

CS 537 - Lecture 13

1. Final exam - Dec 15th (Fri) - 5:30 pm - 7:30 pm

Only concurrency & Persistence.

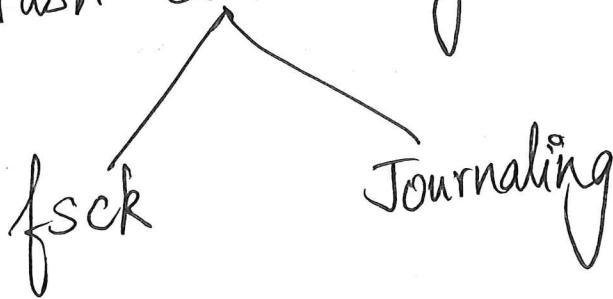
2. P5 - 2 & late days (max).

3. 2 worksheets - dropped.

10 " - graded.

1. Review of FS $\left\{ \begin{array}{l} \text{data structures} \\ \text{access methods} \end{array} \right.$

2. Crash consistency.

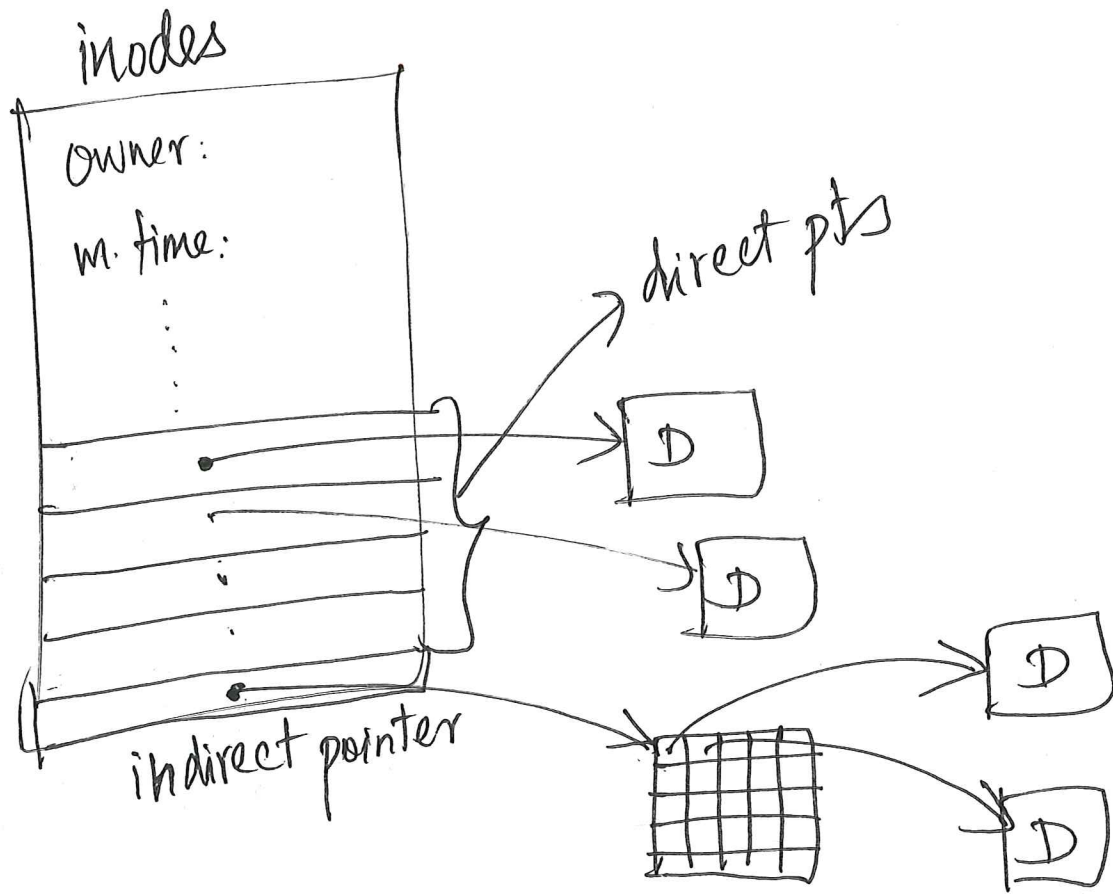
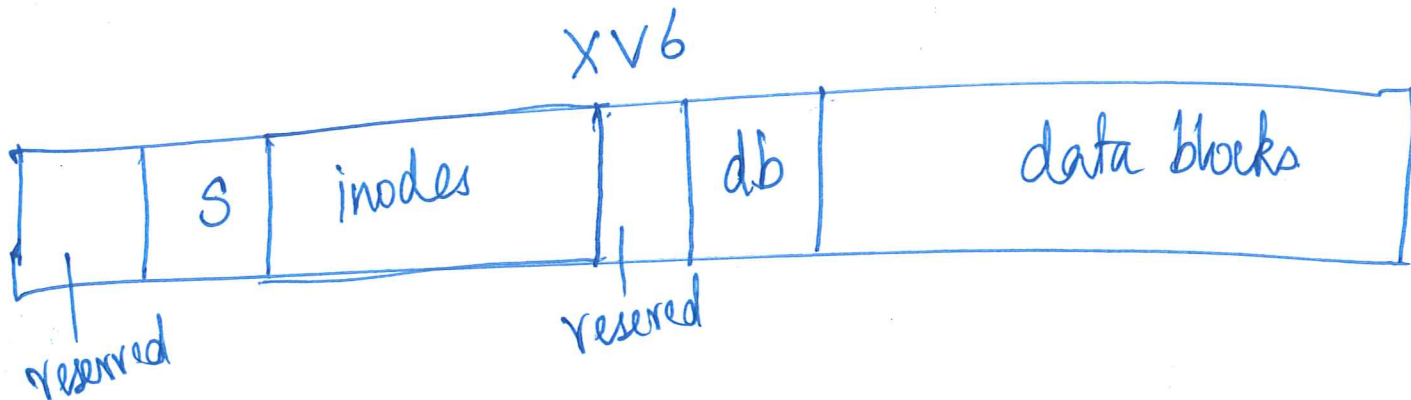


3. RAID.



File System

VFS ✓



contents of root directory

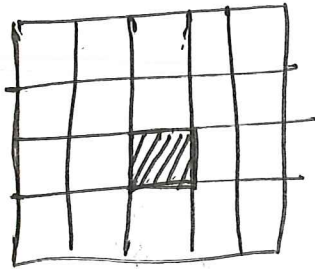
name	inode
.	2
..	2
foo	3

inode # of root = 2.

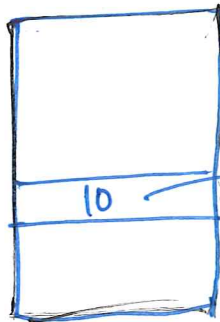
foo's dir contents

.	3
..	2
bar.txt	4

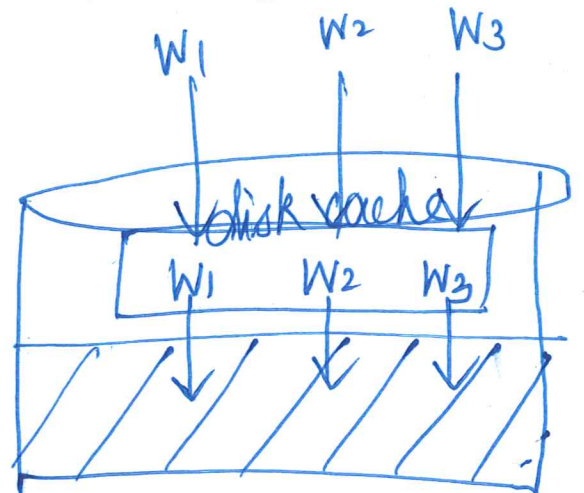
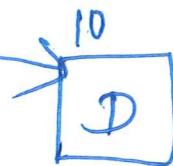
1 block



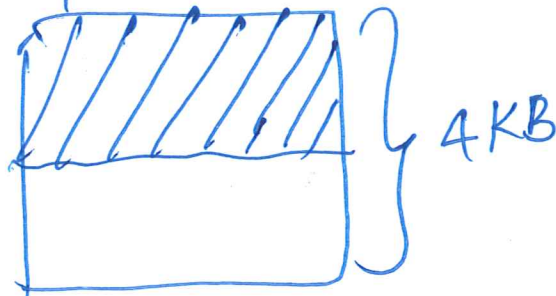
bar's inode



direct ptr



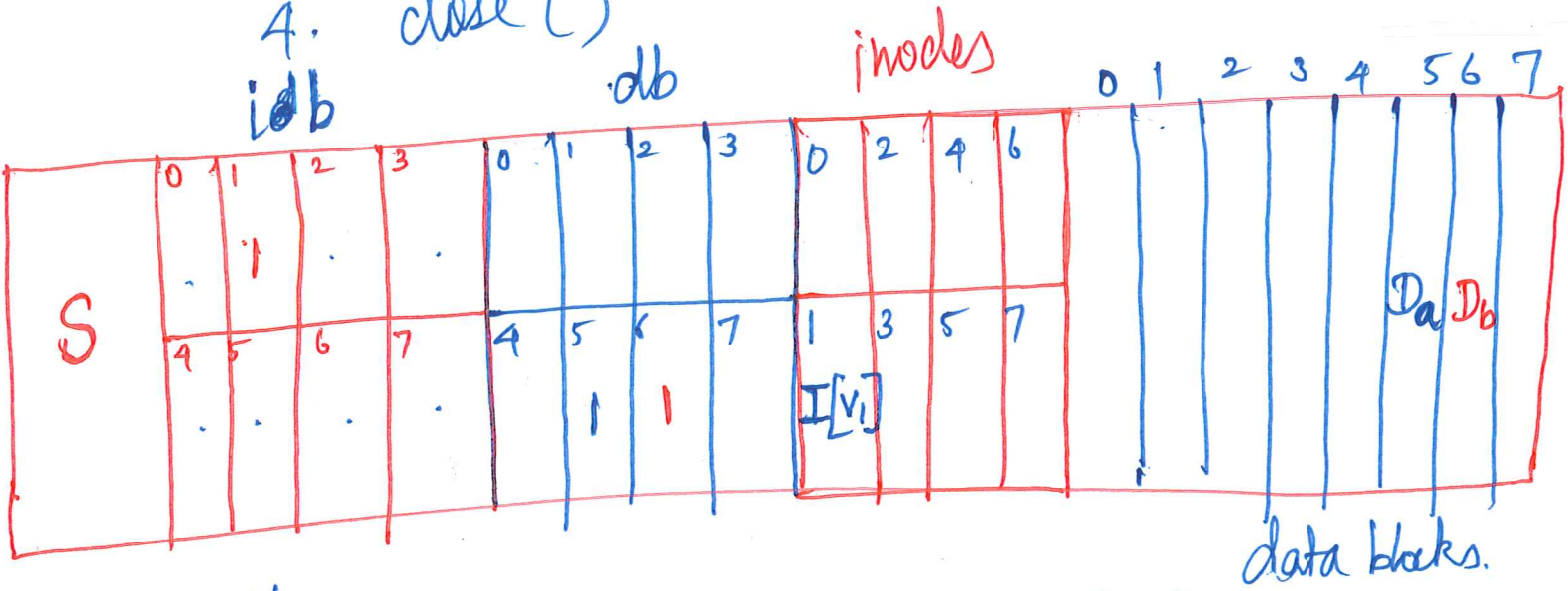
data block



Crash Consistency

Append

1. open ()
2. lseek () to end of file
3. issue a 4KB write
4. close ()



$$ib = \frac{v_1}{01000000}$$

$$db = 00000100$$

Blocks to update?

1. New data block (Db)
2. Data bitmap.

$$00000100 \rightarrow 00000110$$

$$db[v_1] \quad db[v_2]$$

3. Inode of the file.

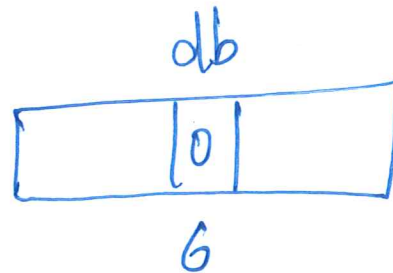
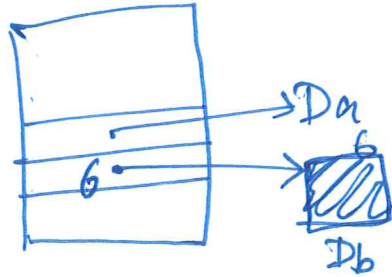
$$I[v_1] \rightarrow I[v_2]$$

→ 2nd direct ptr
 → # blocks: 1 → 2
 → size: 4KB → 8KB
 → modification time

Crash Scenarios

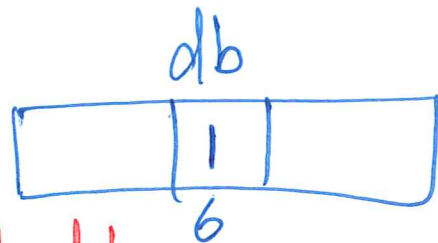
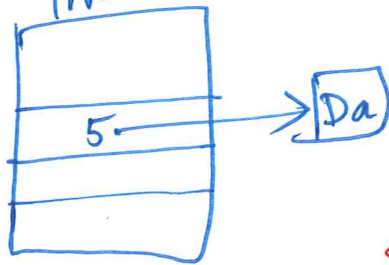
1. only Db is written.
— consistent w.r.t. FS.

2. only Inode is written.



— inconsistent.

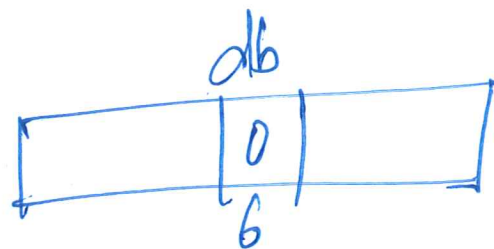
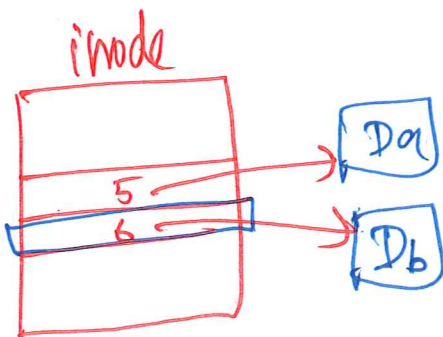
3. only the db[v₂] is written.



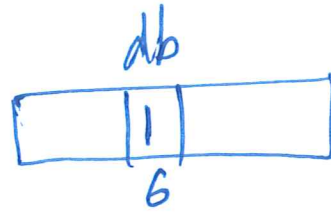
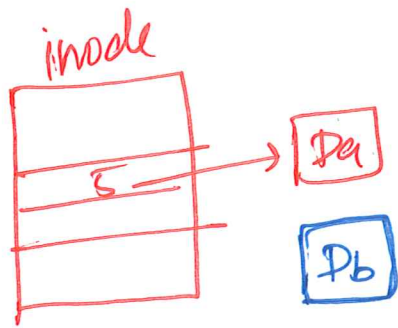
space leak!

— inconsistent state

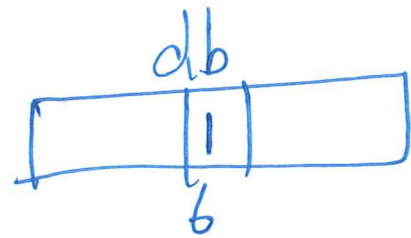
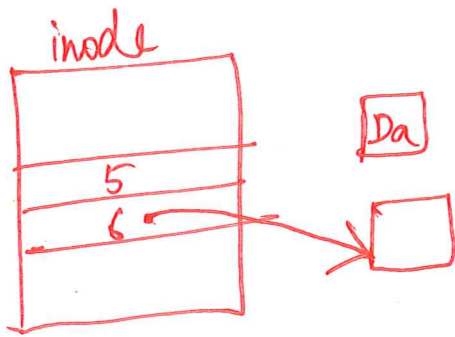
4. Db and I[v₂] are written.



5. D_b and $db[v_2]$.



6. $I[v_2]$ and $db[v_2]$.



Consistent-update problem.

FS $\xrightarrow{\text{atomically}}$ FS
consistent state 1

FS
consistent state 2

FSCK

- detect & recover.
- no other FS activity allowed.
- fsck - run before mounting the FS.

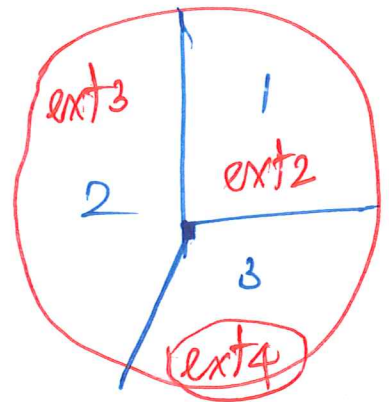
ASIDE:

~~mkfs~~ - make FS.

mkfs -t ext2 /dev/sda1

↓
device

sd → SCSI Device type.

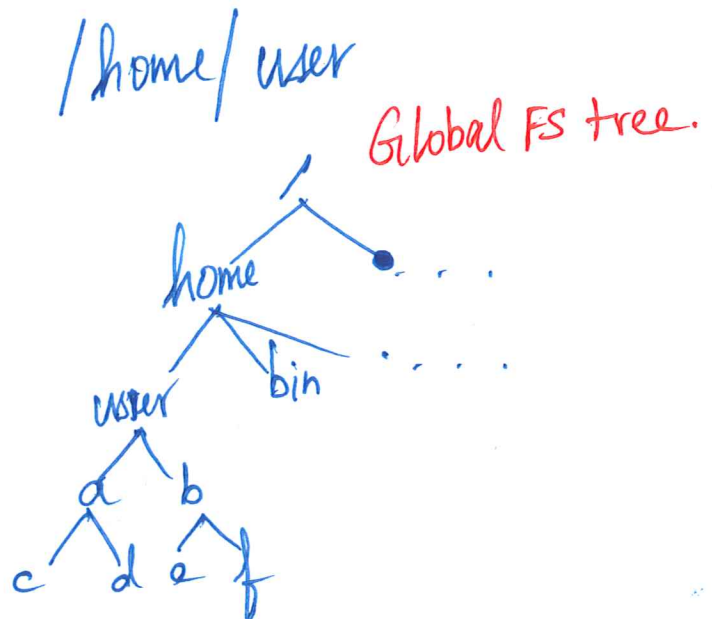
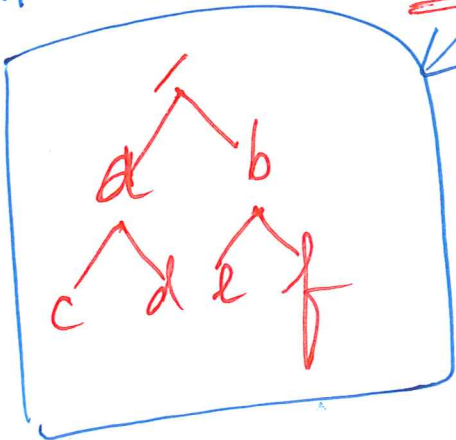


a, b, c, ..., z, Aa, Ab, ..., Az → disk # (order).

1, 2, 3, ..., n → partition # with the disk.

mount()

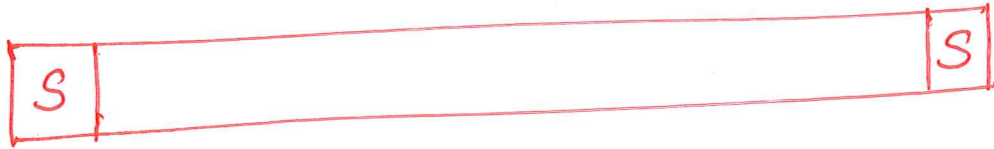
mount -t ext2 /dev/sda1 /home/user



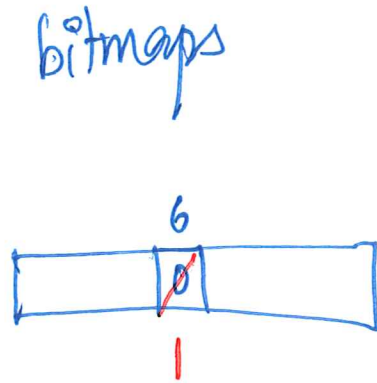
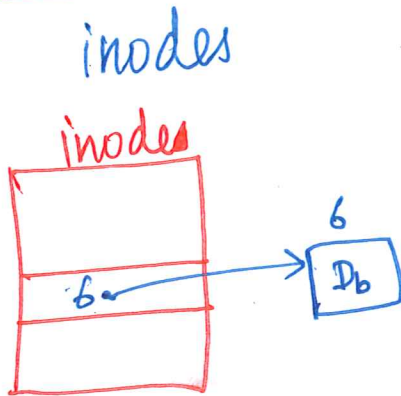
FSCK

1. Superblock.

sanity checks: FS size > # allocated blocks.
- alternative copy of Superblock.



2. Freeblocks



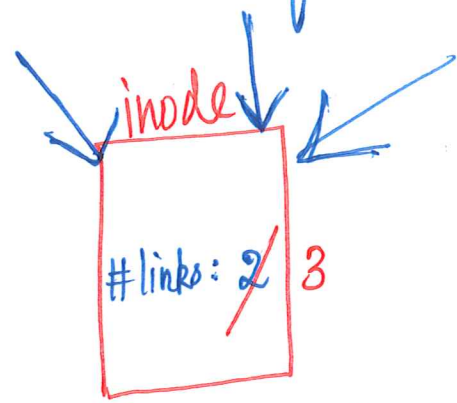
→ trust the inode.

3. Inode state

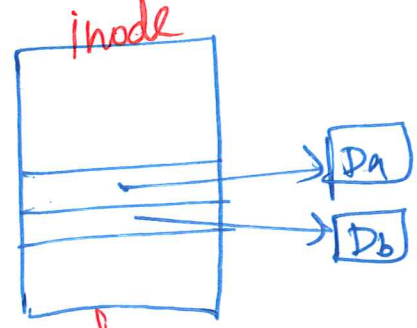
type: dir, reg file, sym. link.
⇒ free the inode
update the inode bitmap.

4. Inode links

— scan the entire FS and update the link counts for each inode.

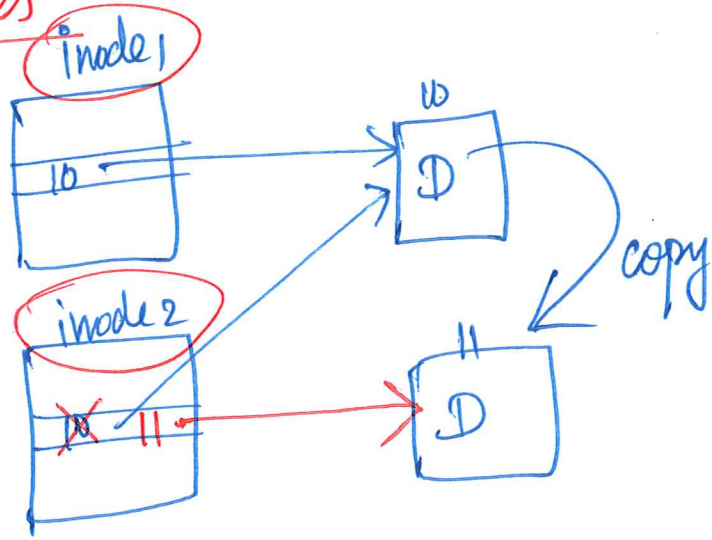


— inode is allocated but no dir is referring to this inode.

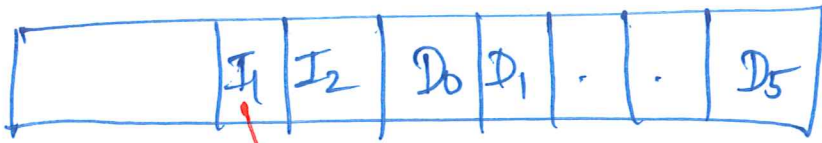


lost + found

5. Duplicates



6. Bad blocks



non-existent data block.

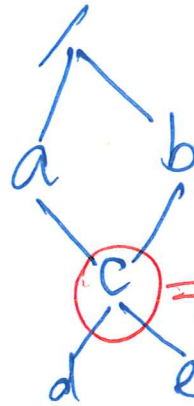
- clear the pointer to the data block.

7. Directory checks

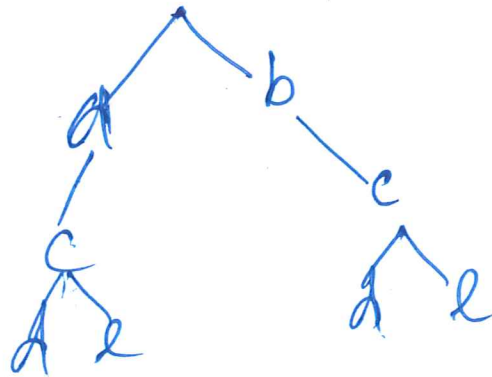
.	i_1
..	i_2
f_1	i_3
f_2	i_4

- first two entries are . and ..

EPIC



ref. count = 2

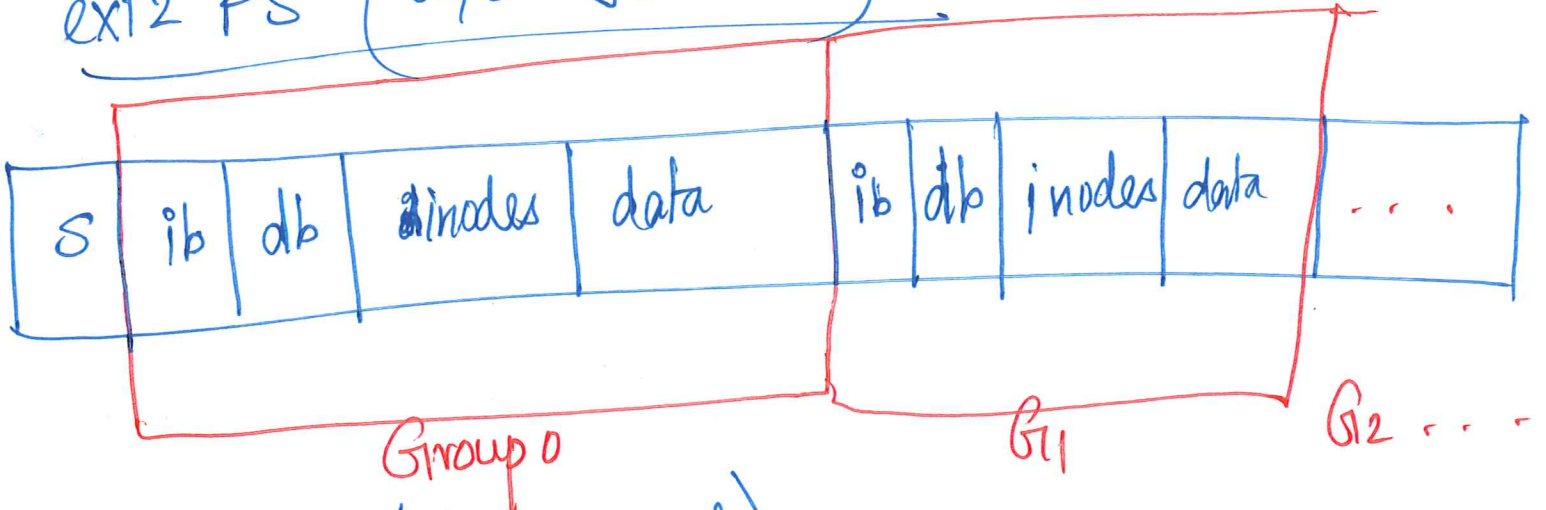


Data Journaling

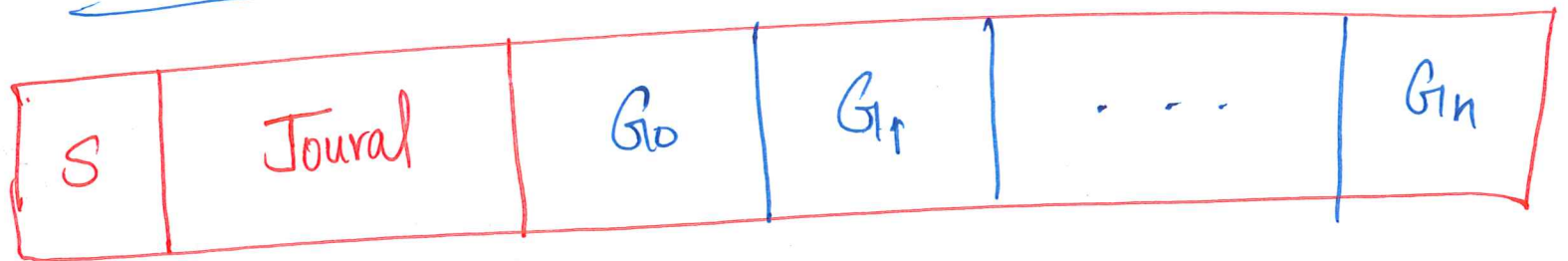
note: write to the journal → also on disk.

Write to the FS

ext2 FS (w/o Journal)



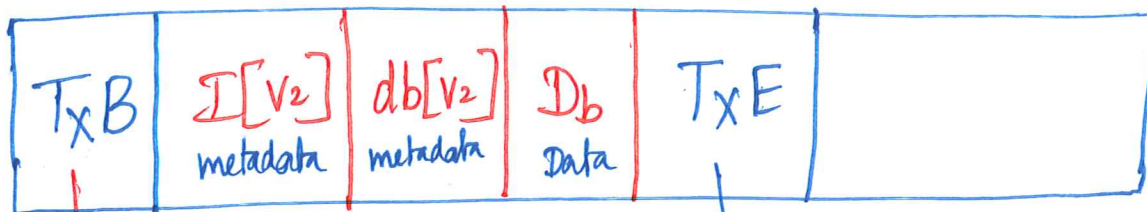
ext3 FS (w/ Journal)



Workload

1. Write $I[v_2]$
2. " $db[v_2]$
3. " Data Db.

Journal



final addr.
of the 3 write
+
TID

Journal

TID

1. Journal Write ✓
2. Checkpoint

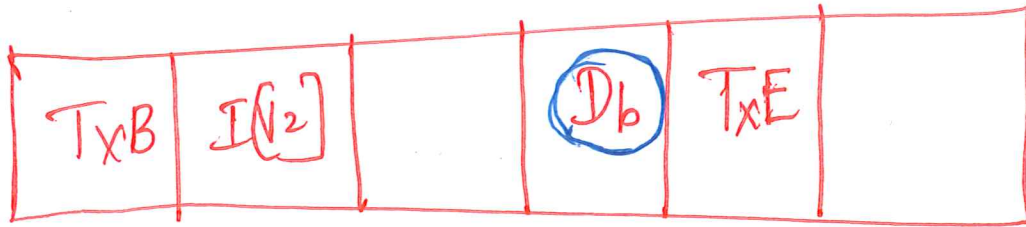
write the metadata & data
to the FS proper.

Crash

- 1) Write I[v₂]
- 2) Write db[v₂]
- 3) " Db

Crash recovery
replaying the tx.

Crash during Journal Write.



- db[V2] is NOT written.
- crash happen.

Protocol

1. Journal Write (TxB, metadata, data).
I[V2], db[V2] Db. WAIT.

WRITE BARRIER.

2. Journal Commit (Write TxE). WAIT.

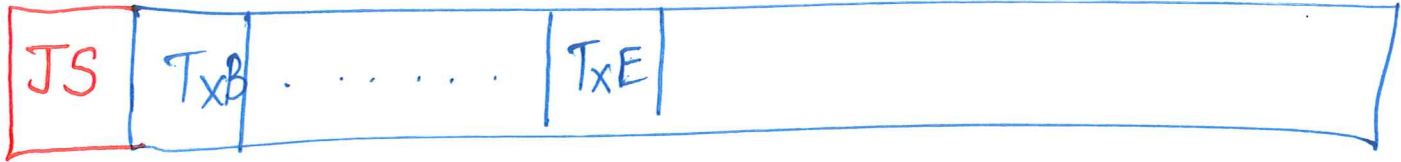
WRITE BARRIER.

3. Checkpoint

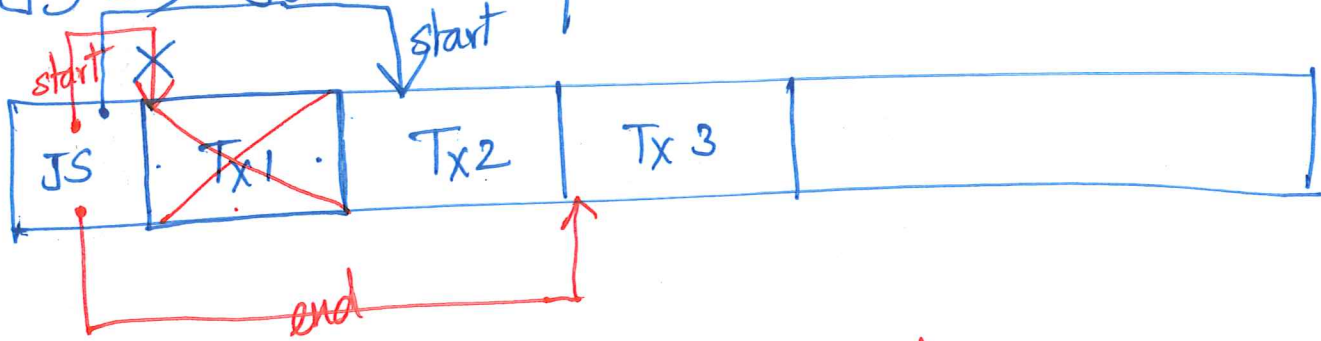
NO WRITE BARRIER

4. Free the transaction on the journal.
- updating the journal super block.

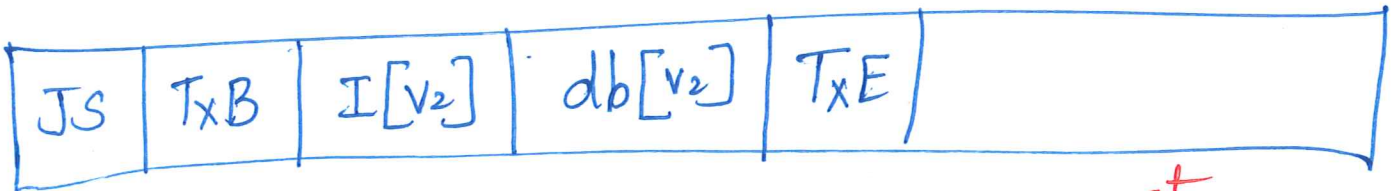
How to make the journal finite?



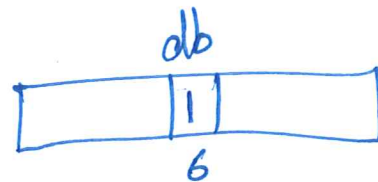
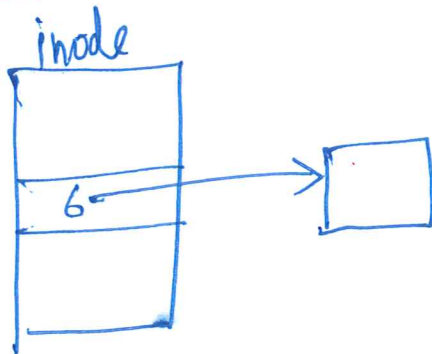
JS → Journal Super block



Metadata Journaling



→ Data block (Db) is written only ^{at} its actual location.



Metadata Journaling

1. Journal Write (TxB, metadata)

2. Journal Commit (Tx E)

3. Checkpoint

4. Free

When to write the data block?

1. Write data block. (Db). WAIT (optional)

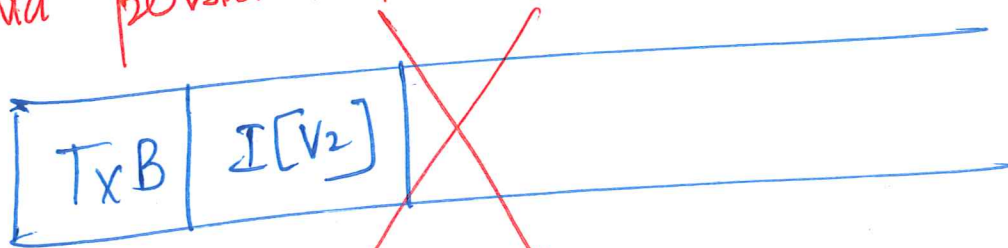
2. Journal metadata write (Tx B, metadata). WAIT

3. Journal commit. (Tx E). WAIT

4. Checkpoint the metadata.

5. Free the Tx from the journal.

→ Data persisted to disk.



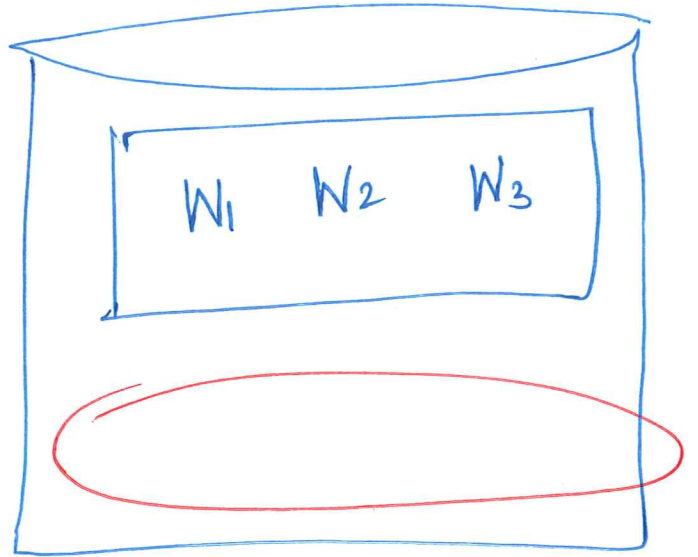
fsync(fd);

fd = open();

→ write(); W1

→ write(); W2

write(); W3



write(fd); W1

fsync(fd);

→ write(fd); W2

fsync(fd);

1. create a new file
2. write to the file

3. fsync.

→ CRASH

DIRECTORY.

foo

..	3
..	2
bar.txt	4

/foo/bar.txt

