Computer Science 162 David E. Culler University of California, Berkeley Practice Final Exam December XX, 2019

Name	
TA Name and Section Time	

This is **practice** for a closed book exam with two 2-sided pages of notes permitted. It is intended to be a 100 minute exam. You have 170 minutes to complete it. The number at the beginning of each question indicates the points for that question. Write all of your answers directly on this exam. Make your answers as concise as possible. If there is something in a question that you believe is open to interpretation, please raise your hand to request clarification. When told to open the exam, put your Student ID on every page and check that you have them all. The final pages are for reference. [* You can expect the topical coverage, length, and format to be similar. The places where we ask you to demonstrate working knowledge through innovation versus recall may differ, but the mix will be similar. *

By my signature below, I swear that this exam is my own work. I have not obtained answers or partial answers from anyone. Furthermore, if I am taking the exam early, I promise not to discuss it with anyone prior to completion of the regular exam, and otherwise I have not discussed it with anyone who took the early alternate exam.

X _____

Grade Table (for instructor use only)

			•				· ,		
Question:	1	2	3	4	5	6	7	8	Total
Points:	1	20	20	10	10	20	20	0	101
Score:									

1. (1 point) Write your Student ID on every page of the exam.

2. (20 points) Operating System Concepts: True/False

Choose True or False for 20 of the 25 questions below. Clearly cross out the five that should not be graded. You do not need to justify your answers.

1100	oc gradea	. Tot do not need to justify your answers.
(a)	and exter	If a process writes past the end of an open file the write system call will succeed ad the file with no reported error. ○ False
(b)	(1 point) by giving excessive	If a fair scheduler, like Linux CFS, strives to meet a responsiveness target T each task a quanta of T divided by the number of tasks, the overhead may get with a large number of tasks.
	•	○ False
(c)	hit rate,	Changing a cache from direct-mapped to fully-associative will <i>never</i> decrease its under any size, workload, and eviction policy.
		√ False
(d)	` - /	A single memory access may result in a TLB miss and a cache hit. ○ False
(e)	typically	As a system's workload is increased, below the half-power point, response time increases due to queuing delays. $\sqrt{\text{False}}$
(f)	` - /	LRU is an appropriate policy to evict page frames for demand paging. $\sqrt{\text{False}}$
(g)	have not	In the $x86$ architecture, the MMU clears the use bit of page table entries that been accessed in a long time.
	O True	$\sqrt{\text{False}}$
(h)	addresses	A device driver may interact with I/O devices by reading and writing to memory reserved for memory-mapped I/O. () False
(•)	· .	
(i)	` - /	In a typical file system, large files account for most in-use disk space. O False
(j)	(1 point) √ True	The compiled kernel code contains the file system root (/) at a known in umber. \bigcirc False
(k)	(1 point) file or its	In the FAT file system, file attributes are kept in the directory and not in the header.
	$\sqrt{\text{True}}$	○ False
(1)	` - /	In a Unix FFS-style inode-based file system, there is no limit on how large a file except for the size of the underlying storage device. V False
(100)	•	
(m)	enabled b	In the Windows NTFS file system, unlike FFS, sequential allocation on disk is by variable extents.
	√ True	○ False

(n) (1 point) Low-level file API (read, write) is accelerated by the kernel's buffer cache.
 √ True ○ False

(o)	(1 point) The Transmission Control Protocol (TCP) allows each endpoint of a connection to always have multiple outstanding unACKed packets at a time.
	\bigcirc True $\sqrt{\text{False}}$
(p)	(1 point) In the Network File System (NFS), a user's modifications to a file <i>might</i> become visible to other users in the system as soon as the write operation completes. √ True ○ False
(q)	
(r)	function on the key and shares it with the other nodes to partition the key space into equal-sized blocks.
(s)	 ○ True
(t)	
(u)	Containers handle package dependences for collections of processes, whereas Kubernetes introduces Pods as a mechanism for co-locating a related set of containers on a set of (virtualized) reasources. Very True O False
(v)	(1 point) Containers provide the illusion of a dedicated file system for a well-defined collection of executables and libraries and the processes that instantiate them. √ True ○ False
(w)	 (1 point) The socket abstraction is dictated by the transport layer of the network stack to allow a client and a server to communicate over HTTP. ○ True ✓ False
(x)	total number of resident virtual pages (assuming the size of physical and virtual pages are the same).
	○ True ✓ False
	Solution: This one is a little tricky. Every resident virtual page is placed in a page frame, so it would be true in the absence of shared pages. But many virtual pages may share the same physical page frame.
(y)	(1 point) A key/value store can provide some, but not all, aspects of a file system by

structuring the keys as pathnames, (e.g., /usr/bin/ps) and the values as file contents.

√ True ○ False

GET acts as reading the whole file.

3. (20	points)	Operating	System	Concepts:	Short	Answer

Do 20 of the following 25 points and cross out the ones that are not to be graded.

Fill in the blanks correctly for the questions below. You do not need to provide justifications.

- (a) (1 point) A process may request the operating system to perform a privileged action on its behalf by issuing a/an _____system call____.
- (b) (1 point) A/An ______ is a programming pattern in which a lock is used with one or more condition variables.
- (c) (1 point) The subset of its virtual address space that a process is actively using over a given time interval is called a/an <u>working set</u>.
- (d) (1 point) Running a virtual machine requires support in the host operating system in the form of a kernel module or driver called a/an <u>Virtual Machine Monitor</u>.
- (e) (1 point) An operation is said to be <u>idempotent</u> if performing it multiple times in a row has the same effect as performing it exactly once.
- (f) (1 point) A group of n machines can maintain safety and liveness in the face of fail-stop failures as long as at most $\frac{n-1}{2}$ node(s) are faulty, whereas with Byzantine failures a $\frac{n-1}{2}$ fraction of nodes must be non-faulty.
- (g) (2 points) What is the difference between a directory and an index?

Solution: A directory maps file names to inumbers. An index maps inumbers to inodes. A file description stores the inumber, allowing operations on an open file descriptor to access the index directly.

(h) (2 points) Why can't distributed file systems, like AFS and NFS, keep client-side caches consistent using techniques from multiprocessor cache coherence (e.g., snoopy caches on a bus)? What do they do instead?

Solution: Snoopy caches rely on a shared broadcast medium shared by all caches. No analogous shared broadcast medium exists among clients of a distributed file system. The server notifies clients when data they may have cached has been overwritten by another client.

(i) (2 points) The concept of a *namespace* arises in a variety of places in operating systems. Give two examples.

Solution: File paths, Process IDs, File handles, virtual address space, physical memory address space, thread IDs, User IDs, keys, disk block address.

(j) (2 points) Google's Internet search service receives approximately 64000 search queries per second, on average. Let's assume that the average latency for Google to process each query is approximately 400 ms. How many search queries, on average, are being processed by Google's system at any particular instant? Show your work.

Solution: $64000 \frac{\text{queries}}{\text{s}} \cdot 0.400 \text{ s} = 25600 \text{ queries}$

(k) (2 points) What is the difference between a virtual machine and a container?

Solution: A virtual machine provides software the illusion of having its own dedicated hardware where as a container provides software only its own dedicated operating system, including the set of processes and the file system.

(l) (2 points) Assume that a service is able to perform an operation in 100ms. If it has no internal parallelism, what is its rate of performing operation? If the service is 50% utilized and it takes 50ms to get to and from the services, how long should you expect your request to take.

Solution: At 100ms per operation, it can services 10 operations per second. The expected queue length is is 0.5/(1 - 0.5) = 1. So we expect the total response time to be 50 ms + 1*100 ms + 100 ms, or 250 ms.

(m) (2 points) For a server listening on a TCP socket, what must take place for it to complete the accept call? What does accept return?

Solution: A client performs a *connect* operation which causes TCP to carry out a 3-way handshake. A Syn is sent from client to server proposing a sequence number, a Syn-Ack is replied from the server recognizing the proposal and making a possibly different one, and an Ack is sent from client to server accepting the server-proposed starting sequence number.

After a successful handshake to establish the connection, the Accept call returns a file descriptor for the connection socket.

(n) (3 points) With a two level page table on a 32-bit machine with 4kb pages and 4-byte page table entries, list every page table and memory access involved in fetching and executing an instruction that loads a value into a register assuming no page faults occur. How many TLB accesses occur during this instruction? [* You can expect a more detailed work-em-out variant of this question on the exam. Try to give a complete answer here. *]

Solution: Fetch the Instruction:

Top level page table is accessed using PC bits [31:22] are used to get a level-2 PT page (concatenate 2 bits and add to PTBR); PC bits [21:12] are used to index into that page to get a page table entry; PC bits [11:0] are concatenated with the memory page frame number in the PTE to get the physical address of the instruction. That instruction is fetched, decoded and executed.

Execute the instruction:

The effective address A of the word to be read is computed. Bits [31:22] of A are used to index into the top level PT page to get the 2nd level page for A; Bits [21:12] of A are used to index into this to get the PTE. Bits [11:0] are concatenated with the frame number in the PTE to get the physical address of the value to be read.

That address is read and its contents placed in the destination register.

2 TLB accesses occur (not counting inserting into the TLB on a miss). Assuming all TLB accesses are misses, then four page table accesses are made.

(o) (2 points) Consider an IPC framework used to allow one process to provide machine learning services to other processes. The client has the following object that it wants to pass to the server for it to make a prediction.

```
struct predictor_obj {
    char *str;
    int num_weights;
    int weights[];
};
```

Explain why it cannot pass this object directly. What does it need to do instead?

Solution: It cannot pass any reference to an object in its virtual address space, i.e., a pointer. Instead is must serialize it by sending the string and the sequence of weights. It would format these into a buffer using an agreed upon syntax, for example, 4 bytes as #chars, that many chars, 4 bytes as #weights, that many weights, and send the contents of the buffer to the prediction server process.

4. (10 points) Pintos

Choose True or False for the questions below according to Pintos, the operating system you used in the projects for this class. You do not need to justify your answers.

- (a) (1 point) Accesses to kernel memory (PHYS_BASE and above) in the system call handler and interrupt handlers operate directly on physical addresses, bypassing all page tables.
 - True √ False
- (b) (1 point) It is inherently unsafe to call $sema_down$ on a zero-value semaphore in external interrupt context.
 - √ True False
- (c) (1 point) It is inherently unsafe to call sema_down on a zero-value semaphore in the page fault handler.
 - \bigcirc True $\sqrt{\text{False}}$
- (d) (1 point) Interrupts are disabled while Pintos switches threads (i.e., while the switch_threads function executes).
- (e) (1 point) On a successful call to pagedir_set_page, Pintos flushes the TLB.
 - True √ False
- (f) (1 point) When a page is allocated using palloc_get_page, Pintos must modify the page table of every user process so that the page is accessible in kernel mode.
 - True √ False
- (g) (1 point) The set of free blocks is stored on disk as a linked list.

Answer the question below on the provided line according to Pintos, the operating system you used in the projects for this class. You do not need to justify your answers.

(h) (1 point) What scheduling algorithm did you implement in Pintos in Task 2 of Project 2?

(h) Priority Scheduling

Answer the questions below in the box provided according to Pintos, the operating system you used in the projects for this class.

(i) (2 points) Suppose the scheduling timer expires while interrupts are disabled. When will the timer_interrupt function (i.e., the external interrupt handler) be invoked next?

Solution: It is invoked as soon as interrupts are re-enabled.

Practice	Final	Exam.	- Page	Q	$\circ f$	25
1 lactice	rinai	Daam .	- I age	J	OI	Δv

Student ID:	

5. (10 points) Inode Design

};

Alice, Benjamin, Catherine, and David are in a group for their CS 162 project. For Project 3, they choose an inode design with the following property: the first 162 bytes of a file are stored directly within the inode.

(a) (1 point) What does the "i" in "inode" stand for?

(b) (2 points) List one advantage of storing some file data directly in the inode.

Solution: Fewer disk accesses are needed to read/write very small files.

(c) (2 points) List one disadvantage of storing some file data directly in the inode.

Solution: It decreases number of disk pointers that fit in the inode, either limiting maximum potential file size or requiring more indirections for medium-sized files. Alternatively, requires a larger inode, meaning fewer inodes total.

(d) (3 points) Write out a definition of struct inode_disk with this property. Be sure to include 10 direct pointers, 1 singly indirect pointer, and 1 doubly indirect pointer, in addition to the first 162 bytes. Your inode structure must support directories, as required in Task 3 of the project. Do not assume useful data is stored in the magic or unused fields. struct inode_disk {

```
off_t length;
uint32_t isdir;
uint8_t data[162];
block_sector_t direct[10];
block_sector_t singly_indirect;
block_sector_t doubly_indirect;
unsigned_magic;
uint8_t unused[BLOCK_SECTOR_SIZE - (4 + 4 + 162 + 2 + 40 + 4 + 4 + 4)];
```

(e) (2 points) If the group decides to implement hard links in their file system, how would their struct inode_disk have to change?

Solution: It would need to contain a reference count recording how many entries in the directory structure reference the inode.

6. (20 points) Footex

Assume that shared pages are implemented in Pintos—that is, a physical page may be mapped into the address spaces of multiple processes. This provides the opportunity to synchronize two processes in *userspace* by placing a lock on a shared page.

In class, we studied the futex ("fast userspace mutex") system call. In this problem, you will implement footex, a simplified version of futex, for Pintos. It consists of two system calls, which each return true on success and false on failure. [* You can expect a more familiar variant of this question. *]

```
/* Checks if *ADDR == VAL and atomically blocks the calling process if so. */
bool footex_wait(int* addr, int val);
/* Wakes up NUM threads waiting on the provided address ADDR. */
bool footex_wake(int* addr, int num);
```

You may call the following function from userspace in Pintos:

```
/* Atomically sets *ADDR = VAL and returns the old value of *ADDR. */
int test_and_set(int* addr, int val);
```

As a warm-up, let's do a programming exercise with footex in a Pintos user program.

(a) (4 points) Implement a userspace lock in Pintos using footex_wait, footex_wake, and test_and_set. This code should run in a user program (i.e., not in the Pintos kernel). Your implementation should be syscall-free if the lock is uncontended. You may assume that footex_wait and footex_wait always return true for this question.

```
struct user_lock {
    int state;

    bool maybe_waiting;
};

void user_lock_init(struct user_lock* lock) {
    lock->state = 0; // unlocked state

    lock->maybe_waiting = false;
}

void user_lock_acquire(struct user_lock* lock) {
    while (test and set(&lock->state, 1) != 0) {
        lock->maybe_waiting = true;
        footex_wait(&lock->state, 1);
        lock->maybe_waiting = true;
}
}
```

```
void user_lock_release(struct user_lock* lock) {
  ASSERT(lock->state == 1); // locked state
  lock->state = 0;
  if (lock->maybe_waiting) {
   lock->maybe_waiting = false;
   footex_wake(&lock->state, 1);
}
```

Now, we will implement the footex_wait and footex_wake system calls in the Pintos kernel.

(b) (4 points) Extend the system call handler to support the two new system calls. You may call the following two functions:

```
void validate_user_buffer(void* pointer, size_t length);
void validate_user_string(const char* string);
These functions check if the provided buffer (or string) exists entirely within valid user-
accessible memory. If not, they terminate the calling process with exit code -1.
static void syscall_handler (struct intr_frame* f) {
   uint32_t* args = ((uint32_t*) f->esp);
   validate_user_buffer(args, sizeof(uint32_t));
    switch (args[0]) {
        /* Pre-existing cases (not shown) */
   case SYS_FOOTEX_WAIT:
       validate_user_buffer(&args[1], 2 * sizeof(uint32_t));
       validate_user_buffer((int*) args[1], sizeof(int));
        f->eax = (uint32_t) syscall_footex_wait((int*) args[1], (int) args[2]);
        break;
    case SYS_FOOTEX_WAKE:
        validate_user_buffer(&args[1], 2 * sizeof(uint32_t));
       validate_user_buffer((int*) args[1], sizeof(int));
        f->eax = (uint32_t) syscall_footex_wake((int*) args[1], (int) args[2]);
        break;
```

```
/* Additional pre-existing cases (not shown) */
...
}
```

(c) (1 point) Suppose process P calls footex_wait on the address addr1 and is put to sleep as a result. Then process Q calls footex_wake on the address addr2. How must addr1 and addr2 be related, for Q's footex_wake to wake up P? (Hint: consider that a single shared page may be mapped at different virtual addresses in different address spaces.)

Solution: addr1 must be mapped in P's address space, addr2 must be mapped in Q's address space, and they must both correspond to the same physical address (equivalently, the same kernel virtual address).

lently, the same kernel virtual address).

(d) (1 point) Keep track of the threads waiting on a particular address using struct footex.
struct footex {
 int* addr;
 struct list_elem elem;

 struct condition cond;
 int num waiters;

};

(e) (1 point) Determine what additional global variables must be added to syscall.c, and extend syscall_init to do any necessary initialization work.

/* Add additional global variables on the lines below. */
static struct list footex list;
static struct lock footex lock;

void syscall_init(void) {
 intr_register_int(0x30, 3, INTR_ON, syscall_handler, "syscall");

/* Perform any necessary initialization work. */

```
/* Perform any necessary initialization work. *
    list_init(&footex_list);
    lock_init(&footex_lock);
}
```

(f) (2 points) Complete the footex_get helper function below. Its intended behavior is to find the struct footex for the given address. If none exists and the create flag is true, then one is created. The provided addr is the virtual address in the calling process' address space.

```
static struct footex* footex_get(int* addr, bool create) {
   struct thread* t = thread_current();
   addr = pagedir_get_page(t->pagedir, addr)
   for (struct list_elem* e = list_begin(&footex_list);
                     e != list_end(&footex_list); e = list_next(e)) {
       struct footex* curr = list_entry(e, struct footex, elem);
       if (curr->addr == addr) {
           return curr;
       }
   }
   struct footex* f = NULL;
   if (create) {
       f = malloc(sizeof(struct footex));
       if (f != NULL) {
           f->addr = addr;
           list_push_back(&footex_list, &f->elem);
           cond_init(&f->cond);
           f->num_waiters = 0;
       }
   }
   return f;
}
```

Question continues on the next page \longrightarrow

}

(g) (4 points) Implement the system call handlers for footex_wait.
bool syscall_footex_wait(int* addr, int val) {
 lock acquire(&footex lock);

 if (*addr == val) {
 struct footex* f = footex_get(addr, true);
 if (f == NULL) {
 lock release(&footex lock);
 return false;

 }

 f->num_waiters++;
 cond_wait(&f->cond, &footex_lock);
 f->num_waiters--;
 if (f->num_waiters == 0) {

lock_release(&footex_lock);
return true;

list_remove(&f->elem);

free(f);

}

(h) (4 points) Implement the system call handlers for footex_wake. bool syscall_footex_wake(int* addr, int num) {

```
lock acquire(&footex_lock);

struct footex* f = footex_get(addr, false);
if (f == NULL) {
    lock release(&footex lock);
    return true;

}
for (int i = 0; i < num; i++) {
    cond signal(&f->cond, &footex_lock);

}
lock release(&footex_lock);

return true;
```

Solution: If we want checking *ADDR == VAL and putting the thread to sleep to be truly atomic, then we would have to disable interrupts around that code. This isn't strictly needed for the userspace lock to work, so we didn't include it above (and we wouldn't take off points on an exam for not including it).

An alternative solution is to not rely on a condition variable at all. Instead, one can maintain a list of threads (struct list) in struct footex, and use thread_block and thread_unblock directly. With this solution, one would disable and restore interrupts to create a critical section.

7. (20 points) **Journaling**

You are writing code late into the night to finish Project 3 for 162 when your laptop's battery dies. You plug in and restart your computer, and the operating system consults the file system's journal to determine how to recover from the unexpected shutdown. Fortunately, you're running a journaling file system that you wrote so there is a log of file system operations in persistent storage. The running state of the kernel is lost on reboot (of course) so a consistent file system state must be reestablished from the log on recovery. Its head pointer gives the next available entry and tail is the last completed operation. The contents of the journal are as follows. All prior log entries have been freed. The directory for /home/oski has been set to include an entry with name filesys.c and inode number 40.

Entry #	Transaction ID	Operation
101	10	Begin
102	10	Initialize inode 18
103	10	Add a new entry to the directory /home/oski with name
		inode.c and inode number 18.
104	11	Begin
105	11	Mark blocks 4, 8, and 15 as allocated in the free map.
106	11	Write <data buf=""> to blocks 4, 8, and 15</data>
107	10	Сомміт
108	11	Update inode 18 and set its first three direct pointers to refer to
		blocks 4, 8, and 15.
109	12	Begin
110	12	Mark blocks 16, 23, 42, and 108 as allocated in the free map.
111	11	Сомміт
112	12	Update inode 40 and set its first four direct pointers to refer to
		blocks 16, 23, 42, and 108.
113	13	Begin
114	13	Mark block 108 as free in the free map.
115	12	Commit
116	13	Clear the direct pointer to block 108 in inode 40.

(a) (2 points) Knowing that file system operations take a long time, in your implementation you decided to use threads to do all the updates to the disk for a transaction concurrently. What would the implementation of the COMMIT operation need to do?

Solution: It needs to ensure that all prior disk operations with its Transaction ID have completed successfully and then advance the tail. It is incrementing that pointer that performs the commit.

Practice	Final	Exam	- Page	17	of 25

Student ID:	

(b)	(4 points) What is the status of the file /home/oski/inode.c on disk and what happens to the log on recovery?
	Ŭ V
	1. Blocks 4, 8, and 15 allocated.
	$\sqrt{\text{Completed}}$, free log entry
	Occupied but log entry cannot be freed
	Maybe completed, could be replayed
	Maybe completed, must be replayed
	O Not complete, replay
	Maybe completed, discard
	O Not completed, discard
	○ Completed, discard
	2. Direct pointers in inode 18 are set to blocks 4, 8, and 15. A directory entry for inode.c
	added to /home/oski and refers to inode 18.
	$\sqrt{\text{Completed}}$, free log entry
	Completed but log entry cannot be freed
	Maybe completed, could be replayed
	Maybe completed, must be replayed
	Not complete, replay

Explain your answers.

Maybe completed, discardNot completed, discardCompleted, discard

Solution: All changes to inode.c are made within Transaction 1. There is a COMMIT entry for this transaction present in the journal, which means its individual operations were ready to be applied to the file system on disk, and the commit operation was performed, which means those disk operations completed. All the log entries for the transaction can be freed.

Question continues on the next page \longrightarrow

(c) (7 points) What is the status of the file /home/oski/filesys.c on disk, what happens to

the log on recovery, and what is the status of the file after recover	y?
1. Blocks 16, 23 and 42 allocated	
 Completed and free log entry 	
O Completed but log entry cannot be freed	
 Maybe completed, could be replayed 	
$\sqrt{\text{Maybe completed, must be replayed}}$	
○ Not complete, replay	
Maybe completed, discard	
○ Not completed, discard	
○ Completed, discard	
2. inode 40 initialized, its first four direct pointers set O Completed and log entry can be freed	
 Completed but log entry cannot be freed 	
 Maybe completed, could be replayed 	
$\sqrt{\text{Maybe completed, must be replayed}}$	
○ Not complete, replay	
Maybe completed, discard	
○ Not completed, discard	
○ Completed, discard	
3. Commit Transaction 12.	
O Completed and log entry can be freed	
Ompleted but log entry cannot be freed	
 Maybe completed, could be replayed 	
$\sqrt{\text{Maybe completed, must be replayed}}$	
○ Not complete, replay	
Maybe completed, discard	
○ Not completed, discard	
O Completed, discard	
4. Block 108 is marked as free	
Completed and log entry can be freed	
Ompleted but log entry cannot be freed	
Maybe completed, could be replayed	
○ Maybe completed, must be replayed	
○ Not complete, replay	
○ Maybe completed, discard	
$\sqrt{\text{ Not completed, discard}}$	
Completed discard	

Student ID:	

- 5. Direct pointer to block 108 cleared in inode 40
 - Ompleted and log entry can be freed
 - Ompleted but log entry cannot be freed
 - O Maybe completed, could be replayed
 - O Maybe completed, must be replayed
 - O Not complete, replay
 - O Maybe completed, discard
 - $\sqrt{}$ Not completed, discard
 - O Completed, discard

Explain your answers.

Solution: Transaction #12 has also been committed, but it is not known whether the COMMIT has completed its work, since the tail has not advanced past it. Each of its operations must be replayed to establish the state of the disk to reflect the commit of the transaction.

Transaction #3 has no COMMIT - it is not committed. Its operations, including removing block 108 from the file, must be eliminated in the recovery, effectively rolling back to just after transaction #12 where that block was present. The file system should not initiate these operations prior to the commit.

(d) (1 point) How would your answers change if item 115 were not present in the log?

Solution: Now, Transaction #12, the write of 4 blocks to /home/oski/filesys.c, would not be committed and its operations would need be discarded, rather than performed and committed in recovery.

(e) (1 point) Instead assume there as an additional entry in the log, as below. What does recovery do?

Entry #	Transaction ID	Operation
117	13	COMMIT

Solution: Now, Transaction #3 is in the same situation as Transactions #2. The truncation of file /home/oski/filesys.c would be performed as part of recovery.

Student ID:	

(f) (5 points) Logs to provide proper recovery show up in many places, besides file systems. For each of the following, describe the proper action for a worker to take in a two-phase commit protocol.

A GLOBAL-COMMIT Message is received when log contains:

Transaction ID: 17, State: Ready, Last Sent Msg: Vote-Commit. What action do you take?

Solution: Transition to Transaction ID: 17, State: Commit.

You have just rebooted and are in recovery. The log contains:

Transaction ID: 93, State: Init, Last Sent Msg: <na>.

What action do you take?

Solution: Set log to

Transaction ID: 93, State: Abort, Last Sent Msg: Vote-Abort.

Send Vote-Abort for Tx 93 to the coordinator.

You don't actually know that a VOTE-REQ has been set for 93, but this is a safe and conservative step. The coordinator can simply move on to a later transaction ID.

You have just rebooted and are in recovery. The log contains:

Transaction ID: 117, State: Ready, Last Sent Msg: Vote-Commit.

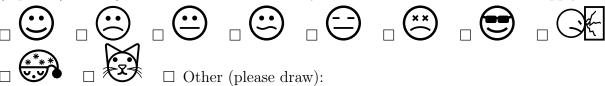
What action do you take?

Solution: It is possible that you died after making the log entry but before sending the message. Thus, you need to wait for either a Global-Abort or Global-Commit from the coordinator to determine the next state. It is safe to send (or resend) the Vote-Commit for #117 (you don't know whether it was sent before. (We didn't do this in HW6 because of the way the protocol was implemented.)

If you receive a message for a later transaction, you will need to query the coordinator for the outcome of #117. Optionally, you could query other workers to determine if any of them have VOTE-ABORT on #117. This would that you will eventually transition to ABORT for #117.

8. (0 points) Optional Questions

(a) (0 points) Having finished the exam, how do you feel about it? Check all that apply:



(b) (0 points) If there's anything you'd like to tell the course staff (e.g., feedback about the class or exam, suspicious activity during the exam, new logo suggestions, etc.) you can write it on this page.

This page is intentionally left blank.

Do not write any answers on this page. You may use it for *ungraded* scratch space.

```
int pthread_create(pthread_t *thread, const pthread_attr_t *attr,
                        void *(*start_routine) (void *), void *arg);
int pthread_join(pthread_t thread, void **retval);
int pthread_mutex_init(pthread_mutex_t *restrict mutex,
   const pthread_mutexattr_t *restrict attr);
int pthread_mutex_lock(pthread_mutex_t *mutex);
int pthread_mutex_unlock(pthread_mutex_t *mutex);
int pthread_cond_init(pthread_cond_t *cond, pthread_condattr_t *cond_attr);
int pthread_cond_signal(pthread_cond_t *cond);
int pthread_cond_broadcast(pthread_cond_t *cond);
int pthread_cond_wait(pthread_cond_t *cond, pthread_mutex_t *mutex);
int sem_init(sem_t *sem, int pshared, unsigned int value);
int sem_post(sem_t *sem);
int sem_wait(sem_t *sem);
pid_t wait(int *status);
pid_t fork(void);
pid_t waitpid(pid_t pid, int *status, int options);
int execv(const char *path, char *const argv[]);
void exit(int status);
/**********************************/
FILE *fopen(const char *path, const char *mode);
FILE *fdopen(int fd, const char *mode);
size_t fread(void *ptr, size_t size, size_t nmemb, FILE *stream);
size_t fwrite(const void *ptr, size_t size, size_t nmemb, FILE *stream);
int fclose(FILE *stream);
/************************************/Ow-Level I/O *****************************
int open(const char *pathname, int flags); (O_APPEND|O_CREAT|O_TMPFILE|O_TRUNC)
ssize_t read(int fd, void *buf, size_t count);
ssize_t write(int fd, const void *buf, size_t count);
int dup(int oldfd);
int dup2(int oldfd, int newfd);
int pipe(int pipefd[2]);
int close(int fd);
/***********************************/intos Lists *****************************/
void list_init(struct list *list);
struct list_elem *list_begin(struct list *list);
struct list_elem *list_next(struct list_elem *elem);
struct list_elem *list_end(struct list *list);
void list_insert(struct list_elem *before, struct list_elem *elem);
void list_push_front(struct list *list, struct list_elem *elem);
void list_push_back(struct list *list, struct list_elem *elem);
struct list_elem *list_remove(struct list_elem *elem);
struct list_elem *list_pop_front(struct list *list);
struct list_elem *list_pop_back(struct list *list);
bool list_empty(struct list *list);
#define list_entry(LIST_ELEM, STRUCT, MEMBER) ...
```

```
void sema_init(struct semaphore *sema, unsigned value);
void sema_down(struct semaphore *sema);
void sema_up(struct semaphore *sema);
void lock_init(struct lock *lock);
void lock_acquire(struct lock *lock);
void lock_release(struct lock *lock);
void cond_init(struct condition *cond);
void cond_wait(struct condition *cond, struct lock *lock);
void cond_signal(struct condition *cond, struct lock *lock);
void cond_broadcast(struct condition *cond, struct lock *lock);
enum intr_level intr_get_level(void);
enum intr_level intr_set_level(enum intr_level);
enum intr_level intr_enable(void);
enum intr_level intr_disable(void);
bool intr_context(void);
void intr_yield_on_return(void);
tid_t thread_create(const char *name, int priority, void (*fn)(void *), void *aux);
void thread_block(void);
void thread_unblock(struct thread *t);
struct thread *thread_current(void);
void thread_exit(void) NO_RETURN;
void thread_yield(void);
struct thread {
   /* Owned by thread.c. */
                                     /* Thread identifier. */
   tid_t tid;
                                    /* Thread state. */
   enum thread_status status;
   char name[16];
                                    /* Name (for debugging purposes). */
                                    /* Saved stack pointer. */
   uint8_t *stack;
                                    /* Priority. */
   int priority;
                                     /* List element for all threads list. */
   struct list_elem allelem;
   /* Shared between thread.c and synch.c. */
                                     /* List element. */
   struct list_elem elem;
#ifdef USERPROG
   /* Owned by userprog/process.c. */
   uint32_t *pagedir;
                                     /* Page directory. */
#endif
   /* Owned by thread.c. */
                                     /* Detects stack overflow. */
   unsigned magic;
void *palloc_get_page(enum palloc_flags); (PAL_ASSERT|PAL_ZERO|PAL_USER)
void *palloc_get_multiple(enum palloc_flags, size_t page_cnt);
void palloc_free_page(void *page);
```

```
void palloc_free_multiple(void *pages, size_t page_cnt);
#define PGBITS 12 /* Number of offset bits. */
#define PGSIZE (1 << PGBITS) /* Bytes in a page. */</pre>
unsigned pg_ofs(const void *va);
                                       uintptr_t pg_no(const void *va);
void *pg_round_up(const void *va);
void *pg_round_down(const void *va);
#define PHYS_BASE 0xc0000000
uint32_t *pagedir_create (void);
                                      void pagedir_destroy(uint32_t *pd);
bool pagedir_set_page(uint32_t *pd, void *upage, void *kpage, bool rw);
void *pagedir_get_page(uint32_t *pd, const void *upage);
void pagedir_clear_page(uint32_t *pd, void *upage);
bool pagedir_is_dirty(uint32_t *pd, const void *upage);
void pagedir_set_dirty(uint32_t *pd, const void *upage, bool dirty);
bool pagedir_is_accessed(uint32_t *pd, const void *upage);
void pagedir_set_accessed(uint32_t *pd, const void *upage, bool accessed);
void pagedir_activate(uint32_t *pd);
bool filesys_create(const char *name, off_t initial_size);
struct file *filesys_open(const char *name);
bool filesys_remove(const char *name);
struct file *file_open(struct inode *inode);
struct file *file_reopen(struct file *file);
void file_close(struct file *file);
struct inode *file_get_inode(struct file *file);
off_t file_read(struct file *file, void *buffer, off_t size);
off_t file_write(struct file *file, const void *buffer, off_t size);
bool inode_create(block_sector_t sector, off_t length);
struct inode *inode_open(block_sector_t sector);
block_sector_t inode_get_inumber(const struct inode *inode);
void inode_close(struct inode *inode);
void inode_remove(struct inode *inode);
off_t inode_read_at(struct inode *inode, void *buffer, off_t size, off_t offset);
off_t inode_write_at(struct inode *inode, const void *buffer, off_t size, off_t offset);
off_t inode_length(const struct inode *inode);
bool free_map_allocate(size_t cnt, block_sector_t *sectorp);
void free_map_release(block_sector_t sector, size_t cnt);
size_t bytes_to_sectors(off_t size);
struct inode_disk {
                                    /* First data sector. */
   block_sector_t start;
                                     /* File size in bytes. */
   off_t length;
   unsigned magic;
                                     /* Magic number. */
   uint32_t unused[125];
                                     /* Not used. */
};
/***********************************/
typedef uint32_t block_sector_t;
#define BLOCK_SECTOR_SIZE 512
void block_read(struct block *block, block_sector_t sector, void *buffer);
void block_write(struct block *block, block_sector_t sector, const void *buffer);
```