

# Embedded Operating System



# Agenda

- **File Management**

1. Directory
2. Links
3. File System Architecture
4. File IO syscalls
5. Disk allocation & Free space management
6. Linux Ext2/3 FileSystems
7. Journaling
8. Disk scheduling algorithms

- ❖ **Reading**

1. Galvin slides (File System & IO subsystem)
2. Professional Linux Kernel Architecture (Virtual File System, Extended File System)
3. Beginning Linux Programming (File SysCalls Programming)



# File Management

- **File = Data + Metadata**
  - Data --> Data blocks
  - Metadata --> Inode (FCB)
- **File System = Boot block + Super block + Inode list + Data blocks**
- **Types of Files**
  - User perspective
    - Text files
    - Archive files
    - Media files
    - Document files
    - Executable files
    - etc.



## • **Kernel perspective**

- Regular files (-) (All user perspective file are regular file)
- Special files
  - Directory files (d)
  - Link files (l)
  - Pipe files (p)
  - Socket files (s)
  - Device files
    - Char device files (c)
    - Block device files (b)



## • **Directory**

- From end user perspective, directory is a container which contains sub-directories and files.
- However, OS treats directory as a special file.
- The directory file contains one entry for each subdirectory or file in it.
- Each directory entry contains i-node number and name of sub-dir / file.
- terminal> ls -a -i -l /home/sunbeam.

## • **Directory Listing**

- terminal> ls dirpath
- Directory access library functions (man section 3)
  - opendir()
  - readdir()
  - closedir()



- opendir()
  - Open the directory file for reading.
  - DIR \*dp = opendir("dir-path");
    - arg1: dir path to be opened
    - returns: DIR pointer if dir opened successfully, otherwise NULL
- closedir()
  - Close the directory file.
  - closedir(dp);
    - arg1: DIR pointer.



## • **readdir()**

- Read the next dirent from the directory file.
- `struct dirent *ent = readdir(dp);`
  - arg1: DIR pointer.
  - returns: Pointer to struct dirent, if next entry is available.
    - Returns NULL if end of dir file is reached.
- struct dirent
  - d\_name --> name of file or sub-directory.
  - d\_ino --> inode number of file or sub-directory



- Symbolic Link

- A symbolic link, also known as a symlink or soft link, is a special type of file that points to another file or directory.
- terminal> ln -s /path/of/target/file linkpath
- Internally use
  - symlink() syscall.
  - man symlink
  - int symlink(const char \*target\_path, const char \*link\_path)

- symlink()

- syscall A new link file is created (new inode and new data block is allocated), which contains info about the target file (absolute or relative path).
- Link count is not incremented.
- If target file is deleted, the link becomes useless.
- Can create symlinks for directories also





# • Hard Link

- A hard link to a file points to the inode of the file instead of pointing to the file itself.
- This way the hard link gets all the attributes of the original file and points to the same data block as the original file.
- terminal> In targetfilepath
- Internally use link() syscall.
  - man link
  - `int link(const char *target_path, const char *link_path);`
- link() syscall
  - A new directory entry is created, which has a new name and same inode number.
  - No new file (inode and data blocks) is created.
  - Link count in the inode of the file is incremented.
  - If directory entry of target file is deleted (rm command), file can be still accessed by link directory entry.
  - Cannot create hard link for directories, because it may lead to infinite recursion (while traversing directories recursively e.g. ls -R)



- rm command
  - The 'rm' means remove.
  - This command is used to remove a file.
  - The rm command in Linux, internally calls unlink() system call.
  - `int unlink(const char *filepath);`
- unlink() syscall
  - It deletes directory entry of the file.
  - It decrements link count in the inode by 1.
  - If link count = 0, the inode is considered to be deleted/free (updated into super-block).
  - It can be reused for any new file.
  - When inode is marked free, data blocks are also made free, so that they can also be reused for some new file



## • Directory

- Directory permissions/mode
  - r -- can read from dir data block -- list directory contents.
  - w -- can write into dir data block -- create new files & sub-directories, remove file/sub-directory, rename file.
  - x -- enable browsing the directory -- "cd" command

SUNBEAM



- File System Architecture

- Virtual File System:

- This layer redirect file system request to the appropriate file system manager.

- File system manager:

- File system manager enables access to repective file system on the disk.
    - OS can see all partitions whose file system managers are installed in that OS.

- IO subsystem:

- Implement buffer cache and other mechanisms to speed up disk IO.

SUNBEAM



## • **Windows vs Linux**

- Linux have FS mgrs for ext3/4, reiserfs, xfs, fat, ntfs, cdfs, etc.
- Hence Linux support many FS.
- Windows have FS mgrs for FAT, NTFS, CDFS.
- Hence Windows do not support Linux FS.
- However, third-party FS managers can be added into Windows to support Linux FS e.g. ext2fsd.

SUNBEAM



## • VFS Structures

- struct inode unsigned long i\_ino; // inode number
- loff\_t i\_size; // file size unsigned
- int i\_nlink; // number of hard links
- umode\_t i\_mode; // file mode (permissions)
- atomic\_t i\_count; // reference count
- struct list\_head i\_list; // inode cache
- Device driver related
  - struct list\_head i\_devices;
  - dev\_t i\_rdev;
  - union {
    - struct pipe\_inode\_info \*i\_pipe;
    - struct block\_device \*i\_bdev;
    - struct cdev \*i\_cdev;
  - };
  - struct file\_operations \*i\_fop;



# Buffer Cache

- Buffer cache is used to speed up disk IO. Each buffer is a memory area and have a buffer header associated with it.
- This buffer header maintains metadata about the buffer i.e. disk block number, dirty flag, pointer to buffer, etc.
- Buffer cache is a doubly circular linked list of buffer headers.
- Whenever a disk block is read into a buffer, its dirty flag is set to false.
- It indicates that data in the buffer is same as data on disk. When write operation is done on the buffer its dirty flag is set to true.
- Then the buffer is scheduled for the write operation (in disk request queue -- which in turn scheduled as per disk scheduling algorithm).
- For read operation from the disk, an empty buffer is allocated in the buffer cache (for intended disk block).
- Then read request is scheduled (in disk request queue -- which in turn scheduled as per disk scheduling algorithm).
- If requested buffer (for read) is already present in buffer cache (non-dirty), then disk read operation is skipped -- to speed-up disk IO



## • File IO syscalls

- open() syscall
  - fd = open("/home/kiran/abc.txt", O\_RDONLY);
  - step 1. Convert given file path into its inode number. This is called as path name translation and is done by a kernel in file from the disk into inode table in memory.
  - step 2. Inodes of all recently accessed files are kept in inode table.
  - step 3. A file position is initialized to 0 and is stored in the open file table. It also stores mode in which file is opened and pointer to the in-memory inode. Information of all files opened in the system, is maintained in this table.
  - step 4. Each process is associated with a open file descriptor table. It keeps info of all files opened by that process. This entry stores pointer to the Open FileTable entry.
  - step 5. Finally index to file desc table entry is returned, which is called as "file descriptor". All further read(), write(), lseek(), close() operations will be using this file desc.





## • **struct file**

- unsigned int f\_flags; // open() arg2
- loff\_t f\_pos; // current file position
- struct path f\_path; // pointer to dentry
- #define f\_dentry f\_path.dentry
- struct list\_head fu\_list; // open file table
- atomic\_t f\_count; // reference count
- Device driver related
  - struct file\_operations \*f\_op;



- **struct dentry**

- struct qstr d\_name; // name of file/sub-directory
- struct inode \*d\_inode; // pointer to the inode
- struct list\_head d\_lru; // dentry cache
- atomic\_t d\_count; // reference count

- **struct fs\_struct**

- struct dentry \* root; // stores "root directory" of the process --> used for absolute path
- struct dentry \* pwd; // stores "current directory" of the process --> used for relative path
- int umask; // user file mode mask -- while creating new file this mask is used.



- **struct files\_struct**

- struct file \* fd\_array[NR\_OPEN\_DEFAULT];

- **struct task\_struct**

- struct fs\_struct \*fs; // current & root directory
  - struct files\_struct \*files; // open file desc tables

SUNBEAM



# Linux open() system call

- **dentry cache**

- Directory entries are in data blocks of directory file (on hard disk).
- Accessing same dentry repeatedly from disk is slower process.
- If dentry is already loaded then it can be accessed quickly using dentry cache.

## **rename() library functions**

- rename() internally use link() call to create new hard link and call unlink() to delete the old name.



## • Reference counting

- Used to manage life-time of any object (in complex systems e.g. Linux kernel, ...).
- Object has a member called as "reference count".
- The count is incremented everytime new pointer points to the same object and decremented everytime the pointer to the object is no more used/required.
- At any moment, reference count is number of pointers referring to the object.
- When reference count become zero, it means no pointer is referring to the object and the object can be deleted safely



# read() system call

- read() syscall
  - count = read(fd, buf, length); -- syscall api
    - sys\_read(fd, buffer, length) -- syscall implementation
    - vfs\_read(file, buffer, length, inode) -- Virtual file system
      - Logical FS considers file as sequential set of bytes and set of blocks.
      - Example: block size = 4096 and file size = 20000 bytes, then number of blocks = 5 (0 to 4).
        - If current file position = 10000, then reading file block = 2
    - ext3\_read(file, inode, file\_block) -- File system manager
      - Refers inode and file disk block corresponding to the file block.
      - check buffer cache
        - if disk block found.
        - if found, return it; otherwise call disk driver to read that disk block from disk
    - disk\_read(disk\_device, disk\_block)
      - device driver Read appropriate sectors and made it available into buffer cache.
      - The current process is blocked/sleep while disk read operation is in progress.



# Write() system call

SUNBEAM



# **lseek() syscall**

- **lseek() syscall**

- Change the file position in open file table entry (struct file --> f\_pos).
- And returns new file position (from the beginning of the file)

- **Examples:**

- `lseek(fd, 0, SEEK_SET);`
  - `filp->f_pos = 0;`
- `lseek(fd, offset, SEEK_SET);`
  - `filp->f_pos = offset;`
- `lseek(fd, 0, SEEK_END);`
  - `filp->f_pos = size; // file size (from the inode)`
- `lseek(fd, offset, SEEK_END);`
  - `filp->f_pos = size + offset; // note that offset will be negative`
- `lseek(fd, offset, SEEK_CUR);`
  - `filp->f_pos = filp->f_pos + offset`





# close() syscall

- **close() syscall**

- Decrement ref count in open file table entry (struct file). If ref count drops to zero, OFT entry is deleted (from OFT).

SUNBEAM



# Operating Systems Concepts

## ➤ Disk space allocation methods:

- When a file is requesting for free data blocks, then in which manner free data blocks gets allocated for that file and how its information can be kept inside inode of that file is referred as disk space allocation method.
- The file data blocks are allocated on the file system (data blocks region) on the disk.
- The data blocks can be allocated in various ways (depending on file system)

## ❖ Three disk space allocation methods are there:

1. Contiguous Allocation
2. Linked Allocation
3. Indexed Allocation



# Operating Systems Concepts

## 1. Contiguous Allocation : free data blocks gets allocated for a file in a contiguous manner.

### ❑ Advantages :

1. Sequential access.
2. Random access .
3. Simple to implement.

### ❑ Disadvantages :

1. File may not grow(Limitations).
2. **External fragmentation**  
-Number of blocks required are available but not contiguous.

### ❖ Defragmentation :

- Moves files on disk so that maximum contiguous free space is available.

Disk Space Allocation Method:

### 1. Contiguous Allocation

india.txt

- inode number: 101

- addr of starting data block=0

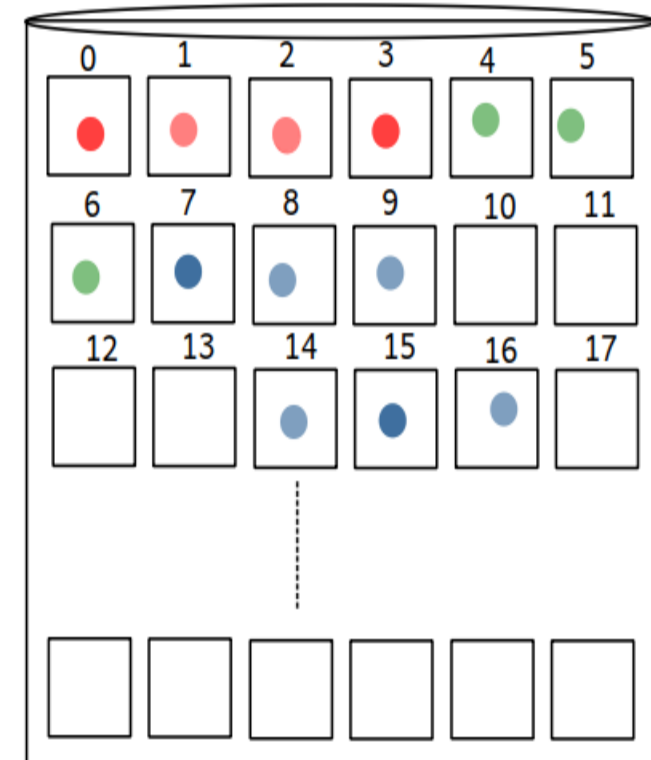
- count=4

pakistan.txt

- inode number: 201

- addr of starting data block=4

- count=3



# Contiguous Allocation

- Number of blocks required for the file are allocated contiguously.
- inode of the file contains starting block address and number of data blocks.
- Faster sequential and random access. Number of blocks required for the file may not be available contiguously.
- This is called as "External Fragmentation". To solve this problem, data blocks of the files can be shifted/moved so that max contiguous free space will be available.
- This is called as "defragmentation".



# Operating Systems Concepts

2. **Linked Allocation**: any free data blocks gets allocated for a file in a linked list manner.

## ❑ Advantages :

1. Sequential access.
2. No file grow limit.
3. No external fragmentation.

## ❑ Disadvantages :

1. Slow random access.

## ❑ Example : FAT

SUN

Disk Space Allocation Method:

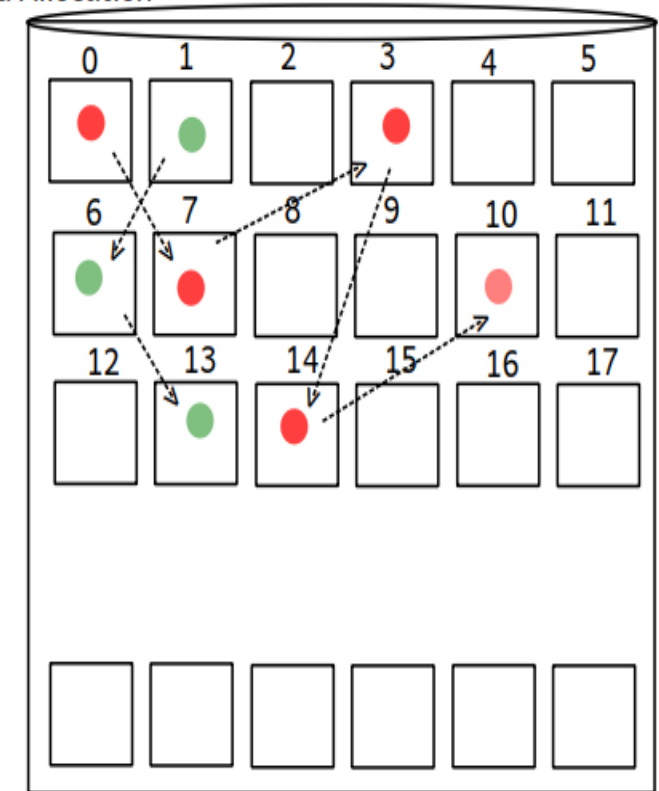
2. Linked Allocation

india.txt

- inode number: 101  
- addr of starting data block=0  
- addr of end data block=10

pakistan.txt

- inode number: 201  
- addr of starting data block=1  
- addr of end data block=13



# Linked Allocation

- Each data block of the file contains data/contents and address of next data block of that file.
- inode contains address of starting and ending data block.
- No external fragmentation, faster sequential access.
- Slower random access.
- e.g. FAT

SUNBEAM



# Operating Systems Concepts

3. **Indexed Allocation** : any free data blocks gets allocated for a file, as by maintaining an index data block information about allocated data blocks can be kept inside it.

❑ **Advantages :**

1. Sequential Access.
2. Random Access.
3. No External Fragmentation

❑ **Disadvantages :**

1. File cant grow up to some limit.

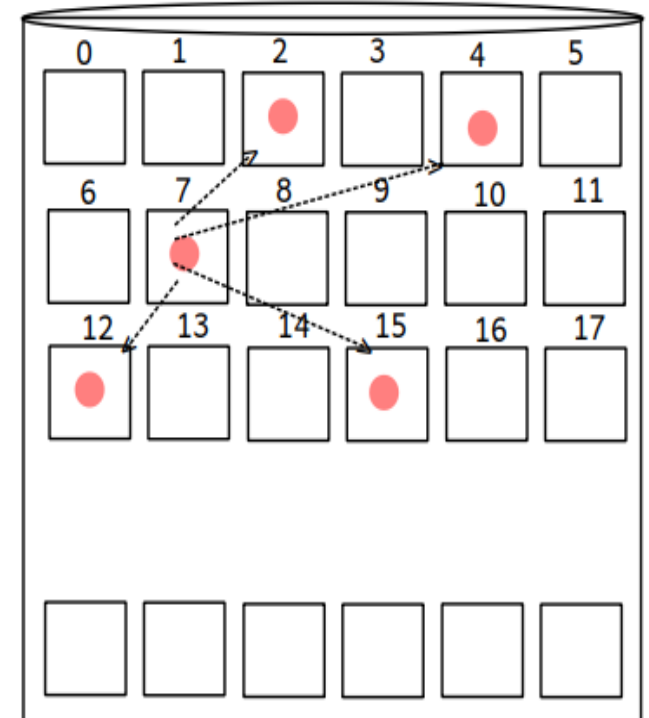
❑ **Example :** UFS , EXT 2 , EXT 3

Disk Space Allocation Method:

3. Indexed Allocation

india.txt

- inode number: 101  
- addr of index data block=7



# Indexed Allocation

- A special data block contain addresses of data blocks of the file.
- This block is called as "index block".
- The address of index block is stored in the inode of the file.
- No external fragmentation, faster random and sequential access.
- File cannot grow beyond certain limit

SUNBEAM





# Free Space Management

- The free data blocks information is maintained in the super-block.
- Super block use some data structures to maintain this information.

## 1. Bit Vector

- Super block maintains an array of bits to keep track of used and free data blocks.
- Number of bits = Number of data blocks. If nth block is used, nth bit in array should be 1.
- If nth block is free, nth bit in array should be 0.
- When allocating a data block for a file, find first free (0) bit in the bitmap.
- Corresponding free block is allocated to the file and then the bit is updated as used (1).

## 2. Linked List

- Super block keep address of first free data block.
- Each free data block keeps address of next free data block.



### 3. Grouping

- Super block keep address of an free data block.
- That free block acts like index block i.e. keeps addresses of other free data blocks

### 4. Counting

- Super block keep address of an free data block.
- That free block acts like index block i.e. keeps addresses of free data block and number of free data blocks after it



# File Systems

- File System is way of organizing files (data and metadata) on the disk.
- File systems differ in metadata (attributes), supported data block sizes, file system layout (**boot block + super block + inode list + data blocks**).
- disk allocation mechanisms, free space management mechanisms, etc.
- File system layout is managed by File system driver/manager.

SUNBEAM



- UNIX File System

Modified indexed allocation

Maximum file size = 16 GB (with 4 KB data block size)

Inode has array of 13 members

10 --> Direct data blocks

11 --> Single indirect data blocks

12 --> Double indirect data blocks

13 --> Triple indirect data blocks

# FAT File System

- Maximum file size = 4 GB.
- FAT -- File Allocation Table
  - A special data structure in the file system on the disk.
  - Array implementation of a linked list.
  - Each element in the FAT table, keep address of next block.
  - This table is used for disk allocation as well as for free space mgmt.
- Data blocks stores only data. Size of data block can be fixed while formatting
  - e.g. 2 kb, 4 kb (default), 8 kb, 16 kb.
- FAT directory entries stores all info about file (no separate inode struct).
  - file name
  - size attributes (read-only, hidden, )
  - start block
  - Boot record (1 sector = 512 bytes) -- information about file system
    - label
    - FAT table information
    - FS block size
    - etc. Does
- **Not support Journaling.**



# NTFS

- NTFS Windows platform
- Much sophisticated than FAT
- Use B-tree for disk allocation
- Max file size > 100 GB
- Also implements Journaling mechanism
- Refer: wiki

SUNBEAM



## ➤ Disk Scheduling Algorithms:

- When system want to access data from a disk, request can sent to disk controller and disk controller accepts one request at a time and complete it.
- There are chances that at a time more than one requests for accessing data from the disk can be made by the processes running in a system, in that case all the requests can **be kept in a waiting queue of the disk maintained by an OS**, and there is need to schedule/select only one request at a time and sent it to the disk controller, to do this there are certain algorithms referred as **disk scheduling algorithms**.

## ➤ Memory Technologies

❖ There are four methods by which data can be accessed from the computer memory:

1. Sequential Access: e.g. Magnetic Tape
2. Direct Access: e.g. Magnetic Disk
3. Random Access: e.g. RAM Memory
4. Associative Access: e.g. Cache Memory

## ❖ Magnetic Disk : Hard Disk Drive Structure

- HDD is made up of one or more circular platters arranged like CD rack.
- A Circular platter is made up of non-magnetic substance like aluminum or aluminum alloy, which is coated with a magnetic substance.



## ➤ Magnetic Disk : Hard Disk Drive Structure

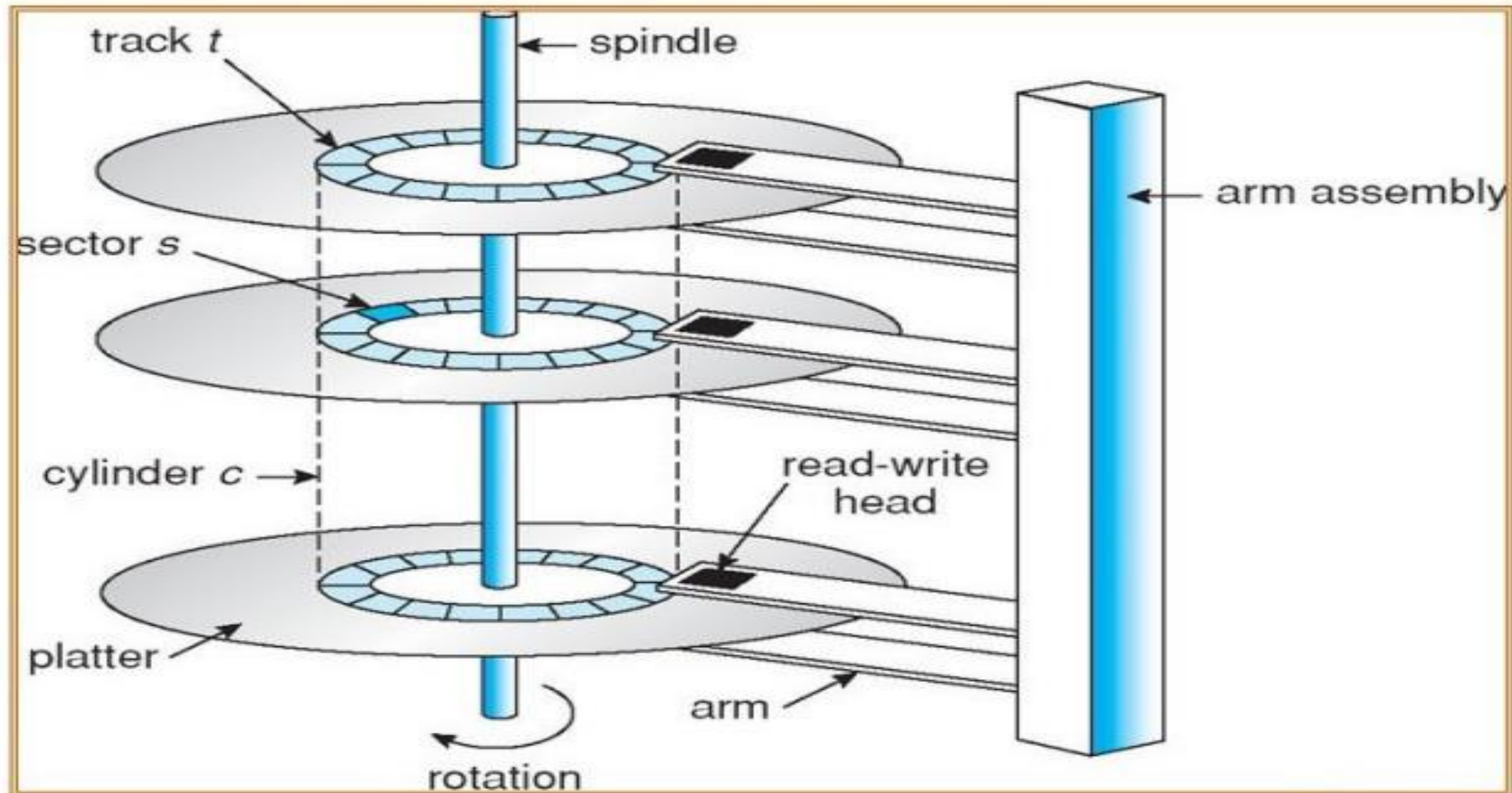
- **HDD** is made up of **one or more circular platters** arranged like CD rack.
- A Circular platter is made up of **non-magnetic substance like aluminum or aluminum alloy**, which is coated with a magnetic substance.
- Coating of magnetic substance is either from one side to the platter or from both the sides (for increasing its capacity) and hence platter in a magnetic disk may be either single sided platter or double sided platter.
- **Circular platter is divided into the hundred's of concentric rings called as tracks** whereas **each track is divided into thousands of same size of blocks called as sectors**.
- Usually the **size of each sector is 512 bytes**
- There is **one conducting coil referred as head** which is used to access data from the sector i.e. head can read and write data from and into a sector at a time.
- **Head writes and read data sector by sector i.e. block by block**, and magnetic disk is also called as **block device**.

# Computer Fundamentals and Operating Systems

- All the operations like read, write, control etc... in a HDD are controlled by **disk controller**, and hence movement of the head also controlled by it.
- **Seek Time:** time required for the disk controller to move head from its current position to the desired track.
- **Rotational Latency:** after reaching head at desired track, circular platter gets rotated till the head does not comes aligned with the desired sector, and time required for this rotation is referred as rotational latency.
- **Access Time = Seek Time + Rotational Latency.**



## Moving-head Disk Mechanism



- **FCFS (First Come First Served):**
  - request which is arrived first gets accepted and completed.
- **SSTF (Shortest Seek Time First):**
  - request which is closed to the current position of the head gets accepted and completed.
- **SCAN:**
  - head keeps scanning the disk from starting cylinder to end cylinder and whichever request came across gets accepted and completed.

- **C-SCAN (Circular SCAN):**

- head scans the disk only in a one direction

- **LOOK:**

- this policy can be used either with SCAN/C-SCAN, in this, if there is no request in a waiting queue then movement of the head gets stopped

SUNBEAM





# Thank you!

Kiran Jaybhavne

email – [kiran.jaybhavne@sunbeaminfo.com](mailto:kiran.jaybhavne@sunbeaminfo.com)

