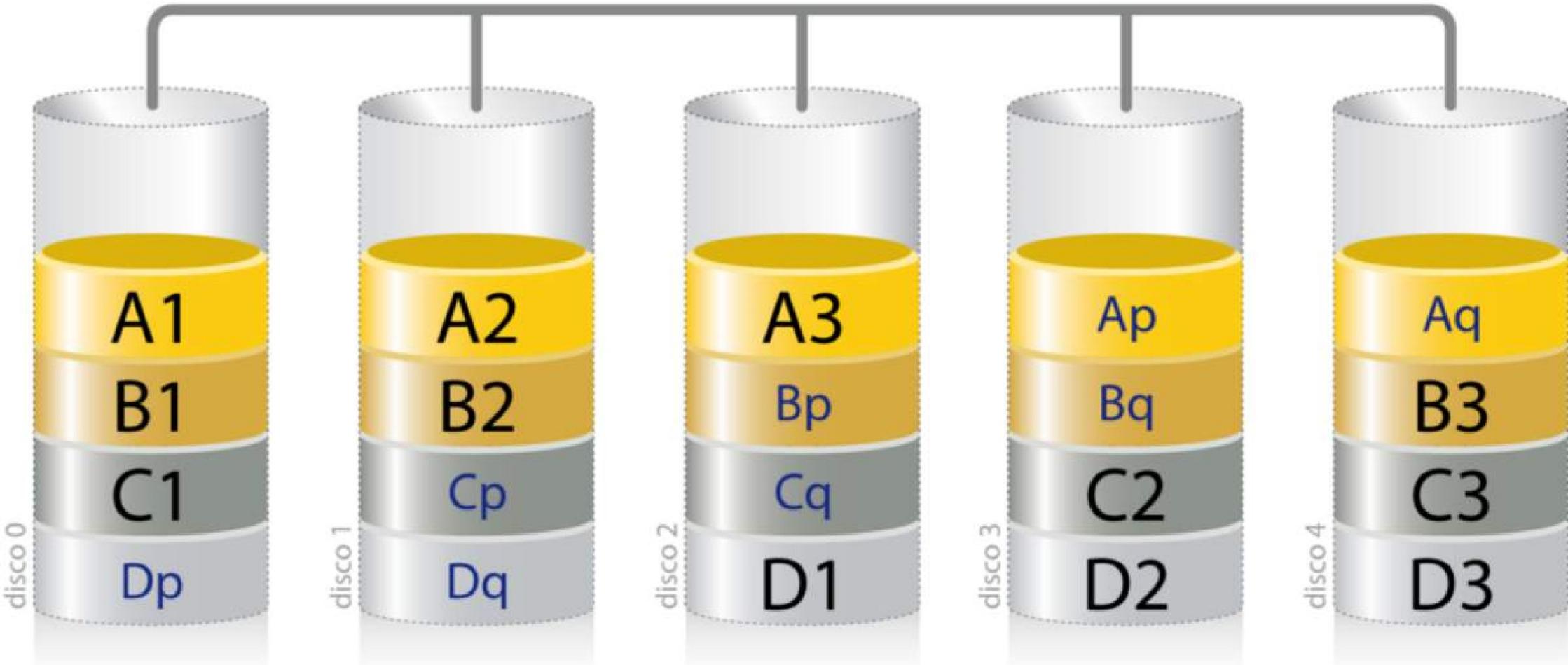
RAID 5



RAID6:

- Widely used;
- Block-level striping;
- Distributed parities (2);
- **Benefits:** fast IO,
secure for 2 devices
failure;
- **Drawbacks:** even
slower rebuild time, 2
extra devices;

RAID 6



Special RAID levels:

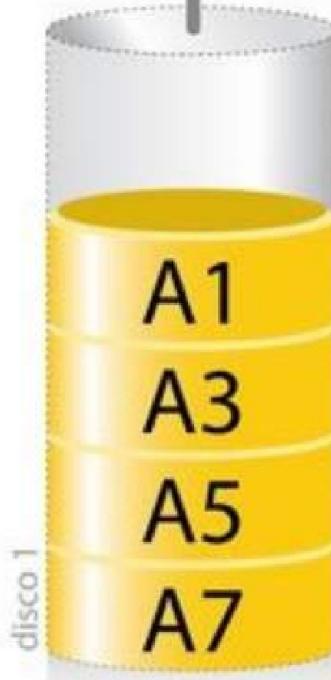
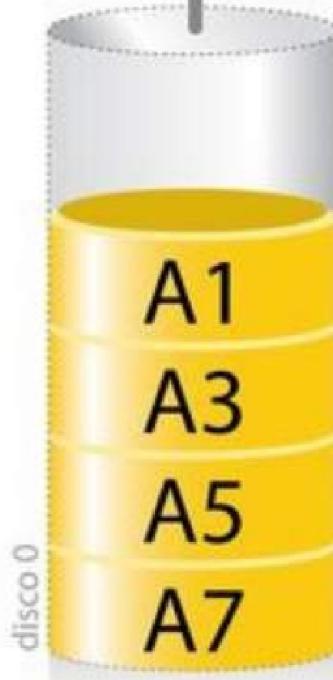
- RAID 1+0: mirrors first, then concatenate;
- RAID 0+1: concatenates first, then mirrors;

RAID 10

RAID 0

RAID 1

RAID 1



Hot spare!



RAID

Where configuration is stored?

- **configuration files** (/etc/m dadm.conf);
- **superblocks**
 - 256 Byte;
 - stores metadata (RAID level, e.g.)

https://raid.wiki.kernel.org/index.php/RAID_superblock_for

Storage architectures

the way we construct network of elementary storage devices and storage structures

Direct Attached Storage (DAS)

- Not a network!
- Storage device directly connected to and exclusively used by the host;
- Most of the protocols are used (SCSI, SAS, ATA, SATA, FC, USB, etc.);
- It is not necessary to be a small device! E.g. SUN 10000 + DAS;
- External or internal;
- Smallest latency, closest to server, but limited complexity.

↓ SAN



Storage Area Networks (SAN)

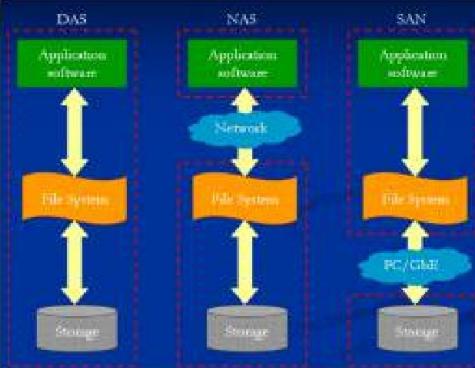
a network that provides access to block level devices



Network Attached Storage (NAS)

file level data storage

- latency larger in 1000x with the system of storage devices, NFS, protocol
- Network appliances with similar internal structure as SANs (SAN protocols)
- Client talks to it as file servers
- Allow file systems access (e.g. SMB/CIFS, Network File, AFP, etc.)
- Requests range from direct read/write
- CPU, memory and I/O bandwidth, power supply, connectors, redundancy, etc.
- Standard version: distributed data and metastable storage on multiple servers
- extra processing power can be added,
- storage failures will not affect file system.
- Can be local mounted on workstation



Unified Storage

offers iSCSI, NFS, CIFS in the same box



Practice#4: Take a look at SAN devices

show enclosure-status -type all
show hardware-information
show front-end-type ctrl
show disks
show raid-groups
show esp-mode
show volumes
set volume volume-name DEV0101
show fc-parameters
show fc-configuration
show host-mac-names
show host-ip-addresses
test net ping port 10 tg 192.168.2.10 -count 2

Clear study: measuring
storage D3000
The controller... The results...

show network
show ems
show storage-system-name
show raid-tuning
show cache-parameters
show extreme-cache

More detailed information about the D3000
can be found here: <http://www.ictnews.de/contimous/universal-unified-storage/>
<http://www.ictnews.de/contimous/universal-unified-storage/>

Direct Attached Storage (DAS)

- Not a network!
- Storage device **directly connected to** and **exclusively used by** the host;
- Most of the protocols are used (SCSI, SAS, ATA, SATA, FC, USB, etc.);
- It is not necessary to be a small device! E.g. SUN 10000 + DAS;
- External or internal;
- Smallest latency, closest to server, but limited complexity.

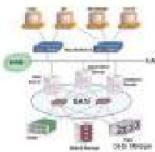
↓
SAN



Storage Area Networks (SAN)

a network that provides access to block level devices

- Basic purpose: **connects servers to disk and tape devices**;
 - OS senses devices as if they were directly attached;
 - Offers **only block level** operations, no higher structures;
 - From 2000s, developed from the mainframe world where multiple servers were connected with storage devices, eliminating Single Point of Failure;
 - **Separate network** from LAN: own topology and devices;
 - Hosts do not own resources! **Conflicts**.
 - As shared device QoS is needed:
 - bandwidth;
 - latency.



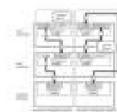
Major SAN types:

- FC-SAN;
 - Ethernet-SAN;

Logical Unit Number (LUN): individually addressable R/W device, volume

- logical disk units;
 - tapes;

Controller, Target, Disk, Slice nomenclature:



SAN variations:

- in-band: data + control;
 - out-of-band: data || control (Ethernet)

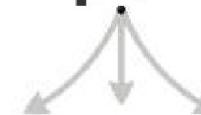


Zoning: creating logical groups of devices

- **soft** - implemented by software;
 - hard - implemented by hardware;
 - WWN - address based;
 - port - port based;



Components



- initiator;
- interface to servers;
- run firmware or software
- make use of storage
- Services.



+

Management
software
g.o.

- active networking layer, such as switches, routers, gateway devices and cables;
- switch:
 - gives dedicated port-to-port connection among the devices;
 - non-blocking;
- copper vs optical cables:

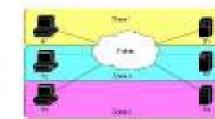
Storage:

- target;
- disks, tapes;
- each device has a unique identifier: WWN;
- Logical Unit Number (LUN) = volumes;
- LUN masking and zoning;
- disk controller, disk array, disk enclosure
- redundant architecture;



NO SPOTS!

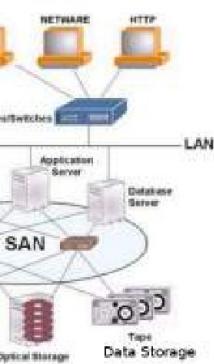
* Redundancy



<https://arkit.co.in/san-switch-basic/>

https://en.wikipedia.org/wiki/Storage_area_network

<https://slideplayer.com/slide/1517247/>

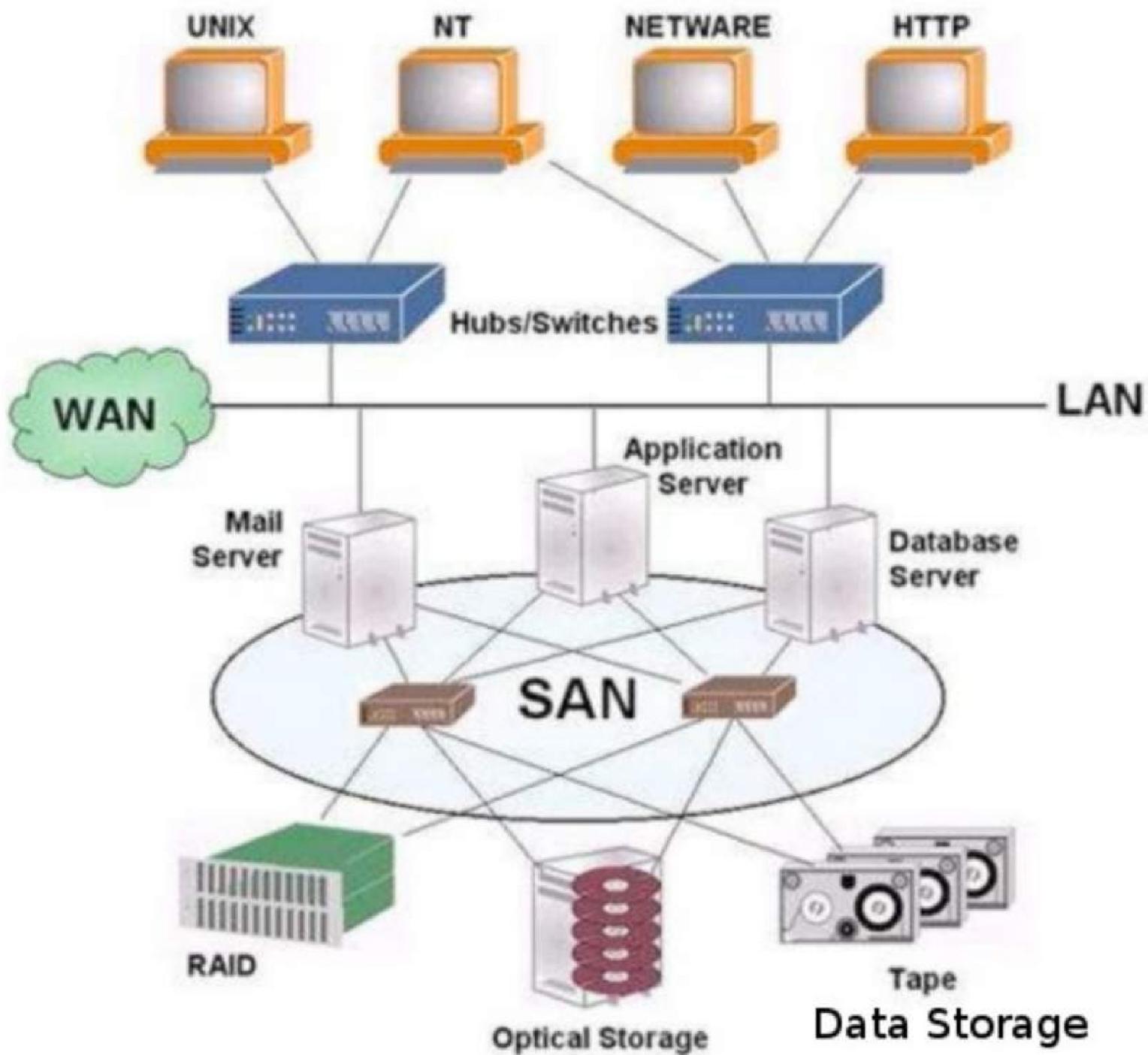


- Basic purpose: **connects servers to disk and tape devices**;
 - OS senses devices as if they were directly attached;
 - Offers **only block level** operations, no higher structures;
 - From 2000s, developed from the mainframe world where multiple servers were connected with storage devices, eliminating Single Point of Failure;
 - **Separate network** from LAN: own topology and devices;
 - Hosts do not own resources! **Conflicts**.
 - As shared device QoS is needed:
 - bandwidth;
 - latency.
- Major SAN types:**
- FC-SAN;

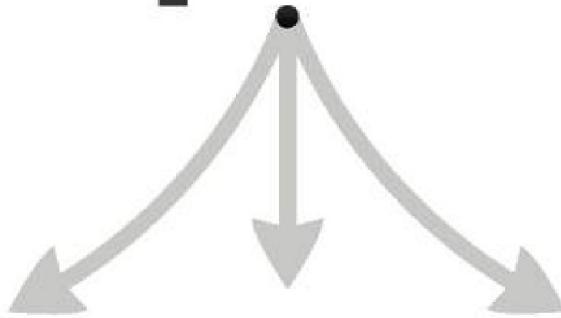
SAN variations

- in-band: c
- out-of-band

(Ethernet)



Components



Host Bus Adapter (HBA):

- initiator;
- interface to servers;
- run firmware or software to make use of storage devices;



+
Management
software
e.g.
SMI

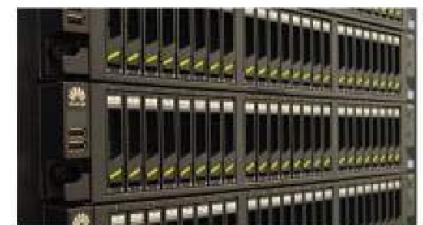
Fabric:

- active networking layer, such as switches, routers, gateway devices and cables;
- **switch:**
 - gives dedicated port-to-port connection among the devices;
 - non-blocking;
- **copper** vs. **optical** cables;



Storage:

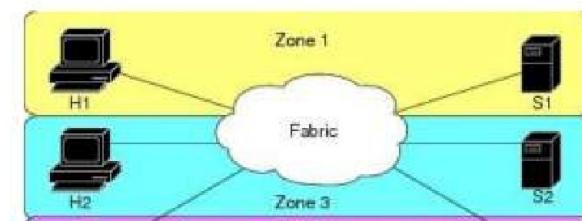
- target;
- disks, tapes;
- each device has a unique identifier: WWN;
- Logical Unit Number (LUN) = volumes;
- LUN masking and zoning;
- disk controller, **disk array**, **disk enclosure**
 - **redundant** architecture;



NO SPOF!
+
hotswap

- Chassis;
- Backplane;
- Connectors;
- Power supply;

- HBAs;
- fabrics;
- disks + structures;



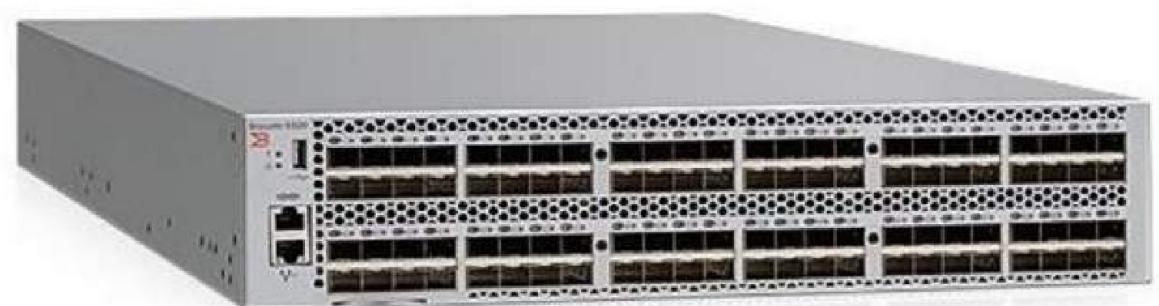
Host Bus Adapter (HBA):

- initiator;
- interface to servers;
- run firmware or software to make use of storage devices;



Fabric:

- active networking layer, such as switches, routers, gateway devices and cables;
- **switch:**
 - gives dedicated port-to-port connection among the devices;
 - non-blocking;
- **copper** vs.
optical cables;



Storage:

- target;
- disks, tapes;
- each device has a unique identifier: WWN;
- Logical Unit Number (LUN) = volumes;
- LUN masking and zoning;
- disk controller, **disk array**, **disk enclosure**
 - **redundant** architecture;



s,

on



NO SPOTS!

- Chassis;
- Backplane;
- Connectors;
- Power supply;
- Cache: RAM + NV;
- Enclosures;
- Power supply;

NO SPOF!

+

hotswap

/ and

e.g.
SMI

Copper vs
optical ca

SAN variations:

- in-band: data + control;
- out-of-band: data || control
(Ethernet)

Intelligent Storage System: Front End

route data in and out

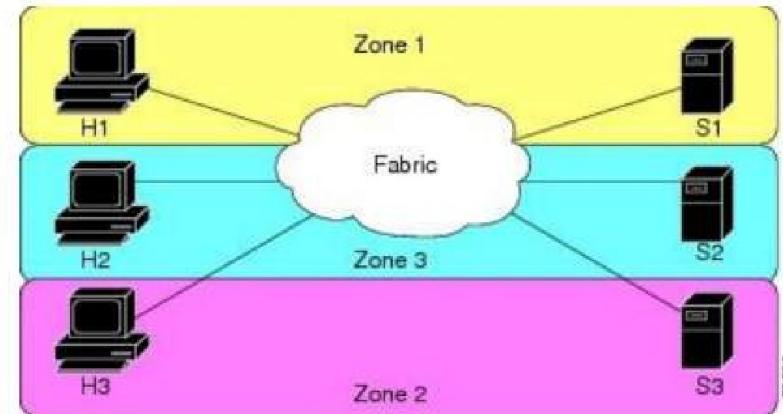
- structure functionality;
↳ Mechanism

Major SAN types:

- FC-SAN;
- Ethernet-SAN;

Zoning:

- HBAs;
 - fabrics;
 - disks + structures;
- creating logical groups of devices



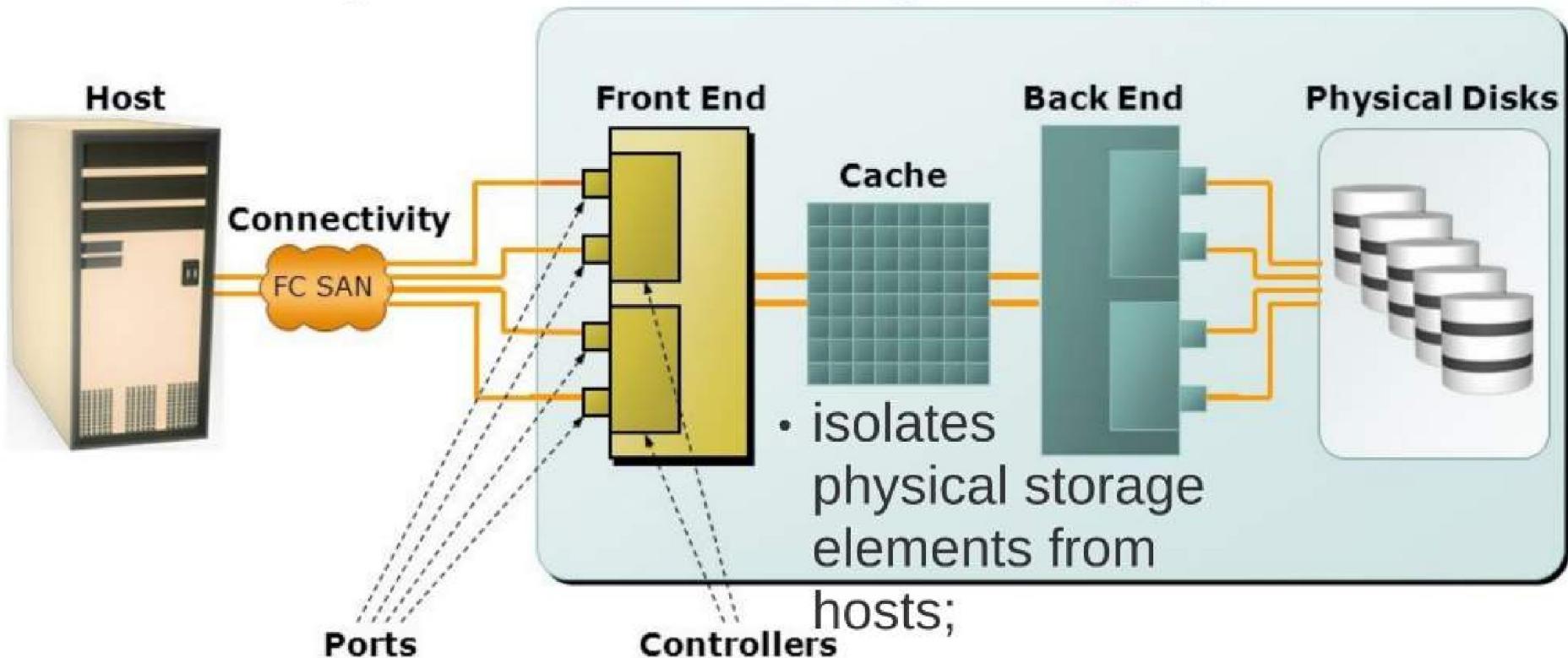
- **soft** - implemented by software;
- hard - implemented by hardware;
- WWN - address based;
- port - port based;

Intelligent Storage System: Front End

- route data in and out;
- optimize traffic;

- structure functionality;
- load balancing;

Intelligent Storage System



Network Attached Storage (NAS)

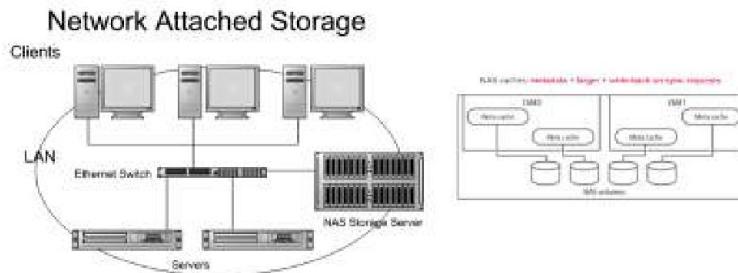
file level data storage

- History began in 1980s with file system sharing protocols: NFS, Novell;
- Network appliances with similar internal structure as SANs (disk pools, storage structures, volumes, etc.);
- Client sees it as a **file server**;
- Allow file system access: NFS, SMB/CIFS, Novell FS, SFTP, etc.;
- Requires simple IP/Ethernet network;
- Offer **some sort of redundancy**: power supply, controllers, connections, etc;
- **Clustered versions**: distributed data and metadata storage on multiple NAS'es, because
 - extra processing power can be added;
 - device failures will not disrupt the system.
- Can be rack mounted or standalone;



NAS management:

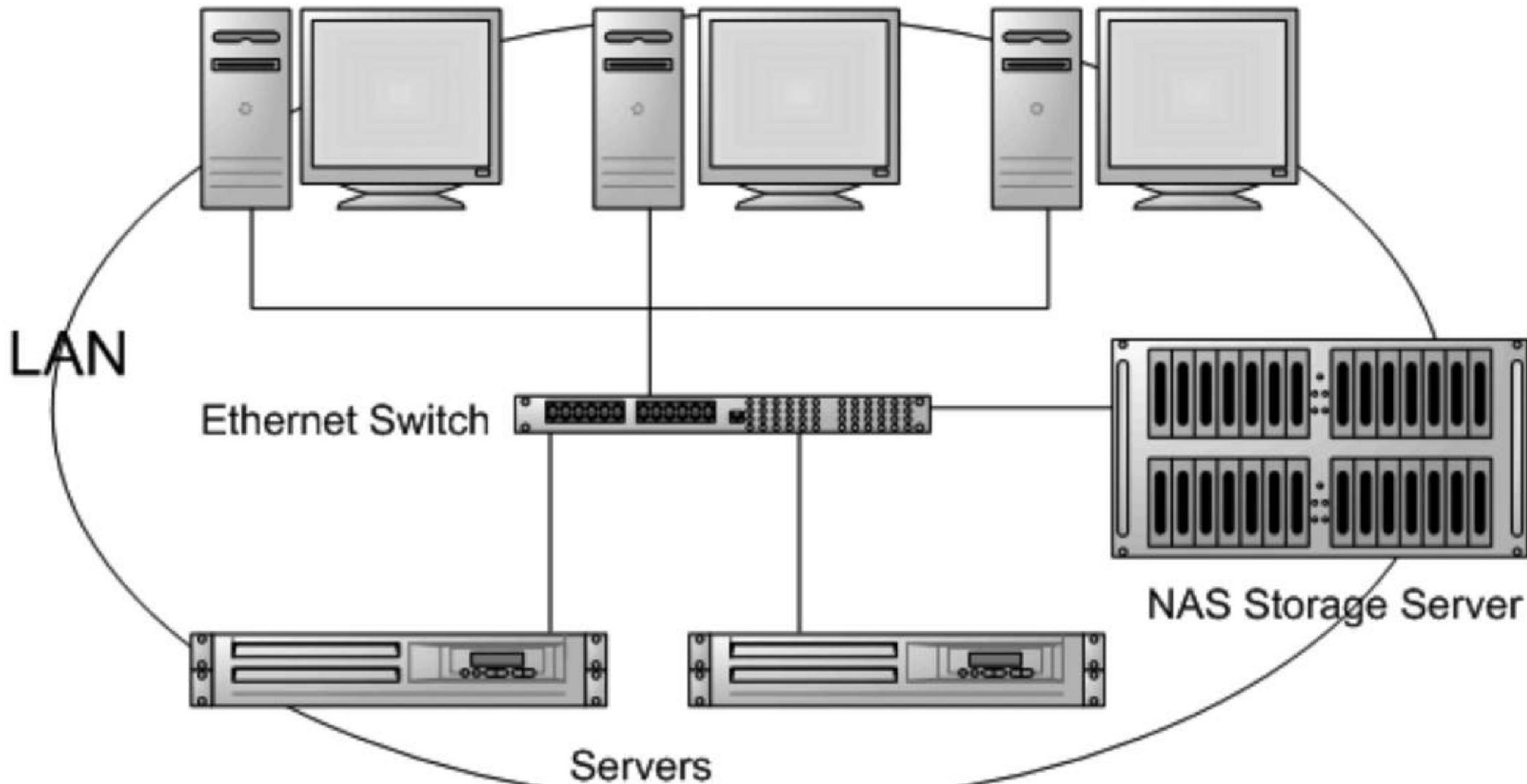
- Similar to that of SAN's: create structures, such as RAID groups, RAIDs;
- Create **NAS volumes**: user or backup (creating FS structures take longer time);
- Access control:
 - users, groups, access lists, shares, (can be taken from auth servers);
 - standards based access control (POSIX, Windows ACL);





Network Attached Storage

Clients



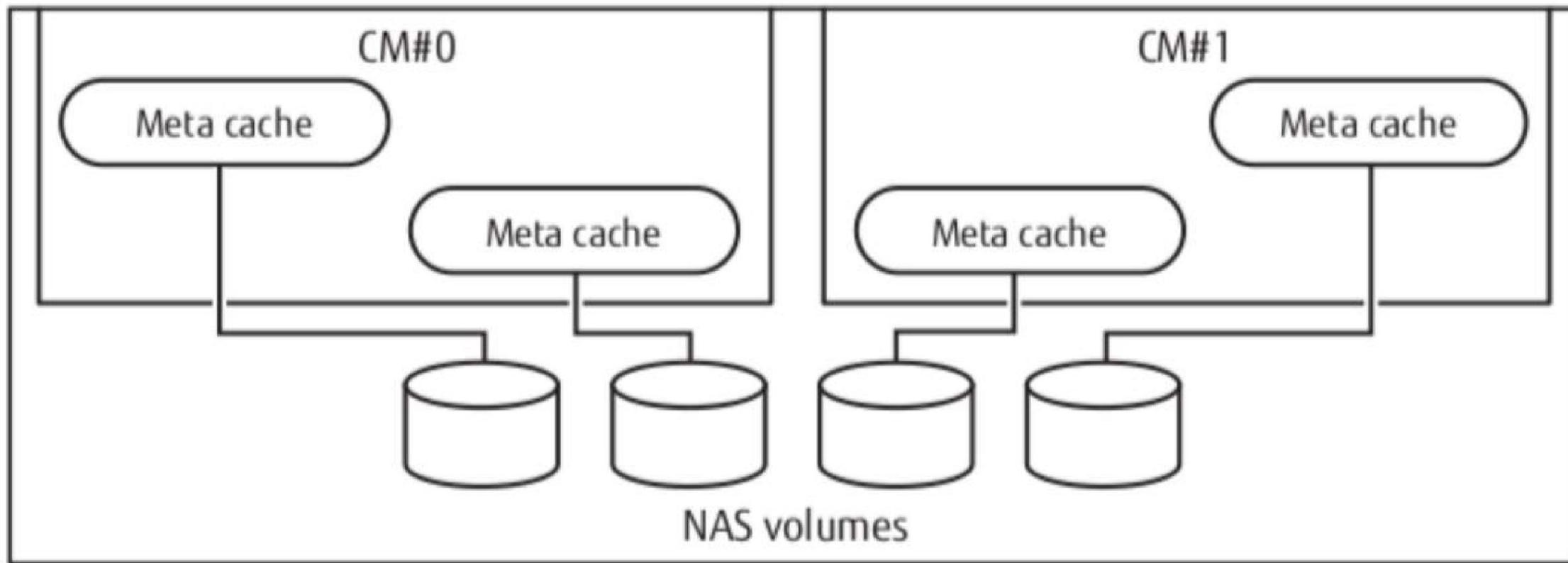
NAS management:

- Similar to that of SAN's: create structures, such as RAID groups, RAIDs;
- Create **NAS volumes**: user or backup (creating FS structures take longer time);
- Access control:
 - **users, groups, access lists, shares,** (can be taken from auth servers);
 - standards based access control (**POSIX, Windows ACL**);

NAS management:

- Similar to that of SAN's: create structures, such as RAID groups, RAIDs;
- Create **NAS volumes**: user or backup (creating FS structures take longer time);
- Access control:
 - **users, groups, access lists, shares,** (can be taken from auth servers);
 - standards based access control (**POSIX, Windows ACL**);

NAS caches: **metadata + larger + write-back on sync requests.**



Network Attached Storage (NAS)

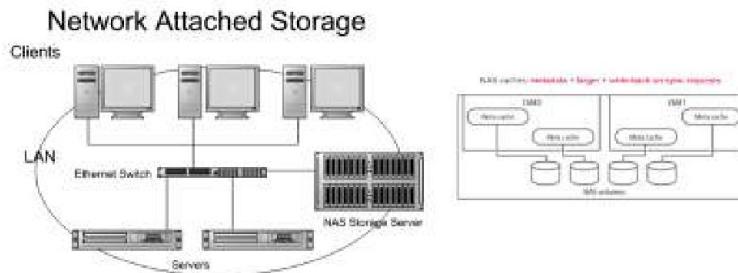
file level data storage

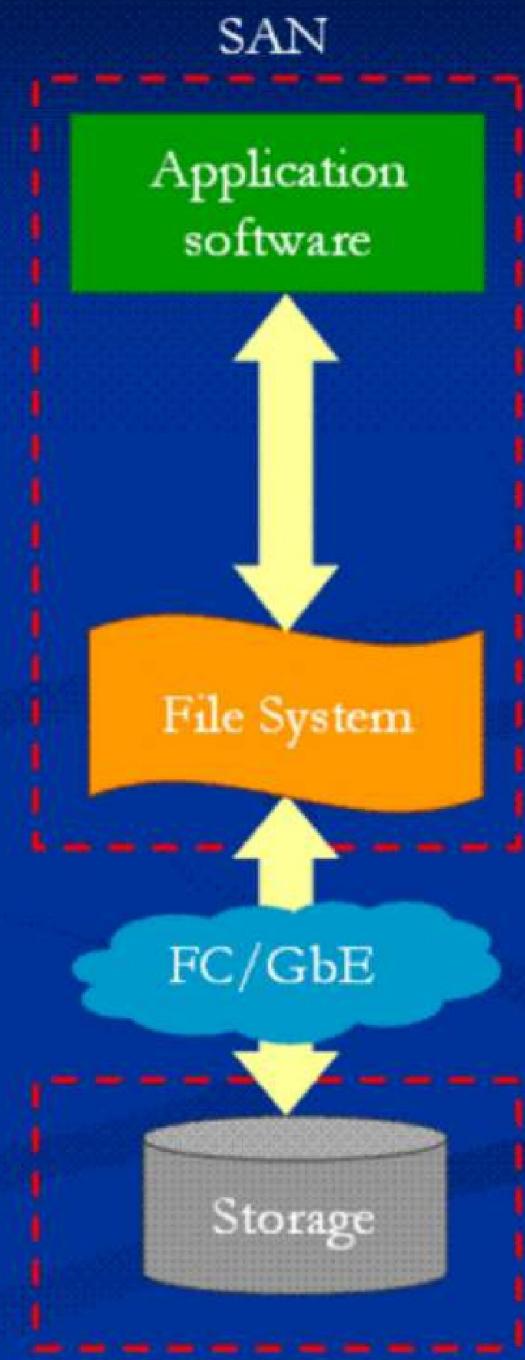
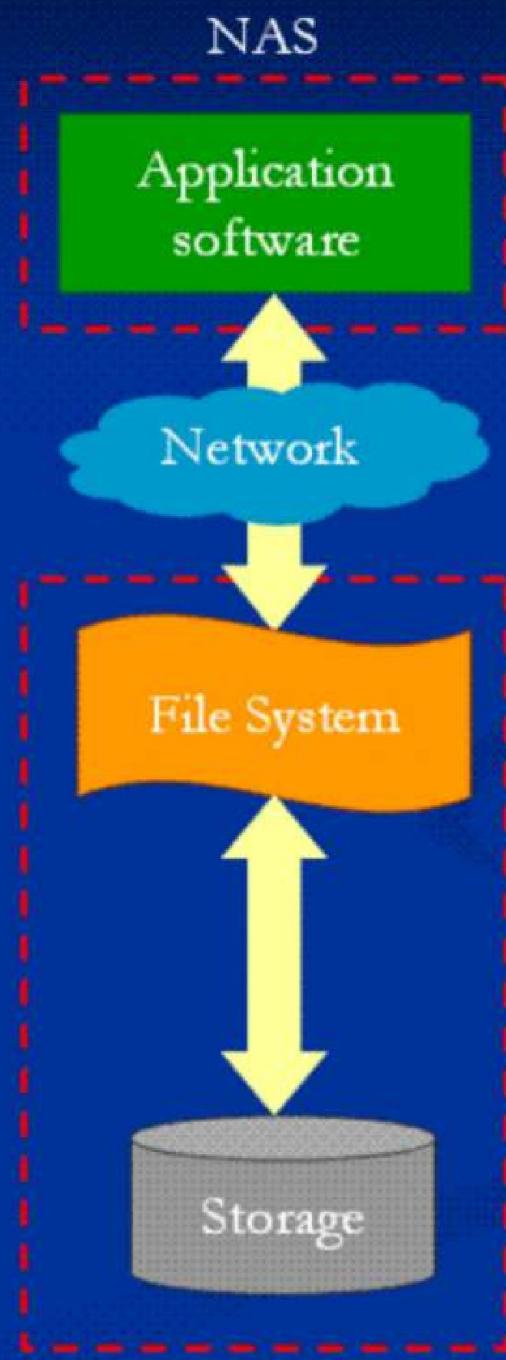
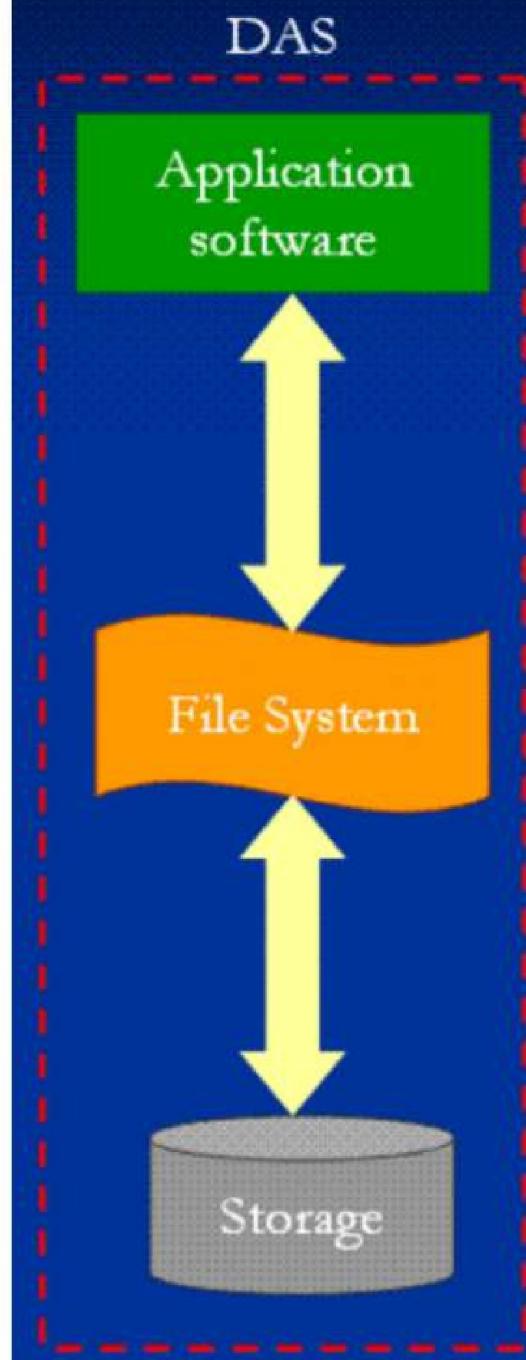
- History began in 1980s with file system sharing protocols: NFS, Novell;
- Network appliances with similar internal structure as SANs (disk pools, storage structures, volumes, etc.);
- Client sees it as a **file server**;
- Allow file system access: NFS, SMB/CIFS, Novell FS, SFTP, etc.;
- Requires simple IP/Ethernet network;
- Offer **some sort of redundancy**: power supply, controllers, connections, etc;
- **Clustered versions**: distributed data and metadata storage on multiple NAS'es, because
 - extra processing power can be added;
 - device failures will not disrupt the system.
- Can be rack mounted or standalone;



NAS management:

- Similar to that of SAN's: create structures, such as RAID groups, RAIDs;
- Create **NAS volumes**: user or backup (creating FS structures take longer time);
- Access control:
 - users, groups, access lists, shares, (can be taken from auth servers);
 - standards based access control (POSIX, Windows ACL);





Unified Storage

offers FCoE, iSCSI, NFS, CIFS in the same box

Conceptual difference:

- **block:**

- physical data concepts referring to the organization of data on disk drives, i.e. tracks, sectors;
- shorter access times, shorter data messages;

- **files:**

- logical data concepts made up of many blocks;
- more random data;
- longer access times, longer data messages;

Transformation is done by file systems!

Performance
difference!



Variations:

- SAN blocks sit on top of shared devices. Translation is done by storage controller on shared Thin Provisioning Volumes, TPVs.
- NAS gateway attached to SANs.
- Separate controllers for SAN and NAS functions.
- **Shared controllers to manage SAN and NAS:**
 - Ethernet as carrier;
 - FCoE, iSCSI, NFS, CIFS protocols;
- STaaS cloud storage, combines on-premises and cloud storage.

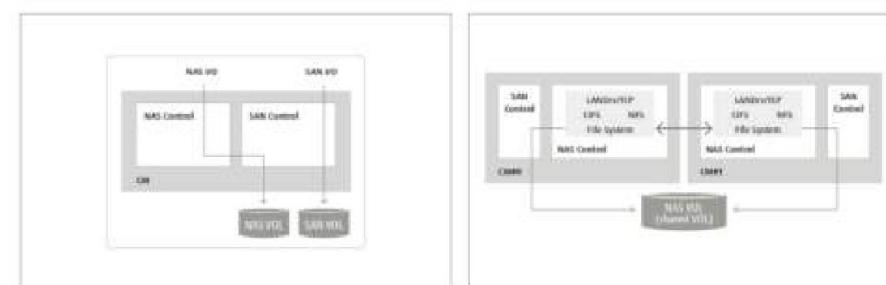


Advantages:

- Save space, device, management efforts, network costs;
- leverages utilization between NAS and SAN;
- **Conform to STaaS providers;**

Disadvantages:

- Performance penalties;



<https://www.computerweekly.com/feature/Unified-Storage-FAQ>

<http://www.ternus-dx.com/ternus-overview/unified-storage/>

Performance difference

Conceptual difference:

- **block:**

- physical data concepts referring to the organization of data on disk drives, i.e. tracks, sectors;
- shorter access times, shorter data messages;

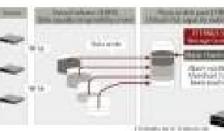
- **files:**

- logical data concepts made up of many blocks;
- more random data;
- longer access times, longer data messages;

Transformation is done by file systems!

Thin Provisioning Vol

- storage virtualisation technology;
- "Normally" 20-30% of capacity is not used;
- helps to improve storage capacity utilization;
- gives the appearance that **there is more physical space than there is**;
- dynamic allocation of resources;
- **allocates data blocks as they are written**, (thick provisioning) at initial formatting;
- intelligent mapping: e.g. if zeroes come in, no space is allocated;

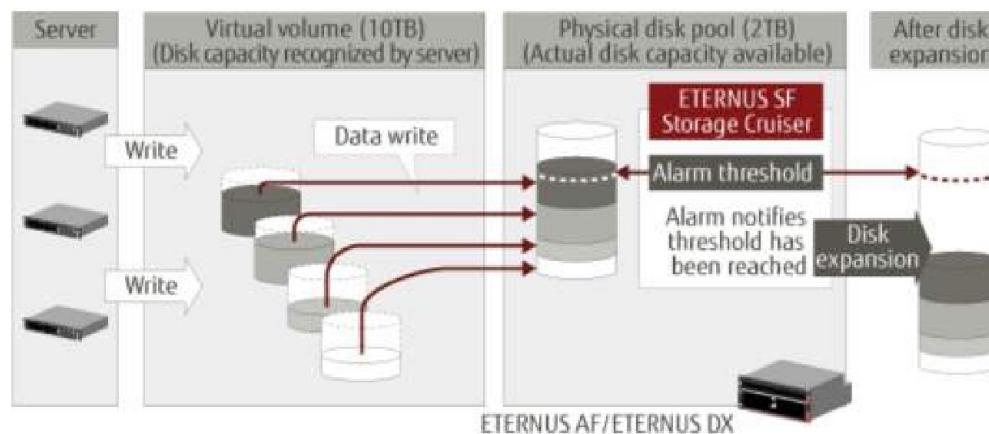


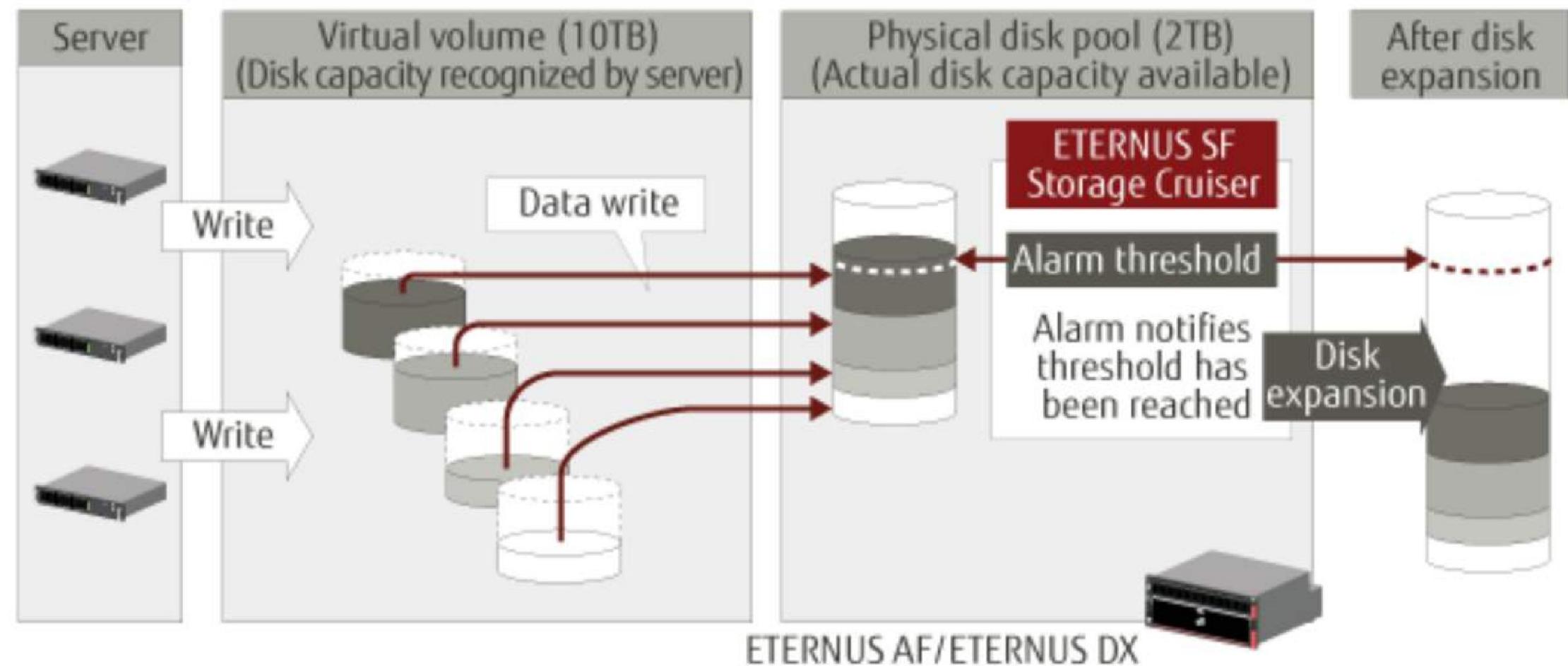
Variations:

- SAN blocks sit on top of shared devices.
Translation is done by storage controller on shared Thin Provisioning Volumes, **TPVs**.
- NAS gateway attached to SANs.
- Separate controllers for SAN and NAS functions.
- **Shared controllers to manage SAN and NAS:**
 - Ethernet as carrier;
 - FCoE, iSCSI, NFS, CIFS protocols;
 - STaaS cloud storage, combines on-premises and cloud storage.

Thin Provisioning Volumes

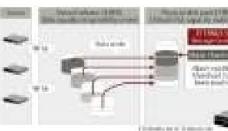
- storage virtualisation technology;
- "Normally" 20-30% of capacity is not used;
- helps to improve storage capacity utilization;
- gives the appearance that **there is more physical resource** as there is;
- dynamic allocation of resources;
- **allocates data blocks as they are written**, not in advance (i.e. thick provisioning) at initial formatting;
- intelligent mapping: e.g. if zeroes come in, no space is allocated;





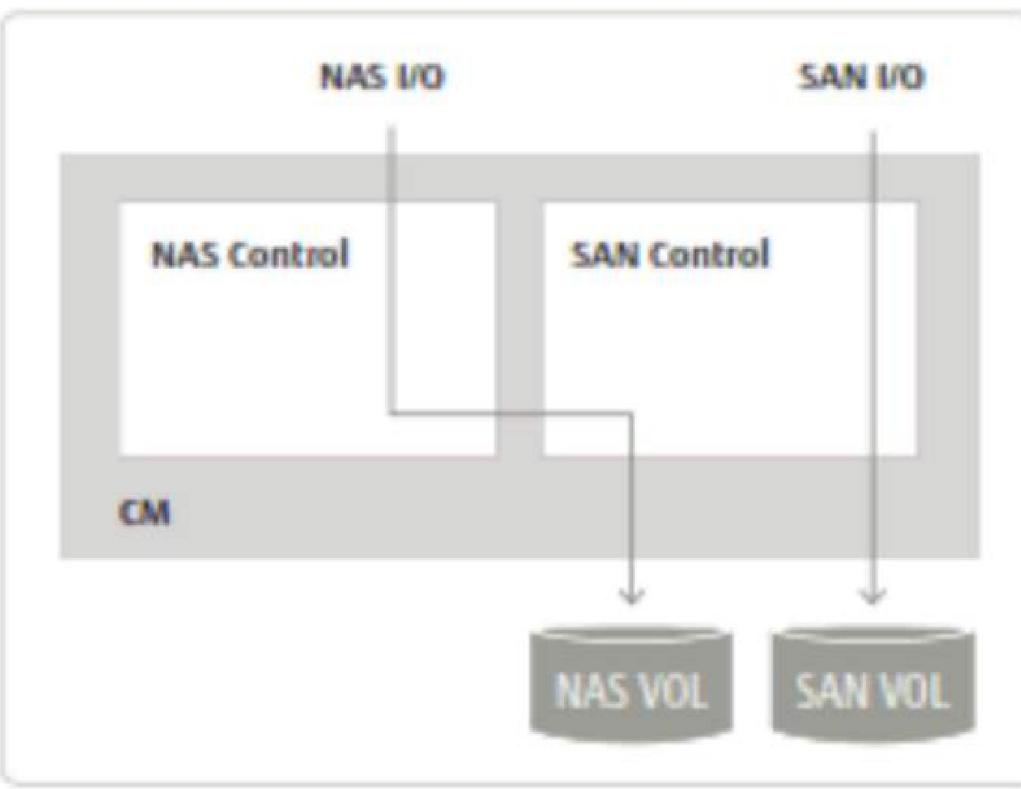
Thin Provisioning Vol

- storage virtualisation technology;
- "Normally" 20-30% of capacity is not used;
- helps to improve storage capacity utilization;
- gives the appearance that **there is more physical space than there is**;
- dynamic allocation of resources;
- **allocates data blocks as they are written**, (thick provisioning) at initial formatting;
- intelligent mapping: e.g. if zeroes come in, no space is allocated;

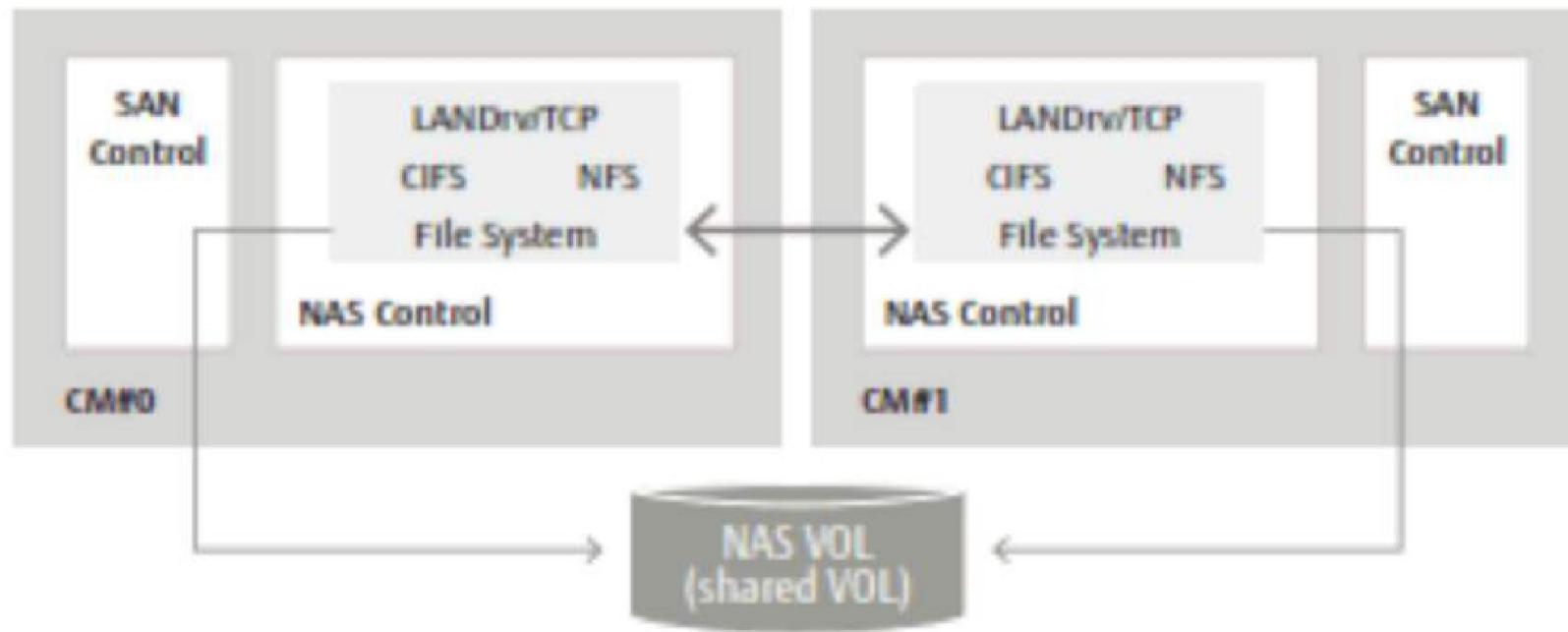


Variations:

- SAN blocks sit on top of shared devices.
Translation is done by storage controller on shared Thin Provisioning Volumes, **TPVs**.
- NAS gateway attached to SANs.
- Separate controllers for SAN and NAS functions.
- **Shared controllers to manage SAN and NAS:**
 - Ethernet as carrier;
 - FCoE, iSCSI, NFS, CIFS protocols;
 - STaaS cloud storage, combines on-premises and cloud storage.



SA
Cont
CMIO



Advantages:

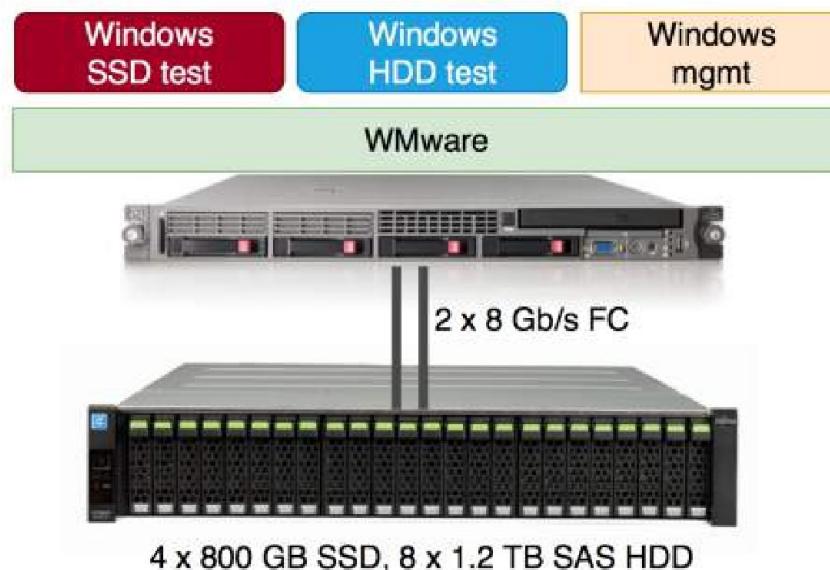
- Save space, device, management efforts, network costs;
- leverages utilization between NAS and SAN;
- **Conform to STaaS providers;**

Disadvantages:

- Performance penalties;

Case study: measuring Eternus DX200

The configuration...



The results...

Test#1:

- 2 x **SSD** in RAID1;
- 2 x **SSD** in RAID0 cache;
- 6 x **HDD** in RAID5;
- 1 x hotspare;



Test#2:

- 4 x **SSD** in RAID5;
- 8 x **HDD** in RAID6;
- no hotspare + no cache;

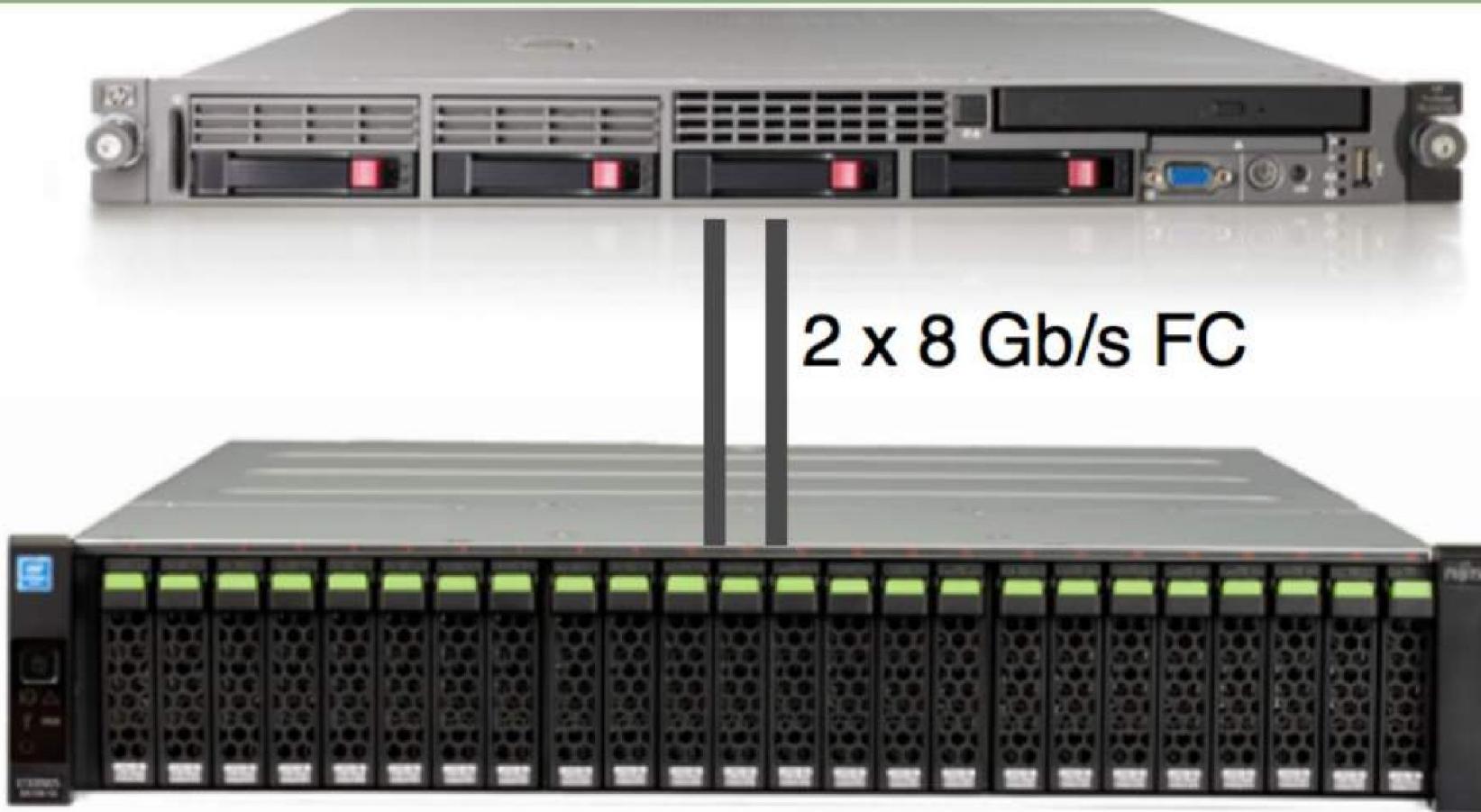


Windows
SSD test

Windows
HDD test

Windows
mgmt

VMware



4 x 800 GB SSD, 8 x 1.2 TB SAS HDD

Test#1:

- 2 x **SSD** in RAID1;
- 2 x **SSD** in RAID0 cache;
- 6 x **HDD** in RAID5;
- 1 x hotspare;

CrystalDiskMark 6.0.2 x64

File Settings Kinézet Súgó Language

All

9

32GiB

C: 14% (14/99GiB)

Read [MB/s]

Write [MB/s]

Seq
Q32T1

493.7

393.5

4KiB
Q8T8

136.5

125.7

4KiB
Q32T1

141.5

129.3

4KiB
Q1T1

22.52

16.80

CrystalDiskMark 6.0.2 x64

File Settings Kinézet Súgó Language

All

9

32GiB

C: 14% (14/99GiB)

Read [MB/s]

Write [MB/s]

Seq
Q32T1

528.9

707.5

4KiB
Q8T8

134.3

11.64

4KiB
Q32T1

140.8

12.48

4KiB
Q1T1

22.24

9.088

Test#1:

- 2 x **SSD** in RAID1;
- 2 x **SSD** in RAID0 cache;
- 6 x **HDD** in RAID5;
- 1 x hotspare;

Test#2:

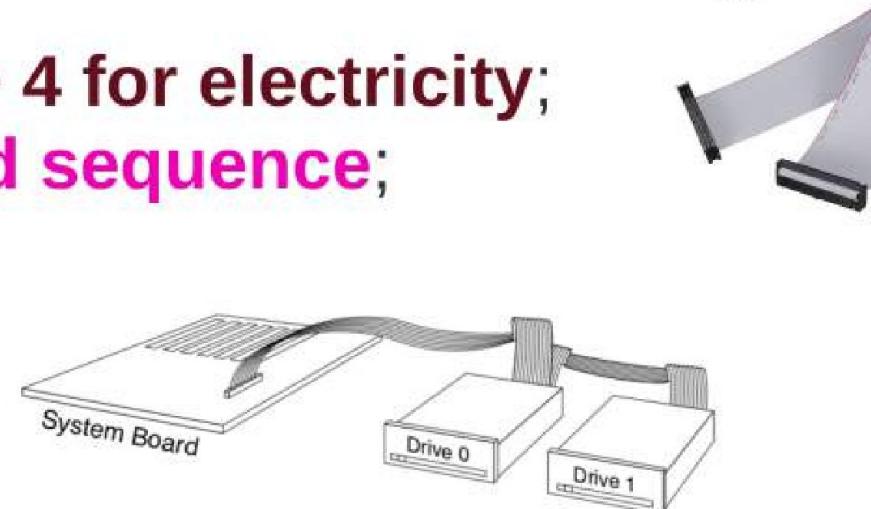
- 4 x **SSD** in RAID5;
- 8 x **HDD** in RAID6;
- no hotspare + no cache;



Attachment (PATA)

- An old PC standard initiated in 1986;
- Made disk controllers independent from drives;
- ATA 0 [**3.3 MB/s**] - 7 [**133 MB/s**];
- Integrated Device Electronics, EIDE at Western Digital;
- Originally **16-bits** wide, then **22, 28**;
- Maximum drive capacity $2^{28} \times 512 = 128$ GBytes;
- **2 drives attached to parallel bus** (master, slave, cable select, OS arbitration);
- Traditional **40-pin ribbon cable + 4 for electricity**;
- Traditionally **serialized command sequence**;

40Pin IDE Extension
Data Cable



PATA commands:

- 8 bit commands, from \$00 to \$FF;

40Pin IDE Extension

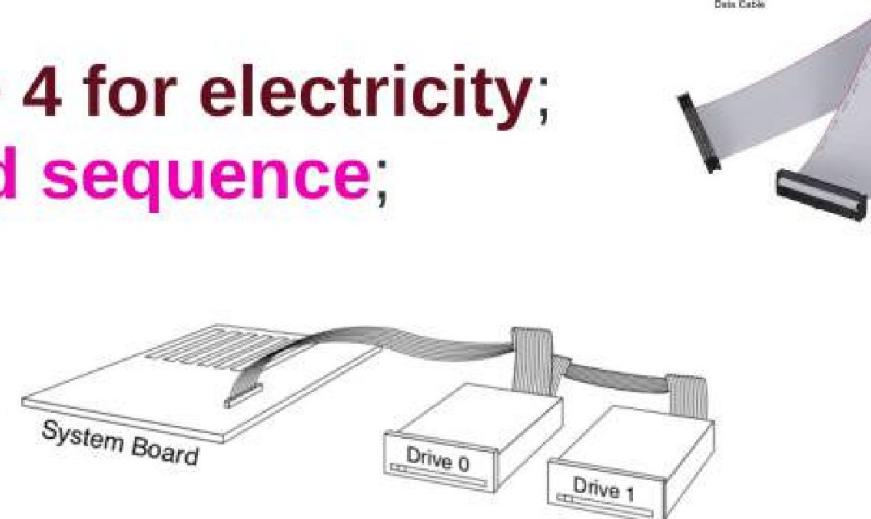
Data Cable



Attachment (PATA)

- An old PC standard initiated in 1986;
- Made disk controllers independent from drives;
- ATA 0 [**3.3 MB/s**] - 7 [**133 MB/s**];
- Integrated Device Electronics, EIDE at Western Digital;
- Originally **16-bits** wide, then **22, 28**;
- Maximum drive capacity $2^{28} \times 512 = 128$ GBytes;
- **2 drives attached to parallel bus** (master, slave, cable select, OS arbitration);
- Traditional **40-pin ribbon cable + 4 for electricity**;
- Traditionally **serialized command sequence**;

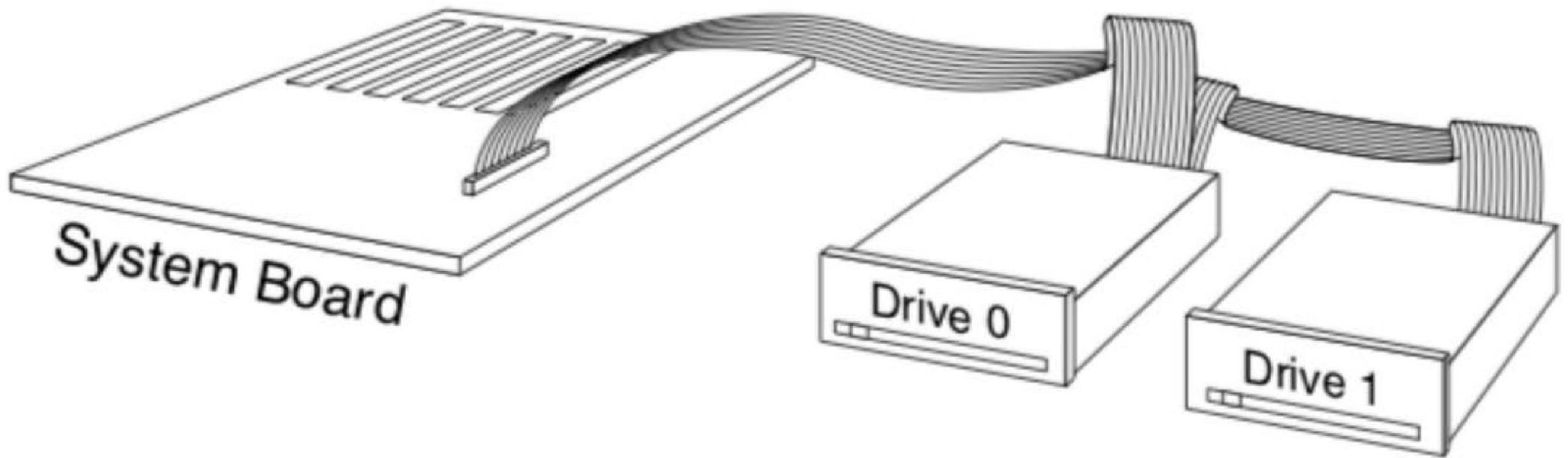
40Pin IDE Extension
Data Cable



PATA commands:

- 8 bit commands, from \$00 to \$FF;

d sequence,

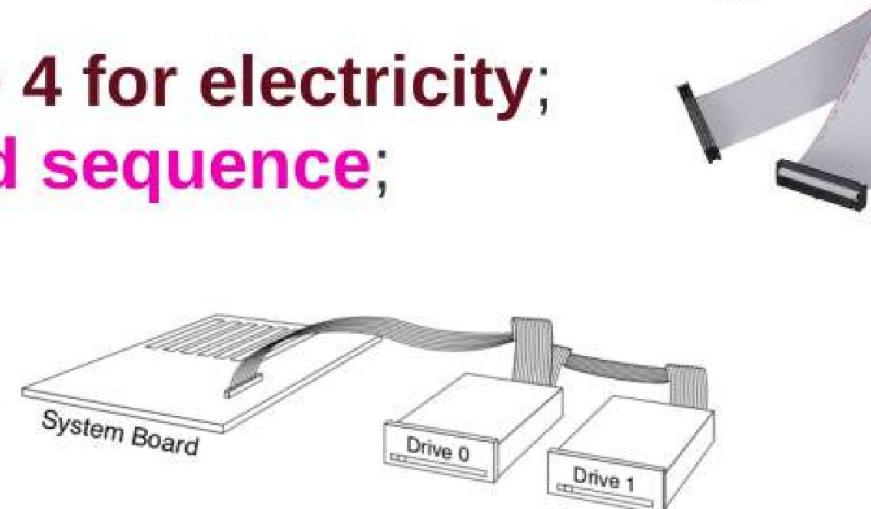


tus, Data;
ATA commands

Attachment (PATA)

- An old PC standard initiated in 1986;
- Made disk controllers independent from drives;
- ATA 0 [**3.3 MB/s**] - 7 [**133 MB/s**];
- Integrated Device Electronics, EIDE at Western Digital;
- Originally **16-bits** wide, then **22, 28**;
- Maximum drive capacity $2^{28} \times 512 = 128$ GBytes;
- **2 drives attached to parallel bus** (master, slave, cable select, OS arbitration);
- Traditional **40-pin ribbon cable + 4 for electricity**;
- Traditionally **serialized command sequence**;

40Pin IDE Extension
Data Cable



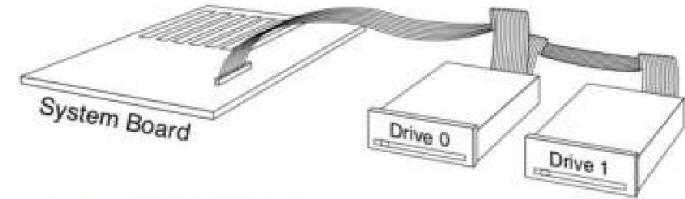
PATA commands:

- 8 bit commands, from \$00 to \$FF;

- Traditional **40-pin ribbon cable + 4 for electricity**;
- Traditionally **serialized command sequence**;

PATA commands:

- 8 bit commands, from \$00 to \$FF;
- IO registers: **Command, Control, Status, Data**;
- Command registers are used to send ATA commands over the parallel ports;
- Status is received in status registers;
- ATA-8 contains a much broader advanced feature set;
- **Protocols:**
 - PIO read, write;
 - DMA transfer;
 - Non-data commands.



Command Block registers							
1	0	0	0	0	Data register		
1	0	0	0	1	Error register		
1	0	0	1	0	Sector Count register		
1	0	0	1	1	Sector Number register		
1	0	1	0	0	ylinder Low register		
1	0	1	0	1	ylinder High register		
1	0	1	1	0	Disk/Head register		
1	0	1	1	1	Status register		
					Control register		

<http://www.t13.org/documents/uploadeddocuments/docs2006/d1699r3f-ata8-ac8.pdf>
<ftp://ftp.seagate.com/acrobat/reference/111-1c.pdf>

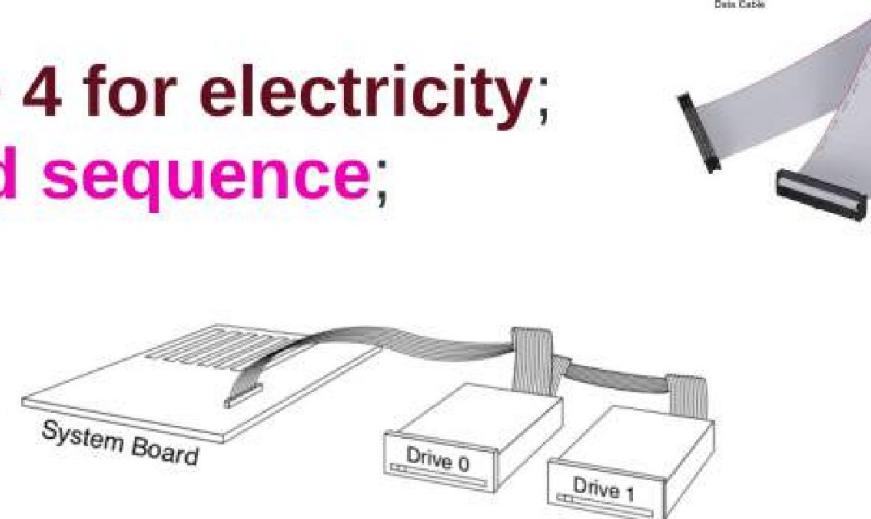
General command set:

- EXECUTE DEVICE DIAGNOSTIC
- FLUSH CACHE
- IDENTIFY DEVICE
- READ DMA
- READ MULTIPLE
- **READ SECTOR(S)**
- READ VERIFY SECTOR(S)
- SET FEATURES
- SET MULTIPLE MODE
- WRITE DMA
- WRITE MULTIPLE
- **WRITE SECTOR(S)**

Attachment (PATA)

- An old PC standard initiated in 1986;
- Made disk controllers independent from drives;
- ATA 0 [**3.3 MB/s**] - 7 [**133 MB/s**];
- Integrated Device Electronics, EIDE at Western Digital;
- Originally **16-bits** wide, then **22, 28**;
- Maximum drive capacity $2^{28} \times 512 = 128$ GBytes;
- **2 drives attached to parallel bus** (master, slave, cable select, OS arbitration);
- Traditional **40-pin ribbon cable + 4 for electricity**;
- Traditionally **serialized command sequence**;

40Pin IDE Extension
Data Cable



PATA commands:

- 8 bit commands, from \$00 to \$FF;

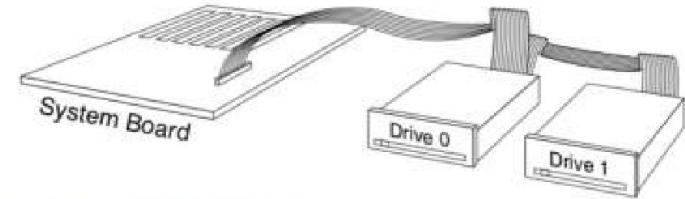
Tagged Command Queuing:

- One request at a time to drives;
- As in case of SCSI, **sending multiple commands at the same time;**
- Due to the high number of interrupts this solution required lot of CPU power and was not really efficient;

- Traditional **40-pin ribbon cable + 4 for electricity**;
- Traditionally **serialized command sequence**;

PATA commands:

- 8 bit commands, from \$00 to \$FF;
- IO registers: **Command, Control, Status, Data**;
- Command registers are used to send ATA commands over the parallel ports;
- Status is received in status registers;
- ATA-8 contains a much broader advanced feature set;
- **Protocols:**
 - PIO read, write;
 - DMA transfer;
 - Non-data commands.



Command Block registers							
1	0	0	0	0	0	0	Data register
1	0	0	0	1	0	0	Error register
1	0	0	0	1	0	0	Sector Count register
1	0	0	0	1	1	0	Sector Number register
1	0	1	0	0	0	0	Sector Low register
1	0	1	0	0	0	1	Cylinder High register
1	0	1	1	0	0	0	Disk/Head register
1	0	1	1	1	0	0	Status register
1	0	1	1	1	1	0	Control register

<http://www.t13.org/documents/uploadeddocuments/docs2006/d1699r3f-ata8-ac8.pdf>
<ftp://ftp.seagate.com/acrobat/reference/111-1c.pdf>

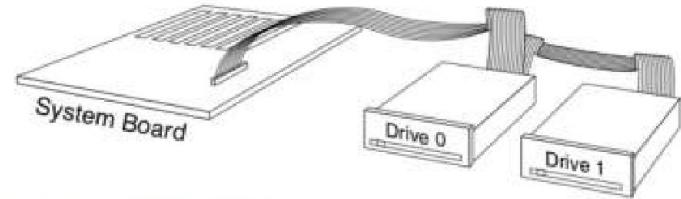
Command Block registers

1	0	0	0	0	Data register	1F0	
1	0	0	0	1	Error register	Feature register	1F1
1	0	0	1	0	Sector Count register	1F2	
1	0	0	1	1	Sector Number register	1F3	
1	0	1	0	0	Cylinder Low register	1F4	
1	0	1	0	1	Cylinder High register	1F5	
1	0	1	1	0	Drive/Head register	1F6	
1	0	1	1	1	Status register	Command register	1F7

- Traditional **40-pin ribbon cable + 4 for electricity**;
- Traditionally **serialized command sequence**;

PATA commands:

- 8 bit commands, from \$00 to \$FF;
- IO registers: **Command, Control, Status, Data**;
- Command registers are used to send ATA commands over the parallel ports;
- Status is received in status registers;
- ATA-8 contains a much broader advanced feature set;
- **Protocols:**
 - PIO read, write;
 - DMA transfer;
 - Non-data commands.



Command Block registers							
1	0	0	0	0	Data register		
1	0	0	0	1	Error register		
1	0	0	1	0	Sector Count register		
1	0	0	1	1	Sector Number register		
1	0	1	0	0	ylinder Low register		
1	0	1	0	1	ylinder High register		
1	0	1	1	0	Disk/Head register		
1	0	1	1	1	Status register		
					Control register		

<http://www.t13.org/documents/uploadeddocuments/docs2006/d1699r3f-ata8-ac8.pdf>
<ftp://ftp.seagate.com/acrobat/reference/111-1c.pdf>

Serial ATA



- It is **serialized version** of Parallel ATA standard;
- Serialization helps reducing cabling structure and increasing clock speed -> increased transfer speeds:
SATA1 [150 MB/s], SATA2 [300 MB/s], SATA3 [600 MB/s], SATA Express, 3.2 [**2 GB/s**];
- 100% software compatible with PATA standard, enhanced command set;
- Dedicated **p2p crossbar** connection among the host bus and drives, eliminating the bus (like data network);
- Frame based protocol;
- Hot-plug, hot-swap;
- Layered communication structure;

SATA Port Multipliers:

- Up to 15 drives;
- Cost efficient connection of drives;
- Share the HBA bandwidth;



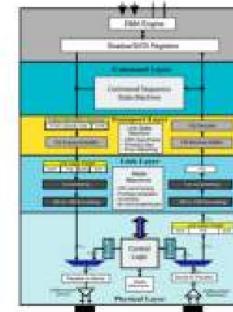
Frame Information Structures:

- Generated in the transport layer;
- Data movement setup;
- Read and write data;
- Used to control IO;
- Carries **contents of ATA command registers**;

	+3	+2	+1	+0
DW 0	Features	Command	C/R# Received	FIS Type (2bit)
DW 1	DevLd	Cy High	Cy Low	Sector Number
DW 2	Features (cont)	Cy High (cont)	Cy Low (cont)	Sect Num (cont)
DW 3	Control	Reserved (0)	Sec Count (cont)	Sector Count
DW 4	Reserved (0)	Reserved (0)	Reserved (0)	Reserved (0)

Benefits:

- lower pin count;
- larger performance, full duplex;
- simple configuration;
- lower voltage;



Link layer:

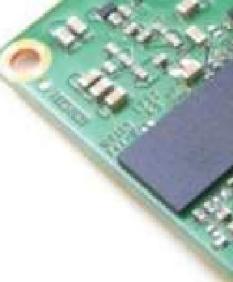
- 8b/10b encoding;
- converts data streams into frames;
- CRC;
- flow control;

Transport layer:

- managing FIS;

Native Command Queuing:

- **Multiple IO commands** in parallel (like SCSI);
- Due to mechanical constraints it is not efficient to serve IO requests in the order of their issue = **elevator scheduling problem**;
- Eliminate unnecessary travels = significantly improve performance;
- MCQ allows **out of order** IO command completion;
- Reordering of 32 IO commands;
- Command order **received** vs. command order optimized for shortest completion time;
- minimize **rotation** and **head movement**;



- It is **serialized version** of Parallel ATA standard;
- Serialization helps reducing cabling structure and increasing clock speed -> increased transfer speeds: SATA1 [**150 MB/s**], SATA2 [**300 MB/s**], SATA3 [**600 MB/s**], SATA Express, 3.2 [**2 GB/s**];
- 100% software compatible with PATA standard, enhanced command set;
- Dedicated **p2p crossbar** connection among the host bus and drives, eliminating the bus (like data network);
- Frame based protocol;
- Hot-plug, hot-swap;
- Layered communication structure;

SATA Port Multipliers:

- Up to 15 drives;
- Cost efficient connection of drives;
- Share the HBA bandwidth;

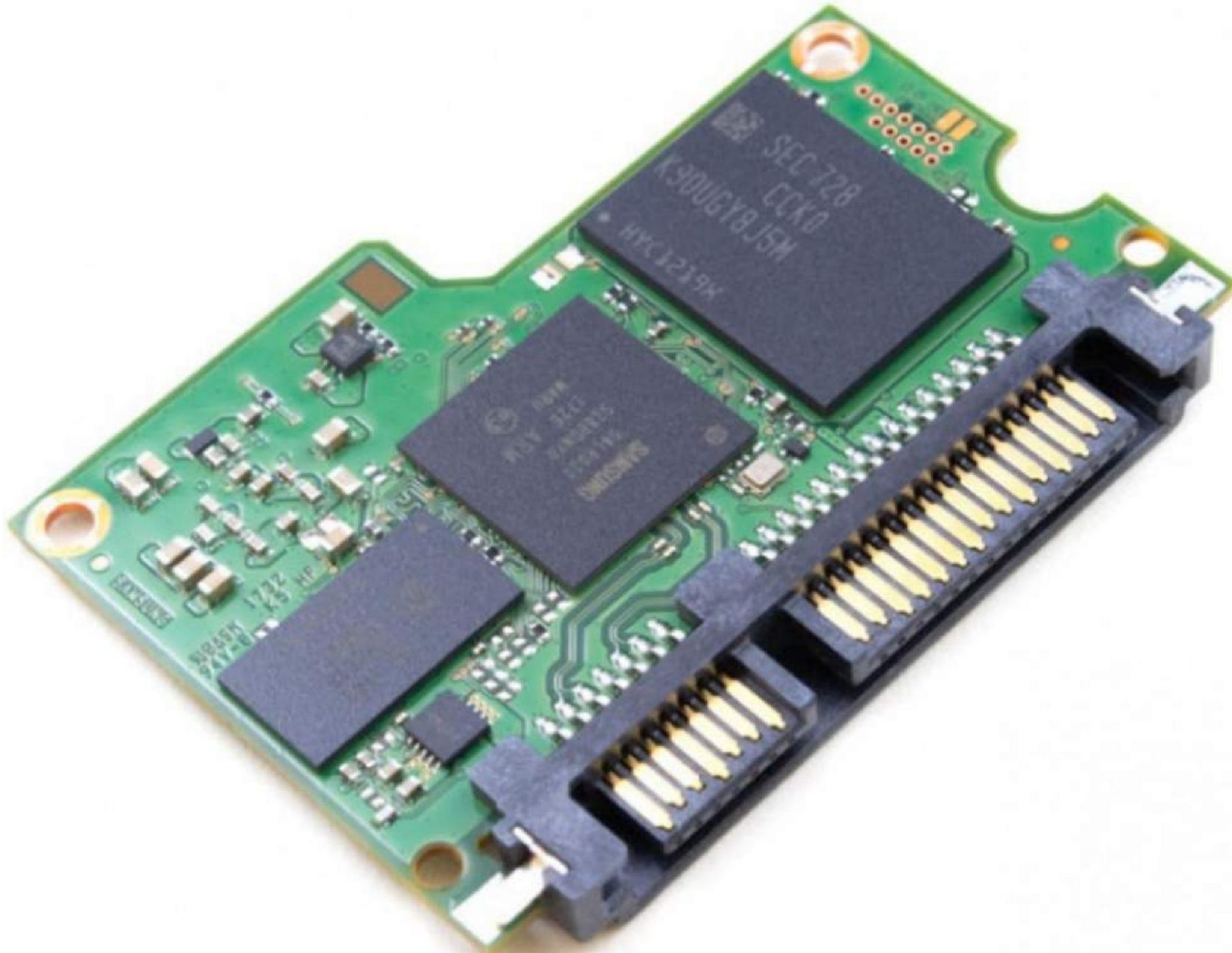


SATA DM compatibility: one port, one cable, four drives; 4-ports = 16 drives

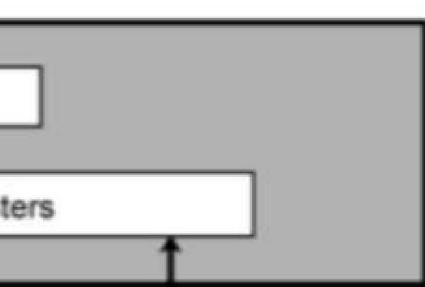
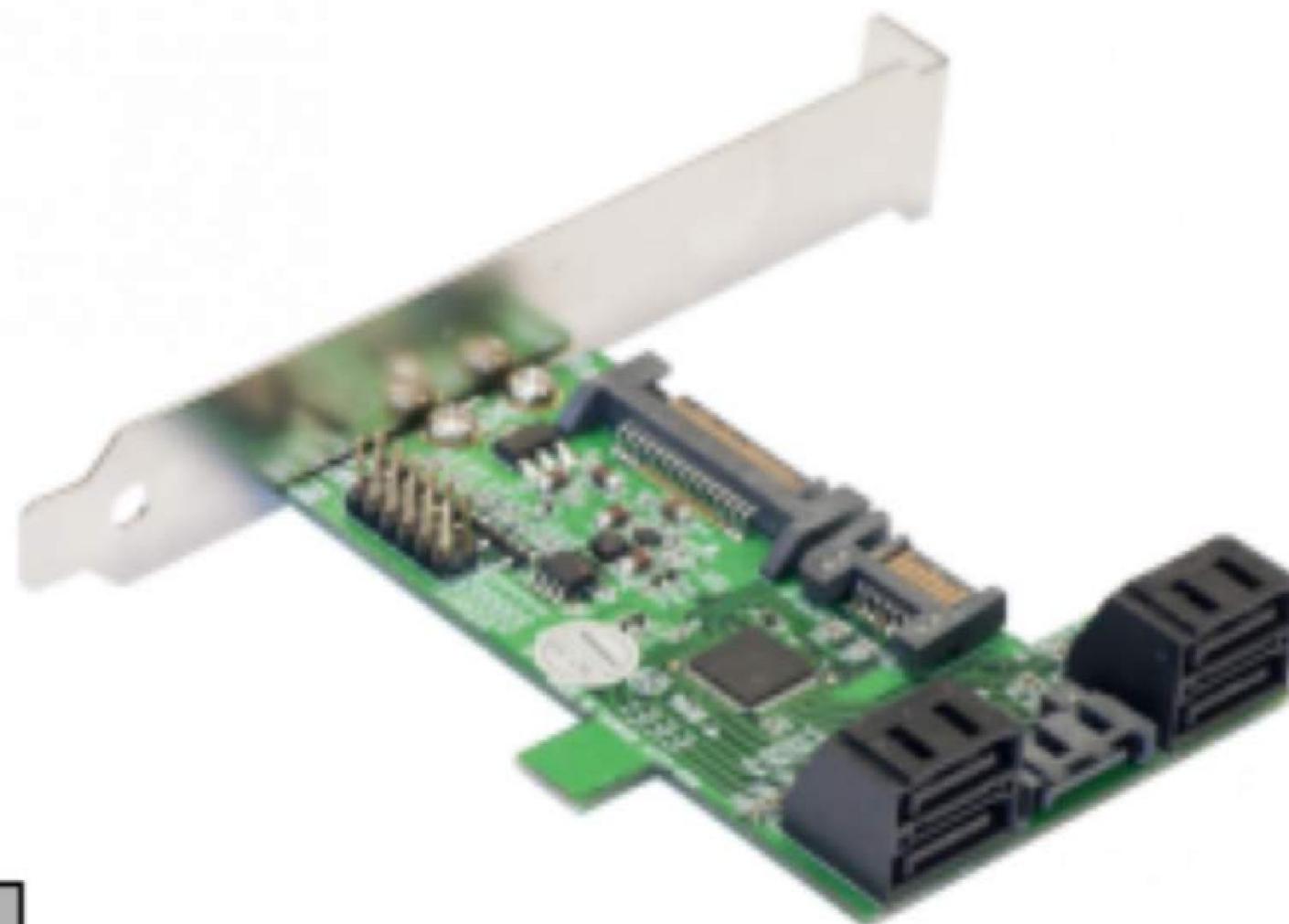
Benefits:

- Lower pin count





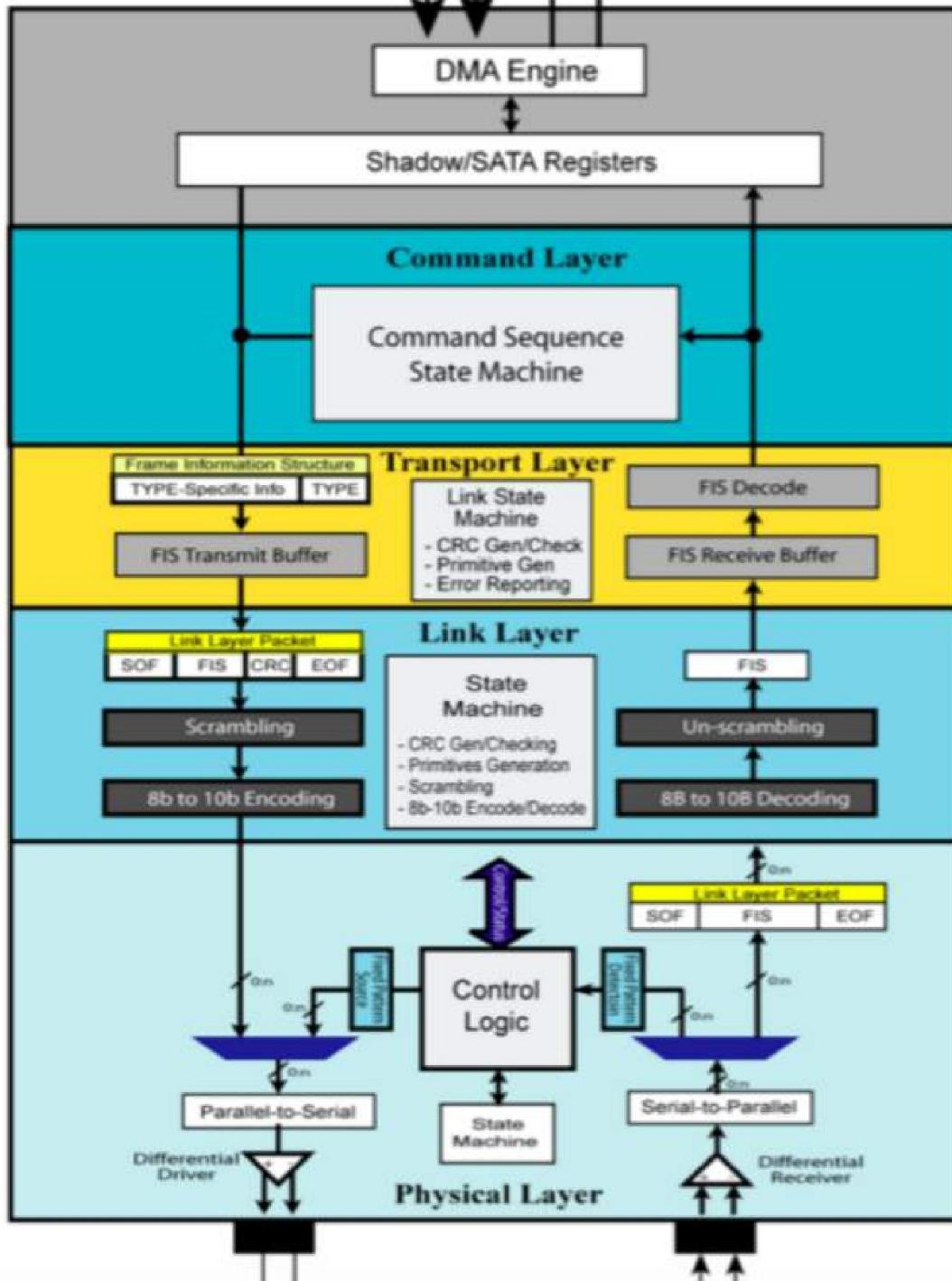




Benefits:

- lower pin count;
- larger performance, full duplex;
- simple configuration;
- lower voltage;

host
work);



Link layer

- 8b/10b encoding
- concatenation of frames
- CRC generation
- flow control

Transp

Link layer:

- 8b/10b encoding;
- converts data streams into frames;
- CRC;
- flow control;

Transport layer

Transport layer:

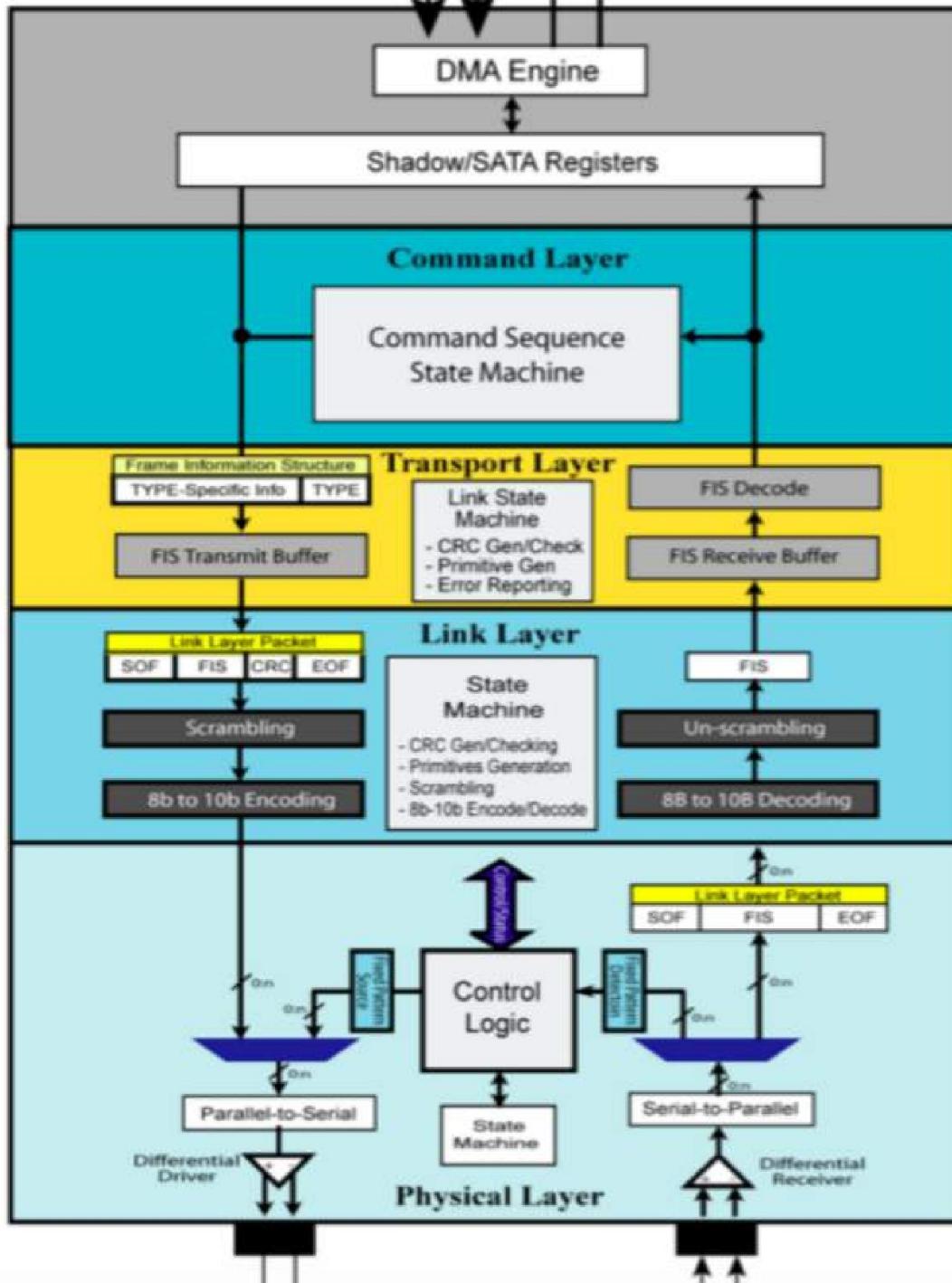
- managing FIS;

Frame Information Structures:

- Generated in the transport layer;
- Data movement setup;
- Read and write data;
- Used to control IO;
- Carries **contents of ATA command registers**;

	+3	+2	+1	+0	
DW 0	Features	Command	C R R	Reserved	FIS Type (27h)
DW 1	Dev/Head	Cyl High	Cyl Low		Sector Number
DW 2	Features (exp)	Cyl High (exp)	Cyl Low (exp)		Sec Num (exp)
DW 3	Control	Reserved (0)	Sec Count (exp)		Sector Count
DW 4	Reserved (0)	Reserved (0)	Reserved (0)		Reserved (0)

host
work);



Link layer

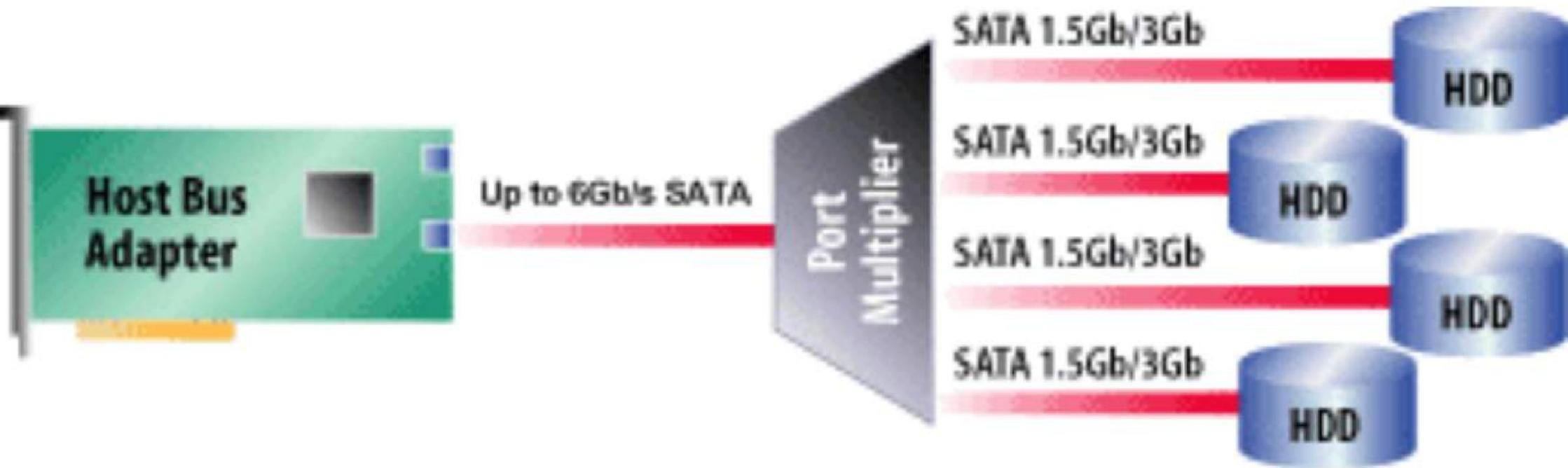
- 8b/10b encoding
- concatenation of frames
- CRC generation
- flow control

Transp

layered communication str

SATA Port Multipliers:

- Up to 15 drives;
- Cost efficient connection of drives;
- Share the HBA bandwidth;



SATA PM connectivity: one port, one cable, four drives; 4-ports = 16 drives

ull

Native Command Queuing:

- **Multiple IO commands** in parallel (like SCSI);
- Due to mechanical constraints it is not efficient to serve IO requests in the order of their issue
= **elevator scheduling problem**;
- Eliminate unnecessary travels = significantly improve performance;
- MCQ allows **out of order** IO command completion;
- Reordering of 32 IO commands;
- Command order **received** vs. command order **optimized for shortest completion time**;
- minimize **rotation** and **head movement**;

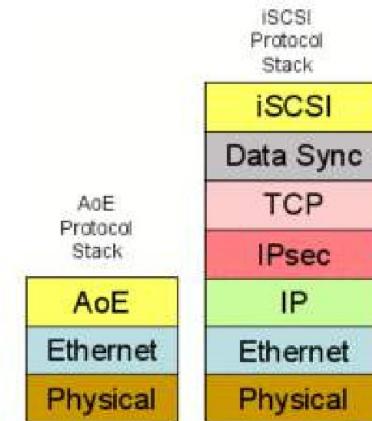
ATA over Ethernet

Octet #	Octet Value	Description
0	0x00	Ethernet Destination MAC Address
4	0x00	Ethernet Destination (most)
5	0x00	Ethernet Source (MAC Address)
8	0x00	Ethernet Source MAC Address (least)
12	0x00	Ethernet Type (0x0d45) + Ver + Flags + Error
13	0x00	Error
14	0x00	Short
15	0x00	Command
20	0x00	Tag
24	0x00	ATA

- Simplicity attracts! :)
- Encapsulates **ATA commands into Ethernet carrier**;
- Carries ATA commands information to be put into registers;
- Data is sent in 8 KBytes = 16 blocks;
- Replaces traditional ATA transport layer (bus and cabling) to Ethernet;
- No need for multipathing;
- Lightweight protocol stack;
- Non-routing protocol;
- **AoE initiator - AoE target**;

Practice #5: Running AoE devices

```
aoe-firmware+vhbae installation  
root@aoe: ~# aoe,devlist=10  
forsid  
mpt2sas test  
vhbae 0.0 is initialized &  
aoe-discovered  
aoe-10M  
forsid  
vhbae 0.0 is initialized &  
mpt2sas test  
vhbae 0.0 is initialized &  
aoe-discovered  
aoe-10M  
forsid  
vhbae 0.0 is initialized &  
mpt2sas test  
vhbae 0.0 is initialized &
```



<https://www.linuxjournal.com/content/mastering-ata-over-ethernet>

<http://storagegaga.com/aoe-all-about-etherne/>

<http://download.alyseo.com:81/pub/partners/Coraid/Docs/AoE/AoEDescription.pdf>

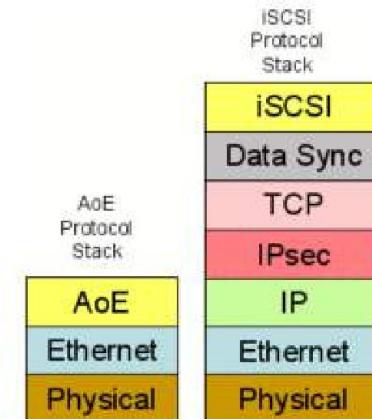
ATA over Ethernet

Octet #	Octet Value	Description
0	0x00	Ethernet Destination MAC Address
4	0x00	Ethernet Destination (most)
5	0x00	Ethernet Source (MAC Address)
8	0x00	Ethernet Source MAC Address (least)
12	0x00	Ethernet Type (0x0d45) Ver Flags Error
13	0x00	Flags
14	0x00	Short
15	0x00	Compressed
20	0x00	Tag
24	0x00	ATA

- Simplicity attracts! :)
- Encapsulates **ATA commands into Ethernet carrier**;
- Carries ATA commands information to be put into registers;
- Data is sent in 8 KBytes = 16 blocks;
- Replaces traditional ATA transport layer (bus and cabling) to Ethernet;
- No need for multipathing;
- Lightweight protocol stack;
- Non-routing protocol;
- **AoE initiator - AoE target**;

Practice #5: Running AoE devices

```
aoe-firmware + vblade installation  
root@aoe: ~# aoe,devlist=10  
banned  
maxtargets=80  
vblade 0.0 is initialized @  
aoe-discovered  
aoe-10M  
fibre IDENTIFIER:0x0  
mdadm --create /dev/mapper/aoe-1 --level=1 --raid-devices=2 /dev/nvme0n1 0 /dev/nvme0n1 1
```



<https://www.linuxjournal.com/content/mastering-ata-over-ethernet>

<http://storagegaga.com/aoe-all-about-etherne/>

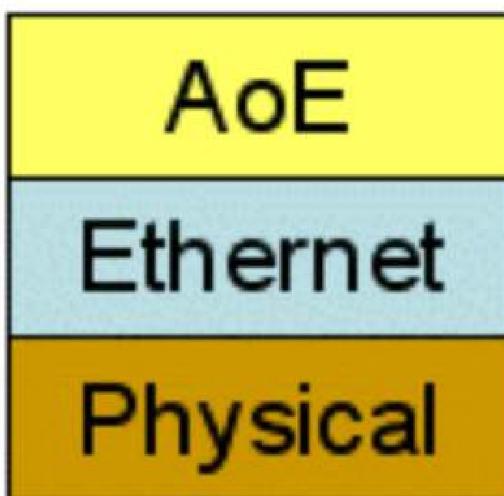
<http://download.alyseo.com:81/pub/partners/Coraid/Docs/AoE/AoEDescription.pdf>

AoE devices

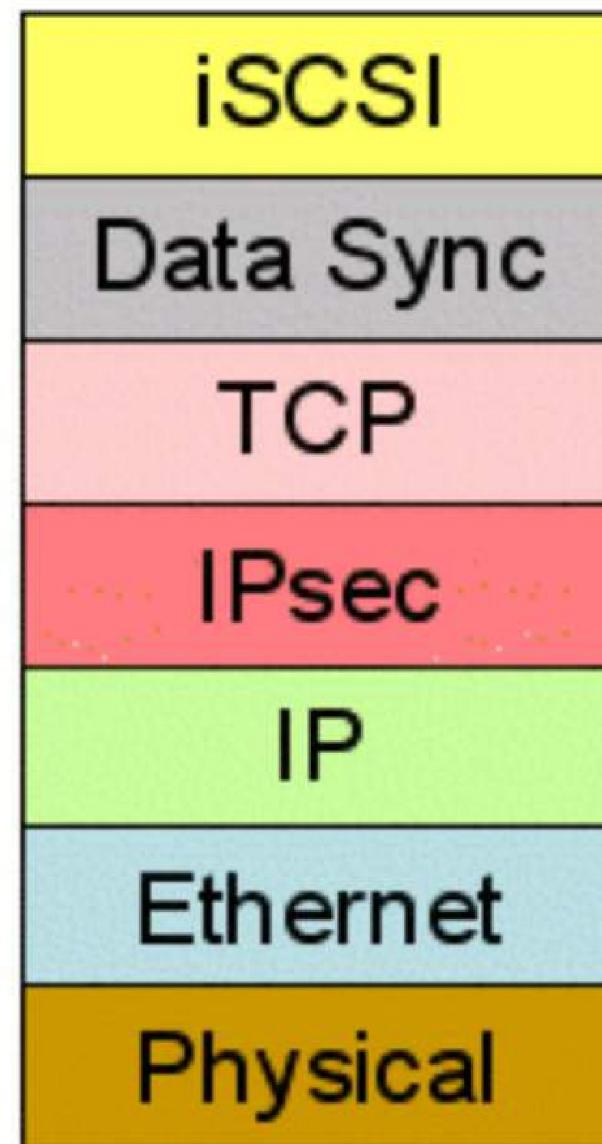
=10

vel=1 --raid-
ev/etherd/e0.1

AoE
Protocol
Stack



iSCSI
Protocol
Stack



SCSI Protocol



Practice #6: SCSI commands

Answer: 1-4: B; 5-8: C; 9-12: D; 13-16: A; 17-20: B; 21-24: C; 25-28: D; 29-32: B; 33-36: A; 37-40: C; 41-44: B; 45-48: D; 49-52: A; 53-56: C; 57-60: B; 61-64: D; 65-68: A; 69-72: C; 73-76: B; 77-80: D; 81-84: A; 85-88: C; 89-92: B; 93-96: D; 97-100: A

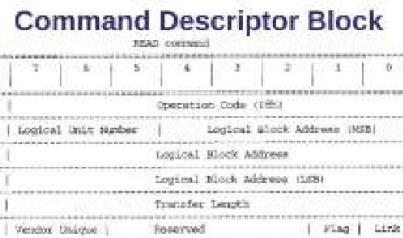
- History starts at 1978 (Alan Shugart, SASI), continues at 1986 (Small Computer System Interface, SCSI-1, first standardization);
 - **Parallel bus standard** to ease the communication between the host and the low level disk drives;
 - Uses 8/16 data + parity + 9 control signals;
 - Originally:
 - 8 bit + parity wide;
 - 8 or 16 drives could be attached;
 - 10 MHz.
 - **Interface names:** SCSI-1 [5 MB/s], Fast SCSI [10 MB/s], Fast Wide SCSI [20 MB/s], Ultra Wide SCSI [40 MB/s], Ultra-320 SCSI [320 MB/s], and Ultra-640 SCSI [640 MB/s];
 - **SCSI initiator - SCSI target(s):**
 - Bus used as **state machine**: bus-free, arbitration, device selection, command, reselection (target disconnect in long operations), data, message, status;

Device types:

- 5-bit field sent in Inquiry;
 - broad range: scanners, printers, block disks, CD-ROMS, floppies...
\$00-\$1F

Device identification:

- manually set: 0-7, 0 - bootable disk,
7 - initiator;
 - BIOS set;
 - slot-set in storage shelves;
 - ID discovery for SAS.



Difference between SCSI and ATA

- On non-protocol level: **Enterprise Storage vs Personal Storage;**
 - Internal CPU driven IO vs **Programmed IO** (lack of DMA);
 - Variable block size vs **Fixed block size;**
 - Parallel queuing vs **Serial command set;**
 - 15 drives vs **2 x 2 drives;**
 - **larger buffers/speed vs smaller buffers/speed;**

- Test unit ready
- Inquiry
- Request sense
- Send/Receive diagnostic
- Start/Stop unit
- Read capacity
- Format unit
- Read (four variants)
- Write (four variants)
- Log sense
- Mode sense
- Mode select

<https://en.wikipedia.org/wiki/SCSI>

<http://fundasbykrishna.blogspot.com/2013/04/scsi-fundamentals-2.html>

- History starts at 1978 (Alan Shugart, SASI), continues at 1986 (Small Computer System Interface, SCSI-1, first standardization);
- **Parallel bus standard** to ease the communication between the host and the low level disk drives;
- Uses 8/16 data + parity + 9 control signals;
- Originally:
 - 8 bit + parity wide;
 - 8 or 16 drives could be attached;
 - 10 MHz.
- **Interface names:** SCSI-1 [5 MB/s], Fast SCSI [10 MB/s], Fast Wide SCSI [20 MB/s], Ultra Wide SCSI [40 MB/s], Ultra-320 SCSI [320 MB/s], and Ultra-640 SCSI [640 MB/s];
- **SCSI initiator - SCSI target(s);**
- Bus used as **state machine**: bus-free, arbitration, device selection, command, reselection (target disconnect in long operations), data, message, status;

Device types:

Command Des

READ command



Device types:

- 5-bit field sent in **Inquiry**;
- broad range: scanners, printers, block disks, CD-ROMS, floppies...
\$00-\$1F

Device identification:

- manually set: 0-7, 0 - bootable disk, 7 - initiator;
- BIOS set;
- slot-set in storage shelves;
- ID discovery for SAS.



2A 00 00 00 1F F0 00 00 01 0
sg_raw -r 512 /dev/sdb 28 0

<http://sg.danny.cz/s>

Connectors:

- SCSI Parallel Interface;
- Fibre Channel;
- Serial Attached SCSI;
- HD Mini SAS;







2A 00 00 00 1F F0 00 00 01 0
sg_raw -r 512 /dev/sdb 28 0

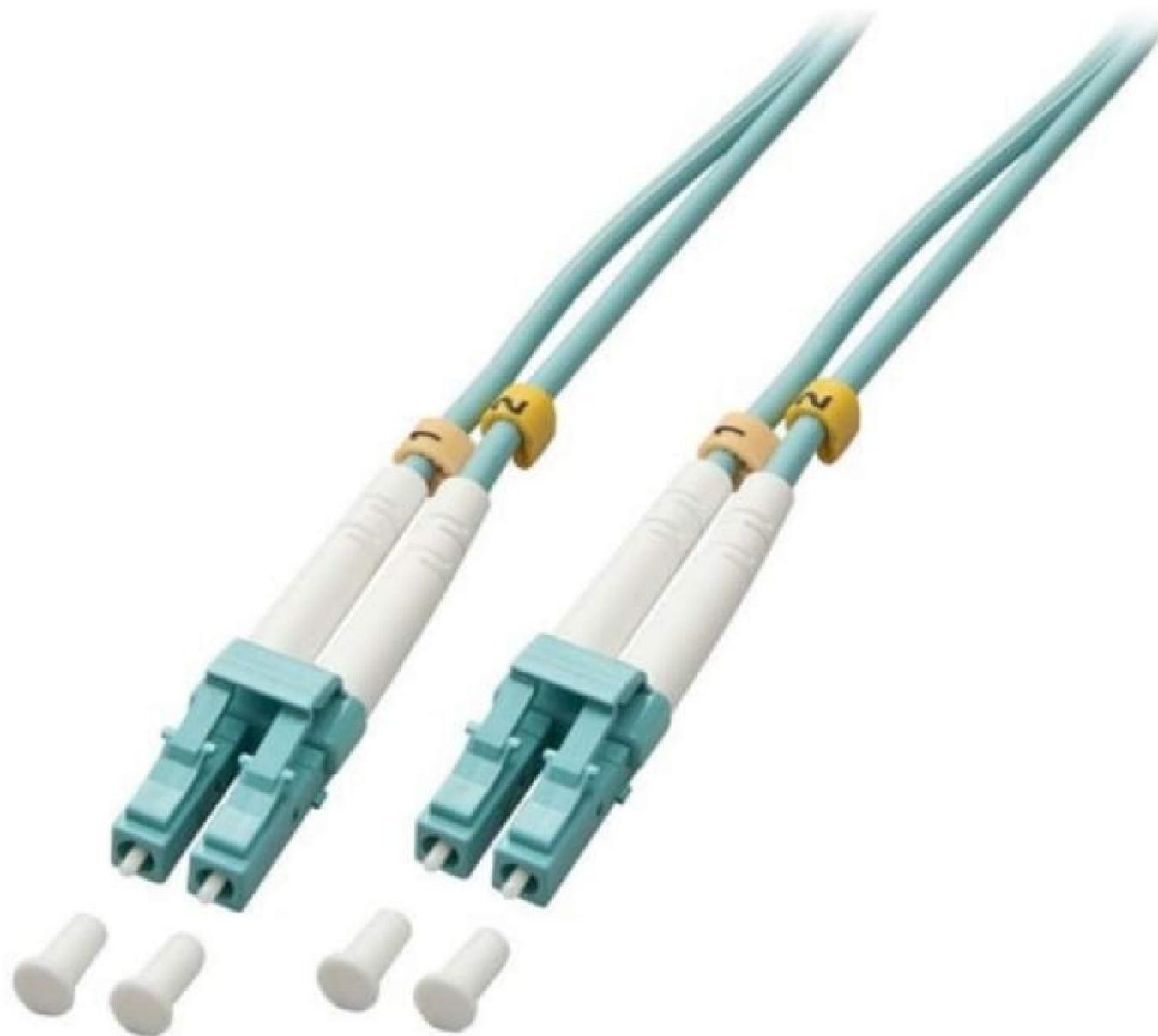
<http://sg.danny.cz/s>

Connectors:

- SCSI Parallel Interface;
- Fibre Channel;
- Serial Attached SCSI;
- HD Mini SAS;



ce;





2A 00 00 00 1F F0 00 00 01 0
sg_raw -r 512 /dev/sdb 28 0

<http://sg.danny.cz/s>

Connectors:

- SCSI Parallel Interface;
- Fibre Channel;
- Serial Attached SCSI;
- HD Mini SAS;



• HD Mini SAS,



DSISEDIMS





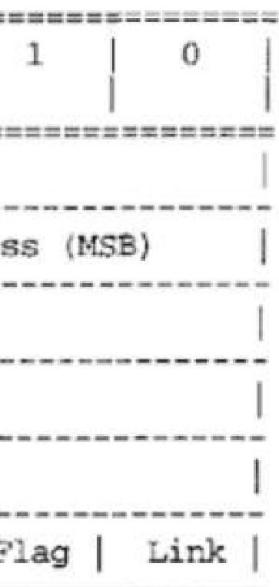
SCSI command families:

- **non-data** commands;
- **reading** data from target;
- **writing** data from the initiator to target;
- **bidirectional**.

Descriptor Block

Test unit ready

lock



e vs **Personal**

PIO (lack of DMA);

- Test unit ready
- Inquiry
- Request sense
- Send/Receive diagnostic
- Start/Stop unit
- Read capacity
- Format unit
- Read (four variants)
- Write (four variants)
- Log sense
- Mode sense
- Mode select

Command Descriptor Block

READ command								
Bit	7	6	5	4	3	2	1	0
Byte								
0								
1	Logical Unit Number				Logical Block Address (MSB)			
2					Logical Block Address			
3					Logical Block Address (LSB)			
4					Transfer Length			
5	Vendor Unique				Reserved		Flag	Link

Difference between **SCSI** and **ATA**:

3	Logical Block Address (LSB)				
4	Transfer Length				
5	Vendor Unique	Reserved		Flag	Link

Difference between **SCSI** and **ATA**:

- On non-protocol level: **Enterprise Storage** vs **Personal Storage**;
- **Internal CPU driven IO** vs **Programmed IO** (lack of DMA);
- **Variable block size** vs **Fixed block size**;
- **Parallel queuing** vs **Serial command set**;
- **15 drives** vs **2 x 2 drives**;
- **larger buffers/speed** vs **smaller buffers/speed**;

<http://pages.cs.wisc.edu/~remzi/Classes/838/Fall2001/Papers/scsi-ata.pdf>

SCSI

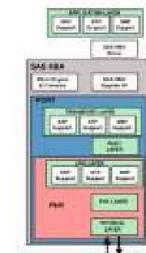
pot.com/2013/04/scsi-fundamentals-2/

Serial Attached SCSI

- P2P crossbar attached version of SCSI;
- Serializes transport between host adapter expanders, and drives;
- Certain compatibility with SATA;
- SAS1 [300 MB/s], SAS2 [600 MB/s], SAS3 [1.5 GB/s], SAS4 [3 GB/s];
- Identifies ports, HAs, devices with unique WWN;
- narrow and wide ports;

SAS and SCSI:

- P2P vs multidrop bus;
- no termination vs terminated bus;
- 65535 devices vs 8/16 devices;
- dedicated full bandwidth vs. shared bus bandwidth;



SAS and SATA:

- multiple initiators vs single initiator;
- TCQ vs NCQ;
- SCSI command set vs ATA command set;
- multipath IO vs single path IO;
- 1.6V signaling vs 0.6V signaling;
- up to 10 meters vs 1 meter cabling;

Serial SCSI Protocol:

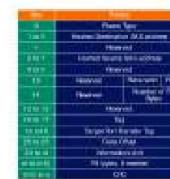
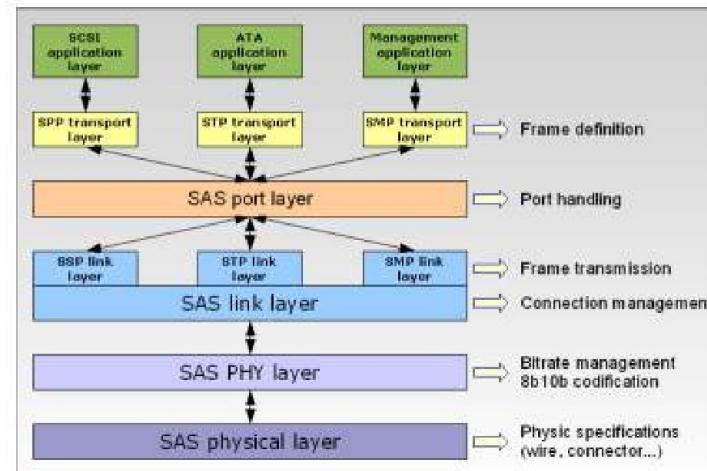
- Wraps SCSI commands into serialized frames;
- Read Command, Data from target, Response;
- Write Command, XFER_RDY, Data to Target, Response;
- Nondata, Command, Response.

SCSI Management Protocol:

- management frames to set up and inquire about SAS topology;
- set up and manage SAS routing;
- like RIP, BGP, or OSPF in case of Ethernet;

SATA Transport Protocol:

- Wraps SATA commands;
- Tunneling;
- Uses underlying SAS layers for transport;
- FIS.



Serial Attached

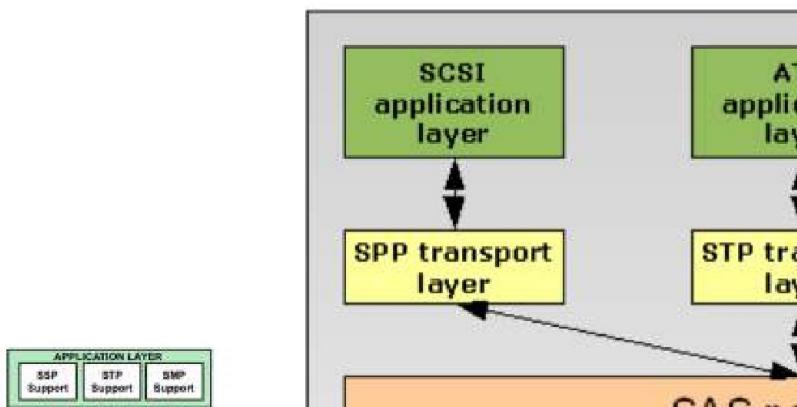
Serial SCSI

- **P2P crossbar** attached version of SCSI;
- **Serializes transport** between host adapter **expanders**, and drives;
- Certain compatibility with SATA;
- SAS1 [300 MB/s], SAS2 [600 MB/s], SAS3 [1.5 GB/s], SAS4 [3 GB/s];
- Identifies ports, HAs, devices with unique WWN;
- **narrow** and **wide** ports;

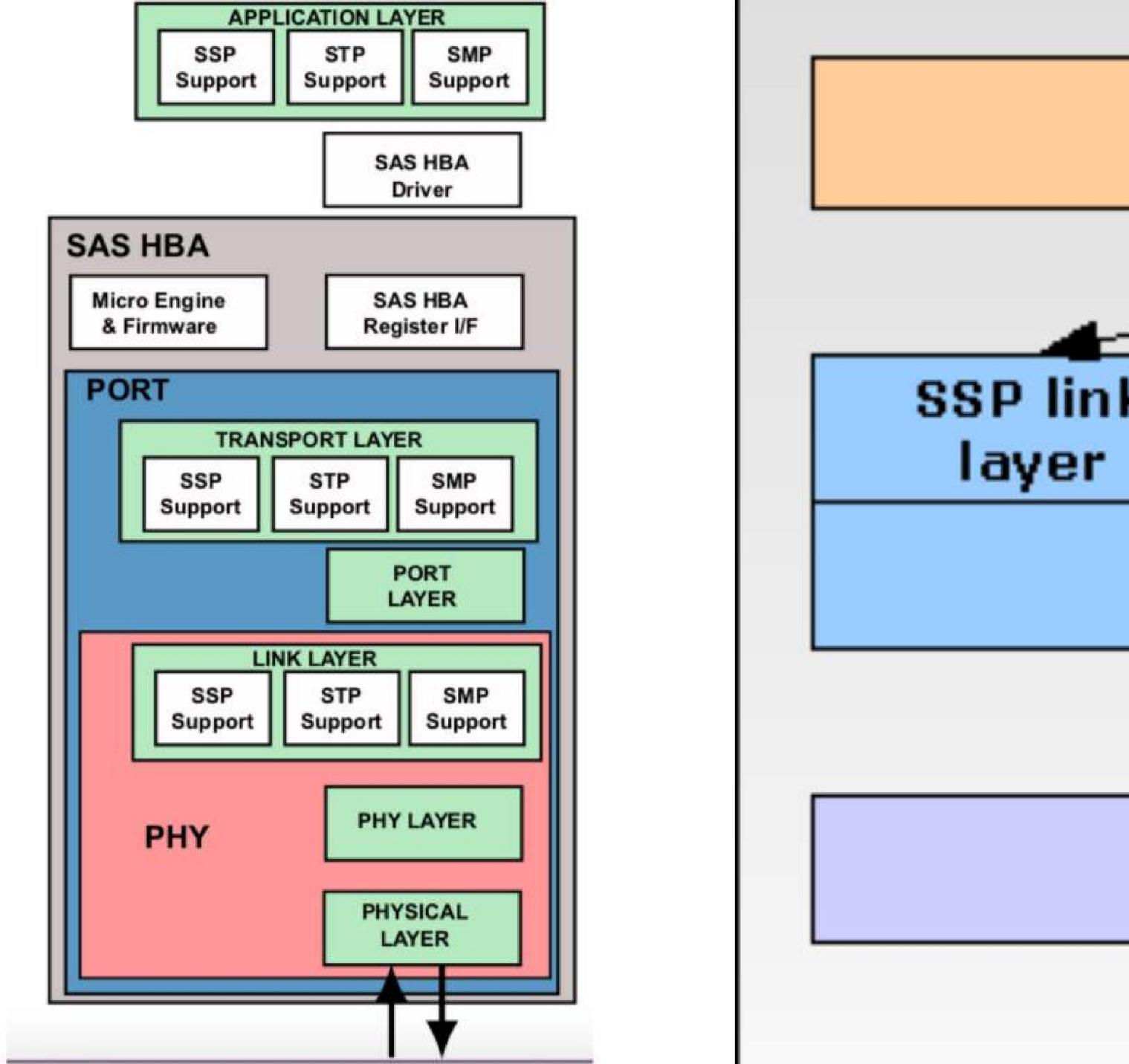
- Wraps SCSI serialized
- Read Command Response
- Write Command to Target
- Nondata

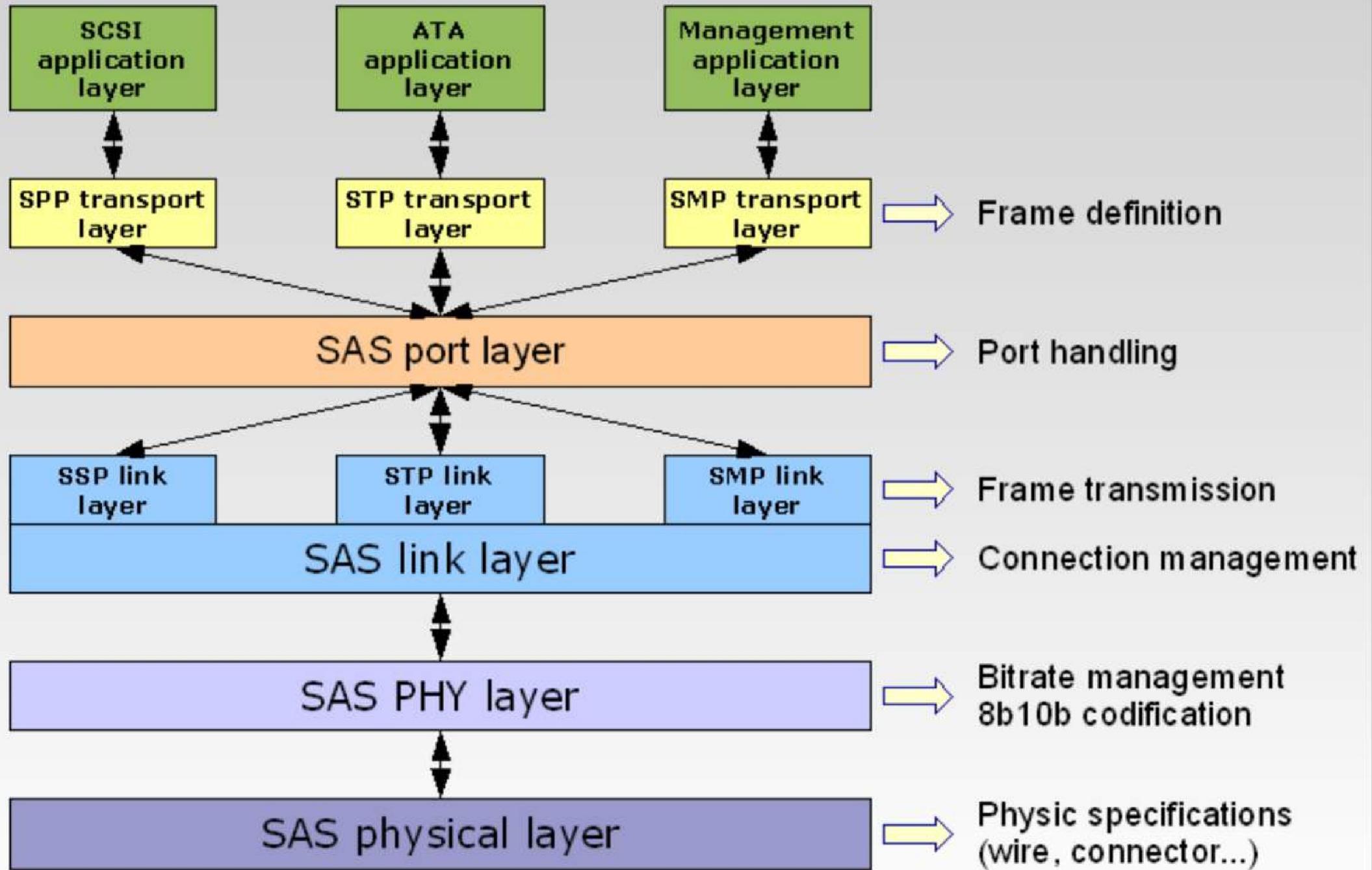
SAS and SCSI:

- P2P vs multidrop bus;

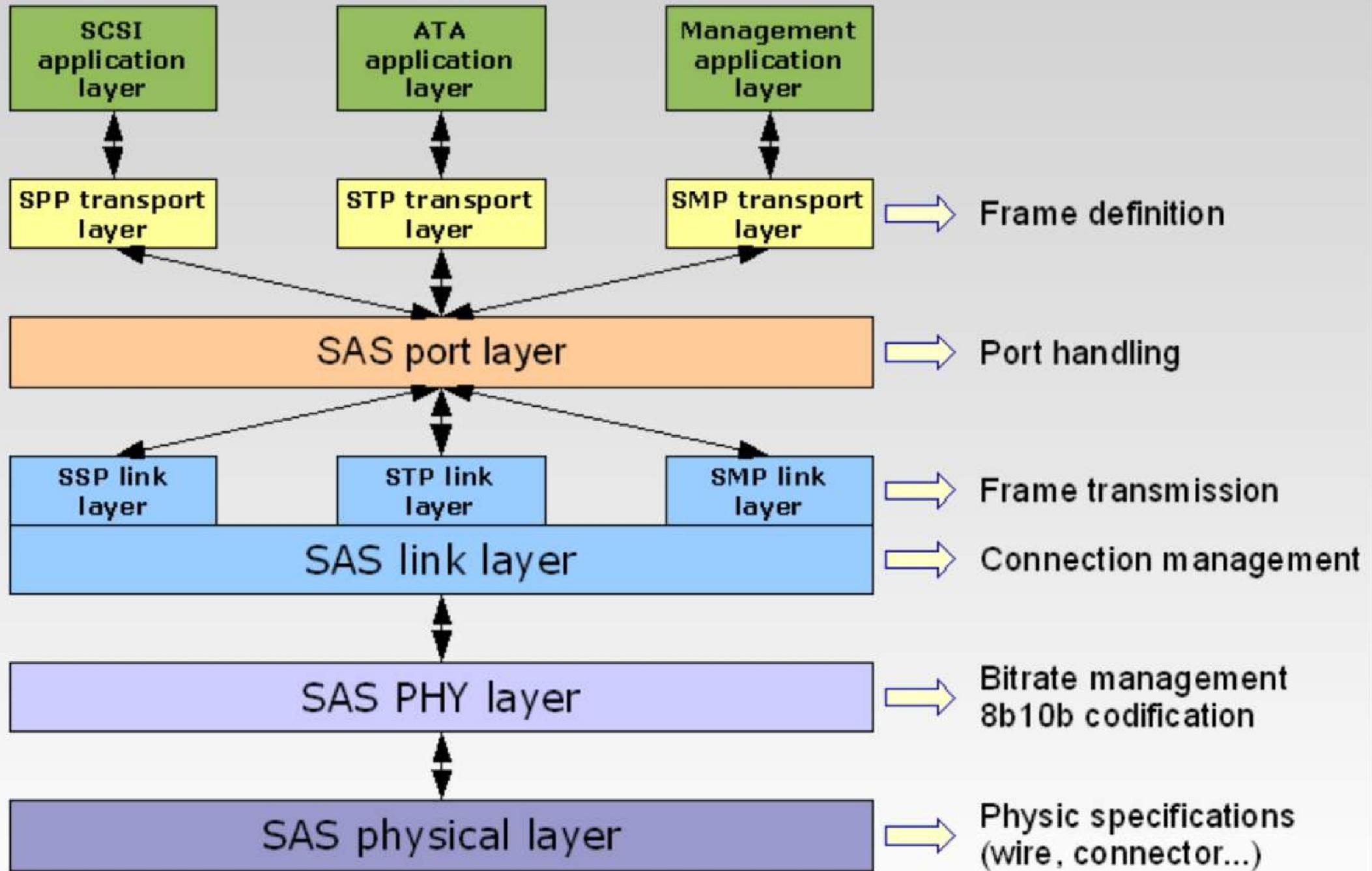


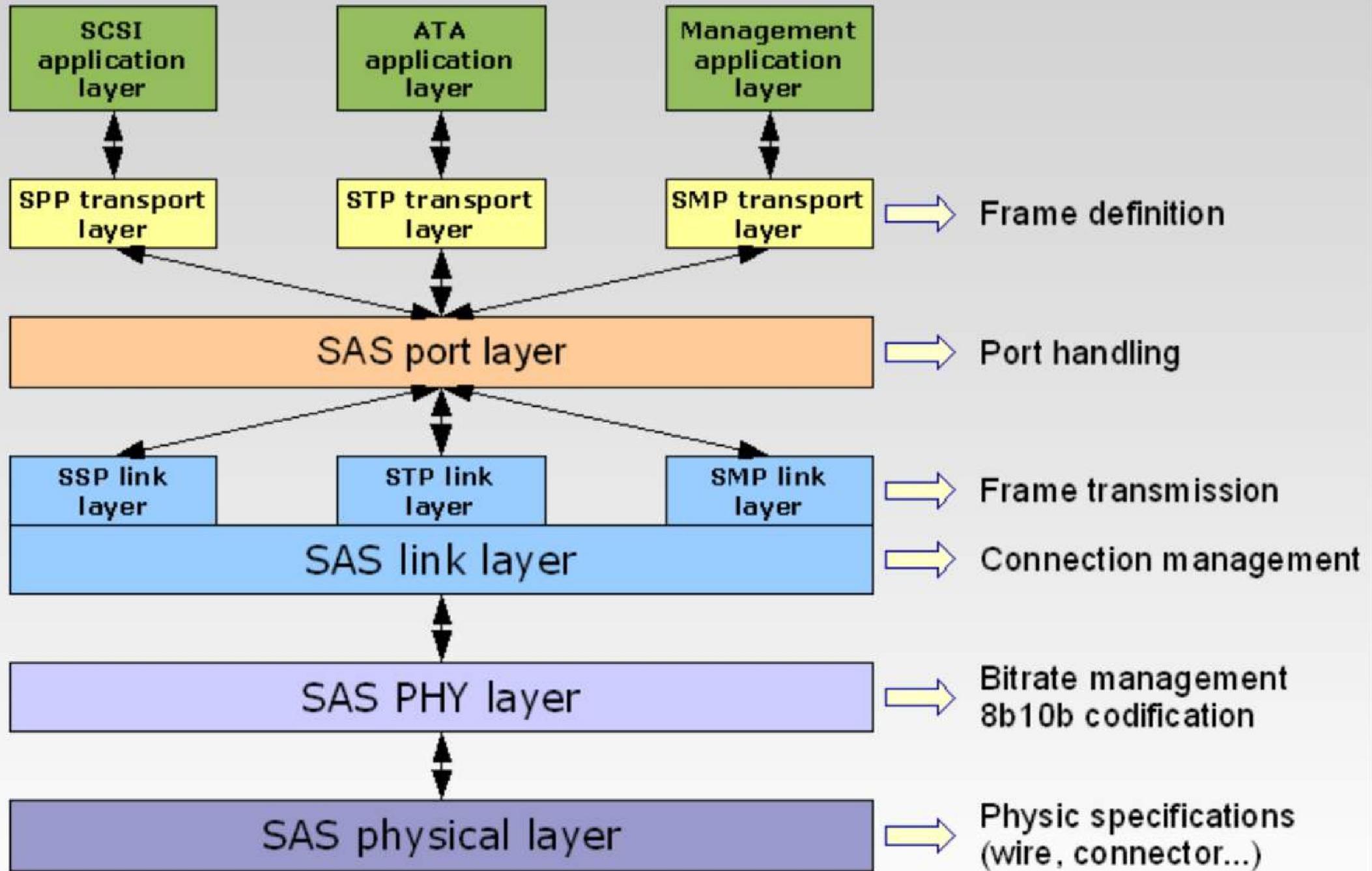
ed bus;
es;



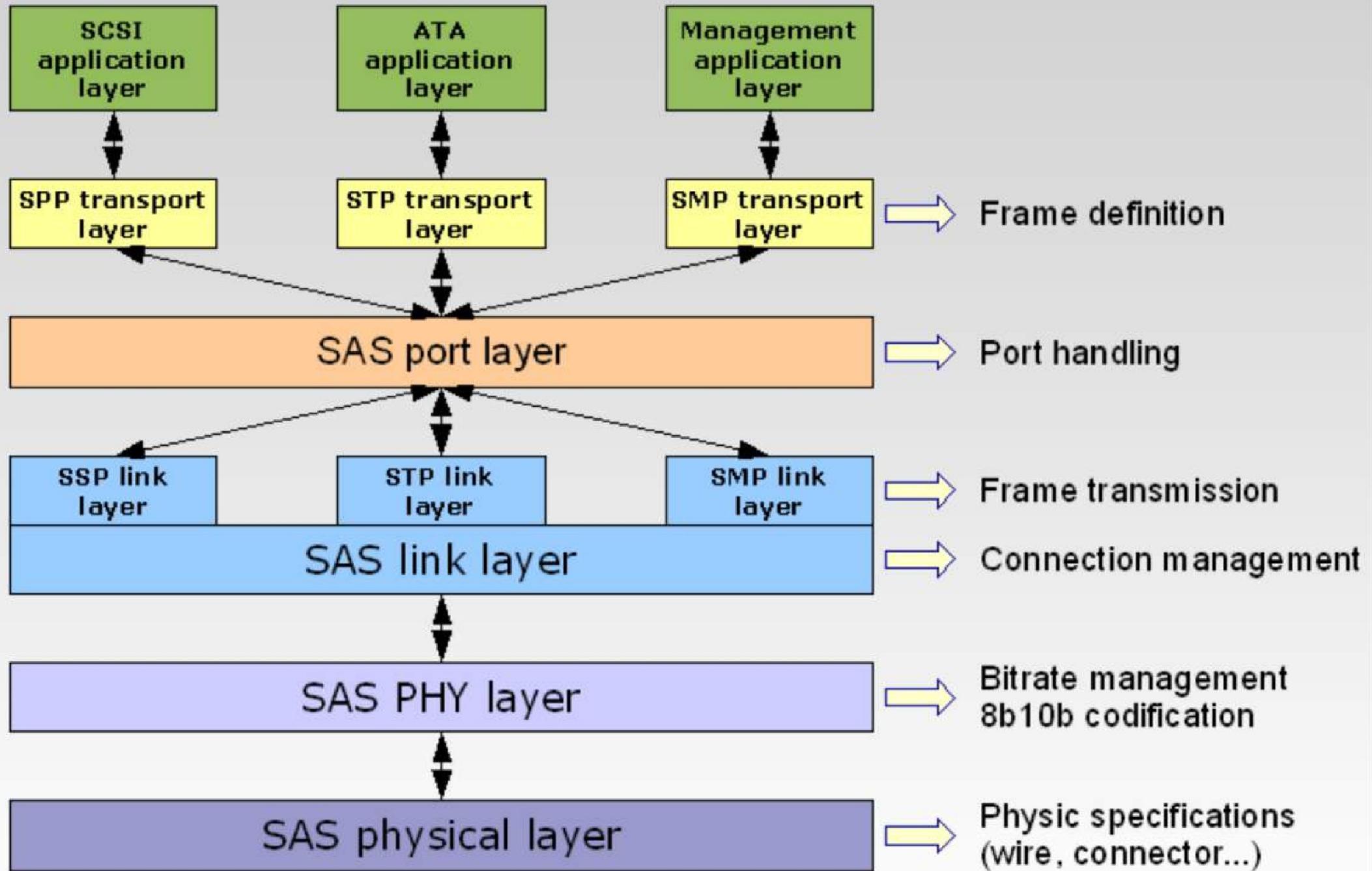


T1	Link	Address Frame Type	Protocol	Initiator Port (H)	Connection Rate
1 707 213 (ms)	1	0x1 : Open	0x1 : SSP	1	0x9 : 3.0 Gbps
Features (H)		Initiator Connection Tag (H)		Destination SAS Address (H)	
0		FFFF		5000C50000103F91	
Compatible Features (H)		Pathway Blocked Count (H)		Arbitration Wait Time (H)	
00		00		03E8	
More Compatible Features (H)				00000000	
CRC (H)	Link Data (H) ↴	Relative Time	Duration		
383F2C91		0 (ns)	133 (ns)		
T1	Link	3 G	Link Data (H) ↵	Relative Time	Duration
1 707 506 (ms)	2		Data	293 (ns)	53 (ns)
Target		RD			
OPEN ACCEPT		+---			
R_RDY NORMAL (x2)		-----			
DONE NORMAL		-----			





Byte	Field(s)				
0	Frame Type				
1 to 3	Hashed Destination SAS address				
4	Reserved				
5 to 7	Hashed Source SAS address				
8 to 9	Reserved				
10	Reserved	Retransmit	Rsvd		
11	Reserved	Number of Fill Bytes			
12 to 15	Reserved				
16 to 17	Tag				
18 to 19	Target Port Transfer Tag				
20 to 23	Data Offset				
24 to m	Information Unit				
m to (n-3)	Fill bytes, if needed				
(n-3) to n	CRC				



Serial SCSI Protocol:

- Wraps SCSI commands into serialized frames;
- Read Command, Data from target, Response;
- Write Command, XFER_RDY, Data to Target, Response;
- Nondata, Command, Response.

SATA Transport Protocol:

- Wraps SATA commands;
- Tunneling;
- Uses underlying SAS layers for transport;
- FIS.

- narrow and wide ports;

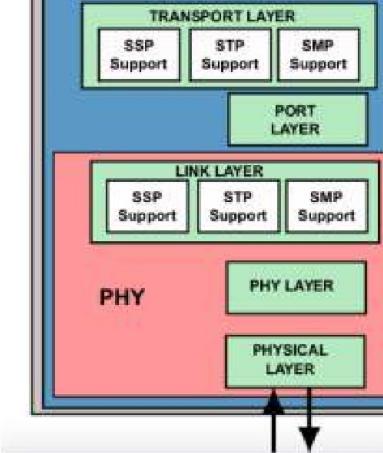
SAS and SCSI:

- P2P vs multidrop bus;
- no termination vs terminated bus;
- 65535 devices vs 8/16 devices;
- dedicated full bandwidth vs. shared bus bandwidth;

SAS and SATA:

• SAS is faster than SATA

- dedicated **full bandwidth** vs.
shared bus bandwidth;



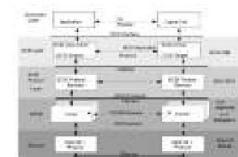
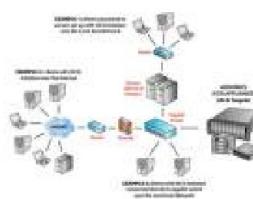
SAS and **SATA**:

- **multiple initiators** vs **single initiator**;
- **TCQ** vs **NCQ**;
- **SCSI command set** vs **ATA command set**;
- **multipath** IO vs **single path** IO;
- **1.6V** signaling vs **0.6V** signaling;
- up to **10 meters** vs **1 meter** cabling;

<https://www.also.com/pub/pdf/hp>
<https://www.mindshare.com/files>

iSCSI Protocol

- Internet SCSI;
- **Block level access over TCP/IP network;**
- Ports 860/3260;
- Purposes:
 - storage consolidation;
 - disaster recovery;
 - iSCSI **initiator** - iSCSI **target**;
- Types:
 - hardware;
 - dedicated;
 - offload engines;
 - software solutions;
- RFC 3720;



Addressing:

- iSCSI Qualified Name;
- iqn-yyyy-mm-reversed_domain:opt
- iqn-2019-04-hu.uni-obuda:storage1

Security:

- Challenge Handshake Authentication Protocol;
- filtering;
- multipathing;
- logical vs physical isolation;
- SSL over IP;
- **no zoning!**

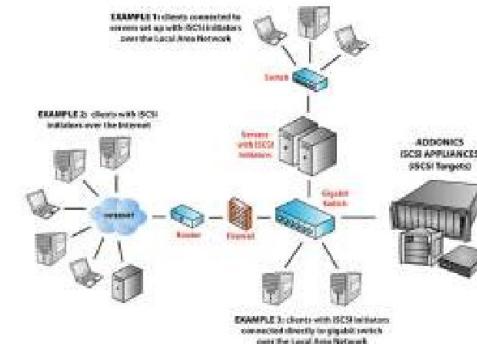
Protocol extensions:

- **MC/S:** Allows multiple TCP/IP connections to set up an iSCSI session;
- Needs additional flow control in the protocol stack!
- Data layer + Basic Header Segments;

Seq	A	B	C	D	E
0	DATA				
1		DATA			
2			DATA		
3				DATA	
4					DATA
5					

https://www.snia.org/sites/default/education/tutorials/2011/spring/networking/HufferdJohn-IP_Storage_Proocols-iSCSI.pdf

- Internet SCSI;
- **Block level access over TCP/IP network;**
- Ports 860/3260;
- Purposes:
 - storage consolidation;
 - disaster recovery;
 - iSCSI **initiator** - iSCSI **target**;
- Types:
 - hardware;
 - dedicated;
 - offload engines;
 - software solutions;
- RFC 3720;



Addressing

- ISCS
- iqn-yy
- iqn-20

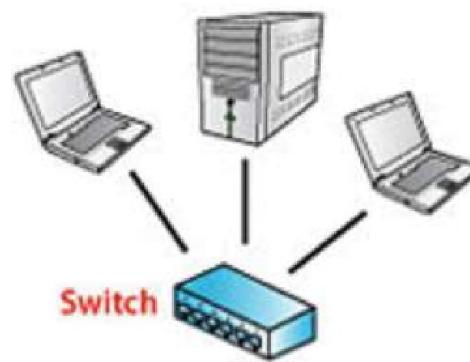
Security

- Challenge
- filtering
- multipath
- logical
- SSL
- **no zoning**

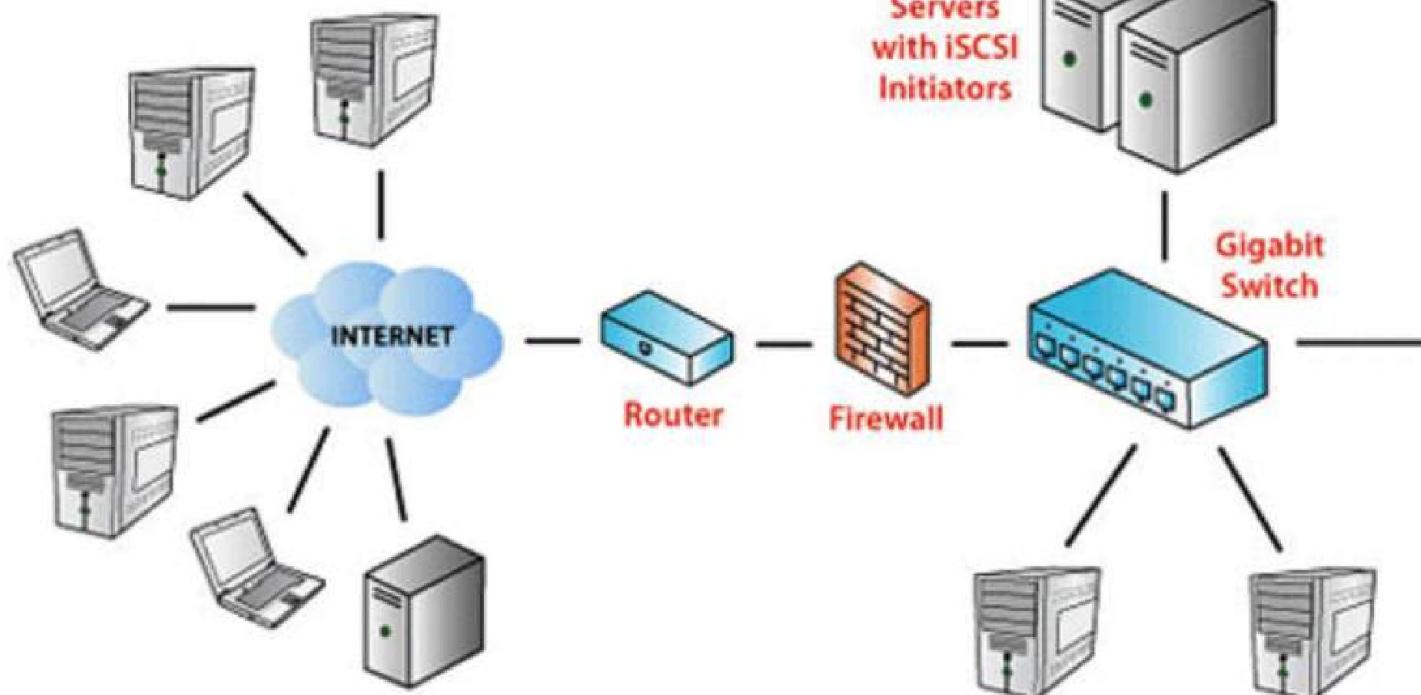
Protocol

- **MC/S**
- connection
- iSCSI
- Need
- in the

EXAMPLE 1: clients connected to servers set up with iSCSI initiators over the Local Area Network



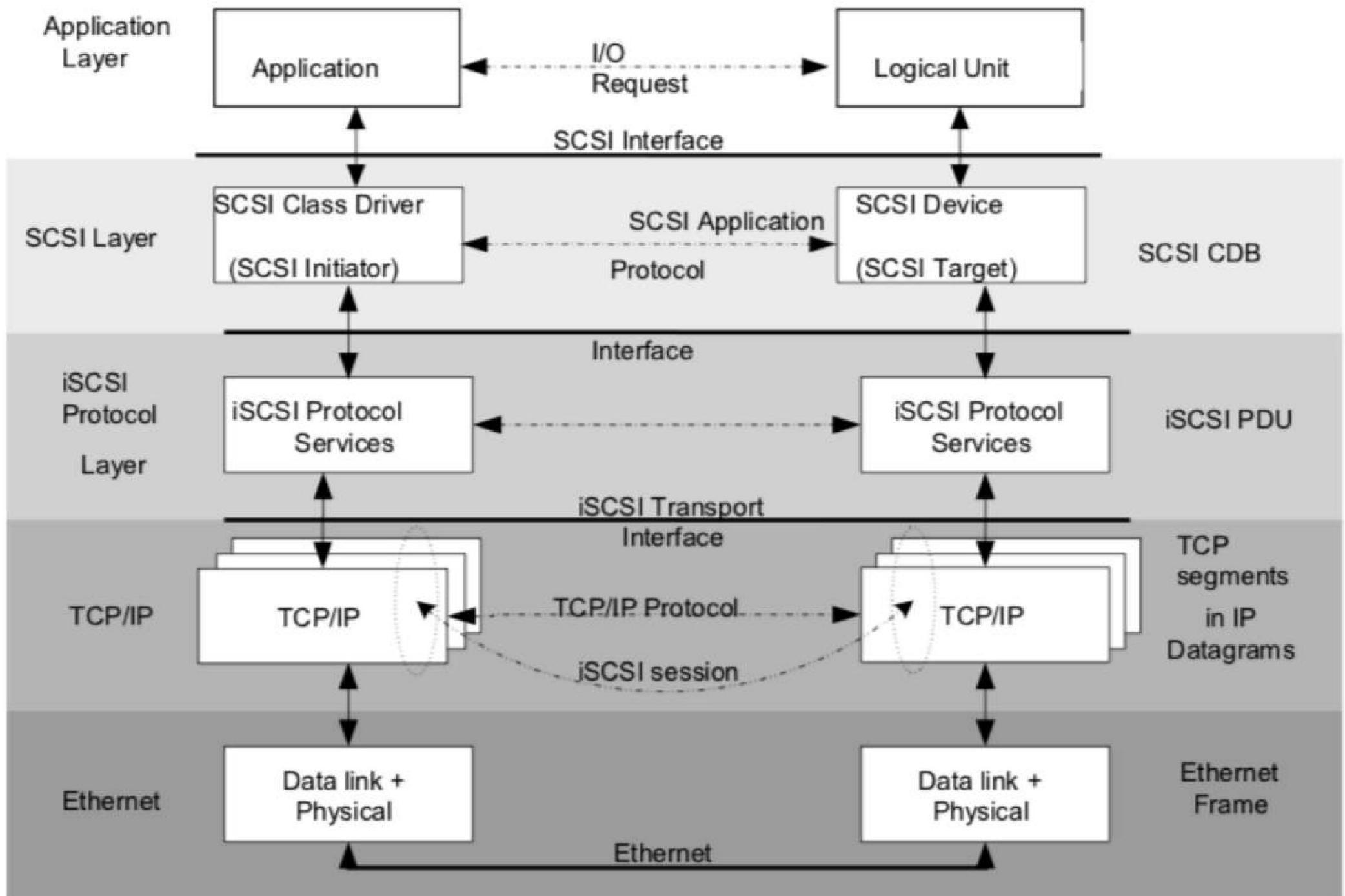
EXAMPLE 2: clients with iSCSI initiators over the Internet



**ADDONICS
iSCSI APPLIANCES
(iSCSI Targets)**



EXAMPLE 3: clients with iSCSI initiators connected directly to gigabit switch over the Local Area Network



Addressing:

- ISCSI Qualified Name;
- iqn-yyyy-mm-reversed_domain:opt
- iqn-**2019-04-hu.uni-obuda:storage1**

Security:

- Challenge Handshake Authentication Proto
filtering:

- iqn-yyyy-mm-reversed_domain:opt
- iqn-**2019-04-hu.uni-obuda:storage1**

Security:

- Challenge Handshake Authentication Protocol;
- filtering;
- multipathing;
- logical vs physical isolation;
- SSL over IP;
- **no zoning!**

Protocol extensions:

- **MC/S**: Allows multiple TCP/IP connections to set up an

Protocol extensions:

- **MC/S**: Allows multiple TCP/IP connections to set up an iSCSI session;
- Needs additional flow control in the protocol stack!
- Data layer + Basic Header Segments;

Byte	0	1	2	3
Bit	0 1 2 3 4 5 6 7	0 1 2 3 4 5 6 7	0 1 2 3 4 5 6 7	0 1 2 3 4 5 6 7
0 to 3	. I Opcode	Opcode-Specific Fields		
4 to 7	Total AHS Length	Data Segment Length		
8 to 15	Logical Unit Number (LUN) or Opcode-Specific Fields			
16 to 19	Initiator Task Tag (ITT) or Opcode-Specific Fields			
20 to 31	Opcode-Specific Fields			
32 to 47	Command Descriptor Block (CDB) or Opcode-Specific Fields			

Infiniband



- A storage data network optimized for high bandwidth and low latency;
- Infiniband Trade Association including major companies;
- SDR [2Gbps, 5μs], DDR [4Gbps, 2.5μs], QDR [8Gbps, 1.3μs], FDR [14Gbps, 700ns], EDR [24Gbps, 500ns], HDR [48Gbps, <500ns];
- Ports can be **trunked to multiply bandwidth**: 1x, 4x, 12x;
- **Non-blocking** switching;
- Switched network, with **Host Channel Adapter**, and **Target Channel Adapter**;
- Physical connections: CX4, SFP, CXP;
- **4k** packets;

Messages:

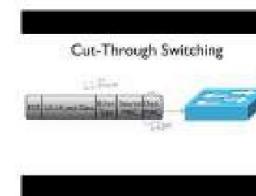
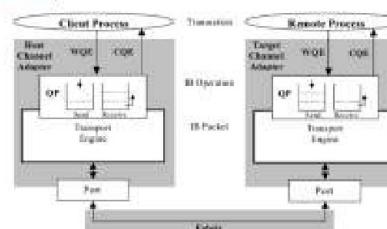
- RDMA;
- Channel send and receive;
- Transaction;
- Multicast message;
- Atomic operation.

Layers:

- **Physical**: cables + interfaces;
- **Link**: Local IDs, Virtual Lanes (VL0-15);
- **Network**: Global Route Headers between subnets, IPv6 headers are used;
- **Transport**: in-order delivery, partitioning, channel multiplexing;

Virtual Interface Architecture:

- distributed messaging technology;
- it offloads traffic control from the client to dedicated execution queues;
- **WQP** is assigned to the transmission and reception side;
- Work Queue Entries -> Completion Queue Entries;



Why IB is faster?

- **Store and forward** switch: stores entire packets and makes decision then;
- **Cut through** switch: stores only headers, not full data frames -> **low latency**;

https://en.wikipedia.org/wiki/Cut-through_switching

https://www.mellanox.com/pdf/whitepapers/IB_Intro_WP_190.pdf

<https://www.arista.com/assets/data/pdf/Infiniband-vs-etherne.pdf>

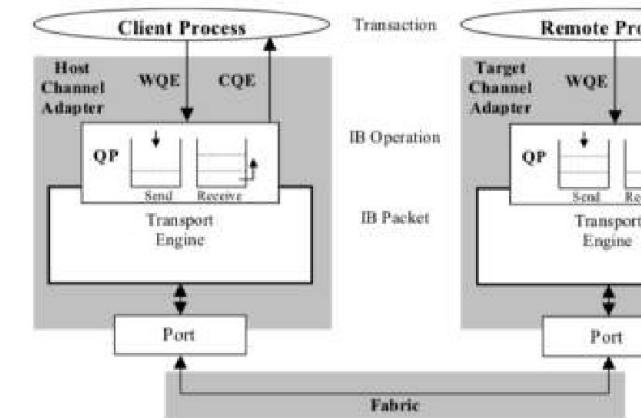
Infiniband

- A storage data network optimized for high bandwidth and low latency;
- Infiniband Trade Association including major companies;
- SDR [2Gbps, 5μs], DDR [4Gbps, 2.5μs], QDR [8Gbps, 1.3μs], FDR [14Gbps, 700ns], EDR [24Gbps, 500ns], HDR [48Gbps, <500ns];
- Ports can be **trunked to multiply bandwidth**: 1x, 4x, 12x;
- **Non-blocking** switching;
- Switched network, with **Host Channel Adapter**, and **Target Channel Adapter**;
- Physical connections: CX4, SFP, CXP;
- **4k** packets;

Messages:

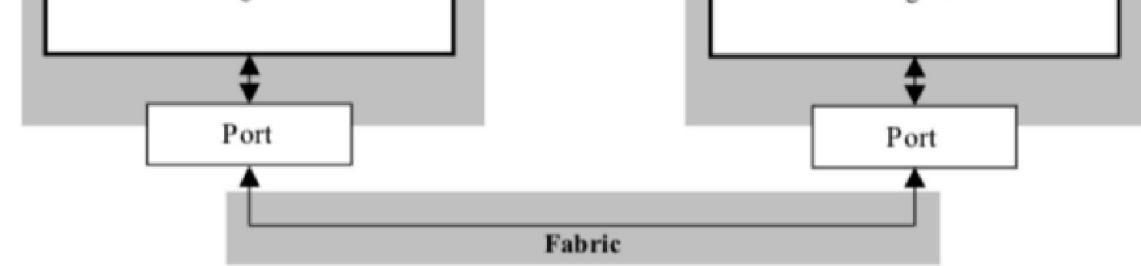
- RDMA;

Layers:



Messages:

- RDMA;
- Channel send and receive;
- Transaction;
- Multicast message;
- Atomic operation.



Layers:

- **Physical:** cables + interfaces;
- **Link:** Local IDs, Virtual Lanes (VL0-15);
- **Network:** Global Route Headers between subnets, IPv6 headers are used;
- **Transport:** in-order delivery, partitioning, channel multiplexing;

Virtual Interface Architecture:

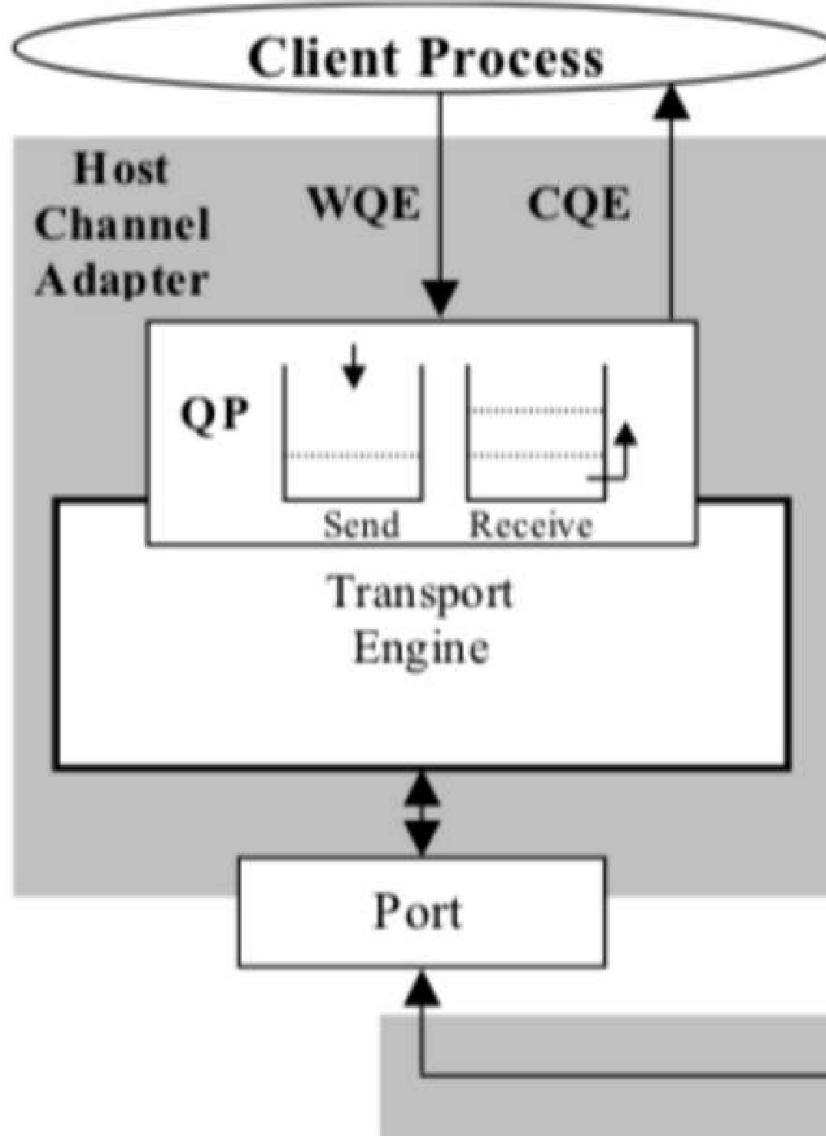
- distributed messaging technology;
- it offloads traffic control from the client to dedicated execution queues;
- **WQP** is assigned to the transmission and reception side;
- Work Queue Entries -> Completion Queue Entries;

Why IB is faster?

- **Store and forward** switch:
stores entire packets and makes decision then;
- **Cut through** switch: stores only headers, not full data frames -> **low latency**;

https://en.wikipedia.org/wiki/Cut-through_switching

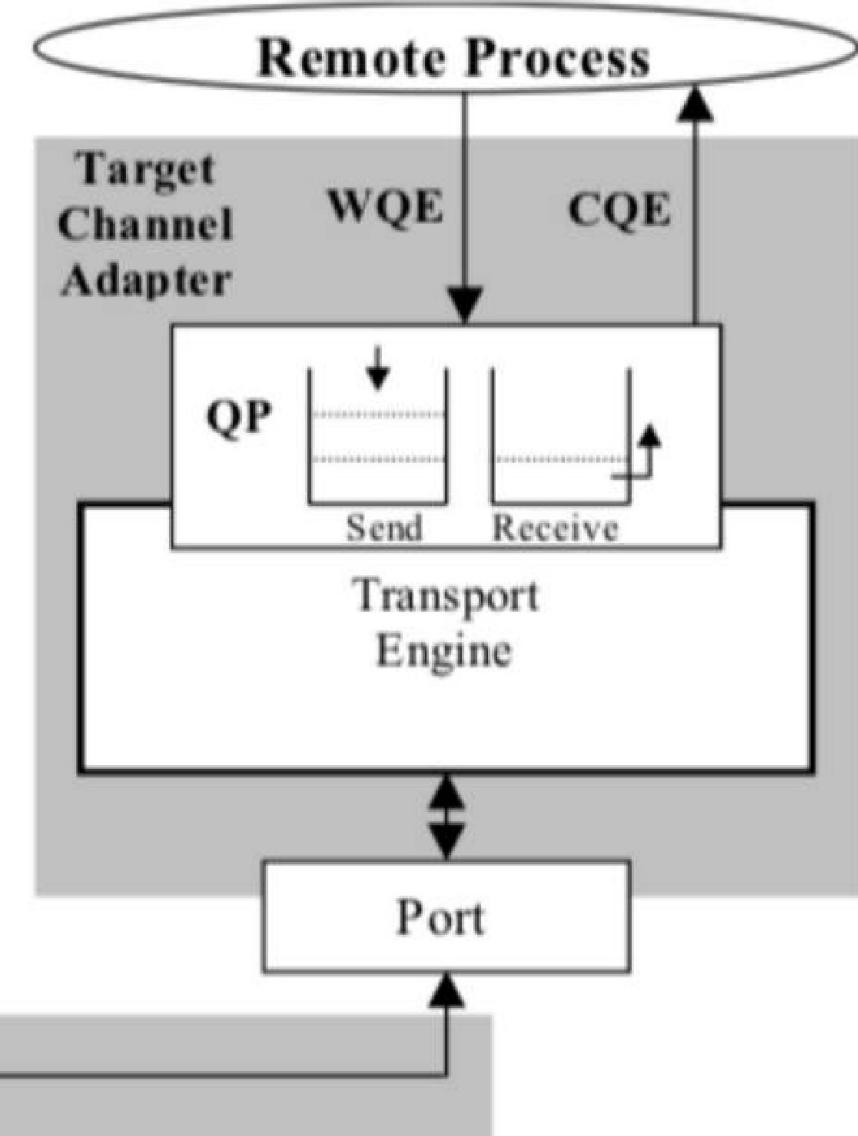
apter , and



Transaction

IB Operation

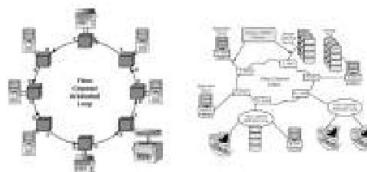
IB Packet



Fabric

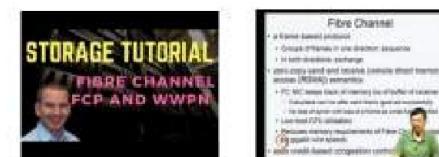
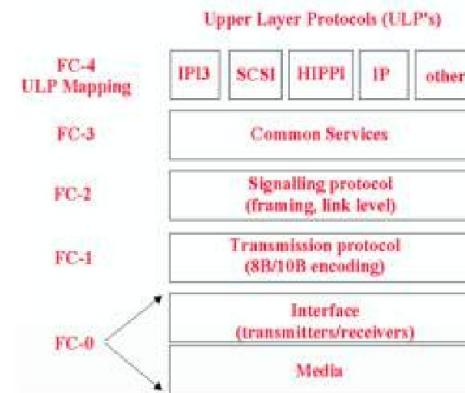
Fibre Channel Protocol

- Delivering **SCSI command set** over FC network.
- Besides SCSI: HIPPI, ATM, IP, NVMe.
- Bandwidth: 1-32 Gb/s. We have 4 and 8 Gb/s.
- **Addressing:**
 - 3 byte addresses are used
 - FFFFFA - fabric manager service
 - FFFFFC - fabric name service
 - FFFFFD - fabric controller service
 - FFFFFE - fabric login service
- **Principal switch:** responsible for managing and distributing IDs within the domain.
- **Flow control:** throttling traffic.
- **Topologies:**
 - FC-PP: point-to-point;
 - FC-AL: arbitrated loop, ring;
 - FC-SW: switched fabric.



Protocol exchanges:

- **Fabric Login (FLOGI):** node enters the fabric;
- **Port Login (PLOGI):** session between initiator port and target port;
- **Process login (PRI):** sending SCSI commands.

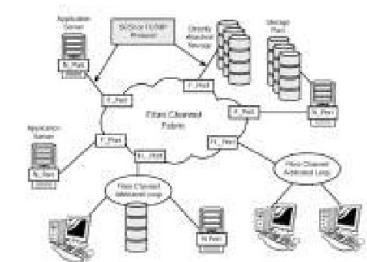
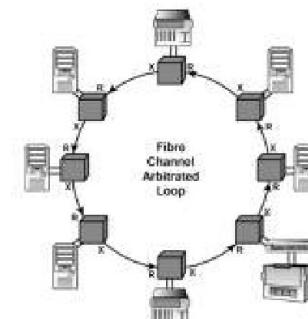


Frame structure:

- Exchange: a session;
- Sequence: a set of frames;
- Frame:
 - data;
 - control (e.g. FSPF).

<http://www.tsmtutorials.com/2016/08/fc-san-protocols.html>
https://en.wikipedia.org/wiki/Fibre_Channel

- Delivering **SCSI command set** over FC network.
- Besides SCSI: HIPPI, ATM, IP, NVMe.
- Bandwidth: **1-32 Gb/s**. We have 4 and 8 Gb/s.
- **Addressing:**
 - 3 byte addresses are used
 - FFFFFA - fabric manager service
 - FFFFFC - fabric name service
 - FFFFFD - fabric controller service
 - FFFFFE - fabric login service
- **Principal switch:** responsible for managing and distributing IDs within the domain.
- **Flow control:** throttling traffic.
- **Topologies:**
 - FC-PP: point-to-point;
 - FC-AL: arbitrated loop, ring;
 - FC-SW: switched fabric.



Domain ID

Bits (23-16)

Area ID

Bits (15-08)

Port ID

Bits (07-00)

Upper Layer Protocols (ULP's)

FC-4
ULP Mapping



FC-3

Common Services

FC-2

Signalling protocol
(framing, link level)

FC-1

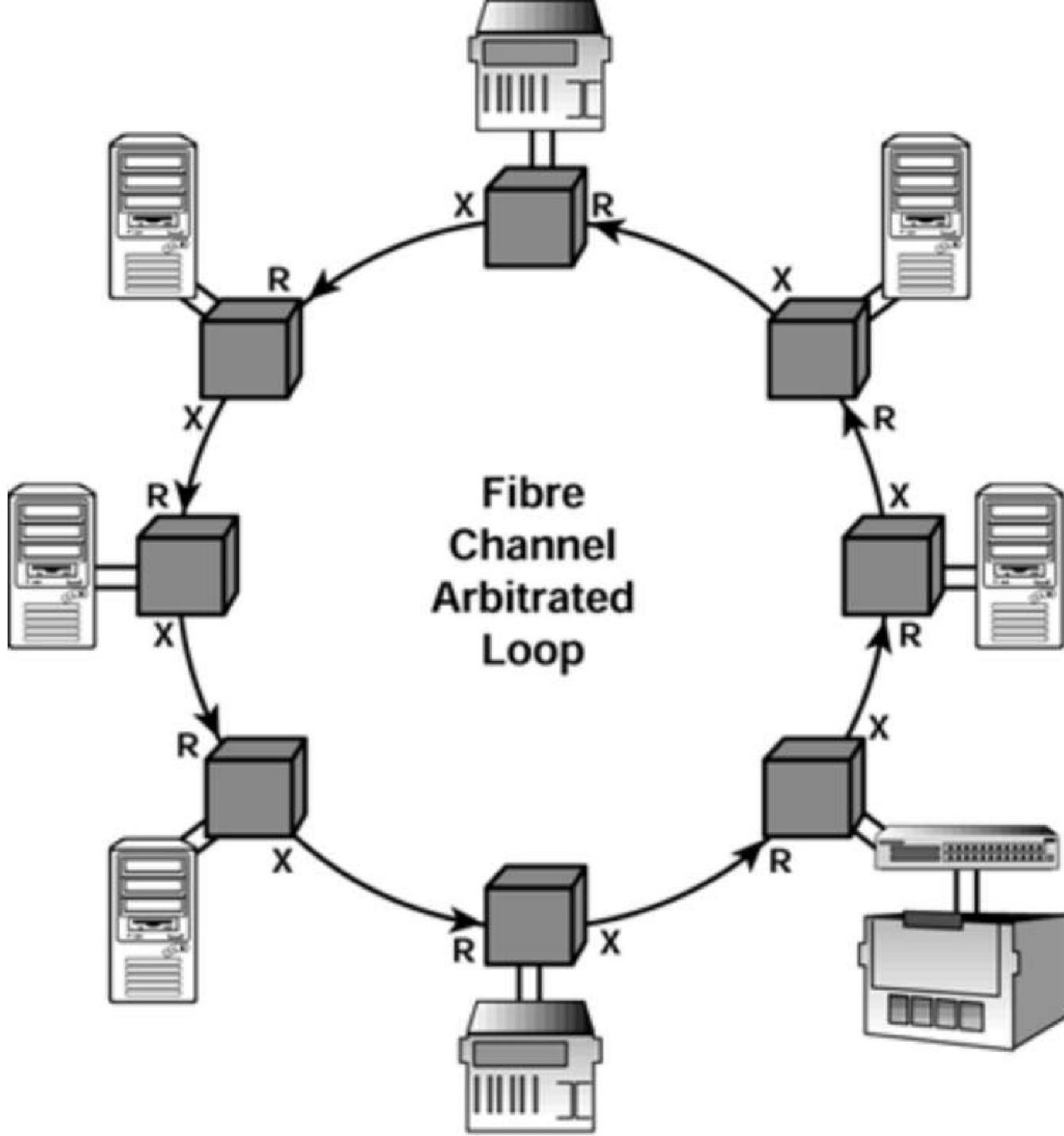
Transmission protocol
(8B/10B encoding)

FC-0

Interface
(transmitters/receivers)

Media



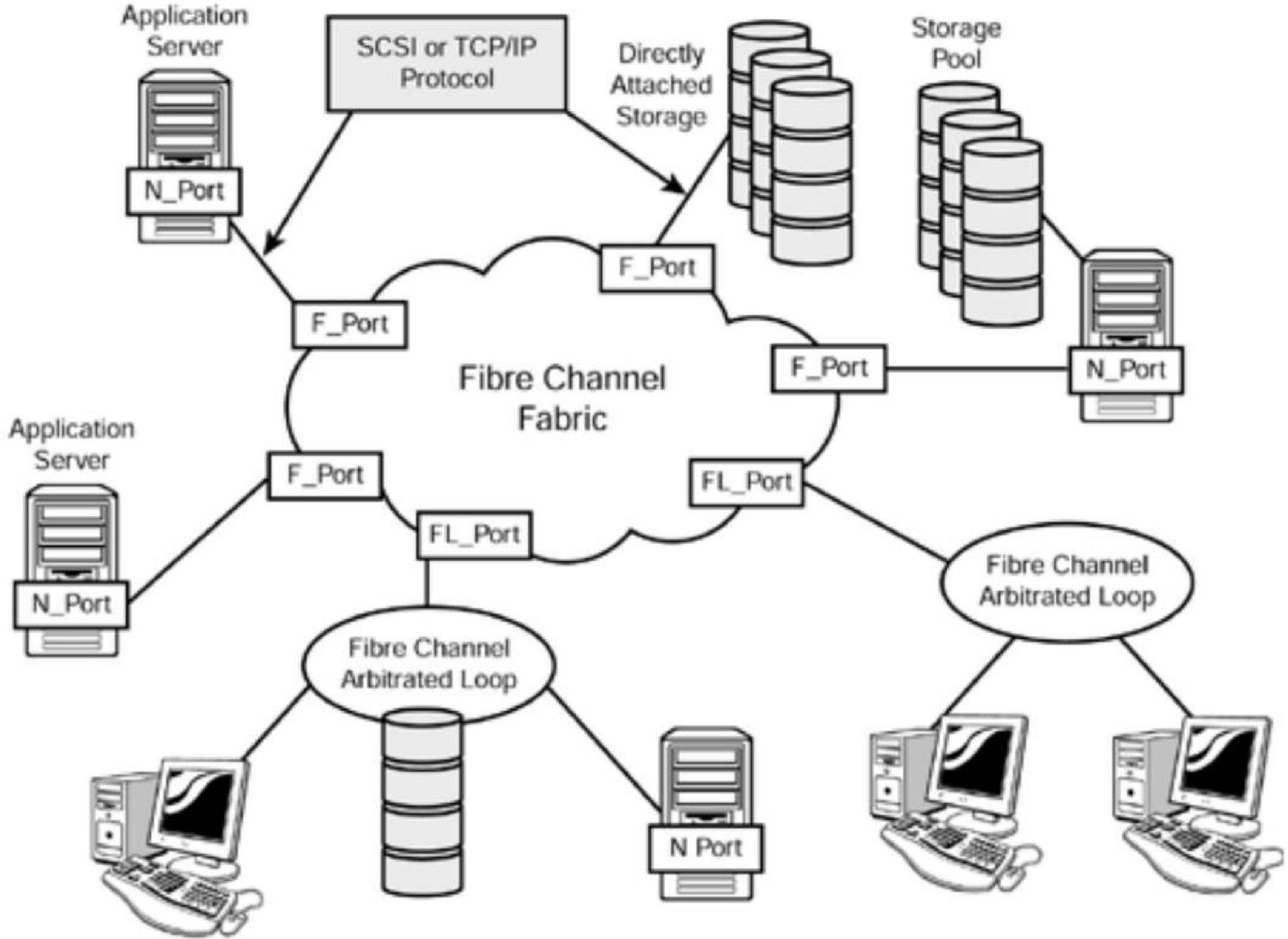


Application Server



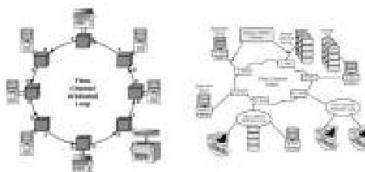
N_Port





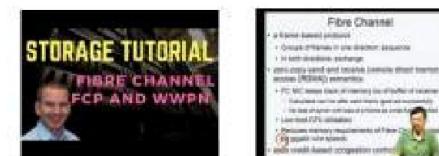
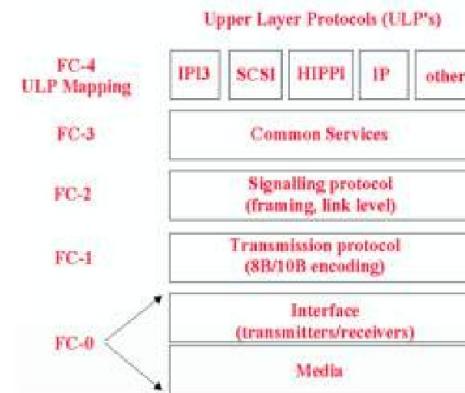
Fibre Channel Protocol

- Delivering **SCSI command set** over FC network.
- Besides SCSI: HIPPI, ATM, IP, NVMe.
- Bandwidth: 1-32 Gb/s. We have 4 and 8 Gb/s.
- **Addressing:**
 - 3 byte addresses are used
 - FFFFFA - fabric manager service
 - FFFFFC - fabric name service
 - FFFFFD - fabric controller service
 - FFFFFE - fabric login service
- **Principal switch:** responsible for managing and distributing IDs within the domain.
- **Flow control:** throttling traffic.
- **Topologies:**
 - FC-PP: point-to-point;
 - FC-AL: arbitrated loop, ring;
 - FC-SW: switched fabric.



Protocol exchanges:

- **Fabric Login (FLOGI):** node enters the fabric;
- **Port Login (PLOGI):** session between initiator port and target port;
- **Process login (PRI):** sending SCSI commands.



Frame structure:

- Exchange: a session;
- Sequence: a set of frames;
- Frame:
 - data;
 - control (e.g. FSPF).

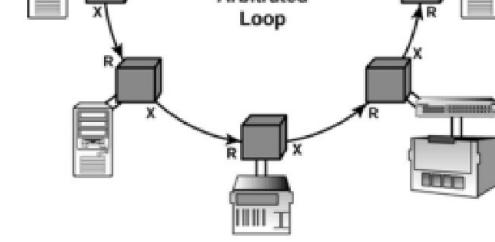
<http://www.tsmtutorials.com/2016/08/fc-san-protocols.html>
https://en.wikipedia.org/wiki/Fibre_Channel



Frame structure:

- Exchange: a session;
- Sequence: a set of frames;
- Frame:
 - data;
 - control (e.g. FSPF).

- FC-Sw: switched fabric.



Protocol exchanges:

- **Fabric Login (FLOGI)**: node enters the fabric;
- **Port Login (PLOGI)**: session between initiator port and target port;
- **Process login (PRI)**: sending SCSI commands.

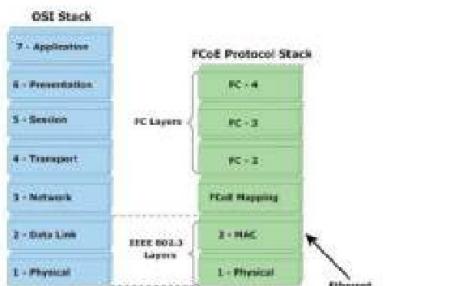
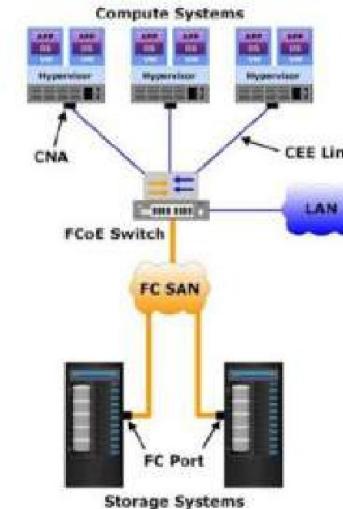
Do
Bits

<http://www.tsmtutorials.com/>

<https://en.wikipedia.org/wiki/>

FC over Ethernet Protocol

- Encapsulates **FC data into Ethernet frames**;
- Uses **multi-functional HBAs** and **switches**;
- Adapters, cables, switches:
 - **Converged Network Adapter**: combines Ethernet and FC functionality in the same device;
 - **Software FCoE Adapter**: uses kernel modules to process FCoE traffic;
- **Fibre Channel Forwarder**:
 - Resides in a combined Ethernet and FC switch;
 - Provides FC services: zoning, name resolving, etc.
- The same port types: **VN_Port**, **VF_Port**, **VE_Port**.



FCoE addressing:

- MAC addresses are used;
- jumbo frames: 2112 bytes vs. 1500 bytes;
- **SPMA**: compute servers provide MAC addresses;
- **FPMA**: fabrics provide MAC addresses (MAC:FC addresses concatenated);
- **Discovery (FIP) - Login - Data forwarding**.

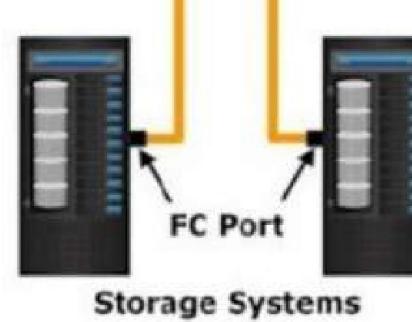
Discovery and Login:

- Multicast Frame Initializing
- Packet messages to find FCF;
- FCFs send FIP Advertisement frames;
- Nodes send FLOGI requests to the closest FCF;

Practice #7: FCoE and iSCSI target creation



ed Ethernet and FC switch;
s: zoning, name resolving, etc.
N_Port, VF_Port, VE_Port.



FCoE addressing:

- MAC addresses are used;
 - jumbo frames: 2112 bytes vs. 1500 bytes;
 - **SPMA**: compute servers provide MAC addresses;
 - **FPMA**: fabrics provide MAC addresses (MAC:FC addresses concatenated);
 - **Discovery (FIP)** - **Login** - **Data forwarding**.

Discovery and Login:

- Multicast Frame Initializing
Packet messages to find FCF;
 - FCFs send FIP Advertisement frames;
 - Nodes send FLOGI requests to the closest FCF;

Practice #7: FCoE and iSCSI target creation

Discovery and Login:

- Multicast Frame Initializing
Packet messages to find FCF;
- FCFs send FIP Advertisement frames;
- Nodes send FLOGI requests to the closest FCF;

Volume management

create flexible volumes and mappings on nodes

Logical Volume Management

a more flexible way of allocating block device resources than in traditional way, i.e. partitioning

- Divides physical space to chunks called Physical Extents (PEs);
- Maps PEs to Logical Extents (LEs);
- Elementary terms:
 - Physical Volume(s) (PV): a set of devices formulating the physical space;
 - Volume Group (VG): a pool of resources sum of all PEs serving the same purpose;
 - Logical Volume(s) (LV): a set of dynamically allocated LEs;
- PE + LE have fixed size - 4MB;
- Structure: first Megabyte of each PV shows the same LVM metadata;

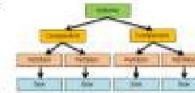


- Operations on LVs:
 - Concatenation of PVs;
 - Mirroring map multiple PEs to a single LE;
 - Growing and shrinking LVs;
 - Growing VGs by adding PVs;
 - Creating snapshots (copy on write LE mapping);
 - Hybrid volumes:
 - beache, dm-cache, Fusion Drive, etc.

https://en.wikipedia.org/wiki/Logical_volume_management
[https://en.wikipedia.org/wiki/Logical_Volume_Manager_\(Linux\)](https://en.wikipedia.org/wiki/Logical_Volume_Manager_(Linux))
<http://www.cyberphoton.com/questions/question/what-is-the-difference-between-lvm-and-raid>

Logical Disk Manager

- Implements Logical Volume Management in MS Windows environment;
- By Veritas and MS;
- Volume types:**
 - Basic volumes: data storage is limited to a physical disk + no partitioning;
 - Dynamic volumes: flexible volumes on the same physical disk, or on multiple disks:
 - Partitionable;
 - Striped volumes;
 - SPAN volumes;
- Limited to **32 dynamic volumes**;
- Uses special partition tables + special partition bounds, i.e. 1MB bound.



https://en.wikipedia.org/wiki/Logical_Disk_Manager

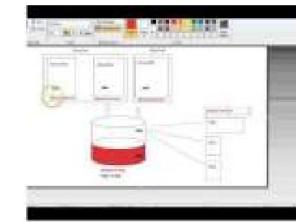
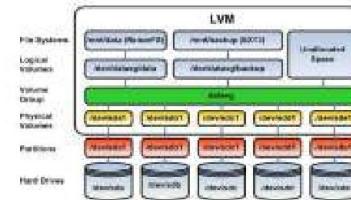
Logical Volume Management

a more flexible way of allocating block device resources than in traditional way, i.e. partitioning

- Divides physical space to chunks called Physical Extents (**PE**);
- Map PEs to Logical Extents (**LE**);
- Elementary terms:
 - **Physical Volume(s) (PV)**: a set of devices formulating the physical space;
 - **Volume Group (VG)**: a pool of resources, sum of all PEs serving the same purpose;
 - **Logical Volume(s) (LV)**: a set of dynamically allocated LEs;
- PE + LE have fixed size: 4MB.
- Structure: first MegaByte of each PVs stores the same LVM metadata;

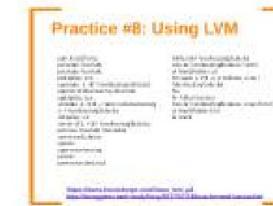
Operations on LVs:

- Concatenation of PVs;
- Mirroring: map multiple PEs to a single LE;
- **Growing and shrinking LVs**;
- Growing VGs by adding PVs;
- creating **snapshots** (copy on write LE mapping);
- Hybrid volumes:
 - similar to hierarchical storage systems;
 - bcache, dm-cache, Fusion Drive, etc.



LE allocation policies:

- contiguous: LEs are adjacent and ordered;
- cling: LEs are on same physical devices;
- normal: indiscriminate allocation;
- anywhere: random allocation.



- Originally LVM is planned for **on host use**;
- **Clustered LVM**:
 - PVs are located on different hosts;
 - A shared and lock-managed metadata service is used to store shared metadata;
- **CLVM**: metadata communication flows through the lock manager;
- **HA-LVM**: uses the default file locking mechanism, and avoids metadata contention, no concurrent accesses, useful in master-slave configs.

https://en.wikipedia.org/wiki/Logical_volume_management

[https://en.wikipedia.org/wiki/Logical_Volume_Manager_\(Linux\)](https://en.wikipedia.org/wiki/Logical_Volume_Manager_(Linux))

<http://www.cyberphoton.com/questions/question/what-is-the-difference-between-lvm-and-raid>

- Divides physical space to chunks called Physical Extents (**PE**);
- Map PEs to Logical Extents (**LE**);
- Elementary terms:
 - **Physical Volume(s) (PV)**: a set of devices formulating the physical space;
 - **Volume Group (VG)**: a pool of resources, sum of all PEs serving the same purpose;
 - **Logical Volume(s) (LV)**: a set of dynamically allocated LEs;
- PE + LE have fixed size: **4MB**.
- Structure: first MegaByte of each PVs stores the same LVM metadata;

LVM

File Systems

/mnt/data (ReiserFS)

/mnt/backup (EXT3)

Unallocated
Space

Logical
Volumes

/dev/datavg/data

/dev/datavg/backup

Volume
Group

datavg

Physical
Volumes

/dev/sda1

/dev/sdb1

/dev/sdc1

/dev/sdd1

/dev/sde1

Partitions

/dev/sda1

/dev/sdb1

/dev/sdc1

/dev/sdd1

/dev/sde1

Hard Drives

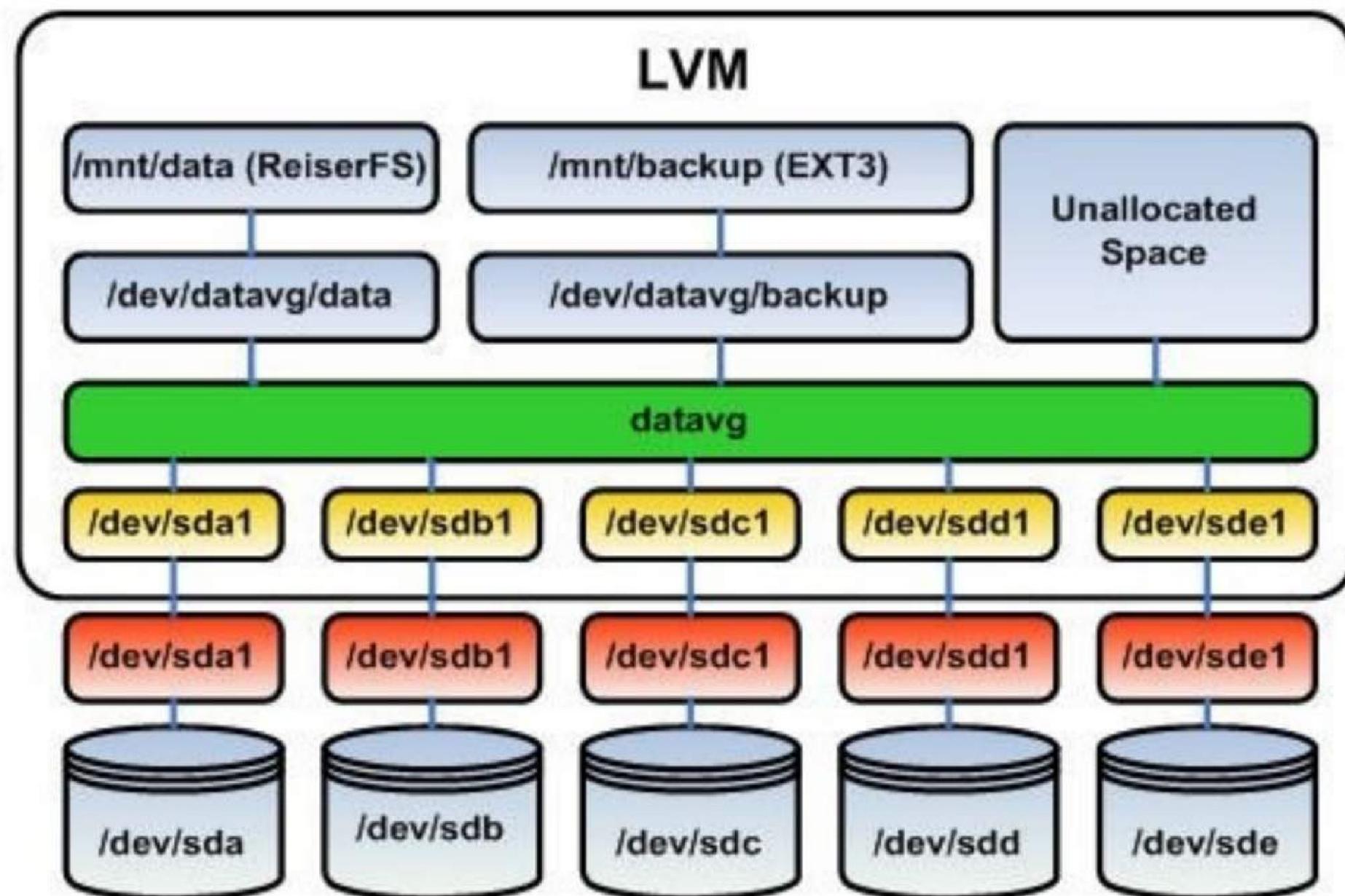
/dev/sda

/dev/sdb

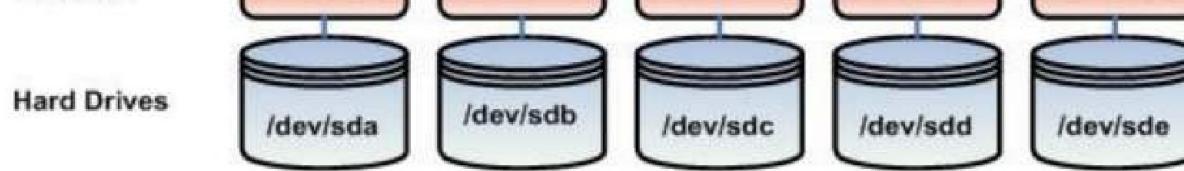
/dev/sdc

/dev/sdd

/dev/sde



- Divides physical space to chunks called Physical Extents (**PE**);
- Map PEs to Logical Extents (**LE**);
- Elementary terms:
 - **Physical Volume(s) (PV)**: a set of devices formulating the physical space;
 - **Volume Group (VG)**: a pool of resources, sum of all PEs serving the same purpose;
 - **Logical Volume(s) (LV)**: a set of dynamically allocated LEs;
- PE + LE have fixed size: **4MB**.
- Structure: first MegaByte of each PVs stores the same LVM metadata;



LE allocation policies:

- **contiguous**: LEs are adjacent and ordered;
- **cling**: LEs are on same physical devices;
- **normal**: indiscriminate allocation;
- **anywhere**: random allocation.

- Originally LVM is planned for

Clustered LVM

Pr

yum
pver
pvcr
pvdi
vgcr
vge
vgdi
lvcre
ls -l
lvdis
lvext
pvm
lvren
vgsc
vgre
pvsc
pvre

11

same LVM metadata;

Operations on LVs:

- Concatenation of PVs;
- Mirroring: map multiple PEs to a single LE;
- **Growing and shrinking LVs;**
- Growing VGs by adding PVs;
- creating **snapshots** (copy on write LE mapping);
- **Hybrid volumes:**
 - similar to hierarchical storage systems;
 - bcache, dm-cache, Fusion Drive, etc.

[https://en.wikipedia.org/wiki/Logical_Volume_Manager_\(Linux\)](https://en.wikipedia.org/wiki/Logical_Volume_Manager_(Linux))
[https://en.wikipedia.org/wiki/Logical_Volume_Manager_\(Windows\)](https://en.wikipedia.org/wiki/Logical_Volume_Manager_(Windows))

- Divides physical space to chunks called Physical Extents (**PE**);
- Map PEs to Logical Extents (**LE**);
- Elementary terms:
 - **Physical Volume(s) (PV)**: a set of devices formulating the physical space;
 - **Volume Group (VG)**: a pool of resources, sum of all PEs serving the same purpose;
 - **Logical Volume(s) (LV)**: a set of dynamically allocated LEs;
- PE + LE have fixed size: **4MB**.
- Structure: first MegaByte of each PVs stores the same LVM metadata;

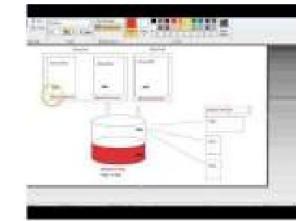
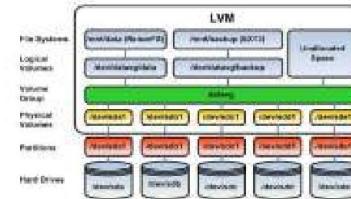
Logical Volume Management

a more flexible way of allocating block device resources than in traditional way, i.e. partitioning

- Divides physical space to chunks called Physical Extents (**PE**);
- Map PEs to Logical Extents (**LE**);
- Elementary terms:
 - **Physical Volume(s) (PV)**: a set of devices formulating the physical space;
 - **Volume Group (VG)**: a pool of resources, sum of all PEs serving the same purpose;
 - **Logical Volume(s) (LV)**: a set of dynamically allocated LEs;
- PE + LE have fixed size: 4MB.
- Structure: first MegaByte of each PVs stores the same LVM metadata;

Operations on LVs:

- Concatenation of PVs;
- Mirroring: map multiple PEs to a single LE;
- **Growing and shrinking LVs**;
- Growing VGs by adding PVs;
- creating **snapshots** (copy on write LE mapping);
- Hybrid volumes:
 - similar to hierarchical storage systems;
 - bcache, dm-cache, Fusion Drive, etc.



LE allocation policies:

- contiguous: LEs are adjacent and ordered;
- cling: LEs are on same physical devices;
- normal: indiscriminate allocation;
- anywhere: random allocation.



- Originally LVM is planned for **on host use**;
- **Clustered LVM**:
 - PVs are located on different hosts;
 - A shared and lock-managed metadata service is used to store shared metadata;
- **CLVM**: metadata communication flows through the lock manager;
- **HA-LVM**: uses the default file locking mechanism, and avoids metadata contention, no concurrent accesses, useful in master-slave configs.

https://en.wikipedia.org/wiki/Logical_volume_management

[https://en.wikipedia.org/wiki/Logical_Volume_Manager_\(Linux\)](https://en.wikipedia.org/wiki/Logical_Volume_Manager_(Linux))

<http://www.cyberphoton.com/questions/question/what-is-the-difference-between-lvm-and-raid>

ere: random allocation.

https://www.howtoforge.com/linux_lvm_p2
<http://strugglers.net/~andy/blog/2017/07/19/bcache-and-lvmpolice/>

- Originally LVM is planned for **on host use**;
- **Clustered LVM**:
 - PVs are located on different hosts;
 - A shared and lock-managed metadata service is used to store shared metadata;
- **CLVM**: metadata communication flows through the lock manager;
- **HA-LVM**: uses the default file locking mechanism, and avoids metadata contention, no concurrent accesses, useful in master-slave configs.

Volume management

Volume Manager (Linux)

Logical Disk Manager

- Implements Logical Volume Management in MS Windows environment;
- By Veritas and MS;
- **Volume types:**
 - **Basic** volumes: data storage is limited to a physical disk + no partitioning;
 - **Dynamic** volumes: flexible volumes on the same physical disk, or on multiple disks:
 - Partitionable;
 - Striped volumes;
 - SPAN volumes;
 - Limited to **32 dynamic volumes**;
- Uses special partition tables + special partition bounds, i.e. 1MB bound.

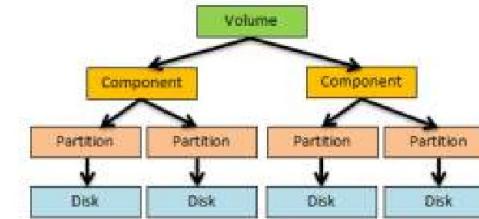
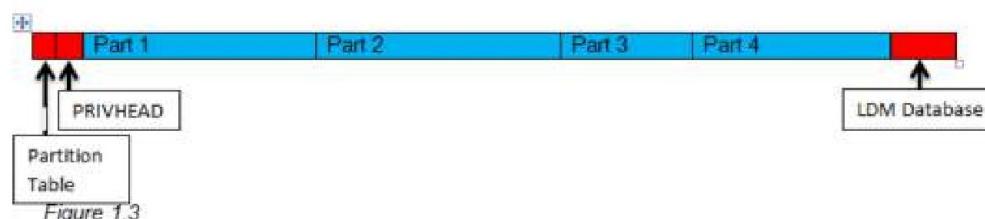


Figure 6.3



- Implements Logical Volume Management in MS Windows environment;
- By Veritas and MS;
- **Volume types:**
 - **Basic** volumes: data storage is limited to a physical disk + no partitioning;
 - **Dynamic** volumes: flexible volumes on the same physical disk, or on multiple disks:
 - Partitionable;
 - Striped volumes;
 - SPAN volumes;
 - Limited to **32 dynamic volumes**;
- Uses special partition tables + special partition bounds, i.e. 1MB bound.

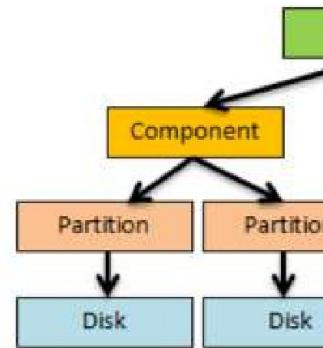


Figure 6.3



same

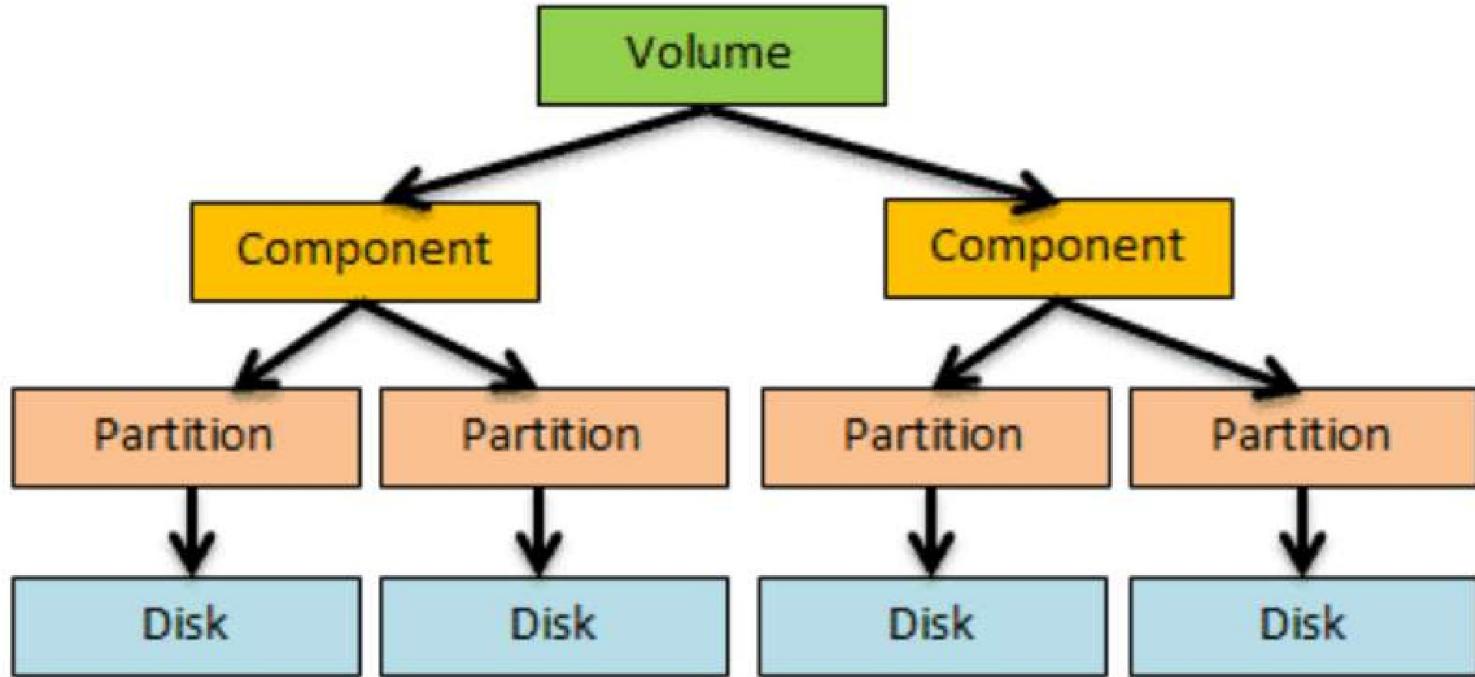


Figure 6.3

bounds,

Names:

- caps sensitive;
- caps insensitive.

Based on the underlying media:

- **disk** file systems (optical, HDD, SSD);
- **tape** file systems (**linear, mixed data and metadata**);
- **RAM** file systems;
- **database** file systems (**stored in RDBM**);
- **device** file systems (**/proc, /sys**);
- **nested** file systems (**disk images**);
- etc.

Arch

- lo
in
W

- There are plenty of them (>100), **it would worth a full semester course.**
- Stores the data organized into **files** (the analogy of paper documents), and **directories** (the analogy of paper folders).
- Since 1961.
- **Classification...**
- **Metadata...**
- **Slack space:**
 - block allocation;
 - files rarely end on block limits;
 - the average amount of space left between the file end and the last block end;
 - on average: $\text{block size} / 2 / \text{file, file number} \times \text{block size} / 2$;
- **Fragmentation:**
 - **free space** fragmentation;
 - **file** fragmentation;
- File versioning (e.g. VMS);
- **File System journaling** (e.g. ext3);

Base
• di
• ta
m
• R
• da
• de
• ne
• et

FS journali
• protecti
corrupt
• metadat
buffer, i
• they are
operatio
• damage
back int

Metadata:

- creation, modification time;
- creator;
- size: block allocation vs byte count;
- etc.

- There are plenty of them (>100), **it would worth a full semester course.**
- Stores the data organized into **files** (the analogy of paper documents), and **directories** (the analogy of paper folders).
- Since 1961.
- **Classification...**
- **Metadata...**
- **Slack space:**
 - block allocation;
 - files rarely end on block limits;
 - the average amount of space left between the file end and the last block end;
 - on average: $\text{block size} / 2 / \text{file, file number} \times \text{block size} / 2$;
- **Fragmentation:**
 - **free space** fragmentation;
 - **file** fragmentation;
- File versioning (e.g. VMS);
- **File System journaling** (e.g. ext3);

Base
• di
• ta
m
• R
• da
• de
• ne
• et

FS journali
• protecti
corrupt
• metadat
buffer, i
• they are
operatio
• damage
back int

FS journaling:

- **protecting metadata against data corruption**;
- metadata changes first go into a **circular buffer**, i.e. a **journal**;
- they are committed to media as atomic operations;
- damaged journal entries cause a **roll back** into last consistent state;



FS maintenance:

- multiple IOs -> different portion change;
- incomplete write operations -> neglected structural parts (e.g. orphan inodes, allocated blocks, etc.);
- **consistency checks, reparations;**
- **defragmentation;**

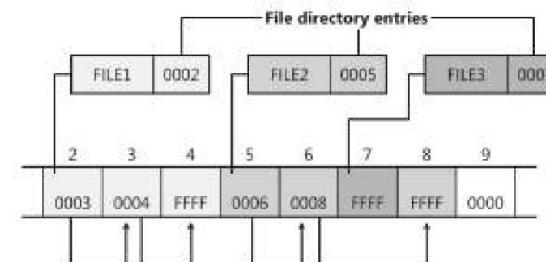
File Allocation Table (FAT) File System

- One of the oldest FS'es from the 1970s;
- **Simple**, focuses only on storing data;
- **Little endian** byte order;
- **Still wide-spread** due to its simplicity;
- Versions: FAT12, FAT16, FAT32, i.e. the width of an entry in the FAT table;
- Structures and **divides space into clusters**: i.e. contiguous space of blocks:
 - cluster size: 4 kByte as default;
 - stores file data in **cluster chains**;
 - that are not necessarily adjacent;
- **File Allocation Tables**: maps/indices of data region (cluster index or delimiter, 0xFFFF);
 - 0: FAT identifier;
 - 1: End-of-Chain (EOC) delimiter, 0xFFFFFFFF;
 - value 0: free cluster;
 - value EOC: end of cluster chain;
 - value **any**: cluster index;
- **Directory entries**:
 - 32-byte long entries;
 - root FS in a special region;
 - original **8.3** file name notation, long names trick;

Device layout structure:

Boot sector	File allocation table 1	File allocation table 2 (duplicate)	Root directory	Other directories and all files
-------------	-------------------------	-------------------------------------	----------------	---------------------------------

Cluster chains:



Directory entries:

Second (and last) long entry											
0x42	w	n	-	t	o	0x0F	0x00	Check sum	x		
0x0000	0xFFFF	0xFFFF	0xFFFF	0xFFFF	0x0000	0xFFFF	0x00	0xFFFF			
0x01	T	h	a	l	q	0x0F	0x00	Check sum	u		
	c	b		b	0x0000	r	0x20	NT	o		
T	H	E	O	U	-	I	O	X	0x20	NT	Create time
Create date	Last access date	0x0000	Last modified time	Last modified date	First cluster						File size

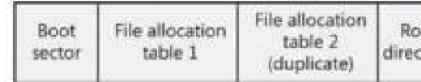
Short entry
First long entry

<https://social.technet.microsoft.com/wiki/contents/articles/6771.the-fat-file-system.aspx>
https://en.wikipedia.org/wiki/Design_of_the_FAT_file_system#FAT

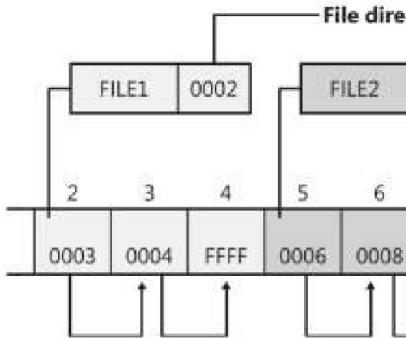
File System

- One of the oldest FS'es from the 1970s;
- **Simple**, focuses only on storing data;
- **Little endian** byte order;
- **Still wide-spread** due to its simplicity;
- Versions: FAT12, FAT16, FAT32, i.e. the width of an entry in the FAT table;
- Structures and **divides space into clusters**: i.e. contiguous space of blocks:
 - cluster size: 4 kByte as default;
 - stores file data in **cluster chains**;
 - that are not necessarily adjacent;
- **File Allocation Tables**: maps/indices of data region (cluster index or delimiter, 0xFFFF);
 - 0: FAT identifier;
 - 1: End-of-Chain (EOC) delimiter, 0xFFFFFFFF;
 - value **0**: free cluster;
 - value **EOC**: end of cluster chain;
 - value **any**: cluster index;
- **Directory entries**:
 - 32-byte long entries;
 - root FS in a special region;
 - original **8.3** file name notation, long names trick;

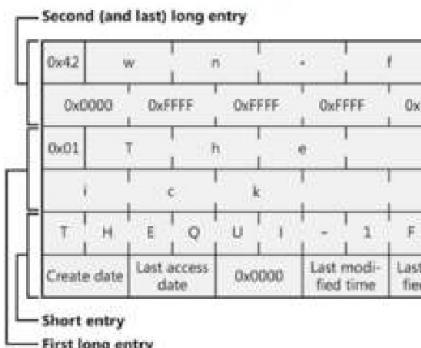
Device layout



Cluster chain



Directory entry



<https://social.technet.microsoft.com/wiki/contents/articles/6771.the-fat-file-system.aspx>
https://en.wikipedia.org/wiki/Design_of_the_FAT_file_system#FAT

Device layout structure:

Boot sector	File allocation table 1	File allocation table 2 (duplicate)	Root directory	Other directories and all files
-------------	-------------------------	-------------------------------------	----------------	---------------------------------

Cluster chains:

File directory entries

Directory entries:

Second (and last) long entry

0x42	w	n	.	f	o	0x0F	0x00	Check sum	x
0x0000	0xFFFF	0xFFFF	0xFFFF	0xFFFF	0x0000	0xFFFF	0xFFFF		
0x01	T	h	e		q	0x0F	0x00	Check sum	u
i	c	k		b	0x0000	r			o
T	H	E	Q	U	I	-	1	F	O
Create date	Last access date	0x0000	Last modified time	Last modified date	First cluster	NT			Create time
									File size

Short entry

First long entry

New Technologies File System (NTFS)

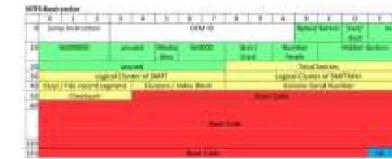
- Since 1993, as a replacement of FAT; more flexible.
- New features:
 - **Variable cluster size**, default to 4 kBytes, up to 2 MBytes;
 - **Journaling** (\$LogFile);
 - Hard links;
 - **Streams**: a logical sequence of file records;
 - **File compression and sparse files**;
 - Volume Shadow Copy based on Copy on Write;
 - ACLs: Directory ACL, Security ACL;
 - Encryption;
 - Quotas;
 - **Resizing**;
- NTFS structure...

File	Name	Content
00	File Record	File record header
01	File Record	File record header
02	File Record	File record header
03	File Record	File record header

File	Name	Content
00	File Record	File record header
01	File Record	File record header
02	File Record	File record header
03	File Record	File record header

File	Name	Content
00	File Record	File record header
01	File Record	File record header
02	File Record	File record header
03	File Record	File record header

File	Name	Content
00	File Record	File record header
01	File Record	File record header
02	File Record	File record header
03	File Record	File record header



File	Name	Content	Description
00	File Record	File record header	Contains one file file record for each file and folder on an NTFS volume. It also contains information for a file or folder to be large so it contains multiple records, either file records are allocated as such.
01	MFT mirror	0	Provides an exact copy of the MFT in case of a single-sector failure. It is a standard file system structure.
02	Log file	0	Contains information used by NTFS for Delta recording. The log file is used by Windows Server 2003 for volume metadata consistency in NTFS volumes. It is a standard file system structure.
03	Volume	0	Contains information about the volume, such as the volume label and the file allocation table.
04	Bitmap	0	Contains attribute names, numbers, and descriptions.
05	File name index	0	The root folder.
06	Cluster bitmap	0	Reserves free clusters by subtracting them from unused clusters.
07	Root sector	0	Contains the root folder of the volume.
08	Bad sector file	0	Contains bad sectors for the volume.
09	Security file	0	Contains security settings for all files within a volume.
10	Unused file	0	Contains unallocated clusters for reading to make room for new clusters.
11	MFT extension	12	Used for various optional extensions such as quotas, roaming point files, and object identifiers.
12	Unused	12	Reserved for future use.

- Metadata: **Master File Table (MFT) + MFT mirror**;
- Stores ACLs, creation times, size info, data blocks, etc.
- 1 kB entries;
- NTFS metafiles (0-26, 27- regular files);
- **MFT structure:**
 - header (block size, cluster size, type, etc.);
 - attribute headers (standard, short name, long name, etc.);
 - attributes (length, type, value, etc.);
 - resident vs. non-resident files/attributes;
- Data stored in **either MFT** (up to 900 bytes) or **clusters**.
- **Runs**: cluster intervals.
- Uses **B-trees** instead of tables and lists.

http://www.ntfs.com/ntfs_basics.htm

<https://en.wikipedia.org/wiki/NTFS>

<https://www.yumpu.com/en/document/read/11722944/ntfs-cheat-sheet-writeblocked>

\$Uppercase	Uppercase table	10	Converts lowercase characters to matching Unicode uppercase characters.
\$Extend	NTFS extension file	11	Used for various optional extensions such as quotas, reparse point data, and object identifiers.
		12-15	Reserved for future use.

- Metadata: **Master File Table (MFT) + MFT mirror**;
- Stores ACLs, creation times, size info, data blocks, etc.
- 1 kB entries;
- NTFS metafiles (0-26, 27- regular files);
- **MFT structure:**
 - header (block size, cluster size, type, etc.);
 - attribute headers (standard, short name, long name, etc.);
 - attributes (length, type, value, etc);
 - resident vs. non-resident files/attributes;
- Data stored in **either MFT** (up to 900 bytes) or **clusters**.
- **Runs**: cluster intervals.
- Uses **B-trees** instead of tables and lists.

NTFS Boot sector

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F																																			
0	Jump Instruction	OEM ID						Bytes/ Sector			Sect/ clust		res																																						
10	0x000000	unused		Media desc	0x0000	Sect / track		Number heads		Hidden Sectors																																									
20	unused						Total Sectors																																												
30	Logical Cluster of \$MFT						Logical Cluster of \$MFTMirr																																												
40	Clust / File record segment	Clusters / Index Block					Volume Serial Number																																												
50	Checksum	Boot Code																																																	
60																																																			
.																																																			
.																																																			
1E0																																																			
1F0	Boot Code													55	AA																																				

\$Uppercase	Uppercase table	10	Converts lowercase characters to matching Unicode uppercase characters.
\$Extend	NTFS extension file	11	Used for various optional extensions such as quotas, reparse point data, and object identifiers.
		12-15	Reserved for future use.

- Metadata: **Master File Table (MFT) + MFT mirror**;
- Stores ACLs, creation times, size info, data blocks, etc.
- 1 kB entries;
- NTFS metafiles (0-26, 27- regular files);
- **MFT structure:**
 - header (block size, cluster size, type, etc.);
 - attribute headers (standard, short name, long name, etc.);
 - attributes (length, type, value, etc);
 - resident vs. non-resident files/attributes;
- Data stored in **either MFT** (up to 900 bytes) or **clusters**.
- **Runs**: cluster intervals.
- Uses **B-trees** instead of tables and lists.

NTFS files

File	Name	\$MFT record #	Description
\$Mft	Master File Table	0	Contains one base file record for each file and folder on an NTFS volume. If the allocation information for a file or folder is too large to fit within a single record, other file records are allocated as well.
\$MftMirr	MFT mirror	1	Guarantees access to the MFT in case of a single-sector failure. It is a duplicate image of the first four records of the MFT.
\$LogFile	Log file	2	Contains information used by NTFS for faster recoverability. The log file is used by Windows Server 2003 to restore metadata consistency to NTFS after a system failure. The size of the log file depends on the size of the volume, but you can increase the size of the log file by using the Chkdsk command.
\$Volume	Volume	3	Contains information about the volume, such as the volume label and the volume version.
\$AttrDef	Attribute definitions	4	Lists attribute names, numbers, and descriptions.
.	Root file name index	5	The root folder.
\$Bitmap	Cluster bitmap	6	Represents the volume by showing free and unused clusters.
\$Boot	Boot sector	7	Includes the BPB used to mount the volume and additional bootstrap loader code used if the volume is bootable.
\$BadClus	Bad cluster file	8	Contains bad clusters for a volume.
\$Secure	Security File	9	Contains unique security descriptors for all files within a volume.
\$Upcase	Upcase table	10	Converts lowercase characters to matching Unicode uppercase characters.
\$Extend	NTFS extension file	11	Used for various optional extensions such as quotas, reparse point data, and object identifiers.
		12-15	Reserved for future use.

File Record Segment Header

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	F	I	L	E	Update Seq array offset	Update Seq array size										\$LogFile Sequence Number
1	Seq no	Hard Link Count		1 st attrib offset		Flags										Allocated size of file record
2					File reference to base file record				Next attrib ID							MFT Record No
3					default location of update seq array (size determined by seq size)											Reserved for update sequence array?
					Reserved for sequence array?											Common location of 1 st attrib

Resident Attribute Header

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	Type ID			Attribute Length			Form code	name len	Name offset	flags			Attrib ID			
1	Content length			Content offset		unused										

Form code

0x00 = Resident

0x01 = Non resident

Flags

0x00FF = Compressed

0x8000 = Sparse

0x4000 = Encrypted

\$File_Name

0A01 - Non Resident

0x4000 = Encrypted

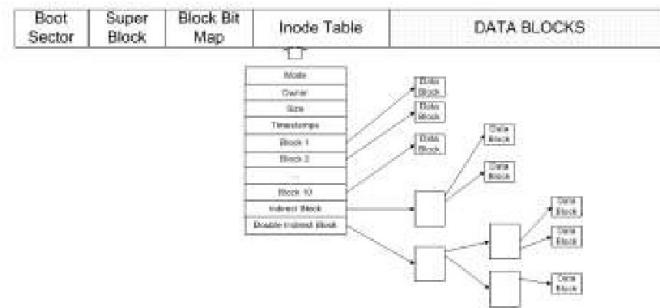
\$Data (Standard Header with data run, may be resident or non resident, non resident shown here)

Unix File Systems (UFS)

- Multiple variations: UFS, EXT2, EXT3, ReiserFS, etc.
- Represents a different principle:
 - No letter tagged devices;
 - **single tree-structure**;
 - each partition, device is mounted under the / tree;
- Metadata stored in **i-nodes**.
- Directory entries are text files containing file names and i-node assignments.
- Multiple references to same i-nodes are hard links.
- **"Everything is a file:"**
 - directories and files;
 - device files, sockets, pipes, links;

UFS structure:

- boot block: serves for OS boot;
- **superblock**: describe FS geometry, version, parameters;
- **cylinder groups**: superblock backup, cylinder group header, number of i-nodes, and data blocks;
- **i-node bitmap**: allocation of i-nodes;
- **data bitmap**: allocation of data nodes;
- **i-node table**;
- **data blocks**;



Linux Ext4 file system features:

- journaling;
- data stored in extents, i.e. contiguous blocks, like cluster runs in NTFS;
- dynamic resizing;
- B-tree like directory structure;
- compression, encryption;

Linux File System

- Multiple variations: UFS, EXT2, EXT3, ReiserFS, etc.
- Represents a different principle:
 - No letter tagged devices;
 - **single tree-structure**;
 - each partition, device is mounted under the / tree;
- Metadata stored in **i-nodes**.
- Directory entries are text files containing file names and i-node assignments.
- Multiple references to same i-nodes are hard links.
- **"Everything is a file:"**
 - directories and files;
 - device files, sockets, pipes, links;

UFS structure:

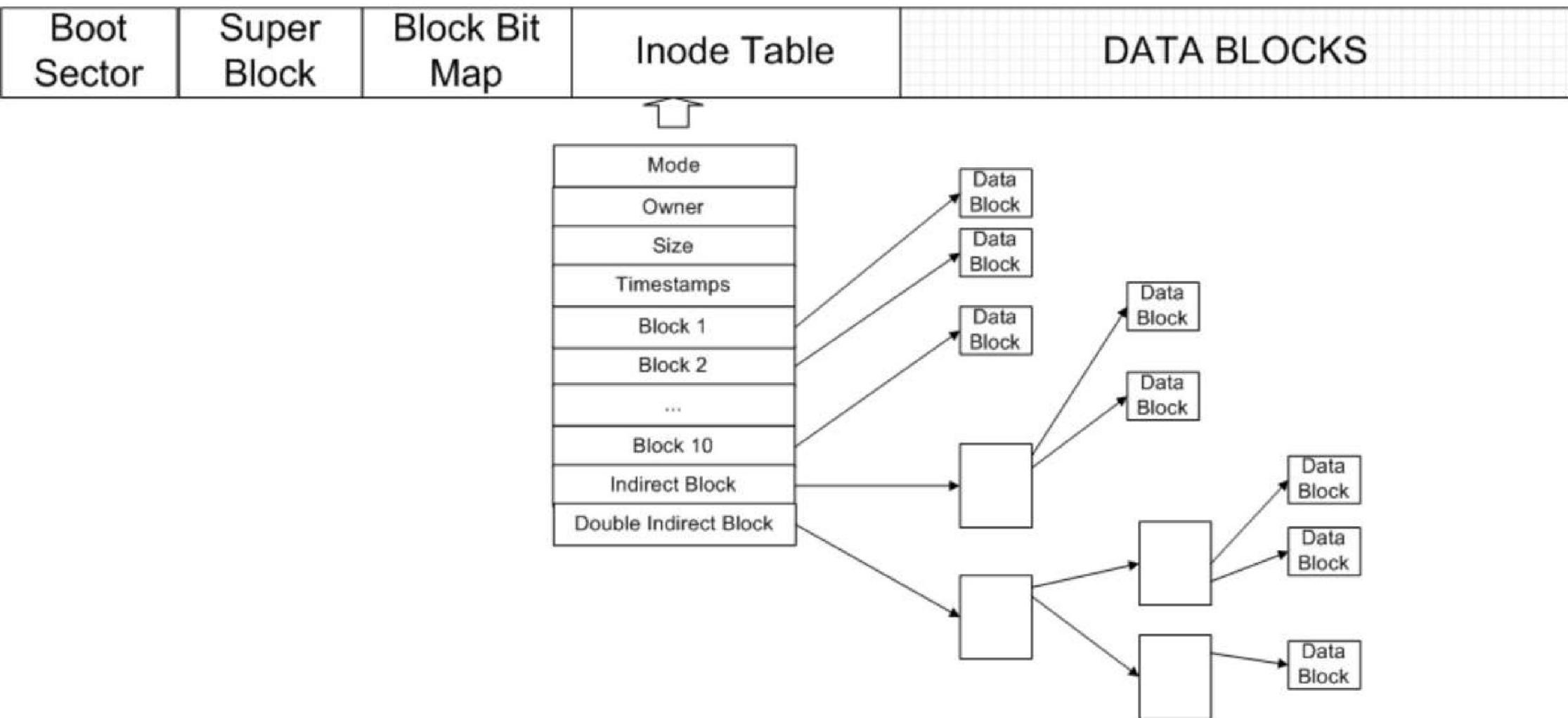
Bo
Sec

Lin

•
•

UFS structure:

- boot block: serves for OS boot;
- **superblock**: describe FS geometry, version, parameters;
- **cylinder groups**: superblock backup, cylinder group header, number of i-nodes, and data blocks;
- **i-node bitmap**: allocation of i-nodes;
- **data bitmap**: allocation of data nodes;
- **i-node table**;
- **data blocks**;



Linux Ext4 file system features:

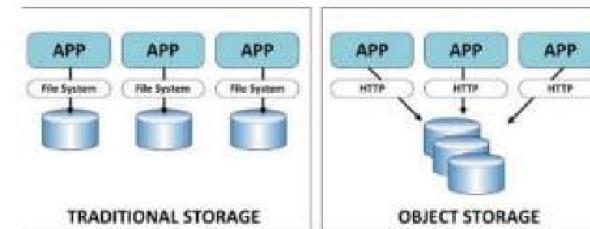
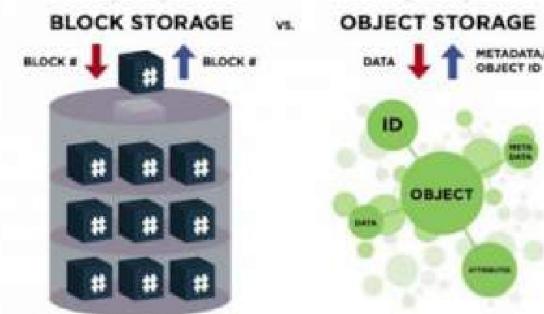
- journaling;
- data stored in extents, i.e. contiguous blocks, like cluster runs in NTFS;
- dynamic resizing;
- B-tree like directory structure;
- compression, encryption;

Distributed, Network and Global File Systems

- Huge areas on their own;
- Eliminate single points of failure;
- Allow **multiple clients to access the same file data structure**;
- Metadata protection, lock and sharing;
- **Virtual and shared name spaces** across local file systems;
- **Variants:**
 - **shared disk** file systems: sharing a common area on disks (NFS, Lustre);
 - **distributed** file systems: use network protocols to share objects (GFS, GPFS);

Object Storage

- block storage: data on physical device;
- file storage: data in human readable structures;
- **object storage:** some sort of **associative storage** over **HTTP**;
- Design principles:
 - simplicity;
 - robustness;
 - partial associativity;



Terms:

- **Bucket:** A data storage pool, e.g. a device;
- **Objects:** Data units to be stored, like files, but have properties:
 - GID;
 - **partially associative GID matching;**
 - any parameters that help indexing;
- optimized for large amounts of data: large number of large files;

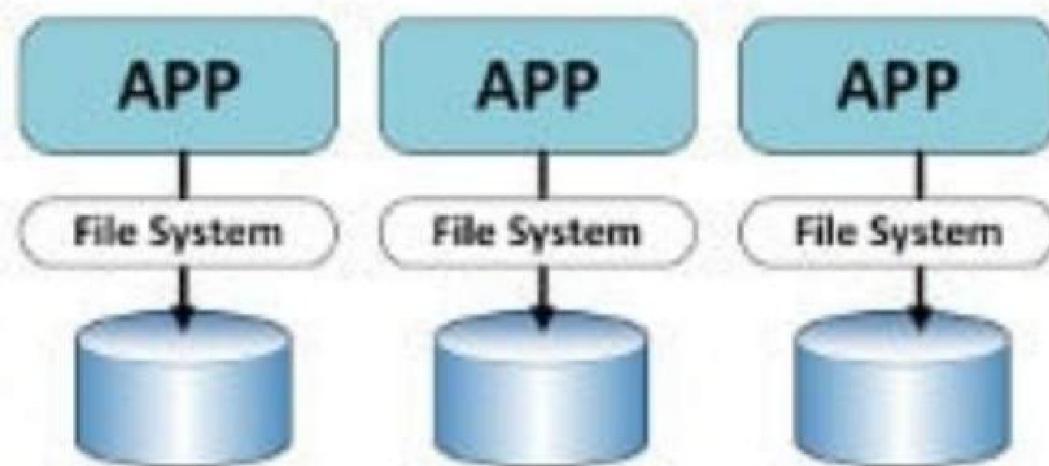
Object storage

- block storage: data on physical device;
- file storage: data in human readable structures;
- **object storage**: some sort of **associative storage** over **HTTP**;
- Design principles:
 - simplicity;
 - robustness;
 - partial associativity;

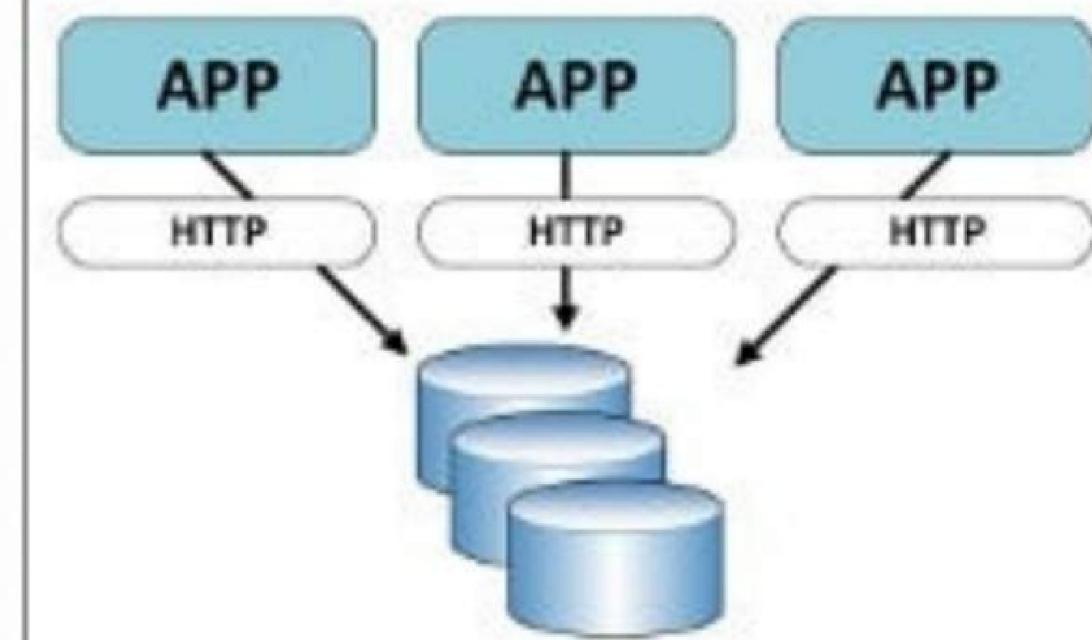


Terms:

- **Bucket**: A data device;
- **Objects**: Data files, but have



TRADITIONAL STORAGE

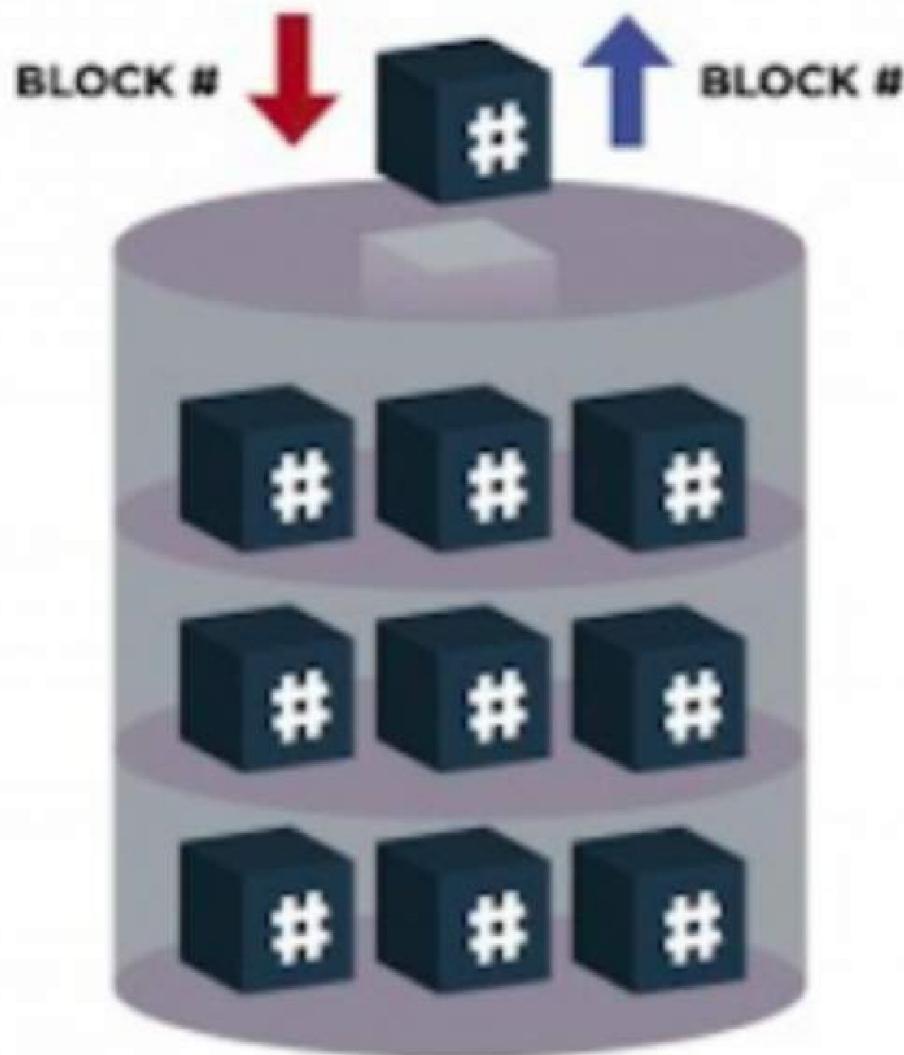


OBJECT STORAGE

Terms:

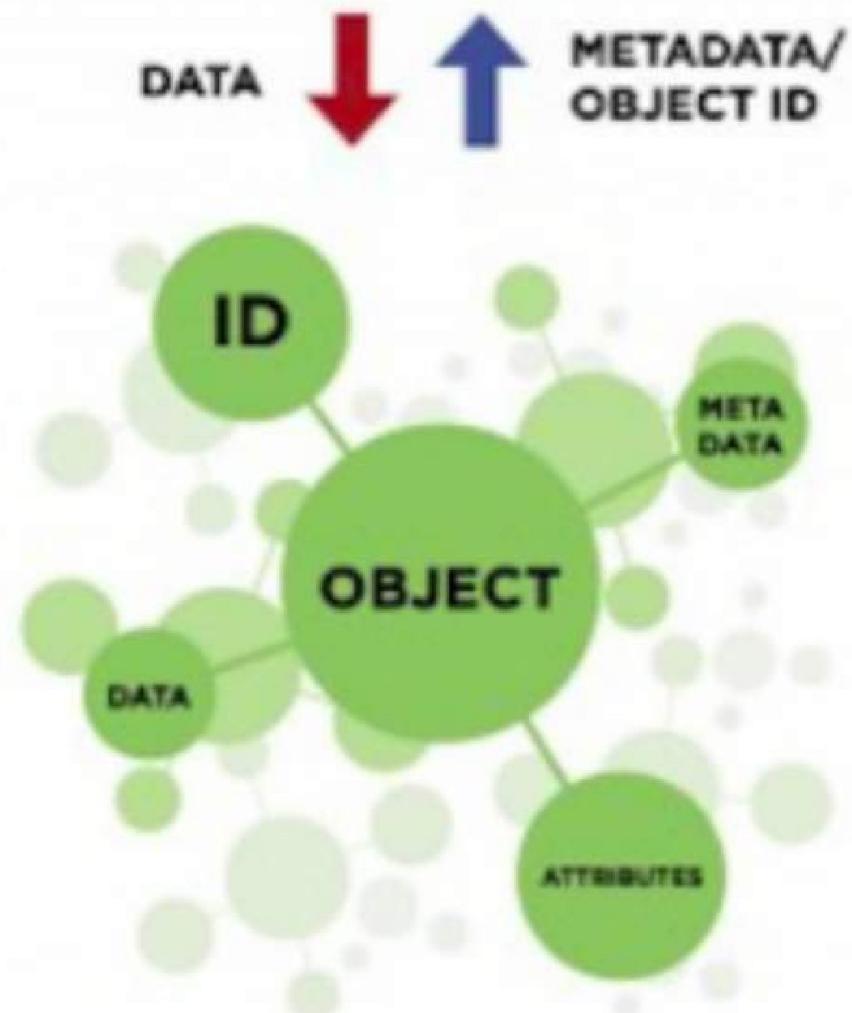
- **Bucket:** A data storage pool, e.g. a device;
- **Objects:** Data units to be stored, like files, but have properties:
 - GID;
 - **partially associative GID matching;**
 - any parameters that help indexing;
- optimized for large amounts of data: large number of large files;

BLOCK STORAGE



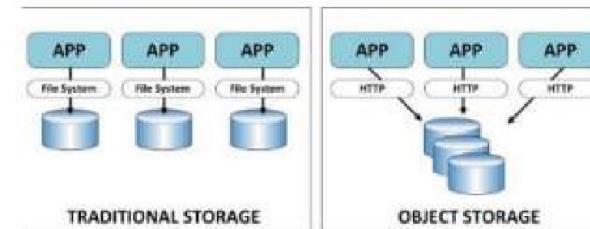
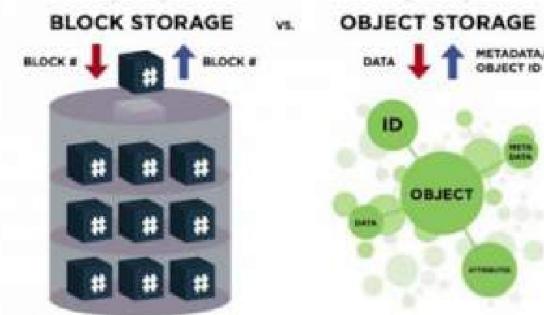
vs.

OBJECT STORAGE



Object Storage

- block storage: data on physical device;
- file storage: data in human readable structures;
- **object storage:** some sort of **associative storage** over **HTTP**;
- Design principles:
 - simplicity;
 - robustness;
 - partial associativity;



Terms:

- **Bucket:** A data storage pool, e.g. a device;
- **Objects:** Data units to be stored, like files, but have properties:
 - GID;
 - **partially associative GID matching;**
 - any parameters that help indexing;
- optimized for large amounts of data: large number of large files;