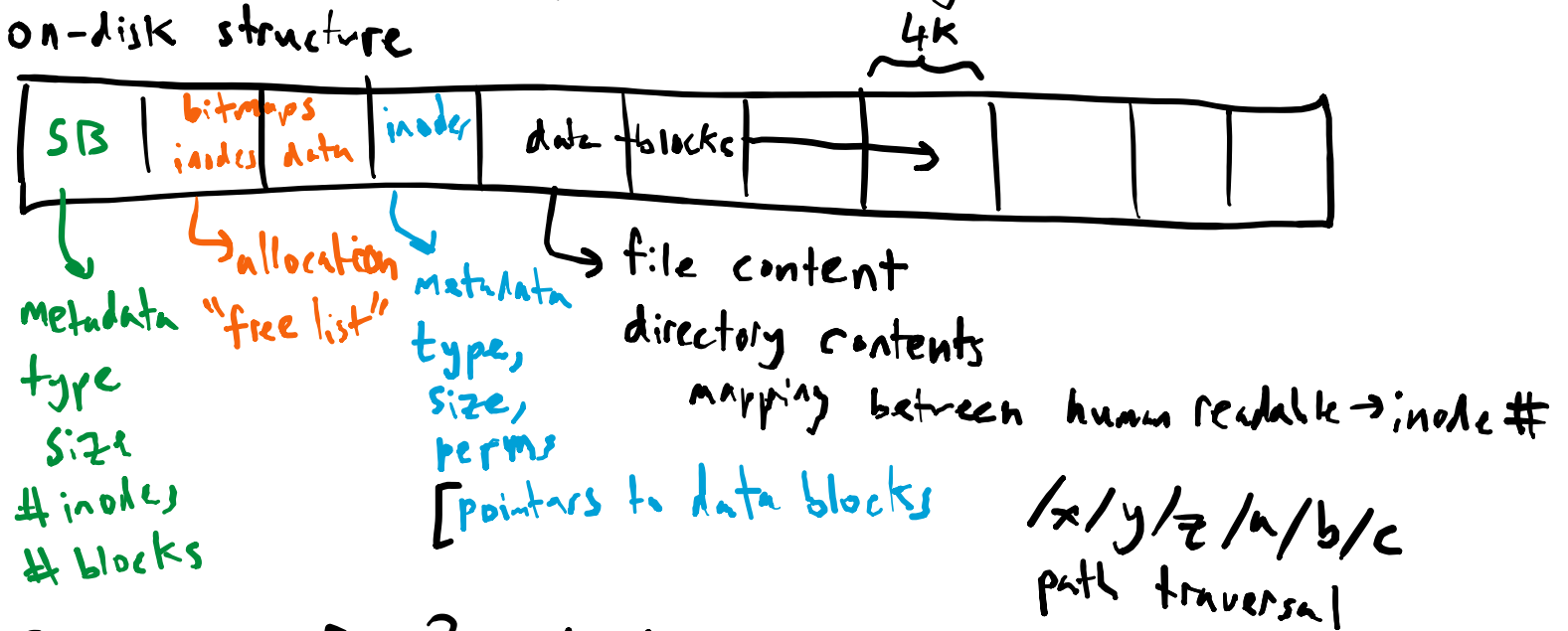


Last time: VSFS (very simple file system)

today: efficiency Fast File System

on-disk structure

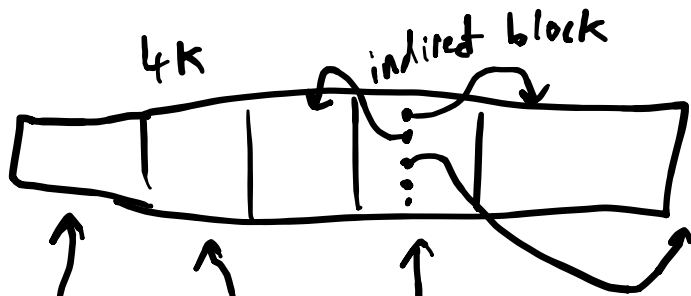


Q's: large files? inode design

what about efficiency? caching... FFs

inode:

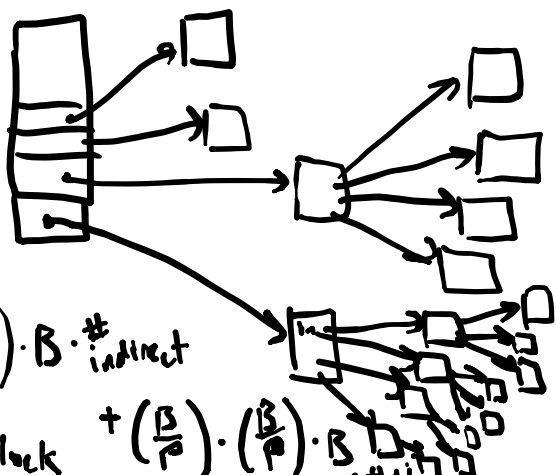
- type
- size
- perms
- pointers to data blocks



direct pointer

indirect pointer

double indirect pointer

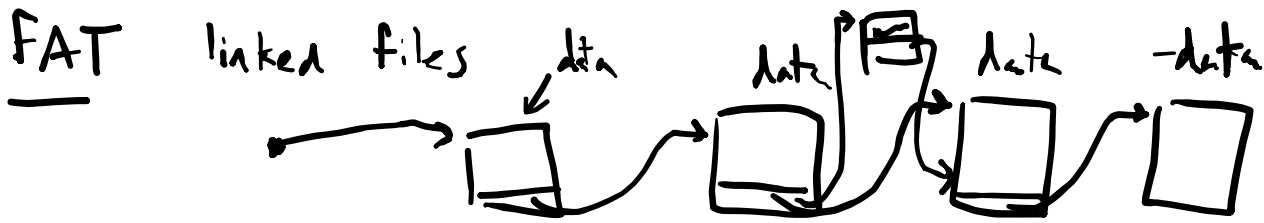


blk size B

Size file we can handle

$$\# \text{ direct} \cdot B + \left(\frac{B}{P}\right) \cdot B \cdot \# \text{ indirect} + \left(\frac{B}{P}\right) \cdot \left(\frac{B}{P}\right) \cdot B \cdot \# \text{ triple indirect}$$

pointers per block



Extents: (ptr + length) contiguous blocks

Efficiency: /a/b/c/d/e/f.txt

root, inode
/ data
a inode
data
b inode
⋮

works well for reads

Keep in memory

buffer-cache: just FS

page cache: unified virtual memory

what about writes?

in background, write to disk

[batching, scheduling]

how frequently should you write

Old UNIX FS

1. API everyone loved

2. performance was terrible

2% of peak & getting worse

problems: seeking too much



{ → inodes & data not "close"

{ → related files? not "close"

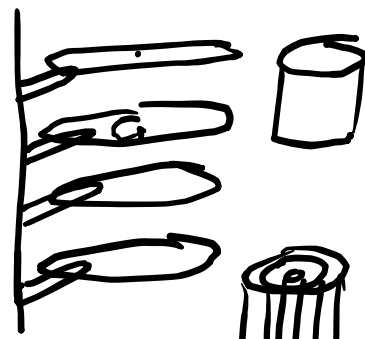
→ fragmentation (internal)

overall: FS assumed disk was like RAM

→ treat disk like a disk ⇒ FFS

FFS: treat disk like a disk

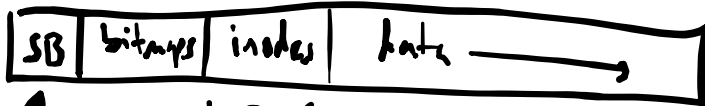
- 1) chop disk up into groups
- 2) put related things in same group



"cylinder group" same cylinder = no seek

today: device doesn't expose geometry
but locality does have perf differences
"block groups"

Each group:



copy
for
reliability

piece
of
inode
table

Allocation policies
related?

- inode + its data
- files within same directory (namespace-based)

ex
mkdir: pick group
lots of free inodes, blocks
low # other dir.

Create
file: put data in same group
as inode

Large files: fill up entire group w/ 1 file

FFS rule: once $>$ size
file
place next block in another group

splitting a large file
between groups is
NOT bad

groups $\left\{ \begin{array}{l} 0 / a \quad aaaaaaaaaa \\ 1 / b \\ 2 \end{array} \right.$

Amortization

seek cost is amortized
over the transfer

0 $1/6$ | aaaaaa bb
1 | aaaa
2 | aaaa

xfer: $\left\{ \begin{array}{l} 40 \text{ MB/s} \\ 10 \text{ ms (avg seek)} \end{array} \right.$

How much data
should be transferred
to amortize seek
time

target
efficiency \rightarrow 90% efficiency

$$\frac{40 \text{ MB}}{\text{sect}} \cdot \frac{1 \text{ sec}}{1000 \text{ ms}} \cdot 90 \text{ ms} \approx \underline{\underline{3.6 \text{ MB}}} \frac{90 \text{ ms}}{10 \text{ ms}} = 90\%$$

Other FFS stuff

Changed impl. without changing API

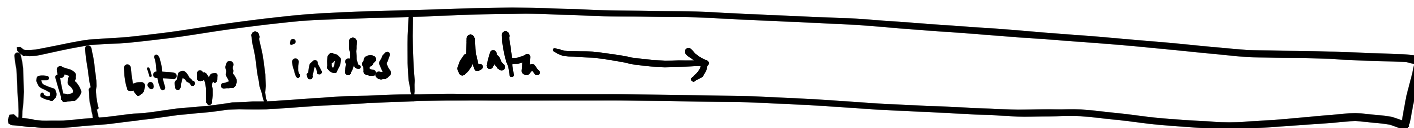
\rightarrow symbolic links

\rightarrow long file names $>$ 8 char

• parameterization (e.g. rotation-aware)

• sub-blocks reduce internal fragmentation (buffer writes)

Understanding structure



mmap

(struct sb *) mmap

sb → size

Crash consistency

