File System: buffer cache

Youjip Won

KAIST EE

Contents

- Overview
- Buffer cache layer
- Code: Buffer cache
- Logging layer
- Log design
- Code: Logging

File System Layout of xv6



- Block 0: Boot sector, what we covered at Chapter 2
- Block 1: Superblock. (metadata about file system)
- Blocks starting at 2
 - Log
 - inodes
 - bitmap
 - data

Superblock in xv6 - Implementation

• A block of metadata describing the file system



(gdb) p sb \$1 = {size = 1000, nblocks = 941, ninodes = 200, nlog = 30, logstart = 2, inodestart = 32, bmapstart = 58}

Log (Write-Ahead-Log)



In-memory representation of the log area



```
struct log {
  struct spinlock lock;
  int start;
  int size;
  int outstanding; // how many FS sys calls are executing.
  int committing; // in commit(), please wait.
  int dev;
  struct logheader lh;
};
```

Inode structure





Bitmap

Bitmap represents the state of the usage of blocks in filesystem



- Each bit in bitmap represents the state of the block allocation.
- 512byte can represent 4096 blocks (512 Bytes = 512 * 8 bit = 4096 bit).

Buffer cache

• Use part of the memory as storage.





Buffer Cache Location at the Memory – Version 3

Rest are mapped for device.



Buffer Cache Entry - Flags

- Flags of buffer head indicates the state of the buffer cache entry
 - B_VALID: Buffer contains a copy of the block
 - B_DIRTY: Buffer content has been modified and needs to be written to the disk
 - If not set, buffer cache entry is unused



Buffer cache layer

- Buffer cache
 - Cache disk block in memory.
 - Ensure that only one copy of a block is in memory.
 - Ensure that only one kernel thread at a time accesses that copy.
- API's
 - bget : Scan the buffer cache and find the matching buffer .
 - bread: Obtain a buffer containing a copy of a block, lock the buffer.
 - bwrite: Write the modified buffer to the appropriate location on the disk.
 - brelse: unlock the buffer.

Buffer



Buffer cache: bcache



```
struct {
  struct spinlock lock;
  struct buf buf[NBUF];
  // Linked list of all buffers, through prev/next.
  // head.next is most recently used.
  struct buf head;
} bcache;
```

Buffer cache mechanism

- The buffer cache uses a per-buffer sleep-lock.
 - Only one thread at a time uses each buffer.
 - bread returns a locked buffer.
 - bget locks the buffer.
 - brelse releases the lock.
- Buffer cache has a fixed number of buffer.
 - If the file system asks for a block that is not in cache, buffer cache recycles a buffer currently holding some other block.
 - LRU scheme for victim selection

Buffer Cache - Initialization

- Buffer cache is already allocated when kernel is loaded (at data section).
- binit() is called at main() and initializes the list.



binit() - initialize the list of buffers





next —	next	→ next	next —	—→ next
buf (head)	buf[NBUF]	buf[NBUF-1]	 buf[1]	buf[0]
prev 🔶	prev ←	prev	prev 🛶	prev

buffer state

- Buffer has two state bits
 - B_VALID: Buffer contains a copy of the block
 - B_DIRTY: Buffer content has been modified and needs to be written to the disk



bread()

- read 'blockno' from the 'dev'.
- bget: scan the buffer cache and find the matching buffer with 'blockno' from the 'dev'.
- If buffer does not exist, read the disk block of blockno' from the 'dev' into buffer pointed by b.





- Find the matching buffer by scanning.
 - [Cache Hit] Return the locked Buffer.
 - [Cache Miss] Return the locked Buffer which is newly allocated.

Scan the array until the matching buffer is found



struct buf * bget(uint dev, uint blockno)

• Find the matching buffer by scanning.

```
bget(uint dev, uint blockno)
{
  struct buf *b;
  acquire(&bcache.lock);
  // Is the block already cached?
  for(b = bcache.head.next; b != &bcache.head; b = b->next) {
    if (b->dev == dev && b->blockno == blockno) {
      b->refcnt++;
                                     Note that we scan from
      release(&bcache.lock);
                                  the beginning of the MRU List.
      acquiresleep(&b->lock);
      return b;
```

• Find the matching buffer by scanning.

```
// Not cached; recycle an unused buffer.
// Even if refcnt==0, B DIRTY indicates a buffer is in use
// because log.c has modified it but not yet committed it.
for (b = bcache.head.prev; b != &bcache.head; b = b->prev) {
  if (b->refcnt == 0 && (b->flags & B DIRTY) == 0) {
    b \rightarrow dev = dev;
    b->blockno = blockno;
                                   Note that we scan from
    b \rightarrow flags = 0;
                                   the end of the MRU List
    b->refcnt = 1;
                                (the beginning of the LRU LIst).
    release(&bcache.lock);
    acquiresleep(&b->lock);
    return b;
panic("bget: no buffers");
```

void bwrite(struct buf *b)

• Write the modified buffer to the appropriate location on the disk.



• Write the modified buffer to the appropriate location on the disk.

```
bwrite(struct buf *b)
{
    if(!holdingsleep(&b->lock))
        panic("bwrite");
    b->flags |= B_DIRTY;
    iderw(b);
}
```

brelse()

- Decrease the reference counter.
- When reference counter is 0, remove it from the buffer cache and add it to the free buffer list.



Summary

- State of the buffer: B_VALID, B_DIRTY
- Locks
 - Buffer lock
 - Buffer cache lock
- API's
 - bread()
 - bget()
 - bwrite()
 - brelse()