

# CS202 ° FILE SYSTEMS III ° CRASH RECOVERY

## ° Logistics Reminder:

- No class on THU (11/23) OR NEXT TUESDAY (11/28)
- LAB 5 HAS NEW DUE DATE: 12/8
- REVIEW SESSION: 11/27 (MONDAY)  
7-8 PM

## ° LAST CLASS

- DIRECTORIES
- FFS

## ° THIS CLASS

- DEALING WITH CRASHES

## CRASH RECOVERY

### CORE PROBLEM

- DISK WRITES PERSIST  
↳ By design
- DISK OPERATION ONLY ATOMIC AT  
SECTOR OR BLOCK GRANULARITY

• WHY BAD?

> ls > /tmp/out

① Allocate inode for out  
↳ Read inode bitmap

→ Update bitmap

② Modify /tmp to point to allocated inode

- Read /tmp inode to find datablock

- Read datablock to find empty index

- Add link to out

③ Allocate datablock for /tmp/out

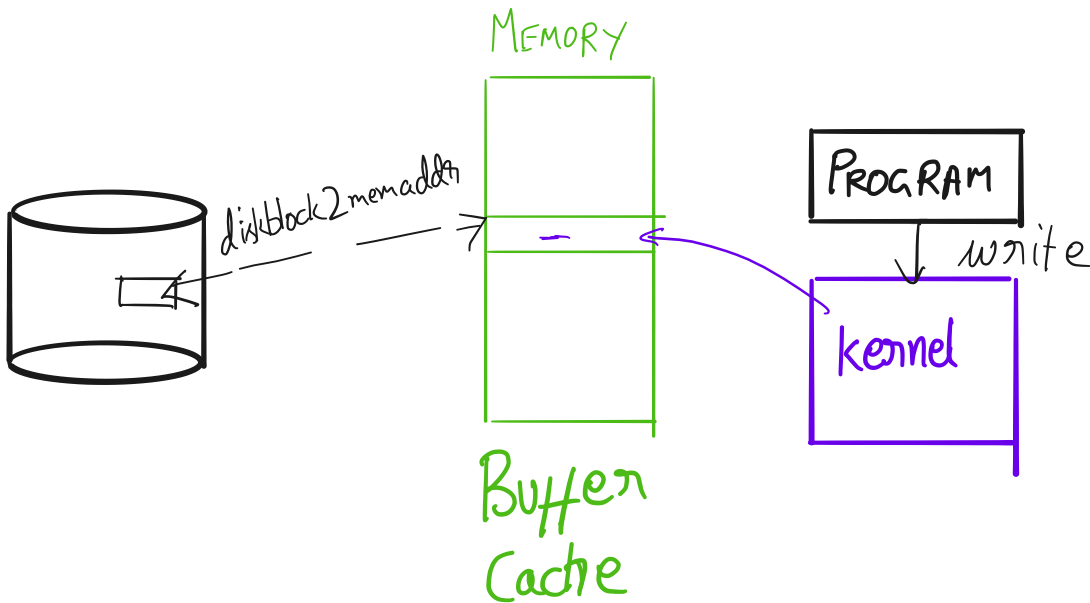
④ Write to datablock

OBSERVATION: ILL-TIMED CRASHES CAN LEAD TO

- RESOURCE LEAKS
- INCONSISTENT FILE METADATA
- DATA LOSS
- ...

GOAL: COME UP WITH MECHANISMS TO ENSURE  
DISK CONSISTENCY AFTER CRASH

BUT FIRST, THINGS GET A BIT WORSE



MODIFICATIONS ARE NOT IMMEDIATELY WRITTEN  
TO DISK.

CAN FORCE THEM (fsync (2)) BUT EXPENSIVE.

WHY IS THIS A PROBLEM FOR CRASH RECOVERY?

RECOVERY MECHANISMS

- ① AD HOC / FSCK
- ② COPY ON WRITE
- ③ JOURNALING

AD HOC / FSCK

ASSUMPTION: RECOVERY LOGIC WILL EXECUTE WHEN  
A COMPUTER RESTARTS AFTER A CRASH

USUALLY IMPLEMENTED IN A PROGRAM

CALLED FSCK

GOAL: ENSURE METADATA REMAINS CONSISTENT

→ NO LEAKED DATA BLOCKS OR

INODES

NON-GOAL: DATA CONSISTENCY

↳ A FILE'S CONTENTS MIGHT

NOT MAKE SENSE

APPROACH: DESIGN FS LOGIC SO THAT FSCK  
CAN REPAIR METADATA AFTER CRASH

> ls > /tmp/out

① WRITE DATA TO DATA BLOCK

---

② WRITE DATA BLOCK REF TO INODE  
FOR /tmp/out

---

③ UPDATE INODE BITMAP

---

④ UPDATE DATA BLOCK BITMAP

---

⑤ UPDATE /tmp TO LINK TO  
/tmp/out

---

D  
U. Juc. Anonymu?

# PROBLEMS WITH THIS APPROACH:

## COPY-ON-WRITE

- CORE IDEA: DO NOT MODIFY\* BLOCKS

COPY INSTEAD

(ONE EXCEPTION: UBERBLOCK/SUPERBLOCK)

- CAN ACHIEVE BOTH METADATA &

DATA CONSISTENCY

- How? [SWITCH TO HANDOUT]

- WHY THIS WORKS

- UBERBLOCK UPDATED ATOMICALLY  
(ONE BLOCK)

- OPERATION  
- COPIES

---

→ UBER BLOCK UPDATE

---

COMMIT  
POINT

OTHER QUESTIONS

(a) PROBLEMS WITH THIS APPROACH?

(b) BENEFITS BEYOND CRASH RECOVERY

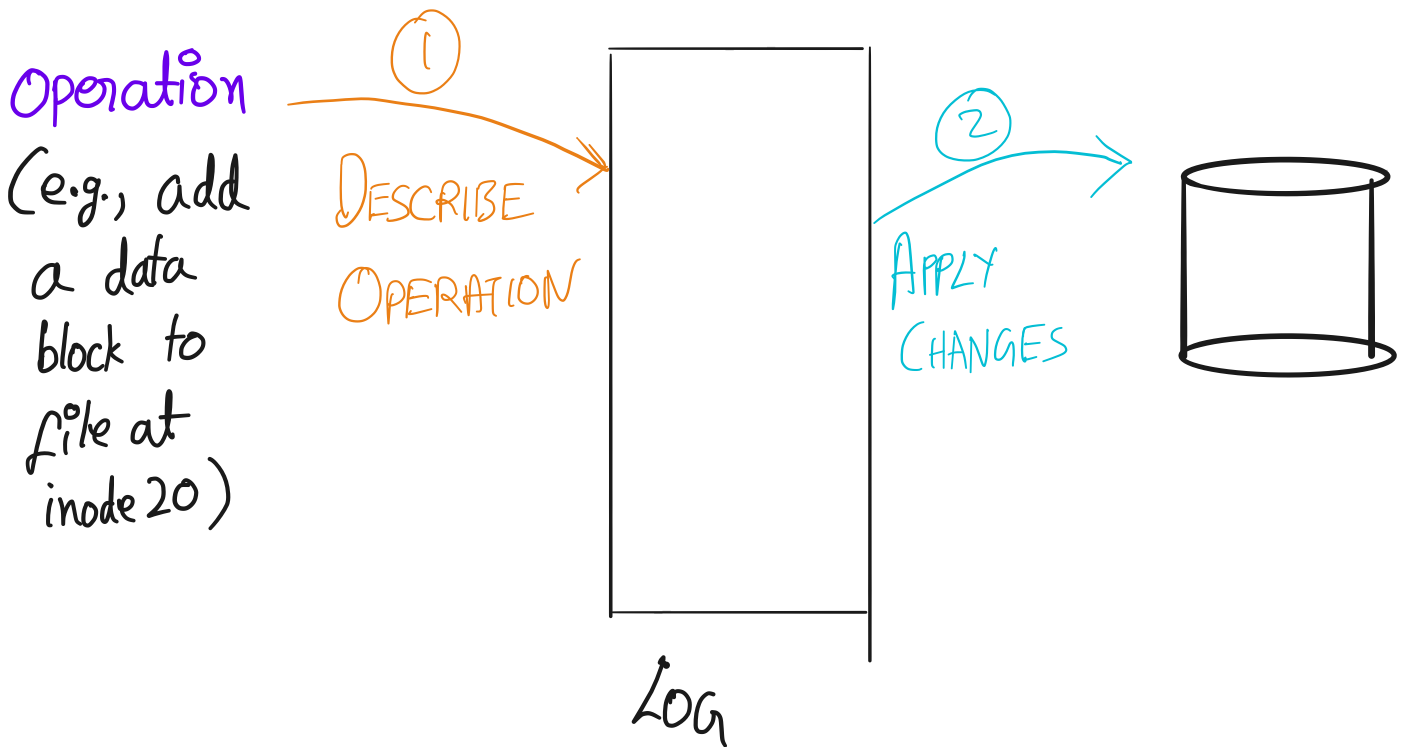
# JOURNALING

- BUILDS ON IDEAS THAT ARE SIMILAR TO COPY ON WRITE

DO NOT MODIFY THE ONLY COPY OF THE  
FILESYSTEM

⇒ NEED A COMMIT POINT AFTER WHICH OPERATION  
IS GUARANTEED TO APPEAR

- DIFFERENCE ◦ REDUCE OVERHEADS



## LOG REQUIREMENTS



NEED TO BE ABLE TO DISTINGUISH  
B/W OPS THAT HAVE BEEN  
COMPLETELY LOGGED & THOSE WHO  
HAVE NOT.

WHY MIGHT LOG FOR AN OP BE  
INCOMPLETE?

How? [SWITCH TO]

GENERALLY, A COMPLETELY LOGGED OP IS  
COMMITTED

## LOG CONTENTS

DEPENDS ON TYPE OF LOG

(a) REDO LOGGING

LOG RECORDS HOW OP CHANGES

DISK

↳ > &lt;img alt="arrow pointing right" data-bbox="650 910 680 940"/> &lt;img alt="arrow pointing left" data-bbox="700 910 730 940"/> &lt;img alt="arrow pointing right" data-bbox="750 910 780 940"/> &lt;img alt="arrow pointing left" data-bbox="800 910 830 940"/> &lt;img alt="arrow pointing right" data-bbox="850 910 880 940"/> &lt;img alt="arrow pointing left" data-bbox="900 910 930 940"/>

(i) Allocate inode for out

- (i) Allocate inode
    - ↳ Read inode bitmap
    - Update bitmap
  - (ii) Modify fmp to point to allocated inode
  - (iii) Allocate data block
- o o o

RECOVERY LOGIC APPLIES CHANGES FROM ANY COMPLETELY LOGGED OP

Q: WHEN IS IT SAFE TO APPLY CHANGES TO DISK (OTHER THAN DURING RECOVERY)?

(b) UNDO LOGGING

LOG RECORDS HOW TO UNDO EACH STEP IN AN OP

- (i) Allocate inode for out
  - ↳ VALUE OF BITMAP BEFORE

ALLOCATION  
ii) Modify /tmp to point to  
allocated inode  
↳ OLD VALUE OF /tmp

o o o

RECOVERY LOGIC: UNDO STEPS FROM ANY  
Op THAT IS NOT COMPLETELY LOGGED

Q: WHEN IS IT SAFE TO MARK THAT AN  
Op IS COMPLETELY LOGGED?

Q: WHEN IS IT SAFE TO APPLY CHANGES  
TO THE DISK?

REDO VS UNDO LOGGING

REDO

UNDO

PRO

CON

REDO + UNDO LOGGING

- SEE NOTES

WHAT TO LOG?

- METADATA ◦ INODES, BITMAPS, ETC.

- DATA ◦

- ...

# CS202 Handout 11

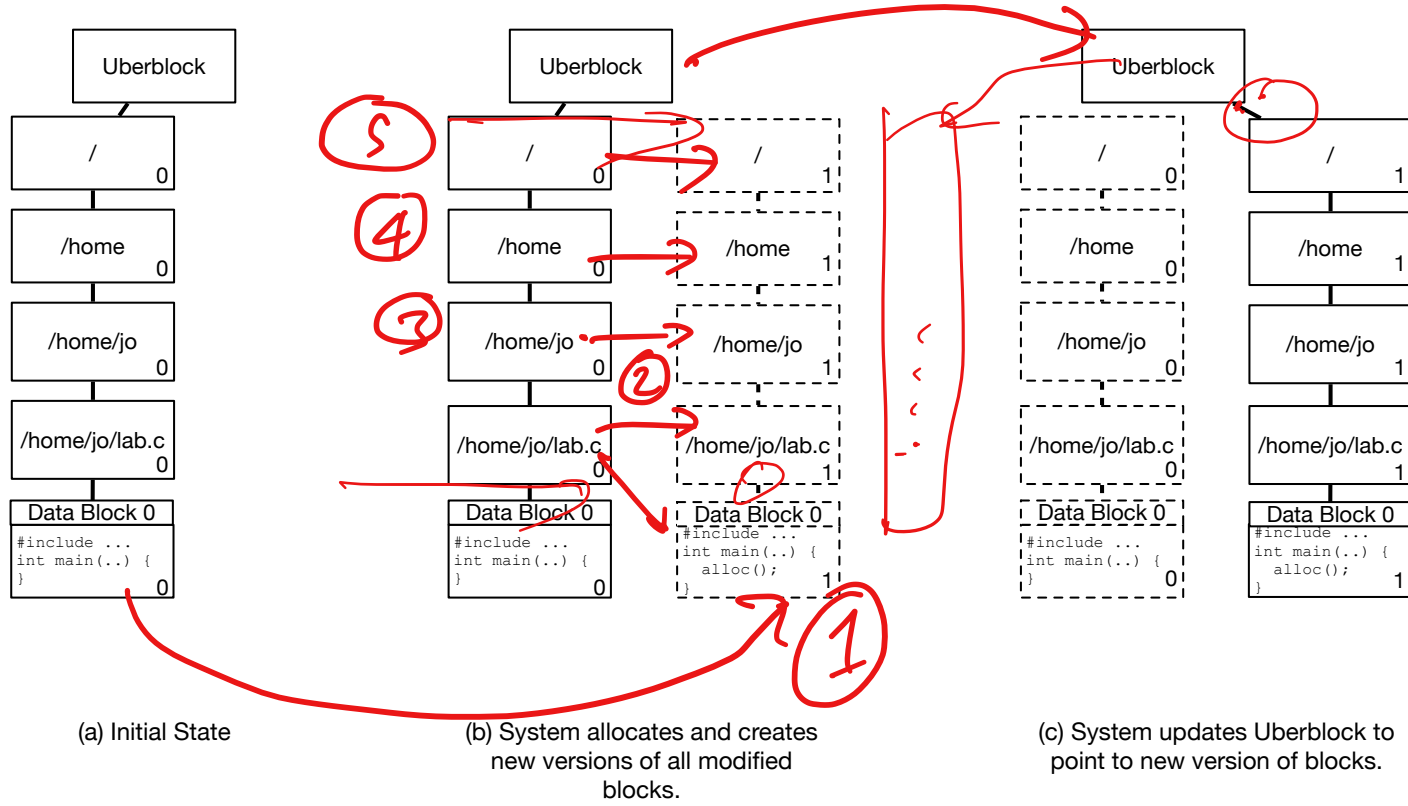
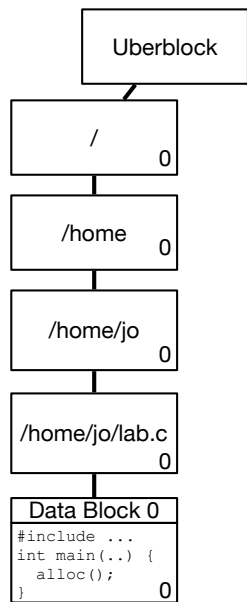
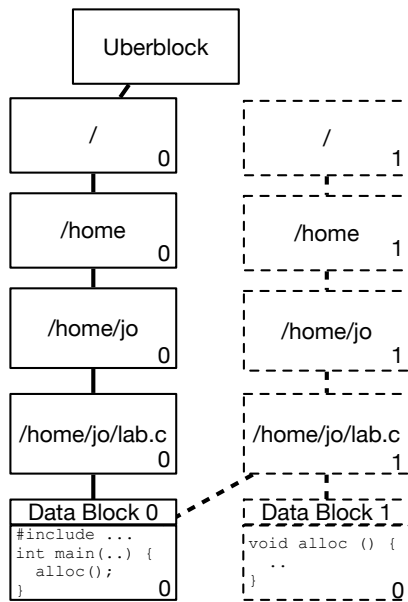


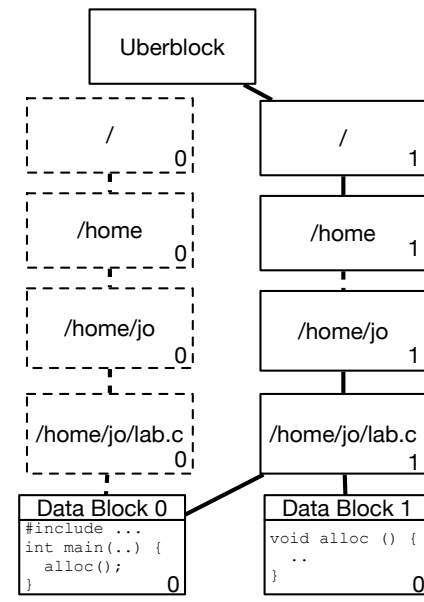
Figure 1: Copy-on-write filesystem: modifying a data block



(a) Initial State

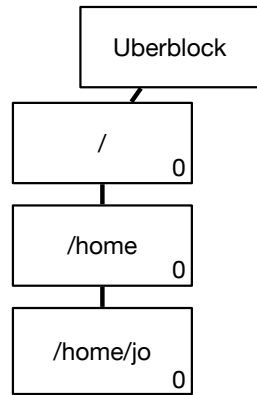


(b) System allocates and creates new versions of all modified blocks.

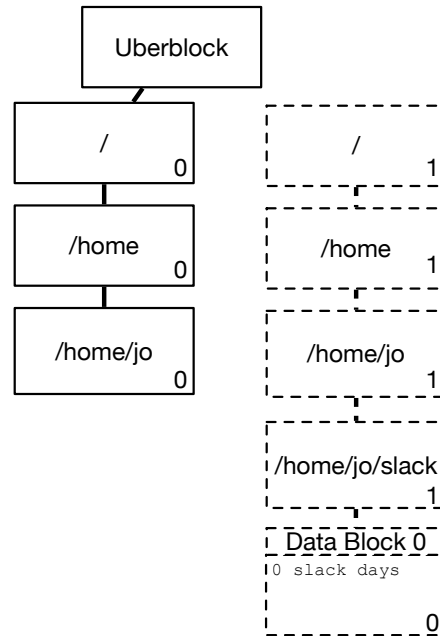


(c) System updates Uberblock to point to new version of blocks.

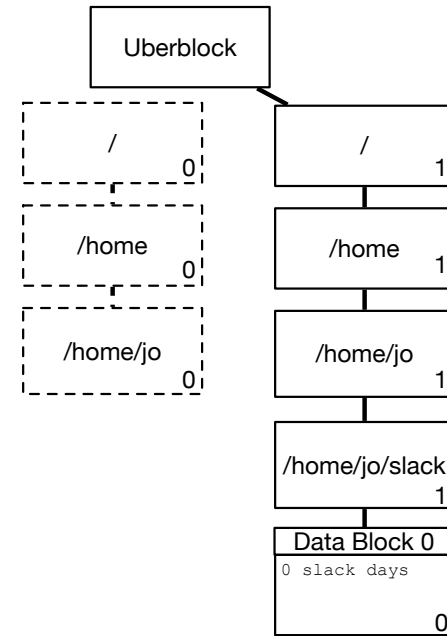
Figure 2: Copy-on-write filesystem: adding a data block



(a) Initial State



(b) System allocates and creates new versions of all modified blocks.

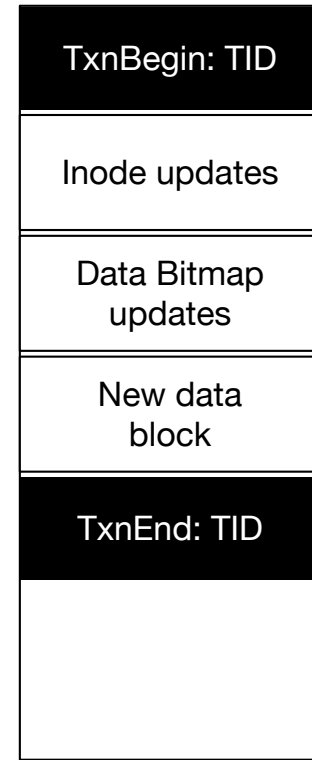


(c) System updates Uberblock to point to new version of blocks.

Figure 3: Copy-on-write filesystem: creating a file



ext3 disk layout



ext3 journal layout

Figure 4: Redo logging in a filesystem