

537 Final Exam,  
Old School Style

"The Job Interview"

Name: \_\_\_\_\_

This exam contains 13 old-school pages.

You are on a job interview. Unfortunately, on job interviews like this one, all they really want to see is that you can write code. So on this exam, you'll have to write some code. Even worse, you'll have to write code that shows that you understand how an operating system works. Yikes!

The good news: if you answer the questions well enough, you'll get the job!

The bad news: this is not true; it is just an exam.

Good luck!

The company is old school, and hence the exam is brought to you in 100% pure old-school style.

1. \_\_\_\_

2. \_\_\_\_

3. \_\_\_\_

4. \_\_\_\_

5. \_\_\_\_

6. \_\_\_\_

7. \_\_\_\_

8. \_\_\_\_

9. \_\_\_\_

10. \_\_\_\_

Total: \_\_\_\_ / 100

What is the maximum file size on a typical Unix system that uses an inode, an indirect pointer, and a double-indirect pointer? Assume a 4KB block size, and 12 direct pointers within the inode, and disk addresses that are 32 bits. Explain your work.

$$12 \text{ direct} + \left( \frac{4 \text{KB}^{\text{block size}}}{4 \text{ bytes}^{\text{disk addr}}} \right) + (1024 \times 1024)$$

indirect
double indirect

1024

$$\text{ptrs} = 12 + 1024 + (1024)^2$$

$$\text{size} = \text{ptrs} * 4 \text{KB}$$

Now, write a program that figures out the maximum file size on the file system you are running. What is the most efficient program you could write to do this? You are to use the following subset of the classic file-system API, i.e., `open()`, `read()`, `write()`, `close()`, `lseek()`.

// assume : file system has room for  
\* big file

```
int fd = open("file", O_RDWR | O_CREAT | O_TRUNC);
assert(fd >= 0);
char buf[4096];
int i = 0;
while(1) {
    int rc = write(fd, buf, 4096);
    if(rc < 0)
        break;
    i++;
}
printf("file is %d bytes\n", i * 4096);
close(fd);
```

You are to write the journal-recovery code. Actually, just a small piece of it: you are to write the piece of code that replays one journal entry. Assume that for this single transaction, all relevant blocks are already in memory. Here is the structure you should use to figure out what is in the transaction.

```
-----  
struct transaction {  
    int numblocks;  
    unsigned int destinations[MAX_BLOCKS];  
    unsigned char *blockarray[MAX_BLOCKS];  
};  
  
struct transaction *t; // use this pointer to access the struct  
-----
```

In it:

- 'numblocks' includes the number of blocks in the transaction.
- 'destinations' is an array that includes, for each block to be replayed, the on-disk block address to which that block should be written.
- 'blockarray' is an array of pointers to the blocks, which are all in memory right now.

Your task: Go through the transaction structure, and write each block to its final destination. You can use the following primitive to write to the disk:

- WRITE(int block, char \*data);

which takes a block number and a pointer to the data and writes the data to that block number.

```
for (i=0; i < t->numblocks; i++)  
    WRITE(t->destinations[i], t->blockarray[i]);  
// assume this never fails;  
// if it does, retry or exit
```

In this question, we have a highly-simplified version of the log-structured file system (LFS). You are supposed to write some code that emulates the cleaner. Specifically, go through each update in a segment and figure out whether the update has a \*live inode\* in it. If you find a live inode, print "LIVE", otherwise print "DEAD".

Here are some data structure definitions to help you out:

```
// the inode map: records (inumber) -> (disk address) mapping
unsigned int imap[MAX_INODES];

typedef struct __inode_t {
    int    direct[10]; // just 10 direct pointers
} inode_t;

typedef struct __update_t {
    int    inumber; // inode number of the inode in this update
    inode_t inode; // the inode
    int    offset; // offset of data block in file, from 0 ... 9
    char   data[4096]; // the data block
} update_t;

typedef struct __segment_t {
    int    disk_addr; // disk address of this segment (in bytes)
    update_t updates[MAX_UPDATES]; // the updates in this segment
    // (assume all MAX_UPDATES are used)
} segment_t;
```

// assume this is also in bytes for simplicity

(just like this)

segment\_t \*segment; // start with this pointer to the segment in question

Assume you are given a pointer to a segment\_t (called 'segment'). Then go through each update in the segment, figure out whether the inode referred to in that update is live or not, printing "LIVE" or "DEAD" as you go:

```
for (i = 0; i < MAX_UPDATES; i++) {
    update_t *u = segment -> updates[i];
```

```
    if (imap[u -> inumber] ==
        segment -> disk_addr + (i * sizeof(update_t))
        + sizeof(int))
        printf("LIVE\n");
    else
        printf("DEAD\n");
}
```

// key: does imap point to this inode?

skip over other updates

start at disk addr of seg

skip over the inumber in the update to get to the inode

}

TLB misses can be nasty. The following code can cause a lot of TLB misses, depending on the values of STRIDE and MAX. Assume that your system has a 32-entry TLB with a 8KB page size.

```
int value = 0;
int data[MAX]; // a big array

for (int j = 0; j < 1000; j++ {
    for (int i = 0; i < MAX; i += STRIDE) {
        value = value + data[i];
    }
}
```

What should you set MAX and STRIDE to so that you can achieve a TLB miss (but *\*not\** a page fault) upon pretty much every access to the array 'data'? (describe)

STRIDE must be  $\geq$  PAGESIZE (to TLB miss every time)  
 $\Rightarrow$  but, incrementing by sizeof(int) each time  
 $\Rightarrow \frac{PAGESIZE}{sizeof(int)} = \frac{8K}{4} = 2K$

MAX must include  $> 32$  pages  $\Rightarrow$  STRIDE \* 33

What happens if MAX is too high? Too low?

$MAX < STRIDE$ TLB hits (just access data[0])	$MAX \leq STRIDE * 32$ still just TLB hits	$MAX \geq STRIDE * 33$ and $MAX < MEMORY SIZE$ TLB misses	$MAX \geq MEMORY SIZE$ Page Faults
---	---	---	---------------------------------------

What happens if STRIDE is too high? Too low?

$STRIDE < PAGESIZE$ some hits + some TLB misses e.g. if $STRIDE = PAGESIZE/4$ , (TLB miss, TLB hit, TLB hit, TLB hit) repeat	$STRIDE \geq PAGESIZE$ TLB miss every time
---	---

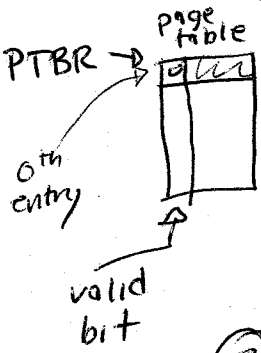
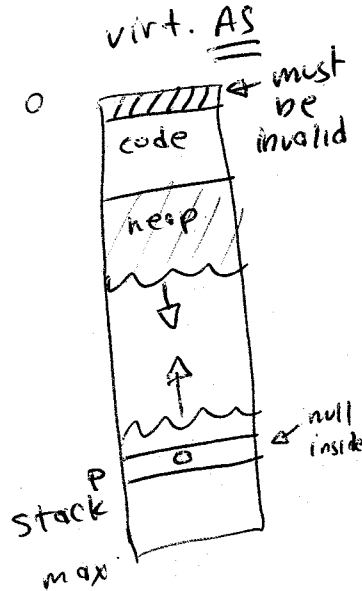
Sometimes badly-written C code dereferences a null-pointer, causing the program in question to crash. Write some code that does this; be brief!

```
int *p = NULL; // assign p to NULL
*p = 10; // deref p to attempt to set it to 10
```

Then, explain, in as much detail as needed, what the OS does in reaction to the null-pointer dereference, i.e., why it causes a program to fault, and how all the machinery behind it works. Be as detailed as you need to be!

Assume a system with linear page tables and 4KB pages.

①  $P$  is a valid memory location somewhere on the stack inside  $p$  we have stored an address, which is set to "NULL" or 0.

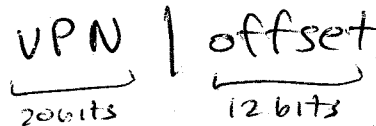


② when we dereference " $p$ " (e.g.  $*p = 10$ )

we are first loading the contents of  $p$  (which works, giving us a "0") and then accessing the virtual address within  $p$  or  $VA = \underline{0}$  in this case

⑥ an invalid access exception is raised; OS kills offending process

③ when  $VA = 0$  is accessed, the H/W splits it into



which here is 

0	0
VPN	off

 (null VA)

⑦ KEY: first page of AS must be invalid

for null deref to always seg fault!

④ the H/W checks  $VPN = 0$  in TLB, which misses  
 ⑤ the H/W or S/W checks the Page Table for  $VPN = 0$  (the first entry), and finds it invalid ( $valid = 0$ )



Some RAID code has been lost. You have to write it!

Assume you have a RAID-4 (parity-based RAID + a single parity disk), with a 4KB chunk size, and 5 disks total as follows:

DISK-0	DISK-1	DISK-2	DISK-3	DISK-4
block0	block1	block2	block3	parity(0..3)
block4	block5	block6	block7	parity(4..7)
.....	.....	.....	.....	.....

Fill in the routine SMALLWRITE() below:

```
// SMALLWRITE()
//
// This routine takes a logical block number 'block' and writes
// the single block of 4KB referred to by 'data' to it.
//
// It may have to use the underlying primitives:
// READ(int disk, int offset, char *data)
// WRITE(int disk, int offset, char *data)
// XOR(char *d1, char *d2) (xors one block with another)
```

```
void SMALLWRITE(int block, char *data) {
```

```
    int disk = block % 4;
    int off = block / 4;
```

```
    char olddata [4096], oldparity [4096];
```

```
could do in parallel { READ (disk, off, olddata);
                       READ (disk 4, off, oldparity); } read old data + parity
```

```
    char *newparity = XOR (oldparity, XOR (olddata, data));
```

```
could write new data + new parity in parallel too { WRITE (disk, off, data);
                                                       WRITE (4, off, newparity); }
    ② then flip bits in old parity to get new parity
    ① use this to see what data has changed
```

```
}
```

You are given the following code, which adds two vectors together, and does so in a multithread-safe way.

```
void
vector_add(vector *v1, vector *v2) {
    mutex_lock(v1->lock);
    mutex_lock(v2->lock);
    for (i = 0; i < v1->size; i++) {
        v1[i] = v1[i] + v2[i];
    }
    mutex_unlock(v1->lock);
    mutex_unlock(v2->lock);
}
```

Then you are told that two different concurrently-executing threads, 1 and 2, call this code as follows:

```
Thread 1:      Thread 2:
vector_add(&vectorA, &vectorB);      vector_add(&vectorB, &vectorA);
```

Unfortunately, this can lead to a DEADLOCK\*, in which the program gets stuck, with each thread waiting for the other to make progress.

- 1)- Why does this happen? (describe, or show with a picture)
- 2)- How could you write vector\_add() so that this deadlock never happens?

1) Thread 1:

```
mutex_lock(&vectorA->lock);
```

~~~~~>  
interrupt

Thread 2:

```
mutex_lock(&vectorB->lock);
//now this thread tries to
grab A's lock but
can't
```

//and same over here

2) Lots of solutions possible:

1) general lock around (avoid hold+wait)  
lock acquisition

```
@ 0 above: mutex_lock(&generalLock);
@ 1 above: mutex_unlock(&generalLock);
```

(avoid cycle)  
2) order locks, always require in same order  
e.g.  
if (v1 > v2) {  
 //grab v2, then v1  
} else {  
 //grab v1, then v2  
}

\* Sorry, a deadlock question. But it should not be too hard\*\*, should it?  
\*\* Probably not true.

Your co-worker implements the following code for condition variables (and specifically, the cond\_wait() and cond\_signal() routines) using \*semaphores\*:

```
typedef struct __cond_t {
    sem_t s;
} cond_t;

void cond_init(cond_t *c) {
    sem_init(&c->s, 0);
}

// cond_wait(): assumes that the lock 'm' is held
void cond_wait(cond_t *c, mutex_t *m) {
    mutex_unlock(&m); // release lock and go to sleep
    sem_wait(&c->s);
    mutex_lock(&m); // grab lock before returning
}

void cond_signal(cond_t *c) {
    sem_post(&c->s); // wake up one sleeping waiter (if there is one)
}
```

Unfortunately, it is buggy. Why doesn't this code work properly?

10

→ w/ "real" conditions, the wait would still get stuck waiting  
 (cond-signal only wakes a waiter if one is already waiting)

→ here, a sem-post would happen on (cond-signal), and a subsequent wait would not get stuck waiting

the problem is that the semaphore has state, whereas a condition var. does not

You are given a new atomic primitive, called FetchAndSubtract(). It executes as a single atomic instruction, and is defined as follows:

```
int FetchAndSubtract(int *location) {
    int value = *location; // read the value pointed to by location
    *location = value - 1; // decrement it, and store result back
    return value;          // return old value
}
```

You are given the task: write the lock\_init(), lock(), and unlock() routines (and define a lock\_t structure) that use FetchAndSubtract() to implement a working lock.

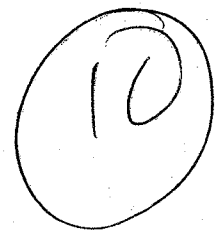
```
typedef struct lock_t {
    int ticket;
    int turn;
} lock_t;
```

```
lock_init(lock_t *lock) {
    lock->ticket = lock->turn = 0;
}
```

// ok to start @ 0 and go negative right away

```
lock(lock_t *lock) {
    int myturn = FAS(&lock->ticket);
    while (myturn != lock->turn)
        ; // spin
}
```

```
unlock(lock_t *lock) {
    FAS(&lock->turn);
}
```



// this is a ticket lock -> simpler solutions were ok too

// e.g. init lock to something      acquire by while FAS(&lock) < something ; // spin      release by setting back to something

NFS and AFS are two famous distributed file systems, yet they each have cases where one performs noticeably different than the other.

Assuming you can only access a single file (and using only the limited API: open(), read(), write(), close(), and lseek() calls), write a program that runs MUCH MUCH faster when run upon NFS than AFS.

open("file"); ← very fast on NFS //

many answers possible!

on AFS fetches WHOLE file

(even if file isn't accessed, or only some of it was read)

5

Using the same assumptions as above, write a program that runs MUCH MUCH faster when run upon AFS than NFS.

```
fd=open("file2", O_RDONLY);
while (1) // well, not forever, but a loop would do
  int rc=read(fd, ...);
  if (rc < 0)
    break;
```

```
close(fd)
```

key: if so, AFS → local NFS → goes to server

// assuming file > local mem } but < local disk }

(and (not needed) file has been accessed before)

5