

Chapter 4: File System Management and I/O System

4.1. Mass-Storage Systems



GV: Nguyễn Thị Thanh Vân

Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018



Mass-Storage Systems

- Overview of Mass Storage Structure
- HDD Scheduling
- NVM Scheduling
- Error Detection and Correction
- Storage Device Management
- Swap-Space Management
- Storage Attachment
- RAID Structure



Operating System Concepts – 10th Edition

11.2

Silberschatz, Galvin and Gagne ©2018



Objectives

- Describe the physical structure of secondary storage devices and the effect of a device's structure on its uses
- Explain the performance characteristics of mass-storage devices
- Evaluate I/O scheduling algorithms
- Discuss operating-system services provided for mass storage, including RAID



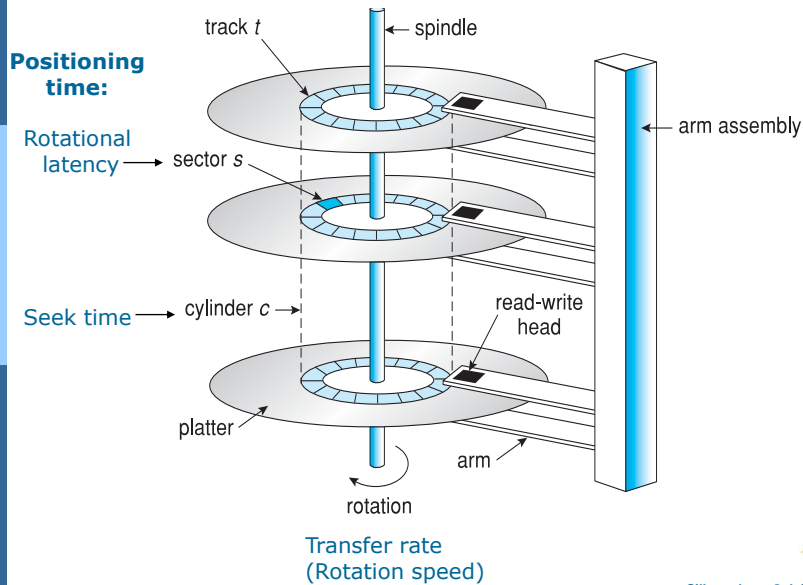
Overview of Mass Storage Structure

- Bulk of secondary storage for modern computers is **hard disk drives (HDDs)** and **nonvolatile memory (NVM)** devices
- **HDDs** spin platters of magnetically-coated material under moving r/w heads
 - The two surfaces of a platter are covered with a magnetic material.
 - store information by recording it magnetically on the platters,
 - read information by detecting the magnetic pattern on the platters
 - Drives rotate at 60 to 250 times per second
 - **Transfer rate (Rotation speed)** is rate at which data flow between drive and computer
 - **Positioning time (random-access time)** is:
 - ▶ time to move disk arm to desired cylinder (**seek time**) and
 - ▶ time for desired sector to rotate under the disk head (**rotational latency**)
 - **Head crash** results from disk head making contact with the disk surface
 - That' s bad
- Disks can be removable





Moving-head Disk Mechanism



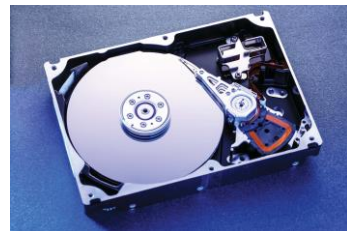
Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018



Hard Disk Drives

- Platters range from .85" to 14" (historically)
 - Commonly 3.5", 2.5", and 1.8"
- Range from 30GB to 3TB per drive
- Performance
 - Transfer Rate – theoretical – 6 Gb/sec
 - Effective Transfer Rate – real – 1Gb/sec
 - Seek time from 3ms to 12ms – 9ms common for desktop drives
 - Average seek time measured or calculated based on 1/3 of tracks
 - Latency based on spindle speed
 - ▶ $1 / (\text{RPM} / 60) = 60 / \text{RPM}$
 - Average latency = 1/2 latency



* RPM: Revolutions per minute, số vòng quay trên 1 phút.

Operating System Concepts – 10th Edition

11.6

Silberschatz, Galvin and Gagne ©2018





Hard Disk Performance

- **Access Latency = Average access time** = average seek time + average latency (milliseconds)
 - For fastest disk $3\text{ms} + 2\text{ms} = 5\text{ms}$
 - For slow disk $9\text{ms} + 5.56\text{ms} = 14.56\text{ms}$
- Average I/O time = average access time + (amount to transfer / transfer rate) + controller overhead
- For example to transfer a 4KB block on a 7200 RPM disk with a 5ms average seek time, 1Gb/sec transfer rate with a 0.1ms controller overhead
 - **Average latency** = $(60 \times 1000 \div 7200) \div 2 = 4.17\text{ms}$
 - Transfer time = $4\text{KB} \div 1\text{Gb/s} = (4 \times 1024 \times 8) \text{b} \div 10^9 = 0.033\text{ms}$
 - Average I/O time = $(5\text{ms} + 4.17\text{ms}) + 0.1\text{ms} + 0.033\text{ms} = 9.303\text{ms}$
 - Average I/O time for 4KB block = 9.303ms

* RPM: Revolutions per minute, số vòng quay trên 1 phút.



The First Commercial Disk Drive



1956
IBM RAMDAC computer included
the IBM Model 350 disk storage
system

5M (7 bit) characters
50 x 24" platters
Access time = < 1 second





Nonvolatile Memory (NVM) Devices

- NVM devices are electrical, is composed of a controller and **flash NAND** die semiconductor chips, which are used to store data.
 - It is used in a disk-drive-like container, called **solid-state disks (SSDs)**
- Other forms include: **USB drives** (thumb drive, flash drive), surface-mounted on motherboards, and main storage in devices like smartphones
- NVM
 - Can be more reliable than HDDs
 - More expensive per MB
 - Maybe have shorter life span – need careful management
 - Less capacity
 - But much faster
 - Busses can be too slow -> connect directly to PCI for example
 - No moving parts, so no seek time or rotational latency



Flash NAND vs NOR

- Bộ nhớ flash NAND (Not And) – on NVM devices:
 - là ổ nhớ flash thể rắn, dung lượng cao, rất nhanh.
 - sử dụng rất ít năng lượng, thuận tiện cho máy tính xách tay,
 - các chip nhỏ hơn, cho phép nhiều bộ nhớ được đóng gói vào cùng một không gian kích thước.
 - xóa nhanh hơn, cho nên NAND ghi nhanh hơn (vì bộ nhớ flash phải xóa các khối bộ nhớ trước khi ghi vào chúng).
 - NAND có tuổi thọ dài hơn với nhiều chu kỳ xóa hơn.
 - ưu điểm hơn so với các loại bộ nhớ flash khác, ex flash NOR (Not Or)
- Bộ nhớ flash NOR
 - thường được sử dụng trong điện thoại di động, PDA và các thiết bị khác, trong đó các chương trình nhỏ được thực thi tại chỗ thay vì ghi vào RAM.
 - Bộ nhớ flash NOR là một lựa chọn tốt cho các ứng dụng yêu cầu ít bộ nhớ và thực hiện hầu hết công việc của chúng trên mã không thay đổi, chẳng hạn như hướng khởi động hoặc hệ điều hành.





Nonvolatile Memory Devices

- Have characteristics that present challenges
- Read and written in “page” increments (think sector) but can’t overwrite in place
 - Must first be erased, and erases happen in larger “block” increments
 - Can only be erased a limited number of times before worn out – ~ 100,000
 - Life span measured in **drive writes per day (DWPD)**
 - ▶ A 1TB NAND drive with rating of 5DWPD is expected to have 5TB per day written within warranty period without failing
- Ex, A 3.5-inch SSD circuit board



SSD

- Chi phí của NAND flash đã giảm đáng kể => các thiết bị lưu trữ chính mới như ổ cứng thể rắn SSD trở thành phổ biến thay thế cho HDD truyền thống
- Ổ cứng SSD
 - lợi thế đáng kể về hiệu năng và độ bền so với ổ cứng tiêu chuẩn.
 - không có bộ phận chuyển động, nên không gặp phải tình trạng trễ cơ học như ổ cứng, có thể chịu được va đập và rung lắc tốt hơn
 - thiết bị này chỉ có các thiết bị bán dẫn NAND.





NAND Flash Controller Algorithms

- With no overwrite, pages end up with mix of valid and invalid data
- To track which logical blocks are valid, controller maintains **flash translation layer (FTL)** table
- Also implements **garbage collection** to free invalid page space
- Allocates **overprovisioning** to provide working space for GC
- Each cell has lifespan, so **wear leveling** needed to write equally to all cells

valid page	valid page	invalid page	invalid page
invalid page	valid page	invalid page	valid page

NAND block with valid and invalid pages



Volatile Memory

- DRAM frequently used as mass-storage device
 - Not technically secondary storage because volatile, but can have file systems, be used like very fast secondary storage
- **RAM drives** (with many names, including RAM disks) present as raw block devices, commonly file system formatted
- Computers have buffering, caching via RAM, so why RAM drives?
 - Caches / buffers allocated / managed by programmer, operating system, hardware
 - RAM drives under user control
 - Found in all major operating systems
 - Linux `/dev/ram`, macOS `diskutil` to create them, Linux `/tmp` of file system type `tmpfs`
- Used as high speed temporary storage
 - Programs could share bulk data, quickly, by reading/writing to RAM drive





Magnetic Tape

Magnetic tape was used as an early secondary-storage medium. Although it is nonvolatile and can hold large quantities of data, its access time is slow compared with that of main memory and drives. In addition, random access to magnetic tape is about a thousand times slower than random access to HDDs and about a hundred thousand times slower than random access to SSDs so tapes are not very useful for secondary storage. Tapes are used mainly for backup, for storage of infrequently used information, and as a medium for transferring information from one system to another.

A tape is kept in a spool and is wound or rewound past a read-write head. Moving to the correct spot on a tape can take minutes, but once positioned, tape drives can read and write data at speeds comparable to HDDs. Tape capacities vary greatly, depending on the particular kind of tape drive, with current capacities exceeding several terabytes. Some tapes have built-in compression that can more than double the effective storage. Tapes and their drivers are usually categorized by width, including 4, 8, and 19 millimeters and 1/4 and 1/2 inch. Some are named according to technology, such as LTO-6 (Figure 11.5) and SDLT.



Figure 11.5 An LTO-6 Tape drive with tape cartridge inserted.



Disk Attachment

- A secondary storage device is attached to a computer by the system bus or an I/O bus
- Several busses available, including **advanced technology attachment (ATA)**, **serial ATA (SATA)**, **eSATA**, **serial attached SCSI (SAS)**, **universal serial bus (USB)**, and **fibre channel (FC)**.
- Most common is SATA
- Because NVM much faster than HDD, new fast interface for NVM called **NVM express (NVMe)**, connecting directly to PCI bus
- Data transfers on a bus carried out by special electronic processors called **controllers** (or **host-bus adapters, HBAs**)
 - Host controller on the computer end of the bus, while device controller on device end
 - Computer places command on host controller, using memory-mapped I/O ports
 - ▶ Host controller sends messages to device controller
 - ▶ Data transferred via DMA between device and computer DRAM





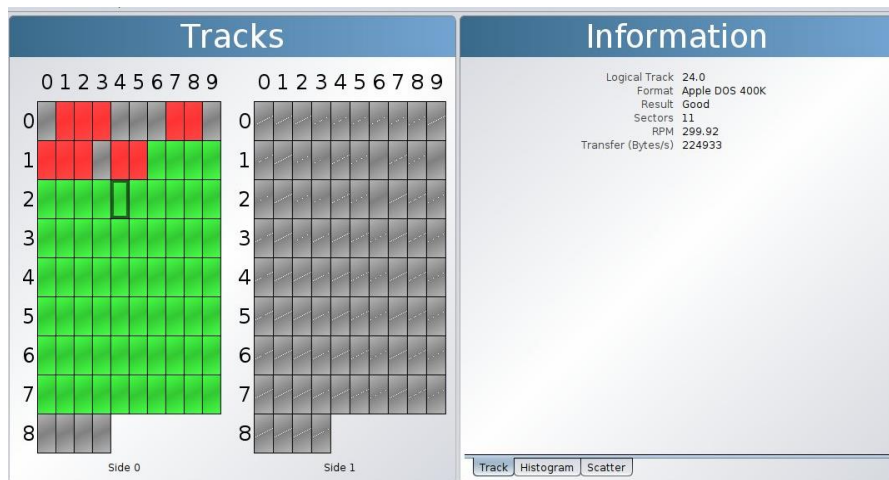
Address Mapping

- Disk drives are addressed as large 1-dimensional arrays of **logical blocks**, where the logical block is the smallest unit of transfer
 - Low-level formatting creates **logical blocks** on physical media
- The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially
 - Sector 0 is the first sector of the first track on the outermost cylinder
 - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost
 - Logical to physical address should be easy
 - Except for bad sectors
 - Non-constant # of sectors per track via constant angular velocity



Ex

- 1-dimensional array of logical blocks is mapped into the sectors





HDD Scheduling

- The operating system is responsible for using hardware efficiently — for the disk drives, this means having a fast access time and disk bandwidth
 - Minimize seek time
 - Seek time \approx seek distance
- Disk **bandwidth** is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer
- Whenever a process needs I/O to or from the drive, it issues a system call to the operating system.
 - If the desired drive and controller are available, It can be serviced
 - If the drive or controller is busy, It will be placed in the queue
- OS maintains queue of requests, per disk or device
 - Optimization algorithms only make sense when a queue exists



Disk Scheduling (Cont.)

- In the past, operating system responsible for queue management, disk drive head scheduling
 - Now, built into the storage devices, controllers
 - Just provide Logical Block Addressing (LBA), handle sorting of requests
 - ▶ Some of the algorithms they use described next
- Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying “depth”)
- Several algorithms exist to schedule the servicing of disk I/O requests
- The analysis is true for one or many platters
- We illustrate scheduling algorithms with a request queue (0-199)

98, 183, 37, 122, 14, 124, 65, 67

Head pointer 53



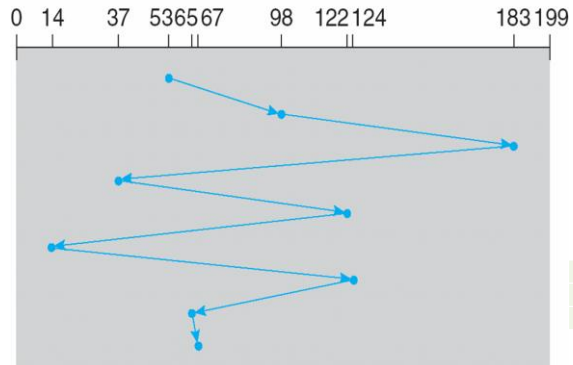


FCFS - first-come, first-served

- Illustration shows total head movement of 640 cylinders

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53



53	98	183	37	122	14	124	65	67
45	85	-146	85	-108	110	-59	2	
45	85	146	85	108	110	59	2	640

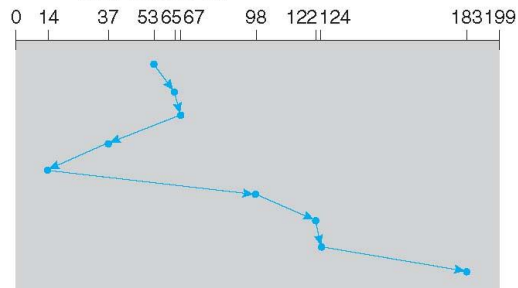


SSTF

- Shortest Seek Time First selects the request with the minimum seek time from the current head position
- SSTF scheduling is a form of SJF scheduling; may cause starvation of some requests
- Illustration shows total head movement of 236 cylinders

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53



53	65	67	37	14	98	122	124	183
12	2	-30	-23	84	24	2	59	
12	2	30	23	84	24	2	59	236

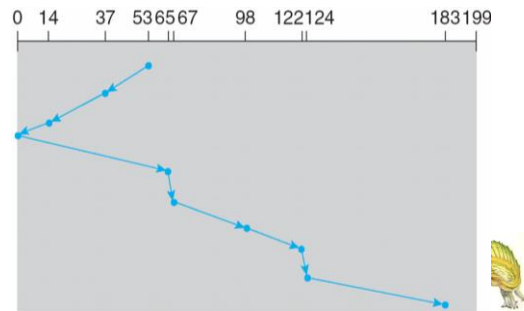




SCAN

- The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues.
- **SCAN algorithm** Sometimes called the **elevator algorithm**
- Illustration shows total head movement of 208 cylinders
- But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest

queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53



Operating System Concepts – 10th Edition

11.23

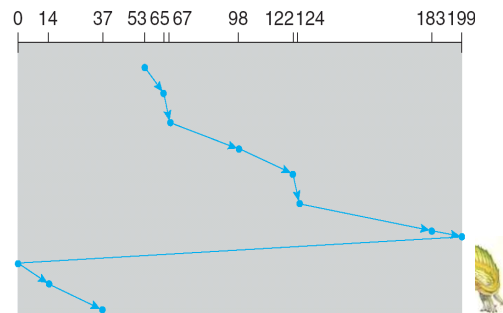
Silberschatz, Galvin and Gagne ©2018



Circular-SCAN (C-SCAN)

- Provides a more uniform wait time than SCAN
- The head moves from one end of the disk to the other, servicing requests as it goes
 - When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip
- Treats the cylinders as a circular list that wraps around from the last cylinder to the first one
- Total number of cylinders?

queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53



Operating System Concepts – 10th Edition

11.24

Silberschatz, Galvin and Gagne ©2018



Selecting a Disk-Scheduling Algorithm

- SSTF is common and has a natural appeal
- SCAN and C-SCAN perform better for systems that place a heavy load on the disk
 - Less starvation, but still possible
- To avoid starvation Linux implements **deadline** scheduler
 - Maintains separate read and write queues, gives read priority
 - Because processes more likely to block on read than write
 - Implements four queues: 2 x read and 2 x write
 - 1 read and 1 write queue sorted in LBA order, essentially implementing C-SCAN
 - 1 read and 1 write queue sorted in FCFS order
 - All I/O requests sent in batch sorted in that queue's order
 - After each batch, checks if any requests in FCFS older than configured age (default 500ms)
 - If so, Logical Block Addressing (LBA) queue containing that request is selected for next batch of I/O
- In Red Hat (RHEL 7) also **NOOP** and **completely fair queueing** scheduler (**CFQ**) also available, defaults vary by storage device



NVM Scheduling

- NVM devices do not contain moving disk heads and commonly use a simple FCFS policy.
- For example, the Linux NOOP scheduler uses an FCFS policy but modifies it to merge adjacent requests.
- In NVM devices:
 - the time required to service reads is uniform
 - write service time is not uniform - slower than reading, decreasing the adv
 - I/O can occur sequentially or randomly, HDD at sequential
 - Throughput can be similar
 - Random-access I/O (is measured in Input/Output operations per second - **IOPS**): much higher with NVM than HDD
 - **Write amplification** (one write, causing garbage collection and many read/writes) can decrease the performance advantage





Error Detection and Correction

- **Error detection** determines if there a problem has occurred (for example a bit flipping)
 - If detected, can halt the operation
 - Detection frequently done via parity bit
- Parity is one form of **checksum** – uses modular arithmetic to compute, store, compare values of fixed-length words
- Another error-detection method common in networking is **cyclic redundancy check (CRC)** which uses hash function to detect multiple-bit errors
- **Error-correction code (ECC)** not only detects, but can correct some errors
 - Soft errors correctable, while hard errors detected but not corrected



Storage Device Management

- **Low-level formatting**, or **physical formatting** — Dividing a disk into sectors that the disk controller can read and write
 - Each sector can hold header information, plus data, plus error correction code (**ECC**)
 - Usually 512 bytes of data but can be selectable
- To use a disk to hold files, the operating system still needs to record its own data structures on the disk
 - **Partition** the disk into one or more groups of cylinders, each treated as a logical disk
 - **Logical formatting** or “making a file system”
 - To increase efficiency most file systems group blocks into **clusters**
 - ▶ Disk I/O done in blocks
 - ▶ File I/O done in clusters





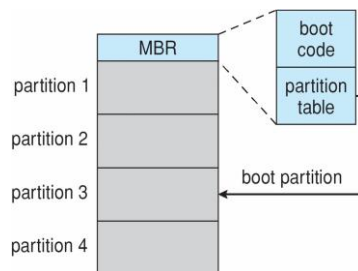
Storage Device Management (cont.)

- **Root partition** contains the OS, other partitions can hold other Oses, other file systems, or be raw
 - **Mounted** at boot time
 - Other partitions can mount automatically or manually
- At mount time, file system consistency checked
 - Is all metadata correct?
 - ▶ If not, fix it, try again
 - ▶ If yes, add to mount table, allow access
- Boot block can point to boot volume or boot loader set of blocks that contain enough code to know how to load the kernel from the file system
 - Or a boot management program for multi-os booting



Device Storage Management (Cont.)

- Raw disk access for apps that want to do their own block management, keep OS out of the way (databases for example)
- Boot block initializes system
 - The bootstrap is stored in ROM, firmware
 - **Bootstrap loader** program stored in boot blocks of boot partition
- Methods such as **sector sparing** used to handle bad blocks



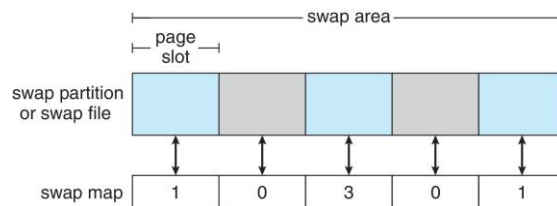
Booting from secondary storage in Windows





Swap-Space Management

- Used for moving entire processes (swapping), or pages (paging), from DRAM to secondary storage when DRAM not large enough for all processes
- Operating system provides **swap space management**
 - Secondary storage slower than DRAM, so important to optimize performance
 - Usually multiple swap spaces possible – decreasing I/O load on any given device
 - Best to have dedicated devices
 - Can be in raw partition or a file within a file system (convenience of adding)
 - Data structures for swapping on Linux systems:



Operating System Concepts – 10th Edition

11.31

Silberschatz, Galvin and Gagne ©2018



Storage Attachment

- Computers access storage in three ways
 - Host-attached
 - Network-attached
 - Cloud
- Host attached access through local I/O ports, using one of several technologies
 - To attach many devices, use storage busses such as USB, firewire, thunderbolt
 - High-end systems use **fibre channel (FC)**
 - ▶ High-speed serial architecture using fibre or copper cables
 - ▶ Multiple hosts and storage devices can connect to the FC fabric

Operating System Concepts – 10th Edition

11.32

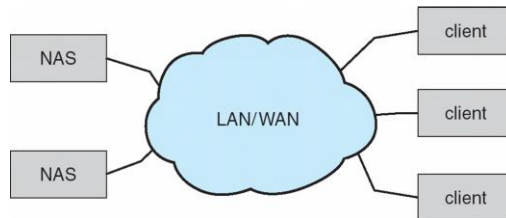
Silberschatz, Galvin and Gagne ©2018





Network-Attached Storage

- Network-attached storage (**NAS**) is storage made available over a network rather than over a local connection (such as a bus)
 - Remotely attaching to file systems
- NFS and CIFS are common protocols
- Implemented via remote procedure calls (RPCs) between host and storage over typically TCP or UDP on IP network
- **iSCSI** protocol uses IP network to carry the SCSI protocol
 - Remotely attaching to devices (blocks)



Cloud Storage

- Similar to NAS, provides access to storage across a network
 - Unlike NAS, accessed over the Internet or a WAN to remote data center
- NAS presented as just another file system, while cloud storage is API based, with programs using the APIs to provide access
 - Examples include Dropbox, Amazon S3, Microsoft OneDrive, Apple iCloud
 - Use APIs because of latency and failure scenarios (NAS protocols wouldn't work well)





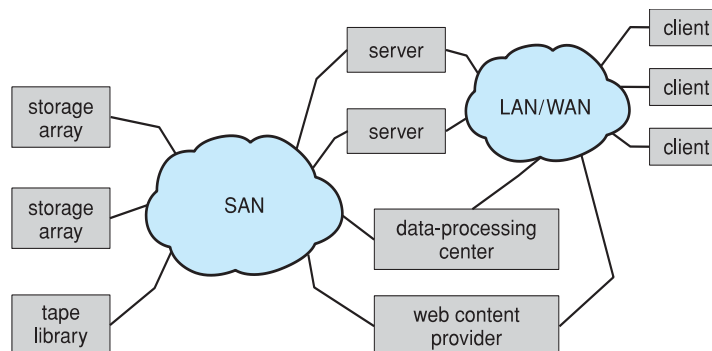
Storage Array

- Can just attach disks, or arrays of disks
- Avoids the NAS drawback of using network bandwidth
- Storage Array has controller(s), provides features to attached host(s)
 - Ports to connect hosts to array
 - Memory, controlling software (sometimes NVRAM, etc.)
 - A few to thousands of disks
 - RAID, hot spares, hot swap (discussed later)
 - Shared storage -> more efficiency
 - Features found in some file systems
 - ▶ Snapshots, clones, thin provisioning, replication, deduplication, etc



Storage Area Network

- Common in large storage environments
- Multiple hosts attached to multiple storage arrays – flexible





Storage Area Network (Cont.)

- SAN is one or more storage arrays
 - Connected to one or more Fibre Channel switches or **InfiniBand (IB)** network
- Hosts also attach to the switches
- Storage made available via **LUN Masking** from specific arrays to specific servers
- Easy to add or remove storage, add new host and allocate it storage
- Why have separate storage networks and communications networks?
 - Consider iSCSI, FCOE



A Storage Array



RAID Structure

- **RAID – redundant array of inexpensive disks**
 - multiple disk drives provides reliability via **redundancy**
- Increases the **mean time to failure**
- **Mean time to repair** – exposure time when another failure could cause data loss
- **Mean time to data loss** based on above factors
- If mirrored disks fail independently, consider disk with 1300,000 **mean time to failure** and 10 hour mean time to repair
 - Mean time to data loss is $100,000^2 / (2 * 10) = 500 * 10^6$ hours, or 57,000 years!
- Frequently combined with **NVRAM** to improve write performance
- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively



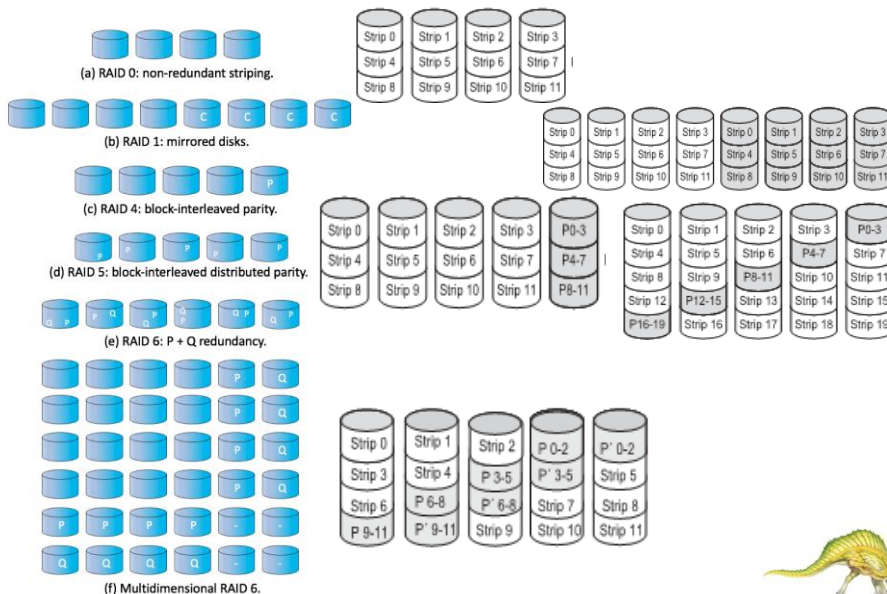


RAID (Cont.)

- Disk **striping** uses a group of disks as one storage unit
- RAID is arranged into six different levels
- RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
 - **Mirroring** or **shadowing** (**RAID 1**) keeps duplicate of each disk
 - Striped mirrors (**RAID 1+0**) or mirrored stripes (**RAID 0+1**) provides high performance and high reliability
 - **Block interleaved parity** (**RAID 4, 5, 6**) uses much less redundancy
- RAID within a storage array can still fail if the array fails, so automatic **replication** of the data between arrays is common
- Frequently, a small number of **hot-spare** disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them

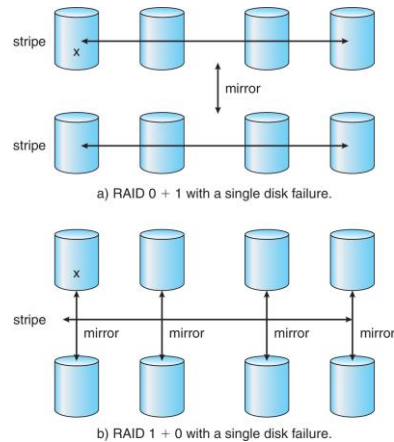


RAID Levels





RAID (0 + 1) and (1 + 0)



Other Features

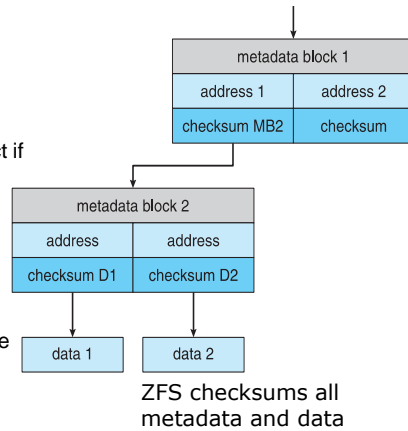
- Regardless of where RAID implemented, other useful features can be added
- **Snapshot** is a view of file system before a set of changes take place (i.e., at a point in time)
- Replication is automatic duplication of writes between separate sites
 - For redundancy and disaster recovery
 - Can be synchronous or asynchronous
- Hot spare disk is unused, automatically used by RAID production if a disk fails to replace the failed disk and rebuild the RAID set if possible
 - Decreases average of the time to repair



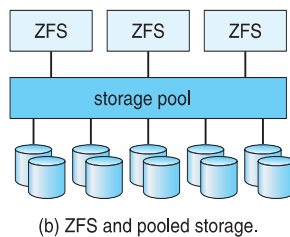
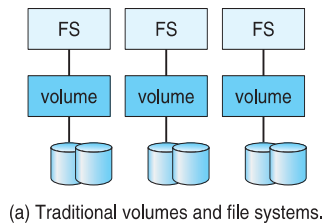


Extensions

- RAID alone does not prevent or detect data corruption or other errors, just disk failures
- Solaris ZFS adds **checksums** of all data and metadata
- Checksums kept with pointer to object, to detect if object is the right one and whether it changed
- Can detect and correct data and metadata corruption
- ZFS also removes volumes, partitions
 - Disks allocated in **pools**
 - Filesystems with a pool share that pool, use and release space like `malloc()` and `free()` memory allocate / release calls



Traditional and Pooled Storage





Object Storage

- General-purpose computing, file systems not sufficient for very large scale
- Another approach – start with a storage pool and place objects in it
 - Object just a container of **data**
 - No way to navigate the pool to find objects (no directory structures, few services)
 - Computer-oriented, not user-oriented
- Typical sequence
 - Create an object within the pool, receive an object ID
 - Access object via that ID
 - Delete object via that ID
- Object storage management software like **Hadoop file system (HDFS)** and **Ceph** determine where to store objects, manages protection
 - Typically by storing N copies, across N systems, in the object storage cluster
 - **Horizontally scalable**
 - **Content addressable, unstructured**



End of Chapter 4.1

