Hello! 4/17

Welcome to the
CS Picnic!

# File Systems : Implementation

-> Simple : Very Simple
File System
(VSFS)

-> Locality : Fast File
System
(FFS)

-> Crash Consistency :
Journaling (JFS)
(write Ahead . ext3/4

Logging ) Linux

Prof $\longleftrightarrow$ Student

$\downarrow$     $\rightarrow$ minimize   $\downarrow$

$\rightarrow$ <u>Write Performance</u>:

$\Big\{$ Logging $\rightarrow$

Log-Structured

File System (<u>LFS</u>)

$\rightarrow$ Flash-based

Solid State Drives

(<u>SSDs</u>)

$\rightarrow \big[$ <u>The End</u> $\big] \Leftarrow$ future:

sadness

Abstraction:

$\rightarrow$ Files     ?

→ Directories )

File: An array of
underline{bytes}

has ⟹
low-
level
name
(number)
{
Operations:
open, read/write,
unlink, close
}

Directory:
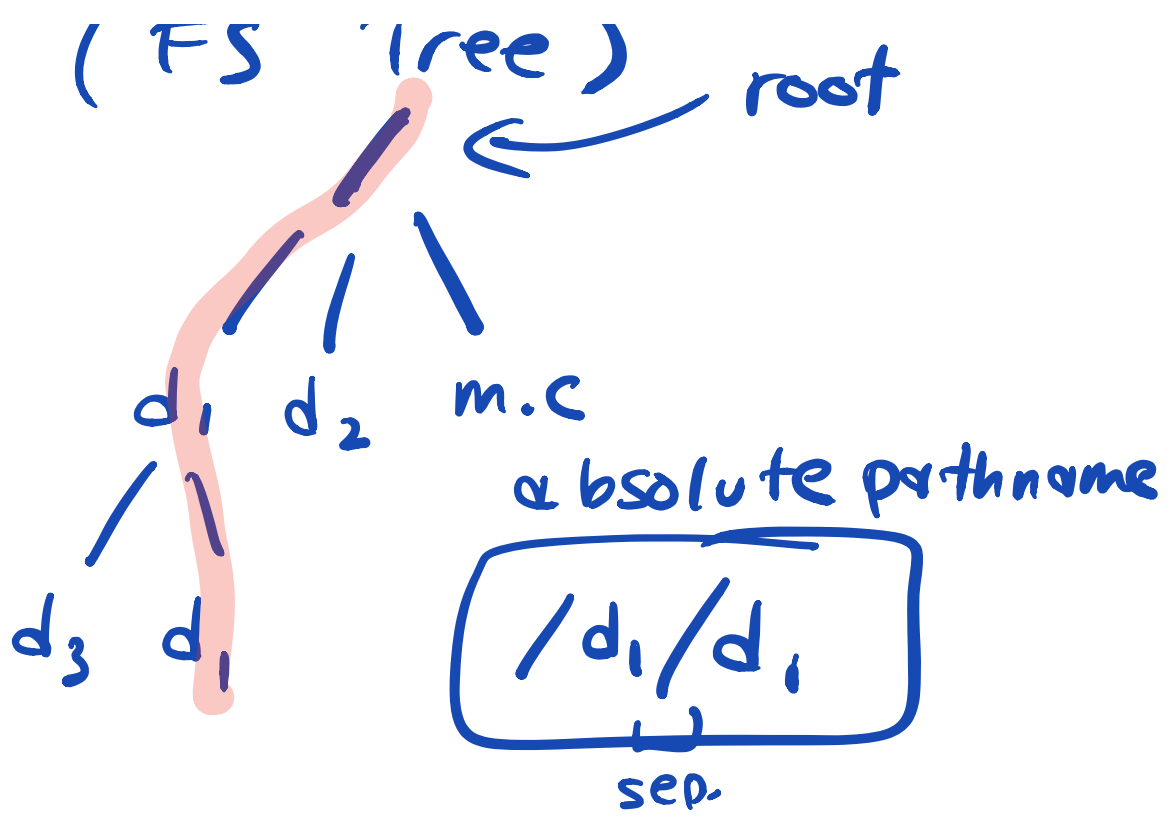list of ⟨names,
(of files,)  low-level name⟩
dirs

one file:
FS ⟶ file $\underline{3}$

"main.c"

dir: maps
"main.c" ⟶ 3

Directory Hierarchy:

( FS Tree ) ⟵ root

d₁  d₂  m.c

d₃  d₁

absolute pathname

$$\boxed{\;/d_1/d_i\;}$$

sep.

## Implementation :

→ On-disk Data Structures

→ Access Methods
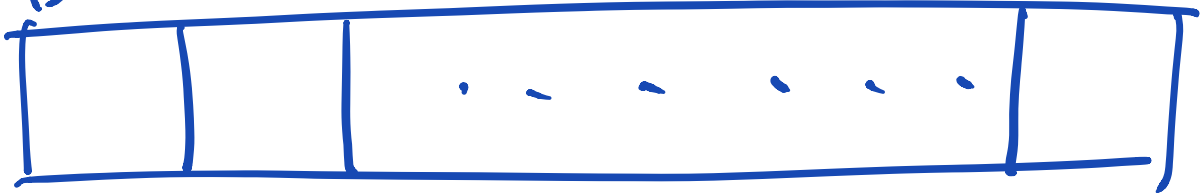open, read/write, etc.

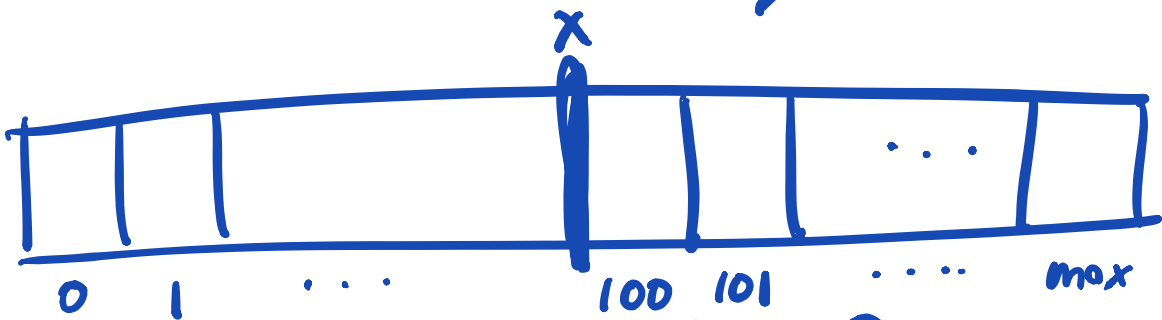## Very Simple File System

(VSFS)

FS : API   open   read ....

read / write to
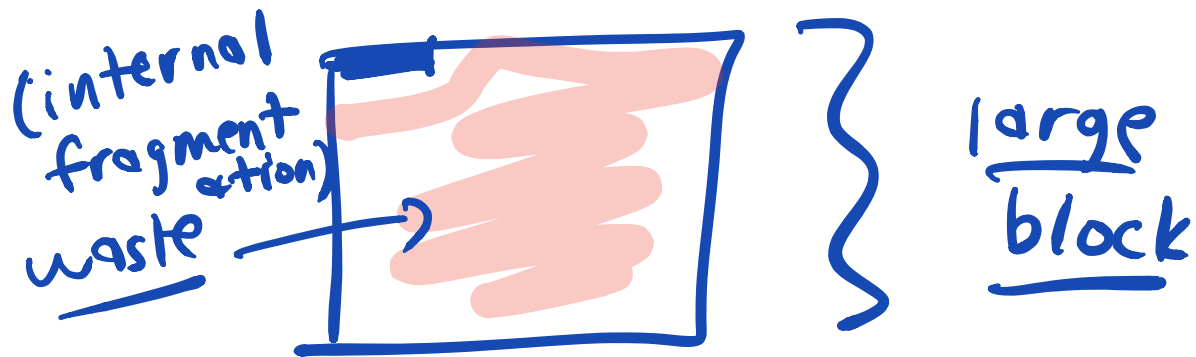   disk blocks

Disk : Bunch of Blocks

what needs to be
   on disk?

→ Data   "data blocks"

x

0   1   ...      100  101   ....   max

(Block) size: ⇒ (4KB)
   must be multiple
   of sector size (512b)

↳ minimal unit of
    allocation → file
                  ↘ directory

main.C :
    #include <stdio.h>

(internal
fragment-
    ation)
waste ————→               } large
                                           block

small blocks:
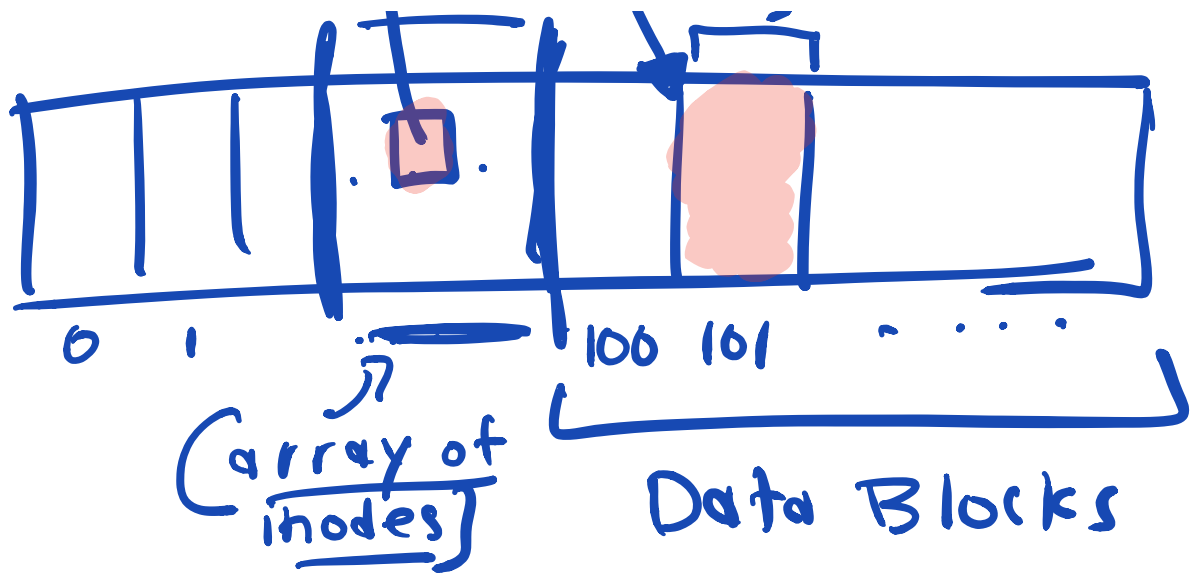    better : less internal
                          waste

    but : lots + lots of
                blocks in
                large file

Blocks : (4KB) ⟹ also:
                             matches
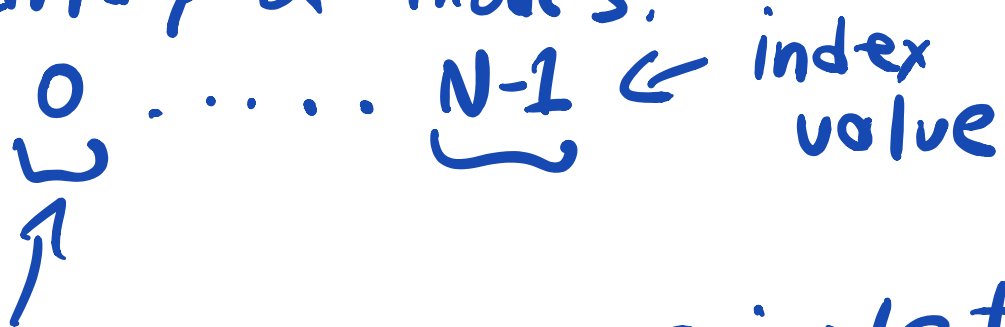                             page size
                     4KB

0    1    ...    100    101    ...

(array of inodes)

Data Blocks

1) Data

2) Info about each file
per file: structure on disk
[inode]    (index node)

↳ what's inside?
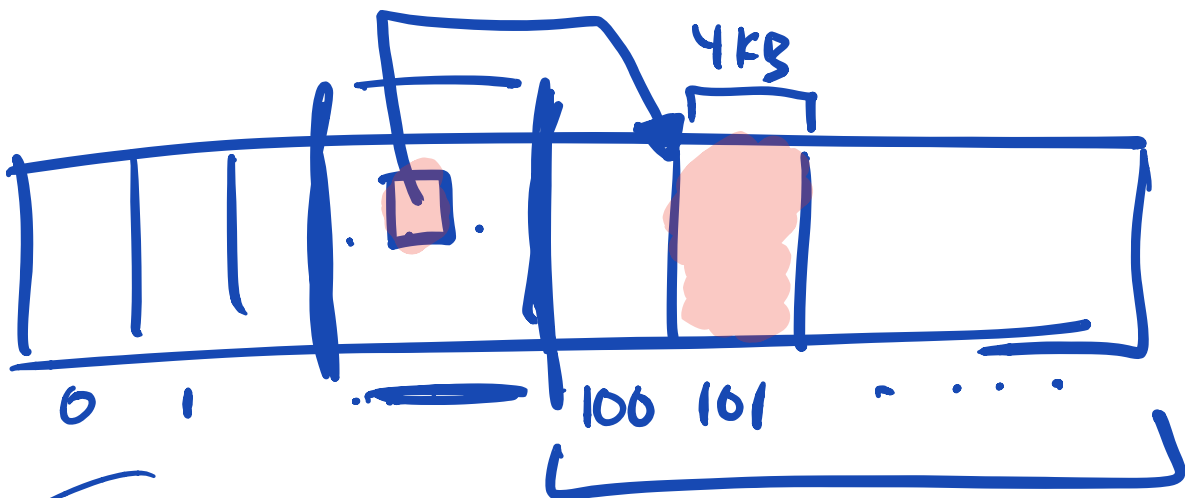
→ type: regular file
         directory

→ ownership: remzi

→ permission:
    who can read/
            write/exec.

→ timestamps

→ size, # of bits
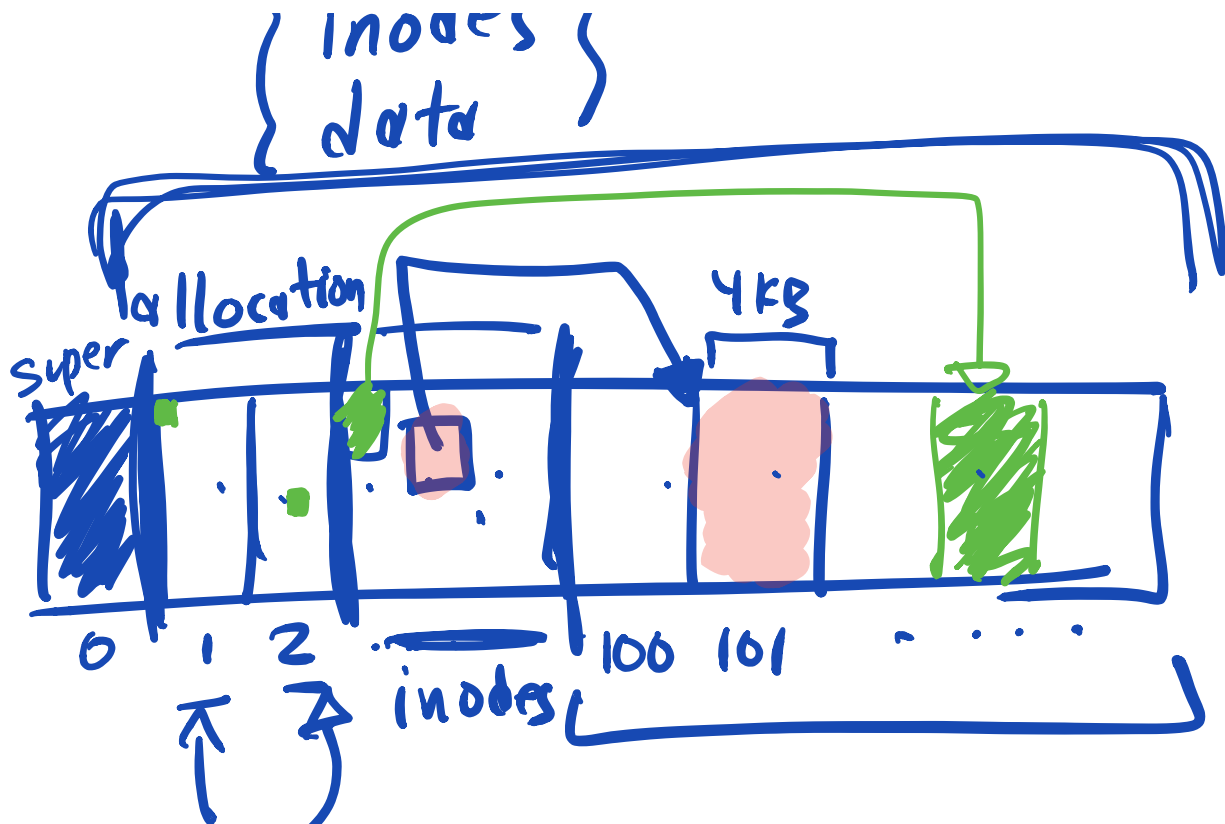
→ location of
blocks of
this file/dir

array of inodes:

0 . . . . . . N-1 ← index
value

low-level name ⟹ inode #
(index
value)



4kB

0   1   . .   100  101   . . . . .

given inode #:

=> <u>size</u>　　　　128 b, 256 b

=> <u>base</u>
　　　(where array starts)

1) <u>data</u> (file, dir)　　a.k.a
2) <u>inode array</u>　(inode table)

3) <u>directories</u> :
　　　treated just like files
　　　=> inode —→ data
　　　　　　　　↘ blocks

data

all dirs:
　.. <u>parent</u>
　. <u>this dir</u>
　:

| .. | 10 |
| . | 100 |
| foo: | 101 |

name　　　low-
(human)　level
　　　　　(inode
　　　　　　#)

4) track free/allocated:

inodes
data

allocation

super

4kB

0  1  2  ...  100  101  . . . .
inodes

1: inode bitmap

2: data bitmap

↓

0 or 1 → whether the corresp. entity is free / in use
(0)  (1)

4KB

1) data    2) inode

3) dir    4) bitmaps

5) superblock :
    per file system info.
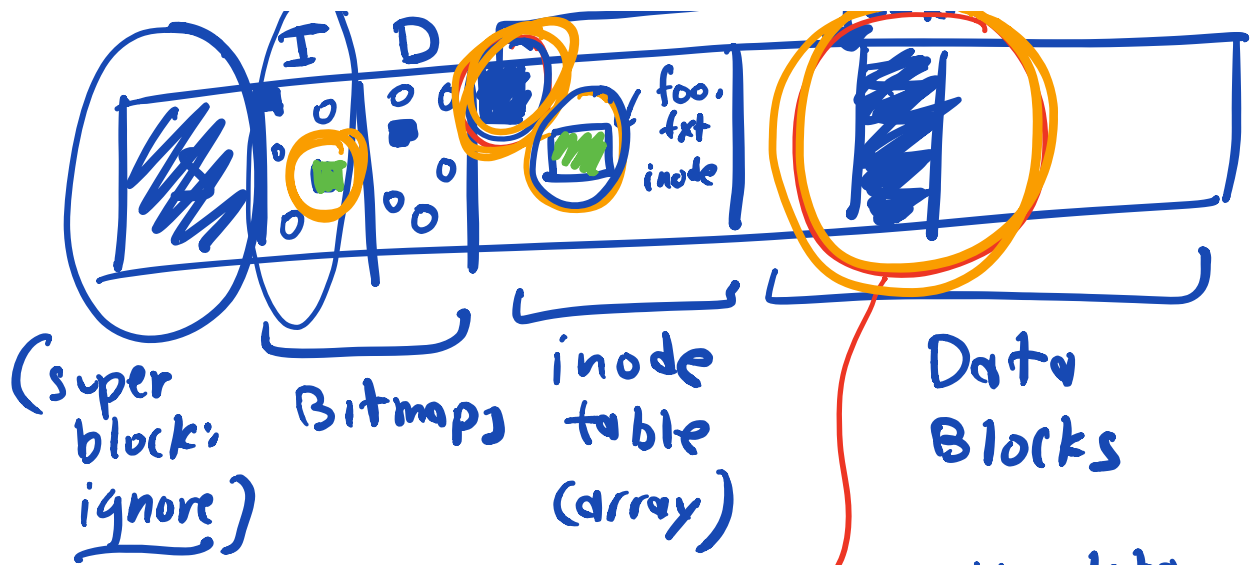    → block size,    type: ext4
      [where is inode table,
       where are data blocks?
       etc. ]

## Access Methods :

API: → what it does
        to disk?

creat ("/foo.txt")

Assume: Empty File System

(super block: ignore)  Bitmaps  inode table (array)  Data Blocks

dir. data block

| | |
|---|---|
| .. | 0 |
| . | 0 |
| foo.txt | 1 |

create /foo.txt

=> read inode of root directory

know: root inode # => 0

1) READ: block containing inode #0 (root inode)

(permissions check)

2) make sure filename is unique:

READ(s) data blocks

(~READ~ (s)) : ... of __root__
__directory__

3) look for space for inode:

( __READ__ ) inode bitmap

4) mark __bit__ (in use) [in memory]

5) ( __WRITE__ ) inode bitmap

6) ~update inode of __foo.txt__

( __READ__ ) inode __block__
__table__

update in memory:
[. . . . .]

( __WRITE__ ) block (inode)
of
foo.txt

7) link file into directory

(READ) root data

update [foo.txt → 1]
↗
directory
entry

(WRITE) dir block

⇒ creation: finished

───────────────

⇒ write ( )

{ ⇒ allocation
or
not? }

⟺) open
("/a/b/c/...")

[path
traversal]

open ( "/d₁/d₂/d₃ (foo.txt) )

root dir

[ read root inode
read root dir →  [          ]

read root's data

read $d_1$'s inode

read $d_1$'s data

read $d_2$'s inode

read $d_2$'s data

read $d_3$'s inode

read $d_3$'s data

read foo.txt's
        inode

| $d_1$ | 10 |
|------|----|

| | |
|------|----|
| $d_2$ | 12 |

| foo.txt | 99 |
|---------|----|

per-process:
   file descriptor (int)

open file table

( in-memory <u>inode</u> )

=) small <u>files</u> ( <u>so far</u> )

what about <u>large</u> <u>files</u>?

inode structure: 128 bytes

| type |
| <u>size</u> |
| <u>blocks</u> |
| perm |
| owner |
| time |
| ÷ |
| = |
| 100 |
| 101 |
| 102 |
| . |
| ∴ |

12 direct ptis

( address ) pointers ( to disk blocks ) ( indirect block )

indirect pointer

full of pointers

double ind. ptr

1024
more
blocks

```
1000
1001
  .
  .
  .
```
} 4KB

double ind
block

4 Bytes
ind. blks    → D

                 → D
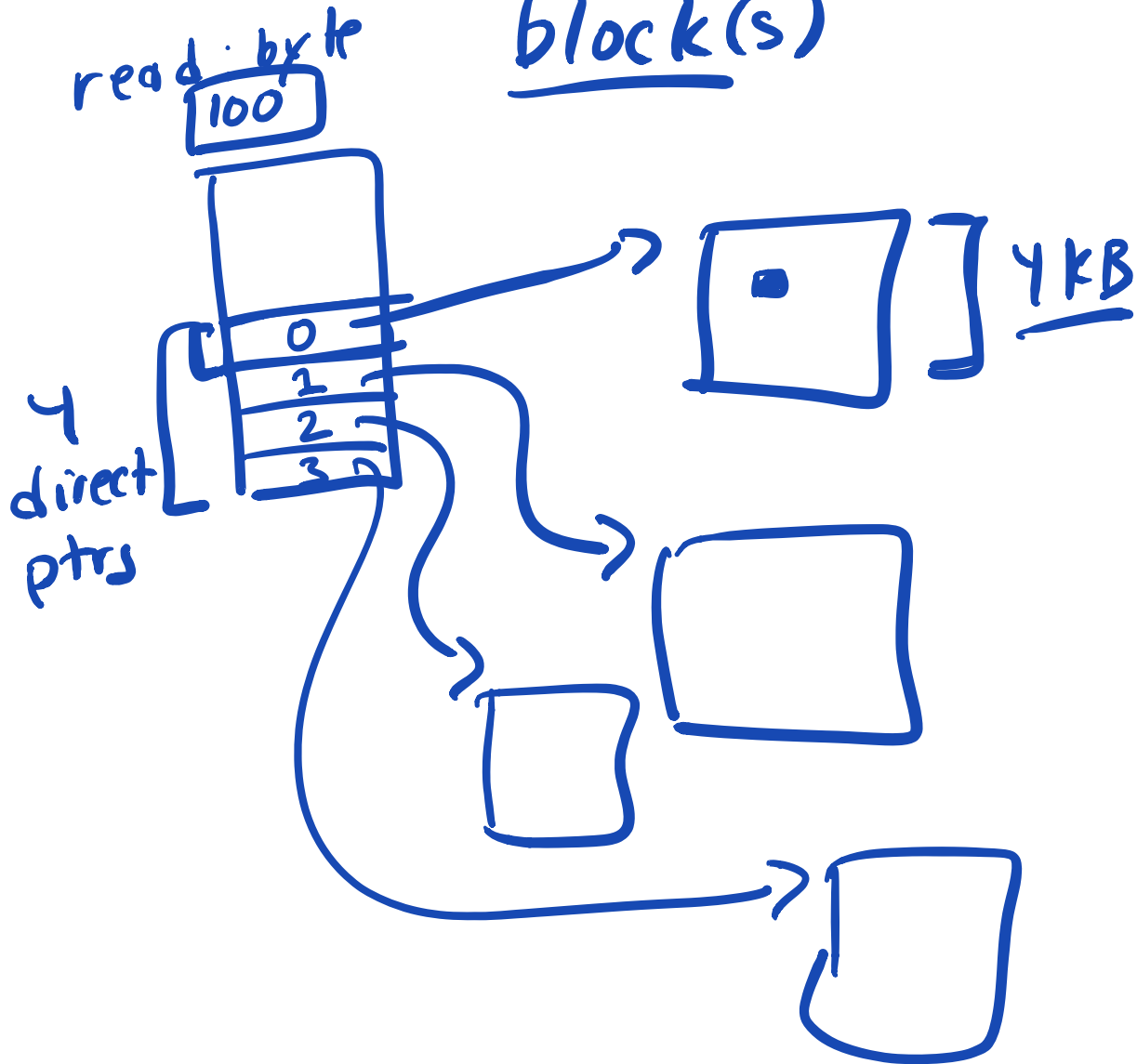
ind blks    → D
                → D

# Mapping : File ⟹ Block

read / write
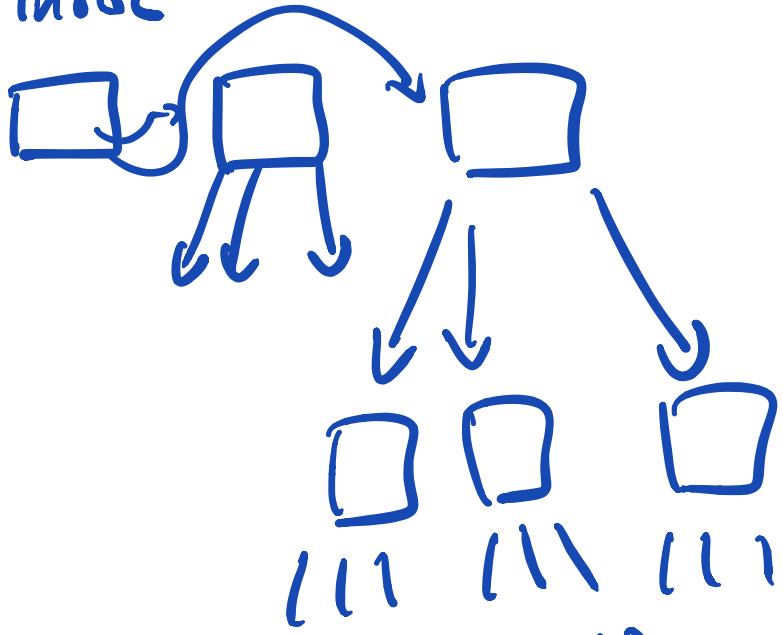                        starts @

$\Rightarrow$ (current offset)    $\underline{O}$

(byte offset)

FS : translating

byte offset $\Rightarrow$

block(s)

read byte

| 100 |

4
direct
ptrs

] 4 KB

FS inode



(imbalanced)

[most files are small]

→ Caching (write Buffering)

→ (Project 4a
=> )

P5 yes
(I-O)

Everyone gets 100!
here

Sign here

→ *[signature]*

Caching / Write Buffering

machine



VM
pages

pretty
fast

slow

=) Use memory
(for FS stuff)

Important on reads:

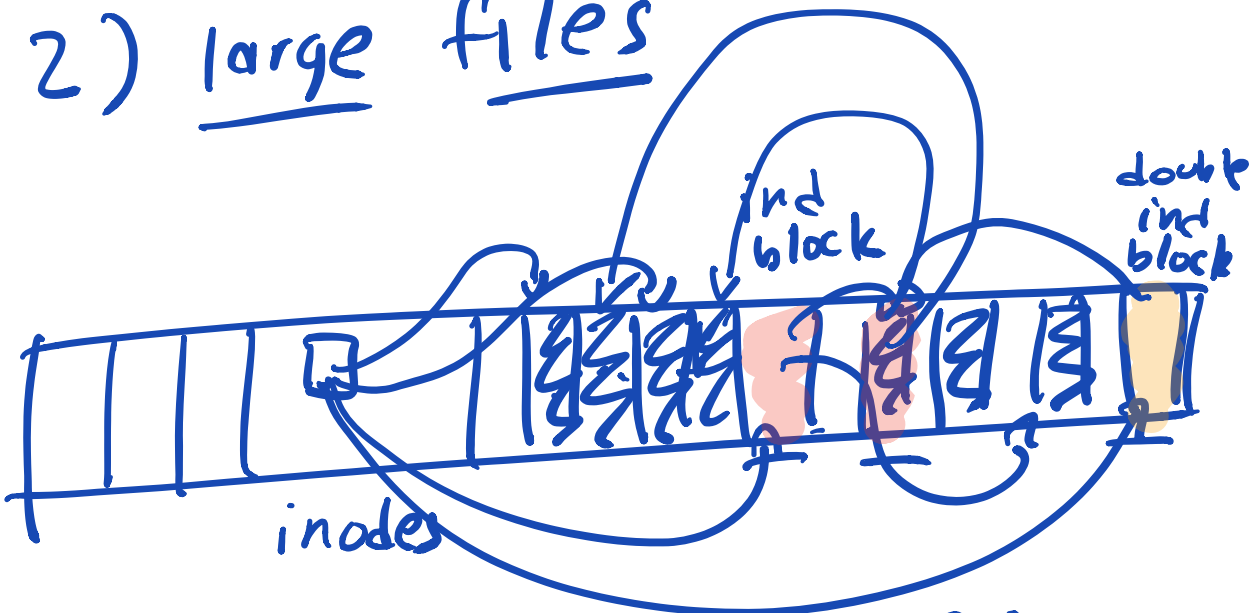1) open ("/a/b/c/d/e.txt"

=>slow                                    O_RDONLY)

many reads

cache  inodes,  (meta data)
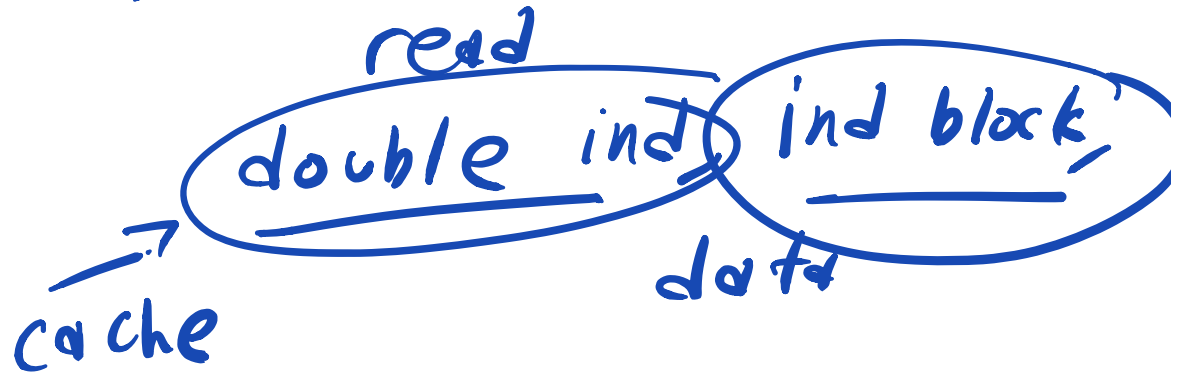
dir data

[in  memory]

re-
open file again => fast


2) large files



ind block

double ind block

inodes

read from large file

=) mighT ~~~ 10
read
(double ind) (ind block,)
→⌐
cache
data

3) cache __data__ too

__writing:__

  open ("/foo", O_WRONLY)
  write ( . . . )
    ⋮
  write ( . . . . )
  close

FS has done __no__ I/O
      to disk

$\Rightarrow$ buffering data in
memory

Later, in background,
OS writes out
data $\rightarrow$ disk

why useful?

$\Rightarrow$ Seems faster (decreases perceived latency)

$\Rightarrow$ leave it in memory : might be read later

$\Rightarrow$ aggregate many "small" writes $\Rightarrow$ 1 block write
(batching)

$\Rightarrow$ overwrite same block repeatedly

(reduces I/O)

=> disk scheduling
   (better w/
        more writes!)

=> open
   write
   ⋮
   write
   close
   ⋮
   unlink
   (delete)

if file
   deleted,
   never write
     to disk
   at all!

why bad?
   => Crash => Data
                Loss!

if you really care
about immediate writes:
open ( )
 write ( )
  :
 write ( )
 fsync ( ) => forces
                data to
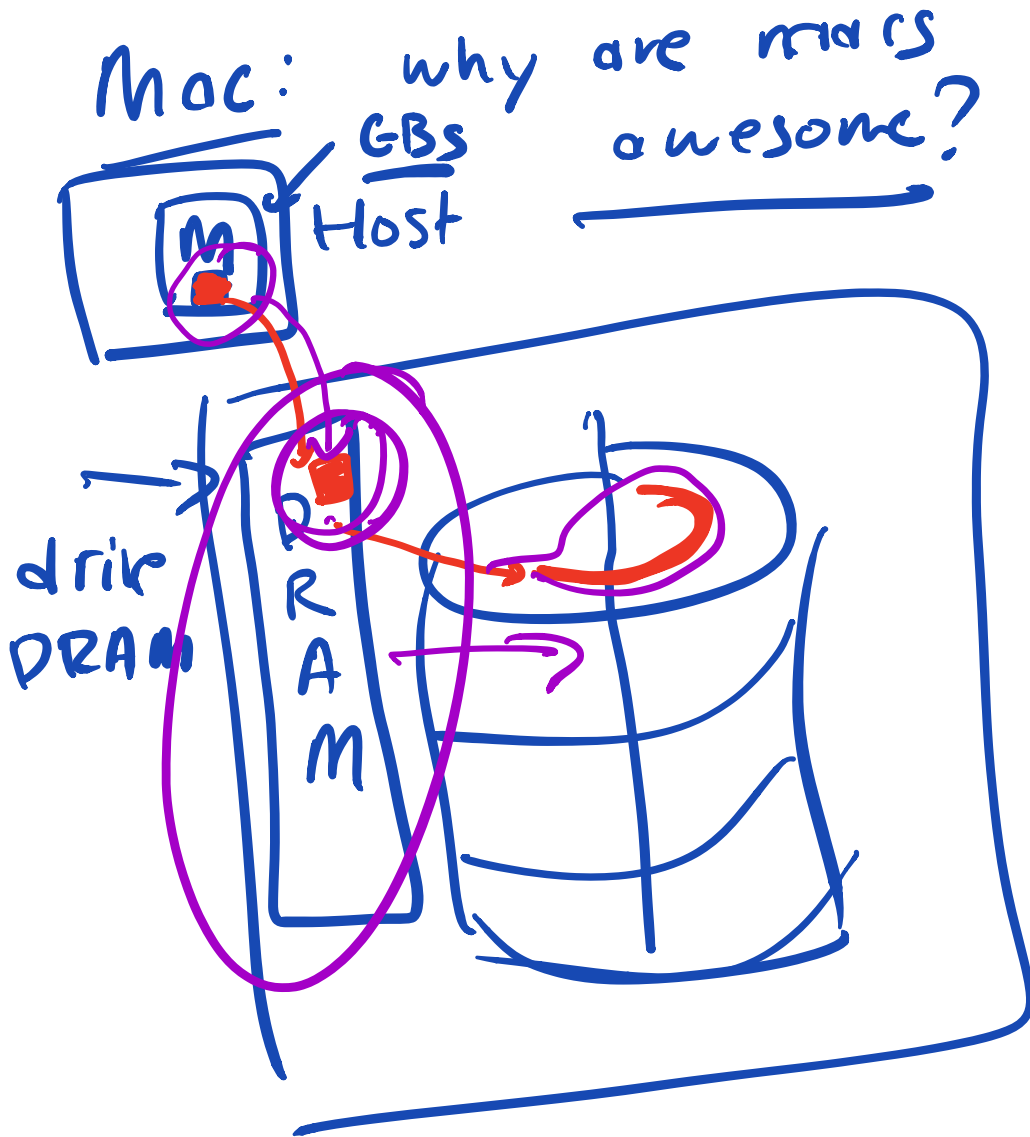                disk
              ( SLOW )
              but safe

P₁                    P₂

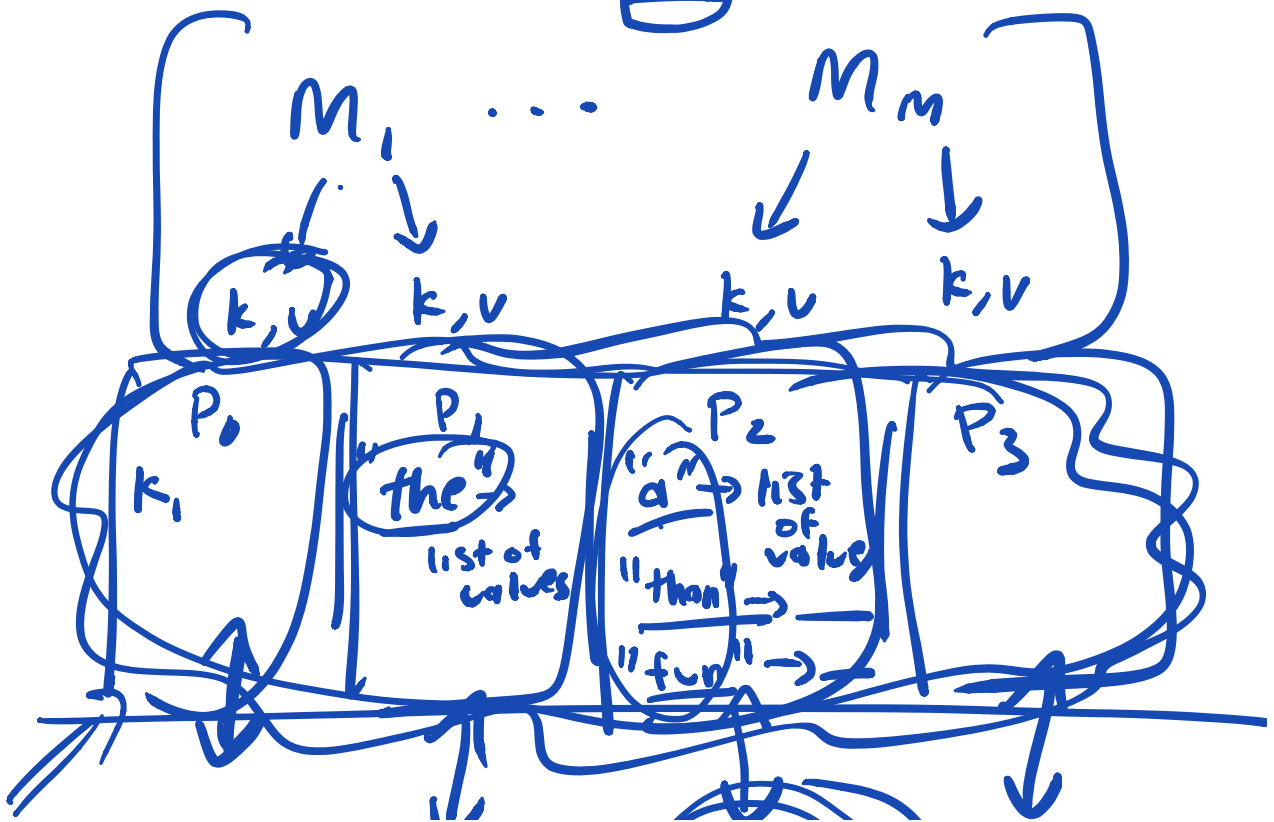W



→R
←R

FS cache:
( Shared )

across all
processes

W

W

Mac: why are macs awesome?

GBs
M   Host

drive
DRAM

R
A
M

Discussion :

Map Reduce Q/A

=> Project reason:
    <u>citations</u>

=> Structure of map-
      reduce problems
          more
      =) examples

=) <u>Perf</u> : [~<u>25%</u>]

=) <u>Sort</u> : 🗒

$M_1$  ...      $M_m$

$k,v$      $k,v$      $k,v$      $k,v$

$P_0$       $P_1$       $P_2$       $P_3$

$k_1$       (the)      (a) → list
            list of         of
            values     values

            "than" → ___

            "fun" →

$R_0$    $R_1$    $R_2$    $R_3$

$\downarrow$

≡        ≡

for all keys k:
    reduce(k, list of
          values
          of k)

alphabetical
order

wait for all reducers

→ clean up

→ return