# Journaling
## and
## Log-structured file systems

Johan Montelius

KTH

2019

A file system is the user space implementation of *persistent storage*.

# The file system

A file system is the user space implementation of *persistent storage*.

- a *file* is persistent i.e. it survives the termination of a process
- a *file* can be access by several processes i.e. a shared resource
- a *file* can be located given a *path* name

## let's write to a file

Assume we want to write to a file `bar.txt`, that requires a new block to be allocated.

We need to:

- update the block bitmap - we have allocated one more data block

Assume we want to write to a file `bar.txt`, that requires a new block to be allocated.

We need to:

- update the block bitmap - we have allocated one more data block
- update the inode of `bar.txt` - a new data block, size and access time

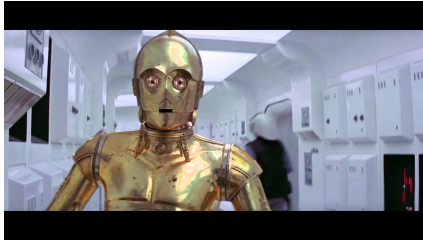Assume we want to write to a file bar.txt, that requires a new block to be allocated.

We need to:

- update the block bitmap - we have allocated one more data block
- update the inode of bar.txt - a new data block, size and access time
- update the block - the new data (it might contain old data).

Assume we want to write to a file bar.txt, that requires a new block to be allocated.

We need to:

- update the block bitmap - we have allocated one more data block
- update the inode of bar.txt - a new data block, size and access time
- update the block - the new data (it might contain old data).

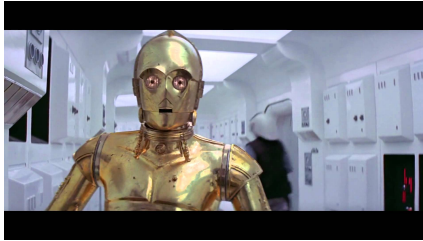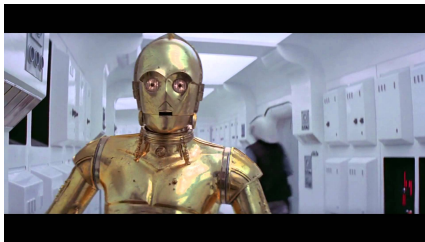*In what order should we perform these operations?*

We're doomed!

We're doomed!
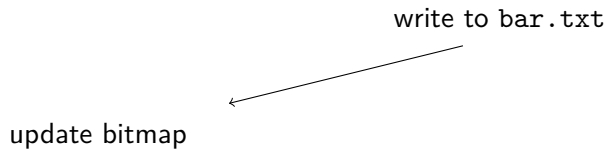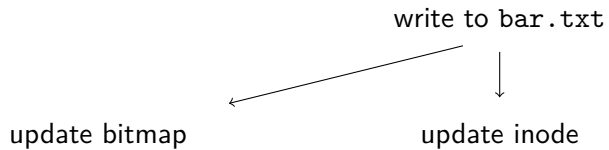
How do we cope with crashing drives?

We're doomed!

How do we cope with crashing drives?

How do we cope with the operating system crashing?

write to `bar.txt`

write to `bar.txt`

update bitmap

write to `bar.txt`

update bitmap          update inode

write to `bar.txt`

update bitmap          update inode          update data block

write to bar.txt

update bitmap        update inode        update data block

bitmap and inode                    inode and data
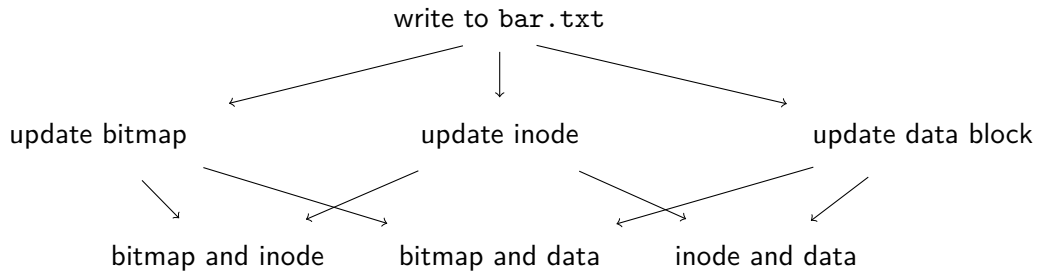
Two out of three <u>is</u> - when it comes to file systems - bad.

Approaches:

Approaches:

- file system check - recover as much as possible

Approaches:

- file system check - recover as much as possible
- journal - write down what you want to do, before you do it

Approaches:

- file system check - recover as much as possible
- journal - write down what you want to do, before you do it
- log - the file system is a log of changes

Approaches:

- file system check - recover as much as possible
- journal - write down what you want to do, before you do it
- log - the file system is a log of changes
- copy on write - create a perfect copy and flip a pointer

Approaches:

- file system check - recover as much as possible
- journal - write down what you want to do, before you do it
- log - the file system is a log of changes
- copy on write - create a perfect copy and flip a pointer

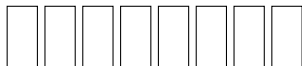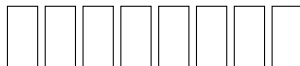Remember the Very Simple File System:
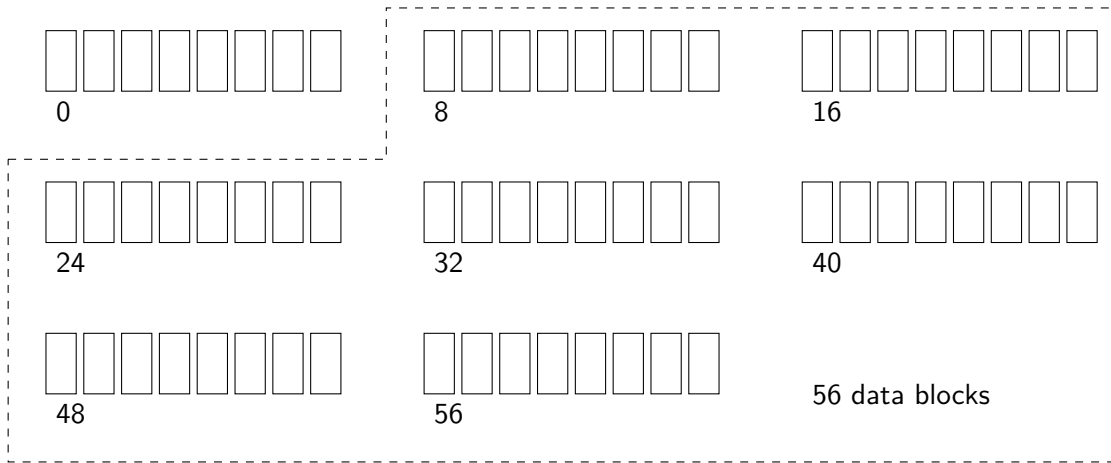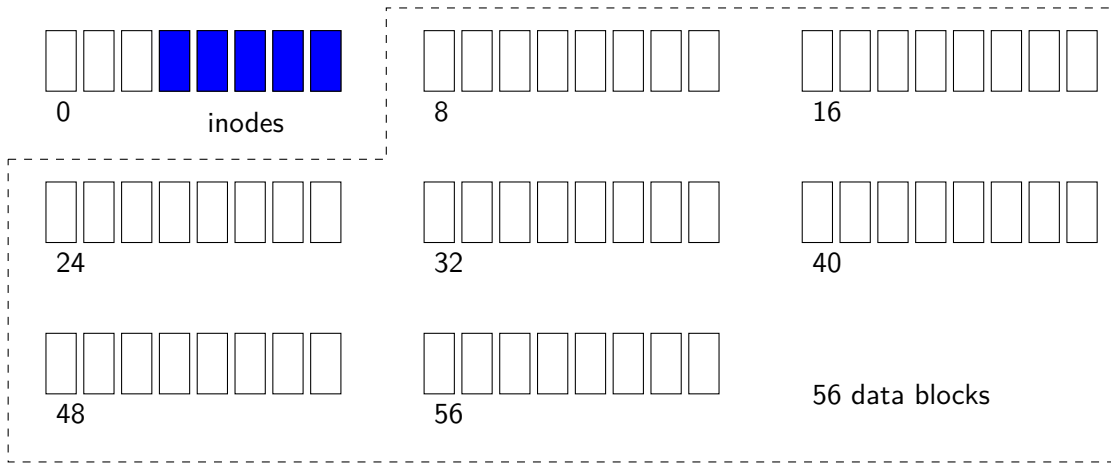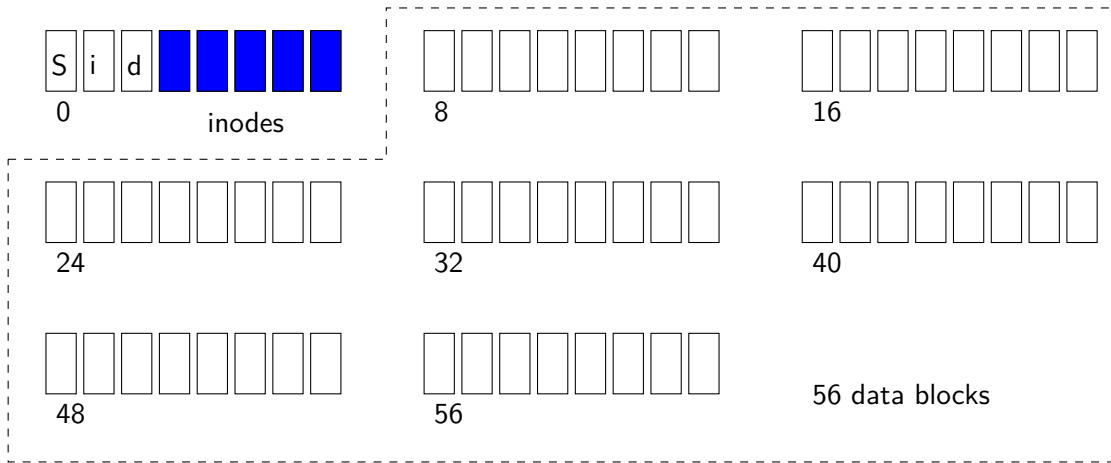

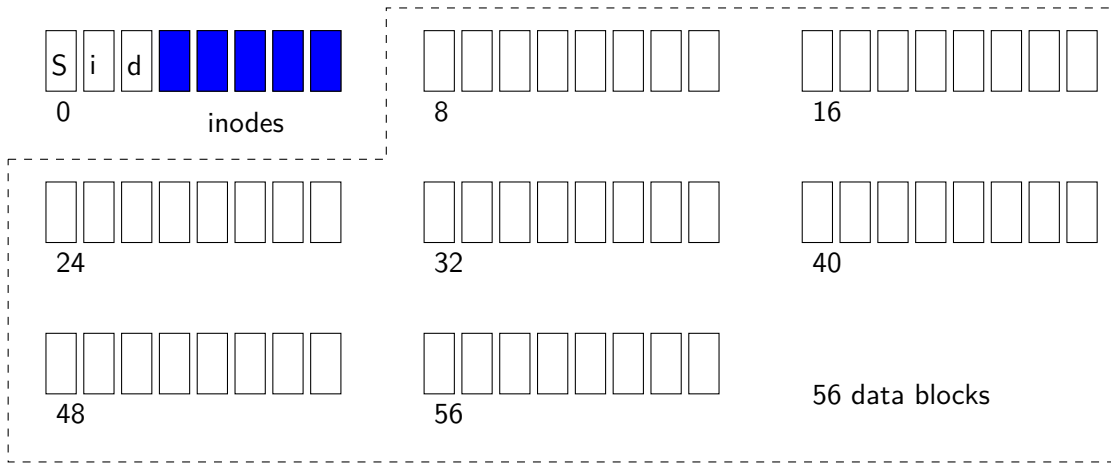
0

8

16

24

32

40

48

56

Remember the Very Simple File System:



56 data blocks

Remember the Very Simple File System:

Remember the Very Simple File System:

Remember the Very Simple File System:

How would we rebuild a file system?

## fsck /dev/sdb1
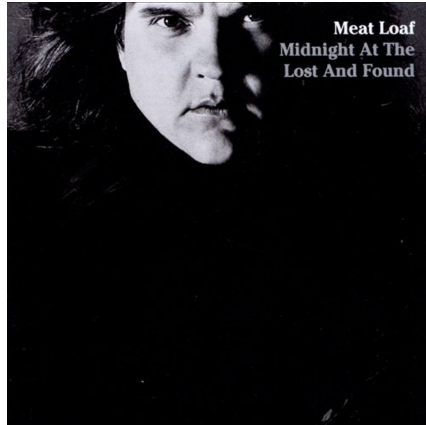
```
$ sudo fsck  -f /dev/sdb1
fsck from util-linux 2.27.1
e2fsck 1.42.13 (17-May-2015)
Pass 1: Checking inodes, blocks, and sizes
Pass 2: Checking directory structure
Pass 3: Checking directory connectivity
Pass 4: Checking reference counts
Pass 5: Checking group summary information
/dev/sdb1: 3339/125952 files (0.1% non-contiguous), 318256/503808 blocks
```

```
> ls -il /
    :
    :

11010049 drwxr-xr-x   2 root root 12288 nov 28 17:49 libx32

      11 drwx------   2 root root 16384 maj  8  2016 lost+found

14155777 drwxr-xr-x   3 root root  4096 jun 29 14:13 media

  262145 drwxr-xr-x   3 root root  4096 okt 22 10:17 mnt

    :
    :
```

We need to move from *a consistent state* to a *consistent state*.

We need to move from *a consistent state* to a *consistent state*.

Let's keep a *journal* of things we are about to do.

We need to move from *a consistent state* to a *consistent state*.

Let's keep a *journal* of things we are about to do.

# Journal or Write-Ahead Logging

We need to move from *a consistent state* to a *consistent state*.

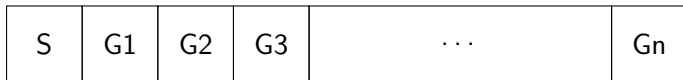Let's keep a *journal* of things we are about to do.

## Journal or Write-Ahead Logging

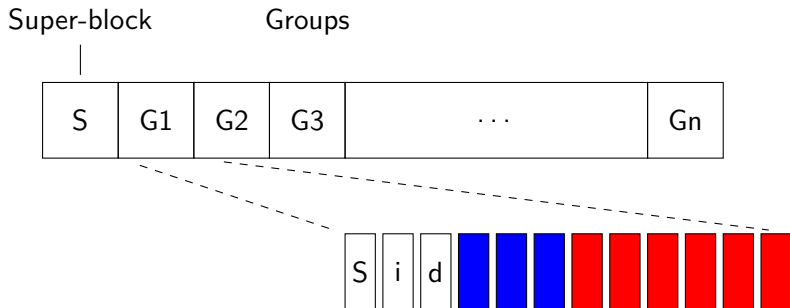If we crash we can look at the journal to repeat the last sequence of operations.

Super-block          Groups

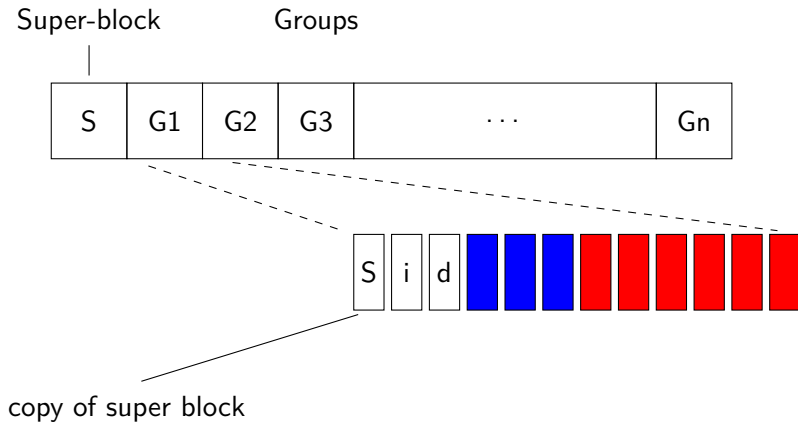| S | G1 | G2 | G3 | $\cdots$ | Gn |
|---|----|----|----|----------|----|

Super-block · · · Groups

| S | G1 | G2 | G3 | · · · | Gn |

| S | i | d | | | | | | | | | |

Super-block      Groups

| S | G1 | G2 | G3 | $\cdots$ | Gn |

S   i   d

copy of super block

Super-block        Groups

| S | G1 | G2 | G3 | $\cdots$ | Gn |

| S | i | d | | | | | | | | | |

copy of super block

inode bitmap

Super-block          Groups

| S | G1 | G2 | G3 | $\cdots$ | Gn |

S | i | d

copy of super block                block bitmap
                inode bitmap

Super-block · Groups

| S | G1 | G2 | G3 | · · · | Gn |

| S | i | d | inodes | data blocks |

copy of super block

inode bitmap

block bitmap

Journal

| S | J | G1 | G2 | $\cdots$ | Gn |

# Linux ext3 - journaling

Journal

| S | J | G1 | G2 | $\cdots$ | Gn |

Journal

| S | J | G1 | G2 | $\cdots$ | | Gn |

| $TxB$ | $I_{v2}$ | $B_{v2}$ | $D_{v2}$ | |

transaction begin          inode      bitmaps   data block

transaction begin inode bitmaps data block transaction end

- Commit: write the transaction

- Commit: write the transaction
  - *TxB* : transaction id, inode id, bit map id, data block id

- Commit: write the transaction
  - *TxB* : transaction id, inode id, bit map id, data block id
  - $I_{v2}$ : the updated inode

- Commit: write the transaction
    - $TxB$ : transaction id, inode id, bit map id, data block id
    - $I_{v2}$ : the updated inode
    - $B_{v2}$ : the updated bitmaps

- Commit: write the transaction
    - $TxB$ : transaction id, inode id, bit map id, data block id
    - $I_{v2}$ : the updated inode
    - $B_{v2}$ : the updated bitmaps
    - $D_{v2}$ : the updated data block

- Commit: write the transaction
    - $TxB$ : transaction id, inode id, bit map id, data block id
    - $I_{v2}$ : the updated inode
    - $B_{v2}$ : the updated bitmaps
    - $D_{v2}$ : the updated data block
    - $TxE$ : transaction id

- Commit: write the transaction
    - $TxB$ : transaction id, inode id, bit map id, data block id
    - $I_{v2}$ : the updated inode
    - $B_{v2}$ : the updated bitmaps
    - $D_{v2}$ : the updated data block
    - $TxE$ : transaction id
- Checkpoint: perform the changes

- Commit: write the transaction
    - $TxB$ : transaction id, inode id, bit map id, data block id
    - $I_{v2}$ : the updated inode
    - $B_{v2}$ : the updated bitmaps
    - $D_{v2}$ : the updated data block
    - $TxE$ : transaction id
- Checkpoint: perform the changes
    - update the blocks: inode, bit maps and data block

- Commit: write the transaction
    - $TxB$ : transaction id, inode id, bit map id, data block id
    - $I_{v2}$ : the updated inode
    - $B_{v2}$ : the updated bitmaps
    - $D_{v2}$ : the updated data block
    - $TxE$ : transaction id
- Checkpoint: perform the changes
    - update the blocks: inode, bit maps and data block
    - remove transaction

- Commit: write the transaction
  - $TxB$ : transaction id, inode id, bit map id, data block id
  - $I_{v2}$ : the updated inode
  - $B_{v2}$ : the updated bitmaps
  - $D_{v2}$ : the updated data block
  - $TxE$ : transaction id
- Checkpoint: perform the changes
  - update the blocks: inode, bit maps and data block
  - remove transaction

- Commit: write the transaction
  - $TxB$ : transaction id, inode id, bit map id, data block id
  - $I_{v2}$ : the updated inode
  - $B_{v2}$ : the updated bitmaps
  - $D_{v2}$ : the updated data block
  - $TxE$ : transaction id
- Checkpoint: perform the changes
  - update the blocks: inode, bit maps and data block
  - remove transaction

We manage to write half of the transaction.

We manage to write half of the transaction.

We manage to write the whole transaction but not updating the blocks.

We manage to write half of the transaction.

We manage to write the whole transaction but not updating the blocks.

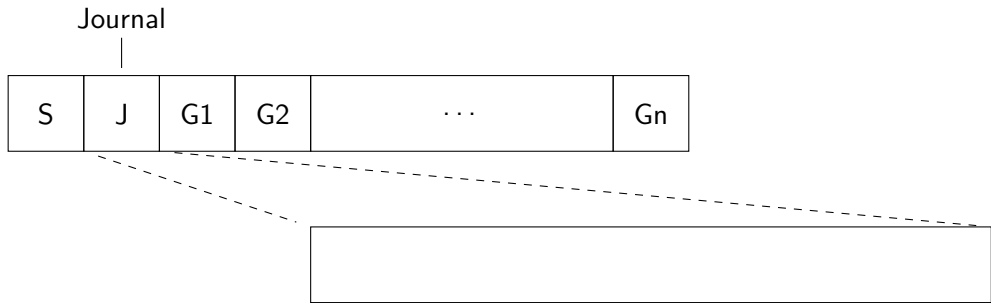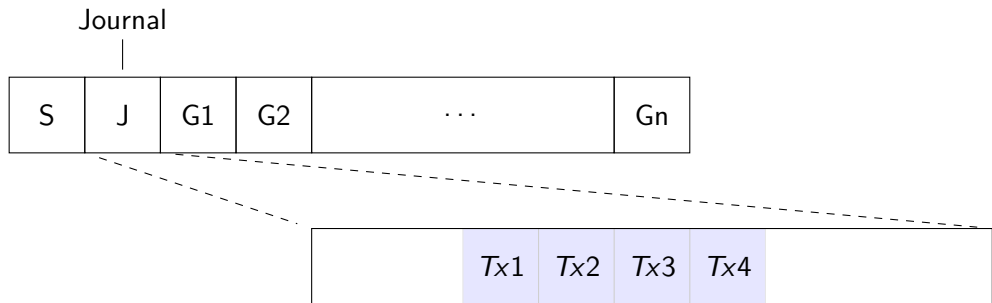We manage to write the whole transaction, updating the blocks but not remove the transaction.

We manage to write half of the transaction.

We manage to write the whole transaction but not updating the blocks.

We manage to write the whole transaction, updating the blocks but not remove the transaction.

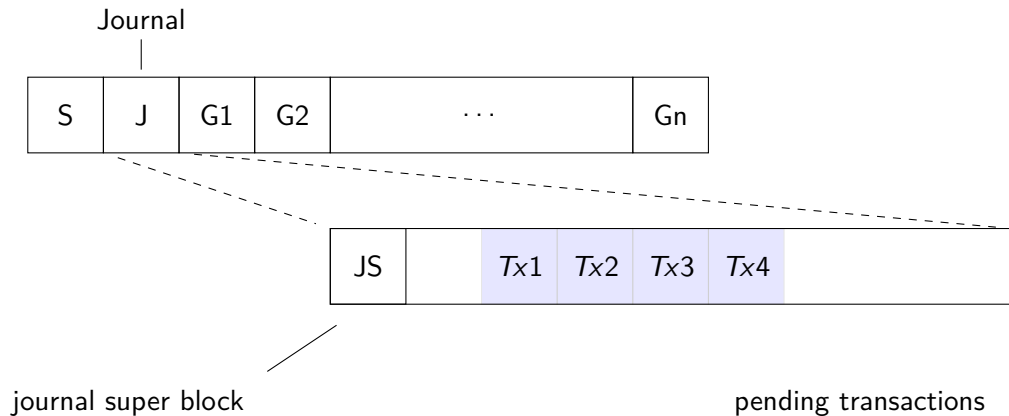We manage to write $TxB$, $I_{v2}$ and $TxE$ and then crash.

We manage to write half of the transaction.

We manage to write the whole transaction but not updating the blocks.

We manage to write the whole transaction, updating the blocks but not remove the transaction.

We manage to write $TxB$, $I_{v2}$ and $TxE$ and then crash.
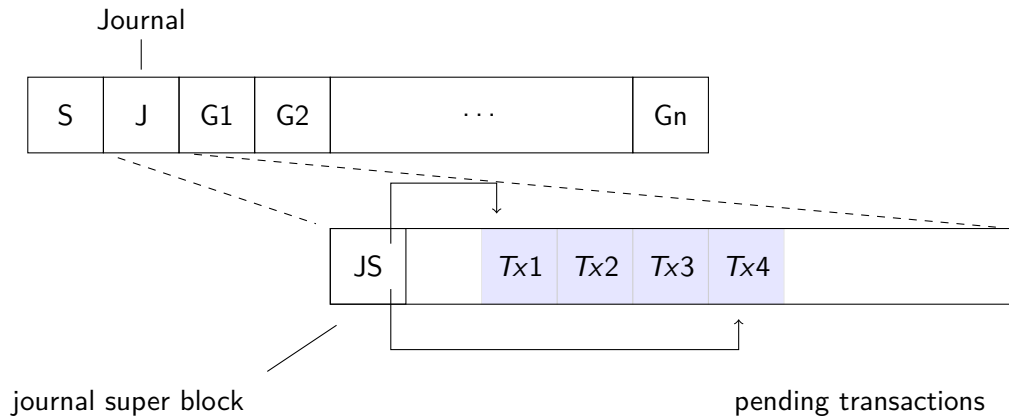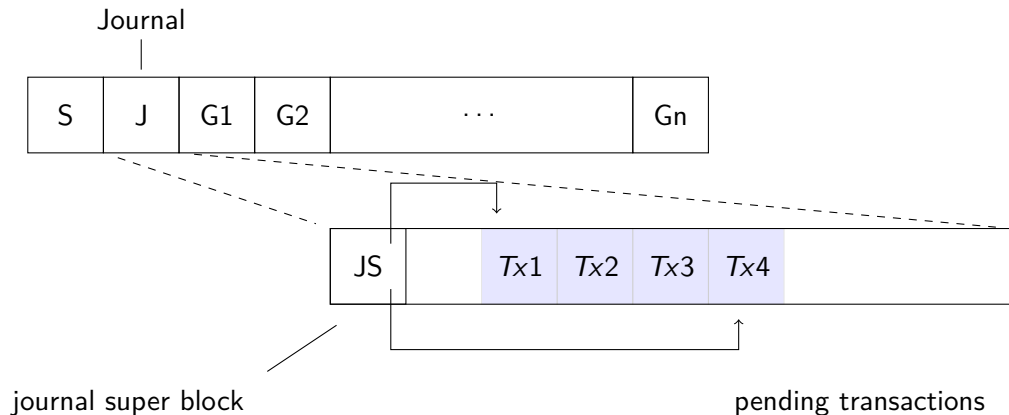
pending transactions

Journal

| S | J | G1 | G2 | $\cdots$ | Gn |

| JS | | Tx1 | Tx2 | Tx3 | Tx4 | |

journal super block                              pending transactions

## pending transactions



journal super block

pending transactions

Journal

| S | J | G1 | G2 | $\cdots$ | | Gn |

| JS | | Tx1 | Tx2 | Tx3 | Tx4 | |

journal super block                                        pending transactions

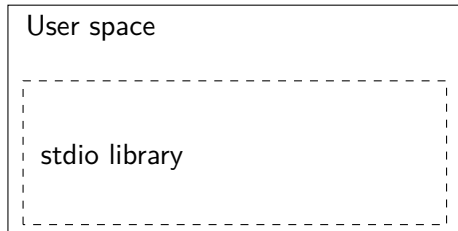What is the state of the file system?

journal super block

pending transactions

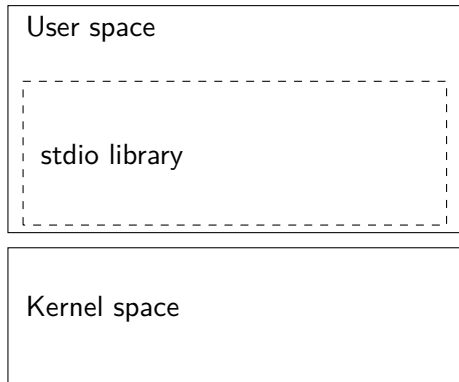What is the state of the file system?

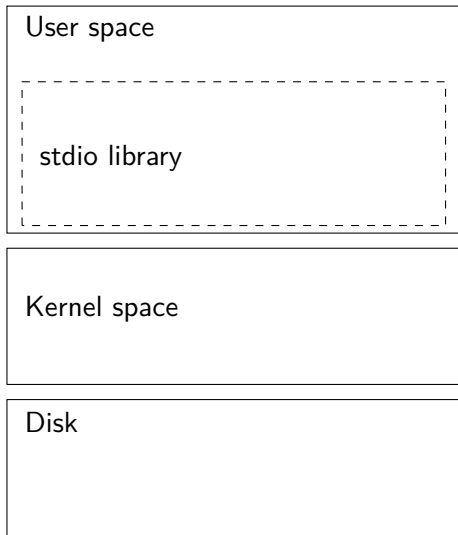Can we read from the file system?

User space

User space

stdio library

User space

stdio library

Kernel space

User space

stdio library

Kernel space

Disk

User space

fwrite()/fread()

stdio library

write buffer

Kernel space

Disk

## Layers of caches

User space

| | |
|---|---|
| | fwrite()/fread() |
| stdio library | |
| | write buffer |

Kernel space

Disk

- flush(): changes in buffer to kernel

User space

stdio library

fwrite()/fread()

write buffer

- flush(): changes in buffer to kernel

write()/read()

Kernel space

file blocks in memory

Disk

User space

stdio library

fwrite()/fread()

write buffer

write()/read()

Kernel space

file blocks in memory

Disk

- flush(): changes in buffer to kernel
- sync(): changes to file system journal/checkpoint

# Layers of caches

User space

- - - - - - - - - - - - - - - - - - - - - - - - -

fwrite()/fread()

stdio library

write buffer

- - - - - - - - - - - - - - - - - - - - - - - - -

write()/read()

Kernel space

file blocks in memory

Disk

pending transactions

checkpoint

- flush(): changes in buffer to kernel
- sync(): changes to file system journal/checkpoint
- checkpointing: from journal to inodes, maps and blocks

Journal is slow:

- Commit: write meta-data and data in a transaction (make sure it's a complete transaction).

Journal is slow:

- Commit: write meta-data and data in a transaction (make sure it's a complete transaction).
- Checkpointing: update the inode, bitmap and data blocks given the transaction.

Journal is slow:

- Commit: write meta-data and data in a transaction (make sure it's a complete transaction).
- Checkpointing: update the inode, bitmap and data blocks given the transaction.

Everything is written twice to disk!

Journal is slow:

- Commit: write meta-data and data in a transaction (make sure it's a complete transaction).
- Checkpointing: update the inode, bitmap and data blocks given the transaction.

Everything is written twice to disk!

*Idea - do the wrong thing and pray for the best.*

Faster:

- Commit data : write data <u>directly to block</u>.

Faster:

- Commit data : write data directly to block.
- Commit meta-data: when data is in block, write meta-data in transaction.

Faster:

- Commit data : write data <u>directly to block</u>.
- Commit meta-data: <u>when data is in block</u>, write meta-data in transaction.
- Checkpointing: update the inode and bitmap given transaction.

Faster:

- Commit data : write data <u>directly to block</u>.
- Commit meta-data: <u>when data is in block</u>, write meta-data in transaction.
- Checkpointing: update the inode and bitmap given transaction.

Even faster:

- Commit data : write data directly to block... eventually, hopefully.

Faster:

- Commit data : write data directly to block.
- Commit meta-data: when data is in block, write meta-data in transaction.
- Checkpointing: update the inode and bitmap given transaction.

Even faster:

- Commit data : write data directly to block... eventually, hopefully.
- Commit meta-data: write meta-data in transaction.

Faster:

- Commit data : write data directly to block.
- Commit meta-data: when data is in block, write meta-data in transaction.
- Checkpointing: update the inode and bitmap given transaction.

Even faster:

- Commit data : write data directly to block... eventually, hopefully.
- Commit meta-data: write meta-data in transaction.
- Checkpointing: update the inode and bitmap given transaction (let's hope the data is there)

Faster:

- Commit data : write data directly to block.
- Commit meta-data: when data is in block, write meta-data in transaction.
- Checkpointing: update the inode and bitmap given transaction.

Even faster:

- Commit data : write data directly to block... eventually, hopefully.
- Commit meta-data: write meta-data in transaction.
- Checkpointing: update the inode and bitmap given transaction (let's hope the data is there)

- `journal`: all data and meta-data is written through journal

- `journal`: all data and meta-data is written through journal
- `ordered` (default): data is written immediately to block, meta-data through journal

## ext4/jdb2

- `journal`: all data and meta-data is written through journal
- `ordered` (default): data is written immediately to block, meta-data through journal
- `write-back` : data is not guaranteed to be written before meta-data

# inode - is everything important

```
> sudo istat /dev/sda1 2236582
inode: 2236582                                    Group: 273
Generation Id: 3805640679
uid/gid: 1000/1000      mode: rrw-rw-r--   Flags: Extents,

size: 43  num of links: 1

Inode Times:
Accessed: 2016-12-06 14:51:17.003254544 (CET)
File Modified: 2016-12-06 15:46:55.667041193 (CET)
Inode Modified: 2016-12-06 15:46:55.667041193 (CET)
File Created: 2016-12-06 13:39:15.084806928 (CET)

Direct Blocks:  6946002
```

*This album has nothing to do with the following material.*

Reading is mostly done from cached copies in memory.

Reading is mostly done from cached copies in memory.

Focus on write operations, try to avoid moving the arm.

Reading is mostly done from cached copies in memory.

Focus on write operations, try to avoid moving the arm.

Writing is best done in large consecutive segments.

Reading is mostly done from cached copies in memory.

Focus on write operations, try to avoid moving the arm.

Writing is best done in large consecutive segments.

The state of the file system is *a log of events*.

D

| D | D |
|---|---|

| D | D | D |

| D | D | D | i7 |

| D | D | D | i7 |

D  D  D  i7  D  D

| D | D | D | i7 | D | D | i9 |

| D | D | D | i7 | D | D | i9 |

| D | D | D | i7 | D | D | i9 |

| D | D | D | i7 | D | D | i9 |

| D | D | D | i7 | D | D | i9 | D |

| D | D | D | i7 | D | D | i9 | D | i7 |

| D | D | D | i7 | D | D | i9 | D | i7 |

| D | D | D | i7 | D | D | i9 | D | i7 |

How do we find the inodes?

D

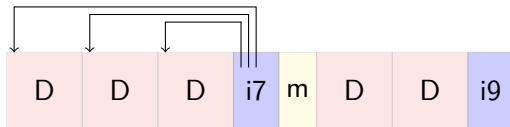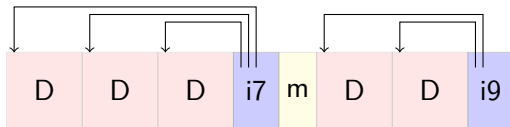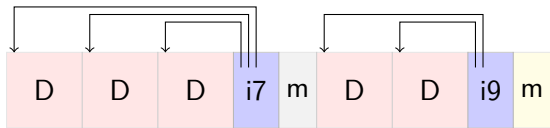| D | D | D |
|---|---|---|

The inode map holds mapping from inode number to block addresses.

The inode map holds mapping from inode number to block addresses.

The inode map holds mapping from inode number to block addresses.
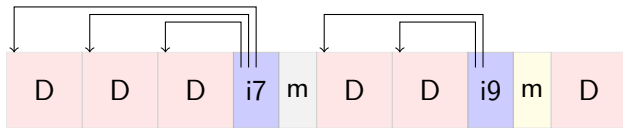
The inode map holds mapping from inode number to block addresses.

The inode map holds mapping from inode number to block addresses.
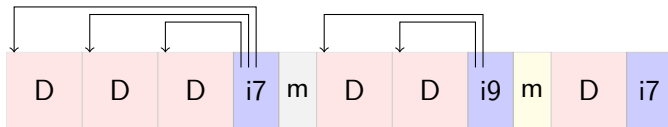
The inode map holds mapping from inode number to block addresses.

The inode map holds mapping from inode number to block addresses.
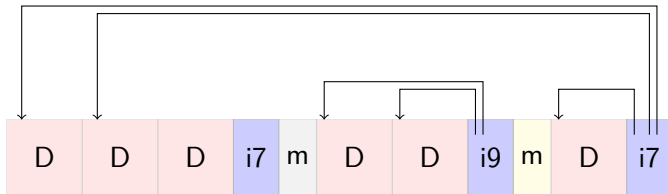
The inode map holds mapping from inode number to block addresses.

The inode map holds mapping from inode number to block addresses.

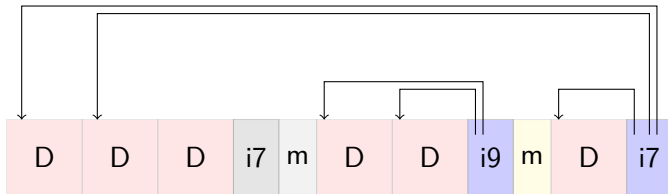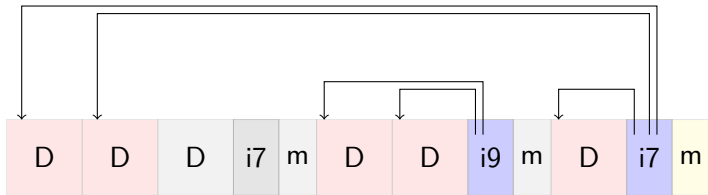The inode map holds mapping from inode number to block addresses.

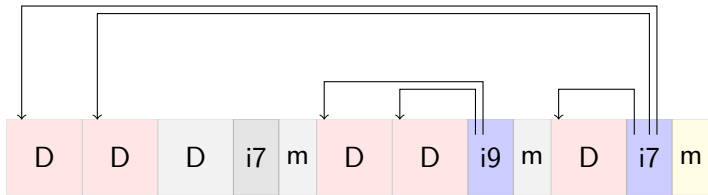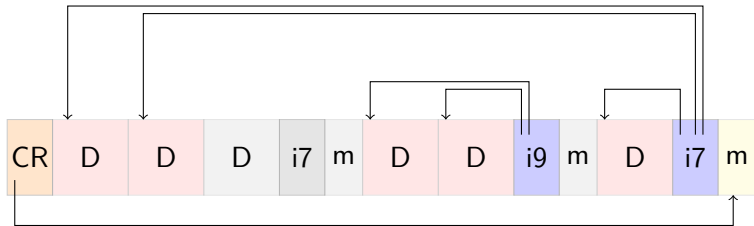The inode map holds mapping from inode number to block addresses.

How do we find the last inode map?

The inode map holds mapping from inode number to block addresses.

How do we find the last inode map?

The inode map holds mapping from inode number to block addresses.

How do we find the last inode map?

reading a file

- read the check region
- find the location of the inode map
- find inode
- read data block

## pros and cons

reading a file

- read the check region
- find the location of the inode map
- find inode
- read data block

writing a file

- write data block
- write new copy of inode
- write new copy of inode map
- update check region

## pros and cons

reading a file

- read the check region
- find the location of the inode map
- find inode
- read data block

writing a file

- write data block
- write new copy of inode
- write new copy of inode map
- update check region

How much can we cache in memory?

## pros and cons

reading a file

- read the check region
- find the location of the inode map
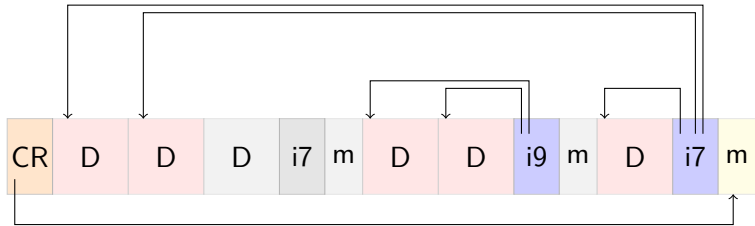- find inode
- read data block

writing a file

- write data block
- write new copy of inode
- write new copy of inode map
- update check region

How much can we cache in memory?

Can we delay updating the check region?

CR | D | D | D | i7 | m | D | D | i9 | m | D | i7 | m

CR D D D i7 m D D i9 m D i7 m

Where is the bit map that keeps track of available blocks?

Do we want to know where to find blocks ..

Do we want to know where to find blocks ..

if they are scattered around the disk?

SS

segment summary
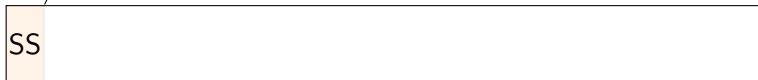
SS

segment summary

SS    A i m B B C i m A i m

segment summary

SS          A i m B B C i m A i m

Segment summary keeps a mapping from block to inode.

segment summary



Segment summary keeps a mapping from block to inode.

segment summary



Segment summary keeps a mapping from block to inode.

segment summary



Segment summary keeps a mapping from block to inode.

segment summary



Segment summary keeps a mapping from block to inode.

segment summary

SS       B C i m A i m DDD i m E i m AA i m

Segment summary keeps a mapping from block to inode.

segment summary

SS    C i m A i m DDD i m E i m AA i m
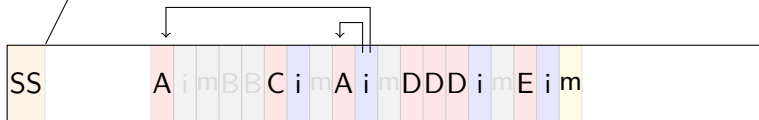
Segment summary keeps a mapping from block to inode.

segment summary



Segment summary keeps a mapping from block to inode.

Segment summary keeps a mapping from block to inode.

*The file system UDF used a log structure to do updates on a write-once CD/DVD*

- ext4 : default Linux system, journaling

- `ext4` : default Linux system, journaling
- `F2FS` : by Samsung, log-structured, optimised for SSD

- `ext4` : default Linux system, journaling
- `F2FS` : by Samsung, log-structured, optimised for SSD
- `NILFS` : Nipon Telecom, log-structured

## Some file systems

- `ext4` : default Linux system, journaling
- `F2FS` : by Samsung, log-structured, optimised for SSD
- `NILFS` : Nipon Telecom, log-structured
- `btrfs` : originally by Oracle, a copy-on-write system

- `ext4` : default Linux system, journaling
- `F2FS` : by Samsung, log-structured, optimised for SSD
- `NILFS` : Nipon Telecom, log-structured
- `btrfs` : originally by Oracle, a copy-on-write system
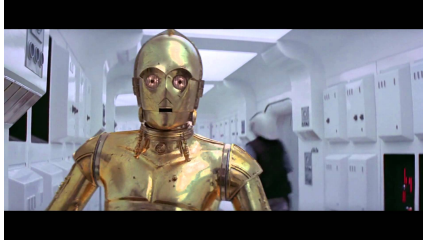- `APFS` : next generation for OSX (Sierra 2017), copy-on-write

## Some file systems

- `ext4` : default Linux system, journaling
- `F2FS` : by Samsung, log-structured, optimised for SSD
- `NILFS` : Nipon Telecom, log-structured
- `btrfs` : originally by Oracle, a copy-on-write system
- `APFS` : next generation for OSX (Sierra 2017), copy-on-write
- `ReFS` : latest file system for Windows servers, copy-on-write

# Some file systems

- `ext4` : default Linux system, journaling
- `F2FS` : by Samsung, log-structured, optimised for SSD
- `NILFS` : Nipon Telecom, log-structured
- `btrfs` : originally by Oracle, a copy-on-write system
- `APFS` : next generation for OSX (Sierra 2017), copy-on-write
- `ReFS` : latest file system for Windows servers, copy-on-write
- `exFAT` : Microsoft system used by SD cards
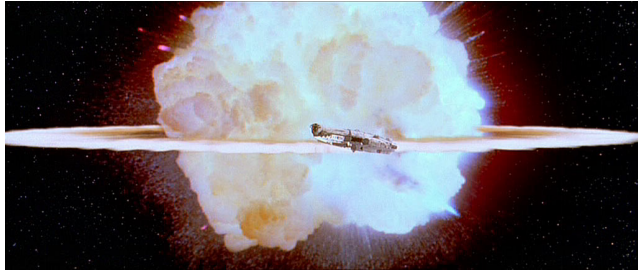
## Some file systems

- `ext4` : default Linux system, journaling
- `F2FS` : by Samsung, log-structured, optimised for SSD
- `NILFS` : Nipon Telecom, log-structured
- `btrfs` : originally by Oracle, a copy-on-write system
- `APFS` : next generation for OSX (Sierra 2017), copy-on-write
- `ReFS` : latest file system for Windows servers, copy-on-write
- `exFAT` : Microsoft system used by SD cards

We're doomed!

Saved!