```
1   CS5600, Handout Week12.b
2
3   // 1. Read/write disk in Lab4 (CS5600 file system)
4   // borrowed from Lab4, fs5600.c
5
6       /* disk access.
7        * All access is in terms of 4KB blocks; read and
8        * write functions return 0 (success) or –EIO.
9        *
10       * read/write "nblks" blocks of data
11       *   starting from block id "lba"
12       *   to/from memory "buf".
13       *     (see implementations in misc.c)
14       */
15      extern int block_read(void *buf, int lba, int nblks);
16      extern int block_write(void *buf, int lba, int nblks);
17
18
19      /*
20       * below are two toy examples of
21       * reading from and writing to a disk block
22       */
23
24      // Reading one block
25
26      char buf[FS_BLOCK_SIZE];   // FS_BLOCK_SIZE=4096; see panel 2
27      int inum = 100;            // block number to read from
28      int ret = block_read(&buf, inum, 1);
29      if (ret < 0) {  // error; ret should be –EIO
30          return ret;
31      }
32
33
34      // Writing one block
35
36      char buf[FS_BLOCK_SIZE];  // again, 4KB buffer
37      ...                       // update buf
38
39      int ret = block_write(&buf, 99, 1);  // update the 99th block
40      if (ret < 0) {  // error; ret should be –EIO
41          return ret;
42      }
43
44
```

```
45
46
47   // 2. CS5600 file system data structures
48   // borrowed from fs5600.h with minor changes
49
50       /*
51        * file:        fs5600.h
52        * description: Data structures for CS5600 file system.
53        *
54        * CS 5600, Computer Systems, Northeastern CCIS
55        * Peter Desnoyers, November 2016
56        *
57        * Modified by CS5600 staff in fall 2021.
58        */
59
60       #define FS_BLOCK_SIZE 4096
61       #define FS_MAGIC 0x30303635
62
63       #define INODE_NUM_PTRS (FS_BLOCK_SIZE/4 – 5)
64
65       /* Superblock – holds file system parameters.
66        */
67       struct fs_super {
68           uint32_t magic;
69           uint32_t disk_size;          /* in blocks */
70
71           /* pad out to an entire block */
72           char pad[FS_BLOCK_SIZE – 2 * sizeof(uint32_t)];
73       };
74
75
76       /* inode = 4096 bytes */
77       struct fs_inode {
78           uint16_t uid;
79           uint16_t gid;
80           uint32_t mode;
81           uint32_t ctime;
82           uint32_t mtime;
83           int32_t  size;
84           uint32_t ptrs[INODE_NUM_PTRS];
85       };
86
87
88       /* Entry in a directory
89        */
90       struct fs_dirent {
91           uint32_t valid : 1;
92           uint32_t inode : 31;
93           char name[28];        /* with trailing NUL */
94       };
```

```
1   CS5600, Handout Week 12.1
2
3   //   Code snippets borrowed from WeensyOS boot loader
4   //   They illustrate how kernel "talks" to a disk through programmed I/O:
5   //   the bootloader reads in the kernel from the disk.
6   //
7   //   See the functions boot_waitdisk() and boot_readsect(). Compare to
8   //   Figures 36.5 and 36.6 in OSTEP.
9
10
11  //   WeensyOS boot loader loads the kernel at address 0x40000 from
12  //   the first IDE hard disk.
13  //
14  //   A BOOT LOADER is a tiny program that loads an operating system into
15  //   memory. It has to be tiny because it can contain no more than 510 bytes
16  //   of instructions: it is stored in the disk's first 512-byte sector.
17
18
19      #define SECTORSIZE      512
20
21      // boot_readsect(dst, src_sect)
22      //    Read disk sector number `src_sect` into address `dst`.
23      static void boot_readsect(uintptr_t dst, uint32_t src_sect) {
24          // programmed I/O for "read sector"
25          boot_waitdisk();
26          outb(0x1F2, 1);            // send `count = 1` as an ATA argument
27          outb(0x1F3, src_sect);     // send `src_sect`, the sector number
28          outb(0x1F4, src_sect >> 8);
29          outb(0x1F5, src_sect >> 16);
30          outb(0x1F6, (src_sect >> 24) | 0xE0);
31          outb(0x1F7, 0x20);         // send the command: 0x20 = read sectors
32
33          // then move the data into memory
34          boot_waitdisk();
35          insl(0x1F0, (void*) dst, SECTORSIZE/4); // read 128 words from the disk
36      }
37
38
39      // boot_waitdisk
40      //    Wait for the disk to be ready.
41      static void boot_waitdisk(void) {
42          // Wait until the ATA status register says ready (0x40 is on)
43          // & not busy (0x80 is off)
44          while ((inb(0x1F7) & 0xC0) != 0x40)
45              /* do nothing */
46          }
47      }
```

---

Handwritten annotations:

1 TB ≠ 1 TB
↓ (disk)    (mem)
1 Billion Bytes < $2^{40}$ B

1MB = 1 million Bytes
(disk)

VA    LBA
32 bits

32

× 4B = 512B

I/O polling

interface of disk.

OS → disk
· platter
· track
· sector

LBA

0 1 2

Sectors → disk

· Capacity: 100GB ~ 10TB
· #Platters: 8~10
· #tracks: 10s ~ 1000s
· #Sectors/track: ~1000
· RPM: ~10 000
· transfer: 50 ~ 150 MB/s
· mean time between failure

Q: 8TB 512B.  #bits for LBA?

$$\frac{8TB}{512B} = \frac{2^3 \cdot 2^{40}}{2^9} = 2^{34}$$

48 bit → $2^{48}$ Byte → 128 PB (max)
34 bit
1 million TB

# Hard disk drive organization

Sector

Platter 1

Platter 2

Track X

Platter 3

Platter 4

Cylinder X

Direction of Rotation

Track

Sector

Read/Write
Head

Fig 5.11, 5.12 from
http://www.ccs.neu.edu/~pjd/x600-book-v0903.pdf

# A Flash-based SSD (logical diagram)

(128KB, 256KB)
blocks
↓
pages.
(4KB.
?
16 KB)

bank

Interface Logic

Memory

Flash
Controller

Flash  Flash  Flash

Flash  Flash  Flash

FTL

Fig 44.3 from
https://pages.cs.wisc.edu/~remzi/OSTEP/file-ssd.pdf

1. Last time ↙
2. SSD
3. Intro to file systems ↙
4. Files
5. Directories

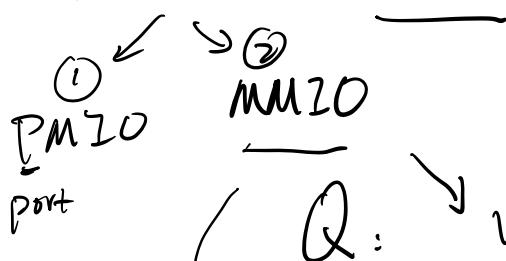------------------------------------------------

I/O.

| OS | ←——→ | device |

Sync vs. async.

User-level threading

{ polling, interrupt }

×

{ PIO, DMA }

① PMIO    ② MMIO

PMIO
port

Q:  ↓ vs. ↓ ?

Py Mem Address space.

| DRAM |

1 MB — | | —— | | devices
        | |

DRAM
buffer ← device

disk

SSD.  ——→ banks ——→ blocks ——→ pages.

| Operations. | granularity. | perform. |
|---|---|---|
| • read. | 1 page | ~10s MS |
| • erase | 1 block., resttig all bits ⇒ **1** | ~1 ms |
| • program. | 1 page., setting some bits ⇒ 0 | ~100s MS |

( Cannot __program__ the same page twice without erase )

Q: page A. update A. HOW?



block   copy   mem

↑
erase X 10,000   ___wear-out___

```
for (i=0 ; i<10000 ; i++){
    update A.
}
```

• FTL

Logical page $\xrightarrow{\text{FTL}}$ physical SSD
SSD                                    page

• log-structure FTL
write ( LP X ) → append., update FTL
read ( l Dx ) ← check mapping return

PX.

```
** an example:
   --Given a flash bank has three blocks; each has two pages.
   --there are four writes to pages:
     wirte(logic_page_1)    [short as LP1]
     wirte(logic_page_10)   [short as LP10]
     wirte(logic_page_99)   [short as LP99]

   --what will happen:                              GC.

           +----------------------------+
   blocks  | block 0 | block 1 | block 2 |
           +---------+---------+---------+
   pages   | P1 | P2 | P3 | P4 | P5 | P6 |
           +----+----+----+----+----+----+
   data    |LP1'''|LP10|LP99|    |    |    |
           +----+----+----+----+----+----+

   mapping:
     LP1 => P1, LP10 => P2, LP99 => P3

Question:
   what will happen if the following op is write(logic_page_1¹)?
```

$> 85\%$

$< 50\%$

write( LP1" )

• FS.

① persistent storage.          ROM

② give named bytes. (file)

③ friendly way to find    ( directories

⬆                     named bytes