LECTURE

Flash

# (NAND) Flash

(refs 7, 16, 17)

→ USCD paper (Micro 2009)
→ DFTL (ASPLOS '09)

**data:** charge trapped in a gate
each gate can store 1/more bits
(SLC) (MLC)

**internally:**

plane/bank

flash chip

**bank:** divided into <u>blocks</u> (size 128k → 256k)

~~block~~ divided into pages (size ~2K → 8K)
(data + OOB)

**three primary operations:**

erase (block) → sets all bits to 1

program (page) → can change (some) 1's → 0's

⟹ must <u>erase</u> block before <u>programming</u> page w/in block
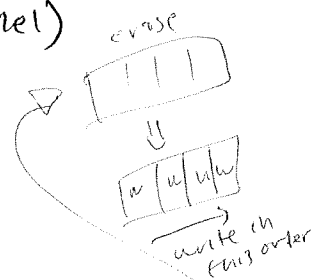
read (page) → read whole page (in parallel)

**performance:** (per page)

read latency ~10 us (10-50)    ~200 MB/s

prog latency ~200 → 500 us    ~

erase        ~2 <u>ms</u>

→ <u>no moving parts</u> !
(rare catastrophic failure)
→ relatively <u>high bit-error rates</u>
(disturb)
→ can't tell difference between 1/0

**reliability:**

primary problem: <u>wearout</u>

erase/program cycle → trapped charge →

ugly
→ partial-page prog
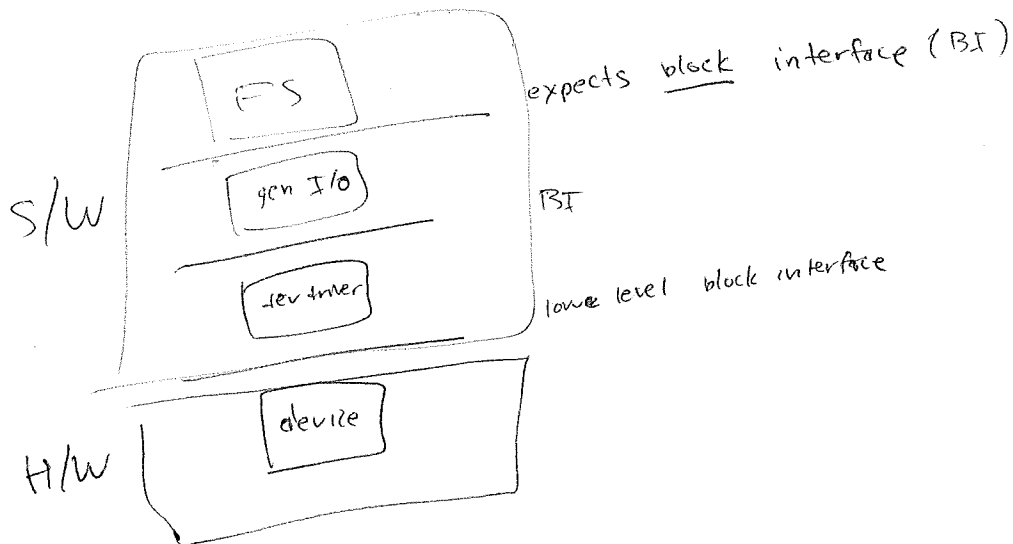→ data retention @ high temp

for each <u>block</u>: 10K P/E cycles (MLC)
100K P/E cycles (SLC)

also, program disturb ⟹ when prog one page, neighbor may be affected
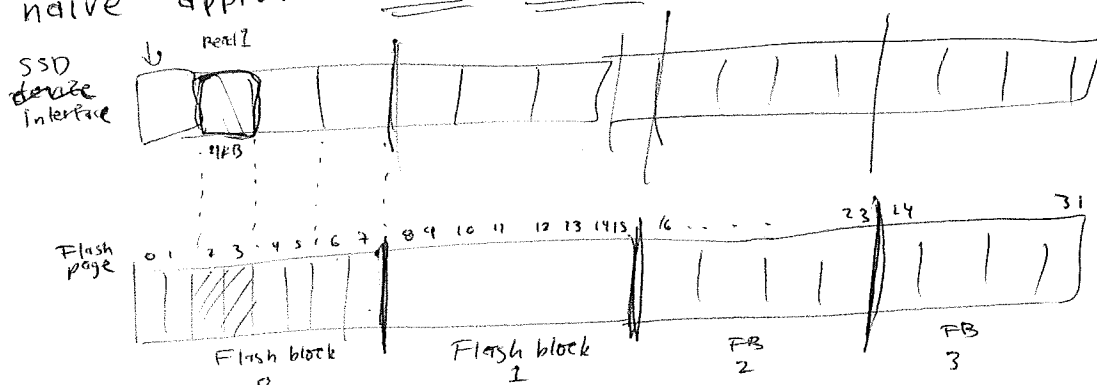⟹ program in <u>order</u> (once written, only one strong disturb will happen)

how to use in system?

S/W

- FS — expects block interface (BI)
- gen I/O — BI
- dev driver — lower level block interface

H/W
- device

---

naive approach: direct mapping

SSD
device
interface

read1

4KB

Flash
page

0 1 2 3 4 5 6 7  8 9 10 11  12 13 14 15  16 . . . . . .  23 24  . . . 31

Flash block 0     Flash block 1     FB 2     FB 3

read $d_1$ => read (2,3) ←p => works!

write $d_1$ => erase (b0), program (2,3)

bad why? => performance (each erase ~ 2 ms)
=> reliability (wear out!)
(prog. disturb)

---

more sophisticated/real: flash translation layer (FTL)

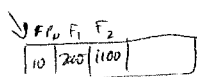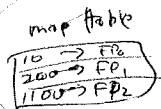SSD: looks like disk   takes read/write blocks

FTL | state: mapping table

⇓

e, p, r of blocks/page

( simple page mapping )

write   10
write   200
write   1100

⇓

map table
10 → FP₀
200 → FP₁
1100 → FP₂

⇓ FP₀ FP₁ FP₂
10 | 200 | 1100

Implications of FTL: (page-mapped)
   → can erase blocks in background          performance ⇒ good! ⎫
   → fill in log-structured style ⇒ helps w/ wear                    ⎬ perf
                                        leveling (reliability)        ⎭

⇒ generates garbage                                                    reliability
   → device must find old blocks, clean, make available again
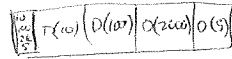   → device must also make sure to __wear level__ blocks   erase/program

device
tracks:
   live/dead : how?      ⇒      [OOB area]  | T(10) | D(100) | X(200) | O(5) |
   erasure count / block : how?  ⇒
   ⎵⎵⎵⎵⎵⎵⎵⎵⎵⎵                        examine block
   (must periodically                check map table
    clear block for
    wear too!)

Problem: mapping table is too big! (costly, energy)
   compute size : 260 GB device w/ 2K pages
                  100 m entries * 4 MB ⇒ [ 400 MB ]
                                          in device!

          solutions:
            zoning (low-end) ⇒   swap table per zone in/out mem
                                  (keep most tables in flash)

Block size=4
Page Num
0 1 2 3 | 4 5 6 7      (larger mappings ⇒   makes table smaller:
BN=0    BN=1              (not used)          @ what cost?                 erase save new block,
BN= PN / block size                           [ program of small page ⇒ read block, write block
Boff= PN % block                                                          (+ new page)
       size                                     (RMW cycle) ⇒ costly

          common: __hybrid__  (both page-level +
                                      block-level)
            most of SSD: Data blocks  (block mapped)  ≥ 95% %
            small, active
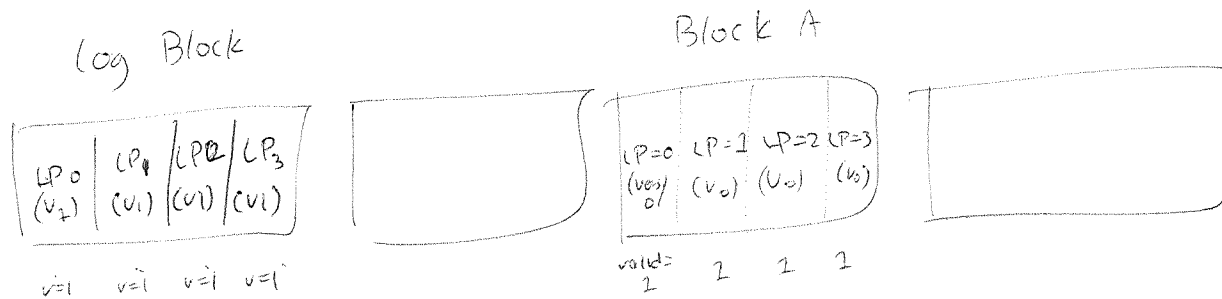               write area : Log blocks  (page mapped)  < 5%

          tricky part: __GC__
            switch merge: after seq write ⇒ log,
                          switch log + dest block

# switch merge

### log Block

| LP 0 (v7) | LP1 (v1) | LP2 (v1) | LP3 (v1) |
|---|---|---|---|
| v=1 | v=1 | v=1 | v=1 |

### Block A

| LP=0 (ver) 0 | LP=1 (v0) | LP=2 (v0) | LP=3 (v0) |
|---|---|---|---|
| valid: 2 | 2 | 2 | 2 |

make A => log block (erase it)

make L => data block (update mapping)

# partial merge

### Log block

| LP=0 | LP=2 | . . . |
|---|---|---|
| v=1 | v=2 | |

### Block A

| LP=0 | v=1 | LP=2 | LP=3 |
|---|---|---|---|
| v=0 | v=0 | vano | |

read 2,3 from A

write to L (prog)

call L → A
call A → L

# full merge

### L

| LP=4 | LP=0 | MMM |
|---|---|---|

### Block A

| 4 | 5 | 6 | 7 |
|---|---|---|---|

### Block B

| 0 | 1 | 2 | 3 |
|---|---|---|---|

invalidated

=> find new erased C

=> merge 0(L) + (1,2,3)B => C
(erase B) => L

=> merge (5,6,7)A => L
(erase A)

end result:

SSDs
     seq r/w => good (100s of MB/s)
     rand read => good (nearly seq)
     rand write => bad (lots of GC + merging)

( FS implications? )