# Lecture 31: Device driver and block I/O in xv6

Mythili Vutukuru

IIT Bombay

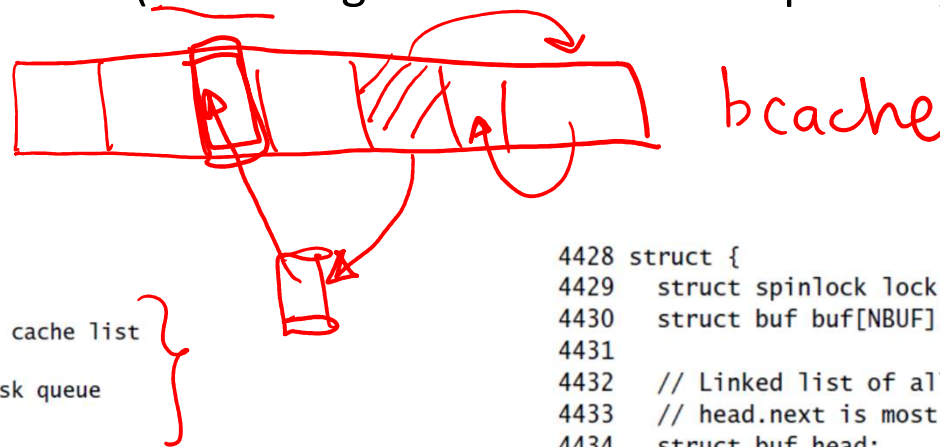https://www.cse.iitb.ac.in/~mythili/os/

# File systems and I/O in xv6

- Multiple layers of abstraction in file systems
  - System call implementations (open, read, write)
  - Operations on file system data structures (inodes, files, directories)
  - Block I/O layer (in-memory cache of disk blocks)
  - Device driver (communicates with hard disk to read/write blocks)
- This lecture and next: overview of these various layers in the xv6 file system

# Disk blocks and buffers

- Disk maintains data as 512-byte blocks
- Disk buffer (struct buf) = copy of disk block in memory
- Buffer cache (bcache) is an array of disk buffers
    - Pointers across buffers create a linked list, most recently used buffers at head
- Reading a block from disk: assign buffer for the block number in buffer cache, device driver sends read request to disk controller, disk controller raises interrupt when data is ready, data copied from disk controller into buffer cache(VALID flag set after data is read)
- Writing a block to disk: first write into buffer in buffer cache, device driver copies data from buffer to disk controller, disk controller raises interrupt when write completes (DIRTY flag is set until disk is updated)

bcache

```
3850 struct buf {
3851    int flags;
3852    uint dev;
3853    uint blockno;
3854    struct sleeplock lock;
3855    uint refcnt;
3856    struct buf *prev; // LRU cache list
3857    struct buf *next;
3858    struct buf *qnext; // disk queue
3859    uchar data[BSIZE];
3860 };
3861 #define B_VALID 0x2  // buffer has been read from disk
3862 #define B_DIRTY 0x4  // buffer needs to be written to disk
```

```
4428 struct {
4429    struct spinlock lock;
4430    struct buf buf[NBUF];
4431
4432    // Linked list of all buffers, through prev/next.
4433    // head.next is most recently used.
4434    struct buf head;
4435 } bcache;
```

# Device driver (1)

- Process that wishes to read/write calls iderw function, buffer as argument
  - If buffer is <u>dirty</u>, <u>write</u> request. If buffer is <u>not valid</u>, <u>read request</u>
  - Requests added to queue, function idestart issues requests one after another
  - Process sleeps until request completes
- Communication with disk controller <u>registers</u> via in/out instructions

```
4350 // Sync buf with disk.
4351 // If B_DIRTY is set, write buf to disk, clear B_DIRTY, set B_VALID.
4352 // Else if B_VALID is not set, read buf from disk, set B_VALID.
4353 void
4354 iderw(struct buf *b)
4355 {
4356   struct buf **pp;
4357
4358   if(!holdingsleep(&b->lock))
4359     panic("iderw: buf not locked");
4360   if((b->flags & (B_VALID|B_DIRTY)) == B_VALID)
4361     panic("iderw: nothing to do");
4362   if(b->dev != 0 && !havedisk1)
4363     panic("iderw: ide disk 1 not present");
4364
4365   acquire(&idelock);
4366
4367   // Append b to idequeue.
4368   b->qnext = 0;
4369   for(pp=&idequeue; *pp; pp=&(*pp)->qnext)
4370     ;
4371   *pp = b;
4372
4373   // Start disk if necessary.
4374   if(idequeue == b)
4375     idestart(b);
4376
4377   // Wait for request to finish.
4378   while((b->flags & (B_VALID|B_DIRTY)) != B_VALID){
4379     sleep(b, &idelock);
4380   }
4381
4382
4383   release(&idelock);
4384 }
```

```
4272 // Start the request for b.  Caller must hold idelock.
4273 static void
4274 idestart(struct buf *b)
4275 {
4276   if(b == 0)
4277     panic("idestart");
4278   if(b->blockno >= FSSIZE)
4279     panic("incorrect blockno");
4280   int sector_per_block =  BSIZE/SECTOR_SIZE;
4281   int sector = b->blockno * sector_per_block;
4282   int read_cmd = (sector_per_block == 1) ? IDE_CMD_READ :   IDE_CMD_RDMUL;
4283   int write_cmd = (sector_per_block == 1) ? IDE_CMD_WRITE : IDE_CMD_WRMUL;
4284
4285   if (sector_per_block > 7) panic("idestart");
4286
4287   idewait(0);
4288   outb(0x3f6, 0);  // generate interrupt
4289   outb(0x1f2, sector_per_block);  // number of sectors
4290   outb(0x1f3, sector & 0xff);
4291   outb(0x1f4, (sector >> 8) & 0xff);
4292   outb(0x1f5, (sector >> 16) & 0xff);
4293   outb(0x1f6, 0xe0 | ((b->dev&1)<<4) | ((sector>>24)&0x0f));
4294   if(b->flags & B_DIRTY){
4295     outb(0x1f7, write_cmd);
4296     outsl(0x1f0, b->data, BSIZE/4);
4297   } else {
4298     outb(0x1f7, read_cmd);
4299   }
```

4

# Device driver (2)

- When disk controller completes read/write operation, it raises an interrupt
  - Data is read from disk controller into buffer using "in" instruction
  - Process sleeping for data is woken up
  - Next request from queue is issued
- No support for DMA in x86. With DMA, data is copied by disk controller into memory buffers directly before raising interrupt
  - Interrupt handler need not copy data using I/O instructions

```
4302 // Interrupt handler.
4303 void
4304 ideintr(void)
4305 {
4306   struct buf *b;
4307
4308   // First queued buffer is the active request.
4309   acquire(&idelock);
4310
4311   if((b = idequeue) == 0){
4312     release(&idelock);
4313     return;
4314   }
4315   idequeue = b->qnext;
4316
4317   // Read data if needed.
4318   if(!(b->flags & B_DIRTY) && idewait(1) >= 0)
4319     insl(0x1f0, b->data, BSIZE/4);
4320
4321   // Wake process waiting for this buf.
4322   b->flags |= B_VALID;
4323   b->flags &= ~B_DIRTY;
4324   wakeup(b);
4325
4326   // Start disk on next buf in queue.
4327   if(idequeue != 0)
4328     idestart(idequeue);
4329
4330   release(&idelock);
4331 }
```

# Disk buffer cache: block read/write (1)

- All processes access disk via buffer cache only
- Only copy of disk block in cache, only one process can access it at a time
- Process calls "bread" to read a disk block, which calls function bget
  - Function bget returns buffer if it already exists in cache and no other process using it
  - If valid buffer not returned by bget, read from disk
- Process calls "bwrite" to write a block to disk: set dirty bit and request device driver to write
- When done with block, process calls brelse to release block, moves to head of list (most recently used)

```
4500 // Return a locked buf with the contents of the indicated block.
4501 struct buf*
4502 bread(uint dev, uint blockno)
4503 {
4504   struct buf *b;
4505
4506   b = bget(dev, blockno);
4507   if((b->flags & B_VALID) == 0) {
4508     iderw(b);
4509   }
4510   return b;
4511 }
```

```
4513 // Write b's contents to disk.  Must be locked.
4514 void
4515 bwrite(struct buf *b)
4516 {
4517   if(!holdingsleep(&b->lock))
4518     panic("bwrite");
4519   b->flags |= B_DIRTY;
4520   iderw(b);
4521 }
```

```
4525 void
4526 brelse(struct buf *b)
4527 {
4528   if(!holdingsleep(&b->lock))
4529     panic("brelse");
4530
4531   releasesleep(&b->lock);
4532
4533   acquire(&bcache.lock);
4534   b->refcnt--;
4535   if (b->refcnt == 0) {
4536     // no one is waiting for it.
4537     b->next->prev = b->prev;
4538     b->prev->next = b->next;
4539     b->next = bcache.head.next;
4540     b->prev = &bcache.head;
4541     bcache.head.next->prev = b;
4542     bcache.head.next = b;
4543   }
4544
4545   release(&bcache.lock);
4546 }
```
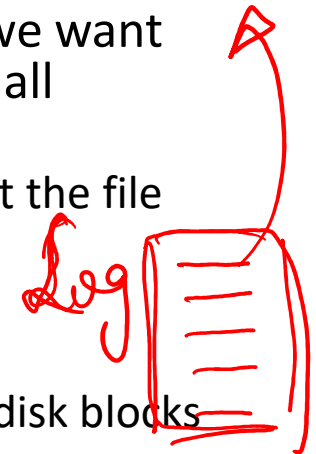
6

# Disk buffer cache: block read/write (2)

- Function bget returns pointer to a disk block if it exists in the cache
  - Ensures only one process at a time accesses a disk buffer
- If block in cache and another process using it, sleep until the block is released by the other process
- If block not in cache, find a least recently used non-dirty buffer and recycle it to use for this block
- Two goals achieved by buffer cache
  - Recently used disk blocks stored in memory for future use
  - Disk blocks modified by one process at a time

```
4465 static struct buf*
4466 bget(uint dev, uint blockno)
4467 {
4468   struct buf *b;
4469
4470   acquire(&bcache.lock);
4471
4472   // Is the block already cached?
4473   for(b = bcache.head.next; b != &bcache.head; b = b->next){
4474     if(b->dev == dev && b->blockno == blockno){
4475       b->refcnt++;
4476       release(&bcache.lock);
4477       acquiresleep(&b->lock);
4478       return b;
4479     }
4480   }
4481
4482   // Not cached; recycle an unused buffer.
4483   // Even if refcnt==0, B_DIRTY indicates a buffer is in use
4484   // because log.c has modified it but not yet committed it.
4485   for(b = bcache.head.prev; b != &bcache.head; b = b->prev){
4486     if(b->refcnt == 0 && (b->flags & B_DIRTY) == 0) {
4487       b->dev = dev;
4488       b->blockno = blockno;
4489       b->flags = 0;
4490       b->refcnt = 1;
4491       release(&bcache.lock);
4492       acquiresleep(&b->lock);
4493       return b;
4494     }
4495   }
4496   panic("bget: no buffers");
4497 }
```

# Logging layer (overview)

*disk blocks* (handwritten annotation)

- A system call can change multiple blocks at a time on disk, and we want atomicity in case the system crashes during a system call. Either all changes are made or none is made
  - Example: we do not want disk block added to the inode of a file but the file data not yet written to it
- Logging ensures atomicity by grouping disk block changes into transactions
  - Every system call starts a transaction in the log, writes all changed disk blocks in the log, and commits the transaction
  - Later, the log installs the changes in the original disk blocks one by one
  - If crash happens before log is written fully, no changes made
  - If crash happens after log entry is committed, log entries are replayed when system restarts after crash
- In xv6, changes of multiple system calls are collected in memory and committed to log together. Actual changes happen to disk blocks only after the group transaction commits
  - Process must call "log_write" instead of "bwrite" during system call

*Log* (handwritten annotation)

# Summary

- Device driver in xv6 communicates with disk controller using in/out instructions to read/write disk blocks
  - Simple driver with no DMA capability
- Buffer cache stores all recently read disk blocks in memory, and synchronizes access to disk blocks across processes
- All blocks changed in a system call are logged on disk and changes are installed atomically
- Next: File system code translates system calls into block read/write operations