

Journaling and Log-structured file systems

Johan Montelius

KTH

2020

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

1 / 35

2 / 35

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
- 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?

what if we crash

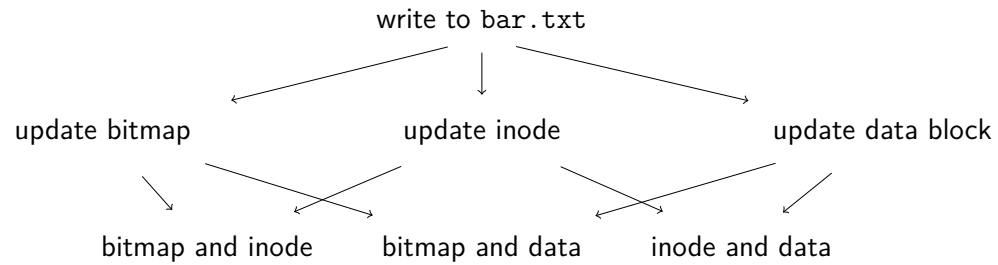
How do we cope with crashing drives?

How do we cope with the operating system crashing?

3 / 35

4 / 35

one or two out of three



5 / 35

two out of three ...

6 / 35

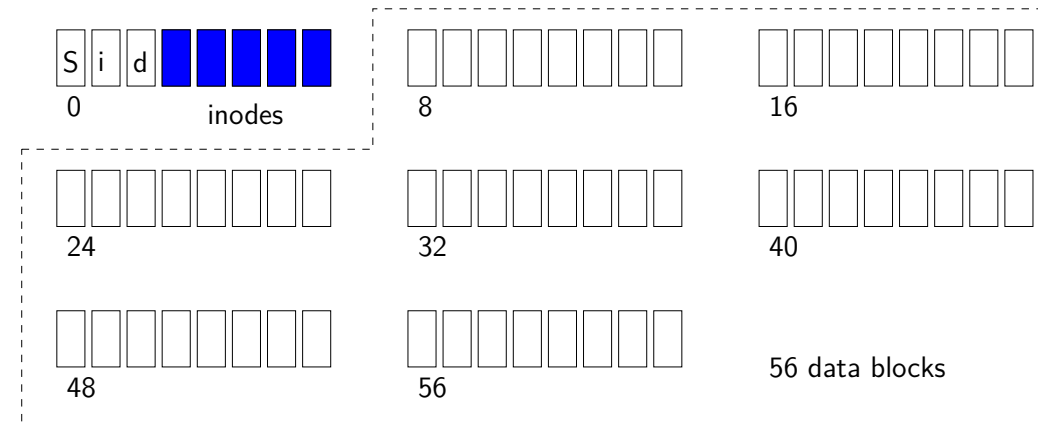
Survive a crash

crash recovery

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:



7 / 35

8 / 35

```
$ 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 blo
```

9 / 35

10 / 35

```
> 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*.

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

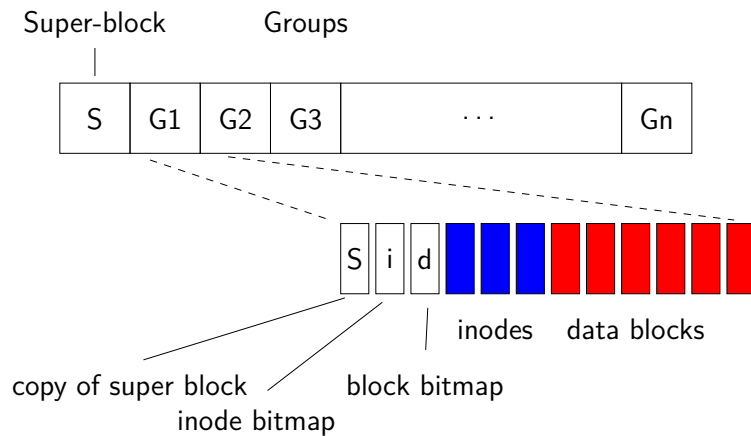
Journal or Write-Ahead Logging

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

11 / 35

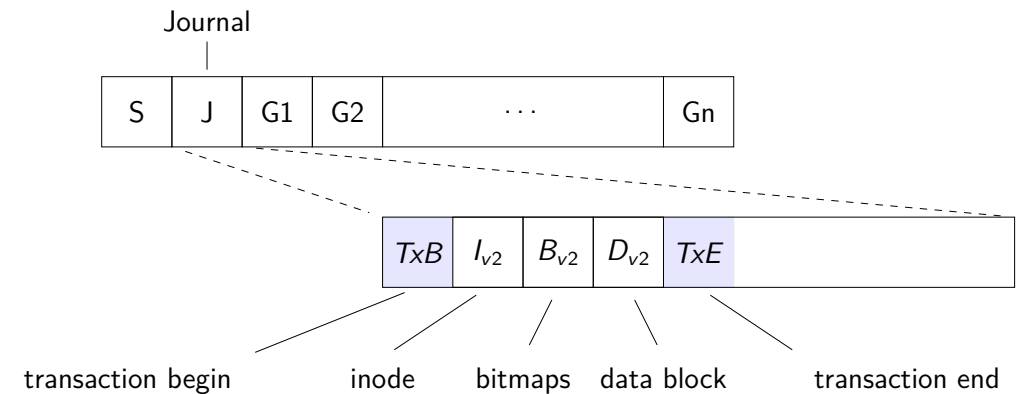
12 / 35

Linux ext2 - no journaling



13 / 35

Linux ext3 - journaling



14 / 35

the safe way

- 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

15 / 35

disaster scenarios

We manage to write half of the transaction.

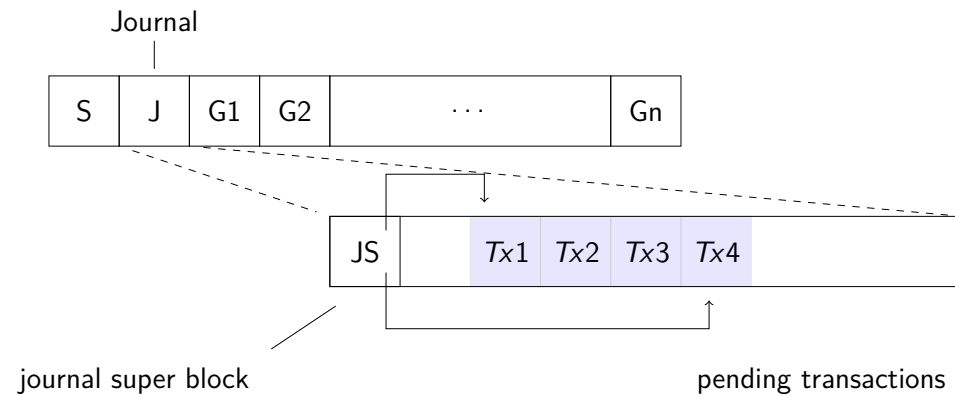
We manage to write the whole transaction but not updating the blocks.

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.

16 / 35

pending transactions

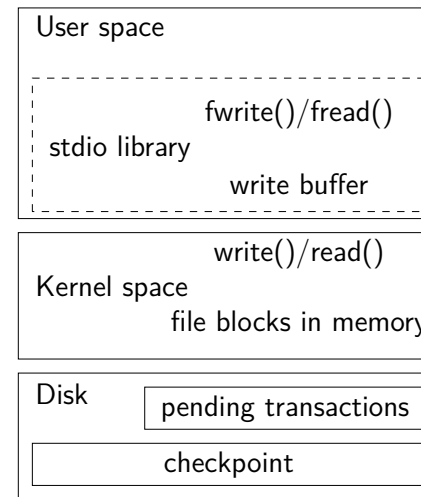


What is the state of the file system?

Can we read from the file system?

17 / 35

Layers of caches



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

18 / 35

do the right thing....?

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.

19 / 35

data only once

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)

20 / 35

- 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

```
> 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)
```

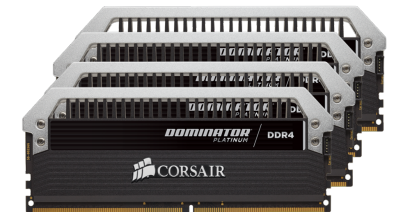
```
Inode Modified: 2016-12-06 15:46:55.667041193 (CET)
```

```
File Created: 2016-12-06 13:39:15.084806928 (CET)
```

```
Direct Blocks: 6946002
```

21 / 35

22 / 35



This album has nothing to do with the following material.

23 / 35

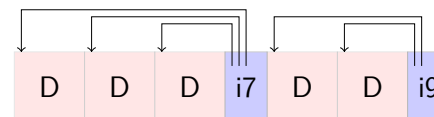
24 / 35

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.

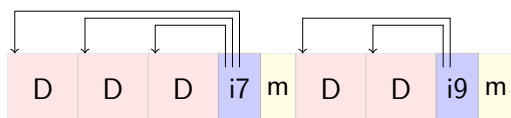
The state of the file system is *a log of events*.



How do we find the inodes?

25 / 35

26 / 35



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

writing a file

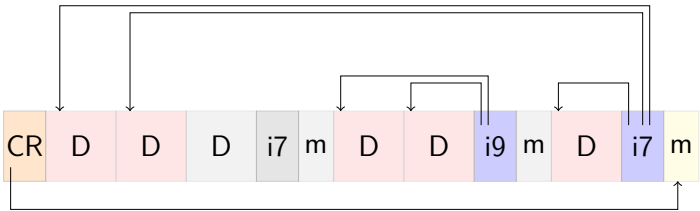
- write data block
- write new copy of inode
- write new copy of inode map
- update check region

How much can we cache in memory?

Can we delay updating the check region?

27 / 35

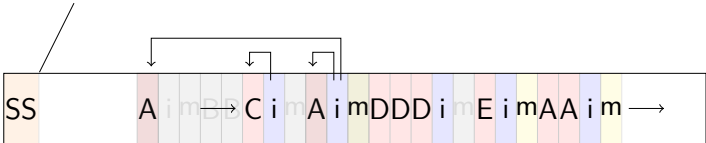
28 / 35



Do we want to know where to find blocks ..
if they are scattered around the disk?

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

segment summary



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

33 / 35

- 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

34 / 35