# UNIT IV

## STORAGE MANAGEMENT



Operating

Systems



The operating system is responsible for using hardware efficiently — for the disk

**HDD Scheduling** 

drives, this means having a fast access time and disk bandwidth

- Minimize seek time
- Seek time ≈ 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



### Disk Scheduling (Cont.)

- There are many sources of disk I/O request OS, System processes & Users processes
- I/O request includes input or output mode, disk address, memory address, number of sectors to transfer
- OS maintains queue of requests, per disk or device
- Idle disk can immediately work on I/O request, busy disk means work must queue
- Optimization algorithms only make sense when a queue exists
- In the past, operating system responsible for queue management, disk drive head scheduling
- Now, built into the storage devices, controllers



## Disk Scheduling (Cont.)

- 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

Illustration shows total head movement of 640 cylinders



INSTITUTIONS





- 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





- 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

**C-SCAN** 

- 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?





#### **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, LBA queue containing that request is selected for next batch of I/O



#### **NVM Scheduling**

- No disk heads or rotational latency but still room for optimization
- In RHEL 7 NOOP (no scheduling) is used but adjacent LBA requests are combined
  - NVM best at random I/O, HDD at sequential
  - Throughput can be similar
  - Input/Output operations per second (IOPS) much higher with NVM (hundreds of thousands vs hundreds)
  - But write amplification (one write, causing garbage collection and many read/writes) can decrease the performance advantage



#### **Error Detection and Correction**

- Fundamental aspect of many parts of computing (memory, networking, storage)
- 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 one form of checksum uses modular arithmetic to compute, store, compare values of fixedlength 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, 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

Dr.B.Anuradha / ASP / CSD/ SEM 4 / OS



#### 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



- 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

Dr.B.Anuradha / ASP / CSD/ SEM 4 / OS



#### **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 (for convenience of adding)
  - Data structures for swapping on Linux systems:



# TEXT BOOK

1. Abraham Silberschatz, Peter B. Galvin, "Operating System Concepts", 10<sup>th</sup> Edition, John Wiley & Sons, Inc., 2018.

2. Andrew S Tanenbaum, Herbert Bos, Modern Operating systems, Pearson, 5th Edition,2022 New Delhi.

#### REFERENCES

- 1. Ramaz Elmasri, A. Gil Carrick, David Levine, "Operating Systems A Spiral Approach", Tata McGraw Hill Edition, 2010.
- 2. William Stallings, Operating Systems: Internals and Design Principles, 7th Edition, Prentice Hall, 2018
- 3. Achyut S.Godbole, Atul Kahate, "Operating Systems", McGraw Hill Education, 2016.

#### **THANK YOU**