

Filesystem  $\rightarrow$  OS subsystem  $\rightarrow$  secondary or smaller systems part of a larger system  
 $\rightarrow$  provides abstraction like files and directories  
 $\rightarrow$  hides complexity of underlying storage devices

### Filesystem Interfacing

$\rightarrow$  Processes identify files through a file handle / file descriptor.

$\rightarrow$  In UNIX  $\rightarrow$  the POSIX API is used to access files, devices, sockets, etc.

Input / Output libraries

fopen, fclose, fread, fprintf, ...



System call API

open, close, read, write



Files

Devices

Sockets

What is the mapping b/w  
library functions and system  
calls?

`fopen()` → `open()`

`fclose()` → `close()`

∴

---

open : getting a handle

`int open (char *path, int flags, mode_t mode)`

→ Access  
permission  
check performed  
by OS (ls -al)

access mode

O - RONLY

O - RDRW

O - WRONLY

O - SEARCH

O - EXEC

if flags  
contains

O\_CREAT

specify file  
creation

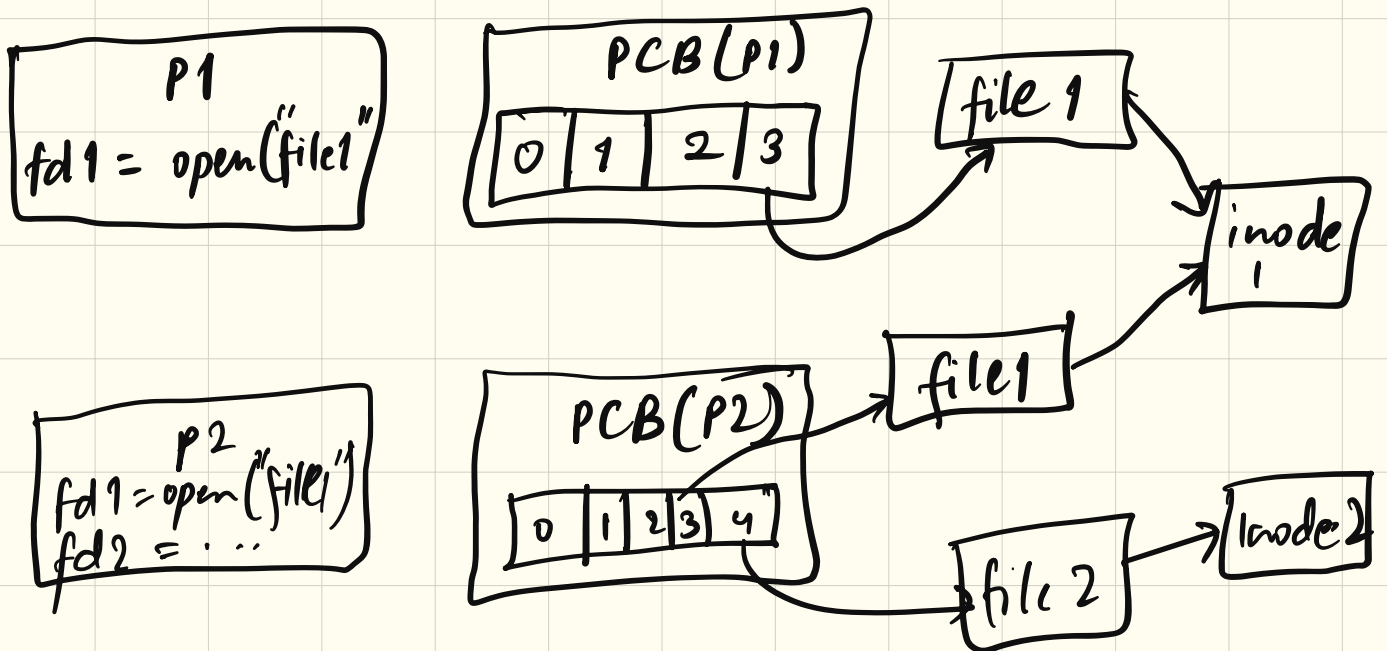
(rwx) using  
mode (owner  
group, other)

⋮

open → returns a non-negative integer  
↓  
file descriptor

-1 on error  
sets errno to indicate error

## Process View of File





→ Per process file descriptors table with pointer to a "file" object.

→ file object → inode  
many  
one

### Questions:

1) What do file descriptors 0, 1 and 2 represent?

2) What happens to the FD and the file objects across fork() ?  
— exec?

3) Can multiple FDs point to the same file object?

# Read and Write

ssize\_t read ( int fd , void \* buf ,  
size\_t count );

file  
handle

pointer  
to the  
buffer where  
read data will  
be stored

# of bytes  
to read

read() → returns # of bytes  
actually read  
→ can be smaller  
than count

0  
⇒  
EOF  
reached  
-1  
error

ssize\_t write ( int fd, void \* buf, size\_t count);

↙  
pointer to  
the data to  
be written

↓  
no. of  
bytes  
to write  
from the  
buffer

→ return no. of  
bytes written

file descriptor :

0 → STDIN

1 → STDOUT

2 → STDERR

lseek

file handle

off\_t lseek (int fd, off\_t offset, int whence);

signed integer type

target offset

SEEK\_SET,  
SEEK\_CUR,  
SEEK\_END

On success, returns offset from the starting of the file

Examples

→ lseek (fd, 100, SEEK\_CUR) → forward the file position by 100 bytes

→ lseek (fd, 0, SEEK\_END)

↓  
file position at EOF  
returns the file size

→ lseek (fd, 0, SEEK\_SET)

↪ file position at the beginning of the file

## File Information (stat, fstat)

```
int stat (const char * path, struct stat *sbuf);
```

→ returns information about the file pointed to by path

→ Information is filled up in the structure stat.

Example:

```
struct stat sbuf;
```

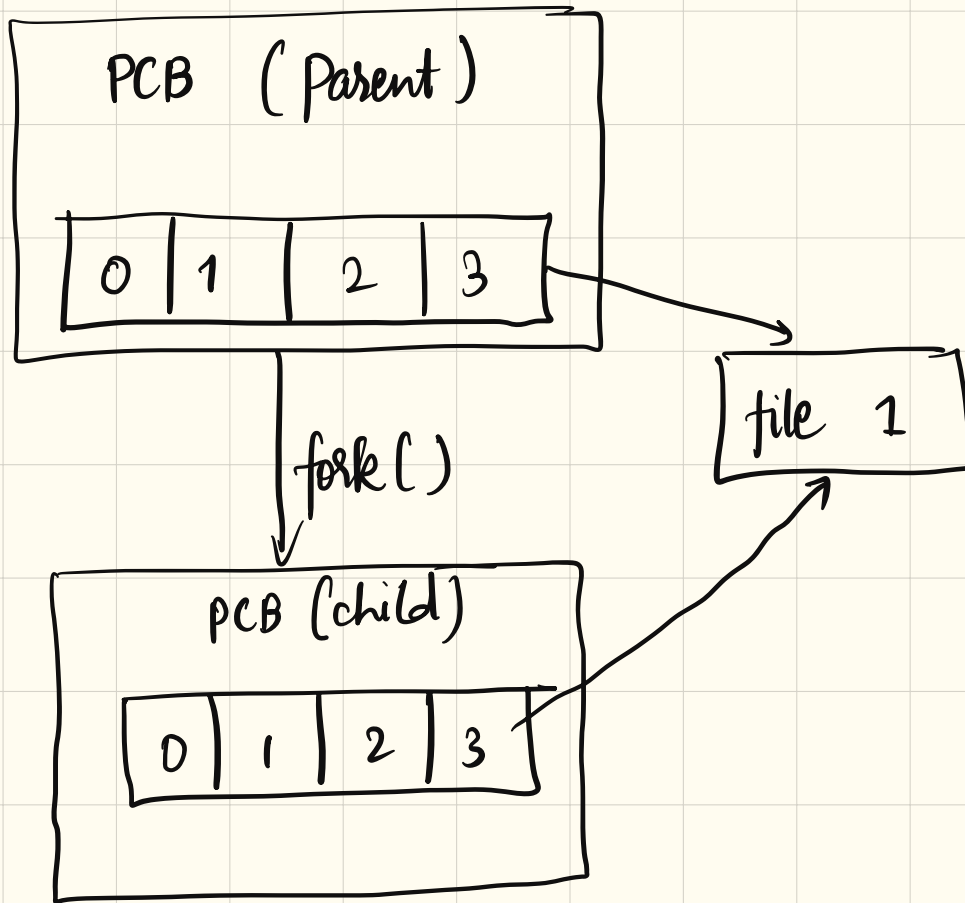
```
stat ("/home/user/tmp.txt", &sbuf);
```

```
printf ("inode = %d, size = %ld \n",  
        sbuf->st_ino, sbuf->st_size);
```

other useful info in stat

st\_uid, st\_mode

## Process view of file after fork()

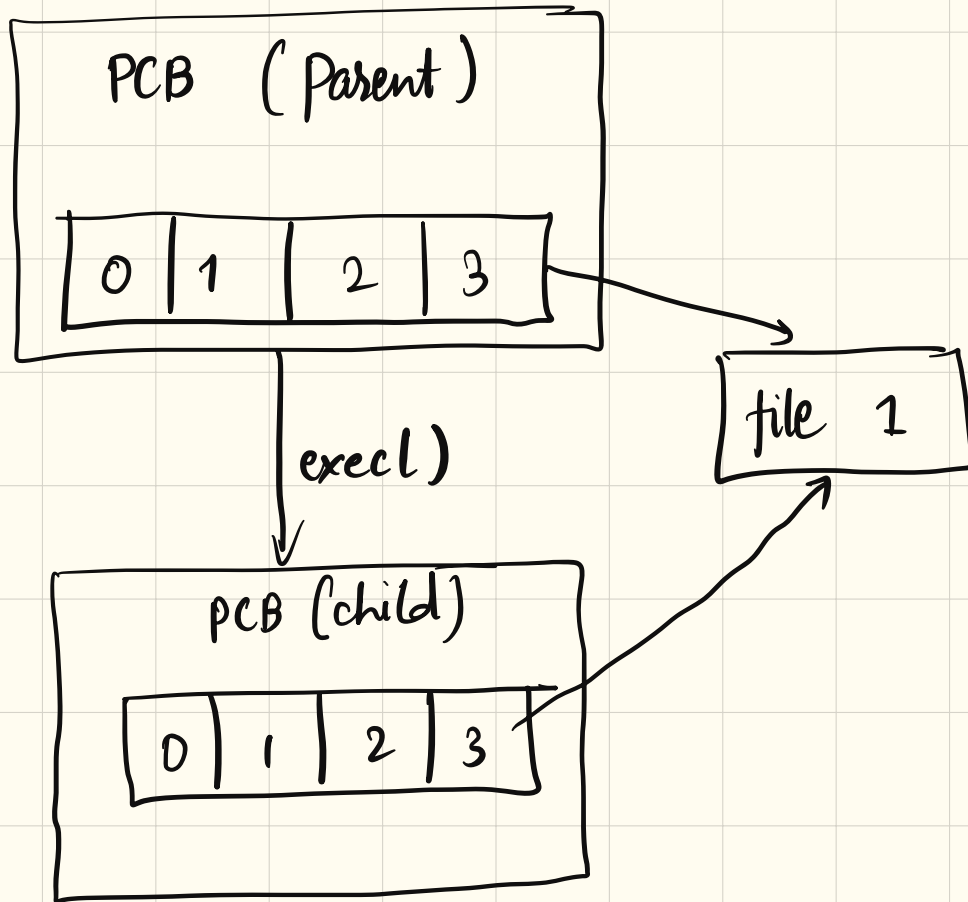


→ On `fork()`, child inherits the parent process state.

→ All file descriptors remain open in the child

→ FDs point to the same file object.

## Process View of file after exec()



→ On `exec()`, destroy memory state of calling process to load new binary

→ By default, it does not destroy FD tables

→ FDs point to the same file objects as pointing earlier, and remains open

→ To close FDs on exec, specify `O_CLOEXEC` flag during open

\* What happens to FD and the file objects across

— fork?

→ the FD table is copied across fork → file objects are shared.

— exec?

→ Open files remain shared by default.



## Duplicate file handle (dup and dup2)

```
int dup (int oldfd);
```

→ The `dup()` system call creates a copy of the file descriptor `oldfd`.

→ Returns the lowest numbered unused descriptors as the new descriptors.

→ Old and New FD → represent the same file

```
int fd, dupfd;  
fd = open ("tmp.txt");
```

```
close (1);
```

```
dupfd = dup (fd);
```

```
printf ("Hello World \n");
```

↘ Will be written to `tmp.txt`

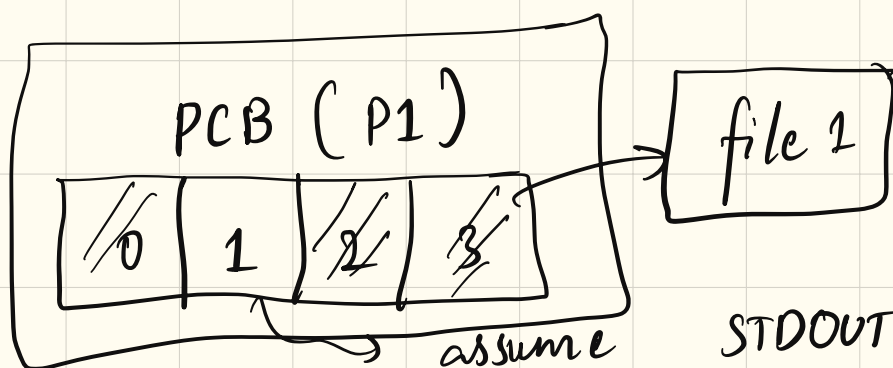
→ `dupfd` will be 1 (assuming `STDIN = 0` is open)

P1

```
fd 1 = open ("file 1");
```

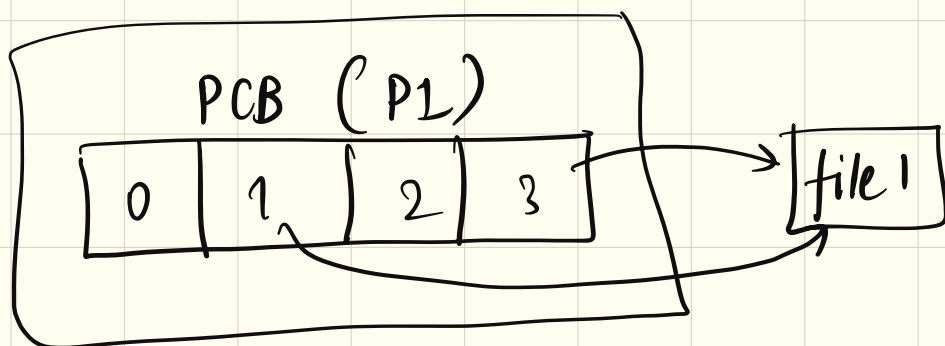
```
+ dup (fd 1)
```

Before dup()



STDOUT is closed  
before

After dup()



→ Duplicate descriptors share the same  
file state

→ Closing one file descriptor does not  
close file 1

```
int dup2 ( int oldfd, int newfd );
```

→ close newfd before duping the file descriptor oldfd

→ dup2 ( fd, 1 ) equivalent to

→ close ( 1 );

→ dup ( fd );

## Use of dup ( ) : Shell redirection

Example : ls > tmp.txt

Implementation :

```
fd = open ( " tmp.txt " );
```

```
close ( 1 ); // close STDOUT
```

```
close ( 2 ); // close STDERR
```

```
dup ( fd ); // 1 → fd
```

```
dup ( fd ); // 2 → fd
```

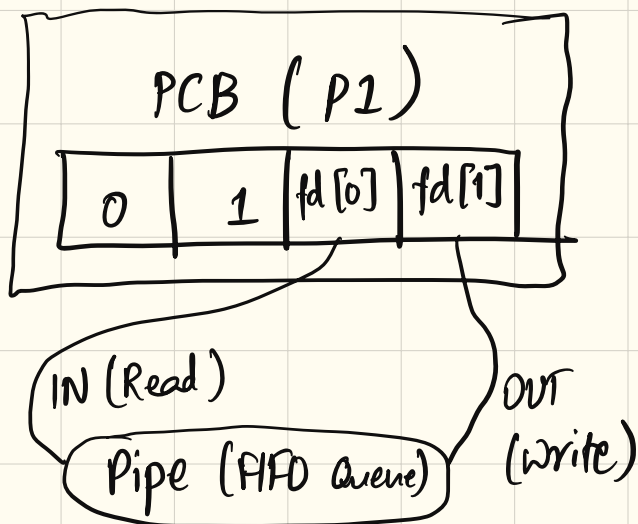
```
exec ("ls");
```

## UNIX pipe system call

P<sub>1</sub>

```
int fd[2];
```

```
pipe (fd);
```



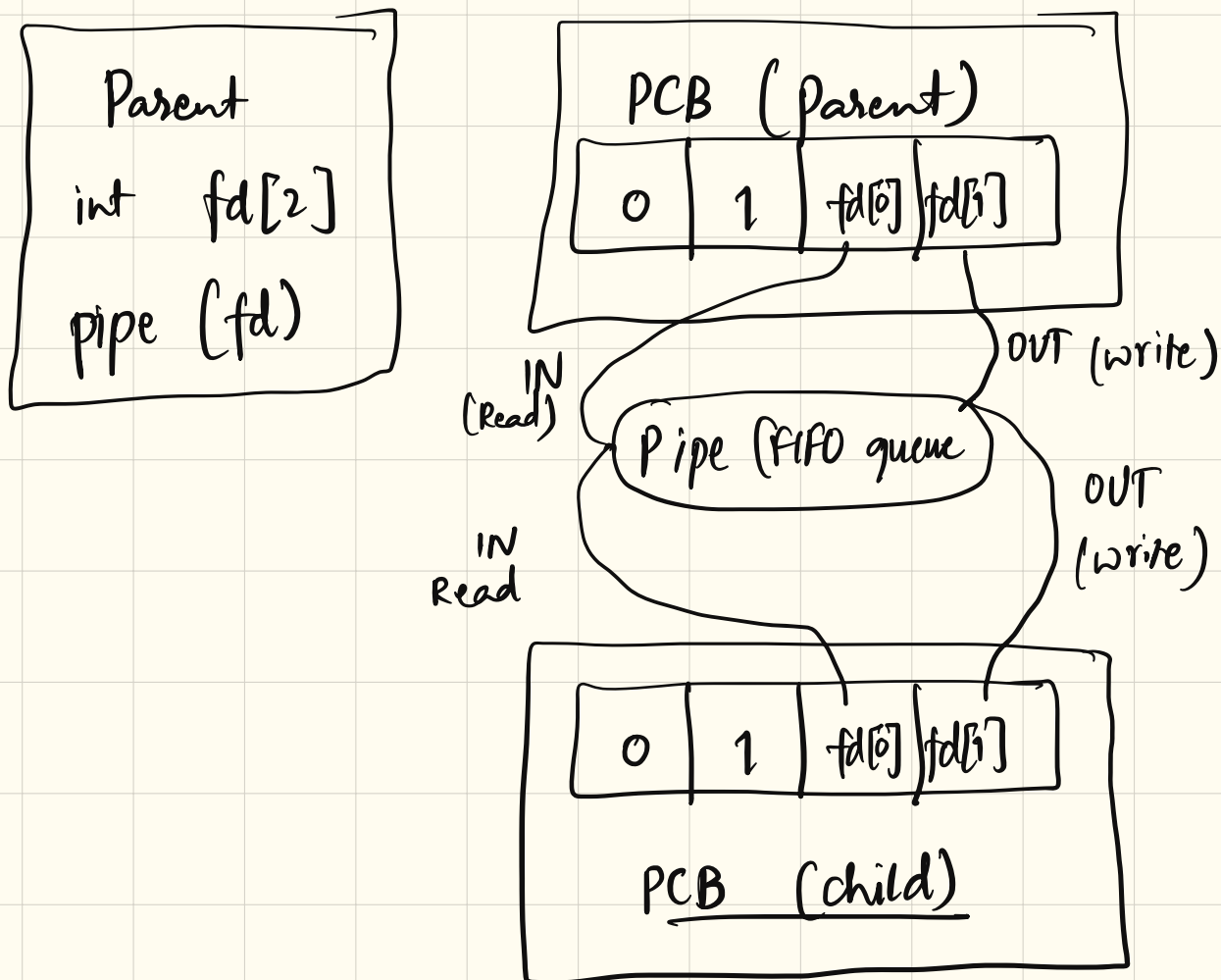
pipe() takes array  
of two FDs as  
input

→ fd[0] ~ read  
end of  
the  
pipe

→ fd[1] ~ write  
end of  
the pipe

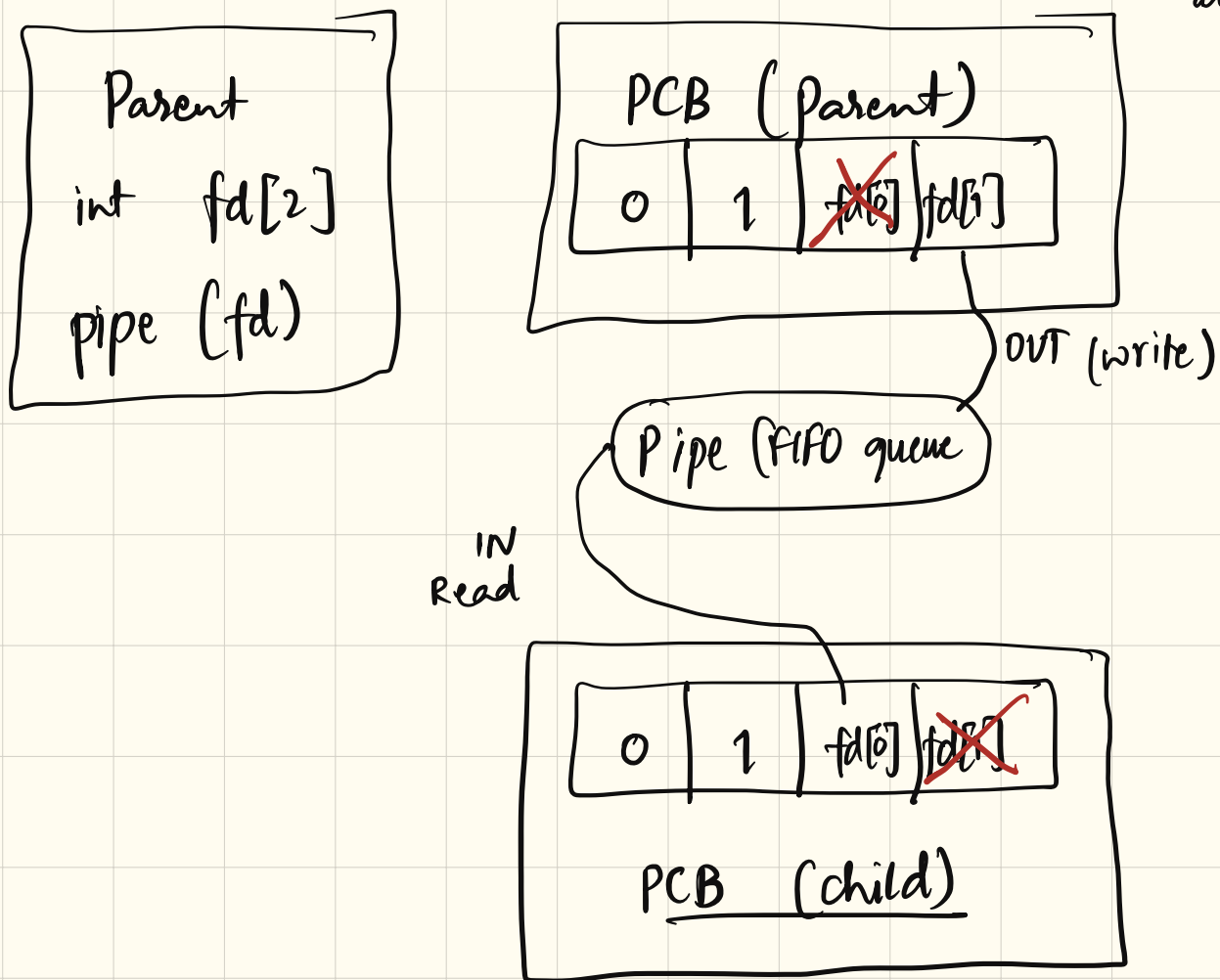
→ Implemented as a  
FIFO queue in  
OS.

→ fork duplicates the file descriptors



→ At this point, both parent and child can read/write to the pipe.

\* close() one end of the pipe, both in parent and child



Result:

→ a queue b/w parent and child.

Shell piping: `ls | wc -l`

→ `pipe()` followed by `fork()`

→ Parent: `exec("ls")` after making  
STDOUT → out fd of the pipe (using dup)

→ Child : `exec("wc")` after closing  
STDIN and duping in fd  
of pipe.

Result : input of `wc` is connected  
to output of `ls`.

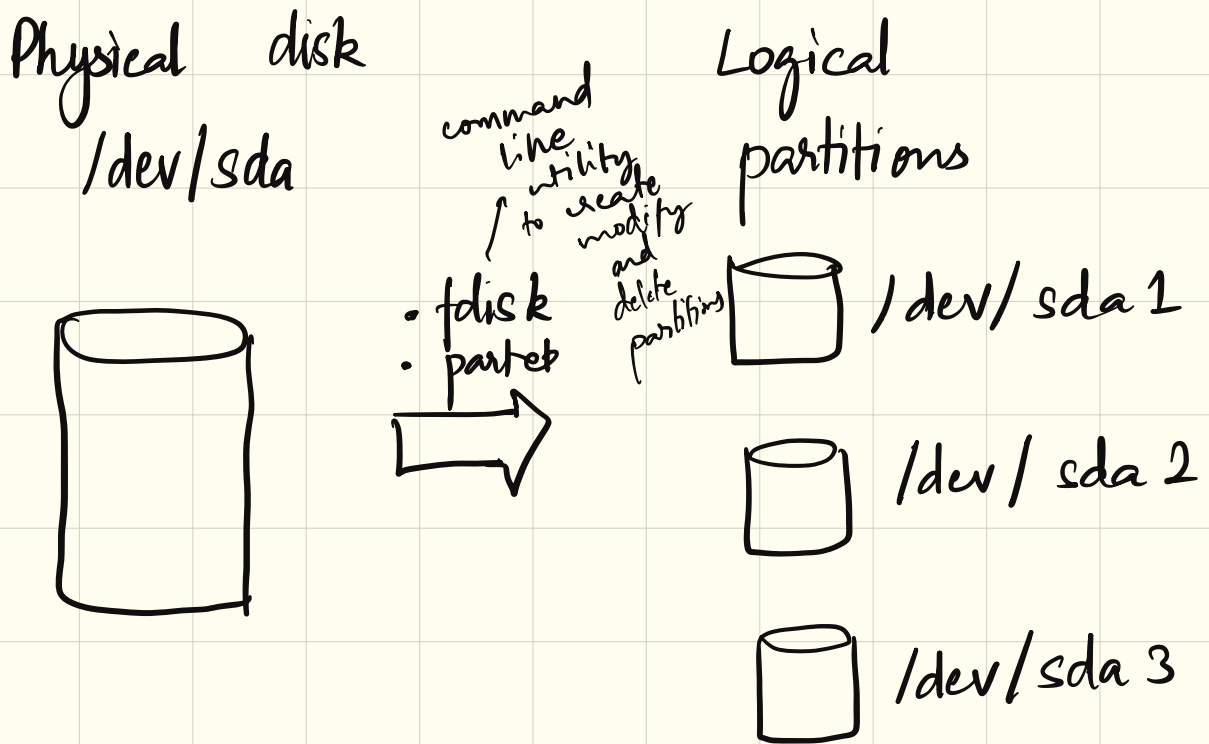
## Filesystem Implementation

### Step 1 : Disk Device Partitioning

→ Partition is stored in the boot  
sector of the disk.

→ Creation of partition is the first  
step → does not create  
filesystem

→ A filesystem → created on a partition to manage the physical device and present the logical view.

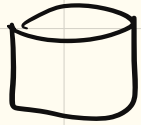


→ All filesystems provide utilities to initialize the filesystem on the partition (e.g. MKFS)

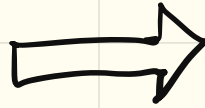


## Step 2: filesystem Creation

Logical partitions



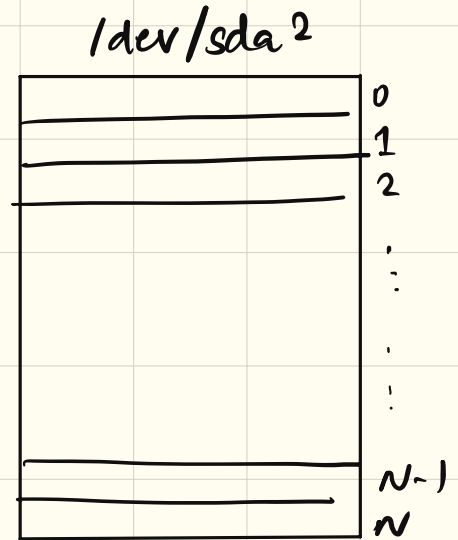
/dev/sda 1 mkfs EXFS /dev/sda2



/dev /sda 2



/dev/sda 3



→ MKFS creates initial structures in the logical partition

→ creates entry point to the FS  
superblock

→ the filesystem is ready to be mounted.

## Step 3: Filesystem mounting

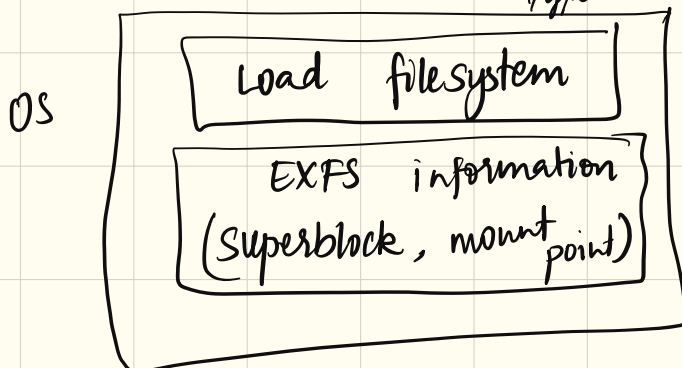
→ `mount()` associates a superblock with the filesystem mount point

→ example: the OS will use the

OS associates a superblock with FS mount point } superblock associated with the  
mount point `"/home"` to reach any file/dir under `"/home"`

USER: `mount -t exfs /dev/sda2 /home`

↓  
`mount [ "/dev/sda2" , "/home" ,`  
device path ↙ ↘ mount point  
"EXFS" , flags,  
↓ ↙ ↘  
filesystem type fs-options )



## Structure of an example superblock

```
struct superblock {  
    u16    block_size;  
    u64    num_blocks;  
    u64    last_mount_time;  
    u64    root_inode_num;  
    u64    max_inodes;  
    disk_off_t    inode_table;  
    disk_off_t    blk_usage_bitmap;  
    ...  
};
```

→ Superblock contains information regarding the device and the filesystem organisation in the disk.

→ Pointer to different metadata related to the filesystem

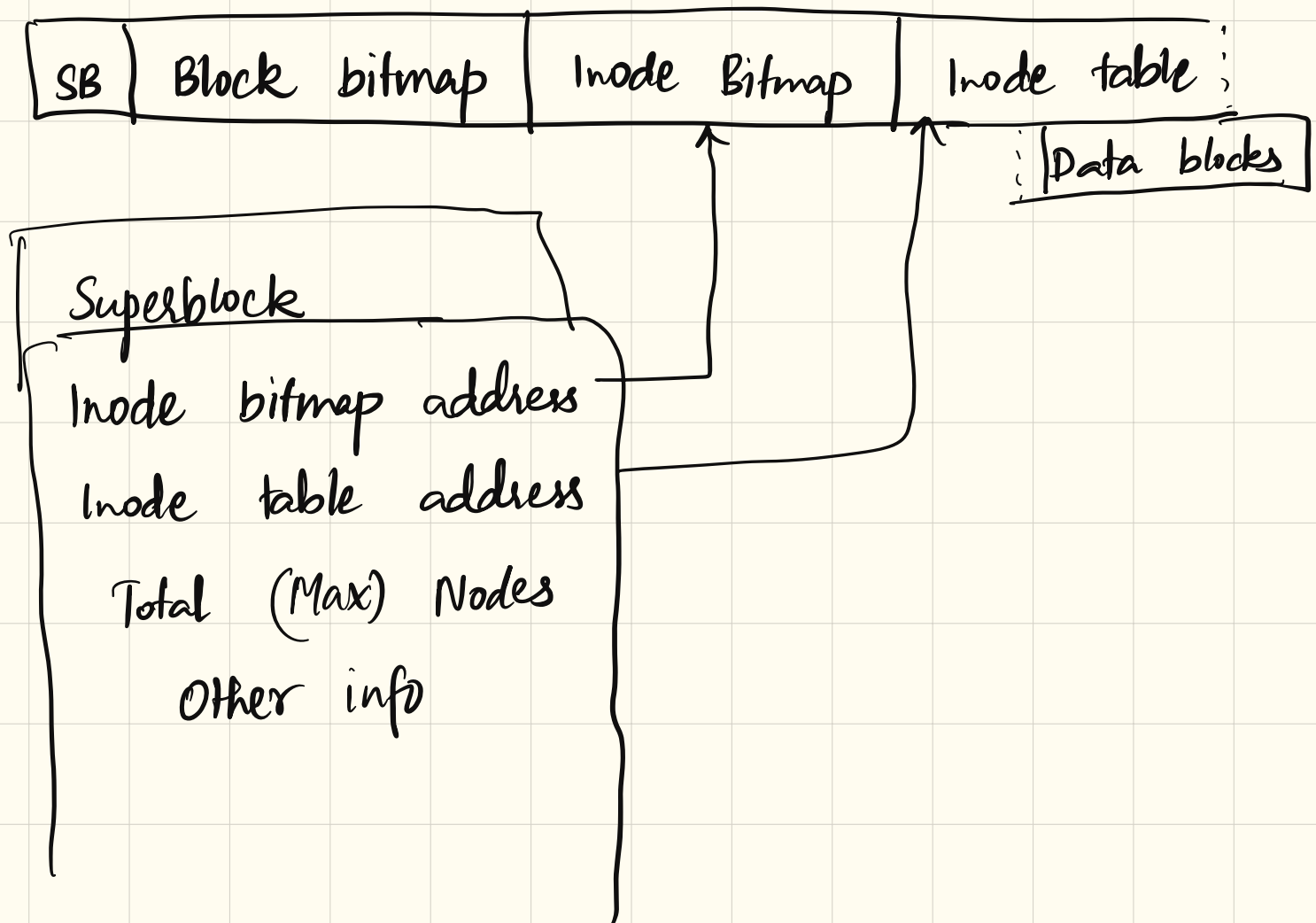
→ ex: list of free

blocks is required

before adding data

to a new file/directory

## Filesystem organization



```

inode_t * get_inode (SB *sb, long info)
{
    inode_t * inode = alloc_mem_inode();
    read_disk (inode, sb → inode_table
               + info * size_of
                 (inode));

    return inode;
}

```

→ given any inode number, loads the inode structure into memory

### Questions:

\* filesystem is mounted, the inode no. for root of the filesystem (i.e., the mount point → known, root node can be accessed

- How to lookup/search files/dir under root inode?
- How to locate the content in disk?
- How to keep track of size, permissions, etc?

## Inode (UNIX)

- On-disk structure
- contains information regarding files/directories in UNIX systems
- unique number in the fs  
"ls -li filename"
- contains access permissions, access time, file size, etc.
- IMPORTANTLY → information regarding file data location on the device

→ Directory inodes → also contain information regarding its content (structure)

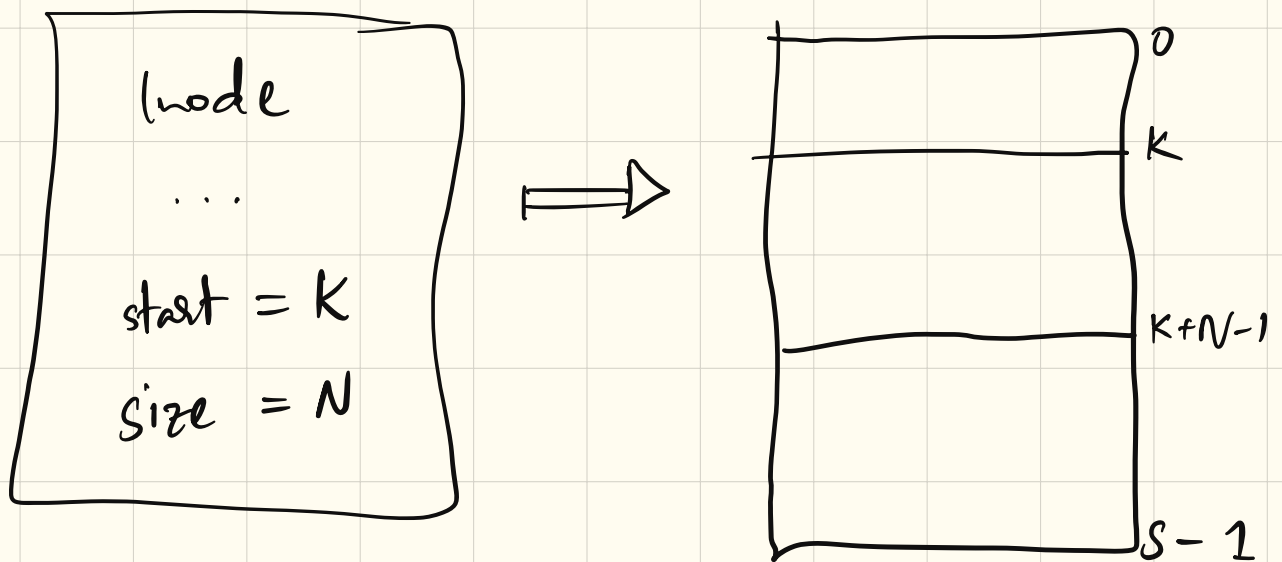
Problem: File Offset to Disk Address

Mapping

How to efficiently translate file offset to device address?

- file size : few bytes to GBs
- can be accessed in sequential or random manner
- How to design mapping structure.

## Contiguous allocation



→ Works nicely for both sequential and random access

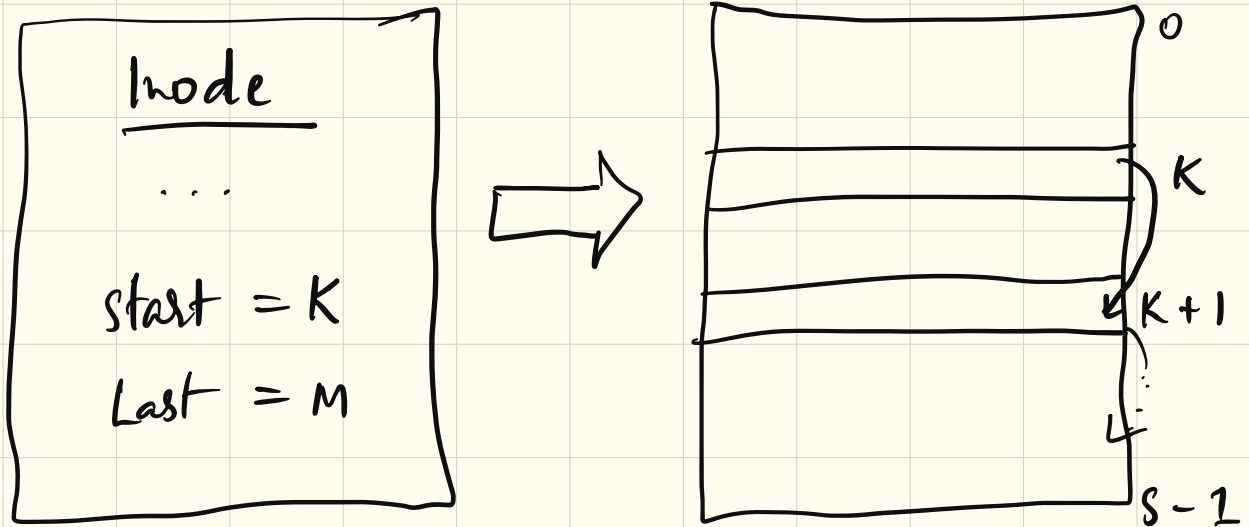
→ Append operation is difficult.  
how to expand files?

↓  
Relocation

→ External fragmentation



# Linked Allocation



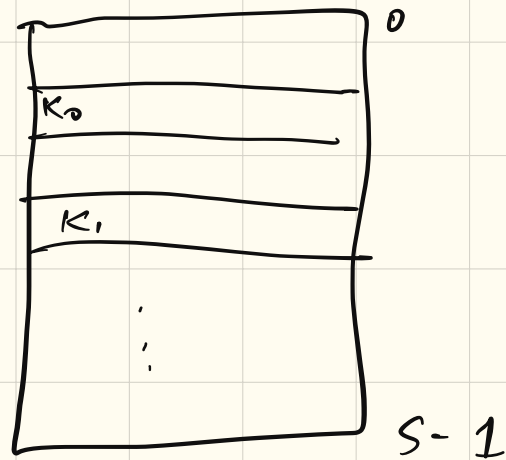
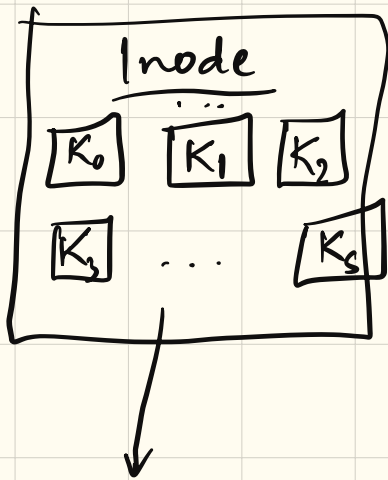
→ Every block contains pointer to next block.

→ Advantage: flexible, easy to grow and shrink

→ Disadvantage: random access

→ Why last? → efficient append operation

# Direct Block Pointers

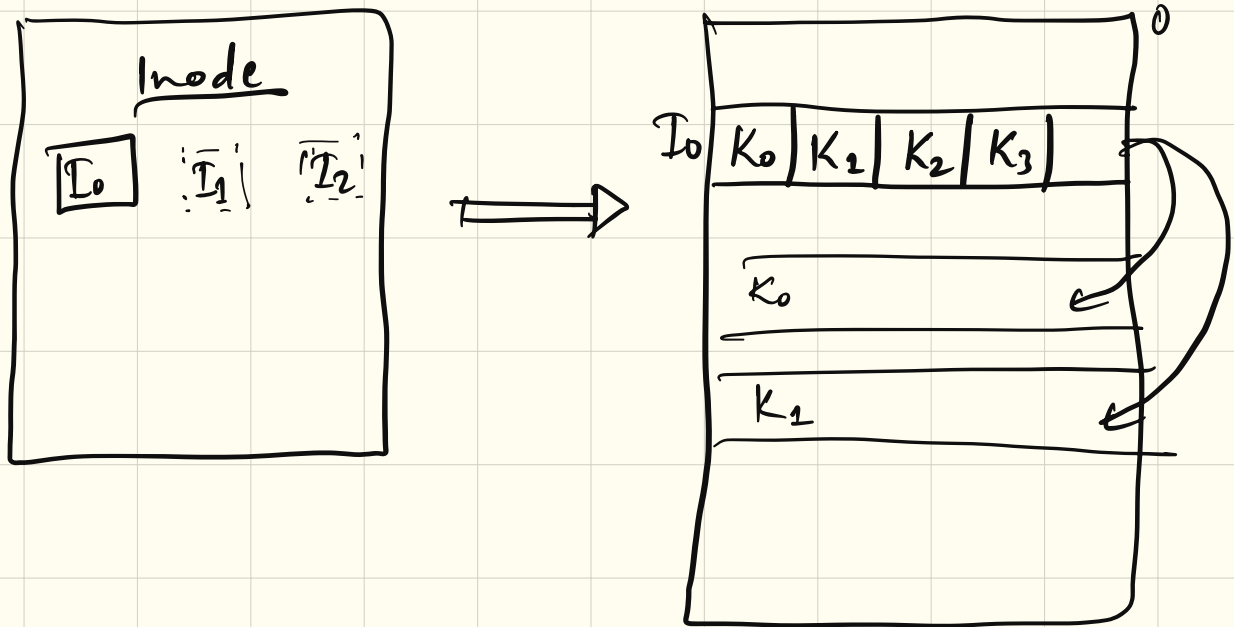


→ direct pointers  
to block

→ flexible growth, shrink, random access  
is good

→ Cannot support files of large size

# Indirect Block Pointers



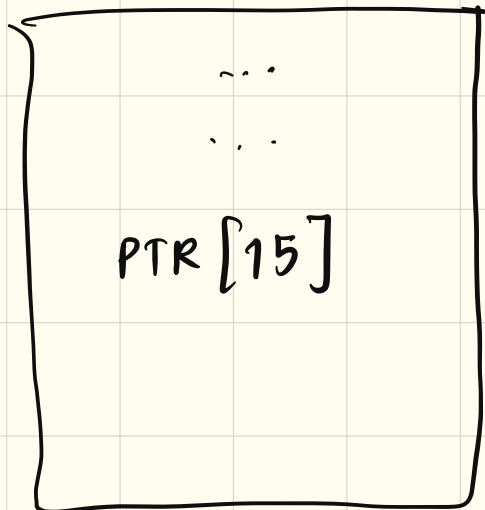
→ Inode contains pointers to a block containing pointers to a data block

→ Advantage: flexible, random access is good

→ Disadvantage: Indirect block access overheads (even for small files)

# Hybrid Block Pointers : Ext 2/3 FS

Ext 2/3 inode



Direct pointers {PTR[0] to PTR[11]}



→ File block address (0 to 11)

Single indirect pointers

{PTR[12]}

file block address

(12 to 1035)

Double indirect

{PTR[13]}

→ 1036 to 1049611

Triple Indirect

\* How to locate content in disk?

→ Inode structures in inode are used to map file offset to disk location

\* How to keep track of size permissions etc?

→ Inode maintains these info

## Organizing directory content

```
struct dir_entry {  
    inode_t  inode_num;  
    char     name[FNAME_MAX];  
};
```

→ fixed size directory entry simple way to organize

→ Advantage : avoid fragmentation,  
rename

→ Disadvantage: space wastage;

Variable size:

```
struct dir_entry {  
    inode_t inode-num;  
    us      entry-len;  
    char    name[nam-len];  
}
```

→ contain length explicitly

→ Advantage: less space wastage (compact)

→ Disadvantage: inefficient rename,  
require compaction

\* How to search/lookup files/dir under root inode?

→ Read the content of the root inode and search the next level dir/file

### Caching and Consistency

Filesystem maintains several metadata structures like superblocks, inodes, directory entries to provide a filesystem abstraction like files, directories

How to lookup/search files/directory under root inode?

- Read the content of the root inode and search the next level dir using the name and find out its inode number
- Read the inode to check permissions and repeat the process.

## Filesystem and Caching

- Accessing data and metadata from disk impacts performance
- Many file operations require multiple block access:  
e.g: opening a file  
normal shell operations: ls



- Executables, config files, library, etc. are accessed frequently.
- Many directories containing executables config files are also accessed frequently.
- Metadata blocks storing inodes, indirect block pointers are also accessed frequently

Can we store frequently accessed disk data in memory?

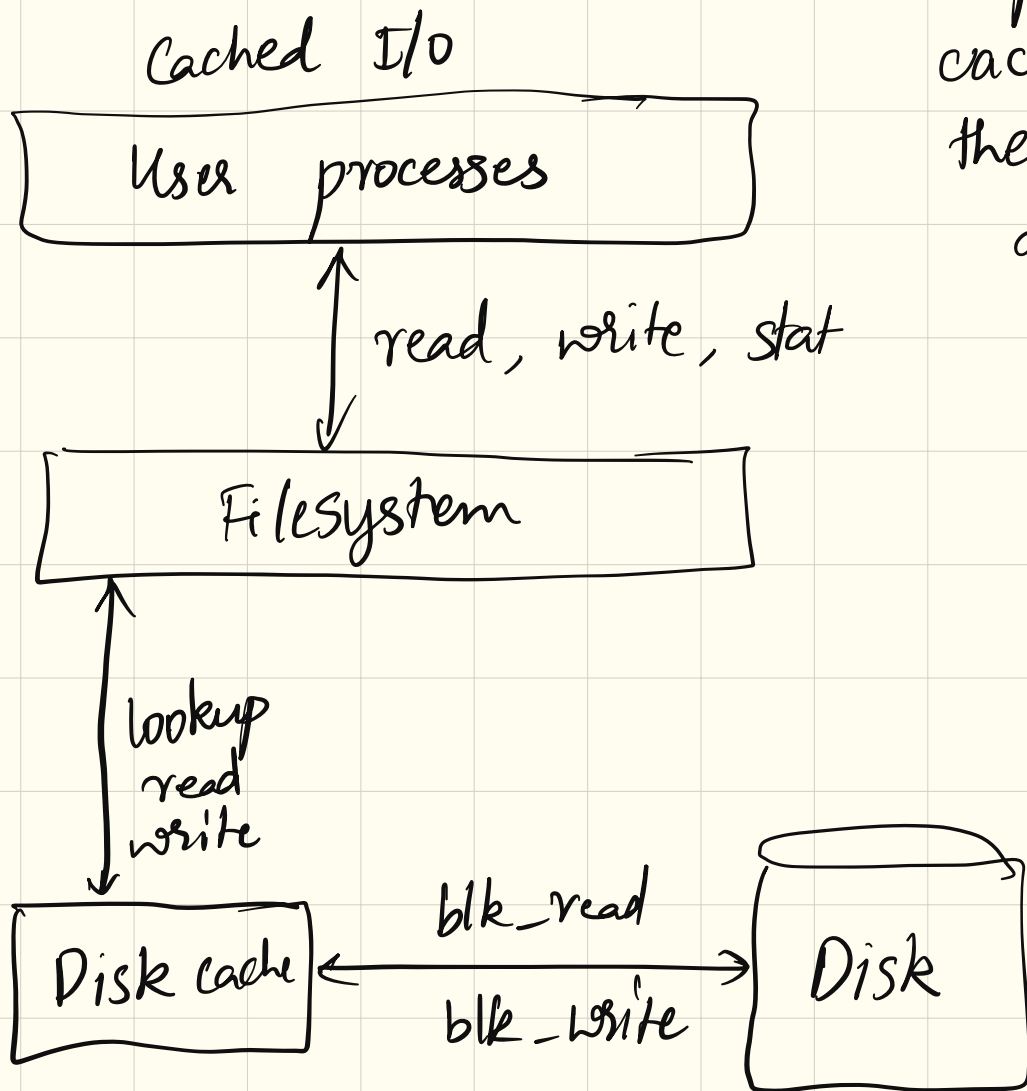
→ What is the storage and lookup mechanism?

→ Are the data and metadata caching same?

- How is the cache managed?
- Eviction policy?
- Complication?

## Block Layer Caching

→ lookup memory cache using the block numbers as the key.



① For data caching : file offset to block address mapping is required before using the cache

② Works fine for metadata as they are addressed using block numbers

### File Layer Caching (Linux Page Cache)

→ Store and lookup memory cache using {inode # , file offset} as the key.

→ For data, index translation is not required for file access.

→ Metadata may not have a  
file association → should be  
handled  
(using a special  
inode maybe)

For storage and lookup mechanism



File Layer caching is desirable  
as it avoids index accesses  
on hit, special mechanism  
required for metadata

∴ How is the cache managed?

→ lookup memory using block #  
or { inode, file offset } as  
the key.

Any eviction policy  
can be used

If page is dirty, write to  
the disk first.

---

Caching may result in inconsistency:

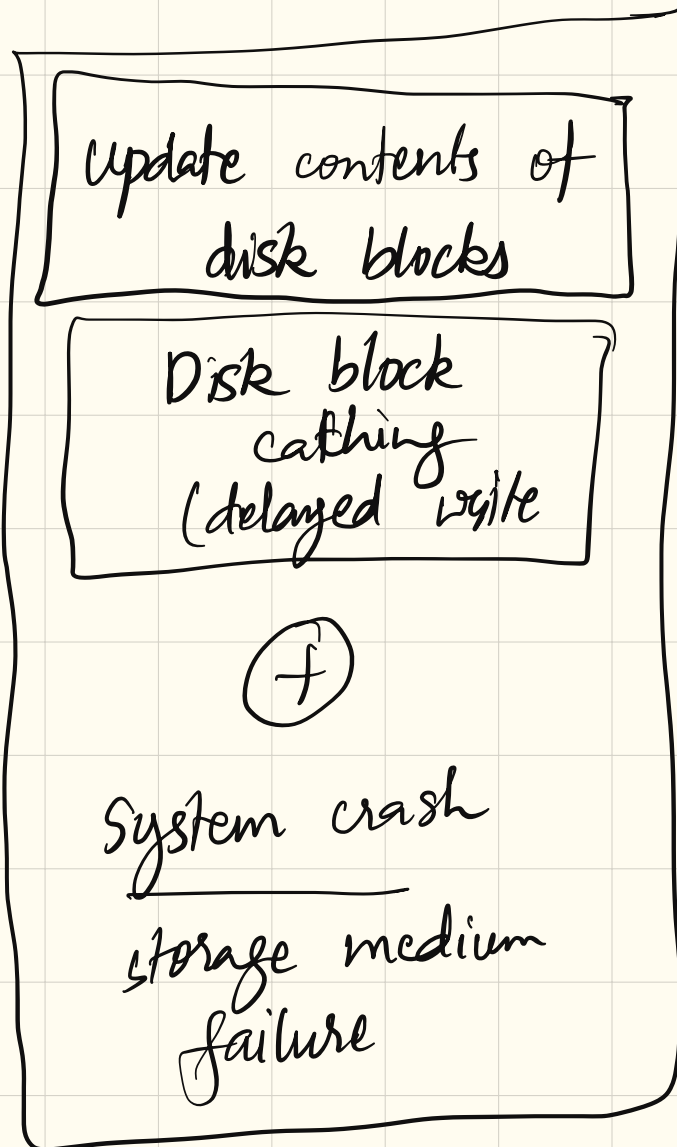
Ex 1: If a write() → successful  
data must be  
written

Ex 2: If a file creation → successful  
→ data must  
be created

block  $\leftarrow$  inode } bitmap must be set

directory entry  
contains  
an inode  $\rightarrow$  inode  
must  
be valid

## Root causes of inconsistency



$\Rightarrow$  possible  
inconsistent  
system

$\rightarrow$  No issues if  
only read  
operations

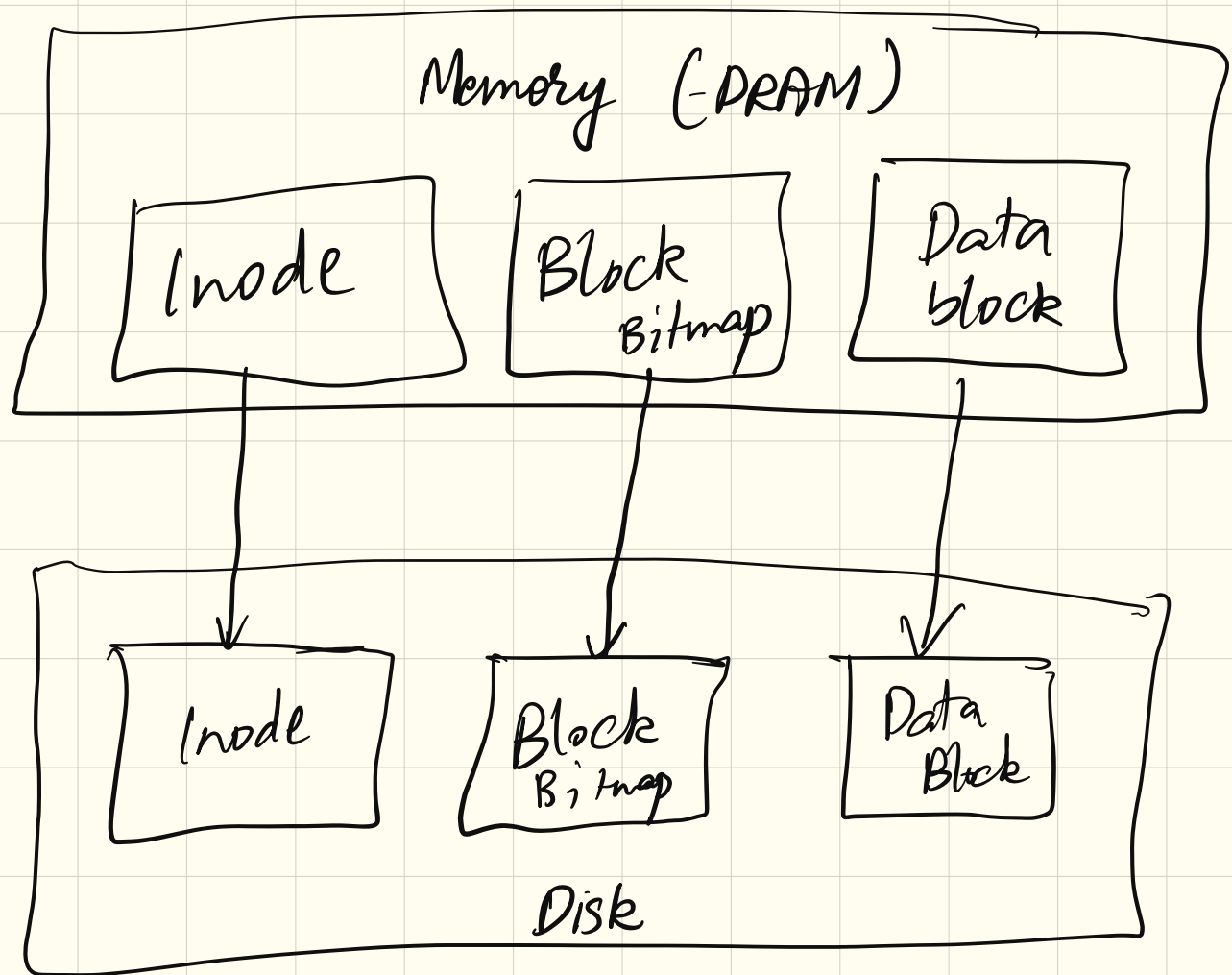
## Example: Append a file

- ① seek to the end of file
- ② allocate new block
- ③ write user data

→ Inode modification: size & block pointers

→ Block bitmap update: set used block bit for newly allocated block(s)

→ Data update: Data block content is updated



Three write operations required to complete operation. What if some of them are incomplete?



Written	Yet to be written	Implications
Data block	Inode, Block bitmap	Filesystem is consistent (lost data)
Inode	Block bitmap, data block	Inconsistent (correctness x)
Block bitmap	Inode, Data block	Inconsistent (space leakage)
Data block Block bitmap	Inode	Inconsistent (space leakage)
Inode Data block	Block bitmap	Inconsistent (correctness)
Inode Block bitmap	Data block	Consistent (Incorrect data)

→ careful ordering of operations  
may reduce risk of inconsistency

## Filesystem consistency with fsck

→ Strategy: Do not worry about  
consistency, recover  
after abrupt failure

→ During FS mount, check if  
it had been cleanly  
unmounted when it was  
last used. How to know?

→ Maintain last unmount  
info on superblock.

→ If the FS was not clearly  
unmounted, perform  
sanity checks at different  
levels: superblock,  
block bitmap,  
inode, directory  
content

→ Sanity checks and verify invariants  
across metadata.

→ Example:

Block bitmap vs. inode block  
ptrs

used inodes vs. dir content

# Filesystem Consistency with Journaling

\* Idea : Before the actual operation, note down the operation in some special blocks journal

- ❑ Journal entry for append operation: [start] [Indoe block] [Block bitmap] [Data block] [End]
- ❑ When the FS is updated (in a delayed manner) and marked the journal entry a completed (a.k.a checkpoint)
- ❑ Recovery mechanism: Journal entries are inspected during the next mount and operations of non-checkpointed entries are re-performed

- Journal write should not only be synchronous but also performed in the specified order (especially the end marker)
- Failure after updating some blocks and rewritten during recovery is not an issue as the data is consistent at the end
- Failure during journal write is not a problem w.r.t file system consistency

Journal entry a completed (a.k.a checkpoint)

# Metadata Journaling: Performance-reliability Tradeoff

- ❑ Journaling comes with a penalty, specially for maintaining data in journal
- ❑ Metadata Journaling: Data block is not part of the journal entry
  - Practical with tolerable performance overheads
- ❑ Example journal entry for append: [start] [Indoe block] [Block bitmap] [End]
- ❑ Strategy: First write the data block (to disk) followed by the journal write and metadata commit afterwards, Why?
  - If the metadata blocks are not written, FS can be recovered
  - If journal write fails, a write is lost (syscall semantic broken)