

Video 1.

File System

File → read/write. { OS ने file system का Deal करता है } .
, डिक्टी बोर्ड डिक्टी

HDD structure / Magnetic Disk

RAM is just a piece of continuous memory.

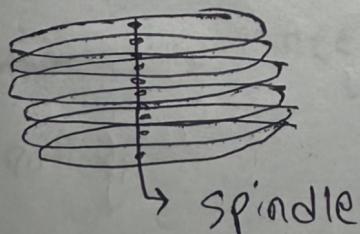
RAM में कोई file नहीं होता।

Secondary memory (HDD) में जारी रखता है,

file is just an imaginary boundary bounded by 0/1 and which stores information

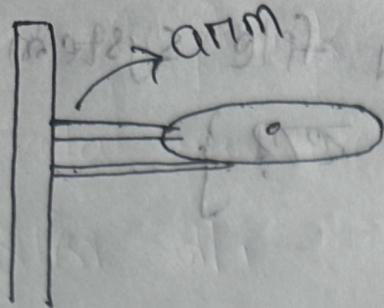
— in HDD is called a platter.

A stack of platters altogether is called a HDD.



Platter का top & bottom

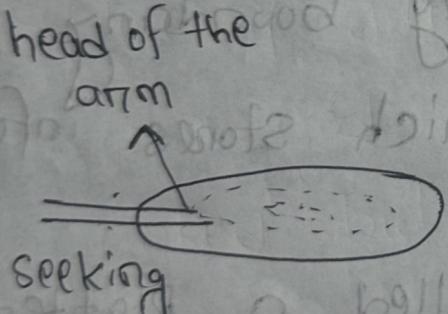
पर्सीडे लेते हैं independently
information store होता है।



Platter route $E78$ bot arm
Same place 92 $2T8$ &
data read $E78$.

Arm assembly. So, a platter should have two arms. And, will work independently.

If we have n platters then we will have $2n$ arms.

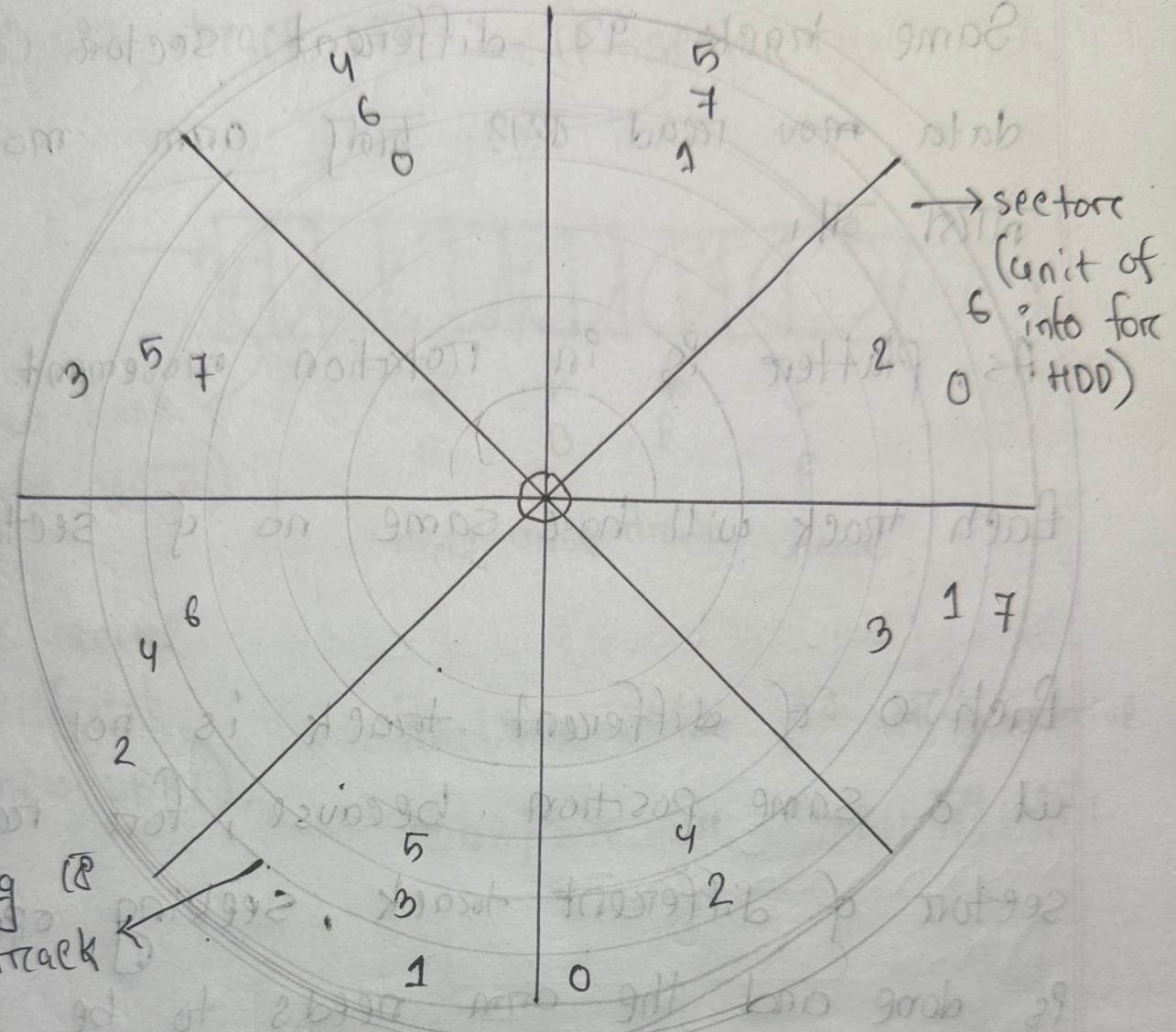


arm mechanism \Rightarrow seeking mechanism

seeking

Arm assembly. Same place \Rightarrow $2T8$, arm just extend





each sector will have equal amount of data always.
So, I have 6 tracks.

When I want to move from one track to another, we need to move the arm. Seek operation.

Same track \neq different sector \neq
data now read \neq bits arm move \neq
not ok.

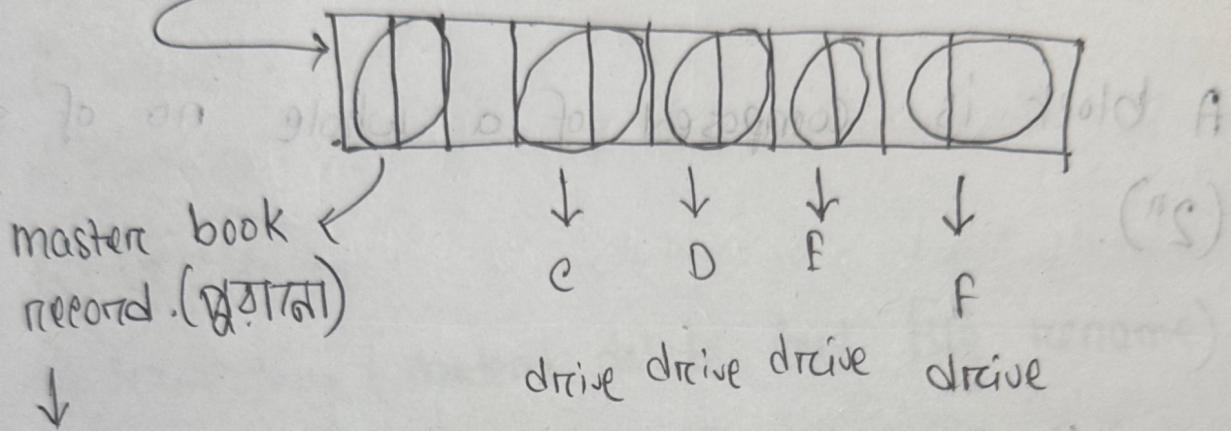
As, platter is in rotation movement.

Each track will have same no of sectors.

Each o of different track is not placed
at a same position because, for reading
sector of different track, seeking operation
is done and the arm needs to be extended/
shrinked but the platter is rotating. So, to
align the arm and the sector, we don't
place the zero at a same position.

* RAM \neq unit 1 byte. $\{$ 1 sector $\}$ \neq
* HDD \neq unit 1 sector. $\} \text{ byte } \neq \text{ 1 byte } \} \text{ 512 bytes }$

HDD वाले track पर 1st sector



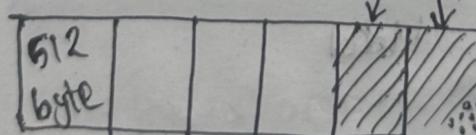
keeps record
of different
OS. new (UEFI)

Partition table वाले सेट, drive
name / size जैसे स्टोर रखते हैं।

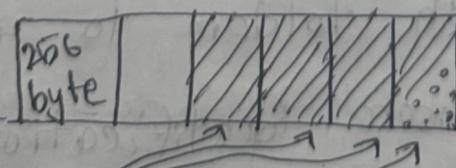
File System Concept

So now
hi.txt
file needs
2 block
to store
in both
HDD.

{ we cannot
change the
size }



Samsung { 1 block = 1 sector }



Toshiba { 1 block = 2 sectors }

hi.txt → 900 bytes

Operating system ने sectors पर फिर किया है।
तो, we make a new unit of OS → a block.

A block is composed of a. whole no of sectors
 (2^n) .
↓ ↓ ↓ ↓
? ? ? ?
↓ ↓ ↓ ↓
quadrants with sub quadrants

sector / physical block / data block = same.

File Attributes

- i) Name : Yo
 - ii) Identifier : unique id found most of the time
 - iii) Type or Extension : .pdf, .mp3
 - iv) Location : where file is stored
 - v) size : 512 byte
 - vi) Protection : which folder to access & not.
 - vii) Time, date & user identification : 1 pm, 11/08/25

File Operations

- i) Create
- ii) Read
- iii) Write
- iv) Delete
- v) Truncating (content delete but file rename)
- vi) Repositioning : changes the position of cursor.

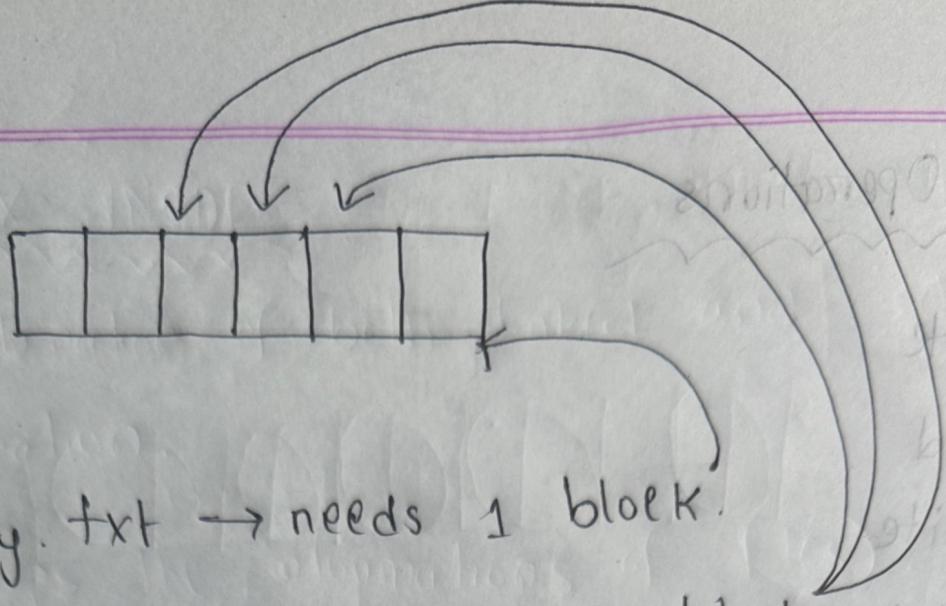
File Allocation Methods

- i) linked allocation
- ii) indexed allocation
- iii) UNIX inode

Allocation methods are of 29 types in zidt

i) Contiguous location

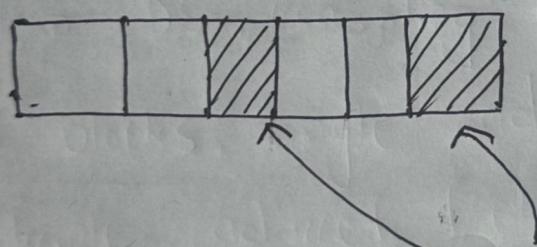
ii) non-contiguous location.



hey.txt → needs 1 block.

hello.txt → needs 3 block.

This is called contiguous (continuous) allocation.



hi.txt → needs 2 block

This is called non-contiguous allocation.

We will study this.

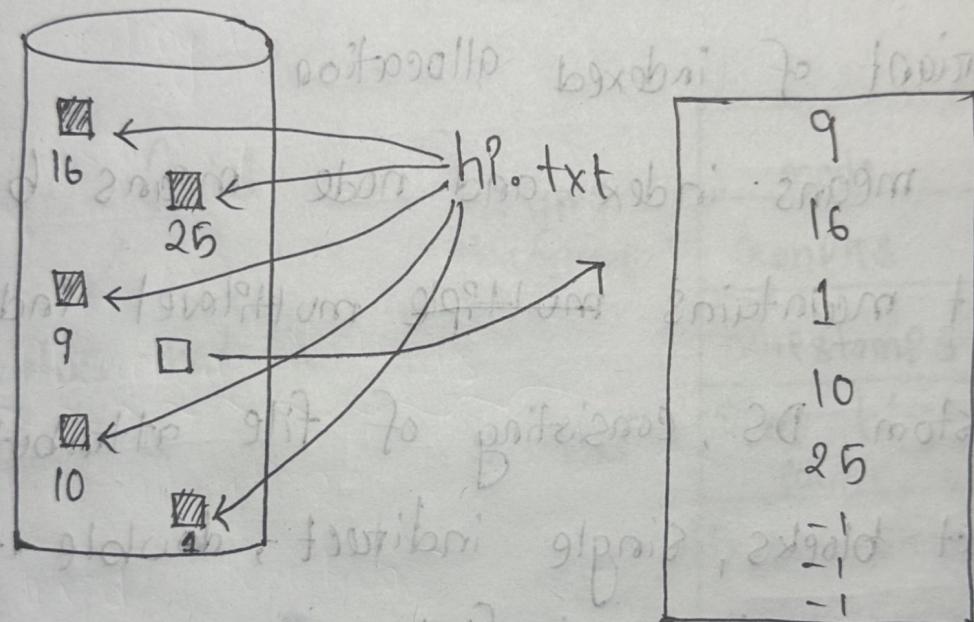
↳ contains a list of numbers to know where are my files are located.

↳ it can be in any order.

↳ -1 caused the list is unused.

↳ my list is fixed.

Non-contiguous allocation (indexed allocation)



This is called index block. Here, highest 8 blocks can be held. We have given -1. So, that we can

Benefits

→ Supports random or direct access.

→ No external fragmentation.

Drawbacks

→ pointer overhead

→ multilevel index when file is too large.

Video: 2

UNIX INODE

Combination of both linked and indexed allocation.

It keeps the informations of files.

- i) Variant of indexed allocation
- ii) I means index and node means block.
- iii) It maintains ~~multiple~~ multilevel index.
- iv) Custom DS, consisting of file attributes, direct blocks, single indirect, double indirect and triple indirect fields.
- v) Direct blocks: Stores such data blocks where each block stores pointer to a block where data is stored.
- vi) Single indirect: Stores such data blocks where each blocks store pointer to a direct block.

vii) Double indirect: Stores such data blocks where each block stores pointers to a single indirect.

viii) Triple indirect: Stores such data blocks where each block stores pointers to a double indirect.

field
↓
metadatas that represent a file.

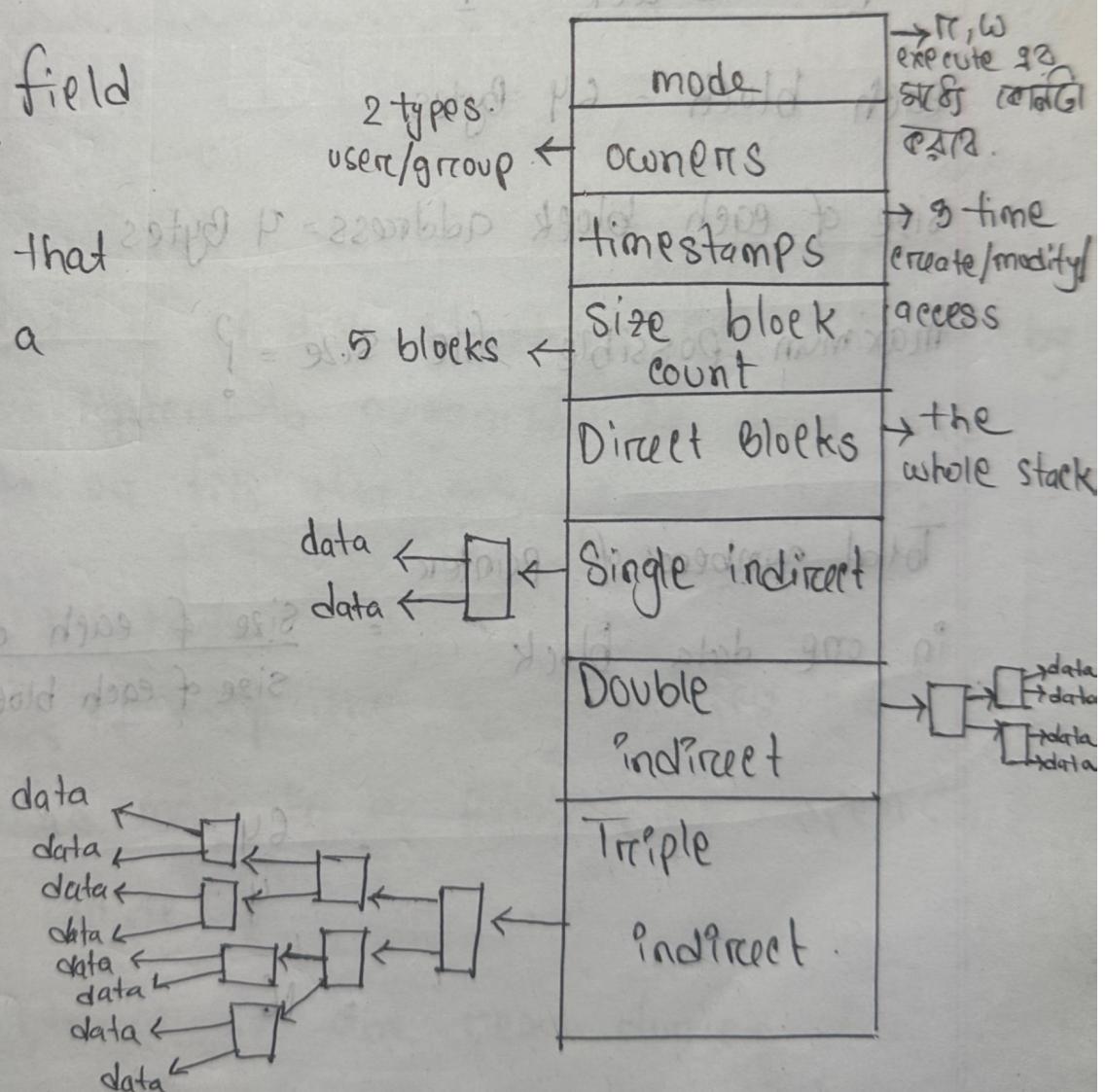


Fig: The UNIX INODE.

Math

UNIX INODE Structure

direct blocks = 4

single indirect = 1

double indirect = 1

triple indirect = 1

each block = 64 Bytes

size of each block address = 4 Bytes

maximum possible file size = ?

Total number of pointers
in one data block

= size of each data block

size of each block address/pointer

$$\begin{aligned} &= \frac{64}{4} \\ &= 16 \end{aligned}$$

$$\text{Total number of pointers} = (4 + 16 + 16^2 + 16^3) = 4372$$

maximum file size = no of total pointers \times block size

$$= 4372 \times 64$$

external disk 279808 Byte

$$\Rightarrow P = 273.25 \text{ kB}$$

at most 632201600 210 87301 911 ←

File System Implementation

2 aspect

→ Data Structures

what type of on-disk structures are utilized by the file system to organize its data and metadata or file attributes.

→ Access methods

how does it map the system calls on operations made by a process as `open()`, `read()`, `write()`

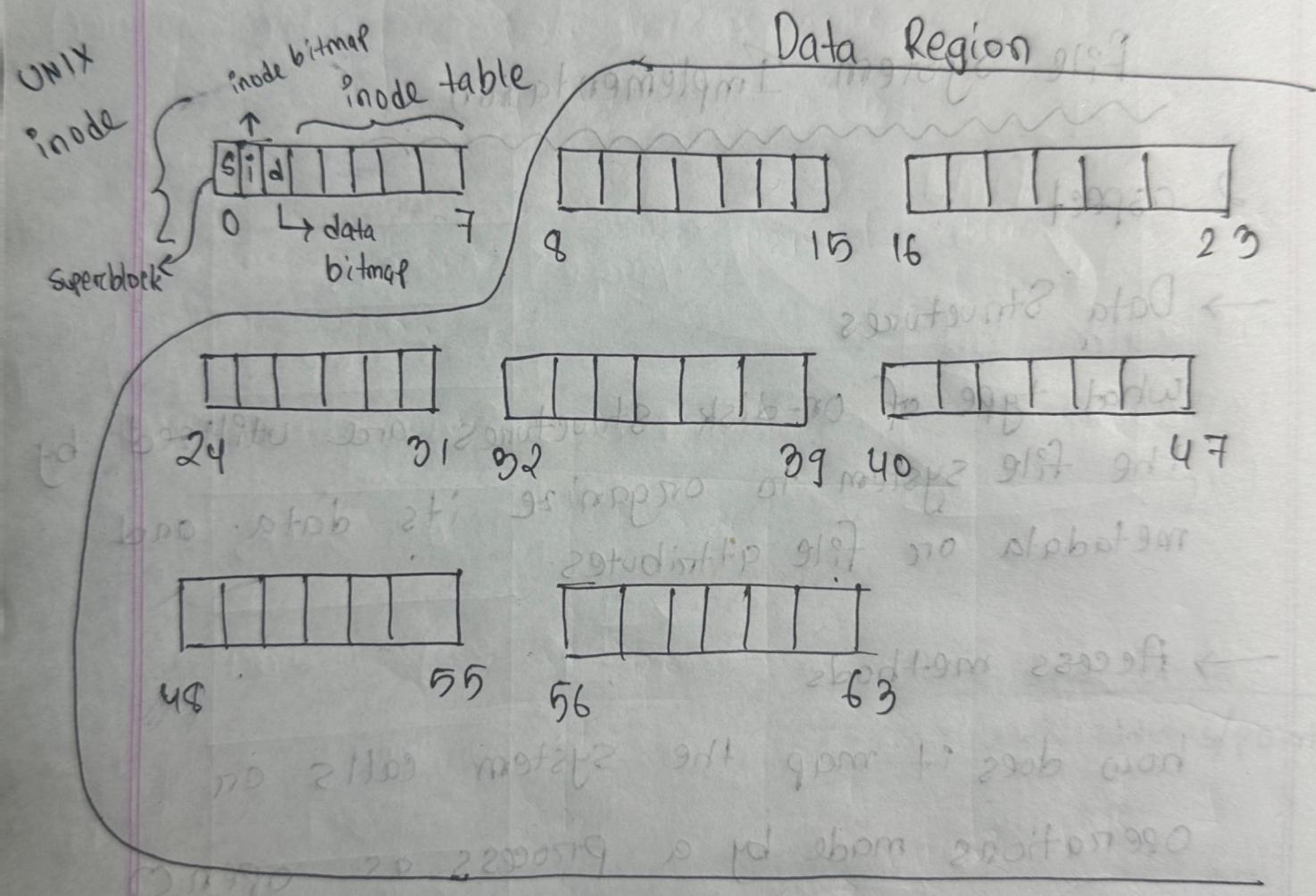
which structures are read during the execution of a particular system call

Overall Organization of VSFS (Very Simple file system)

Divide the disk into data blocks.

→ each block size is 4 kB

→ the blocks are addressed from 0 to N-1.



File organization of the INODE for 256 byte (7)

(*) inode number = 32

Size of inode = 256 bytes

$$\text{offset} = 32 \times 256 = 8192$$

↳ मान आज्ञा jump के क्रमागामी मानों

byte address or location of inode no = the inode table

location of 32 no inode = 12KB + 8192 B

12KB + 8 KB

= 20 KB.

block 0 block 1 block 2

Super	i-bitmap	d-bitmap	0	1	2	3	16	17	18	19	32	33	34	35	48	49	50	51	64	65	66	67
			4	5	6	7	20	21	22	23	36	37	38	39	52	53	54	55	68	69	70	71
			8	9	10	11	24	25	26	27	40	41	42	43	56	57	58	59	72	73	74	75
			12	13	14	15	28	29	30	31	44	45	46	47	60	61	62	63	76	77	78	79

0 8 16 24 32
4 12 20 28

(*) Disk are not byte addressable, sectors addressable.

Disk consist of a large no of addressable sectors,
(512 Bytes)

fetch = inode no 32

size of inode 256 byte

size of block = 4 kB

sector address is the physical address of the
inode block

$$\text{block} = \frac{\text{inode no} \times \text{size of inode}}{\text{size of block size}}$$

$$= \frac{32 \times 256}{4}$$

$$= \frac{8192}{4}$$

= 2048

$$= \frac{8 \text{ kB}}{4 \text{ kB}}$$

= 2

Page 1 block 2.

$$\text{sector address} = \frac{(\text{block no} \times \text{block size}) + \text{inode table start address}}{\text{Sector size}}$$

C orbit

$$= \frac{1}{2} (2 \times 4) + 12 \} \text{ k} \Omega$$

(as above, sum of total no. of teeth in each row) \Rightarrow 512 B

$$= 40 \cdot 2919 \& \text{ last position has } \leftarrow$$

product of different factors not. (i) $\left\{ \begin{array}{l} \text{odd} \\ \text{even} \end{array} \right\}$

Product of factors not even \therefore (ii) $\left\{ \begin{array}{l} \text{odd} \\ \text{even} \end{array} \right\}$

(product mod 100) \Rightarrow last two digits of product

Term	2919	291991	most abstr
.	S	P	Q
..	S	P	Q
ost	P	P	Q
mod	P	P	Q
product	F	P	P

transposition method

add 6 to 6 and 8 to 8 make it 12 and 16 \Rightarrow 100 no. of teeth

and 20 no. of teeth \Rightarrow 100 no. of teeth

Video 3

Directory Organization / Structure

- contains a list of (entry name, inode no) pairs
- each directory has 2 files
 - i) . for current directory
 - ii) .. for parent directory

for example

Directory has 3 files. (foo, bar, foobar)

inode num	name len	stolen	name
5	4	2	.
2	4	3	..
12	4	4	foo
13	4	4	bar
24	8	7	foobar

Free Space Management

- file system tracks which inode and data block are free or not.
- In order to manage free space, we have 2 simple

bitmaps.

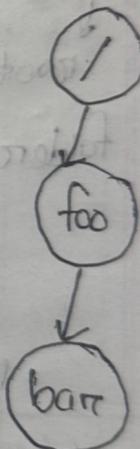
↳ when file is newly created, it's allocated bitmap inode by searching in the inode bitmap and update on disk.

↳ Pre-allocation policy is commonly used for allocate contiguous blocks.

Access Paths: Reading a File from Disk

Issue an open ("/foo/bar", O_RDONLY)

→ root
directory



traverse recursively the path name
until the desired inode 'bar'.

→ issue read() to read from file {system call}.
read the first block of the file, consulting the
inode to find the location of such a block.

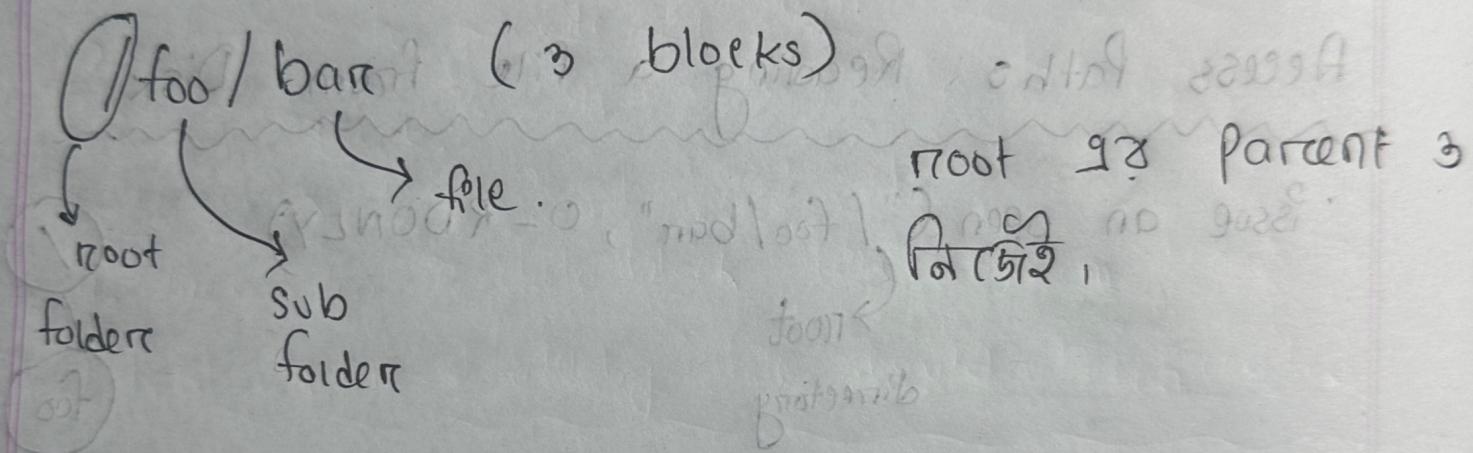
- update the inode with a new last accessed time.
- update in-memory open file table for file descriptor,

→ When file is closed.

file descriptor should be allocated, but for now, that is all the file system really needs to do.

No disk / I/O take place

→ avoid duplicates



read() → call once for each block.

target is to read bar.

file read() করতে প্রথম আবে ওপেন() করতে হবে।

1. open file → read mode একটি ফাইলের মধ্যে

2. to find inode no of bar, read datablock of parent directory.

3. to find data block of x & head inode of x.

- root is pnode no fixed for pibmap of
- /foo/bar → create bar and keep 3 blocks.

① Create

② write

	data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data	bar data	bar data
open (bar)			read 1			read 2				
				read 3				read 4		
					read 5					
read ()					read 6			1st read 7		
read ()						read 8				
					read 9			2nd read 10		
								read 11		
						read 12				
								read 13		

File Read Timeline.

Access Paths : Writing to Disk

1. reading root inode and retrieve into root data block.
2. from root data block it will retrieve into regarding foo inode.
3. from foo inode will learn about foo data block.
4. from foo data block file system learn that on which directory file should be created.
5. then file system need inode bitmap & write "inode" into for new bar file.
6. update directory of foo.

	data inode	inode root	foo inode	bar inode	root	foo data	bar data 0	bar data 1	bar data 2
Create (/foo/bar)		read 1			read 2				
			read 3			read 4			
				write 5		write 6			
					read 7				
write()	read 8				read 9				
	write 10					write 11			
write()	read 11				read 12				
	write 13					write 14			
write()	read 15				read 16				
	write 17					write 18			
							write 19		
								write 20	

File Creation Timeline.

Crash Consistency and Journaling

like hard-disk crash ঘৰছি,

data রাখিব নাই,

* lseek() → file এই pointers indicate কোথা কোথা,

1) Just the data block is written to disk. But

not update bitmap_inode, data bitmap, inode
কেন disk কে point কৰে না,

2) Just the updated inode is written to disk.

But data block update কোথা নাই, file

System 'Inconsistency' happened. inode data
block কে point কৰে কোথা block কে point
কৰে নাই,

3) Just the updated bitmap is written to the
disk.

bitmap indicates that data block is allocated
but inode কোন data block কে point কৰে নাই, space leak

4) Inode & data bitmap are written to disk. But

data block 9 points to block 10 which has been freed.

5) Inode and data block are written. But data bitmap points to block 10 which is free.

6) data bitmap and data block are written. But not the inode block is stored last position of it.

Solution

The file System Checker

* fseck → is a unix tool for finding inconsistency and repair them.

* Super block → super block copy एवं मालिक file size एवं file corrupt एवं असर
alternate copy of super block
fseck use 2nd iteration only (ii)

* free blocks → inode, indirect blocks, double indirect block एवं check करना भी फ्री ब्लॉक की

* inode state → consistency check ~~eg. about 10%~~ (P)

* inode link → check if the reference count of each inode is consistent.

→ check if a block is shared by two inodes.

→ check for 'bad' block pointer is the one that points to the location that lies outside the file partition.

* directory: check if . & .. are properly set up. Make sure that there are only one hard link for a directory.

Drawbacks of FSCK

i) It's slow. It's overhead 2TCP

ii) Not a solution to journaling.

iii) Very complex knowledge to build.

Video: 04

Journaling

A (মেমো) B (ডি.সি) file copy হওতে একে file system
অন্তর্বর্তী জায়গামুখে নিখুঁত সাধারণ নে A point (মেমো)
B point রে copy হব। Direct copy হও না, নিখুঁত

পথ copy হওয়া পাই অন্তর্বর্তী / পথ, This is called structure
that is the destination of the 'write ahead' is
called "log."

* Journaling itself has a super-block.

Super	Journal	inode	Data bitmap	Data block
-------	---------	-------	----------------	---------------

* If update ২টা অন্তর্বর্তী journal রে update হওয়া,
যদি,

TxB → entry in journal. Transaction begin.

TxE → out in journal. Transaction End.
physical logging

Journal	TxB	I [V ₁]	B [V ₂]	Db	TxE
old status					

এখন এখনো legit Disk রে update ২টা,

Replay → when failure happens at the time of TXE. Then, after the failure, we have to update the things of journal in DISK again.

Disk is strict Disk scheduler. Can reorder the operation. And, the block loses.

To avoid being inconsistent.

When, inode and Data bitmap matches with each other, then it is consistent.

Journal doesn't care about data.

→ Journal write → write the contents of the transaction to the log.

→ Journal commit → write the transaction commit block.

→ Checkpoint → writes the contents of the update to their locations.

not serializable in 2nd out group goes to 1st.

Recovery

→ if the crash happens, before the transaction is written to the log then the pending update is skipped.

not present in log appends previous

→ if the crash happens after the transaction is written to the log, but before the checkpoint

recover the updated as follows.

• Scan the log and look for transactions that have committed to the disk.

• transactions are replayed

• fail back to original or to latest log
• hold queue in read format

exit	pxt	exit	sxt	ixt	logcat
					msg

Batching Log Update

- If we create two files in same directory, the same inode and the directory entry block is to the log and committed twice.

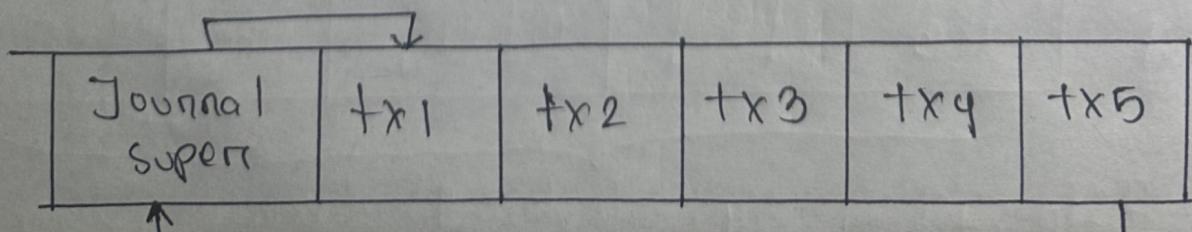
- To reduce excessive write traffic to disk, journaling manage the global transaction.

- Write the content of the global transaction forced by synchronous request.
- Write the content of the global transaction after timeout of 5 seconds.

Making the log infinite

The Journal is a circular linked list.

Journal has a super block.



Journal does not include Data blocks because it doesn't care. So, space reduces. of data blocks.

This is Meta Data Journal. Means information of Journal.

"circular log"

$$= (\text{index} + 1) \% \text{array-length}$$

Metadata Journaling

Because of the high cost of writing every

data block to disk twice.

- commit to log journal

- checkpoint to on-disk location.

→ Data write

→ journal metadata (inode & bitmap) write

→ journal commit

→ checkpoint metadata

→ free

To start writing data Disk is called flushing.

old slab p. about 2092 id. 0109 free

new free block 1000 slab id is 0109

Tricky Case: Block Reuse

directory

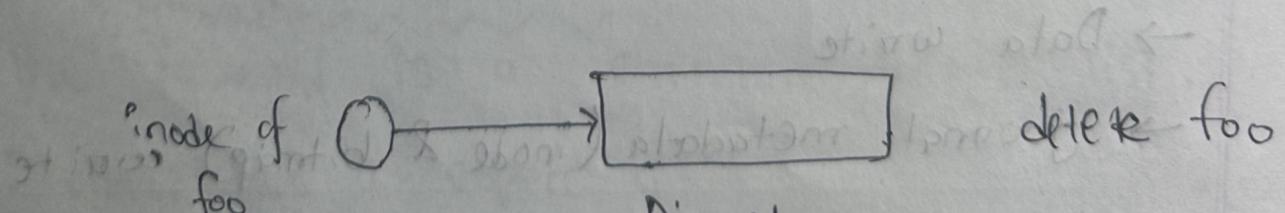
foo inode

direct : 1000

owner : CSE 321

modified :

TxB	inode (foo)	Data (foo)	TxE	TxB
id=1	bitmap 1000	bitmap 1000	id=1	id=2



Directory free ←
block of foo.

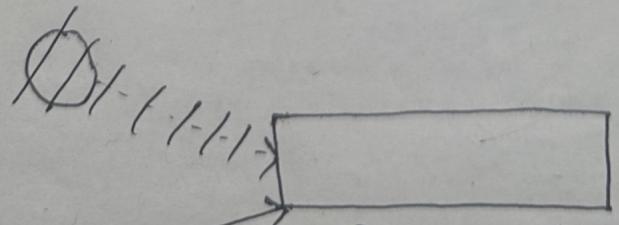
data of foo is deleted, but the inode is still there. After power on, 1000 data bitmap is empty as foo is deleted. But, inode is there. So, we need to remove the inode too, like pointer. Then assign a file to it.

foobarr	TXF	
1000	id: 2	

Now,

inode of
foo

inode of
foobarr



Directory block
of foo.