

Layered Block-Structured File System & Caching

Yu-Ju Huang

Robbert van Renesse

CS4411/5411 Project overview

P0	Queue
P1	Hello World
P2	Cooperative threads
P3	Preemptive threads & MLFQ
P4	Syscall & Memory protection
P5	Disk cache
P6	File system

File System

- Mapping a file name to its content and metadata
 - A file may be stored in HDD, SSD, or RAMDISK
 - Metadata like permission, owner, etc

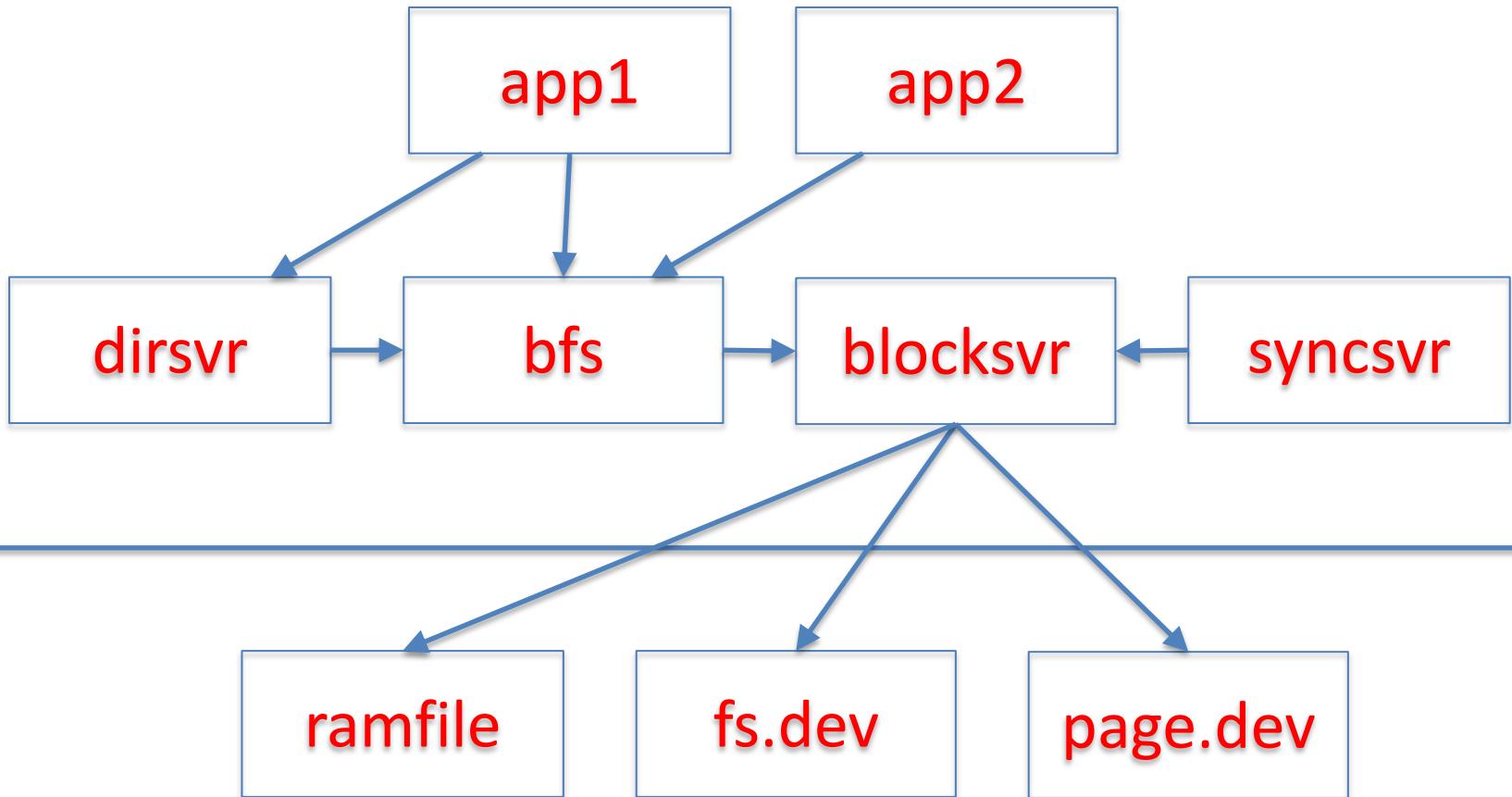
Block store

- Underneath any file system, database system, etc. there are one or more *block stores*
- A block store provides a disk-like interface:
 - a storage object is a sequence of blocks
 - typically, a few kilobytes
 - you can read or write a block at a time
- The block store abstraction doesn't deal with file naming, security, etc., just storage

Inode

- A data structure that stores metadata about a **file**.
- Each file on a filesystem has a unique inode number within that filesystem.
- An Inode contains
 - **File type** (regular file, directory, symbolic link, etc.)
 - **Permissions** (read, write, execute)
 - **Owner (UID) and Group (GID)**
 - **File size**
 - **Timestamps**

EGOS Storage Architecture



bfs: block file server

- Implements a full-featured file system
- Maintains for each file a “stat structure”:
 - size in bytes
 - owner
 - modification time
 - access control information
 - etc.
- files are indexed by i-node numbers

blocksvr: block server

- Manage block device (or block store)
 - HDD, SSD, RAMDISK, etc
- Given (inode, offset), blocksvr read/write the content from/to the underlying block device

BLOCK STORE

Block Store Abstraction

- A block store consists of a collection of *i-nodes*
- Simple interface:
 - `getninode()` → integer
 - returns the number of i-nodes on this block store
 - `getsize(inode number)` → integer
 - returns the number of blocks on the given inode
 - `setsize(inode number, nblocks)`
 - set the number of blocks on the given inode

Block Store Abstraction, cont'd

- `read(inode, block number) → block`
 - returns the contents of the given block number
- `write(inode, block number, block)`
 - writes the block contents at the given block number
- `sync(inode)`
 - make sure all blocks are persistent
 - if `inode == -1`, then all blocks on all inodes

Simple block stores

- “filedisk”: a simulated disk stored on a Posix file
 - `block_if bif = filedisk_init(char *filename, int nblocks)`
- “ramdisk”: a simulated disk in memory
 - `block_if bif = ramdisk_init(block_t *blocks, nblocks)`
 - Fast but volatile
- `block_if` is a pointer to the block interface

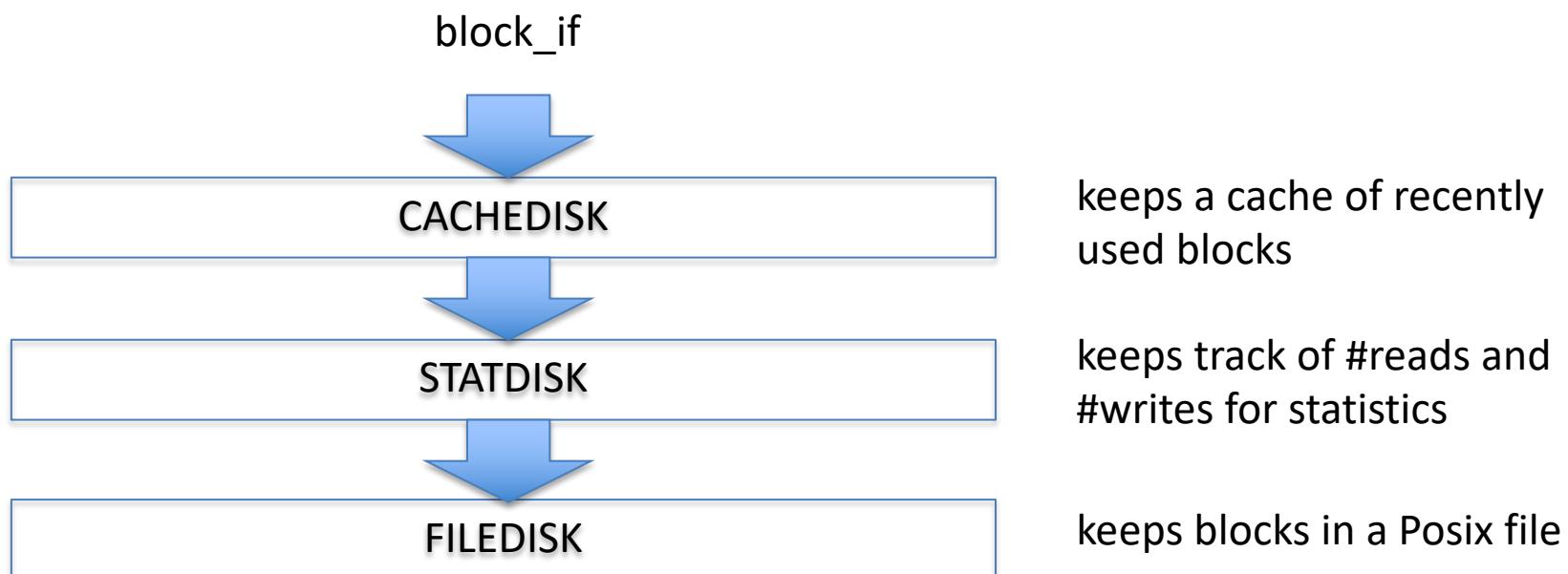
Example code

```
#include ...
#include "egos/block_store.h"

int main(){
    block_if disk = filedisk_init("disk.dev", 1024);
    block_t block;
    strcpy(block.bytes, "Hello World");
    (*disk->write)(disk, 0, 0, &block);
    (*disk->release)(disk);
    return 0;
}
```

Block Stores can be Layered!

Each layer presents a `block_if` abstraction



Example code with layers

```
#define CACHE_SIZE 10          // #blocks in cache

block_t cache[CACHE_SIZE];

int main(){
    block_if disk = filedisk_init("disk.dev", 1024);
    block_if sdisk = statdisk_init(disk);
    block_if cdisk = cachedisk_init(sdisk, cache, CACHE_SIZE);

    block_t block;
    strcpy(block.bytes, "Hello World");
    (*cdisk->write)(cdisk, 0, 0, &block);
    (*cdisk->release)(cdisk);
    (*sdisk->release)(sdisk);
    (*disk->release)(disk);

    return 0;
}
```

Example Layers

```
block_if clockdisk_init(block_if below,
                      block_t *blocks, block_no nblocks);
// implements CLOCK cache allocation / eviction

block_if statdisk_init(block_if below);
// counts all reads and writes

block_if debugdisk_init(block_if below, char *descr);
// prints all reads and writes

block_if checkdisk_init(block_if below);
// checks that what's read is what was written
```

How to write a layer

```
struct statdisk_state {  
    block_if below;           // block store below  
    unsigned int nread, nwrite; // stats  
};  
  
block_if statdisk_init(block_if below){  
    struct statdisk_state *sds = calloc(1, sizeof(*sds));  
    sds->below = below;  
  
    block_if bi = calloc(1, sizeof(*bi));  
    bi->state = sds;  
    bi->getsize = statdisk_nblocks;  
    bi->setsize = statdisk_setsize;  
    bi->read = statdisk_read;  
    bi->write = statdisk_write;  
    bi->release = statdisk_release;  
    return bi;  
}
```

statdisk implementation, cont'd

```
static int statdisk_read(block_if bi, unsigned int ino, block_no offset,
block_t *block){
    struct statdisk_state *sds = bi->state;
    sds->nread++;
    return (*sds->below->read)(sds->below, ino, offset, block);
}

static int statdisk_write(block_if bi, unsigned int ino, block_no offset,
block_t *block){
    struct statdisk_state *sds = bi->state;
    sds->nwrite++;
    return (*sds->below->write)(sds->below, ino, offset, block);
}

static int statdisk_getsize(block_if bi){ ... }
static int statdisk_setsize(block_if bi, block_no nbblocks){ ... }

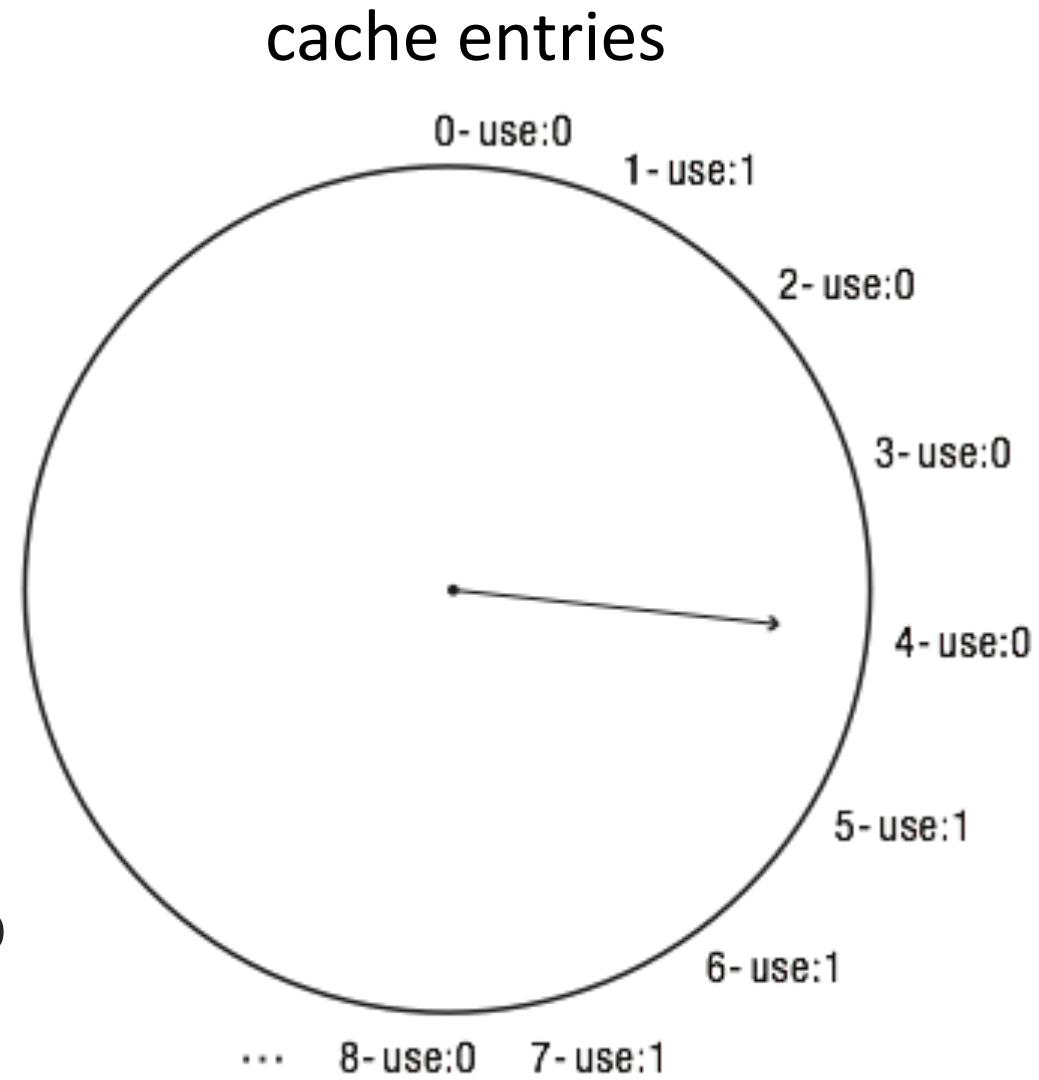
static void statdisk_release(block_if bi){
    free(bi->state);
    free(bi);
}
```

P5: Implement a cache layer

- Based on clock algorithm
- Two versions:
 1. write-through
 2. write-behind *aka* write-back
- Tricky part: what to do if cache is full?

Clock Algorithm

- To evict a block, inspect the *use* bit in the cache entry at clock hand and advance clock hand
- Used? Clear *use* bit and repeat
- Not used -> victim to be evicted



Implement a CLOCK cache

What should be in a cache data structure?

Implement a CLOCK cache

What should be in a cache data structure?

- Content
- Use flag
- Valid flag

Implement a CLOCK cache

What should be in a cache data structure?

- Content
- Use flag
- Valid flag

What should be in the DISK structure?

Implement a CLOCK cache

What should be in a cache data structure?

- Content
- Address
- Use flag
- Valid flag

What should be in the DISK structure?

- Clock hand

Today

- Block store
- P4 due next week
- P5 release this weekend, due in 4 weeks