# Review (537) 5/7

=> Open

General      Old Exam Questions

Suggestions: => focus on now

Read+ => what??
Do => How?

S '16, Q. 7
F '11, 3, Q. 4,6,5
S '18, All !

1/3 RAID

3/5 FFS: Large file + chunk (math) sizes

3/5 SSDs + FTLs + Hybrid Maps

m + F = 50%

1/3 Journaling: why good? why bad?

4 Disk Sched: SATF

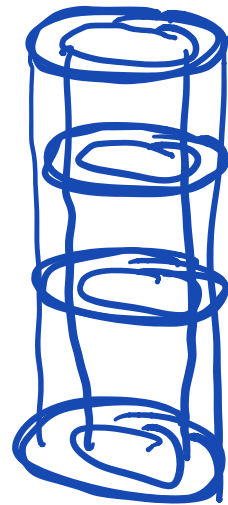**FFS:** Fast File System

=) (Groups) + Locality

esp. on Hard Drives



related

costly seek (+ rotation)

related'

=) Cylinder Groups
(Block)

**FFS:**
place (related) items
into **same** group

(spread out
unrelated item)
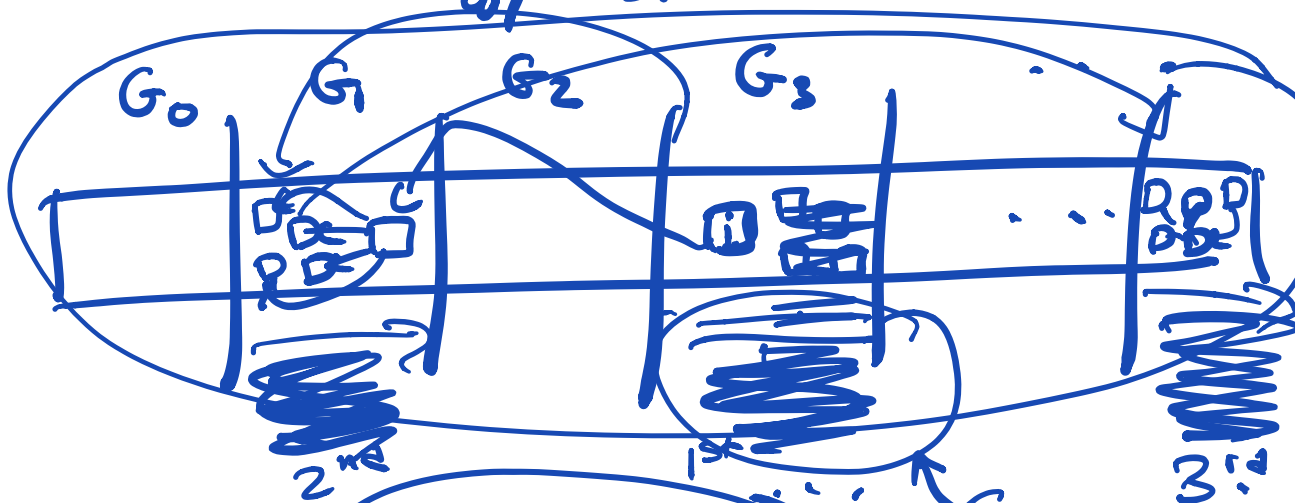
=> related?
 => inode + its data ↗↗ (heuristic)
  => files in same directory
   (make xv6?)

Exception:
(Large Files)

Avoid over-filling a group
         w/ one file



$G_0$    $G_1$    $G_2$    $G_3$

      2nd           1st           3rd

cost : (performance)        (chunk size)
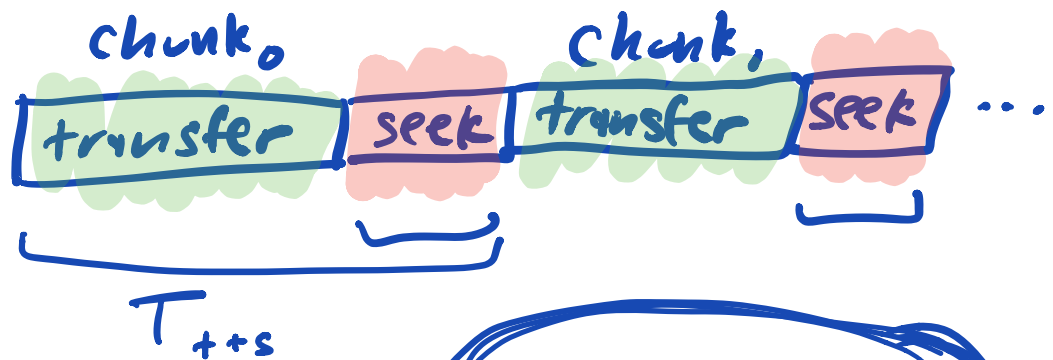  old: one large seq transfer
  new: shorter seq xfer, seek
         (repeat)

$\Rightarrow$ how to get good perf?

A: empirically answer

how to determine
chunk size?

chunk$_0$                              chunk$_1$

| transfer | seek | transfer | seek | ...

$T_{t+s}$

make high: $\boxed{\dfrac{T_{transfer}}{T_{transfer} + T_{seek}}}$

$\Rightarrow .9$
(90%)

Disk:
$\Rightarrow$ Transfer : 100 MB/s
$\Rightarrow$ Avg seek : (20 ms)

Compute    Chunk size ↰ $_{200\,ms}$

$\left( 18\ MB \right) \Rightarrow$ | transfer | seek | $\dfrac{180}{200} =$
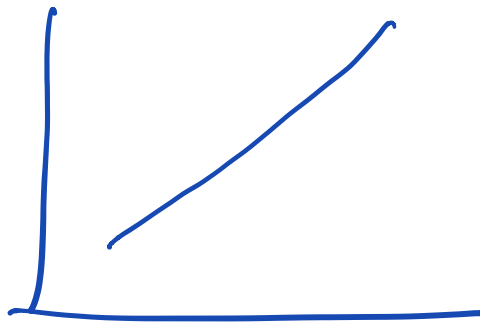
180 ms          20 ms     .9

$$180 \text{ ms} \left( \frac{100 \text{ MB}}{\text{sec}} \right) \left( \frac{1 \text{ sec}}{1000 \text{ ms}} \right) =$$

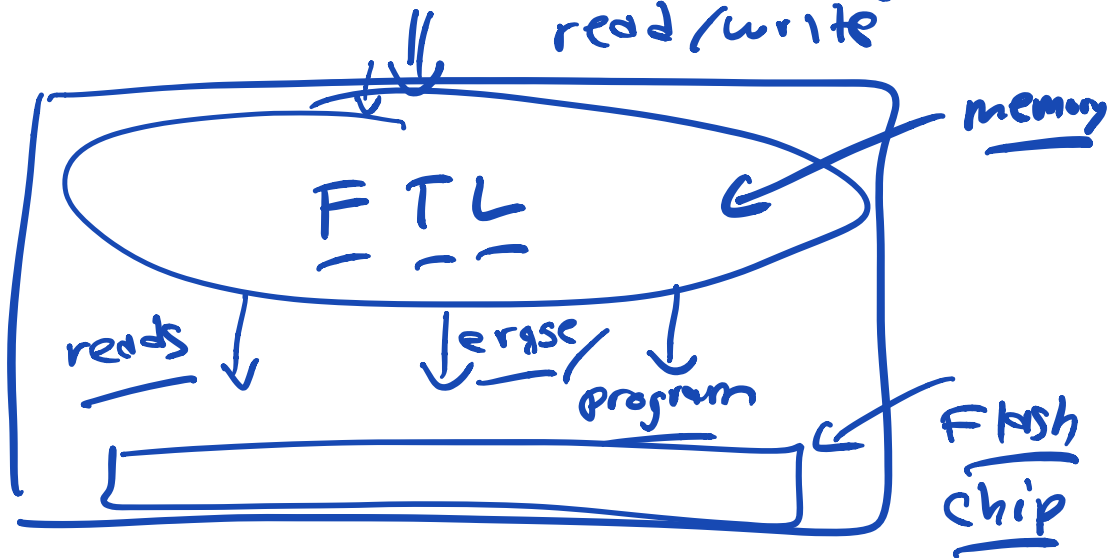$$\frac{180 \text{ MB}}{10} = \boxed{18 \text{ MB}}$$

## SSDs:

pages (YKB)

Blocks (512KB)
(Erase Blocks)

writing: ⎧ => Erase Block
⎪ (makes all pages
⎨ valid)
⎪
⎩ => Program each page once

worries: too many erases => wear
                                    out
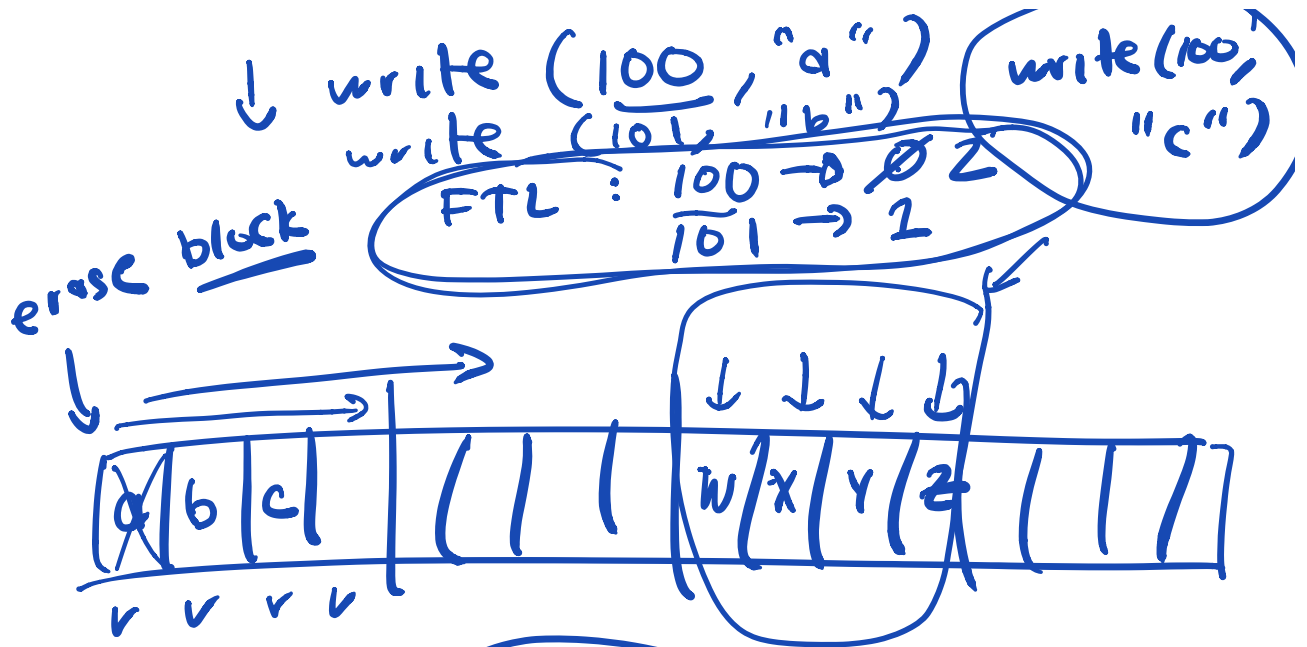same "block" interface
read/write



FTL

memory

reads        erase/program

Flash chip

Flash Translation Layer:

=> Simplest "working" design
Logging of each logical block

+ mapping (logical blocks
           => physical pages)

↓ write (100, "a")
write (101, "b")

write (100, "c")

erase block

FTL : 100 → 0̸ 2
101 → 1

↓ ↓ ↓ ↓

↓

→

→

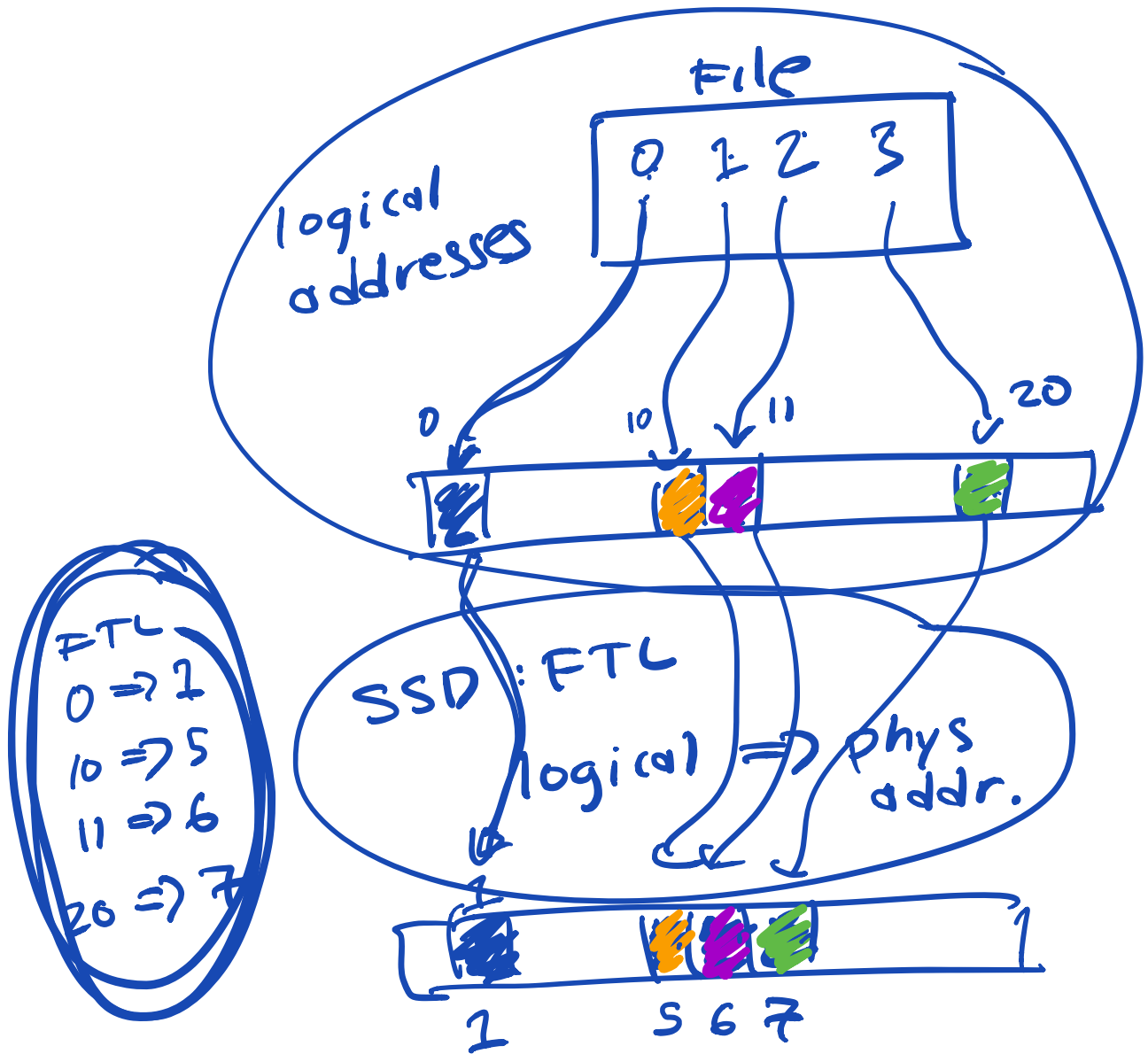| a̶ | b | c | | | | | w | x | y | z | | | | |

✓ ✓ ✓ ✓

concerns ( logging )

=> garbage collection
(background)

=> wear leveling:
move live blocks
elsewhere to
balance erase load
(wear)

Problem: FTL => too big
(one mapping / logical block)

( too many mappings! )

File : inode (+indirect blocks)
       addresses (locations of
                  data blocks)

File

| 0 | 1 | 2 | 3 |

logical addresses

0          10    11         20

SSD : FTL

logical ⇒ phys addr.

FTL
0 ⇒ 1
10 ⇒ 5
11 ⇒ 6
20 ⇒ 7

1          5 6 7

Reduce: FTL size

=> FTL memory : (cache of translations)

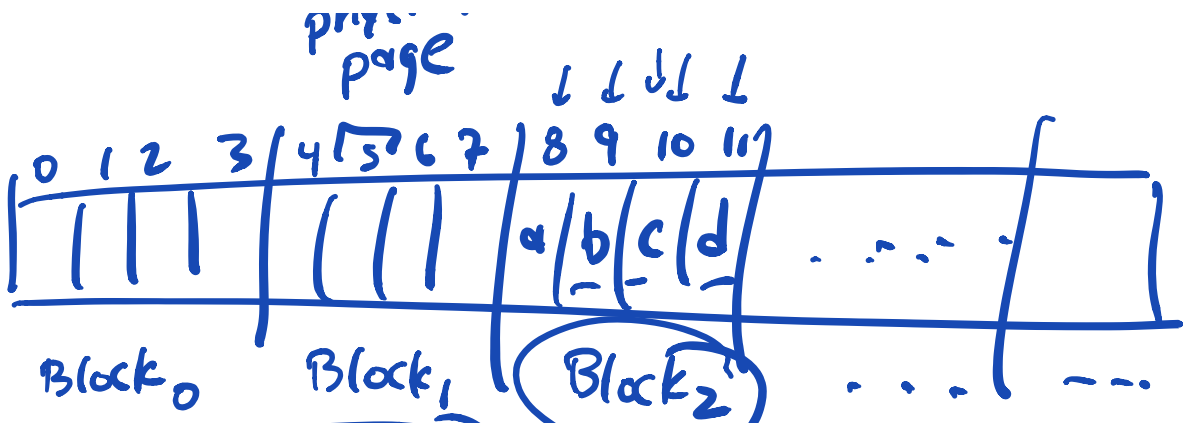=> only keep "most popular" translations in memory of device

=> cost : performance →

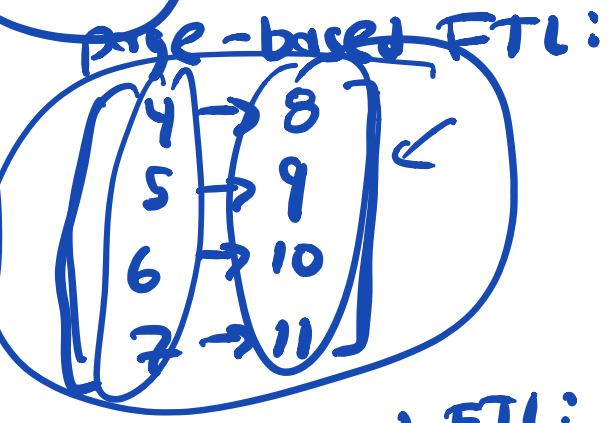e.g. random reads (to more blocks than there are translations that fit in memory)

=> miss per access
=> 2x slower reads

=> More sophisticated mappings

=> Block-based

physical

phys
page

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | | | ... ... | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

↓ ↓ ↓↓ ↓

a | b | c | d  . . . . . .

Block₀        Block₁        Block₂     . . . .     . . .

write (4; a)
write (5, b)
write (6, c)
write (7, d)

page-based FTL:

4 → 8
5 → 9
6 → 10
7 → 11

block-based FTL:
one-mapping
per block

4, 5, ..., → phys block

2

Hybrid FTL : page mappings +
block mappings

write (4; a);
write (5, b);
write (6, c) ];

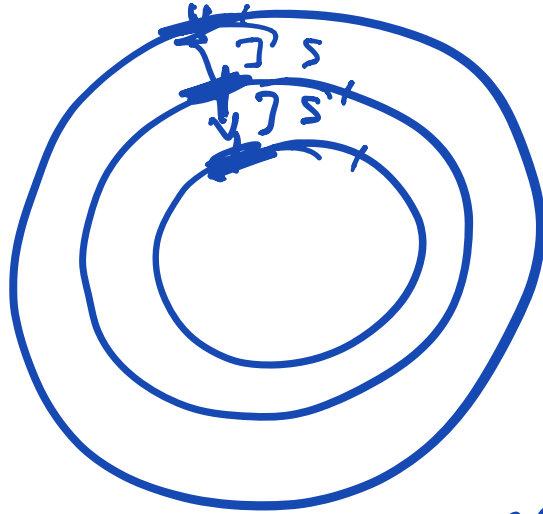$$\begin{pmatrix} \text{write } (6, \underline{\quad}) \\ \text{write } (7, d) \end{pmatrix}$$

FTL : use page mappings
at first

$$\begin{bmatrix} 4 \rightarrow 100 \\ 5 \rightarrow 102 \\ 6 \rightarrow 102 \\ 7 \rightarrow 103 \end{bmatrix}$$

Convert
to
smaller
block-based
mapping

$$[1 \rightarrow 25]$$

$$[ \quad A \quad | \quad AB \uparrow \quad B \uparrow \quad \quad ]$$

⟶

=) Relation of S to R

=) Same request

SATF:

$$(2S + 3R)$$

old:

$$\frac{D \text{ time units}}{3D}$$

new:

$$\frac{\overset{Avg}{Seek} + \overset{Avg}{Rotate}}{S + \frac{R}{2}}$$
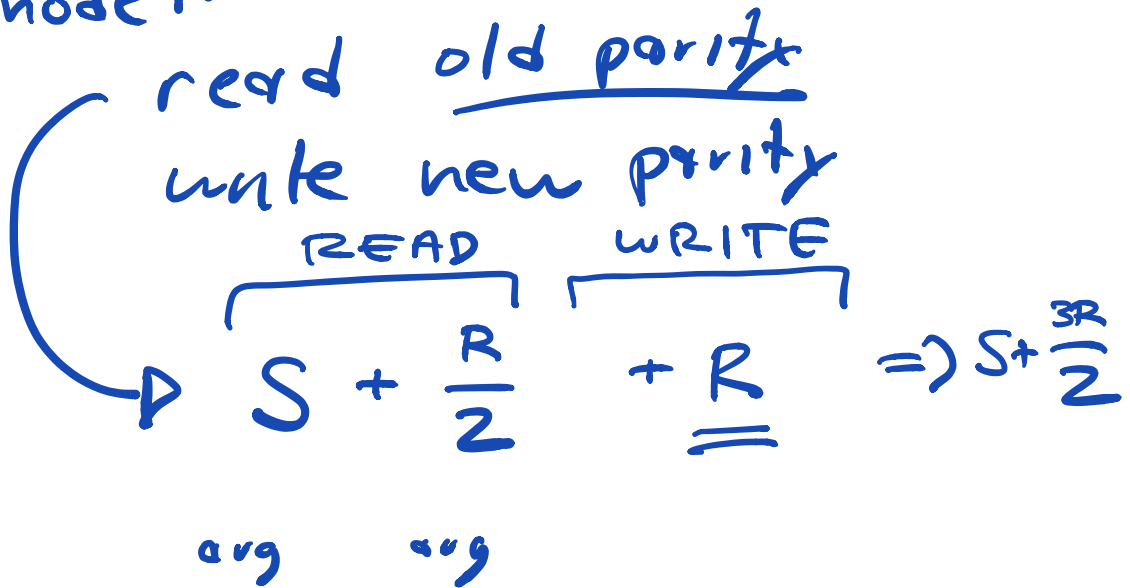
RAID - Y:

$(0,1,2)$

[ 0 ] [ 1 ] [ 2 ] [ P ]

[ Read old data, read old parity ]

→ compute new parity

[ write new data, write new parity ]

model:

( read old parity
  write new parity

READ: $S + \dfrac{R}{2}$   WRITE: $+R$   $\Rightarrow S + \dfrac{3R}{2}$

avg    avg

RAID-5 : 12 writes

0  1  [ 2  P ]    logical 2 writes

4   5   P   3

8   P   6   7

P   9   10   11

logical write:

D $\begin{cases} \text{read} \\ \text{write} \end{cases}$  $\left.\begin{array}{l} \text{read} \\ \text{write} \end{array}\right]$ in parallel / in parallel

$\Rightarrow$ logical write $\Rightarrow$ 2 D
                                    $\times$
                                    6

$\Rightarrow$ 12 D

Append:   write ( data )

ID=1
1000,20
60

data   inode   Bitmap           ID=1

$T_B$   D   I   B   $T_E$

journal

(WAL)

1) write everything to journal
$\Rightarrow T_B$, contents all atonce
wait
$\Rightarrow T_E$ : commits transaction

wait

2) update (overwrite $B, D, I$
in place)