

# CSE 421/521 Midterm Solutions

## —SOLUTION SET—

25 Mar 2015

Please fill out your name and UB ID number above. Also write your UB ID number at the bottom of each page of the exam in case the pages become separated.

This midterm exam consists of three types of questions:

1. **10 multiple choice** questions worth 1 point each. These are drawn directly from lecture slides and intended to be easy.
2. **6 short answer** questions worth 5 points each. You can answer as many as you want, but we will give you credit for your best four answers for a total of up to 20 points. You should be able to answer the short answer questions in four or five sentences.
3. **2 long answer** questions worth 20 points each. **Please answer only one long answer question.** If you answer both, we will only grade one. Your answer to the long answer should span a page or two.

Please answer each question as **clearly** and **succinctly** as possible. Feel free to draw pictures or diagrams if they help you to do so. **No aids of any kind are permitted.**

The point value assigned to each question is intended to suggest how to allocate your time. So you should work on a 5 point question for roughly 5 minutes.

### Statistics:

- **135** students took this exam.
- **19** was the median score. (This exam was way too hard.)
- **19.5** was the average score.
- **8.37** was the standard deviation of the scores.

## Multiple Choice

1. (10 points) Answer all **ten** of the following questions. Each is worth **one** point.

- (a) What singer was *not* played before class this semester?  
 **Tony Rio**    **Amanda Banana**    **Stacy Keys**    **Lisa Lee**
- (b) All of the following are critical section requirements *except*  
 mutual exclusion.    **concurrency.**    progress.    performance.
- (c) Interprocess communication is harder than intra-process communication.  
 **True**    False
- (d) All of the following are private to each process thread *except*  
 stack.    **file handles.**    registers.
- (e) What action does *not* require the kernel?  
 **Switching between two threads.**    Reading from a file.    Creating a new process.    Altering a virtual address mapping.
- (f) What part of the address space is not initialized using information from the ELF file?  
 Code segment.    **The heap.**    Uninitialized static data.    Initialized static data.
- (g) The Rotating Staircase Deadline Scheduler is most similar to which other scheduling algorithm?  
 Lottery scheduling.    Round-robin.    **Multi-level feedback queues.**  
 Random.
- (h) Which of the following is *not* a part of making a system call?  
 Arranging the arguments in registers or on the stack.    Loading the system call number into a register.    Generating a software interrupt.  
 **Allocating memory in the heap using malloc.**
- (i) What information would probably *not* be stored in a page table entry?  
 The location on disk.    Read, write or execute permissions.    **The process ID.**    The physical memory address.
- (j) Paging using fixed-size pages can suffer from internal fragmentation.  
 **True.**    False.

## Short Answer

Choose **4 of the following 6** questions to answer. You may choose to answer additional questions, in which case you will receive credit for your best four answers.

2. (5 points) One way to eliminate deadlock when acquiring multiple locks is to eliminate cycles. Describe another way to avoid deadlock (2 points) and how to implement it (3 points).

Graded by Zihe Chen.

### Solution:

There isn't much you can do about protected access to shared resources, which requires waiting—either passive or active. In some cases adding resources to eliminate unnecessary sharing could be an option, and we'll accept that. (Eliminating sleeping isn't a solution, since replacing deadlocks with livelocks isn't a great idea.) But the other two conditions are potential solutions:

1. **Add resources**, eliminating the need for sharing. As an example from the dining philosophers problem, adding chopsticks eliminates the need to share. We'll accept this answer if you were quite specific about this not working in all circumstances, since in certain cases global resources are required for correctness.
2. **Allow resource preemption**, creating a way for the operating system to forcibly take a resource from a thread that is holding it. One way to implement this is to detect a deadlock by identifying a circular dependency in the wait queue and then interrupt one of the threads holding the locks causing the cycle. This would require some way for the thread to register a handler each time it requests a lock that would be called if that lock was force dropped.
3. **Disallow multiple independent requests**, preventing a thread from holding some resources while requesting others. You could implement this by checking each time a thread requests a lock to make sure that it doesn't hold any other locks. Enforcing this without any new lock acquisition mechanisms would require programmers to change their code to introduce new locks covering cases. Alternatively, you could create a new locking mechanism allowing multiple locks to be acquired simultaneously.

### Statistics:

- **122 out of 135** students answered this question.
- **2** was the median score.
- **2.18** was the average score.
- **1.9** was the standard deviation of the scores.

**Grader Feedback:**

Common mistakes:

- Rephrasing eliminating cycles method stated in the problem, for example, saying that resources need to be ordered for locking.
- Implementing their methods by using `do if cannot acquire a lock`, which is not correct since it is impossible to know if you can immediately acquire a lock successfully.
- Not allowing parallelism.

3. (5 points) Describe two changes to the OS virtual memory system that you would need or might want to make to accommodate 48-bit virtual and physical addresses (1 point each). For each, describe how it would affect the computation or memory overhead of virtual to physical translation (2 points each).

Graded by Jinghao Shi.

**Solution:**

There were a few options here:

1. **Page table entry size.** You clearly need to do something about your page table entries, since whatever you did to bitpack for 32-bit virtual addresses is going to change with 48-bit virtual addresses. Almost inevitably, you're going to end up storing more state.
2. **Address space size.** You don't necessarily have to change the address space size, since you can simply fit more content from 32-bit wide address spaces into a system with 48-bit virtual addresses. You may want to change the break point between userspace and the kernel, however, since there is no point wasting any of the 4 GB wide virtual address space on the kernel at this point. Everything from  $0 \times 0$  to  $0 \times \text{FFFFFFFF}$  should be used for the process, with kernel addresses outside of that range ( $\geq 0 \times 100000000$ ). This doesn't necessarily have any direct effect on the overheads, although it may affect the size of other data structures (PTEs or page tables).
3. **Larger pages.** You may really want to increase the 4K page size at this point, with the concordant tradeoff between TLB faults and spatial locality. Larger pages may reduce the number of translations required.
4. **Different page table structures.** It could be time to add a third level to the existing two-level page tables, particularly if you decide to expand the the address space size. The goal would be to reduce the amount of space required for the page tables for common address space layout patterns.
5. **MMU changes.** Clearly the MMU and TLB need to be modified to help map 48-bit virtual addresses to 48-bit physical addresses. No effect on the kernel overheads here, but any speculation about the affect on hardware (wider addresses might imply fewer entries) would be accepted.

**Statistics:**

- **60 out of 135** students answered this question.
- **3** was the median score.
- **2.7** was the average score.
- **1.84** was the standard deviation of the scores.

4. (5 points) Describe how a scheduling algorithm might improve resource allocation by observing processes communicating using pipes (4 points). What is a tradeoff involved in this decision (1 point)?

Graded by Jinghao Shi.

**Solution:**

Assuming I have the following shell pipeline: `foo | bar`. Due to the pipe between them, before each run we know that `bar` cannot make much progress until `foo` runs and begins to fill the pipe. So we could just run `foo` to completion, buffer all of its output in the pipe, and then run `bar` to completion. This would minimize the context switches between them.

However, the tradeoff here is that this requires a large amount of memory—or *some kind of storage*—to buffer the pipe contents. Switching between the processes can help reduce memory overhead by keeping the pipe size smaller. Ideally, we could run them on separate cores, but might still need to schedule them carefully so that `bar` has enough work to do when it runs to justify the context switch while `foo` doesn't fill the buffer too quickly.

**Statistics:**

- 108 out of 135 students answered this question.
- 2 was the median score.
- 1.82 was the average score.
- 1.6 was the standard deviation of the scores.

**Grader Feedback:**

There are two common mistakes:

- Processes that communicate with a pipe should be allocated with same priority and be scheduled to run together. **False:** writers should be scheduled before readers.
- Using pipe is an indication that the processes are more interactive, therefore should be given higher priority. **False:** there is no absolute correlation between using pipes and being interactive.

```
1 #include <stdio.h>
2 #include <stdlib.h>
3
4 int32_t foo[1024];
5
6 int
7 main(int argc, const char * argv) {
8     buffer = (char *) malloc(10240); // You can assume that malloc succeeds.
9     int32_t bar[1024];
10    // <-- Here -->
11    struct foover[24];
12 }
```

5. (5 points) Examine the simple program above. At the point in its execution indicated, *at minimum* how many 4K pages of memory will it require in each segment: code, stack and heap? Justify your answer (5 points).

Note that you can ignore dynamically-loaded libraries, and assume that `malloc` does not consume any memory for its own data structures (if only).

Graded by Zihe Chen.

**Solution:**

The process has:

- At least **two code pages**: one for the 4K `foo` array, and at least one more for the rest of the code. (We'd accept answers where you pointed out that we didn't ask about the data segment and note that that's where `foo` would actually be allocated.)
- At least **two stack pages**: one for the 4K `bar` array, and at least one more for anything else on the stack. (There is at least one other thing on the stack: `argv`!)
- At least **three heap pages** for the 10K dynamically-allocated `buffer`.

So at least seven.

**Statistics:**

- **81 out of 135** students answered this question.
- **2** was the median score.
- **1.49** was the average score.
- **1.63** was the standard deviation of the scores.

6. (5 points) Some systems provide separate exception handlers for TLB-faults and all other kinds of exceptions, allowing the operating system to handle them separately. First, describe why this might be a useful feature for the hardware to provide (2 points). Second, suggest one way that the operating system might take advantage of this differentiation (3 points).

Graded by Jinghao Shi.

**Solution:**

The reason to have separate exceptions is to allow the OS to establish separate code paths—this is kind of given away by the second part of the question. This is particularly important given that TLB exceptions are extremely common and, if you improve performance on this hot code path, then the system will run faster overall.

One way that the OS could take advantage of this differentiation is by designing a TLB handler that doesn't require saving all of the state required by the system call path. In general, system calls and hardware interrupts can branch off into very long code paths, meaning that it is reasonable to save all of the registers up front. However, the TLB code path is extremely specific, and so might be able to avoid saving as much state—at least initially. If it branches off into a page fault, then more state might need to be saved, but that could be done at that point.

**Statistics:**

- 104 out of 135 students answered this question.
- 1.5 was the median score.
- 1.22 was the average score.
- 1.22 was the standard deviation of the scores.

**Grader Feedback:**

- Some students mistakenly interpreted “separate exception handler for TLB-faults”—they thought the system provides hardware TLB fault handler.
- Most students got the first part, yet missed the second part of the question.



VPN	SLT	VPN	PPN	VPN	PPN	VPN	PPN
0		0	9	0	31	0	94
1	450	1	84	1		1	
2		2		2		2	
3		3	39	3	72	3	10
4		4		4	62	4	13
5	120	5		5		5	
6		6	14	6	40	6	8
7		7	10	7	43	7	
8		8		8		8	
9	320	9		9		9	

(a) TLT                      (b) SLT at 450                      (c) SLT at 320                      (d) SLT at 120

Table 1: Process Page Tables

7. (5 points) The HappyStudent architecture has byte-addressable memory with 10-byte virtual pages. The HappyOS operating system uses two-level pages tables to support 1000-byte virtual address spaces by dividing the virtual page number into a 10-byte second-level index and using the rest as the top-level index.

First, given the virtual address 273, identify the virtual page number, the offset, and the first and second level page table indices (1 point).

Second, given the top-level page table (TLT) and second-level page tables (SLT) above for the currently-running process provided above, indicate the result of the following four virtual to physical page translations (1 point each).

Graded by Zihe Chen.

**Solution:**

For 273, the virtual page number is 27, the offset is 3, and the first and second level page indices are 2 and 7.

Translations follow:

- 165 → 145
- 578 → Invalid (no entry in second-level table).
- 57 → Invalid (no entry in top-level table).
- 900 → 310

**Statistics:**

- 63 out of 135 students answered this question.
- 2 was the median score.
- 1.94 was the average score.
- 1.66 was the standard deviation of the scores.

## Long Answer

Choose 1 of the following 2 questions to answer. **Do not answer both questions.** If you do, we will only read one, most likely the one that looks shorter and more incorrect. If you need additional space, continue and clearly label your answer on other exam sheets.

8. (20 points) Choose one of the following two questions to answer:

1. **Interface Size Tradeoffs.** Linux and other UNIX variants provide a fairly thin system call interface, consisting of roughly 300 system calls. You are now familiar with a subset of these calls, particularly those you implemented for ASST2. In contrast, the Windows kernel interface contains almost an order-of-magnitude more system calls: approximately 3000.

First, describe the tradeoff between small and large system call interfaces. What is good about a thin interface? How about a thick interface (4 points)? Second, provide three examples illustrating the problems with the thin system call interface you are familiar with. For each, describe how adding additional system calls would help, and what their interface would be (4 points each).

Finally, for *one* of your new system calls describe the operating system changes that would be required to implement it. You are free (but not required) to reference OS/161-specific implementation details in your answer (4 points).

2. **Concurrent System Calls.** As core counts have increased, the OS has become a source of potential bottlenecks for multithreaded applications. One source of poor scaling behavior is when two system calls cannot be executed concurrently, and in some cases redesigning the system call interface can improve concurrency and performance on multicore systems.

Consider the code below and a multithreaded application where, on an  $N$  core machine,  $N$  thread run `openclose`<sup>1</sup>. First, describe the OS performance bottleneck that this code might encounter on a multicore system. Note that `open` is implemented to return the lowest available file descriptor, but very few apps rely on this behavior. Be sure to describe the source of the bottleneck clearly (10 points). Second, describe how alter the system call interface to remove this particular performance bottleneck (10 points).

```
1 void openclose() {
2     // Repeatedly open and close a file.
3     while (1) {
4         int fd = open("/tmp/foo");
5         close(fd);
6     }
7 }
```

---

<sup>1</sup>Don't worry about why an app would do anything so seemingly meaningless. If it's helpful, just assume your boss wrote it and is wondering why it doesn't work well.

Graded by Guru Prasad Srinivasa.

**Statistics:**

For (8.1):

- **75 out of 57** students answered this question.
- **5** was the median score.
- **6.85** was the average score.
- **4.7** was the standard deviation of the scores.

For (8.2):

- **57 out of 135** students answered this question.
- **0** was the median score. (This question was too hard.)
- **3.79** was the average score.
- **6.3** was the standard deviation of the scores.

**Grader Feedback:**

Generally these questions were not that hard if you read them carefully, figure out what they were asking, and then *followed instructions* when crafting your answer. Both questions include a fairly clear template for what your solution was intended to look like: "First ... continue ... You will want to consider ... Second ... You will want to think about ..." Sadly some of you seemed to fail to read this part of the question!

It also seemed that many students ran out of time answering the short answer questions and didn't devote enough time to the long answer. Given that it was 20 points, this was not a particularly wise decision.

Individual grader feedback specific to each question is included below.

## Solution: Interface Size Tradeoffs.

**First, describe the tradeoff between large and small system call interfaces.** Large interfaces allow processes to be more specific about exactly what they want the operating system to do, which eventually leads to fewer system calls and fewer transitions into the kernel, which could improve performance. (Large systems call interfaces also *might* expose more features, although you were free to assume that the feature set exposed by both large and small interfaces was the same.)

In contrast, small interfaces are easier for the operating system to implement and maintain and make it more clear to application developers which of multiple ways to do things is the correct way. If there's only one way to accomplish something, it's what the OS developers will work to support; if there are multiple ways, it's not necessarily clear which is the most up-to-date. This was the canonical complaint about the Windows interface by third-party developers, which was (partly due to backwards compatibility), with half-a-dozen options available to accomplish a particular task the right approach wasn't clear. However, if you *worked* for Microsoft (say, developing Internet Explorer), you could just email the guys at the Windows team and ask "Hey, which way of reading data from a file has the best performance?"

**Second, provide three examples illustrating the problems with the thin system call interface you are familiar with.** The best way to do this is to show where multiple system calls are required to accomplish something useful that could be done using a single call. In none of the cases that we provide as solutions does the implementation required much past adding a new system call number and handler, but your solution may be different.

1. `fork` followed by `exec`: this was mentioned in class. Creating a child and loading a new image in one new call isn't just easier on the process, it also allows the OS to avoid some of the memory overhead associated with `fork`. Implementing this requires adding the system call and handling errors properly, which could occur either during `fork` or `exec`.
2. `open` followed by `read` or `write`: if I want to read the entire contents of a file into memory, why not just have a single call that takes a path and a buffer and performs the operation without requiring two calls and the maintenance associated with the process file table? Implementing this requires adding the new system call and handler, but not necessarily much else.
3. `read` or `write` and `lseek`: the fact that `read` and `write` manipulate the file offset implicitly and therefore require separate calls to `lseek` is both annoying and may cause problems when multiple threads try to use the same file handle. A separate system call could perform a `read` or `write` using an offset provided as an argument, eliminating the need to use `lseek`. Implementing this requires adding the new system call.

**Grader Feedback:**

1. Instead of talking about problems and solutions of a thin syscall interface, students have often described what they believed were missing but useful syscalls, such as 'copyfile()'.  
2. They also seem to consider thin interfaces to be faster than thick which is incorrect—more code is necessary to maintain necessary abstractions such as the VFS layer

## Solution: Concurrent System Calls

This question was inspired by this fantastic MIT paper on scalable interfaces. We may look at this in more detail later in the semester.

**First, describe the OS performance bottleneck that this code might encounter on a multicore system.** The key was to process the note that we left you, noting that `open` currently is implemented to return the lowest-available file descriptor. However, processes don't rely on this behavior and really only treat the file descriptor as a reference to the file handle, so you could alter this.

And you want to, because it's the performance bottleneck. Each time `open` is called it has to lock and search the file table from the beginning to find the lowest-available entry. (Note that in this case we *cannot* apply the read-lock-read optimization used in the 2014 midterm, because the file descriptor is supposed to be the lowest available.) The locking and long critical section caused by the required linear search have the effect of unnecessarily serializing concurrent calls to `open` that could proceed independently on multiple cores.

**Second, describe how to alter the system call interface to remove this particular performance bottleneck.** Once you see the bottleneck this is pretty obvious: change `open` to not guarantee that it will return the lowest-available file descriptor. This has two benefits. First, we can start the search at a random location, meaning that we will find an available file descriptor more quickly *even* if decide to grab a lock during the process. And if you were even more clever (and present at our midterm review), you probably noticed that you could apply the read-lock-read optimization here as well and avoid holding a lock outside of the loop.

### Grader Feedback:

Scores on this question were unusually low for the following reasons:

1. A lot of people did not fully understand what was being asked
2. People tried to fix the code provided instead of the underlying syscall implementation
3. Many failed to identify the underlying bottleneck