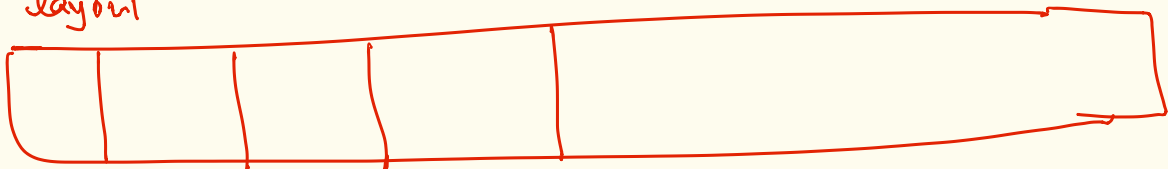


Lecture # 30(?)

CS 347/333

⊕ some more file systems

disk layout



super block

inode bitmap

data bitmap

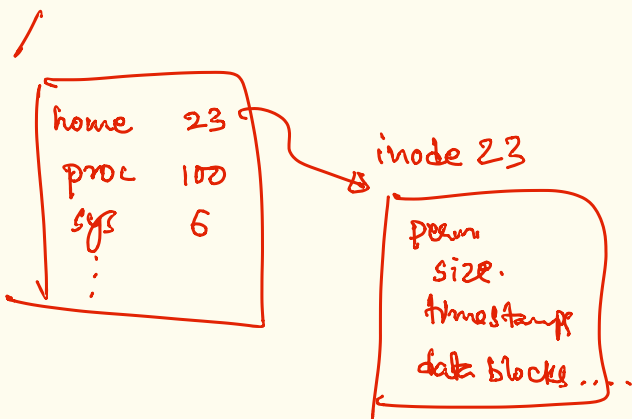
inodes

data blocks

⊕ every file maps/has/represented by an inode.

(HW) whether an inode object/structure has a filename field in it?

Q. What is a directory? — it is a file?

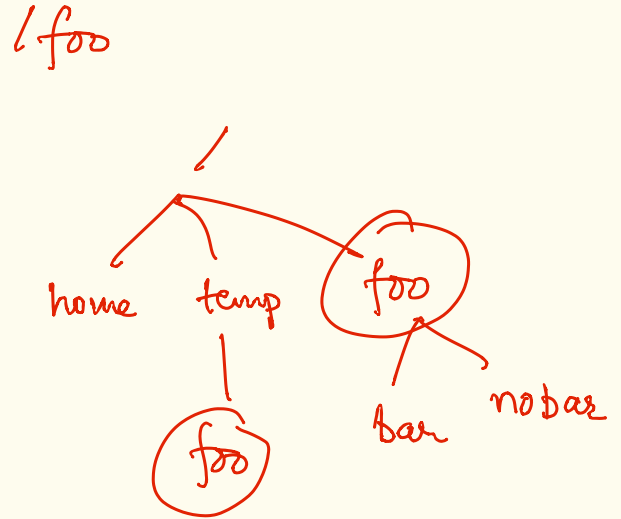


/home/
/sys
/proc

What FS actions are invoked on file accesses via the FILE API.

fd = open("/foo/bar");

(disk) read root inode
read root data block
read foo inode
read foo data
read bar inode



5 disk accesses to open a file.

CPU > Mem > IO

open("/1/2/3/4/5 - .../100/what really.txt");

~ 200+ disk accesses!

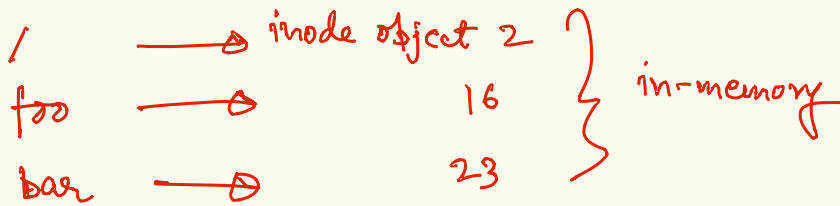
FS caching techniques

- (i) dentry caches (directory entry) } metadata
- (ii) inode caches } metadata
- (iii) page cache (buffer/disk/data cache) } data

* (i) dentry cache (in-memory)

maps directory entries to inode objects.

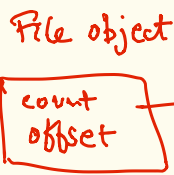
/ foo / bar



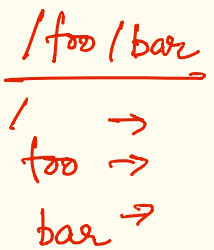
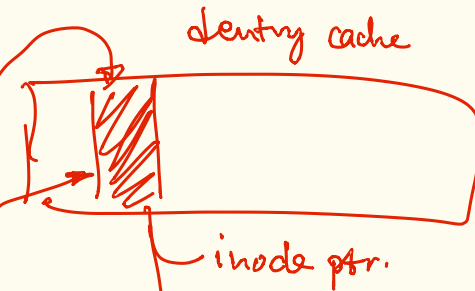
(ii) inode cache (in-memory)

maps inode number to inode objects.

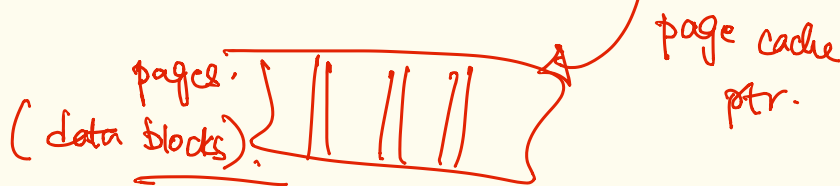
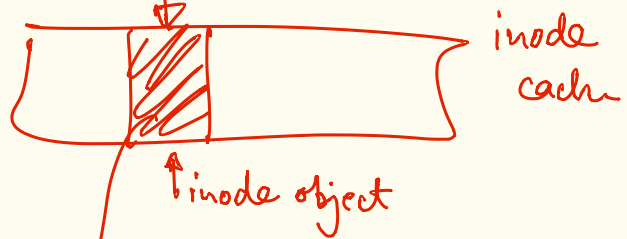
P1



Lecture 31 (?)



P2



* pointer for page cache is stored in the in-memory inode object.

* - in-memory inode object is not the same / but is built using disk inode object.

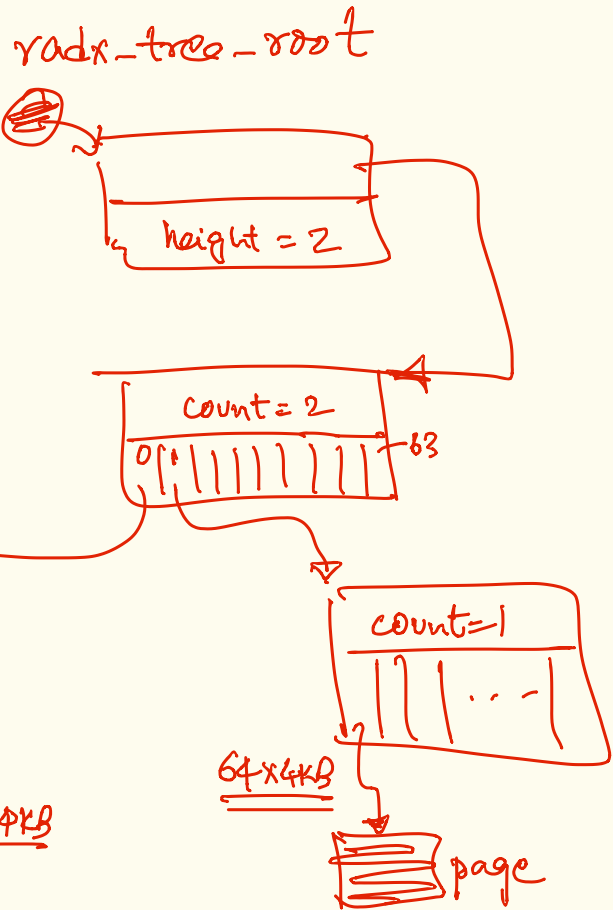
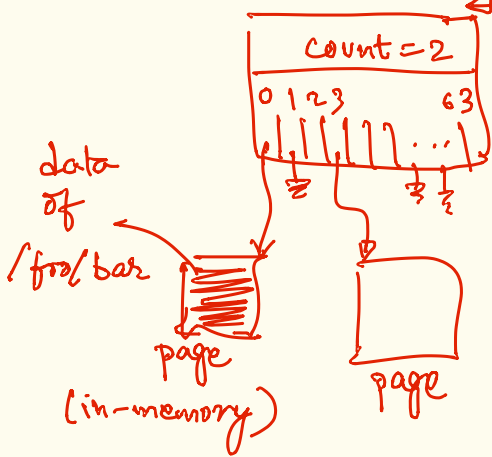
page cache

$\langle \text{fd, offset} \rangle \Rightarrow \boxed{\text{FS}} \Rightarrow \langle \text{disk, blocknum} \rangle$

user / applⁿ
process
view

page cache || — in-memory store
w/ page granularity
objects.

⊛ radix-tree page cache org. in Linux.



(*) cache : great for reads.

↳ not so much for writes.

↳ free (!) access
↳ temporal & spatial locality
↳ can be exploited

issue

⇒ multiple copies of ~~g~~ data being updated.

(i) cache invalidation

- invalidate / removes object in cache
- writes to disk

(ii) write through

- writes to cache & to disk synchronously
- then returns from write.

(iii) write back.

- lazy flush to disk.
- write returns on cache update.

`fsync()`

(*) fsck

Virtualization

Lecture 32(?)

CS 333/CS 347

- ① file consistency invariant \Rightarrow always keep the file system "consistent"
 - ~ all metadata & data are in-sync
 - ~ all or nothing updates
- ② persistent storage \Rightarrow what you write is what you get, across restarts (if write succeeds).

~ crash-consistency problem } ~ many disk writes for
 consistent update problem } a single FS action.
 ~ crashes are non-deterministic.

① pro-active approach
 e.g: journaling

② lazy approach
 e.g: fsck file system checker
 consistency check
 ~ let the world go by
 ~ if crash/restart
 - check FS consistency
 - fix as best as possible.

fsck.

write to a file.	inode	data bitmap	data block	
	C: ✓	✓	✓	
	C: ✗	✗	✗	
<u>cannot identify</u> , no action.	IC: ✗	✗	✓	↔
inconsistency detected: data bitmap: data is garbage	IC: ✓	✗	✗	(*)
	IC: ✗	✓	✗	~
detectable: data bitmap: reset	C: ✓	✓	✗	
appl ⁿ , brief inconsistency: FS not detectable (no action)	C: ✓	✗	✓	
detectable: set data bitmap	IC: ✗	✓	✓	(*) as data bitmap reset.

⑧ fsck. = needs to know file system specifications

+ superblock ~ verify integrity of super block.

+ free data blocks ~ are all ^{data} blocks pointed to by some inode.
that are set

+ inodes ~ verify type of inode
- checksum

+ inode links ~ how many files linked to this inode?

+ duplicates ~ data blocks pointed by more than one inode.

+ bad blocks - r/w checks
- block offsets within valid range.

+ directory checks - loops
└ " ." & " .. " are first two entries of each directory file.

⑨ VM files / images

- vmdk
img
vbx