

CS 0449: Introduction to  
**S Y S T E M S   S O F T W A R E**

---

**Jonathan Misurda**  
Computer Science Department  
University of Pittsburgh  
[jmisurda@cs.pitt.edu](mailto:jmisurda@cs.pitt.edu)  
<http://www.cs.pitt.edu/~jmisurda>

**VERSION 3, REVISION 0**

Last modified: August 12, 2014 at 8:54 P.M.

Copyright © 2014 by Jonathan Misurda

This text is meant to accompany the course CS 0449 at the University of Pittsburgh. Any other use, commercial or otherwise, is prohibited without permission of the author. All rights reserved.

Java is a registered trademark of Oracle Corporation.

*This reference is dedicated to the students of CS 0449, Fall 2007 (2081). Their patience in dealing with a changing course and feedback on the first version of this text was greatly appreciated.*



# Contents

<b>Contents</b>	<b>i</b>
<b>List of Figures</b>	<b>v</b>
<b>List of Code Listings</b>	<b>vii</b>
<b>Preface</b>	<b>ix</b>
<b>1 Pointers</b>	<b>1</b>
1.1 Basic Pointers . . . . .	2
1.1.1 Fundamental Operations . . . . .	2
1.2 Passing Pointers to Functions . . . . .	4
1.3 Pointers, Arrays, and Strings . . . . .	5
1.3.1 Pointer Arithmetic . . . . .	6
1.4 Terms and Definitions . . . . .	7
<b>2 Variables: Scope &amp; Lifetime</b>	<b>8</b>
2.1 Scope and Lifetime in C . . . . .	9
2.1.1 Global Variables . . . . .	10
2.1.2 Automatic Variables . . . . .	10
2.1.3 Register variables . . . . .	11
2.1.4 Static Variables . . . . .	12
2.1.5 Volatile Variables . . . . .	15
2.2 Summary Table . . . . .	16
2.3 Terms and Definitions . . . . .	16

<b>3 Compiling &amp; Linking: From Code to Executable</b>	<b>17</b>
3.1 The Stages of Compilation . . . . .	17
3.1.1 The Preprocessor . . . . .	18
3.1.2 The Compiler . . . . .	19
3.1.3 The Linker . . . . .	20
3.2 Executable File Formats . . . . .	25
3.3 ☕ Linking in Java . . . . .	28
3.4 Terms and Definitions . . . . .	28
<b>4 Function Pointers</b>	<b>30</b>
4.1 Function Pointer Declarations . . . . .	30
4.2 Function Pointer Use . . . . .	32
4.2.1 Function Pointers as Parameters . . . . .	32
4.2.2 Call/Jump Tables . . . . .	34
4.3 Terms and Definitions . . . . .	36
<b>5 Processes &amp; Address Spaces</b>	<b>37</b>
5.1 Pages . . . . .	38
5.2 Terms and Definitions . . . . .	39
<b>6 Stacks, Calling Conventions, &amp; Activation Records</b>	<b>40</b>
6.1 Calling Convention . . . . .	42
6.2 Variadic Functions . . . . .	46
6.3 Buffer Overrun Vulnerabilities . . . . .	48
6.4 Terms and Definitions . . . . .	49
<b>7 Dynamic Memory Allocation Management</b>	<b>50</b>
7.1 Allocation . . . . .	50
7.1.1 Allocation Tracking Using Bitmaps . . . . .	51
7.1.2 Allocation Tracking Using Linked Lists . . . . .	53
7.1.3 Allocation Algorithms . . . . .	55
7.2 Deallocation . . . . .	57
7.2.1 Using Linked Lists . . . . .	58
7.2.2 Using Bitmaps . . . . .	58
7.2.3 ☕ Garbage Collection . . . . .	59
7.3 Linked List Example: <code>malloc()</code> . . . . .	62
7.4 Reducing External Fragmentation: The Buddy Allocator . . . . .	64

7.5	Terms and Definitions . . . . .	65
<b>8</b>	<b>Operating System Interaction</b>	<b>66</b>
8.1	System Calls . . . . .	66
8.1.1	Crossing into Kernel Space from User Space . . . . .	68
8.1.2	Unix File System Calls . . . . .	68
8.1.3	Process Creation with <code>fork()</code> and <code>execv()</code> . . . . .	70
8.2	Signals . . . . .	72
8.2.1	Sending Signals . . . . .	73
8.2.2	Catching Signals . . . . .	74
8.3	Terms and Definitions . . . . .	76
<b>9</b>	<b>Multiprogramming &amp; Threading</b>	<b>78</b>
9.1	Threads . . . . .	81
9.1.1	User Threading . . . . .	81
9.1.2	Kernel Threading . . . . .	83
9.2	Terms and Definitions . . . . .	83
<b>10</b>	<b>Practical Threading, Synchronization, &amp; Deadlocks</b>	<b>85</b>
10.1	Threading with <code>pthreads</code> . . . . .	85
10.2	Synchronization . . . . .	88
10.2.1	Mutexes . . . . .	90
10.2.2	Condition Variables . . . . .	91
10.2.3	Semaphores . . . . .	94
10.3	Deadlocks . . . . .	96
10.4	Terms and Definitions . . . . .	97
<b>11</b>	<b>Networks &amp; Sockets</b>	<b>98</b>
11.1	Introduction . . . . .	98
11.2	Berkeley Sockets . . . . .	101
11.3	Sockets and Threads . . . . .	102
11.4	Terms and Definitions . . . . .	104
<b>A</b>	<b>The Intel x86 32-bit Architecture</b>	<b>106</b>
A.1	AT&T Syntax . . . . .	108
A.2	Intel Syntax . . . . .	108
A.3	Memory Addressing . . . . .	110
A.4	Flags . . . . .	110

A.5	Privilege Levels . . . . .	111
<b>B</b>	<b>Debugging Under <code>gdb</code></b>	<b>113</b>
B.1	Examining Control Flow . . . . .	114
B.2	Examining Data . . . . .	116
B.3	Examining Code . . . . .	117
B.4	<code>gdb</code> Command Quick Reference . . . . .	118
B.5	Terms and Definitions . . . . .	119
<b>References For Further Reading</b>		<b>120</b>
<b>Index</b>		<b>122</b>
<b>Colophon</b>		<b>125</b>

# List of Figures

1.1	Two diagrammatic representations of a pointer pointing to a variable.	3
1.2	The values of variables traced in the <code>swap()</code> function.	5
2.1	Static electricity as an analogy to static data	9
2.2	The nested scopes in a C program.	10
3.1	The phases of the <code>gcc</code> compiler for C programs.	18
3.2	Static linking inserts library code into the executable.	21
3.3	Dynamic linking resolves calls to library functions at load time.	22
3.4	Dynamically loaded libraries are loaded programmatically and on-demand.	24
5.1	A process has code and data in its address space.	38
6.1	The shared boundary between two functions is a logical place to set up parameters.	41
6.2	A comparison of the calling conventions of MIPS and x86	42
6.3	A function with one parameter.	43
6.4	A function with two parameters.	45
6.5	A program with a buffer overrun vulnerability.	48
7.1	A bitmap can store whether a chunk of memory is allocated or free.	52
7.2	A linked list can store allocated and unallocated regions.	53
7.3	Coalescing free nodes on deallocation.	58
7.4	Reference counting can lead to memory leaks.	60
7.5	Copying garbage collectors divide the heap in half and move the in-use data to the reserved half, which has been left empty.	62
7.6	Heap management with <code>malloc()</code> .	63

8.1	A library call usually wraps one or more system calls. . . . .	67
8.2	A “Hello world” program run through <code>strace</code> . . . . .	67
8.3	The standard signals on a modern Linux machine. . . . .	73
9.1	While the CPU sees just one stream of instructions, each process believes that it has exclusive access to the CPU. . . . .	78
9.2	The life cycle of a process. Dashed lines indicate an abnormal termination. . . . .	80
9.3	Thread State versus Process State. . . . .	80
10.1	Two threads independently running the same code can lead to a race condition. . . . .	89
10.2	Synchronizing the execution of two threads using a mutex. . . . .	90
10.3	The producer/consumer problem in pseudocode. . . . .	92
11.1	The Internet Layer Model. . . . .	99
11.2	Layout of an IP packet. . . . .	99
11.3	Each layer adds its own header to store information. . . . .	100
11.4	The functions of Berkeley Sockets divided by role. . . . .	101
A.1	The 32-bit registers. . . . .	107
A.2	<code>%eax</code> , <code>%ebx</code> , <code>%ecx</code> , and <code>%edx</code> have subregister fields. . . . .	107
A.3	The instructions used in this book. . . . .	108
A.4	Privilege rings on an x86 processor. . . . .	112

# List of Code Listings

2.1	Variables in inner scopes can shadow variables of the same name in enclosing scopes. . . . .	11
2.2	A pointer error caused by automatic destruction of variables. . . . .	12
2.3	Unlike automatic variables, pointers to static locals can be used as return values. . . . .	13
2.4	The <code>strtok()</code> function can tokenize a string based on a set of delimiter characters. . . . .	14
3.1	The a.out header section. . . . .	26
3.2	String literals are often deduplicated. . . . .	27
4.1	Function pointer example. . . . .	31
4.2	<code>qsort()</code> definition. . . . .	32
4.3	<code>qsort()</code> ing strings with <code>strcmp()</code> . . . . .	33
4.4	<code>qsort()</code> ing structures with a custom comparator . . . . .	35
6.1	A variadic function to turn its arguments into an array. . . . .	47
8.1	Using the Unix system calls to do a “Hello World” program. . . . .	69
8.2	An example of process creation using <code>fork</code> . . . . .	71
8.3	Launching a child process using <code>fork</code> and <code>execvp</code> . . . . .	72
8.4	Signals can be sent programmatically via <code>kill</code> . . . . .	74
8.5	<code>SIGALRM</code> can be used to notify a program after a given amount of time has elapsed. . . . .	75
10.1	Basic thread creation. . . . .	86
10.2	Inserting a yield to voluntarily give up execution. . . . .	87
10.3	Waiting for the spawned thread to complete. . . . .	88

10.4	Using a <code>pthread_mutex</code> to protect an enqueue operation on a shared queue. . . . .	91
10.5	The producer function using condition variables. The consumer function would be similar. . . . .	93
10.6	The producer function using semaphores. The consumer function would be similar. . . . .	95
11.1	A server using sockets to send “Hello there!”. . . . .	103
A.1	Hello world in AT&T assembler syntax. . . . .	109
A.2	Hello world in Intel assembler syntax. . . . .	109

# Preface

AT SOME POINT late last decade, the core curriculum of many major cs programs, including the one here at the University of Pittsburgh, switched to teaching Java, C#, or python. While there is no mistaking the prevalence or importance of a modern, Object-Oriented, garbage-collected programming language, there is also plenty of room for an “old-fashioned” do-it-yourself language like C. It is still the language of large programs and Operating Systems, and learning it opens the door to doing work with real-life systems like GNU/Linux.

Armed with a C compiler, we can produce executable programs, but how were they made? What does the file contain? What actually happens to make the program run? To be able to answer such questions about a program and its interactions seems to be fundamental to the issue of defining a system. Biologists have long known the benefit of studying life in its natural environment, and Computer Scientists should do no different. If we follow such a model and study a program’s “life” as we run it on the computer, we will begin to learn about each of the parts that work together. Only then can we truly appreciate the abstractions that the Operating System, the hardware, and high-level programming languages provide to us.

This material picks up where high-level programming languages stop, by looking at interactions with memory, libraries, system calls, threading, and networking. This course, however, avoids discussing the implementations of the abstractions the system provides, leaving those topics for more advanced, specialized courses later in the curriculum.

I started out writing this text not as a book, but as a study guide for the second exam in CS 0449: *Introduction to Systems Software*. I did not have a specific textbook except for the C programming portion of the course, and felt that the students could use something to help them tie in the lectures, the course slides, and the projects. So I sat down one Friday afternoon and, through the next four days, wrote a 60-page “pamphlet.” I’ve since decided to stick with my effort for the next term,

as part of a four-book curriculum that includes two other freely available texts on Linux Programming and Linux Device Drivers (links available in the Bibliography Section).

Over time, I hope to continue to add new topics and refine the material presented here. I appreciate any feedback that you, the reader, might have. Any accepted significant improvement is usually worth some extra credit to my students. All correspondence can be directed to the email address found on the front cover.

## Acknowledgments

No book, not even one so quickly written, is produced in a vacuum. I would like to thank all of my students for their patience while we all came to terms with my vision for the course. I especially want to thank Nathaniel Buck, Stacey Crotti, Jeff Kaminski, and Gerald Maloney for taking the time to provide detailed feedback. I'd also like to thank my parents and friends for contributing in ways from just basic support to the artistic. This text would not be the same without all of your help.

## Notes on the Second Edition

In this edition, I have attempted to improve the integration of the material into the modern Java curriculum. Several sections (some old, some new) are marked with the  icon and show up in the PDF outline in *italics*. These indicate that the material relates to the Java language or the Java Virtual Machine.

—Jonathan Misurda  
May 1, 2008

# 1 | Pointers

AS PART OF IMPLEMENTING a sorting algorithm, we often need to exchange the values of two items in an array. Good programming practice suggests that when we have some commonly-reused code we should wrap it in a function:

```
void swap(int a, int b) {  
    int t = a;  
    a = b;  
    b = t;  
}
```

Then we can call it from our program with `swap(a,b)`. However, if we initially set `x=3; y=5;` and run the above swap function, the values of `x` and `y` remain unchanged. This should be unsurprising to us because we know that when we pass parameters to a function they are passed “by value” (also known as passing “by copy”).

This means that `a` and `b` contain the same values as `x` and `y` at the moment that the call to `swap()` occurs because there is an implicit assignment of the form `a=x; b=y;` at the call site. From that point on in the function `swap()`, any changes to `a` and `b` have no effect on the original `x` and `y`. Thus, our swap code works fine inside the scope of `swap()` but once the scope is destroyed when the function returns, the changes are not reflected in the calling function.

We then wonder if there is another way to write our `swap()` so that it succeeds. Our first inclination might be to attempt to return multiple values. In C, like in Java, this is not directly possible. A function may only return one thing. However, we could wrap our values in a structure or array and return the one aggregate object, but then we have to do the work of extracting the values from the object, which is just as much work as doing the swap in the first place. We soon realize that we cannot write a swap function that exchanges the values of its arguments in any reasonable fashion. If we are programming in Java with the primitive data types,



this is where our attempts must stop. However, in C we have a way to rewrite the function to allow it to actually work.

In fact, this should not be surprising because we are essentially asking a function to return multiple pieces of data, and we have already seen one function that can do that: `scanf()`. If we pass `scanf()` a format string like `"%d %d"` it will set two argument variables to the integers that the user inputs. In essence, it is doing what we just said could not be done: It is modifying the values of its parameters. How is that accomplished? The answer is by using pointers.

## 1.1 Basic Pointers

A **pointer** is a variable that holds an address. An **address** is simply the index into memory that a particular value lives at. Memory (specifically RAM) is treated as an array of bytes, and just like our regular arrays, each element has a numerical index. We haven't needed pointers until now because we have given our variables names and used those names to refer to these locations. We don't even know the actual addresses because the compiler and the system automate the layout and management of many of our variables. However, unlike in Java, it is possible in C to ask for the location of a particular variable in memory.

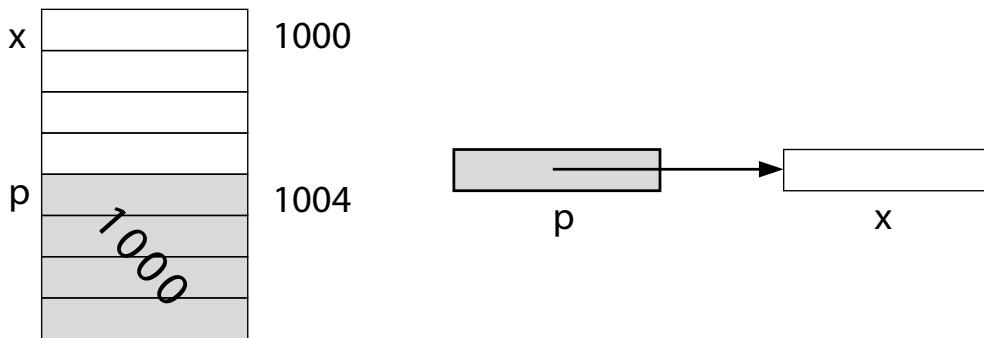
Referring to a location via a name or an address is not something unique to C, or even to computers in general. For example, you can refer to the room in which I work at Pitt as "Jon's Office," which is a name, or as "6203 Sennott Square" which is an address. Both are ways of referring to the same location and this duplication is known as *aliasing*.

### 1.1.1 Fundamental Operations

We can declare a pointer variable using a special syntax. Pointers in C have a type to them just like our variables did, but in the context of pointers, this type indicates that we are pointing to a memory location that stores a particular data type. For instance, if we want to declare a pointer that will hold the address in memory of where an integer lives, we would declare it as:

```
int *p;
```

where the asterisk indicates that p is a pointer. We need to be careful with declaring multiple variables on one line because its behavior in regards to pointers is surprising. If we have the declaration:



**Figure 1.1:** Two diagrammatic representations of a pointer pointing to a variable.

```
int *p, q;
```

or even:

```
int* p, q;
```

we get an *integer pointer* named `p` and an *integer* named `q`. No matter where you place the asterisk, it binds to the next variable. To avoid confusion, it is best to declare every variable on its own separate line.

To set the value of a pointer, we need to be able to get an address from an existing variable name. In C, we can use the address-of operator, which is the unary ampersand (`&`) to take a variable name and return its address:

```
int x;
int *p;
```

```
p = &x;
```

This code listing declares an integer `x` that lives someplace in memory and an integer pointer `p` that also lives somewhere in memory (since pointers are variables too). The assignment sets the pointer `p` to “point to” the variable `x` by taking its address and storing it in `p`.

Figure 1.1 shows two ways of picturing this relationship. On the left, we have a possible layout of RAM, where `x` lives at address 1000 and `p` lives at address 1004. After the assignment, `p` contains the value 1000, the address of `x`. On the right, much more abstractly, is shown the “points-to” relationship.

Now that we have the pointer `p` we can use it as another name for `x`. But in order to accomplish that, we need to be able to traverse the link of the points-to relationship. When we follow the arrow and want to talk about the location a pointer points-to rather than the pointer itself, we are doing a **dereference** operation. The dereference operator in C is also the asterisk (\*). Notice that although we used the asterisk in the pointer definition to declare a variable as a pointer, the dereference operator is different.

When we place the asterisk to the left of a pointer variable or expression that yields a pointer, we chase the pointer link and are now referring to the location it points to. That means that the following two statements are equivalent:

```
x = 4;           *p = 4;
```

Note that it is usually a mistake to assign a pointer variable a value that is not computed from taking the address of something or from a function that returns a pointer. This general rule should remind us that `p = 4;` would not be appropriate because we do not normally know in advance where memory for objects will be reserved.

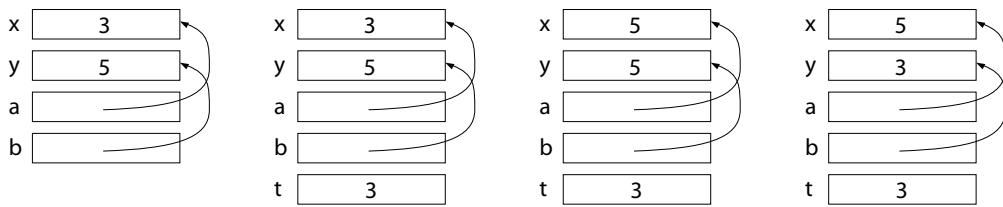
## 1.2 Passing Pointers to Functions

With the basic operations of address-of and dereference, we can begin to use pointers to do new and useful tasks. If we make a slight modification to our previous `swap()` function, we can get it to work:

```
void swap(int *a, int *b) {
    int t = *a;
    *a = *b;
    *b = t;
}
```

We also need to change the way we invoke it. If we have our same variables, `x` and `y`, we would call the function as `swap(&x, &y)`. After `swap()` returns, we now find that `x` is 5 and `y` is 3. In other words, the swap worked.

To better understand what happened here, we can trace the code constructing a picture like before. Figure 1.2 shows the steps. When the `swap()` function is called, there are four variables in memory: `x` and `y` which contain 3 and 5, respectively, and `a` and `b` which get set to the addresses of `x` and `y`. Next, a temporary variable `t` is created and initialized to the value of what `a` points to, i.e., the value of `x`. We then



**Figure 1.2:** The values of variables traced in the `swap()` function.

copy the value of what `b` points to (namely the value of `y`) into the location that `a` points to. Finally, we copy into the location pointed to by `b` our temporary variable.

Our swap function is an example of one of the two ways that pointers as function parameters are used in C. In the case of `swap()`, `scanf()`, and `fread()`, among many others, the parameters that are passed as pointers are actually acting as additional *return values* from the functions.

The other reason (which is not mutually exclusive from the previous reason) that pointers are used as parameters is for time and space efficiency in passing large objects (such as arrays, structs, or arrays of structs). Since we have already established that C is pass-by-value, if we pass a large object to a function, that object would have to be duplicated and that might take a long time. Instead, we can pass a pointer (which is really just integer-sized), thus taking no noticeable time to copy. This is why `fwrite()` takes a pointer parameter even though it does not change the object in memory.

## 1.3 Pointers, Arrays, and Strings

With our knowledge of pointers, we now can understand why we had to do certain things such as prefix variable names with an ampersand for `scanf()`. However, you may recall there passing a string was the exception to that rule. To understand why this is the case, we need to explain a fundamental identity in C:

The name of an array is a (read-only) pointer to its beginning.

Imagine we declare an array: `int a[4];`. If we wish to refer to a particular integer in the array, we can subscript the array and write something along the lines of `a[i]`, which represents the  $i^{th}$  item. An alternative way to view it is that `a[i]` is  $i$  integers away from the start of the array. This is valid because we know that all

elements of an array are laid out consecutively in memory. Thus, in essence, if we knew where the entire array started, we could easily add an offset to that address and find a particular element in the array.

To express it mathematically, if we have an array with address base, then the  $i^{th}$  element lives at  $\text{base} + i \times \text{sizeof(type)}$  where type is the type of the elements in the array. This simple calculation is so convenient that C decides to enable it by making the name of the array be a pointer to the beginning of the array in memory. This means that if we have an array element indexed as  $a[3]$ , it lives at the address  $a + 3 \times \text{sizeof(int)}$ .

However, when we convert this formula to C code, there is a slight adjustment we must make. The appropriate way to rewrite  $a[3]$  using pointers and offsets is  $*(a+3)$ . Where did the `sizeof(int)` term go? The answer lies in how C does **pointer arithmetic**.

### 1.3.1 Pointer Arithmetic

Pointers are allowed three mathematical operations to be performed on them. Addition and subtraction of an integer offset is allowed in order to support the aforementioned address/offset calculations. The third operation is that two pointers of the same type are allowed to be subtracted from each other in order to determine an offset between the pointers.

In the case of addition or subtraction of an integer offset, C recognizes that if you start with a pointer to a particular type, you will still want a pointer to that same type when you do your arithmetic. If the expression `a+1` added one byte to the pointer, we'd now be pointing into the middle of the integer in memory. What C wants is for `a` to point to the next int, so it automatically scales the offset by the size of the type the pointer points to.

Because addition and subtraction are supported, C also allows the pre- and post-increment and -decrement operators to be applied to pointers. These are still equivalent to the `+1` and `-1` operations as on integers.

This means that a particularly sadistic person could write the following function:

```
void f(char *a, char *b) {
    while(*a++ = *b++) ;
}
```

This function is one we have already discussed in the course. It is `strcpy()`. The way that it works is that the post-increment allows us to walk one character at a time through both the source and destination arrays. The dereferences turn the

pointers into character values. The assignment does the copy of a letter, and yields the right-hand side as the result of the expression. The right-hand side is then evaluated in a boolean context to determine if the loop stops or continues. Since every character is non-zero, the loop continues. The loop terminates when the nul-terminator character is copied, because its value is zero and thus false. The loop needs no body, the entirety of the work is done as side-effects of the loop condition.

## 1.4 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Address** The index into memory where a particular value lives.

**Dereference** To follow the link of a pointer to the location it points to. In C, the dereference operator is `*`.

**Pointer** A variable that holds an address.

**Pointer Arithmetic** Adding or subtracting an offset to a pointer value. In C, this offset is automatically scaled by the type the pointer points to.

## 2 | Variables: Scope & Lifetime

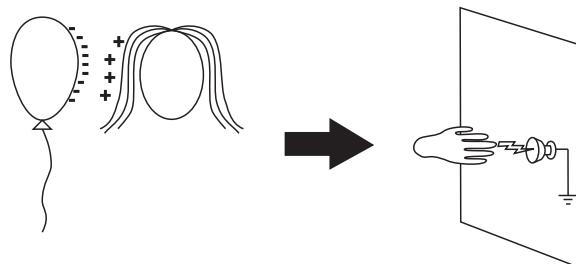
IN A PROGRAMMING LANGUAGE, **scope** is a term that refers to the region in a program where a symbol is legal and meaningful. A **symbol** is a name that represents a constant, literal, or variable. Scope serves a dual purpose. It allows for a symbol to be restricted to some logical division of the program, and it allows for duplicate names to exist in non-overlapping regions.

Restricting the portion of a program where a symbol is legal allows for a simple form of encapsulation, also known as “data-hiding.” By restricting a symbol’s access to a programmer-defined, limited region of the program, a programmer can be sure that minimal code affects the value of a particular variable. Allowing duplicate names is important when there are many programmers working together. Imagine a language where every variable was global. For many people to work together, they would have to avoid using the same variable names as anyone else. This would lead to some cumbersome naming convention and not encourage using the clearest name for a variable.

While scope is a compile-time property of code, when the program is actually executed, variables are created and destroyed. The time from which a particular memory location is allocated until it is deallocated is referred to as that variable’s **lifetime**.

The precise rules governing scope are programming language dependent, and a language construct rather than something enforced by the organization of the computer. If we have a variable that is restricted to a particular function (usually called a **local variable**), there is nothing preventing that variable from being created at program start and not destroyed until program termination.

What determines the lifetime of a variable is whether it is **static** or **dynamic** data. Static means non-moving, just as *static electricity* is an electrical charge separated from an opposite charge, and not moving towards it (as shown in Figure 2.1). In computing, static refers to when a program is compiled or anytime before it is run. If



**Figure 2.1:** Static electricity is a built-up charge that is separated from an opposite charge. When you touch a grounded object, the charges suddenly flow to cancel out, giving you a shock.

the compiler can compute exactly how much space will be needed for variables, the variables are static data. Dynamic data is that which is allocated while the program runs, since its size may depend on input or random values.

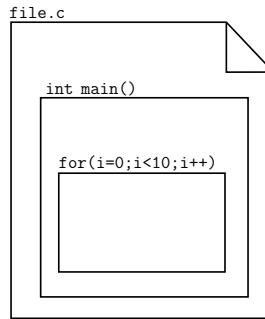
At first glance, local variables would seem to be static data, since the compiler can determine exactly how much space is necessary for them. However, while the compiler may be able to compute the memory needs of a function, there is not always a way to determine how many times that function may be called if, for instance, it is recursive. Since each invocation of the function needs its own copy of the local variables, allocation of their storage must be dealt with at run-time.

Static data can be allocated when the program begins and freed when the program exits. Dynamic data, on the other hand, will need special facilities to handle the fluctuations in the demand for memory while the program runs.

## 2.1 Scope and Lifetime in C

Scopes in C are defined by files and blocks. Remember that a block is a region of code enclosed in curly braces: { and }. In C, local variables are legal in the scope they are declared in, as well as all scopes that are nested inside of that scope. Figure 2.2 illustrates three nested scopes in a C program: the file, the function, and a loop within the function. A variable declared in the file is accessible anywhere, but a variable declared inside the loop can only be used inside that loop.

While two variables with the same name may not occupy the same scope, there is nothing preventing a nested scope from naming a variable the same as one in an enclosing scope. When this occurs, the variable of the outer (larger) scope is



**Figure 2.2:** The nested scopes in a C program.

hidden, and the innermost variable is said to be **shadowing** the outer one.

Listing 2.1 shows an example of a shadowed variable. The new block redeclares the shadowed variable, and when this program is run, the value “6” is displayed on the screen. In general, using shadowed variables is not good practice and can lead to a great deal of confusion to someone reading the code.

### 2.1.1 Global Variables

A **global variable** is a variable that is accessible to all functions and which retains its value throughout the entire execution of the program. In C, any variable declared outside of a function could be considered global. However, C treats a file as a scope, so “global” variables are actually limited to the file they are declared in.

What makes these file-scoped variables special is that they can be imported into the scope of other files by the `extern` keyword. If File A contains a global declaration like `int x`, File B can also refer to that variable by redeclaring it, but with the `extern` qualifier: `extern int x`.

### 2.1.2 Automatic Variables

Variable declarations that occur inside functions are implicitly declared `auto`, meaning an automatic variable. Automatic variables are variables that are created and destroyed automatically by code generated from the compiler. In general, the lifetime of automatic variables is the lifetime of the block they are defined in. However, it may be more convenient for all automatic variables to be created on function entry and destroyed at function return. This is implementation specific and has

```
#include <stdio.h>

int main() {
    int shadowed;
    shadowed = 4;
{
    int shadowed;
    shadowed = 6;
    printf("%d\n", shadowed);
}
return 0;
}
```

**Listing 2.1:** Variables in inner scopes can shadow variables of the same name in enclosing scopes.

no bearing on the correctness of a program since the scope is narrower than the lifetime.

The only problem that could arise with automatic variables comes as an abuse of pointers. Consider the code of Listing 2.2. Here function `f()` creates and returns a pointer to an automatic variable, which function `main()` captures. However, when `main()` goes to use the memory location referenced by the pointer `p`, that variable is “dead” and can no longer safely be used. The `gcc` compiler is kind enough to issue a warning if we do this:

```
(14) thot $ gcc escape.c
escape.c: In function `f':
escape.c:5: warning: function returns address of local variable
```

However, the code actually compiles and a program is produced. Proving again that C is not picky about what you do, no matter if you mean it or not.

### 2.1.3 Register variables

The old C compilers were not always very intelligent when it came to machine code generation. Thus, a programmer was allowed to give a hint to the compiler to indicate that a certain variable was particularly important. Such heavily accessed variables should be stored in architectural registers rather than in memory. Declar-

```

int *f() {
    int x;
    x = 5;
    return &x;
}

int main() {
    int *p;
    p = f();
    *p = 4;
    return 0;
}

```

**Listing 2.2:** A pointer error caused by automatic destruction of variables.

ing a variable as `register` could have a significant improvement on the performance of frequently executed code such as in loops.

Modern compilers support the `register` keyword, but most simply ignore it, due to the fact that the compiler will examine the code and make intelligent decisions about what data to place in registers or in memory. This is done for every variable without having to specify anything more than a compiler option, if even that.

#### 2.1.4 Static Variables

Static is a much-overloaded term in Computer Science. Earlier it was defined as pertaining to when a program is compiled. In Java, it was used to declare a method or field that existed independently of any instances of a class. In C, `static` is a keyword with two new meanings, depending on whether it is applied to a local variable, or to a global variable or function.

##### *Static Local Variables*

So far, the distinction between scope and lifetime seems somewhat unnecessary. Local variables have a lifetime approximately that of their scope, and global variables need to live for the whole execution of a program. However, if we think about all possible combinations of scope and lifetime, there are two we have not addressed. The first, a global scope but a local lifetime, is a recipe for disaster. This is basically what occurs when returning a pointer to an automatic variable, and even the com-

```

char *asctime(const struct tm *timeptr) {
    static char wday_name[7][3] = {
        "Sun", "Mon", "Tue", "Wed", "Thu", "Fri", "Sat"
    };
    static char mon_name[12][3] = {
        "Jan", "Feb", "Mar", "Apr", "May", "Jun",
        "Jul", "Aug", "Sep", "Oct", "Nov", "Dec"
    };
    static char result[26];

    sprintf(result,
            "%3s%3s%3d.%2d:%.2d%2d\n",
            wday_name[timeptr->tm_wday],
            mon_name[timeptr->tm_mon],
            timeptr->tm_mday, timeptr->tm_hour,
            timeptr->tm_min, timeptr->tm_sec,
            1900 + timeptr->tm_year);

    return result;
}

```

**Listing 2.3:** Unlike automatic variables, pointers to static locals can be used as return values.

piler warned that was a bad idea. The other combination is a variable with a local scope but a global lifetime. This combination would imply that the variable could be used only within the block it was declared in but would retain its value between function invocations. Such a variable type might be useful as a way to eliminate the need for a global variable to fulfill this role. Anytime a global variable can be eliminated is usually a good thing.

By declaring a local variable `static`, the variable now will keep its value between function invocations, unlike an automatic variable. Interestingly, an implication of this is that static local variables can safely have pointers to them as return values. Listing 2.3 shows the `man` page for the `asctime()` function, which builds the string in a static local variable. The advantage to this is that the function can handle the allocation but does not require the caller to use `free()` as if `malloc()` had been used.

```

#include <string.h>
#include <stdio.h>

int main() {
    char str[] = "The\u00a0quick\u00a0brown";
    char *tok;

    tok = strtok(str, "\u00a0");

    while(tok != NULL) {
        printf("token:\u00a0%s\n", tok);
        tok = strtok(NULL, "\u00a0");
    }

    return 0;
}

```

**Listing 2.4:** The `strtok()` function can tokenize a string based on a set of delimiter characters.

The quintessential example of static local variables is the standard library function `strtok()`, whose behavior is somewhat atypical. Listing 2.4 shows an example of splitting a string based on the space character. The first time `strtok()` is called, the string to tokenize and the list of delimiters are passed. However, the second and subsequent times the function is invoked, the string parameter should be `NULL`. Additionally, `strtok()` is destructive to the string that is being tokenized. In order to understand this, imagine that the following string is passed to `strtok()` with space as the delimiter:

t	h	e		q	u	i	c	k		b	r	o	w	n	\0
---	---	---	--	---	---	---	---	---	--	---	---	---	---	---	----

On the first call to `strtok()`, the return value should point to a string that contains the word “the.” On the second call, where `NULL` is passed, the return value should point to “quick.” If we now examine the original string in a debugger, we would see the following:

t	h	e	\0	q	u	i	c	k	\0	b	r	o	w	n	\0
---	---	---	----	---	---	---	---	---	----	---	---	---	---	---	----

The delimiter characters have all been replaced by the null terminator! That explains why we cannot pass the original string on the second call, since even `strtok()` will stop processing when it encounters the null, thinking that the string is over. So now we have “lost” the remainder of the string. The `strtok()` function, however, remembers it for us in a static local variable. Passing `NULL` tells the function to use the saved pointer from the last call and to pick up tokenizing the string from the point it left off last.

### *Static Global Variables*

The `static` keyword, when applied to file-scope variables, takes on an entirely different meaning. Coming from an object-oriented language such as Java, the natural comparison is to consider the `static` keyword as equivalent to the `private` keyword. In C, `static` restricts use of a file-scoped variable to only the file it was declared in. No other file may import it into that file’s scope through the `extern` keyword. The variable is hidden and may not be shared.

As a note, this is also true if the `static` keyword is used to prefix a function. The function then may not be called by code from any other files in the program. In this manner, variables and functions can be kept isolated from each other, one of the original points of scoping.

#### **2.1.5 Volatile Variables**

In most cases, the compiler will decide when to put a value into a register and when to keep it in memory. In certain cases, the memory location that data is stored in is special. It could be a shared variable that multiple threads or processes (see Chapter 9) are using. It could be a special memory location that the Operating System maps to a piece of hardware. In these cases, the value of the variable in the register may not match the value that is in memory. Declaring a variable as `volatile` tells the C compiler to always reload the variable’s value from memory before using it. In a way, this could be thought of as the opposite of the `register` keyword.

The `volatile` keyword is somewhat rare to see in normal desktop applications, but it is useful in systems software when interacting with hardware devices. Though you may not encounter it much, it is nonetheless important to remember for that one time you might need it (or more if you develop low-level programs).

## 2.2 Summary Table

The following table summarizes the scope and lifetimes that C provides:

	Scope	Lifetime
<b>Automatic</b>	The block it is defined in	The life of the function
<b>Global</b>	The entire file plus any files that import it using <code>extern</code>	The life of the program
<b>Static Global</b>	The entire file, but may not be imported into any other file using <code>#1extern</code>	The life of the program
<b>Static Local</b>	The block it is defined in	The life of the program

## 2.3 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Dynamic** While a program is being run.

**Lifetime** The time from which a particular memory location is allocated until it is deallocated.

**Local variable** A variable with a scope limited to a function or smaller region.

**Scope** The region in a program where a symbol is legal and meaningful.

**Shadowing** A variable in an inner scope with the same name as one in an enclosing scope. The innermost declaration is the variable that is used; it shadows the outermost declaration.

**Static** While a program is being compiled.

**Symbol** A name that represents a constant, literal, or variable.

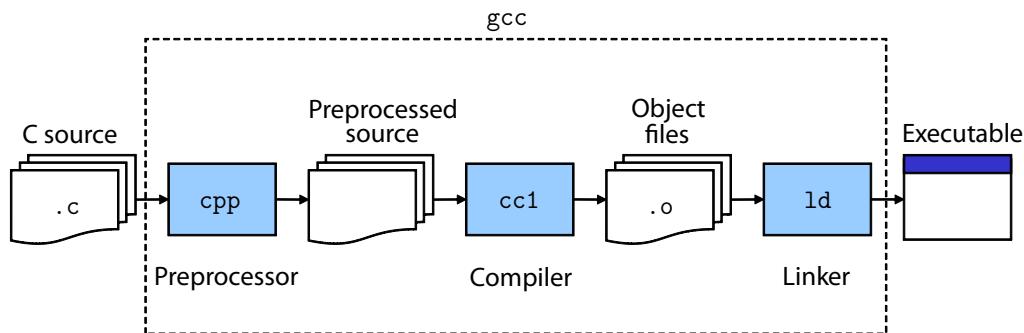
# 3 | Compiling & Linking: From Code to Executable

IN THIS CHAPTER, we begin to look at the computer as a **system**: A group of programs written by many different people, all trying to work together to accomplish useful tasks. No single component exposes the innermost workings of a system quite as well as the compiler. We begin this chapter by briefly describing the process of compilation, linking, and the makeup of executable files.

Certainly programming languages other than C can be compiled and executed, but C's original purpose was for writing Operating Systems. Thus we will discuss the C compiler as the prototypical compiler. Modern Linux systems use the `gcc` compiler, created by Richard Stallman as part of his free `GNU` system. This text will assume the `gcc` compiler, and it has been used to test the code listings found throughout the book.

## 3.1 The Stages of Compilation

Figure 3.1 shows the typical stages that a C program goes through while being compiled under the `gcc` compiler on Unix/Linux. Starting with one or more C code source files, the files are first sent to `cpp`, the **C Preprocessor**. The preprocessor is responsible for doing textual replacement on the input files by following all of the commands that begin with a `#`. After the preprocessor does its job, the fully expanded C code is passed to `cc1`, the **compiler**, which is responsible for syntax checking and code generation. If there are multiple sources of code in your program (which if you are using any of the standard library code, there are), `ld`, the **linker**, needs to resolve how to find this code. Some code will be copied into the executable and some will be left out until the program is loaded. At this point, however, there is a recognizable executable file (assuming there were no errors).



**Figure 3.1:** The phases of the gcc compiler for C programs.

### 3.1.1 The Preprocessor

The two most common preprocessor directives are `#include` and `#define`. When the preprocessor encounters a `#include` command, it seeks out the file whose name comes after the include command and inserts its contents directly into the C code file. The preprocessor knows where to look for these files by the delimiters used around the filename. Here are two examples:

```
#include <stdio.h>
#include "myheader.h"
```

If `<` and `>` are used, the preprocessor looks on the include path<sup>1</sup> where the standard header files for the system are found. A **header file** contains function and data type definitions that are found outside of a particular file. If “ and ” are used instead, the local directory is searched for the named file to include.

The directive `#define` creates a macro. A **macro** is a simple or parameterized symbol that is expanded to some other text. One common use of a macro is to define a constant. For example, we might define the value of  $\pi$  in the following way:

```
#define PI 3.1415926535
```

Now we may use the symbol `PI` whenever we want the value of  $\pi$ :

```
double degrees_to_radians(double degrees) {
    return degrees * (PI / 180.0);
}
```

---

<sup>1</sup> Under Linux this is usually `/usr/include`.

We can actually parameterize our macros to allow for more generic substitutions. For instance, if we frequently wanted to know which of two numbers is larger, we could create a macro called `MAX` that does this for us:

```
#define MAX(a,b) (((a) > (b)) ? (a) : (b))
```

Notice that we do not need to put any type names in our definition. This is because the preprocessor has no understanding of code or types; it is just following simple substitution rules: Whatever is given as `a` and `b` will be inserted into the code. For instance:

```
c = MAX(3,4);
```

will become:

```
c = (((3) > (4)) ? (3) : (4));
```

However,

```
c = MAX("bob","fred");
```

will become:

```
c = (((("bob") > ("fred")) ? ("bob") : ("fred")));
```

which is not legal C syntax. The preprocessor will do anything you tell it to, but there is no guarantee that what it produces is appropriate C.

### 3.1.2 The Compiler

The compiler is a complex and oftentimes daunting piece of software. There is an entire course on it alone. (And another one at the graduate level!) Thankfully, what we need to understand about the compiler is fairly simple. The compiler takes our source code, checks to make sure it is legal syntax, and then generates machine code. This machine code may be **optimized** — meaning that the compiler had some rule by which it transformed your code into something it believed would run faster or take up less memory.

The output of the compiler is an **object file**, which is denoted by the `.o` extension when using `gcc`.<sup>2</sup> There is one object file produced per source code file. An object file contains machine code, but any function call to code that is in a different source file is left unresolved. The address the code should jump to is unknown at this stage.

---

<sup>2</sup> Visual Studio under Windows produces object files with the extension `.obj`

### 3.1.3 The Linker

Code in an executable can come from one of three places:

1. The actual source code
2. Libraries
3. Automatically generated code from the linker

The job of the linker is to assemble the code from these three places and create the final program file.

The source code is, of course, the code the programmer has written. The **libraries** are collections of code that accomplish common tasks provided by a compiler writer, system designer, or other third party.<sup>3</sup> C programs nearly always refer to code provided by the **C Standard Library**. The C Standard Library contains helpful functions to deal with input and output, files, strings, math, and other common programming tasks. For instance, we have made much use of the function `printf()` in our programs. To gain access to this function in our code, two independent steps must be done.

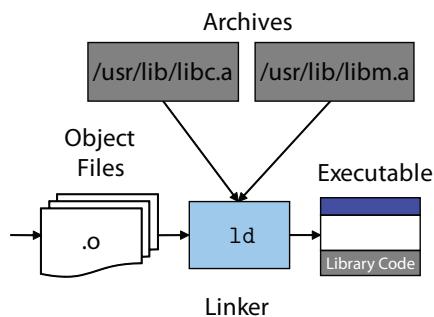
The first step is to inform the compiler that there is a function named `printf()` that takes a format string followed by a variable number of arguments. The compiler needs to know this information to do type checking (to ensure the return value and parameters match their specifications) and to generate appropriate code (which will be elaborated upon in Chapter 6). This information is provided by a function prototype declaration that resides in `<stdio.h>`.

The second step is to tell the linker how to find this code. The simplest way to assemble all three sources of code into a program is to literally put it all into one big file. This is referred to as **static linking**. It is not the only option, however. Remember that static is a term that is often used to describe the time a program is compiled. Its opposite is dynamic — while the program is running. It comes then as no surprise that we also have the option to **dynamically link** libraries, so that the code is not present in the executable program but is inserted later while the program loads (link loading) or executes (dynamic loading).

Each of these techniques has the same goal: put the code necessary for our program to run into RAM. Most common computer architectures follow the **von Neumann Architecture** where both code and data must be loaded into a main memory

---

<sup>3</sup> You can always write your own libraries as well!



**Figure 3.2:** In static linking, the linker inserts the code from library archives directly into the executable file.

before instructions can be fetched and executed. Static linking puts the code into the executable so that when the Operating System loads the program, the code is trivially there. Dynamic linking defers the loading of library code until runtime. We will now discuss the issues and trade-offs for each of these three mechanisms.

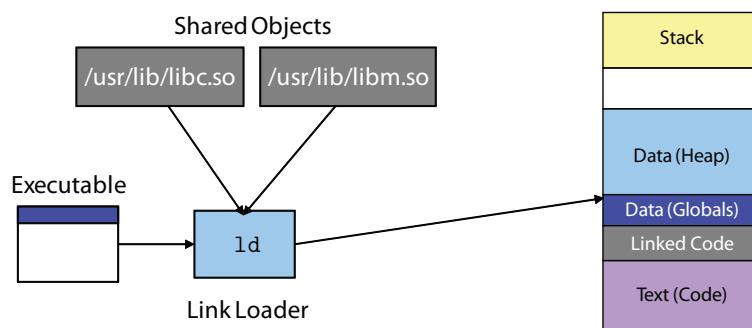
### Static Linking

Static linking occurs during compilation time. Figure 3.2 gives an overview of the linker's function during static linking. When the linker is invoked as part of the compilation process, code is taken from the libraries and inserted into the executable file directly. The code for the libraries in Unix/Linux comes from archive (.a) files. The linker reads in these files and then copies the appropriate code into the executable file, updating the targets of the appropriate `call` instructions.

The advantages of static linking are primarily simplicity and speed. Because all targets of calls are known at compile time, efficient machine code can be generated and executed. Additionally, understanding and debugging the resultant program is straightforward. A statically-linked program contains all of the code that is necessary for it to run. There are no dependencies on external files and installation can be as simple as copying the executable to a new machine.

There are two major disadvantages to static linking, however. The first is an issue of storage. If you examine a typical Unix/Linux machine, you will find hundreds of programs that all make use of the `printf()` function. If every one of these programs had the code for `printf()` statically linked into its executable, we would have megabytes of disk space being used to store just the code for `printf()`.

The second major disadvantage is exposed by examining such programs under



**Figure 3.3:** Dynamic linking resolves calls to library functions at load time.

the light of Software Engineering, where modularity — the ability to break a program up into simple, reusable pieces — is emphasized. Imagine that a bug in `printf()` is subsequently discovered on our system that is entirely statically linked. To fix our bug, we will have to find every program that uses `printf()` and then recompile them to include the fixed library. For programs whose source code we do not have, we would have to wait until the vendor releases an update, if ever.

### *Dynamic Linking*

A better approach for code that will be shared between multiple programs is to use dynamic linking. Figure 3.3 shows the process. Notice that the executable file has already been produced, and that we are about to load and execute the program. In dynamic linking, the linker is invoked twice: once at compile time and once every time the program is loaded into memory. For this reason, the linker is sometimes referred to as a **link loader**.

When the linker is invoked as part of the compiler (`ld` as a part of `gcc`, for instance) the linker knows that the program will eventually be loaded, and any library calls will be resolved. The linker then inserts some extra information into the executable file to help the linker at load time. When the linker runs again, it takes this information, makes sure the dynamically linked library is loaded into memory, and updates all the appropriate addresses to point to the proper locations.

On Unix/Linux, dynamically linked libraries are called **shared objects** and have `.so` as their extension. Windows calls them Dynamically Linked Libraries (appropriately) with the extension `.DLL`.

Dynamic linking allows for a system to have a single copy of a library stored on disk and for the library code to be copied into the process's address space at load time.

The term for the removal of duplicate information is **deduplication**. Deduplication is a common technique that is necessary for systems to work with large amounts of data. Consider an email provider like Google's GMail. If a spammer managed to get a list of one million GMail addresses, they could send a 1MB file to each of them. That would require GMail to naïvely store 1TB of data for this one email message. That obviously wouldn't be possible. Instead, GMail deduplicates email messages and attachments, storing them just once. Each user has a reference to that object that they see. The disadvantage is that GMail cannot delete the file when any individual user clicks delete, but must rather wait for everyone who received the email to delete it before reclaiming that space.

Fixing a bug in a library only requires the library to be updated. Since it is an independent file, the one-and-only copy of the library can be replaced with the updated version and all future references to the library at load time will be resolved with the updated version.

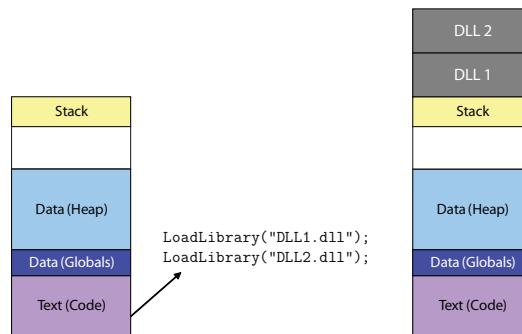
While dynamic linking solves the storage and update problems, it introduces some issues of its own. The first issue is that the extra work to resolve the addresses to their actual location is done at load time and thus slows the program's execution. The solution to this is to let the linking at load time be "lazy."

In the "lazy" approach to dynamic linking, a `call` instruction jumps to a large table in memory that contains the actual address of any function that was dynamically linked. However, this incurs extra penalties in the cost of doing the `call` instruction the first time. To make the second execution of the `call` faster, the linker may rewrite the call to jump directly to the proper location, a technique known as **back-patching**.

A second issue that arises with dynamic linking concerns versioning. A programmer may wish to use a new function included with the most recent version of a library. If the programmer distributes the program without including this library, the end user may have an older version of that library without the important function. To get around this, the developer may wish to distribute the shared libraries with the application. But if every application does this, we are not much better off than with static linking.<sup>4</sup> Unix/Linux solves this with a file naming system that includes the version number, allowing a programmer to specify a particular version if it is needed. Windows does not have an elegant solution to this problem, which is

---

<sup>4</sup> The developer can still fix bugs in either the main application or the library without having to distribute both files. However, this is somewhat minor compared to the advantage of having just one copy of a library on a system.



**Figure 3.4:** Dynamically loaded libraries are loaded programmatically and on-demand.

affectionately referred to as “**DLL Hell**.<sup>5</sup>

A third issue is related to security. A malicious programmer could name their own library the same as one the program expects to find on the system already. As long as it exports the same functions and data, they could write the library to do whatever they please. This seems to be a fatal flaw with dynamic linking. However, the problem is not related directly to linking but rather the security of the file system where the libraries are stored. Thus, the solution to this problem is to restrict the permissions of who can alter the location where shared libraries are stored on disk.

### *Dynamic Loading*

Dynamic loading is a subset of dynamic linking where the linking occurs completely on demand, usually by a program’s explicit request to the Operating System to load a library. Figure 3.4 shows an example of a Windows program making a request to load two DLLs programmatically. In Windows, a programmer can make a call to `LoadLibrary()` to ask for a particular library to be loaded by passing the name as a string parameter. Under Unix/Linux there is an analogous call, `dlopen()`.

One challenge for the Operating System in supporting dynamic loading is where to place the newly loaded code in memory. To handle this and traditional load-time linking, a portion of a process’s address space (see Chapter 5) will typically be reserved for libraries.

Care has to be taken by the programmer to handle the error condition that arises

---

<sup>5</sup> Recent versions of Windows protect core DLLs from being overwritten by older versions, but this does not solve 100% of all problem cases.

if the load fails. It is possible that a library has been deleted or not installed, and the code must be robust to not simply crash. This is one significant issue with dynamic loading, since with compile- or load-time linking the program will not compile or run without all of its dependencies available. A user will not be happy if they lose their work because in the middle of doing something the program terminates, saying that a necessary library was not found.

For this reason, core functionality of the program is probably best included using static or dynamic linking. Dynamic loading is often used for loading *plugins* such as support for additional audio or video formats in a media player or special effect plugins for an image editor. If the plugins are not present, the user may be inconvenienced, but they will likely still be able to use the program to do basic tasks.

## 3.2 Executable File Formats

After the linker runs as part of the compilation process, if there were no errors an executable file is created. The system has a format by which it expects the code and data of a program to be laid out on disk, which we call an **executable file format**.

Each system has its own file format, but the major ones that have been used are outlined here:

**a.out** (Assembler OUTput) — The oldest Unix format, but did not have adequate support for dynamic linking and thus is no longer commonly used.

**COFF** (Common Object File Format) — An older Unix format that is also no longer used, but forms the basis for some other executable file formats used today.

**PE** (Portable Executable) — The Windows executable format, which includes a COFF section as well as extensions to deal with dynamic linking and things like .NET code.

**ELF** (Executable and Linkable Format) — The modern Unix/Linux format.

**Mach-O** — The Mac osx format, based on the Mach research kernel developed at CMU in the mid 1980s.

Though a.out is not used any longer,<sup>6</sup> it serves as a good example of what types of things an executable file might need to contain. The list below lists the seven sections of an a.out executable:

1. exec header
2. text segment
3. data segment
4. text relocations
5. data relocations
6. symbol table
7. string table

The *exec header* contains the size information of each of the other sections as a fixed-size chunk. If we were to define a structure in C to hold this header, it would look like Listing 3.1.

```
struct exec {
    unsigned long    a_midmag; //magic number
    unsigned long    a_text;
    unsigned long    a_data;
    unsigned long    a_bss;
    unsigned long    a_syms;
    unsigned long    a_entry;
    unsigned long    a_trsize;
    unsigned long    a_drsize;
};
```

**Listing 3.1:** The a.out header section.

The *magic number* is an identifying sequence of bytes that indicates the filetype. We saw something similar with ID3 tags, as they all began with the string “TAG”. Word documents begin with “DOC”. The loader knows to interpret this file as an a.out format executable by the value of this magic number.

The *text segment* contains the program’s instructions and the *data segment* contains initialized static data such as the global variables. The header also contains

---

<sup>6</sup> When using gcc without the -o option, you will notice it produces a file named a.out, but this file, somewhat confusingly, is actually in ELF format.

```
public class StringTable {
    public static void main(String[] args) {
        String s = "String\u00dfliteral";

        if(s == "String\u00dfliteral") {
            System.out.println("Equal");
        }
    }
}
```

**Listing 3.2:** String literals are often deduplicated.

the size of the *BSS section* which tells the loader to reserve a portion of the address space for static data initialized to zero.<sup>7</sup> Since this data is initialized to zero, it does not take up space in the executable file, and thus only appears as a value in the header. The two relocation sections allow for the linker to update where external or relocatable symbols are defined (i.e., what addresses they live at).

The *symbol table* contains information about internal and external functions and data, including their names. Since the linker may need to look things up in this table, we want random access of the symbol table to be quick. The quickest way to look something up is to do so by index, as with an array. For this to work, however, we need each record to be a fixed size. Since strings can be variable length, we want some way to have fixed-sized records that contain variable-sized data. To accomplish this, we split the table into two parts. The symbol table with fixed-sized entries, and a *string table* that contains only the strings. Each record in the symbol table will contain a pointer to the appropriate string in the string table.

The string table will also contain any string literals that appear in the program's source code. This is another example of deduplication. In Listing 3.2, we see a common beginner's mistake in Java. The `==` operator tests for equality, but when applied to objects the equality it tests for is that the two objects live at the same address. Obviously, we wanted to use the `.equals` method. But if we compile and run this, what would we see? The output is:

Equal

---

<sup>7</sup> A reputable link on the meaning of bss states that it is “Block Started by Symbol.” See <http://www.faqs.org/faqs/unix-faq/faq/part1/section-3.html>.



Why did we get the right answer? Java class files are executable files too. And to save both storage space and network bandwidth when transferred, Java deduplicates string literals during compilation. When the class file is loaded into memory, the string literal is only loaded once. Thus the references compare as equal since both are pointing to the same object in memory. One small change and it will break. If we change the initialization of `s` to:

```
String s = new String("String\u00dfliteral");
```

We are now constructing a new object in memory that will live in a different memory location and our comparison will fail.

This is not specific to Java. In C, we have similar concerns. If we declare a string variable as:

```
char s[] = "String\u00dfliteral";
```

we get a variable `s` that points to the string literal. It is unsafe (and generally a compiler warning) to modify the string literal by doing something like `s[0] = 's';`. This is prevented in Java because `String` objects are immutable.

### 3.3 Linking in Java

When a Java program is compiled, an executable file called a `.class` file is produced. This has the standard parts including code, data, string table, and references to other classes that it depends on. One interesting thing to note is that Java has no concept of static linking. All references to code that lives in other class files are resolved while the program runs (dynamic loading).

A natural thing to assume then would be that the `import` keyword indicates that you want to link against a particular package. However, this is not the case. The `import` keyword simply indicates to the compiler the fully qualified path to a particular object, such as `ArrayList` being in `java.util`. This is mostly a scope issue, but it does enable the compiler to produce the appropriate name for the imported class, which is then, in turn, used by the dynamic class loader to find it.

### 3.4 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Back-Patching** Rewriting a jump to a dynamically loaded library so that it jumps directly to the code rather than via an intermediate table.

**Compiler** A tool for converting a high-level language (such as C) to a low-level language such as machine code.

**Deduplication** The elimination of duplication, often by storing something once and replacing each copy with a reference or link to it.

**Executable File** or simply “Executable” — A file containing the code and data necessary to run a program.

**Header File** A file containing function and data type definitions that refer to code outside of a given file.

**Library** Collections of code that accomplish common tasks provided by a compiler writer, system designer, or other third party.

**Linker** A tool for combining multiple sources of code and data into a single executable or program.

**Loader** The portion of the Operating System that is responsible for transferring code and data into memory.

**Macro** A simple or parameterized symbol that is expanded to some other text.

**Object File** A file containing machine code, but calls to functions that are outside of the particular source file it was generated from are unresolved.

**Preprocessor** The program responsible for expanding macros.

**von Neumann Architecture** A commonly-used computer architecture where the code and data must both be resident in main memory (RAM) in order to run a program.

# 4 | Function Pointers

BOTH CODE AND DATA live in memory on a computer. We have seen how it is possible to refer to a piece of data both by name and by its address. We called a variable containing such an address a *pointer*. But in addition to being able to point to data, we can also create pointers to functions. A **function pointer** is a pointer that points to the address of loaded code. An example of function pointers is given in Listing 4.1.

This example simply displays “3” upon the console. The interesting thing to note is that there is no direct call to `f()` (there is no `f(...)` in the code), but there is a call to `g()` which seems to have no body. However, we do use `f()` as the right-hand side of an assignment to `g`. The odd declaration of `g` makes it look like a function prototype. However, with experience, it becomes apparent that it is a function pointer because of the fairly unique syntax of having the asterisk inside of the parentheses. The function pointer `g` is now pointing to the location where the function `f()` lives in memory. We can now call `f()` directly as we always could by saying `f(3)`, or we can dereference the pointer `g`.

Remember that with arrays, the name of an array is a pointer to the beginning of that array. There is no need to use the dereference operator (\*) because the subscript notation ([ and ]) does the dereference automatically. The same is true for functions. The name of the function is a pointer to that function. We do not need to dereference it with a \* because the ( and ) do it automatically. Thus as the argument to `printf()`, `g(3)` dereferences the pointer `g` to obtain the actual function `f()` and calls it with the parameter 3.

## 4.1 Function Pointer Declarations

The most complicated part of function pointers is their declaration. Care needs to be taken to distinguish the fact that we have a function pointer, rather than a

```
#include <stdio.h>

int f(int x) {
    return x;
}

int main() {
    int (*g)(int x);

    g = f;
    printf("%d\n", g(3));
    return 0;
}
```

**Listing 4.1:** Function pointer example.

function that returns a pointer. For example:

```
int *f();
```

means that we have a function `f()` that returns a pointer to an integer.

```
int (*f)();
```

means we have a function pointer that can point to any function that has an empty parameter list and returns an integer. The difference is in the parenthesization.

```
int *(*f)();
```

means we have a function pointer that can point to any function that has an empty parameter list and returns a pointer to an integer. If we forget the parentheses:

```
int **f();
```

we are declaring a function that takes no parameters but returns a pointer to an integer pointer.

Why do function pointer declarations need to be so hard? The answer lies in code generation. To correctly set up the arguments to the function, the compiler needs to know exactly how many are required and what size (indicated by the type) they are. This means that there is only one correct way to have declared `g` in Listing 4.1.

```
void qsort( void *base,
            size_t num,
            size_t size,
            int (*comparator)(const void *, const void *)
        );
```

**Listing 4.2:** `qsort()` definition.

## 4.2 Function Pointer Use

Function pointers primarily get used in two particular ways. Both ways depend on the ability to defer knowledge of what function we are calling until the very moment we make the call. The first use, and the most common from a high-level programmer's perspective, is to pass a function pointer as an argument to another function. The second use is to make an array of function pointers to use as a jump or call table. This is a technique that can be used to accomplish dynamic linking.

### 4.2.1 Function Pointers as Parameters

The easiest way to explain the motivation for passing a function a pointer to another function is by example. Let us imagine we are writing a function to sort some data. While we are coding our algorithm, we find some line that requires us to do a comparison. If we are comparing the primitive data types, we can do comparison by simply using the `<` or `>` operators. But what about a more complex data type, such as a structure? Is there any way we could compare them? Additionally, if our function takes the array to sort as a parameter, what type do we define it as?

The C Standard Library includes a function called `qsort()` that seeks to be a generic sorting routine that anyone can call on any array, regardless of what type of data it contains. The declaration for `qsort()` is given in Listing 4.2.

The first parameter is of type `void *`, which means it is a pointer to anything. This solves our problem of how to declare the parameter but presents a new problem. We use typed pointers because C is able to determine the size of the data at the particular address from the size of the type. But a void pointer could be pointing to anything, and thus we need to do something extra to tell `qsort()` the size of each element. We pass that size as the third parameter. The second parameter is the length of the array we want to sort, since there is no way to query the length of an array from its pointer in C.

```

#include <stdio.h>
#include <stdlib.h>
#include <string.h>

#define NUM_NAMES 5
#define MAX_LENGTH 10

int main() {
    char names[NUM_NAMES][MAX_LENGTH] = {
        "Mary", "Bob", "Fred",
        "John", "Carl"};
    int i;

    qsort(names, NUM_NAMES, MAX_LENGTH, strcmp);

    for(i=0;i<NUM_NAMES; i++) {
        printf("%d:\t%s\n", i, names[i]);
    }
    return 0;
}

```

**Listing 4.3:** `qsort()`ing strings with `strcmp()`.

The final parameter looks complex but, based on our earlier discussion, it should be evident that this is a function pointer (the \* inside the parentheses gives it away). The function we pass to `qsort()` should return an integer, and take two parameters that will point to two items in our array which this function is supposed to compare. The return values for `comparator` need to be handled as:

$$\text{comparator} = \begin{cases} < 0, & \text{if the first parameter is } < \text{ the second} \\ 0, & \text{if they are equal} \\ > 0, & \text{if the first parameter is } > \text{ the second} \end{cases}$$

This rule might remind us of the return values for `strcmp()`. In fact, `strcmp()` makes an obvious choice for sorting an array of strings. The only requirement is that all the strings need to be a fixed size for this to work. Since sorting requires swapping elements, `qsort()` must be told as a parameter the size of the elements in the array in order to rearrange the array elements correctly. Listing 4.3 gives an

example. When we run this code, we get the expected output indicating a successful sort:

```
0: Bob    1: Carl    2: Fred    3: John    4: Mary
```

One interesting thing to note is that this code compiles with a warning:

```
(11) thot $ gcc qs.c
qs.c: In function `main':
qs.c:15: warning: passing arg 4 of `qsort' from incompatible pointer type
```

This message is telling us that there is something wrong about passing `strcmp()` to `qsort()`. If we look at the formal declaration of `strcmp()`:

```
int strcmp(const char *str1, const char *str2);
```

we see the parameters are declared as `char *` rather than the `void *` that the function pointer was declared as in Listing 4.2. This is one warning that is all right to ignore. In the old days of C, there was no special `void *` type, and a `char *` was used to point to any type when necessary.

To do something more complicated like sorting an array of structures, we will need to write our own comparator function. Listing 4.4 shows an example. The program sorts the structures first by age and then by name if there is a “tie.” The output is the following:

```
18: Bob    18: Fred    20: John    20: Mary    21: Carl
```

#### 4.2.2 Call/Jump Tables

A **call table**, or **jump table**, is basically an array of function pointers. It can be indexed by a variety of things, for instance, by a choice from a menu to decide what function to call. Compilers may generate such a table to implement `switch` statements.

Of specific interest to us is how a linker might use a call table to do dynamic linking. Each dynamically linked library exports some public functions. If we give these functions a number, referred to as an **ordinal**, we can use this ordinal to index a table created by the linker in order to find a specific function’s implementation in memory.

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>

#define NUM_STUDENTS 5

struct student {
    int age;
    char *name;
};

int mycmp(const void *a, const void *b) {
    struct student *s1 = (struct student *)a;
    struct student *s2 = (struct student *)b;

    if(s1->age < s2->age)
        return -1;
    else if(s1->age > s2->age)
        return 1;
    else
        return strcmp(s1->name, s2->name);
}

int main() {
    struct student s[NUM_STUDENTS] = {
        {20, "Mary"}, {18, "Bob"}, {18, "Fred"},
        {20, "John"}, {21, "Carl"}};
    int i;

    qsort(s, NUM_STUDENTS,
          sizeof(struct student), mycmp);

    for(i=0; i<NUM_STUDENTS; i++) {
        printf("%d: %s\n", s[i].age, s[i].name);
    }
    return 0;
}
```

**Listing 4.4:** qsort()ing structures with a custom comparator

## 4.3 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Call Table** or Jump Table — A table of function pointers.

**Function Pointer** A pointer that points to code rather than data.

**Ordinal** A number used to refer to a particular function in a dynamically linked library.

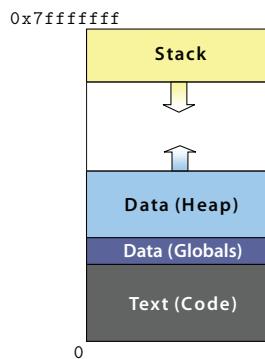
## 5 | Processes & Address Spaces

AFTER THE LOADER has done its job, the program is now occupying space in memory. The program in memory is referred to as a **process**. A process is the Operating System's way of managing a running program and keeping it isolated from other processes on the system. Associated with each process is a range of legal addresses that the program may reference. These addresses form what is known as an **address space**. To make sure that a programmer does not interfere with any other running process (either by accident or maliciously), the Operating System needs to ensure that one program may not change the code or data of another without explicit permission.

One way to solve this problem of protection is to have the computer enforce strict limits on the range that an address in a process may take. At each instruction that references a memory address, that address is checked against this legal range to ensure that the instruction is only affecting the code or data in that process. However, this incurs a performance penalty since the CPU has to do extra checking every time there is a memory operation.

Modern Operating Systems take a different approach. Through a technique referred to as **Virtual Memory**, a process is given the illusion that it is the only one running on the computer. This means that its address space can be all of the memory the process can reference in the native machine word size. On a 32-bit machine with 32-bit pointers, a process can pretend to have all  $2^{32} = 4$  gigabytes to itself. Of course, even the most expensive high-end computers do not have 4GB of memory per process of physical RAM, so the Operating System needs to work some magic to make this happen. It makes use of the hard disk, keeping unneeded portions of your program there until they need to be reloaded into physical memory.

Figure 5.1 shows a diagram of what a typical process's address space contains. The loader has placed all of the code and global variables into the low addresses. But this does not take up all of the space. We also need some dynamic storage



**Figure 5.1:** A process has code and data in its address space.

for function invocations and dynamically generated data. Functions will keep their associated storage in a dynamically managed section called the **stack**. Other dynamic data, which needs to have a lifetime beyond that of a single function call, will be placed in a structure called a **heap**.<sup>1</sup>

Notice that the stack and the heap grow toward each other. In a fixed-sized region, the best way to accommodate two variable sized regions is to have them expand toward each other from the ends. This makes sense for managing program memory as well, since programs that use a large amount of heap space likely will not use much of the stack, and vice versa. Chapter 6 and Chapter 7 will discuss the stack and the heap at more length.

## 5.1 Pages

Memory management by the Operating System is done at a chunk size known as a **page**. A page's size depends on the particulars of a system, but a common size is 4 kilobytes. The Operating System looks at what chunks you are using and those that you do not have allocated. It is clear that smaller programs will have large chunks of unallocated space in the region between the stack and the heap. These unallocated pages do not need to take up physical memory.

---

<sup>1</sup> Unfortunately, the term “heap” has two unrelated meanings in Computer Science. In this context, we mean the portion of the address space managed for dynamic data. It can also refer to a particular data structure that maintains sorted order between inserts and deletes, and is commonly used to implement a priority queue.

Since the Operating System manages the pages of the system, it can do some tricks to make dynamic libraries even more convenient. We already motivated dynamic libraries by saying they can save disk storage space by keeping common code in only one place. However, from our picture of process loading, it would appear that every shared library is copied into every process's address space, wasting RAM by having multiple copies of code loaded into memory. Since physical memory is an even more precious resource than disk storage, this seems to be less of a benefit than we initially thought. The good news is that the Operating System can share the pages in which the libraries are loaded with every process that needs to access them. This way, only one copy of the library's code will be loaded in physical memory, but every program can use the code.

## 5.2 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Address Space** The region of memory associated with a process for its code and data.

**Page** The unit of allocation and management at the Operating System level.

**Process** A program in memory.

**Virtual Memory** The mechanism by which a process appears to have the memory of the computer to itself.

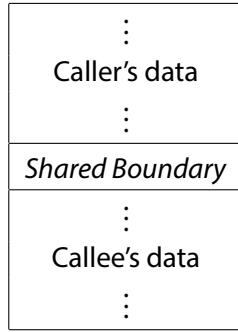
# 6 | Stacks, Calling Conventions, & Activation Records

THE FIRST TYPE of dynamic data we will deal with is local variables. Local variables are associated with function invocations, and since the compiler does not necessarily know how many times a function might be called in a program, there is no way to predict statically how much memory a program needs. As such, a portion of memory needs to be dedicated to holding the local variables and other data associated with function calls. Data in this region should ideally be created on a function's call and destroyed at a function's return. If we look at the effect of having a function call other functions, we see that the local variables of the most recent function call are the ones that are most important. In other words, the local variables created last are used, and are the first ones to be destroyed; those created earlier can be ignored. This behavior is reminiscent of a stack.

To create a stack we need some dedicated storage space and a means to indicate where the top of the stack lives. There is a large amount of unused memory in the address space after loading the code and global data. That unused space can be used for the stack. Since practically every program will be written with functions and local variables, the architecture will usually have a register, called the **stack pointer**, dedicated to storing the top of the stack.

With storage and a stack pointer, we can make great strides in managing the dynamic memory needs of functions. When a function is compiled, the compiler can figure out how many bytes are needed to store all of the local variables in that function and then write an instruction that adjusts the stack pointer by that much on every function call. When we want to deallocate that memory on function return, we could adjust the stack pointer in the opposite direction.

Other than local variables, what information might we want to store on the stack? Since the concept of a stack is so intrinsically linked with function calls, it



**Figure 6.1:** The shared boundary between two functions is a logical place to set up parameters.

seems to make sense that the **return address** of the function should be stored on the stack as well. If we examine the caller/callee relationship, we see that their data will be in adjacent locations on the stack. Figure 6.1 shows the two stack entries and the boundary between them. If the caller function needs to set up arguments for the callee, this boundary seems a natural place to pass them.

On machines with many registers, some registers may be designated for temporary values in the computation of complicated expressions. These temporary registers may be free for any function to use, and thus, if a function wants a particular register to be saved across an intervening call to another function, the calling function must save it on the stack. This is referred to as a **caller-saved** register. Other registers may be counted upon to retain their values across intervening function calls, and if a called function wants to use them, it is responsible for saving it on the stack and restoring them before the function returns. These are **callee-saved** registers. In general, the stack is used to save any architectural state that needs to be preserved.

We now have the following pieces of data that need to be on the stack:

1. The local variables — including temporary storage, such as for saving registers
2. The return address
3. The parameters

All of these together will form a function's **activation record** (sometimes called a **frame**). Not every system will have all of these as part of an activation record. How

	MIPS	x86
<b>Arguments:</b>	First 4 in %a0–%a3, remainder on stack	Generally all on stack
<b>Return values:</b>	%v0–%v1	%eax
<b>Caller-saved registers:</b>	%t0–%t9	%eax, %ecx, & %edx
<b>Callee-saved registers:</b>	%s0–%s9	Usually none

**Figure 6.2:** A comparison of the calling conventions of MIPS and x86

parameters are passed, be it on the stack or in registers, is part of a “contract” the system abides by, called the **calling convention**.

## 6.1 Calling Convention

For one function to call another, they must first agree on the particulars of where parameters will be passed and where return values will be stored. Each system will have a different calling convention depending on its particular details, such as the number of general purpose registers. Figure 6.2 summarizes two different architectures’ calling conventions.

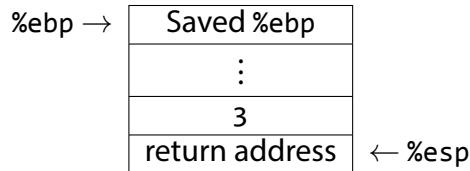
Beyond an architecture or system’s designers specifying the details of a calling convention, individual programming languages may specify how some data is to be exchanged. Let us start with a look at a C program written for an x86 processor running Linux and the assembler output generated by `gcc`. Figure 6.3 shows the source and output for a `main()` function that calls a function `f()` with one parameter.

If we begin to trace the first five instructions of the `main()` function, we notice that they are all related to the management of the two stack-related registers, `%esp`, the stack pointer, and `%ebp`, the base pointer. In this context, `%ebp` is being used to keep track of the beginning of the activation record, and is sometimes referred to as the **frame pointer**. The `andl $-16, %esp` instruction is a way to do data **alignment**. On many architectures, it is faster to access data that begins at addresses that are multiples of some power of two because of memory fetches and caches. The `subl $16, %esp` allocates some memory on the stack for the parameter to `f()`. This is a subtract because the stack starts from a high address and grows towards lower addresses. All figures in this book have been drawn with that detail in mind. Thus, a **push** is a subtract and a **pop** is an add. Whenever a function call is made,

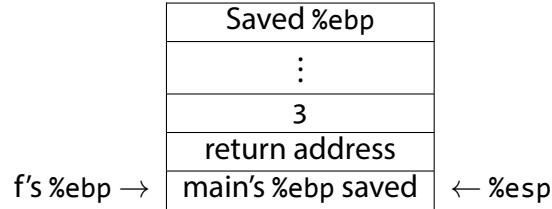
<pre>#include &lt;stdio.h&gt;  int f(int x) {     return x; }  int main() {     int y;      y = f(3);      return 0; }</pre>	<pre>f:    pushl %ebp       movl %esp, %ebp       movl 8(%ebp), %eax       leave       ret  main: pushl %ebp       movl %esp, %ebp       subl \$8, %esp       andl \$-16, %esp       subl \$16, %esp       movl \$3, (%esp)       call f       movl %eax, -4(%ebp)       movl \$0, %eax       leave       ret</pre>
--	---

**Figure 6.3:** A function with one parameter.

the `call` instruction pushes the return address onto the top of the stack and then jumps to the subroutine. Below is a diagram of the current state of the stack (the dots indicate extra space left unused due to alignment).



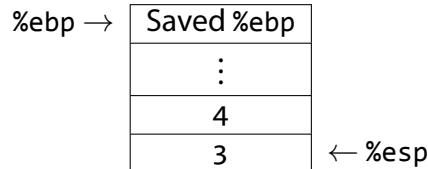
When we enter into `f()`, we again set up an activation record by first saving `main()`'s frame pointer and then adjusting the frame pointer to point to the bottom of `f()`'s activation record.



At this point, `f()` wants the value of its parameter, which is stored above the current base pointer. The instruction `movl 8(%ebp), %eax` accesses it, which is two words (8 bytes) away. There is no need for any calculation in this simple function, and the return value is directly placed into `%eax`. The `leave` instruction restores `main()`'s frame by popping the activation record, setting `%esp` to `%ebp`, and then restoring `%ebp` to the old value. The `ret` instruction pops the return address off the stack and returns back to `main()`.

While a lot is going on in these few instructions, everything seems fairly straightforward. An activation record is created for `main()`, the parameter is placed on the stack, and the function `f()` is called, which sets up its own frame and does the work. Deallocating the frame is simple due to the `leave` and `ret` instructions.

Let us now look at the code in Figure 6.4, which illustrates a function that takes two parameters. Much the same is going on in Figure 6.3; the only difference is we push two parameters instead of one. The stack looks like:



<pre>#include &lt;stdio.h&gt;  int f(int x, int y) {     return x+y; }  int main() {     int y;      y = f(3, 4);      return 0; }</pre>	<pre>f:    pushl %ebp        movl %esp, %ebp        movl 12(%ebp), %eax        addl 8(%ebp), %eax        leave        ret  main: pushl %ebp       movl %esp, %ebp       subl \$8, %esp       andl \$-16, %esp       subl \$16, %esp       movl \$4, 4(%esp)       movl \$3, (%esp)       call f       movl %eax, 4(%esp)       movl \$0, %eax       leave       ret</pre>
--	---

**Figure 6.4:** A function with two parameters.

when we make the call to `f()`.

Something that is peculiar here is that the two parameters seem somehow backwards. They have been pushed in reverse order from the way they are written. This particular quirk in the C calling convention is introduced entirely for the support of one very common function:

```
int printf(const char *format, ...);
```

The `printf()` function takes a variable number of arguments. As such, there is no way for it to know how far down the stack to look for data. The only way it might be able to know is to look inside the format string, which is what it does. The number of scan codes inside the format string tells `printf()` how many arguments to expect. Because of this, you must always be sure to pass as many arguments as you have scan codes, or `printf()` might start printing other values off the stack. By pushing the arguments in reverse order, the closest argument to `%ebp` is the format string, always at a fixed offset in a predictable location. This is why the C calling convention pushes arguments in reverse order.

## 6.2 Variadic Functions

A **variadic function** is a function that takes a variable number of arguments. In C, a variadic function declaration is easy to recognize due to the ellipses (...) in the function declaration. The `stdarg.h` header file provides several macros to help deal with the parameter list. Listing 6.1 shows an example of a variadic function that turns its parameter list into an array. The `va_start` macro takes the last required parameter and initializes the `va_list` variable `ap`. The `va_arg` macro facilitates advancing the pointer to walk the stack by the size of the appropriate type (which must be known by the function). It casts the data at that address back to the appropriate type and advances the pointer automatically.

Java and other languages with built-in support for dynamic arrays often expose the parameters of a variadic function as an array. For example:

```
public static void printArray(Object ... objects) {
    for (Object o : objects)
        System.out.println(o);
}

printArray(3, 4, "abc");
```

displays its parameters on the screen, one per line.



```
#include <stdio.h>
#include <stdarg.h>

int *makearray(int a, ...) {
    va_list ap;
    int *array = (int *)malloc(MAXSIZE*sizeof(int));
    int argno = 0;
    va_start(ap, a);
    while (a > 0 && argno < MAXSIZE)
    {
        array[argno++] = a;
        a = va_arg(ap, int);
    }
    array[argno] = -1;
    va_end(ap);
    return array;
}

int main() {
    int *p;
    int i;
    p = makearray(1,2,3,4,-1);

    for(i=0;i<5;i++)
        printf("%d\n", p[i]);

    return 0;
}
```

**Listing 6.1:** A variadic function to turn its arguments into an array.

```

void f(char *s) {
    gets(s);
}

int main() {
    char input[30];
    f(input);
}

```

---

```

main:
    pushl %ebp
    movl %esp, %ebp
    subl $40, %esp
    andl $-16, %esp
    subl $16, %esp
    leal -40(%ebp), %eax
    movl %eax, (%esp)
    call f
    leave
    ret

```

**Figure 6.5:** A program with a buffer overrun vulnerability.

### 6.3 Buffer Overrun Vulnerabilities

The code in Figure 6.5 has a problem. The code allocates an array of thirty characters, then calls `gets()` to ask the user to enter some input. The `gets()` function takes only one parameter, and from our knowledge of C, we know that there is no way to tell from that one pointer parameter just how big the array is that was passed. This means that a malicious person could enter something much larger than 30 characters. Since `gets()` has no idea when to stop copying input into our array, it keeps going. As the assembly listing shows, the array was allocated on the stack. When `gets()` exceeds that space, it starts writing over the rest of the data on the stack, corrupting it.

If our malicious user notices that the program crashes with long input, he or she may suspect that the input is overwriting the activation record on the stack. By carefully crafting the input string, the user can deliberately overwrite specific offsets on the stack, including the return address of `gets()`. With the return address modified, when `gets()` returns, code will start executing at whatever location the user entered.

With some more careful crafting of the input string, the return address can be modified so that it points to the location in memory where the array is stored. Cleverly, the user can place machine code as the input, and when `gets()` returns, it jumps via the return address to the code that the attacker entered. If this program

had been a system service, running as a superuser on the machine, the attacker could have injected code to make themselves a superuser too. This type of attack is known as a **buffer-overrun vulnerability**.

Great care must always be taken with arrays on the stack. Overrunning the end of the buffer means writing over activation records, and possibly subjecting the system to an attack. *Always check that the destination of a copy is large enough to hold the source!*

## 6.4 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Activation Record** The local variables, parameters, and return address associated with a function's invocation.

**Alignment** Making sure that data starts at a particular address for faster access due to memory fetches or the cache.

**Callee-Saved** A piece of data (e.g., a register) that must be saved by a called function before it is modified and restored to its original value before the function returns.

**Caller-Saved** A piece of data (e.g., a register) that must be explicitly saved if it needs to be preserved across a function call.

**Calling Convention** An agreement, usually created by a system's designers, on how function calls should be implemented — specifically regarding the use of registers and the stack.

**Frame** Another name for an activation record.

**Frame Pointer** A pointer, usually in a register, that stores the address of the beginning of an activation record (frame).

**Return Address** The address of the next instruction to execute after a function call returns.

**Stack** A portion of memory managed in a last-in, first-out (LIFO) fashion.

**Stack Pointer** The architectural register that holds the top of the stack.

# 7 | Dynamic Memory Allocation Management

WHILE THE STACK is useful for maintaining data necessary to support function calls, programs also may want to perform dynamic data allocation. Dynamic allocation is necessary for data that has an unknown size at compile-time, an unknown number at compile-time, and/or whose lifetime must extend beyond that of the function that creates it. The remaining portion of our address space is devoted to the storage of this type of dynamic data, in a region called the **heap**.

As is often the case, there are many ways to track and manage the allocation of memory. There are trade-offs between ease of allocation and deallocation, whether it is done manually or it is automatic, and the speed and memory efficiency need to be considered. Also as usual, the answer to which approach is best depends on many factors.

This chapter starts by describing the major approaches to allocation and deallocation. We first describe the two major ways to track memory allocation. The first is a **bitmap** — an array of bits, one per allocated chunk of memory — that indicates whether or not the corresponding chunk has been allocated. The second management data structure is a **linked list** that stores contiguous regions of free or allocated memory. A third technique, the **Buddy Allocator** attempts to reduce wasted space from many allocations. The chapter also describes an example implementation of `malloc()`, the C Standard Library mechanism for dynamic memory allocation.

## 7.1 Allocation

The two operations we will be concerned with are *allocation*, the request for memory of a particular size, and *deallocation*, returning allocated memory back to the system

for subsequent allocations to use. The use of the stack for function calls led us to create activation records upon function invocation and to remove them from the stack on function return. In essence, we were allocating and deallocating the activation records at runtime — the very operations we are attempting to define for the heap.

The question is then, why is the stack insufficient and what is different about the heap? As the name implies, the stack is managed as a FIFO with allocation corresponding to a push operation and deallocation corresponding to pop. This worked for function calls because the most recently called function, the one whose activation record was allocated most recently and lives at the top of the stack, is the one that returns first. Deallocations always occur in the opposite order from the allocations. New allocations always occur at the top of the stack, and with the stack growing from higher addresses to lower ones by convention, this means that all space above the stack pointer is in-use. All space at lower addresses is free or not part of the stack.

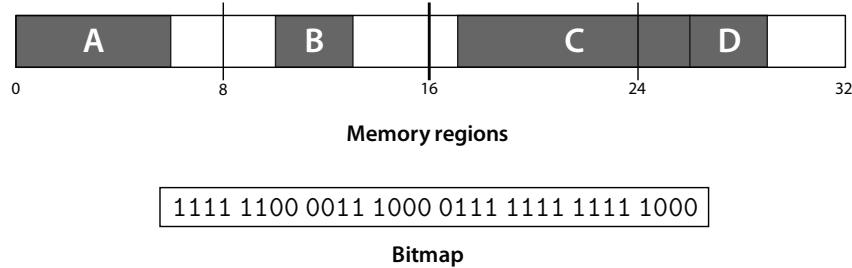
Thus, allocation is simply moving the top of the stack, and deallocation is moving it back. But for objects whose lifetime is not limited to the activation of a particular function, this order requirement is too restrictive. We would like to be able to allocate objects A, B, and C, and then deallocate object B. This leaves an unallocated region in the middle that we may wish to reuse to allocate object D.

In this section, we are considering this more general case of allocation: the possibility that we have free space in between allocated spaces. We need to track that space and to allocate from it. With that in mind, the simple dividing line between free and used space that the stack pointer represented is insufficient and we need to use a more flexible scheme.

### 7.1.1 Allocation Tracking Using Bitmaps

What we wish to do is to ask for an arbitrary piece of memory, is this space in use or is it free? At the heart of it, this question is answered via a single boolean value that is mutually exclusive: yes, it is being used or no, it is free. That information can be stored in a single bit per piece of memory. For an entire region of memory, we can combine all of the individual bits into a single array of bits called a **bitmap**.

When a modern computer user thinks of the term bitmap, he or she invariably recalls the image format. This is appropriate, since if we were to create a format for storing black and white images, we might devote one bit per pixel and treat the image as a large array of pixels. The same concept applies to managing memory



**Figure 7.1:** A bitmap can store whether a chunk of memory is allocated or free.

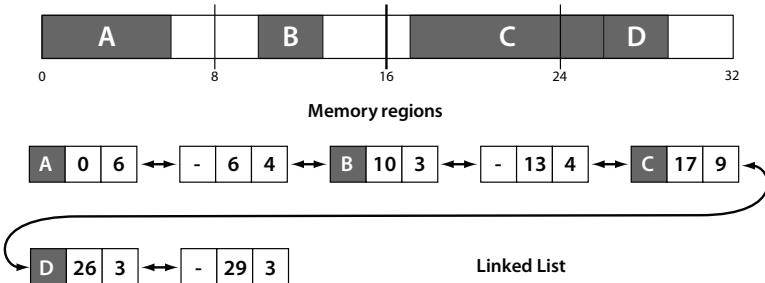
using bitmaps. We store an array of bits where a 0 indicates an unallocated chunk of memory, and a 1 indicates a chunk has been allocated for some purpose. Figure 7.1 shows a region of 32 chunks. Several parts are allocated to A, B, C, and D, with the remainder of the chunks free. Below it is the corresponding bitmap.

Bitmaps have several significant disadvantages that make them undesirable for common use. The first major drawback is the space requirement. A part of the region of memory that we are managing will have to be devoted to storing the bitmap itself. If the unit of allocation is very small, for instance a single byte, the bitmap will need to be large (one bit per byte). This means that for every eight bytes we will need one more byte for the bitmap. One out of nine bytes (11%) will be wasted in management overhead.

To reduce the size of the bitmap then, we could try to reduce the number of bits necessary by having one bit represent more space. This will result in increasing the smallest unit we can allocate. Instead of each bit tracking a single byte being allocated or free, each bit will represent a contiguous chunk of memory. Now, for a memory region of size  $S$  and a chunk size of size  $c$ , we will only need  $\frac{S}{c}$  bits.

While this will reduce the size of the bitmap, a new problem arises: It is impossible to allocate less than a single chunk, since a single bit cannot represent partial allocations — only whole ones. This leads to the problem of **internal fragmentation**, which is wasted space due to an allocation unit being bigger than our need. When we allocate a contiguous region of memory, on average the very last chunk will be half full. This could result in a large amount of wasted space over the lifetime of a program.

A second disadvantage to using bitmaps is the difficulty in finding a large enough free space to hold a given allocation. The challenge lies in how to discern that there are the right number of zeros in a row. If the empty space is at the end, the search



**Figure 7.2:** A linked list can store allocated and unallocated regions.

may be slow. Practically, it would involve a lot of bit shifting and masking. As such, using bitmaps for any significant tracking of dynamic memory allocations is unlikely, but bitmaps do often find a use in tracking disk space allocation.

### 7.1.2 Allocation Tracking Using Linked Lists

The potentially huge size of a bitmap relative to the memory region it was managing drove us to chunk memory and introduce internal fragmentation. However, there are two properties of a bitmap that might allow us to reduce its overall size while avoiding the necessity of chunks. One observation that we may make is that frequently the bitmap may contain many zeros. This would be the case when the region is new and there haven't been many allocations, when there have been a lot of deallocations, or when a region is simply much larger than the current dynamic memory needs. Noticing this, we can take a page from matrices. When a matrix has mostly zero entries, we call it **sparse**. In the CS world, we find that sparseness can apply to various data structures including our bitmap. A space-saving idea with a sparse data structure is to only store the elements that are nonzero. A linked list is a sparse data structure that supports storing a variable number of elements with dynamic inserts and deletes. We can omit the zero entries and infer that anything not in the linked list is free space.

While the choice to use a sparse data structure is likely a good one, there are two issues that prevent us from stopping here. The first is that sparseness is good for size, but we need the linked list to easily support fast allocations that come from the free space. If the information about unallocated space is not directly stored in the linked list, we must infer the necessary space is available and of the proper size. This motivates having additional linked list nodes that represent free space, but in doing so, we have lost the sparseness that made a multi-byte linked list node a reasonable

thing to store over a single bit in a bitmap.

The solution here is that our observation of sparseness in the bitmap was accurate, but it did not go far enough in describing the situation. Not only are many of the bitmap entries zeros, those zeros are also likely to be next to other zeros in long contiguous *runs*. When we have an allocation, that also exists in the bitmap as a run of ones. A key observation is that the number of ones or zeros in a particular run is actually an encoding of that region's length, albeit in unary (base 1). Unary turns out to be the worst of all bases in terms of compactness because to represent the number  $n$ , you need a string of length  $n$ . If we instead consider base two (binary), the length of a string of bits needed to represent the value  $n$  is only  $\log_2(n)$ .

This logarithmic rather than linear growth gives us a better scheme for storing the size of a particular allocation. Instead of denoting a size as a run of  $n$  ones or zeros in a bitmap, we could simply store the size in a normal variable in memory. A 32-bit integer would only be enough space to store information about 32 chunks when used as a bitmap, but because of the slow growth of logarithms, those same 32 bits can store a value of  $2^{32}$ , or about four billion! Even if we decided not to use chunks but track memory at the byte granularity, an often unreasonable choice for a bitmap, we could still represent allocations of up to 4GB in size.

The removal of runs is a simple form of compression known as **run-length encoding** or RLE. Using RLE to compress allocations and free spaces allows us to have one linked list node per allocated or free contiguous region. Figure 7.2 shows the same region and allocations as in Figure 7.1, but with the nodes of a linked list corresponding to the free and allocated spaces. The first number is the starting index of the contiguous space; the second is the size of the allocation.

Like a bitmap, the linked list needs to be stored in the same memory it is managing. With the bitmap, since the size is known ahead of time (it is simply the size of the region divided by the size of the chunk), space can be reserved before any allocations are done. The bitmap will not grow or shrink as long as the region and chunk do not change size. The linked list, on the other hand, has a number of nodes that is proportional to the number of allocations and free spaces. This number changes as the region is used.

To store the linked list, we could reserve the worst-case size in advance, much as we did with the bitmap. The worst case length of the linked list would occur when there is one node per minimum unit of allocation. This could occur for a number of different scenarios, such as with a full region of individual unit-size allocations or where unit allocations are separated by unit-sized free spaces. In this case, we would have  $n$  list nodes just as we had  $n$  bits in the bitmap. However, while each entry

in the bitmap only required a single bit's worth of storage, how large might a list node be? We need to store the size, the start, and links for the linked list. For faster deallocation support, we probably want this to be a doubly-linked list, requiring us to have two node pointers. Assuming all of these fields are four bytes in size, we would need  $4 \times 4 = 16$  bytes or  $16 \times 8 = 128$  bits. Thus, to reserve space for the worst case scenario, we would need  $128n$  bits where the bitmap only needed  $n$  bits. The linked list is 128 times the size!

This is horrible and we may wonder how we began by trying to reduce the size of a data structure but ended up making it 128 times worse. The answer is in the worst case scenarios. They were the worse case because they eliminated the runs that were the bases of our compression. When our assumptions are not valid, our end result is likely to come out worse. The good news is that our assumptions are valid in the *typical* case. Such degenerate linked lists are not likely to result from the normal use of dynamic memory.

While we have convinced ourselves the linked list is still a valid approach, we still need a good solution for where to store the elements of the linked list. Reserving the space in advance is not feasible. A better solution might be to think of the memory-tracking data structure as a “tax” on the region of memory we are tracking. For a bitmap, we pay a fixed-rate tax off the top — before we have even used the region. For paying that tax, we never have to pay again. For the linked list however, we could instead pay a tax on each allocation. Every time we get a request for dynamic memory, we could allocate a bit extra to store the newly-required list node. For instance, if we get a request for 100 bytes, we actually allocate 116 bytes and use the additional space to store one of the nodes we described above.

### 7.1.3 Allocation Algorithms

Searches through the linked list are necessary to find a region to satisfy an allocation request, but the integer size comparisons are easier for the computer compared to the bit matching needed for bitmaps. Whichever technique we use, when an allocation request is made there may be many free spots that could accommodate the request. Which one should we choose? Below is a list of various algorithms from which we could pick:

**First fit** Find the first free block, starting from the beginning, that can accommodate the request.

**Next fit** Find the first free block that can accommodate the request, starting where

the last search left off, wrapping back to the beginning if necessary.

**Best fit** Find the free block that is closest in size to the request.

**Worst fit** Find the free block with the most left over after fulfilling the allocation request.

**Quick fit** Keep several lists of free blocks of common sizes, allocate from the list that nearest matches the request.

**First fit** is the simplest of the algorithms, but suffers from unnecessary repeated traversals of the linked list. Each time first fit runs, it starts at the beginning of the list and must search over space that is unlikely to have any free spaces due to having allocated from the beginning each prior time the function was called. To avoid this cost, we could remember where the last allocation happened and start the search from there. This modification of first fit is called **next fit**.

Both of these algorithms take the first free block they find, which may not be ideal. This strategy may leave uselessly small blocks or prevent a later request from being fulfilled because a large free block was split when a smaller free spot elsewhere in the list might have been a better fit. This wasted space between allocations is **external fragmentation**.

To avoid external fragmentation, we may wish to search for the best fit. The **best fit** algorithm searches the entire list looking for the free space that is closest in size to the request. This means that we will never stop a large future request from being fulfilled because we took a large block and split it unnecessarily. However, this algorithm can turn out to be poor in actual usage because we end up with many uselessly small leftovers when the free space is just slightly larger than the request. This is guaranteed to be as small as possible whenever an exact fit is not found, due to the difference between the free space and the allocation being minimized by our definition of “best”. Additionally, best fit is slow because we must go through the entire linked list, unless we are lucky enough to find a perfect fit.

To avoid having many small pieces remain, we could do the exact opposite from best fit, and find the **worst fit** for a request. This should leave a free chunk after allocation that remains usefully large. As with best fit, the entire list must be searched to find the worst fit, resulting in poor runtime performance. Unlike best fit, which could stop early upon finding a perfect fit, the worst fit cannot be known without examining every free chunk. Despite our intuition, simulation of this algorithm reveals that it is not very good in practice. An insight into why is that

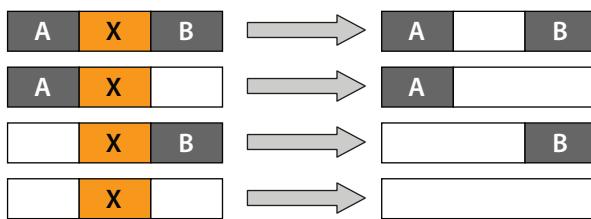
after several allocations, all of the free chunks are around the same, small size. This is bad for big requests and makes looking through the whole list useless as every free chunk is about equal size.

An alternative to the search-based algorithms, **quick fit** acknowledges that most allocations come clustered in certain sizes. To support these common sizes, quick fit uses several lists of free spaces, with each list containing blocks of a predetermined size. When an allocation request is made, quick fit looks at the list most appropriate for the request. Performance is good because searching is eliminated: With a fixed number of lists, determining the right list takes constant time. Leftover space (internal fragmentation) can be bounded since an appropriately-sized piece of memory is allocated. If the lists were selected to match the needs of the program making the allocation, this would leave very little wasted space. Additionally, that leftover space should not harm future large requests because the large requests would be fulfilled from a different list. One issue with quick fit is the question of whether or not to coalesce free nodes on deallocation or to simply return them to their appropriate list. One solution is to provide a configurable parameter to the allocator that says how many adjacent small free nodes are allowed to exist before they are collapsed into one. This ensures large unallocated regions as well as enough of the more common smaller regions.

The two likely “winners” of the allocation battle are next fit and quick fit. They both avoid searching the entire list yet manage to fulfill requests and mostly avoid fragmentation. The GNU glibc implementation of `malloc()` uses a hybrid approach that combines a quick fit scheme with best fit. The writers claim that while it is not the theoretically best performing `malloc()`, it is consistently good in practice.

## 7.2 Deallocation

The other important operation to consider is deallocation. This is where the true distinction against stack allocation is drawn. Whenever we free space on the stack, we reclaim only the most recently allocated data. The heap has no such organization, and thus deallocations may occur regardless of the original allocation order. Since the stack is completely full from bottom to top, the only bookkeeping necessary is an architectural register to store the location of the top. The heap, on the other hand, will inevitably have “holes” — free spaces from past deallocations — that will arise. Keeping track of the locations of these holes motivates the use of a data structure such as a bitmap or linked list. In this section, we look at the various approaches a



**Figure 7.3:** Coalescing free nodes on deallocation.

system might take to deallocation.

### 7.2.1 Using Linked Lists

When we allocate some memory, our linked list changes. The free node is split into two parts: the newly allocated part and the leftover free part. Eventually deallocations happen, and it is time to release a once-used region of memory. Figure 7.3 shows the four scenarios that we might find when doing a free operation. The top-most shows a region being deallocated (indicated with an 'X') that has two allocated neighboring regions; in this case, we simply mark the middle region as free. The second and third cases show when the left or right neighbor is free. In this case, we want to **coalesce** the free nodes into a single node so that we may later allocate this as one large contiguous region. The final case shows both of our neighbors being free, and thus we will need to coalesce them all.

To facilitate coalescing nodes, we may want to use a doubly linked list, which has pointers to the next node as well as the previous node. Note that we do not want to coalesce allocated nodes because we would like to be able to search the linked list for a particular allocation by name (or address, if that is what we are storing as the “name”).

### 7.2.2 Using Bitmaps

For all of the flaws of using bitmaps for dynamic memory management, deallocation of a bitmap-managed region is surprisingly simple: the appropriate bits switch from 1 (allocated) to 0 (deallocated). The beauty of this approach is that free regions are coalesced without any explicit effort. The newly freed region's zeros naturally “melt” into any neighboring free space.

### 7.2.3 Garbage Collection

Up until this point, we have been assuming that requests for deallocation have come directly from the user. Forgetting to explicitly deallocate space can lead to **memory leaks**, where a dynamically allocated region cannot be explicitly freed because all pointers to it have either been overwritten or gone out of scope. It requires the diligence of the programmer to avoid leaking memory, a tedious task that might seem to lend itself to being automated.

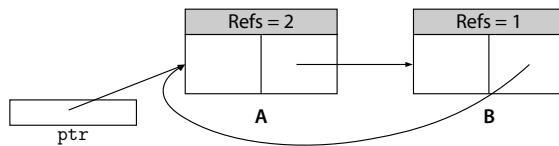
For deallocation to be done automatically, the system needs to know that a chunk of dynamically allocated memory is no longer used. A clue to how this might be determined is in the previous paragraph. If we have leaked the memory — that is, we have no valid pointers to it — then we can reclaim that space. The system is now faced with a fundamental problem: If there are no pointers to a region, how does the system itself (in Java's case, the Virtual Machine) find it? There are several approaches that might work. The JVM could keep an internal list of pointers to every object that is allocated. Another alternative is to walk the stack looking for object references.

Once all of the data items that are unreachable are discovered, the task of freeing their space, called **garbage collection**, starts. Since it takes time to find all of the “garbage,” the collection process is not usually on-demand in the same way that `free()` works in C. Except for the so-called *concurrent collectors*, garbage collectors run only when necessary — usually when the amount of free heap space has dropped below some threshold. When this threshold is hit, the program generally pauses and garbage collection begins. While there is a vast array of techniques by which to free used space, we will discuss three common strategies: *reference counting*, *in-place collectors*, and *copying collectors*.

#### *Reference Counting*

We have already determined that a dynamically-allocated object is garbage and can be collected when it has been leaked and there are no longer any valid references to the memory. Possibly the simplest way to determine this is to count valid links to an object and, when the count reaches zero, automatically free the memory. This strategy is known as reference counting and can be implemented relatively easily even in native code.

Each object needs a reference count variable associated with it. This variable is incremented or decremented as the program runs. It will need to be updated:



**Figure 7.4:** Reference counting can lead to memory leaks. If the pointer `ptr` goes out of scope, the circularly linked list should be collected. However, each object retains one pointer to the other, leading to neither having the requisite zero reference count for deallocation.

1. When a reference goes out of scope.
2. When a reference is copied (explicitly with assignment or implicitly in parameter passing, for example).

When a reference goes out of scope, the reference count on the associated object must be decremented. Copying references affects both sides of the assignment. The left-hand side (often called an *l-value*) might have been referring to an object prior to the assignment. This reference is now going to be lost from the overwrite, so the original object's reference count must be decremented. The right-hand side of the assignment (predictably called the *r-value*) is now going to have one more reference to it and the associated counter must be incremented accordingly.

When an object's reference count reaches zero, the object is garbage and can be collected. This might happen while the program is running (making it a concurrent collector) or at periodic breaks in the program's execution (a stop-the-world collector). The act of garbage collection can be as easy as freeing the object with whatever heap management operation is available.

A problem that can arise in a reference counting garbage collector is, remarkably, that it can leak memory. If a data structure has a cycle, such as in a circularly linked list as shown in Figure 7.4, there can be no way to collect the data structure. With a cycle, there is at least one reference to each object that remains even after all references from the program code are gone. Since the reference count never reaches zero, the objects are not freed. Possible solutions to this problem include detecting that the objects are part of a cycle or by using one of the other garbage collection algorithms.

### In-place Collectors

Another approach to garbage collection is via an in-place collector. The process is comprised of two phases: a *mark* phase and a *sweep* phase. During the mark phase, all of the references found on the stack are followed to find the objects to which they refer. Those objects may contain references themselves. As the algorithm traverses this graph of references, it marks each object it encounters as reachable, thus indicating it is not to be collected. When every reference that can be reached has its associated object marked, the algorithm switches to the sweep phase.

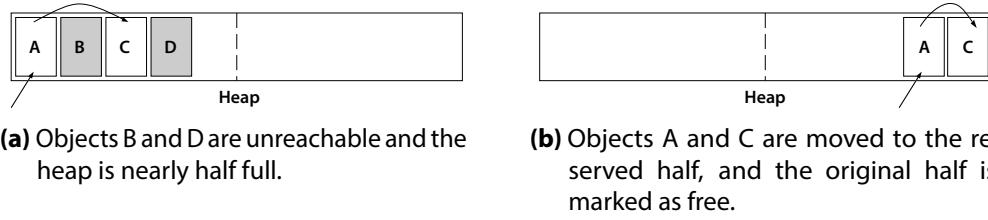
In the sweep phase, all unmarked objects are freed from the heap. All that remains are the reachable objects. When the deallocation is finished, all of the marked objects are reset to unmarked so that the process may begin all over again when the garbage collector is invoked the next time.

This *mark and sweep* approach is simple and relatively fast. It avoids cycles because encountering an object we have already seen can be detected as the object will be already marked as seen. It suffers from a significant problem, however. The newly freed space might be between objects that are still alive and remain in the heap. We now have holes that are small and scattered throughout memory, rather than a big free contiguous chunk of memory from which to allocate new objects. While there might be a significant fraction of space that is free, it might be fragmented to the point of being unusable. This is once again, external fragmentation, and was the motivation behind coalescing the adjacent free nodes in a linked list management scheme.

### Copying Collectors

To fix the fragmentation issue of the in-place collector, a garbage collector could compact the region by moving all of the objects closer together. This would constitute a third *compaction* phase and is actually unnecessary. We can avoid a third phase by combining deallocation and compaction into a single pass through the heap.

Copying garbage collectors such as the *semispace* collector typically divide the heap into two halves and copy from the full half into the reserved, empty half. Figure 7.5 shows an example. In Figure 7.5a, objects B and D have been designated unreachable and should be freed. Rather than explicitly do this freeing and be left with two small holes in memory, objects B and D are left untouched. Objects A and C are referenced and thus alive. A copying collector will move these live objects to the reserved half of the heap, placing them contiguously to avoid wasting space.



**Figure 7.5:** Copying garbage collectors divide the heap in half and move the in-use data to the reserved half, which has been left empty.

Figure 7.5b shows the resultant heap.

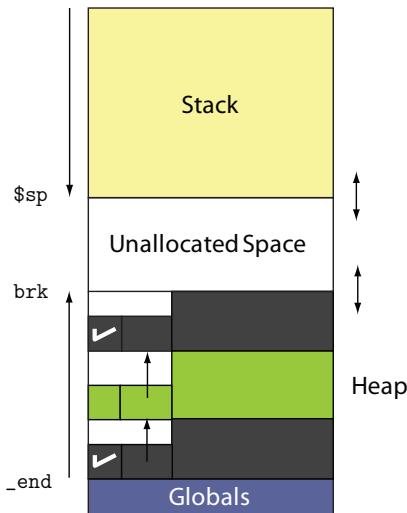
While copying all the live objects seems like it should be more expensive, the reduction of fragmentation in the free space usually negates the cost of copying. An allocation that has to search through a list or bitmap to find space is slow, whereas a contiguous region is simple to dole out.

To accommodate a copying garbage collector, some restrictions on references need to be made. Since the addresses that data items live at can change during the execution of the program, there can be no hard-coded addresses. Additionally, references must be kept distinct from integer types since when the system moves an object, all references must be able to be found before they can be properly updated. Indirection may be employed through a “table of contents,” where references are indices into a table that has real addresses that are updated as needed. This, however, incurs additional cost whenever a dereference occurs. While garbage collection makes life easier for the programmer, it does come at the cost of efficiency.

### 7.3 Linked List Example: `malloc()`

The C Standard Library provides a heap allocator called `malloc()`. When the loader creates an address space for the process, the typical layout is that code and global data start at a low address and extend as needed. The end of this fixed-size portion is represented as the symbol `_end`. The stack, by convention, starts at a high address and grows downward. The space between `_end` and the stack pointer can be used as the heap. `malloc()` and the Operating System denote the maximum space of the heap by the symbol `brk`. Figure 7.6 shows the relation of these symbols to each other.

The break can be set via a system call `brk()` or a library wrapper `sbrk()`. When



**Figure 7.6:** Heap management with `malloc()`.

`malloc()` gets a request for allocation that cannot fit, it extends `brk`. The heap is exhausted if the break gets too near the top of the stack. Likewise, the stack may be exhausted (usually from deep recursion) if it gets too close to `brk`.

Typically `malloc()` uses a linked list allocation strategy to track free and allocated space. One of the issues with linked lists is the question of where to store the list. An implementation of `malloc()` might store the linked list inside the heap, with each node near the allocated region. This allows calls to `free()` to easily access the size field of the node in the list corresponding to the region to reclaim. The drawback to this is that any over- or under-run of a heap-allocated buffer may overwrite the list, resulting in a corrupted heap.

Not all implementations of `malloc()` adjust the value of `brk`. The GNU implementation in `glibc` uses `mmap()` for allocations beyond 128KB. The `mmap()` system call requests pages directly from the Operating System. Some `malloc()`s use only `mmap()` for allocation.

## 7.4 Reducing External Fragmentation: The Buddy Allocator

The concern over external fragmentation has lead to various algorithms being developed beyond those used with the basic bitmap or linked list data structures. One particularly good algorithm for reducing external fragmentation is called the **buddy allocator**. Consider the case that we have a free memory region of 2MB and that we wish to allocate 4KB. The buddy algorithm looks for a free space reasonable to hold the allocation request. Right now, there is a single free space of size 2MB which is too large to reasonably allocate to this space. The buddy allocator takes the region and splits it into two allocations of half of the original size. So in this case, the algorithm turns our space into two allocations of 1MB each. This is still too much, and so one of the 1MB regions is split into two 512KB, and again split to 256KB, and so on, until finally we split an 8KB space into two 4KB regions. We now mark one of them as in-use and return it for the user.

Free regions of a particular size are linked together, using a portion of the region as a list node as in the implementation of `malloc()` described in the previous section. This way the left over regions from previous splits can be handed out quickly. We don't need to keep a list node for an allocated region, however. All that is necessary is a single bit that indicates that the region is free. Practically, we may still wish to reserve space for an entire list node for when the space is later freed.

We don't need to keep a list of allocated space because of a nice property that our scheme has. Every time a region was split into two, the two regions have addresses with a specific relationship. For a block of size  $n$  at address 0, when split it becomes two blocks of size  $\frac{n}{2}$ , one with address 0 and the other at address  $\frac{n}{2}$ . If we started with  $n$  being a power of two, then  $\frac{n}{2}$  is one less power of two. This means the regions that were created have addresses that differ only by the value of a single bit, in the position of the new size.

Given an address  $L$  of a region of size  $2^k$ , we can find its “buddy” (the other node split from the parent node) by inverting the  $k$ th bit. The bitwise XOR (eXclusive-OR) operator computes this for us: `buddy = L ^ size`. Note that this only works if the first region began at address 0 and was a size that is a power of two. If the starting address was something else, as it often will be in a practical application, we can simply subtract the actual starting address from  $L$  before performing the XOR operation.

If a block and its buddy are both free after deallocation, they can be safely

coalesced back into a node of size  $2^{k+1}$ . The free buddy will be removed from the linked list of free spaces of size  $2^k$  and the combined space will be inserted into the free list of size  $2^{k+1}$ . Since this may result in both a region and its buddy being free at size  $2^{k+1}$ , we can repeat this process with progressively larger regions until the newly coalesced region's buddy is not free or only the original entire free space is left.

## 7.5 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Bitmap** An array of bits indicating which chunks are allocated.

**Coalesce** To combine two or more nodes into one.

**External Fragmentation** Free space that is too small to be useful, a result of deallocation without compaction.

**Garbage Collection** Automatic deallocation of dynamic memory that occurs when a memory region is no longer needed.

**Heap** A region of a process's address space dedicated to dynamic data whose lifetime extends beyond that of the function that creates it.

**Internal Fragmentation** Wasted space due to the minimum allocation unit being too large.

**Memory Leak** When a dynamically allocated region cannot be deallocated because all pointers to it have been lost.

**Run-length Encoding** A simple compression scheme where  $n$  values in a row can be replaced by the number  $n$  in any base greater than one. For example, the string "AAAAAA" could be compressed to "5A".

# 8 | Operating System Interaction

THE **Operating System** (os) is a special process on a computer responsible for two major tasks: *managing resources* and *abstracting details*. An os manages the shared resources on the computer. These resources include the CPU, RAM, disk storage, and other input and output devices. The os is also useful in abstracting the specific details of the system away from application programmers. For instance, an os may provide a uniform way to print a document to a printer regardless of its specific make and model.

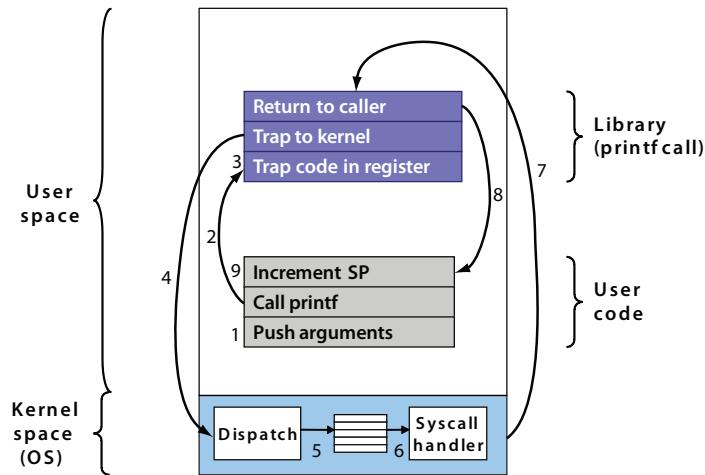
The core process of the os is called the **kernel**. The kernel runs at the highest privilege level that the CPU allows and thus can perform any action. The kernel is responsible for management and protection; it should be the most trusted component on the system. The kernel runs in its own address space that is referred to as **kernel space**. Programs that are not the kernel, referred to as **user programs**, cannot access the memory of the kernel and run at a lower privilege level. This portion of the computer is called **user space**.

Application programmers access the facilities that the Operating System provides via **system calls**. System calls are functions that the Operating System can do, usually related to process management and I/O operations.

## 8.1 System Calls

The particular system calls an Operating System provides depend on the particular os and version. Because of this, most application programmers access system calls via a library **wrapper function**. Wrapper functions provide an os-neutral way to do common tasks such as file and console I/O. The C Standard Library contains many wrapper calls as part of the `stdio` package.

Figure 8.1 shows the steps involved in executing a library call that needs to call the Operating System. If we have a call to `printf()` in our program, it is compiled



**Figure 8.1:** A library call usually wraps one or more system calls.

```
fstat(1, {st_mode=S_IFCHR|0600, st_rdev=makedev(136, 7), ...}) = 0
mmap(NULL, 4096, PROT_READ|PROT_WRITE,
      MAP_PRIVATE|MAP_ANONYMOUS, -1, 0) = 0x2a95557000
write(1, "Hello world!\n", 13Hello world!) = 13
exit_group(0)
```

**Figure 8.2:** A “Hello world” program run through `strace`.

to a `call` to the C Standard Library. In the library, the code to interpolate the arguments is run, and the final output string is prepared. When it is finally time to display the string, the library makes a system call to do the actual I/O operation. The Unix/Linux utility `strace` provides a list of system calls made during the execution of a program. In Figure 8.2 we see the system calls made by a “hello world” program using `printf()`.

The Unix and Linux Operating Systems provide a `write()` system call that interacts with I/O devices. The first parameter being the value `1` indicates that the output should go to the `stdout` device. Section 8.1.2 will detail the I/O calls provided by Unix and Linux systems.

### 8.1.1 Crossing into Kernel Space from User Space

Figure 8.1 also illustrates that a system call forces the CPU to switch modes from dealing with a process running in user space to the kernel running in kernel space. User space applications cannot cross this boundary themselves but instead issue a **trap** instruction that signals to the CPU that a **context switch** into the kernel should occur. A context switch occurs whenever the CPU switches from running one process to another. To resume the suspended process, the state of the machine — called a process's **context**, which includes the registers, open files, and stack — must be saved. To run the new process, its context must be restored. Context switches are very time consuming and all attempts to avoid doing more than absolutely necessary are made.

Whenever the CPU receives a trap or **interrupt**, it transfers control via an array of function pointers indexed by interrupt number called the **interrupt vector**, which is set up by the Operating System when it boots up. Under Linux the system call trap is `int 0x80`.

Upon entering the OS via a trap, the CPU is now in **kernel mode** and can perform the privileged operations of the OS, such as talking to the system I/O devices. The dispatcher selects the appropriate system call routine to execute based on the value of a register when the interrupt was sent. The kernel can now perform the requested operation.

### 8.1.2 Unix File System Calls

Listing 8.1 gives an example of how to write a program that writes text to a file without using the C Standard Library calls. Unix and Linux systems provide the primitive I/O operations of `open()`, `read()`, `write()`, and `close()` to operate on files. These operations end up being used to perform almost every I/O operation due to the traditional Unix paradigm of treating all devices as files in the file system.

While C programs that operate on files make use of a `FILE *` to track open files, the Unix system calls designate an integer as a *file descriptor*, which represents an open file. When a program begins, three file descriptors are automatically opened. They are:

Descriptor	C Name	Usage
0	<code>stdin</code>	Standard Input: Usually the keyboard
1	<code>stdout</code>	Standard Output: Usually the terminal
2	<code>stderr</code>	Standard Error: Usually the terminal

```
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>

int main() {
    int fd;
    char buffer[100];
    strcpy(buffer, "Hello, World!\n");

    fd = open("hello.txt", O_WRONLY | O_CREAT);
    write(fd, buffer, strlen(buffer));
    close(fd);
    exit(0);

    return 0;
}
```

**Listing 8.1:** Using the Unix system calls to do a “Hello World” program.

Notice that `stdout` and `stderr` both display upon the screen by default. They are separate streams, however, and may be redirected or piped independently of each other.

Another thing to notice about Listing 8.1 is the second parameter to `open()`. Two macros are bitwise-ORed together. If we look for the definitions of these in the header files, we see them as:

```
#define O_RDONLY      0
#define O_WRONLY       1
#define O_RDWR         2
#define O_CREAT        16
```

which are all powers of two. This technique is common when we want to send several flags that affect the operation of a function. Each separate bit in an integer can be seen as an independent boolean flag. Bitwise-ORing them together allows the programmer to specify one or more flags simultaneously. In this example the flags to open the file for writing and to create it if it does not already exist are set. The implementation can check whether a particular flag is set by bitwise-ANDing the parameter with the same constants. This technique is also very commonly seen in the functions Microsoft Windows provides.

### 8.1.3 Process Creation with `fork()` and `execv()`

No survey of important system calls under Unix/Linux would be complete without discussing some issues of process management. A programmer may frequently wish to spawn off a new process to do some additional work concurrently. The system call to do this is `fork()`. Listing 8.2 shows an example of `fork()`'s usage. When `fork()` is called, the Operating System creates a clone of the original process that is identical in every way except for the return value of the `fork()` call. In the original process, denoted the parent, the return value is the **process ID** — the number that the os uses to keep track of processes — of the child. In the child process, the return value is zero.

If we run the program, we might try to predict its output. The problem with doing this is that there is usually no guarantee that two separate processes will run in any particular order. At the very least, we can be sure that there will be four lines printed to the screen: the two inside the `if/else` and each process's copy of the “Hi from both” line. The other guarantee is that the lines will print in the proper relative order in terms of a single process. That is, there is no way to see both “Hi from both” lines before “Hi from the child” and “Hi from the parent” have each

```
#include <stdio.h>
#include <unistd.h>

int main() {
    if(fork()==0) {
        printf("Hi from the child!\n");
    }
    else {
        printf("Hi from the parent\n");
    }

    printf("Hi from both\n");
    return 0;
}
```

**Listing 8.2:** An example of process creation using fork.

been displayed. On the test run, the following output was seen:

```
Hi from the child!
Hi from both
Hi from the parent
Hi from both
```

This indicates that the child process ran and completed before the parent process resumed its execution.

The other common use of `fork()` is to launch a separate program entirely. The family of `execv()` functions all wrap around the `execv()` system call, which embodies the loader we described in Chapter 3. The unusual thing about `execv()` is that it needs a process to be created for it. The system call itself will not create a process; rather it will replace the process that called it with the new program to be loaded. This means that `execv()` often comes shortly after a call to `fork()`. Listing 8.3 shows an example of a program that launches the `ls` program. The parent, in the `else`, waits for all child processes to complete before it continues by using the `wait()` function.

```

#include <stdio.h>
#include <unistd.h>

int main() {
    if(fork()==0) {
        char *args[3] = {"ls", "-al", NULL};
        execvp(args[0], args);
    }
    else {
        int status;
        wait(&status);
        printf("Hi from the parent\n");
    }
    return 0;
}

```

**Listing 8.3:** Launching a child process using `fork` and `execvp`.

## 8.2 Signals

A **signal** is a message from the Operating System to a user space program. Signals are generally used to indicate error conditions, in much the same way that Java Exceptions function. A program can register a handler to “catch” a particular signal and will be asynchronously notified of the signal without the need to **poll**. Polling is simply the action of repeatedly querying (e.g., in a loop) whether something is true.

Figure 8.3 shows a list of the os signals on a modern Linux machine. You can generate a complete list for your system by executing the command `kill -1`. Most signals tend to fall into a few major categories. There are the error signals, which indicate something has gone awry:

**SIGILL** The CPU has tried to execute an illegal instruction.

**SIGBUS** A bus error, usually caused by bad data alignment or a bad address.

**SIGFPE** A floating point exception.

**SIGSEGV** A segmentation violation, i.e., a bad address.

There are several ways to tell a program to forcibly exit on a system:

SIGHUP	SIGINT	SIGQUIT	SIGILL	SIGTRAP	SIGABRT
SIGBUS	SIGFPE	SIGKILL	SIGUSR1	SIGSEGV	SIGUSR2
SIGPIPE	SIGNALRM	SIGTERM	SIGCHLD	SIGCONT	SIGSTOP
SIGTSTP	SIGTTIN	SIGTTOU	SIGURG	SIGXCPU	SIGXFSZ
SIGVTALRM	SIGPROF	SIGWINCH	SIGIO	SIGPWR	SIGSYS

**Figure 8.3:** The standard signals on a modern Linux machine.

**SIGINT** Interrupt, or what happens when you hit **CTRL+C**.

**SIGTERM** Ask nicely for a program to end (can be caught).

**SIGKILL** Ask meanly for a program to end (cannot be caught).

**SIGABRT, SIGQUIT** End a program with a core dump.

The remaining signals send information about the state of the os, including things like the terminal window has been resized or a process was paused.

### 8.2.1 Sending Signals

Signals can be sent programmatically by using the `kill()` system call. The name is somewhat misleading, since any of the signals can be sent by using it, but most often the signals that are sent seem to be related to process termination. The code in Listing 8.4 makes the program stop with the `SIGSTOP` signal. The program will not resume until it receives a `SIGCONT` signal. The `getpid()` call asks the os for the current process's id. If you know the process id of another process, you can send it signals as well.

It is useful to be able to send termination signals via the command line in a Unix shell, and the command `kill` allows you to do this. You need to pass the process id, but that can be obtained using the `ps` command. By default, `kill` with a process id argument will send that process a `SIGTERM` signal, which a process can choose to ignore. If a process is truly crashed, it is better to send the `SIGKILL` signal. You can specify the signal to send by saying `-NAME`, where name is the name portion of a signal, like `SIGNAME`. You may also specify a signal's numerical value, and you will often see a forcible termination of a process done like:

```
kill -9 process_to_kill_pid
```

```
#include <unistd.h>
#include <sys/types.h>
#include <signal.h>

int main() {
    pid_t my_pid = getpid();
    kill(my_pid, SIGSTOP);
    return 0;
}
```

**Listing 8.4:** Signals can be sent programmatically via `kill`.

### 8.2.2 Catching Signals

Much like Java Exceptions, signals can be caught by a program. However, since many of the signals indicate something has gone terribly wrong, great care has to be taken in how a caught exception is dealt with. Any of the termination signals, if caught, should only be used as an opportunity to clean up any open or temporary files and exit gracefully. Memory state may be corrupted, so any attempts to continue will just make the program fail further down the line. A few signals are useful to catch, like `SIGALRM`. Listing 8.5 gives an example.

In this program, we set up a signal handler by using the `signal()` function. It takes two parameters, the first is the signal to listen for, the second is a function pointer of a function to call upon receipt of the signal. The alarm signal allows you to specify a timeout upon which the program will be notified that the time has elapsed. In our example, `alarm()` function tells the os to send a `SIGALRM` in one second. This is not a precise time, as it may be beyond one second that the handler is called. The program in Listing 8.5 makes a countdown timer from ten to one and then exits.

One final signal of particular interest is `SIGTRAP`. This is the *Breakpoint Trap* signal. Debuggers such as `gdb` listen for this signal to retake control of an executing process in order to examine it.

The Intel x86 instruction set defines all interrupts with a two-byte instruction encoding. The `int` opcode is `0xCD` followed by the interrupt number. For example, the Linux system call trap (`int 80`) would be `0xCD 0x80`. However, there is a special one-byte encoding for the breakpoint trap, `int 3`. It is the opcode `0xCC`.

Debuggers place breakpoints by overwriting existing instructions, since in-

```
#include <unistd.h>
#include <signal.h>

int timer = 10;

void catch_alarm(int sig_num) {
    printf("%d\n",timer--);
    alarm(1);
}

int main() {
    signal(SIGALRM, catch_alarm);

    alarm(1);
    while(timer > 0) ;
    alarm(0);
    return 0;
}
```

**Listing 8.5:** SIGALRM can be used to notify a program after a given amount of time has elapsed.

serting them would require rewriting the code. With the x86's variable-length instruction set architecture, a two-byte breakpoint might overwrite more than one instruction. This would be problematic if a particular breakpoint was skipped over and the target of a jump was the second byte of the breakpoint. To avoid this problem, the breakpoint trap is given a special one-byte encoding. Remembering this encoding may come in handy if you ever are dealing with low-level code and want to insert a breakpoint by hand.

### 8.3 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Context** The state of a process, necessary to restart it where it left off. Usually includes the state of the registers and other hardware details.

**Context Switch** The act of saving the context of a running process and restoring the context of a suspended process in order to change the currently running program.

**Interrupt** A CPU instruction or signal (the voltage kind) issued by hardware that interrupts the currently executing code and jumps to a handler routine installed by the Operating System. On Intel x86 computers, there is no distinction in name between an interrupt and a trap; both are referred to as interrupts.

**Interrupt Vector** An array of function pointers indexed by interrupt number used to call an os-installed handler routine.

**Kernel** The core process of an Operating System.

**Kernel Space** The Operating System's address space.

**Operating System** A program that manages resources and abstracts details of hardware away from application programmers.

**Poll** The act of repeatedly querying some state in a loop.

**Process ID** The number that the os uses to keep track of processes.

**Signal** A message from the Operating System, delivered asynchronously, that usually indicates an error has occurred.

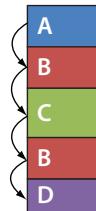
**System Call** A function that the Operating System provides to user programs to interact with the system and/or perform I/O operations.

**Trap** A software interrupt, usually used to signal the CPU to cross into kernel space from user space. **See also:** *interrupt*

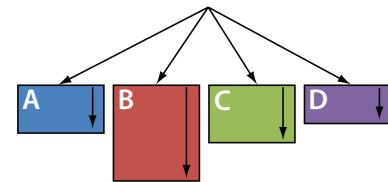
**User Program** Any application that is not part of the Operating System and runs in User Space.

**User Space** The unprivileged portion of the computer in which user programs run.

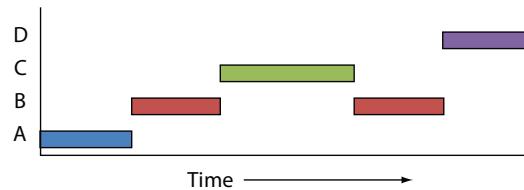
**Wrapper Function** A function, typically part of a library, that provides a generic way to do a system-specific common task such as file and console I/O.



**(a)** The CPU, with only one program counter, sees a contiguous stream of instructions.



**(b)** Each process, with its own isolated address space, appears to have a dedicated program counter and is never interrupted.



**(c)** Multiprogramming as viewed over time on a single CPU.

**Figure 9.1:** While the CPU sees just one stream of instructions, each process believes that it has exclusive access to the CPU.

## 9 | Multiprogramming & Threading

IN CHAPTER 5, we introduced the abstraction of a running program in memory called a process. One of the significant parts of that abstraction was the process's view that it had a large amount of RAM all to itself as part of its address space. Another part of the process abstraction is the idea that a process has the entire CPU to itself. On a large supercomputer or a small embedded device, this might be the

case. In general, however, a modern Operating System has to manage multiple processes all competing for the shared resource of the CPU. Figure 9.1 illustrates the differing perspectives that the CPU and the processes have. The CPU sees one unified stream of instructions, whereas each process believes it has exclusive access to the machine.

An observation about running processes is that they frequently will need to perform some type of I/O operation. Whenever an I/O device is accessed, there is a delay until it is ready to transmit or receive data. During this delay, a process cannot proceed and is said to be **blocked**. For instance, imagine a very fast typist typing at 120 words per minute. Two words per second is fast for a person, but modern computers can do billions of operations per second. Most of the time, the computer is waiting for the user to do something. If during this waiting time the computer could do other work, it could provide the illusion of doing multiple tasks at the same time. This abstraction is called **multiprogramming**. Figure 9.1c shows the progress of processes over time. Note the gap in B's execution while C runs. Process B should be totally unaware that it was not in control of the CPU for a portion of time.

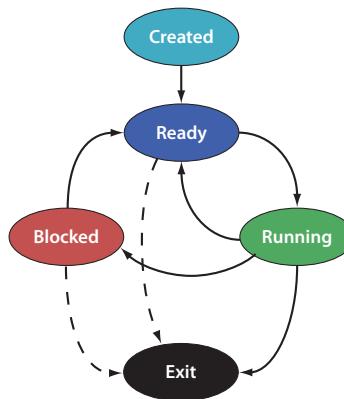
To change from one application to another, we need to save the current application's CPU state, which we defined before as its context. Just as we did a context switch when we wanted to enter the kernel to perform a system call, when one process is paused and another is begun, the Operating System does a context switch.

It is also possible that a program is **CPU-Bound**, meaning that it does not do very much I/O but rather is computing something intensive, using as much of the CPU as it can get. In these cases, a CPU-bound process would stop other processes from being able to run. To prevent such a process from starving out the others, periodically the hardware will issue a timer interrupt. This timer interrupt will cause the Operating System to run, and then it can determine whether the program should continue or should be paused to let another process have the CPU. This process of pausing a running program after a period of time is called **preemption**.

Figure 9.2 shows the life cycle of a process. When a process is created, it enters a queue of processes on the system that are available to run, called the *ready queue*<sup>1</sup>. When a process is in the running state it eventually will yield the CPU, either because it performs an I/O operation and transitions to the blocked state or because it is preempted and goes back into the ready queue. At this point, a component of the Operating System called the **scheduler** chooses a ready process (from the queue) and

---

<sup>1</sup> While this queue may be managed in FIFO order, we will use 'queue' just to imply a set of waiting objects.



**Figure 9.2:** The life cycle of a process. Dashed lines indicate an abnormal termination.

Per Thread	Per Process
Program counter	Address space
Registers	Open files
Stack & stack pointer	Child processes
State	Signals & handlers
	Accounting info
	Global variables

**Figure 9.3:** Thread State versus Process State.

lets it run. Processes eventually finish voluntarily or because of an error (indicated by the dashed transitions in Figure 9.2).

Sometimes, we as programmers know that a process will be **I/O Bound**, and we would explicitly like to do some other task related to the process simultaneously. If we launch another process to do work in parallel, we have two issues. The first issue is that we have no guarantee from the scheduler when a particular process will run, so we may not be doing the other work while the first process is blocked. The second problem is that because of the isolation an address space provides, it is very cumbersome and slow to share information between processes. Ideally, we'd like the ability to run multiple streams of instructions that share a single address space so that sharing data is as easy as loads and stores.

We can have multiple streams of instructions in a single process's address space by having a mechanism called a **thread**. A thread is a stream of instructions and its

associated context. A thread's context should be small, since the Operating System will still manage the process as a whole. For instance, a list of open files is part of a process's context but not an individual thread's. Figure 9.3 shows a list of what might be part of a process's context and what context needs to be stored per thread. If we define a stream of instructions as a function, we can more easily see what state we need to store. A function needs the machine registers and a stack of its own. This will form a minimal thread context.

## 9.1 Threads

Every process has one thread by definition. An application must explicitly create any additional threads it needs. Support for threading can come from two sources. In **user threading**, a user-space library provides threading support with minimal, if any, support from the Operating System. In **kernel threading**, the Operating System has full support for threads and manages them in much the same way as it manages processes.

### 9.1.1 User Threading

With user threading, we assume the Operating System has no explicit support for threads. A library containing helper functions will take care of thread creation and maintenance of the threads' state. There are two major hurdles to making user threading work. Imagine two processes are running on a system, one of which wants 100% of the CPU, the other wants 10%. Since the Operating System can preempt the greedy process, the one that needs only a little bit of CPU time can get it.

Imagine now two threads of a process that have the same characteristics. The greedy thread runs and runs until the process containing both threads is interrupted. When the process resumes execution, the greedy thread continues to run. The process abstraction prevents the process from knowing it was ever stopped. While each process was protected from the others on a system, there is no protection from the Operating System for threads since it does not know they are even there!

This may seem to be the downfall of user threading. However, if we go back to the original motivation for having threads at all, it was that we wanted to do work in a cooperative way and that the isolation afforded to us by an address space was too much. Since threads live within a single application, they can be expected to cooperate. A user threading library could then supply a `yield()` call whereby one thread voluntarily gives up the CPU, and the threading library can pick a different

thread to run. Application programmers writing a multithreaded program only need to be aware that they should call `yield()` at appropriate times. This explicit yielding even gives a bit of extra power to the developer, because it becomes easy to make one thread be more important than the others by having it yield less frequently.

Having solved that problem, let us tackle a second. Imagine that a thread is executing some code and comes upon an I/O operation. The operation traps into the kernel, and finding the data not yet ready, the kernel moves the process into the blocked state. One of the original motivations was that during these times of being blocked, we would like to run some other code to do some useful work, so we might assume that another thread will get to run. But remember, the kernel knows nothing of the threads and has put the entire process to sleep. There is no way any other code in the process can run until the I/O request has finished.

Though we found a reasonable way around the yielding problem, this seems the death knell for user threads. There is no way to avoid the inevitable block that will happen to the process when the I/O operation cannot be completed. The only way around it would be if there was some facility by which going into a blocked state could be prevented. If an Operating System has a facility for non-blocking I/O calls, the user threading library could use them and insert a yield to run a different thread until the requested data was ready.

Unix/Linux systems have a system call named `select()` or `poll()` that tells whether a given I/O operation would block. It has the added side effect of doing the actual I/O request. Since `select()` can be non-blocking, the threading library could provide its own version of the I/O calls. A thread would use the library's routines, and when in the library, the library could make a call to `select()` to see if the operation would block. If it would, it can put that thread to sleep and allow another thread to run. The next time we are in the library, via a `yield()`, a `create()`, or an I/O call for another thread, we can check to see whether the original call is ready, and if so, unblock the requesting thread.

Since an Operating System needs to have non-blocking I/O support to make user threading work, it is arguable that user threads actually require no explicit support. The `select()` call is useful for more than just threads, including checking to see whether any network packets have arrived. While this minimal level of functionality is required, we will see with kernel threads that the level of OS support is far beyond a single system call.

### 9.1.2 Kernel Threading

Kernel threading is the complete opposite of user threading. With kernel threads, the Operating System is completely in charge of managing threads. The OS has a system call that creates threads. The scheduler in the OS knows that when one thread is blocked, another thread from the very same process could still run. A CPU-bound thread can be preempted without any need for the programmer to put in explicit calls to a `yield()` function. In short, kernel threads are everything that user threads were not. However, kernel threads also may be slower to create.

In user threads, any thread operation was a library call, and since no process boundaries were crossed, no context switches needed to be done. But anytime a thread is created in kernel threading, the OS must update its internal record-keeping, and a context switch into the kernel must be done. Context switches are expensive, so creating many threads will be significantly slower. Thread and process switching will both require context switches with kernel threading. These issues may not be a problem for a program using a small number of threads, but a Web server for a high-demand Web site may be spawning new threads for every request a hundred times per second. In that case, kernel threads may have too much overhead.

Ultimately, when it comes time to make a multithreaded program, you will be at the mercy of the system you are developing for. If it has kernel threads, your life may be easier, but performance may not be as good as possible. The good news is that you can always use user-threading libraries on a system with kernel threads. Use the tool that works best for the situation.

## 9.2 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Blocked** A process that is unable to continue because it is waiting for something, usually an I/O request to complete.

**CPU-Bound** A process that primarily needs to do computation and rarely needs to do an I/O operation.

**I/O Bound** A process that spends most of its time blocked.

**Kernel Threading** Operating System support for thread management.

**Multiprogramming** Part of the process abstraction where a process appears to have the CPU entirely to itself, even when there are multiple processes on a single machine.

**Preemption** Interrupting a process because it has had the CPU for some amount of time in order to allow another process to run.

**Scheduler** The portion of the Operating System responsible for choosing which process gets to run.

**Thread** A stream of instructions and its associated context.

**User Threading** Threading done via a user space library that provides thread support with minimal, if any, support from the Operating System.

# 10 | Practical Threading, Synchronization, & Deadlocks

To **FACILITATE** multithreaded programming on a wide range of Operating Systems, a standard library was developed to hide the implementation details of writing multithreaded programs. The **POSIX** group created a standard library called **pthreads** that allows for the creation and management of multithreaded programs without concern for the underlying implementation — that is, whether user threading or kernel threading is supported.

With an abstraction such as the **pthreads** library, it is possible to write portable threaded programs that run on a variety of systems. In this chapter, we will examine the functionality the library provides. Once we are able to create and manage threads, we will consider issues of **synchronization**: The need for many threads that share data to safely modify that data. Likewise, we need to examine the issue of **deadlocks**. Deadlocks arise when two or more threads are unable to make progress while waiting on each other to do some task.

## 10.1 Threading with **pthreads**

The most basic operations for a threading library to provide are functions for creating and destroying threads. The **pthread** library provides several functions to aid us with these tasks. All processes consist of one thread of execution already, and so any additional threads we want to create will be in addition to this original thread. Thread creation is done via the **pthread\_create()** function.

Listing 10.1 shows an example of running the **do\_stuff()** function in two threads. The main thread is “recycled” to run the function in addition to the newly spawned thread. **pthread\_create()** takes four parameters. The first is a pointer to a variable of type **pthread\_t** that is set to a unique identifier for the newly created

```

#include <stdio.h>
#include <pthread.h>

void *do_stuff(void *p) {
    printf("Hello from thread %d\n", *(int *)p);
}

int main() {
    pthread_t thread;
    int id, arg1, arg2;
    arg1 = 1;
    id = pthread_create( &thread, NULL,
                         do_stuff, (void *)&arg1 );
    arg2 = 2;
    do_stuff((void *)&arg2);
    return 0;
}

```

**Listing 10.1:** Basic thread creation.

thread. This identifier is “opaque,” meaning that we do not know what type this identifier is (integer, structure, etc.) and we should not depend on its being any particular type or having any particular value. The second parameter controls how the thread is initialized, and for most simple implementations it can be set to `NULL` to take on the defaults.

The third and fourth parameters specify the stream of instructions to run in the new thread. First comes a pointer to a function containing the code to run. This function can, of course, call other functions, but it could be considered analogous to a “main” function for that particular thread. The signature of the function must be such that it takes and returns a `void *`. Because of the strict type checking done on passing function pointers in terms of return values and parameters, this function needs to be as generic as possible while still having a well-defined prototype. The advantage in using a `void *` is that it can point to anything, even an aggregate data type like an array or structure. In this way, no matter how many arguments are actually needed, the function can receive them. The final parameter is the actual parameters to pass to this function, which can be `NULL` if unnecessary.

```
int main() {
    pthread_t thread;
    int id, arg1, arg2;
    arg1 = 1;
    id = pthread_create(&thread, NULL,
                        do_stuff, (void *)&arg1);
    pthread_yield();
    arg2 = 2;
    do_stuff((void *)&arg2);
    return 0;
}
```

**Listing 10.2:** Inserting a yield to voluntarily give up execution.

Compiling this program requires an additional command-line option to `gcc`. Since the `pthread` library is not part of the C Standard Library, it is not linked against the program automatically. Adding the `-pthread` switch tells the linker to include the appropriate code and data.

Executing the resulting program may lead to some interesting results. On one system, the following output from Listing 10.1 was seen:

```
Hello from thread 2
```

It would appear that the newly created thread is not run. If we did a similar test using `fork()` instead, the output of both would be seen. So what is the difference? When the `main()` function returns, the process is over. Since `main()` terminates relatively quickly, the other thread never gets a chance to run. With two processes (from `fork()`), each will not terminate until its respective `main()` finishes, guaranteeing that each will print its output before completing.

Some control over when a thread runs or when the process terminates is necessary. From the discussion of threading in the previous chapter, we know that it is possible to voluntarily yield control to a separate thread. The `pthread` library exposes this through the `pthread_yield()` function. If we rewrite the `main()` function as in Listing 10.2, we get the following output:

```
Hello from thread 1
Hello from thread 2
```

However, this is no guaranteed solution. While the `pthread_yield()` is a suggestion to let another thread run, there is no way to *force* this to happen. The other

```

int main() {
    pthread_t thread;
    int id, arg1, arg2;
    arg1 = 1;
    id = pthread_create(&thread, NULL,
                        do_stuff, (void *)&arg1);
    arg2 = 2;
    do_stuff((void *)&arg2);
    pthread_join(thread, NULL);
    return 0;
}

```

**Listing 10.3:** Waiting for the spawned thread to complete.

thread(s) might be blocked, or the scheduler might simply ignore the yield. A better solution is to force the process to wait until the other thread completes.

Listing 10.3 illustrates the better way to ensure threads complete. Calling the `pthread_join()` function blocks the thread that issued the call until the thread specified in the parameter finishes. The second parameter to the `pthread_join()` call is a `void **`. When a function needs to change a parameter, we pass a pointer to it. Passing a pointer-to-a-pointer allows a function to alter a pointer parameter, in this case, setting it to the return value of the thread the join is waiting on. We can, of course, choose to ignore this parameter, in which case we can simply pass `NULL`. Note, however, joining the threads still does not guarantee they will run in any particular order before the call to `pthread_join()`.

The moral of the threading story is that, unless explicitly managed, threads run in no guaranteed order. This lesson becomes even more important when we begin to access shared resources in multiple threads concurrently. When we need to manipulate shared objects, we may need to ensure a particular order is preserved, which leads to the next topic: Synchronization.

## 10.2 Synchronization

Imagine that there are two threads, Thread 0 and Thread 1, as in Figure 10.1. At time 3, Thread 0 is preempted and Thread 1 begins to run, accessing the same memory location X. Because Thread 0 did not get to write back its increment to

Thread 0	Thread 1
1    read X	
2    X = X + 1	
	3    read X
	4    X = X + 1
5    write X	
	6    write X

**Figure 10.1:** Two threads independently running the same code can lead to a race condition.

memory, Thread 1 has read an older version. Whichever thread writes last is the one that makes the update, and the other is lost. When the order of operations, including any possible preemptions, results in different values, the program is said to have a **race condition**.

Determining whether code is susceptible to race conditions is an exercise in Murphy's Law.<sup>1</sup> Race conditions occur when code that accesses a shared variable may be interrupted before the change can be written back. The obvious solution to preventing a race condition is to simply forbid the thread from being interrupted during execution of this **critical region** of code. Allowing a user-space process to control whether it can be preempted is a bad idea, however. If a user program were allowed to do this, it could simply monopolize the CPU and not allow any other programs to run. Whatever the solution, it will require the help of the Operating System, as it is the only part of the system we can trust to make sure an action is not interrupted.

A better solution is to allow a thread to designate a portion of its code as a critical region and control whether other threads can enter the region. If a thread has already entered a critical region of code, all other threads should be blocked from entering. This lets other threads still run and do “non-critical” code; we have not given up any parallelism. The marking of a critical region itself must not be interruptible, a trait we refer to as being “atomic.” This atomicity and the ability to make other threads block means that we need the Operating System or the user-thread scheduler's help.

Several different mechanisms for synchronization are in common use. We will focus on three in this text, although a fourth, known as a *Monitor*, forms the basis

---

<sup>1</sup> Anything that can go wrong, will go wrong.”

Thread 0	Thread 1
1 lock mutex	
2 read X	
3 X = X + 1	
4	lock mutex
5 write X	
6 unlock mutex	
7	read X
8	X = X + 1
9	write X
10	unlock mutex

**Figure 10.2:** Synchronizing the execution of two threads using a mutex.

for Java's support for synchronization. The `pthread` library provides support for *Mutexes* and *Condition Variables*. *Semaphores* can be used with the inclusion of a separate header file.

The `pthread` library provides an abstraction layer for synchronization primitives. Regardless of the facilities of the Operating System, with respect to its support for threading or synchronization, mutexes and condition variables will always be available.

### 10.2.1 Mutexes

The first synchronization primitive is a **mutex**. The term comes from the phrase **Mutual Exclusion**. Mutual exclusion is exactly what we are looking for with respect to critical regions. We want each thread's entry into a particular critical region to be exclusive from any other thread's entry. A mutex behaves as a simple lock, and thus we get the two operations `lock()` and `unlock()` to perform.

With a mutex and the lock and unlock operations, we can solve the problem of Figure 10.1. Figure 10.2 assumes a mutex variable named `mutex` that is initially in an unlocked state. Thread 0 comes along first, locks the mutex, and proceeds to do its work up until time 4, when it is preempted and Thread 1 takes over. Thread 1 attempts to acquire the mutex lock but fails and is blocked. With no other threads to run, Thread 0 resumes and finishes its work, unlocking the mutex. With the mutex now unlocked, the next time that Thread 1 is scheduled to run it can, as it is no longer in the blocked state.

The `pthread` library provides a simple and convenient way to use mutexes. A

```
#include <stdio.h>
#include <pthread.h>
int tail = 0;
int A[20];
pthread_mutex_t mutex = PTHREAD_MUTEX_INITIALIZER;

void enqueue(int value) {
    pthread_mutex_lock(&mutex);
    A[tail] = value;
    tail++;
    pthread_mutex_unlock(&mutex);
}
```

**Listing 10.4:** Using a `pthread_mutex` to protect an enqueue operation on a shared queue.

`mutex` is declared of type `pthread_mutex_t` and can easily be initialized to start off unlocked by assigning `PTHREAD_MUTEX_INITIALIZER` to it. Locking and unlocking are done via the `pthread_mutex_lock()` and `pthread_mutex_unlock()` functions. Listing 10.4 shows an example of protecting a shared queue using a mutex.

### 10.2.2 Condition Variables

Sometimes we would like to do more than just protect a region. Imagine two threads working together, one of which is producing items into a fixed-size buffer, the other consuming them from that same buffer. This is a classic problem of synchronization known as the *Producer/Consumer Problem*. Figure 10.3 gives a pseudocode implementation. Since the buffer is fixed-size, we must be careful not to overrun or underrun the bounds. We want to stop the producer when the buffer is full and stop the consumer when it is empty. Let us then assume that we have a `sleep()` call to put the current thread to sleep and a corresponding `wakeup()` that will wake up a particular sleeping thread by notifying the scheduler that the thread is no longer blocked.

If we examine the consumer, we see the conditional `if(counter==0)` and the action `sleep()`. Here, obviously, the code is being guarded from the possibility of underrun. Possibly less obvious is the `if(count==N-1)` followed by the `wakeup(producer)`. If `count` is currently one less than the maximum, we know the

### Shared Variables

```
#define N 10;
int buffer[N];
int in = 0, out = 0, counter = 0;

Consumer
while(1) {
    if(counter == 0)
        sleep();
    ... = buffer[out];
    out = (out+1) % N;
    counter--;
    if(counter == N-1)
        wakeup(producer);
}

Producer
while(1) {
    if(counter == N)
        sleep();
    buffer[in] = ... ;
    in = (in+1) % N;
    counter++;
    if(counter==1)
        wakeup(consumer);
}
```

**Figure 10.3:** A pseudocode implementation of the Producer/Consumer problem. Note that this code has an unresolved synchronization problem.

```

void *producer(void *junk) {
    while(1) {
        pthread_mutex_lock(&mutex);
        if( counter == N )
            pthread_cond_wait(&prod_cond, &mutex);
        buffer[in] = total++;
        printf("Produced: %d\n", buffer[in]);
        in = (in + 1) % N;
        counter++;
        if( counter == 1 )
            pthread_cond_signal(&cons_cond);
        pthread_mutex_unlock(&mutex);
    }
}

```

**Listing 10.5:** The producer function using condition variables. The consumer function would be similar.

buffer was full before we consumed an item, and if the buffer was full, the producer is asleep, so wake it up.

However, there is a subtle problem here. Imagine that the consumer is running, executes the `if(counter==0)` line, and finds the buffer empty. But right before the sleep is executed, the thread is preempted and stops running. Now the producer has a chance to run, and since the buffer is empty, successfully produces an item into it. The producer notices that the count is now one, meaning the buffer was empty just before this item was produced, and so it assumes that the consumer is currently asleep. It sends a wakeup which, since the consumer has not yet actually executed its sleep, has no effect. The producer may continue running and eventually will fill up the buffer, at which point the producer itself will go to sleep. When the consumer regains control of the CPU, it executes the `sleep()` since it has already checked the condition and cannot tell that it has been preempted.<sup>2</sup> Now both the consumer and the producer threads are asleep and no useful work can be done. This is called a **deadlock** and is the subject of Section 10.3.

There are two ways to prevent this problem. The first is to make sure that the check and the sleep are not interrupted. The second is to remember that there was a

---

<sup>2</sup> In fact, it does not need to have been preempted if these two threads were running on separate cores or processors.

wakeup issued while the thread was not sleeping and to immediately wake up when the next sleep is executed in that thread. The first way is implemented via condition variables, the second solution is a semaphore.

A **condition variable** is a way of implementing sleep and wakeup operations in the pthread library. A variable of type `pthread_cond_t` represents a “condition” and acts somewhat like a phone number. If a thread wants to sleep, it can invoke `pthread_cond_wait()` and go to sleep. The first parameter is a condition variable that enables another thread to “phone it” and wake it up. A thread sleeping on a particular condition variable is awoken by calling `pthread_cond_signal()` with the particular condition variable that the sleeping thread is waiting on. Condition variables can be initialized much the same way that mutexes were, by assigning a special initializer value to them appropriately called `PTHREAD_COND_INITIALIZER`.

While this enables us to wait and signal (sleep and wake up), we still have the issue of possibly being interrupted. This is where the second parameter of `pthread_cond_wait()` comes into play. This parameter must be a mutex that protects the condition from being interrupted before the wait can be called. As soon as the thread sleeps, the mutex is unlocked, otherwise deadlock would occur. When a thread wakes back up it waits until it can reacquire the mutex before continuing on with the critical region. Listing 10.5 shows the producer function rewritten to use condition variables.

### 10.2.3 Semaphores

The `semaphore.h` header provides access to a third type of synchronization, a **semaphore**. A semaphore can be thought of as a counter that keeps track of how many more wakeups than sleeps there have been. In this way, if a thread attempts to go to sleep with a wakeup already having been sent, the thread will not go to sleep. Semaphores have two major operations, which fall under a variety of names. In the `semaphore.h` header, the operations are called wait and post, but they can also be known as lock and unlock, down and up, or even P and V. Whenever a wait is performed on a semaphore, the corresponding counter is decremented. If there are no saved wakeups, the thread blocks. If the counter is still positive or zero, the thread can continue on. The post function is an increment to the counter and if the counter remains negative, it means that there is at least one thread waiting that should be woken up.

One way to conceptualize the counter is to consider it as maintaining a count of how many resources there are currently available. In the Producer/Consumer

```
void *producer(void *junk) {
    while(1) {
        sem_wait(&semempty);
        sem_wait(&semmutex);
        buffer[in] = total++;
        printf("Produced: %d\n", buffer[in]);
        in = (in + 1) % N;
        sem_post(&semmutex);
        sem_post(&semfull);
    }
}
```

**Listing 10.6:** The producer function using semaphores. The consumer function would be similar.

example, each array element is a resource. The producer needs free array elements, and when it exhausts them it must wait for more free spaces to be produced by the consumer. We can use a semaphore to count the free spaces. If the counter goes negative, the magnitude of this negative count represents how many more copies of the resource there would need to be to allow all of the threads that want a copy to have one.

Semaphores and mutexes are very closely related. In fact, a mutex is simply a semaphore that only counts up to one. Conceptually, a mutex is a semaphore that represents the resource of the CPU. There can only be one thread in a critical region that may be running, and all other threads must block until it is their turn.

Semaphores can be declared as a `sem_t` type. There is no way to have a fixed initializer, however, because a semaphore can initially take on an integer value rather than being locked or unlocked. A semaphore is initialized via the `sem_init()` function, which takes three parameters: the semaphore variable, the value `0` on all Linux machines, and the initial value for the semaphore. Listing 10.6 shows the producer/consumer problem solved by using semaphores. Notice that they can even replace the mutex, although we could use a mutex if we wanted. The `semempty` (initialized to `N`) and `semfull` (initialized to `0`) semaphores count how many empty and full slots there are in the buffer. When there are no more empty slots, the producer should sleep, and when there are no more full slots, the consumer should sleep.

## 10.3 Deadlocks

The formal definition of a deadlock is that four things must be true:

**Mutual exclusion** Only one thread may access the resource at a time.

**Hold and wait** When trying to acquire a new resource, the requesting thread does not release the ones it already holds.

**No preemption of the resource** The resource cannot be forcibly released from the holding thread.

**Circular wait** A thread is waiting on a resource that is owned by a second thread which, in turn, is waiting on a resource the first thread has.

For our purposes, we will truly worry about the circular wait condition. This means that we must be careful about how and when we acquire resources, including mutexes and semaphores. If we do something as simple as mistakenly alter the order of the semaphores from Listing 10.6 to be:

```
sem_wait(&semmutex);
sem_wait(&semempty);
```

our program will instantly deadlock. If there are no empty slots, the thread does not release the semaphore used for mutual exclusion so that the other thread may run and consume some items.

Ensuring that your code is deadlock-free can sometimes be a difficult task. A simple rule of thumb can help you avoid most deadlocks and produce code that spends as much time unblocked as possible:

*Always place the mutex (or semaphore being used as a mutex) around the absolute smallest amount of code possible.*

This is not a perfect rule, and surely there is a counter-example to defeat it. No rules will ever replace understanding the issues of synchronization and using them to illuminate the potential problems of your own code.

## 10.4 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Critical Region** A region of code that could result in a race condition if interrupted.

**Deadlock** A program that is waiting for events that will never occur.

**Race Condition** A region of code that results in different values depending on the order in which threads are executed and preempted.

**Synchronization** The protection against race conditions in critical regions.

# 11 | Networks & Sockets

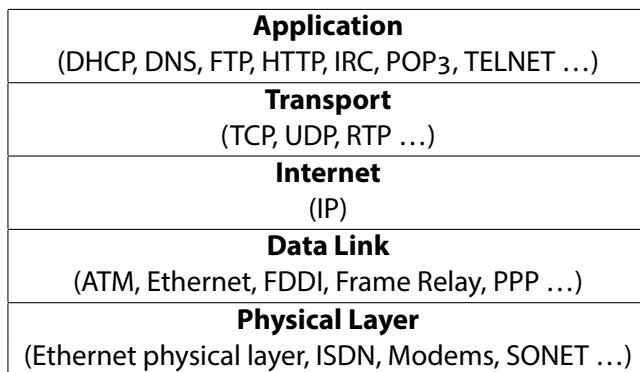
IN THIS CHAPTER we will examine the basics of having two or more computers talk to each other over an electronic or radio-frequency connection. Having computers connected to a **network** is almost taken for granted in this day and age, with the quintessential network being the Internet. There are plenty of other networks, from telephones (both cellular and land-line) to the local-area networks that share data and applications in businesses and homes.

We will start with an introduction to networking basics from a programmer's perspective. We focus on the makeup and potential issues of network communication, and how they affect the performance and reliability of transmitting and receiving data. We then move to the de facto standard for programming network applications: **Berkeley Sockets**. Berkeley Sockets is an **Application Programming Interface (API)**. An API is an abstraction, furnished by an Operating System or library that exposes a set of functions and data structures to do some specific tasks.

## 11.1 Introduction

A network is a connection of two or more computers such that they can share information. While networking is ubiquitous today, some details are important to understand before network-aware applications can be adequately written.

Networks, like Operating Systems, can be broken up into several layers to abstract specific details and to allow a network to be made up of heterogeneous components. There is a formal seven-layer model for networking known as the *OSI model* that serves as a set of logical divisions between the different components into which a network can be subdivided. The Internet uses five of these layers and results in the diagram shown in Figure 11.1. The bottom-most layer represents the actual electronic (physical) connection. While wireless networks are common, most networks will consist of a closed electrical circuit that sends signals and needs to



**Figure 11.1:** The Internet Layer Model.

Bits	0–3	4–7	8–15	16–18	19–31		
<b>0</b>	Version	Header Length	Type of Service	Total Length			
<b>32</b>	Identification			Flags	Fragment Offset		
<b>64</b>	Time to Live		Protocol	Header Checksum			
<b>96</b>	Source Address						
<b>128</b>	Destination Address						
<b>160</b>	Options						
<b>192–</b>	Data						

**Figure 11.2:** Layout of an IP packet.

resolve issues of message collision. On top of this layer comes an agreement on how to send data by way of these electronic signals. The data is organized into discrete chunks called **packets**. How these packets are organized needs to be standardized for communication to be intelligible to the recipient. This standard agreement on how to do something is known as a **protocol**. The protocol governing the Internet is appropriately known as the *Internet Protocol* (IP).

The Internet Protocol defines a particular packet, as illustrated in Figure 11.2. The first 192 bits (24 bytes) form a header that indicates details such as the destination and source of the packet. To identify specific computers in a network, each computer is assigned a unique **IP address**, a 32- or 128-bit number. Because of the way IP addresses are allocated, it often works out that many computers share a single IP address, but the details of how this works are beyond the scope of this text.



**Figure 11.3:** Each layer adds its own header to store information.

The different sizes of IP addresses come about as a result of two different standards. IPv4 is the current system using 32-bit addresses. Addresses are represented in the familiar “dotted decimal” notation such as 127.0.0.1. As networked devices keep growing, an effort to make sure that every device can have a unique address spawned the IPv6 standard. With this standard’s 128-bit addresses, there is little chance of running out anytime in the foreseeable future.

Packets sent via IP make no guarantees about arrival or receipt order. As a theoretical concept, such a guarantee is impossible to make. Imagine that Alice sends a message to Bob and wants to know that Bob receives it, so she asks Bob to send a reply when he gets it. A week passes, and Alice hears nothing from Bob and begins to wonder. However, she is met with an unanswerable question: Did Bob not get her message, or did she not get Bob’s reply?

The good news in regard to this conundrum is that modern networks are usually reliable enough that “dropped” packets are rare. A protocol that does nothing to guarantee receipt is known as a **Datagram** protocol. The term Datagram comes from a play on telegram, which also had no guarantee about receipt.

While mostly reliable communication might be adequate for some uses, the majority of applications want reliable, order-preserving communication. Email, for instance, would be useless if the message arrived garbled and with parts missing. Since we assume a mostly-reliable network, we can do better. A protocol implemented on top of IP called **Transmission Control Protocol** (TCP) attempts to account for the occasional lost or out-of-order packet. It does this through acknowledgment messages and a sequence number attached to each packet to indicate relative order. To do this, TCP needs to add this sequence number to the packet, and so it adds its own header following the IP header. Figure 11.3 illustrates the concatenation of headers done by each layer. (Note that the figure assumes the data link layer is Ethernet.)

Some applications, like streaming audio or video, or data for video games, can tolerate the occasional lost or out-of-order packet. Not worrying about receipt or order allows for larger amounts of data to be sent at faster rates. In fact, no connection is even necessarily made between the sender and the recipient. The most common

<b>Server</b>	<b>Client</b>
socket()	
bind()	connect()
listen()	
accept()	
send() and recv()	
close()	

**Figure 11.4:** The functions of Berkeley Sockets divided by role.

of these so called **connectionless** protocols is **UDP**, or the **User Datagram Protocol**. Using **UDP** provides nearly raw access to the **IP** packet without the overhead, or guarantees, associated with **TCP**.

The topmost layer is the *Application Layer*. This is implemented on top of **TCP** or **UDP** and consists of a protocol for programs to talk to each other. The protocol might be binary, like the Oscar protocol for AOL's Instant Messenger, or it might be text such as the famous Hypertext Transfer Protocol (**HTTP**) used by Web browsers to ask for Web pages from Web servers.

While **IP** addresses are convenient for computers to store and manipulate, they are not generally easily remembered by humans. People would rather have a name or other word to associate with a particular computer. On the Internet, each Web site has a unique *Domain Name* that corresponds to a particular **IP** address of the Web server. The World Wide Web provides a set of servers to facilitate translating a name into an **IP** address, a process known as domain name resolution. A **Domain Name Server (DNS)** provides a way to look up a particular **IP** address based upon the parts of a domain name.

## 11.2 Berkeley Sockets

Unix-like Operating Systems try to interact with devices, files, and networks in a uniform fashion by treating them all as part of the filesystem. By doing this, the programmer's interaction with the **I/O** device is uniform: The device can be opened, read or written, and closed. **Berkeley Sockets** serve to implement this abstraction for network communication.

The functions of the Berkeley Sockets API are listed in Figure 11.4. A socket is an **I/O** device representing a connection to a computer via a network. The **socket()**

call creates a file descriptor representing the connection to be used by the other functions. Berkeley Sockets distinguish between a listening server and a connecting client. Listing 11.1 gives the code for a server that simply replies with “Hello there!” to any program that connects to that particular machine and port pair. A **port** is an application-reserved connection point on a particular IP address.

Due to the fact that a network communication failure is much more likely than failures with many other I/O operations, even a simple program ends up with many lines of code. If something goes wrong, every function will return a negative number and set `errno`, the global error code, to the appropriate value. Using `perror()` converts `errno` into a (sometimes) useful error message and prints it to the screen. For the `send()` and `recv()` functions, the returned value indicates how many bytes were actually sent. If the data is too large, multiple calls may be needed to handle it.

We can connect to the server in Listing 11.1 by using `telnet`, which emulates a terminal and connects to a specified address and port. If the server is running on the local machine, the following output would be seen:

```
(1) thot $ telnet localhost 1100
Trying 127.0.0.1...
Connected to localhost.
Escape character is '^]'.
Hello there!
Connection closed by foreign host.
```

Berkeley Sockets also support connectionless protocols like UDP. The `sendto()` and `recvfrom()` calls take extra parameters that specify the address of the recipient or the sender. There is no need to do anything other than set up a socket to use them.

### 11.3 Sockets and Threads

The theme of Chapter 9 was that threads are useful when a program wants to do tasks in parallel. However, if all of those threads are CPU-bound, we may not see any performance advantage on a single-processor machine. Fortunately, we quickly realized that many modern programs are I/O-bound, and while the I/O operations are blocked, another thread could be scheduled to run.

While the I/O operations on a local machine may seem slow to the CPU, they are nothing compared to the delay incurred by doing network communication. Thus,

```
#include <string.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <stdio.h>
#include <stdlib.h>

#define MYPORt 1100

int main() {
    int sfd, connfd, amt = 0;
    struct sockaddr_in addr;
    char buf[1024];

    if((sfd=socket(PF_INET, SOCK_STREAM, 0)) < 0) {
        perror("Socket\u2014failed");
        exit(EXIT_FAILURE);
    }

    memset(&addr, 0, sizeof(addr));
    addr.sin_family = AF_INET;
    addr.sin_port = htons(MYPORt);
    addr.sin_addr.s_addr = INADDR_ANY;

    if(bind(sfd, (struct sockaddr *)&addr, sizeof(addr)) < 0) {
        perror("Bind\u2014failed");
        exit(EXIT_FAILURE);
    }
    if(listen(sfd, 10) < 0) {
        perror("Listen\u2014failed");
        exit(EXIT_FAILURE);
    }
    if((connfd=accept(sfd, NULL, NULL)) < 0) {
        perror("Accept\u2014failed");
        exit(EXIT_FAILURE);
    }

    strcpy(buf, "Hello\u2014there!\n");
    while(amt < strlen(buf)) {
        int ret = send(connfd, buf+amt, strlen(buf)-amt, 0);
        if(ret < 0) {
            perror("Send\u2014failed");
            exit(EXIT_FAILURE);
        }
        amt += ret;
    }
    close(connfd);
    close(sfd);
    return 0;
}
```

**Listing 11.1:** A server using sockets to send “Hello there!”.

we frequently see programs that access remote machines written in a multithreaded fashion. One such application is the multithreaded Web server.

A Web server does not necessarily meet the traditional idea of a multithreaded application, since each page request is likely to be independent and there is no real advantage to sharing a single address space. However, the real benefit comes from the idea of a main thread that accepts connections and then spawns a worker thread to take care of the I/O operations. The thread creation cost should be much cheaper than the overhead needed to create a full process, and thus the server can utilize CPU time more effectively. Other performance enhancements can also be incorporated. With a single address space, the server is free to make a shared cache in memory of frequently accessed files, reducing the need for disk I/O.

The more obvious marriage of threads and sockets comes from the frequent need to do asynchronous, bidirectional communication. Consider writing a simple instant messaging program that can talk to another program across a network. Which of the two instances of the program is the server and which is the client? Both programs want to send data at the request of the user. If the program was written in a single threaded fashion, it would need to have a sequence of `send()`s and `recv()`s, but their order dictates who can talk and who must listen.

The solution to this problem comes by having two separate threads, one devoted to sending and the other to receiving messages. This way, the receiving thread can remain safely blocked and the user can send any number of messages without delay.

## 11.4 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Application Programming Interface/API** An abstraction, furnished by an Operating System or library, which exposes a set of functions and data structures to do some specific tasks.

**Internet Protocol (IP) Address** A number that represents a particular computer on the Internet.

**Network** A group of computers wired to talk to each other.

**Packet** The unit of data transmission on a network.

**Port** An application-reserved connection point on a particular IP address.

**Protocol** An agreement on how data should be sent.

**Socket** An abstraction representing a connection to another computer over a network.

# A | The Intel x86 32-bit Architecture

THE INTEL X86 32-bit architecture is an example of a **Complex Instruction Set Computer** (CISC). While more recently designed CPUs have a simplified set of minimal operations, a CISC computer has many different instructions, special purpose registers, and complex addressing modes.

Figure A.1 gives a list of the general purpose registers. The first six can be used for most any purpose, although some instructions expect certain values to be in a particular register. %esp and %ebp are used for managing the stack and activation records. The program counter is %eip, which is read-only and can only be set via a jump instruction. The results of comparisons for conditional branches are stored in the register EFLAGS.

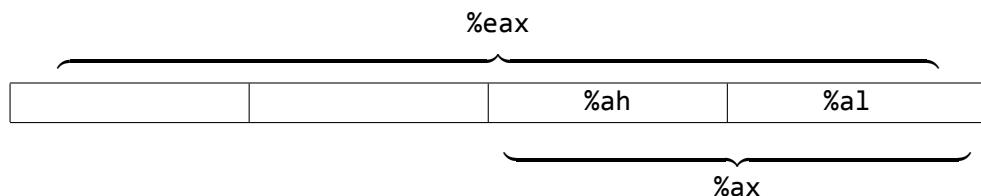
The registers %eax, %ebx, %ecx, and %edx each have subregister fields as shown in Figure A.2. For example, the lower (least-significant) 16 bits of %eax is known as %ax. %ax is further subdivided into two 8-bit registers, %ah (high) and %al (low). There is no name for the upper 16-bits of the registers. Note that these subfields are all part of the %eax register, and not separate registers, so if you load a 32-bit quantity and then read %ax, you will read the lower 16-bits of the value in %eax. The same applies to the other three registers.

Operations in x86 usually have two operands, which are often a *source* and a *destination*. For arithmetic operations like [add](#), these serve as the addends, and the addend in the destination position stores the sum. In mathematical terms, for two registers  $a$  and  $b$ , the result of the operation is  $a = a + b$ , overwriting one of the original values.

---

%eax	Accumulator
%ebx	Base
%ecx	Counter
%edx	Data
%esi	String Source
%edi	String Destination
%esp	Stack Pointer
%ebp	Base or Frame Pointer
%eip	Instruction Pointer
EFLAGS	Flag register

**Figure A.1:** The 32-bit registers.



**Figure A.2:** %eax, %ebx, %ecx, and %edx have subregister fields.

sub	Subtract
add	Add
and	Bitwise AND
push	Push a value onto the stack
pop	Pop a value off of the stack
mov	move a value
call	call a function
leave	clean up a stack frame
ret	return from a function
lea	compute an address (pointer)

**Figure A.3:** The instructions used in this book.

## A.1 AT&T Syntax

AT&T syntax is used by default in `gcc` and `gdb`. In AT&T assembler syntax, every operation code (opcode) is appended with the type of its operands:

b	byte (8-bit)
w	word (16-bit)
l	long (32-bit)
q	quad (64-bit)

After the opcode, the first operand is the source and the second operand is the destination. Memory dereferences are denoted by ( ). Listing A.1 gives an example of a “hello world” program as produced by `gcc`.

## A.2 Intel Syntax

Intel assembler syntax is the default syntax of the Intel documentation, Microsoft’s compilers and assemblers (`MASM`), and `NASM` – the Netwide Assembler. Instead of appending the operand size to the opcode, Intel syntax uses C-like casts to describe the size of operands. The type sizes are spelled out:

```
.file  "asm.c"
      .section .rodata.str1.1,"aMS",@progbits,1
.LC0:
      .string "hello\u0000world!"
      .text
.globl main
      .type  main, @function
main:
      pushl %ebp
      movl %esp, %ebp
      subl $8, %esp
      andl $-16, %esp ;1111 1111 1111 0000
      subl $16, %esp
      movl $.LC0, (%esp)
      call puts
      movl $0, %eax
      leave
      ret
```

**Listing A.1:** Hello world in AT&T assembler syntax.

```
main:
      push  ebp
      mov   ebp, esp
      sub   esp, 8
      and   esp, -16 ;1111 1111 1111 0000
      sub   esp, 16
      mov   DWORD PTR [esp], .LC0
      call  puts
      mov   eax, 0
      leave
      ret
```

**Listing A.2:** Hello world in Intel assembler syntax.

BYTE	1 byte
WORD	2 bytes
DWORD	4 bytes (double word)
QWORD	8 bytes (quad word)

Intel syntax orders the operands completely in reverse from the AT&T convention. The first operand is the *destination*, the second operand is the source. Dereferences are denoted by [ ]. Listing A.2 gives a sample of the same “hello world” program as in Listing A.1 rewritten in Intel syntax.

## A.3 Memory Addressing

One of the biggest surprises in the x86 instruction set for RISC assembly programmers is the memory addressing model. Architectures such as MIPS require that all Arithmetic/Logical Unit (ALU) operations have only registers as operands. Moving to and from memory requires explicit *load* and *store* instructions. However, most x86 instructions may take one operand as a memory location, as long as the other (if necessary) is a register. You may not have both a source and a destination that are in memory.

Memory addresses can be constructed from four parts: a signed offset (constant), a base (register), an index (register), and a scale (constant: 1, 2, 4, or 8). The resulting address is determined as: *base* + *index* × *scale* + *offset*. As an example, we can choose %ebx as the base, %eax as the index, 4 as the scale, and 16 as the offset. In AT&T syntax this would be expressed as: 16(%ebx,%eax,4). In Intel syntax, it would be written as: [ebx+eax\*4+16]. If the offset is zero or the scale is one it can be omitted, as can either of the registers if they are unnecessary.

## A.4 Flags

Suppose we have a conditional statement in C such as `if(x == 0) { ... }` which we could translate into x86 using the compare and jump-if-equals instructions as:

```
cmpl $0, %eax
je .next
; ...
.next:
```

One thing that is not immediately apparent in the code is how the branch “knows” the result of the previous compare instruction. The answer is that the compare instruction has a side-effect: It sets the `%eflags` register based on the result of the comparison.

The `%eflags` register is a collection of single-bit boolean variables that represent various pieces of state beyond the normal result. Many instructions modify `%eflags` as a part of their operation. Some arithmetic instructions like addition set flags if they overflow the bounds of the destination. The conditional jumps consume the state of various flags as the condition on which to branch. In the above example, the jump-equals instruction actually checks the value of the special zero flag (ZF) that is part of `%eflags`. In fact, the `je` instruction is actually a pseudonym for the `jz` instruction: jump if the zero flag is set.

The side-effect of an operation setting flags can lead to confusing code. Consider the listing:

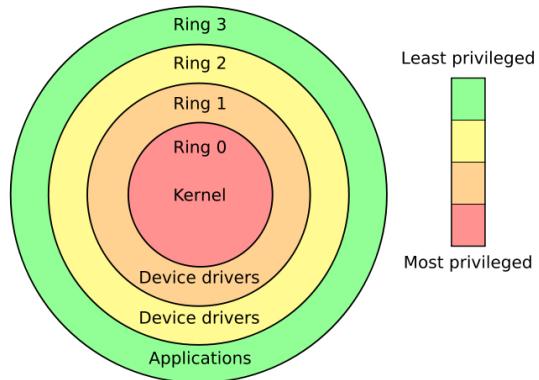
```
test    %eax, %eax
je      .next
; ...
.next:
```

This is functionally equivalent to the version that used `cmpl` above. The `test` instruction computes the bitwise-AND of the two arguments, in this case both are the register `%eax`. Since anything AND itself is going to be itself, this seems to be a no-op. But `test` takes the result of that AND and sets the ZF based on it. To learn about these side-effects, it is always handy to have the instruction set manual nearby.

A good compiler will probably generate the listing that uses `test` rather than `cmp` since the immediate 0 takes up 4 bytes of representation that are not needed in the encoding of `test`. The smaller code is generally preferred for performance (caching) reasons.

## A.5 Privilege Levels

To keep user programs and the kernel separate, x86 processors have four different privilege rings that processes may run in (Figure A.4). The kernel itself runs as the most privileged process in ring 0. Ring 1 is usually reserved for drivers, which can be thought of as kernel processes. Ring 2 is used for either drivers or user libraries. The final ring is for unprivileged user programs. The general protection rule is that a process running in a particular ring may access the data of the rings above but not



**Figure A.4:** Privilege rings on an x86 processor.

below it. This allows the OS to modify anything and to be protected from malicious or buggy user programs.

## B | Debugging Under gdb

OFTEN YOU WILL WANT to be able to examine a program while it executes, usually to find a bug or learn more about what a program is doing. A **debugger** is a program that gives its user control over the execution and data of another program while running.

With the very first programs a new programmer writes come the first problems, known as “bugs.” These are logical errors that the compiler cannot check. To track down these bugs, a concept known as *debugging*, beginners often use print statements. While such statements certainly work, this approach is often tedious and time-consuming due to frequent recompiling and re-execution. Adding print statements can also mask bugs. Since these are additional function calls, they provide legitimate stack frames where before an array out of bounds problem might have been a segmentation violation. With multithreaded programs, print statements can change timings and context switch points resulting in a different execution order that may hide a concurrency issue or race condition.

Print statements generally come in two kinds: The “I’m here” variety, which indicates a particular path of execution, and the “x is 5” variety that examine the contents of variables. The first type is an attempt at understanding the decisions that a program makes in its execution, i.e., what branches were taken. The path of a program through its flowchart representation is known as its **control flow**. The second type explores data values at certain points of execution.

We can do both of these things with a debugger, without the need to modify the source code. We may, however, choose to modify the executable at compile-time to provide the debugger with extra information about the program’s structure and correspondence to the source. Remember that a compiled executable has only memory locations and machine instructions. Gone are symbols like x or count for variables; all that remains are registers and memory addresses. For a debugger to be more helpful, we can choose to add extra information to the executable while

compiling that equates those particular addresses back to the human-readable names they originally had. For `gcc`, the `-g` flag indicates that the compiler should augment the symbol table and executable code with debugging information.

Debuggers come in all shapes and sizes. Some might be part of an Integrated Development Environment (IDE); others might be stand-alone. They can be graphical or text-based. For this chapter, we will use the stand-alone, textual debugger `gdb`. All of the concepts and commands we establish will be portable to other debuggers, usually without many, if any, differences. We will begin our discussion of the role of the debugger by discussing how to examine (and stop) the flow of control in a program.

## B.1 Examining Control Flow

The simplest way to control the execution of a program is to make it pause. We can ask a program to stop by having the debugger insert a **breakpoint** at a specific location. Locations can be specified in many ways, such as:

**Function Name** Execution will stop *before* the specified function is executed.

**Line Number** Execution will stop *before* the specified line of code is executed.

**Absolute Address** Execution will stop *before* the instruction at that address is executed.

The first two specifications require the executable to have additional information that may not necessarily be there. The line number information requires the executable to have been built with debugging information. The function names are usually part of the symbol table even without a special compilation but the symbol table may be “stripped” out after compilation. Specifying an absolute address always works, but gives no high-level language support.

Running `gdb` is as simple as specifying an executable to debug on the command line. If your program requires command line arguments, they can be specified by using the `--args` command line option, or by issuing the `set args` command once `gdb` has started up.

Once in `gdb`, the command to place a breakpoint is `break`, which can be abbreviated just by typing `b`. For example, a breakpoint could be placed at the `main()` function by typing `b main`. When the program is run via `run` or `r`, the program will immediately stop at the `main()` function. If there is debugging information, a

breakpoint can be set at a particular file and line number, separated by a colon: `b main.c:10` sets a breakpoint at line 10 in the `main.c` file.

It is also possible to put a breakpoint at an arbitrary instruction by specifying its address in memory. The syntax is `b *0x8048365`, where the hexadecimal number needs to be the start of an instruction. It is up to you to ensure that this address is valid. If it is not aligned to the start of an instruction, the program might crash. Section 8.2.2 gives some insight on how breakpoints are implemented and why placing a breakpoint in the middle of a multibyte instruction would be catastrophic to the program.

Once the program is stopped, you will probably want to examine some data, the topic of the next section. One useful thing to check, however, is what the call stack contains, i.e., what function calls led to the current place. This is called a **backtrace** and can be seen via the `backtrace` (abbreviated `back` or `bt`) command.

When you are finished with your examination, you have either found your bug and want to stop debugging, or you will need to continue on. Typing `run` or `r` allows you to restart the program from the beginning. The `quit` command exits `gdb`. If you want to `continue`, you can simply issue the command or the shortcut `c`. `Continuing` runs the program until it hits another breakpoint or the program ends, whichever comes first.

You may also want to step through the execution of the program, as if there were a breakpoint at each line. If the program was built with debugging information, you can use the commands `next` and `step`. These two commands are identical except for how they behave when they encounter a function call. The `step` command will go to the next source line inside the called function. The `next` command will skip over the function call and stop at the source code line immediately following the call. In other words, it will not leave the current function. Both `step` and `next` can be abbreviated with their first letter.

If the source code is not available and the program was not built with debugging information, `step` and `next` cannot be used. However, there are parallel commands for operating directly on the machine instructions. The `stepi` command goes to the next machine instruction, even if that is inside a separate function, whereas `nexti` skips to the next instruction following a `call` without leaving the current function. The abbreviation for `stepi` is `si` and for `nexti` is `ni`.

Each of the above control flow instructions (i.e., `continue`, `next`, `step`, `nexti`, and `stepi`) take an optional numerical argument. This number indicates a *repeat count*. The particular operation is performed that many times before control is returned to the debugger.

The technique of last resort is to interrupt the program with a `SIGINT` signal (see Section 8.2) by pressing `CTRL+C`. The signal will be handled by `gdb` and it will return control back to the user. One word of caution, however: it may be somewhat surprising where execution has stopped (it could be deep in a bunch of library calls), so it may be helpful to use `back` to get a backtrace and possibly some locations for regular breakpoints.

## B.2 Examining Data

With execution stopped and control transferred back to the debugger, most likely you will want to examine data values in memory. In `gdb` there are two primary commands for examining data, `print` and `x`. The `print` command will display the value of an expression, which can be written using C-style syntax:

```
(gdb) print 1+2  
$1 = 3
```

If the program was built with debugging symbols, you can write your expressions in terms of actual program variables to see what they contain. If your program is without such symbols, you can always look at register values. For example, on an x86 platform, you could display the contents of the register `%eax` by prefixing it with a dollar sign (\$) like so:

```
(gdb) print $eax  
$1 = 10
```

With the ability to use casts and dereferences, the `print` command is likely all that you need. However, the relative frequency with which you will want to look at the contents of some memory locations is high enough that there is a dedicated examine command, `x`. With the `x` command, the specified argument must be an address. By default, the contents are dumped as a hexadecimal number in the machine's native word size. It is possible to also specify a format, which acts as a typecast for the data. The type is specified by a single letter code following a forward slash. For instance:

```
(gdb) x 0x8048498  
0x8048498 <_IO_stdin_used+4>: 0x6c6c6548  
(gdb) x/d 0x8048498  
0x8048498 <_IO_stdin_used+4>: 1819043144
```

```
(gdb) x/x 0x8048498
0x8048498 <_IO_stdin_used+4>: 0x6c6c6548
(gdb) x/s 0x8048498
0x8048498 <_IO_stdin_used+4>: "Hello, world!"
```

Note that all four `x` commands operate upon the same address, but each has a different data interpretation. With `/d`, the number is printed out in decimal, rather than the hexadecimal (also obtainable via `/x`). Specifying `/s` will treat the address as the start of a C-style string and will attempt to print data until it encounters a null character.

The above example also illustrates a few of the output features of `gdb`. The first column is the address that is being examined, but in between the `<` and `>`, `gdb` attempts to map this address back to the nearest entry in the symbol table. In some cases, this can be quite useful, since with debugging symbols, individual variable names will be identified even when you know only an address (say from the contents of a pointer). Without full debugging symbols, or with a stripped executable, the output might not be correct, so common sense must always be used when interpreting the output.

## B.3 Examining Code

Whenever a breakpoint is encountered and control is returned to `gdb`, the current source line will be displayed if it is available. When the program contains debugging information and the source files are available, the debugger can operate in terms of the source code. You can see the source code around the current instruction pointer by using the `list` command. You can also specify a file and line number, as with breakpoints, or a function name.

When the source is unavailable, you can only see the machine instructions by using the command `disassemble`. Like `list`, `disassemble` will, by default, attempt to disassemble the instructions around `%eip`. You can additionally specify a region of memory addresses to dump. Care must be taken to ensure that the first address is actually the valid start of an instruction; otherwise, on a variable-length instruction set architecture, the disassembler could get confused.

One final trick is especially useful when `CTRL+C` is used to stop the program. Using the `examine` command with the `/i` format, you can disassemble an individual instruction at a certain location. To disassemble the instruction at the current instruction pointer location, do:

```
x/i $eip
```

## B.4 gdb Command Quick Reference

Command	Abbrv.	Description
help		Get help on a command or topic
set args		Set command-line arguments
run	r	Run (or restart) a program
quit	q	Exit gdb
break	b	Place a breakpoint at a given location
continue	c	Continue running the program after hitting a breakpoint
backtrace	bt	Show the function call stack
next	n	Go to the next line of source code without entering a function call
step	s	Go to the next line of source code, possibly entering a new function
nexti	ni	Go to the next instruction without entering a function call
stepi	si	Go to the next instruction, possibly entering a new function
print		Display the value of an expression written in C notation
x		Examine the contents of a memory location
list		List the source code of the program
disassemble	disas	List the machine code of the program

## B.5 Terms and Definitions

The following terms were introduced or defined in this chapter:

**Backtrace** A list of function calls that led to the current call; a stack dump.

**Breakpoint** A location in code where execution should stop or pause, usually used to transfer control to a debugger.

**Control Flow** The path or paths possible through a region of code as a result of decision (control) structures.

**Debugger** A program that controls and examines the execution and data of another program.

# References For Further Reading

- [1] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison Wesley, 1986.
- [2] Jeff Bonwick and Sun Microsystems. The slab allocator: An object-caching kernel memory allocator. In *USENIX Summer*, pages 87–98, 1994.
- [3] Jonathan Corbet, Alessandro Rubini, and Greg Kroah-Hartman. *Linux Device Drivers*. O'Reilly, 3rd edition, 2005. Available from: <http://lwn.net/Kernel/LDD3/>.
- [4] Intel Corporation. *Intel 64 and IA-32 Architectures Software Developer's Manual*, volume 1: Basic Architecture. 2007. Available from: <http://www.intel.com/design/processor/manuals/253665.pdf>.
- [5] Intel Corporation. *Intel 64 and IA-32 Architectures Software Developer's Manual*, volume 2A: Instruction Set Reference, A-M. 2007. Available from: <http://www.intel.com/design/processor/manuals/253666.pdf>.
- [6] Intel Corporation. *Intel 64 and IA-32 Architectures Software Developer's Manual*, volume 2B: Instruction Set Reference, N-Z. 2007. Available from: <http://www.intel.com/design/processor/manuals/253667.pdf>.
- [7] Tim Lindholm and Frank Yellin. *The Java Virtual Machine Specification*. Addison Wesley Longman, 2nd edition, 1999. Available from: <http://java.sun.com/docs/books/jvms/>.
- [8] Sandra Loosemore, Richard M. Stallman, Roland McGrath, Andrew Oram, and Ulrich Drepper. *The GNU C Library Reference Manual*. Free Software Foundation, 0.11 edition, 2007. Available from: <http://www.gnu.org/software/libc/manual/pdf/libc.pdf>.

- [9] Mark Mitchell, Jeffrey Oldham, and Alex Samuel. *Advanced Linux Programming*. New Riders Publishing, 2001. Available from: <http://www.advancedlinuxprogramming.com/alp-folder>.
- [10] David A. Patterson and John L. Hennessy. *Computer Organization and Design: The Hardware/Software Interface*. Morgan Kaufmann, 3rd edition, 2007.
- [11] Richard Stallman, Roland Pesch, Stan Shebs, et al. *Debugging with gdb*. Free Software Foundation, 9th edition, 2006. Available from: <http://sourceware.org/gdb/current/onlinedocs/gdb.html>.
- [12] Andrew S. Tanenbaum. *Modern Operating Systems*. Prentice Hall, 2nd edition, 2001.

# Index

A page number appearing in **bold** indicates an end-of-chapter definition.

## A

Activation record **41, 49**  
Address **2, 7**  
Address space **37, 39, 78**  
Alignment **42, 49**  
API **98, 104**  
AT&T Syntax **108**

## B

Back patching **28**  
Back-patching **23**  
Backtrace **115, 119**  
Berkeley Sockets **98, 101**  
Best Fit **56**  
Bitmap **50, 51, 65**  
Blocked **79, 83**  
Breakpoint **114, 119**  
BSS section **27**  
Buddy Allocator **50, 64**  
Buffer-overrun **49**

## C

C Standard Library **20**  
Call table **34, 36**

Callee-saved **41, 49**  
Caller-saved **41, 49**  
Calling convention **42, 49**  
CISC **106**  
Coalesce **58, 65**  
Compiler **17, 19, 29**  
Condition variable **94**  
Context **68, 76, 79**  
    Switch **68, 76, 79**  
Control flow **113, 119**  
CPU bound **79, 83**  
Critical region **89, 97**

## D

Data segment **26**  
Datagram **100**  
Deadlock **85, 93, 97**  
Debugger **74–76, 113, 119**  
Deduplication **23, 29**  
Dereference **4, 7**  
DLL Hell **24**  
Domain Name Server **101**  
Dynamic **8, 16**

**E**

Exec header **26**  
 Executable file **25, 29**

**F**

File descriptor **68**  
 First Fit **56**  
 Fragmentation  
   External **56, 65**  
   Internal **52, 65**  
 Frame **41, 49**  
   Pointer **42, 49**  
 Function pointer **30, 36, 74**

**G**

Garbage collection **59, 65**  
 gcc **11, 17, 42, 108, 114**  
 gdb **74, 108, 113–119**

**H**

Header file **18, 29**  
 Heap **38, 50, 65**

**I**

I/O bound **80, 83**  
 Intel syntax **108**  
 Interrupt **68, 76, 79**  
   Vector **68, 76**  
 IP Address **99, 104**

**J**

Jump Table **34**

**K**

Kernel **66, 76**  
   Mode **68**  
   Space **66, 76**

**L**

Libraries **20**  
 Library **29**  
 Lifetime **8, 16**  
 Link loader **22**  
 Linked list **50**  
 Linker **17, 20–25, 29**  
 Linking  
   Dynamic **20, 22–24**  
   Dynamic loading **24–25**  
   Static **20–22**  
 Loader **29, 71**

**M**

Macro **18, 29**  
 Magic number **26**  
 Memory leak **59, 65**  
 Multiprogramming **79, 84**  
 Mutex **90**

**N**

Network **98, 104**  
 Next Fit **56**

**O**

Object file **19, 29**  
 Operating System **66, 76**  
 Optimized **19**  
 Ordinal **34, 36**

**P**

Packet **99, 104**  
 Page **38, 39**  
 Plugins **25**  
 Pointer **2, 7**  
 Pointer Arithmetic **6, 7**  
 Poll **72, 76**

- Port **102, 104**  
 Preemption **79, 84**  
 Preprocessor **17–19, 29**  
 Process **37, 39, 78**  
   Identifier **70, 76**  
 Producer/Consumer problem **91**  
 Protocol **99, 105**  
   Connectionless **101**
- Q**  
 Quick Fit **57**
- R**  
 Race condition **89, 97**  
 Return address **41, 49**  
 Run Length Encoding **54**  
 Run-length Encoding **65**
- S**  
 Scheduler **79, 84**  
 Scope **8, 16**  
 Semaphore **94**  
 Shadowing **10, 16**  
 Shared Object **22**  
 Signal **72, 76**  
 Socket **105**  
 Sparse **53**  
 Stack **38, 40, 49**  
   Pointer **40, 49**  
 Static **8, 16**  
 String table **27**  
 Symbol **8, 16**  
 Symbol table **27**  
 Synchronization **85, 97**  
 System **17**  
 System Call **66**  
 System call **77**
- T**  
 Text segment **26**  
 Thread **80, 84**  
   Kernel **81, 83**  
   User **81, 84**  
 Transmission Control Protocol **100**  
 Trap **68, 77**
- U**  
 User Datagram Protocol **101**  
 User program **66, 77**  
 User space **66, 77**
- V**  
 Variable  
   Automatic **10–11**  
   Global **10**  
   Local **8, 16**  
   Register **11–12**  
   Static global **15**  
   Static local **12–15**  
   Volatile **15**  
 Variadic function **46**  
 Virtual Memory **37, 39**  
 von Neumann Architecture **20, 29**
- W**  
 Worst Fit **56**  
 Wrapper function **66, 77**

# Colophon

**col·o·phon**, noun, 1: an inscription placed at the end of a book or manuscript usually with facts relative to its production

—Merriam-Webster's Collegiate Dictionary, Eleventh Ed.

This text was typeset in  $\text{\LaTeX}$  2 $\varepsilon$  using lualatex under the MiK $\text{\TeX}$  2.9 system on Windows 7.  $\text{\TeX}$  is a macro package based upon Donald Knuth's  $\text{\TeX}$  typesetting language.<sup>1</sup>  $\text{\TeX}$  was originally developed by Leslie Lamport in 1985.<sup>2</sup> MiK $\text{\TeX}$  is maintained and developed by Christian Schenk.<sup>3</sup>

The typefaces are Minion Pro, Myriad Pro, and Consolas. Illustrations were edited in Adobe® Illustrator® and the final PDF was touched-up in Adobe® Acrobat®.

## About the Author

Jonathan Misurda is a Lecturer in the Department of Computer Science at the University of Pittsburgh, where he did his Ph.D. in the Software Testing aspect of Software Engineering. His heart lies, however, in Computer Science Education.

Jonathan's hobbies and interests can vary seemingly on the month, but his experiences in preparing this text have sparked an interest in Graphic Design and Typography. He is always on the lookout to learn how to better present the topics he explains and the information he presents. His goal is to make both his teaching and his book look better and to express their content more clearly. To this end he jokes that he is in his "Knuth" phase.

---

<sup>1</sup> <http://www-cs-faculty.stanford.edu/~knuth/>

<sup>2</sup> <http://www.latex-project.org/>

<sup>3</sup> <http://www.miktex.org/>



[Dashboard](#) / [My courses](#) / [Computer Engineering & IT](#) / [CEIT-Even-sem-21-22](#) / [OS-even-sem-21-22](#) / [7 February - 13 February](#)

/ [Quiz-1: 10 AM](#)

**Started on** Saturday, 12 February 2022, 10:00:21 AM

**State** Finished

**Completed on** Saturday, 12 February 2022, 11:25:53 AM

**Time taken** 1 hour 25 mins

**Grade** 4.94 out of 10.00 (49%)

**Question 1**

Complete

Mark 0.00 out of 0.50

Select all the correct statements about code of bootmain() in xv6

```

void
bootmain(void)
{
    struct elfhdr *elf;
    struct proghdr *ph, *eph;
    void (*entry)(void);
    uchar* pa;

    elf = (struct elfhdr*)0x10000; // scratch space

    // Read 1st page off disk
    readseg((uchar*)elf, 4096, 0);

    // Is this an ELF executable?
    if(elf->magic != ELF_MAGIC)
        return; // let bootasm.S handle error

    // Load each program segment (ignores ph flags).
    ph = (struct proghdr*)((uchar*)elf + elf->phoff);
    eph = ph + elf->phnum;
    for(; ph < eph; ph++){
        pa = (uchar*)ph->paddr;
        readseg(pa, ph->filesz, ph->off);
        if(ph->memsz > ph->filesz)
            stobs(pa + ph->filesz, 0, ph->memsz - ph->filesz);
    }

    // Call the entry point from the ELF header.
    // Does not return!
    entry = (void(*)(void))(elf->entry);
    entry();
}

```

Also, inspect the relevant parts of the xv6 code. binary files, etc and run commands as you deem fit to answer this question.

- a. The elf->entry is set by the linker in the kernel file and it's 0x80000000
- b. The condition if(ph->memsz > ph->filesz) is never true.
- c. The readseg finally invokes the disk I/O code using assembly instructions
- d. The elf->entry is set by the linker in the kernel file and it's 0x80000000
- e. The kernel ELF file contains actual physical address where particular sections of 'kernel' file should be loaded
- f. The kernel file gets loaded at the Physical address 0x10000 +0x80000000 in memory.
- g. The kernel file in memory is not necessarily a continuously filled in chunk, it may have holes in it.
- h. The elf->entry is set by the linker in the kernel file and it's 8010000c
- i. The kernel file gets loaded at the Physical address 0x10000 in memory.
- j. The stobs() is used here, to fill in some space in memory with zeroes
- k. The kernel file has only two program headers

The correct answers are: The kernel file gets loaded at the Physical address 0x10000 in memory., The kernel file in memory is not necessarily a continuously filled in chunk, it may have holes in it., The elf->entry is set by the linker in the kernel file and it's 8010000c, The readseg finally invokes the disk I/O code using assembly instructions, The stosb() is used here, to fill in some space in memory with zeroes, The kernel ELF file contains actual physical address where particular sections of 'kernel' file should be loaded, The kernel file has only two program headers

**Question 2**

Complete

Mark 0.50 out of 0.50

What's the trapframe in xv6?

- a. A frame of memory that contains all the trap handler's addresses
- b. A frame of memory that contains all the trap handler code's function pointers
- c. The IDT table
- d. The sequence of values, including saved registers, constructed on the stack when an interrupt occurs, built by hardware + code in trapasm.S
- e. The sequence of values, including saved registers, constructed on the stack when an interrupt occurs, built by code in trapasm.S only
- f. The sequence of values, including saved registers, constructed on the stack when an interrupt occurs, built by hardware only
- g. A frame of memory that contains all the trap handler code

The correct answer is: The sequence of values, including saved registers, constructed on the stack when an interrupt occurs, built by hardware + code in trapasm.S

**Question 3**

Complete

Mark 0.21 out of 0.50

Order the events that occur on a timer interrupt:

- |   |   |
|---|---|
| Jump to a code pointed by IDT                       | 3 |
| Jump to scheduler code                              | 2 |
| Select another process for execution                | 5 |
| Change to kernel stack of currently running process | 4 |
| Set the context of the new process                  | 6 |
| Save the context of the currently running process   | 1 |
| Execute the code of the new process                 | 7 |

The correct answer is: Jump to a code pointed by IDT → 2, Jump to scheduler code → 4, Select another process for execution → 5, Change to kernel stack of currently running process → 1, Set the context of the new process → 6, Save the context of the currently running process → 3, Execute the code of the new process → 7

**Question 4**

Complete

Mark 0.50 out of 0.50

Suppose a program does a scanf() call.

Essentially the scanf does a read() system call.

This call will obviously "block" waiting for the user input.

In terms of OS data structures and execution of code, what does it mean?

Select one:

- a. OS code for read() will move PCB of current process to a wait queue and call scheduler
- b. OS code for read() will call scheduler
- c. OS code for read() will move the PCB of this process to a wait queue and return from the system call
- d. read() will return and process will be taken to a wait queue
- e. read() returns and process calls scheduler()

The correct answer is: OS code for read() will move PCB of current process to a wait queue and call scheduler

**Question 5**

Complete

Mark 0.50 out of 0.50

In bootasm.S, on the line

```
ljmp    $(SEG_KCODE<<3), $start32
```

The SEG\_KCODE << 3, that is shifting of 1 by 3 bits is done because

- a. The value 8 is stored in code segment
- b. The code segment is 16 bit and only upper 13 bits are used for segment number
- c. While indexing the GDT using CS, the value in CS is always divided by 8
- d. The ljmp instruction does a divide by 8 on the first argument
- e. The code segment is 16 bit and only lower 13 bits are used for segment number

The correct answer is: The code segment is 16 bit and only upper 13 bits are used for segment number

**Question 6**

Complete

Mark 0.40 out of 0.50

Select Yes if the mentioned element should be a part of PCB

Select No otherwise.

**Yes      No**

Memory management information about that process

Process context

Function pointers to all system calls

PID

EIP at the time of context switch

List of opened files

PID of Init

Process state

Pointer to the parent process

Pointer to IDT

Memory management information about that process: Yes

Process context: Yes

Function pointers to all system calls: No

PID: Yes

EIP at the time of context switch: Yes

List of opened files: Yes

PID of Init: No

Process state: Yes

Pointer to the parent process: Yes

Pointer to IDT: No

**Question 7**

Complete

Mark 0.40 out of 1.00

Mark the statements, w.r.t. the scheduler of xv6 as True or False

**True      False**

The work of selecting and scheduling a process is done only in `scheduler()` and not in `sched()`

`swtch` is a function that saves old context, loads new context, and returns to last EIP in the new context

The variable `c->scheduler` on first processor uses the stack allocated entry.S

`sched()` and `scheduler()` are co-routines

The function `scheduler()` executes using the kernel-only stack

When a process is scheduled for execution, it resumes execution in `sched()` after the call to `swtch()`

`swtch` is a function that does not return to the caller

the control returns to `mycpu()->intena = intena; ()`; after `swtch(&p->context, mycpu()->scheduler);` in `sched()`

the control returns to `switchkvm();` after `swtch(&(c->scheduler), p->context);` in `scheduler()`

`sched()` calls `scheduler()` and `scheduler()` calls `sched()`

The work of selecting and scheduling a process is done only in `scheduler()` and not in `sched()`: True

`swtch` is a function that saves old context, loads new context, and returns to last EIP in the new context: True

The variable `c->scheduler` on first processor uses the stack allocated entry.S: True

`sched()` and `scheduler()` are co-routines: True

The function `scheduler()` executes using the kernel-only stack: True

When a process is scheduled for execution, it resumes execution in `sched()` after the call to `swtch()`

: True

`swtch` is a function that does not return to the caller: True

the control returns to `mycpu()->intena = intena; ()`; after `swtch(&p->context, mycpu()->scheduler);` in `sched()`:

False

the control returns to `switchkvm();` after `swtch(&(c->scheduler), p->context);` in `scheduler()`: False

`sched()` calls `scheduler()` and `scheduler()` calls `sched()`: False

**Question 8**

Complete

Mark 0.00 out of 0.50

Consider the following programs

**exec1.c**

```
#include <unistd.h>
#include <stdio.h>
int main() {
    exec("./exec2", "./exec2", NULL);
}
```

**exec2.c**

```
#include <unistd.h>
#include <stdio.h>
int main() {
    exec("/bin/ls", "/bin/ls", NULL);
    printf("hello\n");
}
```

Compiled as

```
cc  exec1.c -o exec1
cc  exec2.c -o exec2
```

And run as

```
$ ./exec1
```

Explain the output of the above command (./exec1)

Assume that /bin/ls , i.e. the 'ls' program exists.

Select one:

- a. "ls" runs on current directory
- b. Execution fails as the call to execl() in exec1 fails
- c. Execution fails as one exec can't invoke another exec
- d. Program prints hello
- e. Execution fails as the call to execl() in exec2 fails

The correct answer is: "ls" runs on current directory

**Question 9**

Complete

Mark 0.00 out of 0.50

For each line of code mentioned on the left side, select the location of sp/esp that is in use

**cli****in bootasm.S**

0x7c00 to 0

**readseg((uchar\*)elf, 4096, 0);****in bootmain.c**

Immaterial as the stack is not used here

**ljmp \$(SEG\_KCODE<<3), \$start32****in bootasm.S**

0x10000 to 0x7c00

**call bootmain****in bootasm.S**

The 4KB area in kernel image, loaded in memory, named as 'stack'

**jmp \*%eax****in entry.S**

0x7c00 to 0x10000

The correct answer is: **cli**

**in bootasm.S** → Immateral as the stack is not used here, **readseg((uchar\*)elf, 4096, 0);**

**in bootmain.c** → 0x7c00 to 0, **ljmp \$(SEG\_KCODE<<3), \$start32**

**in bootasm.S** → Immateral as the stack is not used here, **call bootmain**

**in bootasm.S** → 0x7c00 to 0, **jmp \*%eax**

**in entry.S** → The 4KB area in kernel image, loaded in memory, named as 'stack'

**Question 10**

Complete

Mark 0.63 out of 1.00

Select the correct statements about interrupt handling in xv6 code

- a. xv6 uses the 64th entry in IDT for system calls
- b. xv6 uses the 0x64th entry in IDT for system calls
- c. On any interrupt/syscall/exception the control first jumps in vectors.S
- d. The function trap() is called irrespective of hardware interrupt/system-call/exception
- e. Before going to altraps, the kernel stack contains upto 5 entries.
- f. Each entry in IDT essentially gives the values of CS and EIP to be used in handling that interrupt
- g. The trapframe pointer in struct proc, points to a location on kernel stack
- h. The function trap() is called only in case of hardware interrupt
- i. The trapframe pointer in struct proc, points to a location on user stack
- j. On any interrupt/syscall/exception the control first jumps in trapasm.S
- k. The CS and EIP are changed only after pushing user code's SS,ESP on stack
- l. The CS and EIP are changed only immediately on a hardware interrupt
- m. All the 256 entries in the IDT are filled

The correct answers are: All the 256 entries in the IDT are filled, Each entry in IDT essentially gives the values of CS and EIP to be used in handling that interrupt, xv6 uses the 64th entry in IDT for system calls, On any interrupt/syscall/exception the control first jumps in vectors.S, Before going to altraps, the kernel stack contains upto 5 entries., The trapframe pointer in struct proc, points to a location on kernel stack, The function trap() is called irrespective of hardware interrupt/system-call/exception, The CS and EIP are changed only after pushing user code's SS,ESP on stack

**Question 11**

Complete

Mark 0.50 out of 0.50

Order the sequence of events, in scheduling process P1 after process P0

timer interrupt occurs

2

Process P0 is running

1

context of P1 is loaded from P1's PCB

4

Process P1 is running

6

Control is passed to P1

5

context of P0 is saved in P0's PCB

3

The correct answer is: timer interrupt occurs → 2, Process P0 is running → 1, context of P1 is loaded from P1's PCB → 4, Process P1 is running → 6, Control is passed to P1 → 5, context of P0 is saved in P0's PCB → 3

**Question 12**

Complete

Mark 0.00 out of 1.00

Select the sequence of events that are NOT possible, assuming a non-interruptible kernel code

(Note: non-interruptible kernel code means, if the kernel code is executing, then interrupts will be disabled).

Note: A possible sequence may have some missing steps in between. An impossible sequence will have n and n+1th steps such that n+1th step can not follow n'th step.

Select one or more:

a. P1 running

P1 makes system call  
timer interrupt  
Scheduler  
P2 running  
timer interrupt  
Scheduler  
P1 running  
P1's system call return

b. P1 running

P1 makes system call  
system call returns  
P1 running  
timer interrupt  
Scheduler running  
P2 running

c. P1 running

P1 makes system call and blocks  
Scheduler  
P2 running  
P2 makes system call and blocks  
Scheduler  
P3 running  
Hardware interrupt  
Interrupt unblocks P1  
Interrupt returns  
P3 running  
Timer interrupt  
Scheduler  
P1 running

d. P1 running

keyboard hardware interrupt  
keyboard interrupt handler running  
interrupt handler returns  
P1 running  
P1 makes system call  
system call returns  
P1 running  
timer interrupt  
scheduler  
P2 running

e.

P1 running  
P1 makes system call  
Scheduler  
P2 running  
P2 makes system call and blocks  
Scheduler  
P1 running again

f. P1 running

P1 makes system call and blocks  
Scheduler  
P2 running  
P2 makes system call and blocks  
Scheduler  
P1 running again

The correct answers are: P1 running

P1 makes system call and blocks

Scheduler

P2 running

P2 makes system call and blocks

Scheduler

P1 running again, P1 running

P1 makes system call

timer interrupt

Scheduler

P2 running

timer interrupt

Scheuler

P1 running

P1's system call return,

P1 running

P1 makes system call

Scheduler

P2 running

P2 makes system call and blocks

Scheduler

P1 running again

**Question 13**

Complete

Mark 0.50 out of 0.50

Some part of the bootloader of xv6 is written in assembly while some part is written in C. Why is that so?

Select all the appropriate choices

- a. The code in assembly is required for transition to protected mode, from real mode; after that calling convention applies, hence code can be written in C
- b. The setting up of the most essential memory management infrastructure needs assembly code
- c. The code in assembly is required for transition to protected mode, from real mode; but calling convention was applicable all the time
- d. The code for reading ELF file can not be written in assembly

The correct answers are: The code in assembly is required for transition to protected mode, from real mode; after that calling convention applies, hence code can be written in C, The setting up of the most essential memory management infrastructure needs assembly code

**Question 14**

Complete

Mark 0.17 out of 0.50

The bootmain() function has this code

```
elf = (struct elfhdr*)0x10000; // scratch space  
readseg((uchar*)elf, 4096, 0);
```

Mark the statements as True or False with respect to this code.

In these statements 0x1000 is referred to as ADDRESS

**True      False**

- The value ADDRESS is changed to a 0 the program could still work
- If the value of ADDRESS is changed to a higher number (upto a limit), the program could still work
- This line loads the kernel code at ADDRESS
- If the value of ADDRESS is changed to a lower number (upto a limit), the program could still work
- This line effectively loads the ELF header and the program headers at ADDRESS
- If the value of ADDRESS is changed, then the program will not work

The value ADDRESS is changed to a 0 the program could still work: False

If the value of ADDRESS is changed to a higher number (upto a limit), the program could still work: True

This line loads the kernel code at ADDRESS: False

If the value of ADDRESS is changed to a lower number (upto a limit), the program could still work: True

This line effectively loads the ELF header and the program headers at ADDRESS: False

If the value of ADDRESS is changed, then the program will not work: False

**Question 15**

Complete

Mark 0.50 out of 1.00

Which parts of the xv6 code in bootasm.S bootmain.c , entry.S and in the codepath related to scheduler() and trap handling() can also be written in some other way, and still ensure that xv6 works properly?

Writing code is not necessary. You only need to comment on which part of the code could be changed to something else or written in another fashion.

Maximum two points to be written.

We can use a scheduling algorithm. We can use the kernel stack in scheduler function in entry.S and bootmain.c .

**Question 16**

Complete

Mark 0.13 out of 0.50

Select all the correct statements about zombie processes

Select one or more:

- a. If the parent of a process finishes, before the process itself, then after finishing the process is typically attached to 'init' as parent
- b. Zombie processes are harmless even if OS is up for long time
- c. A zombie process occupies space in OS data structures
- d. A process becomes zombie when it finishes, and remains zombie until parent calls wait() on it
- e. A process can become zombie if it finishes, but the parent has finished before it
- f. init() typically keeps calling wait() for zombie processes to get cleaned up
- g. A process becomes zombie when its parent finishes
- h. A zombie process remains zombie forever, as there is no way to clean it up

The correct answers are: A process becomes zombie when it finishes, and remains zombie until parent calls wait() on it, A process can become zombie if it finishes, but the parent has finished before it, A zombie process occupies space in OS data structures, If the parent of a process finishes, before the process itself, then after finishing the process is typically attached to 'init' as parent, init() typically keeps calling wait() for zombie processes to get cleaned up

◀ Extra Reading on Linkers: A writeup by Ian Taylor (keep changing url string from 38 to 39, and so on)

Jump to...



**Started on** Saturday, 20 February 2021, 2:51 PM

**State** Finished

**Completed on** Saturday, 20 February 2021, 3:55 PM

**Time taken** 1 hour 3 mins

**Grade** 7.30 out of 20.00 (37%)

#### Question 1

Partially correct

Mark 0.80 out of 1.00

Select all the correct statements about the state of a process.

- a. A process can self-terminate only when it's running ✓
- b. Typically, it's represented as a number in the PCB ✓
- c. A process that is running is not on the ready queue ✓
- d. Processes in the ready queue are in the ready state ✓
- e. It is not maintained in the data structures by kernel, it is only for conceptual understanding of programmers
- f. Changing from running state to waiting state results in "giving up the CPU" ✓
- g. A process in ready state is ready to receive interrupts
- h. A waiting process starts running after the wait is over ✗
- i. A process changes from running to ready state on a timer interrupt ✓
- j. A process in ready state is ready to be scheduled ✓
- k. A running process may terminate, or go to wait or become ready again ✓
- l. A process waiting for I/O completion is typically woken up by the particular interrupt handler code ✓
- m. A process waiting for any condition is woken up by another process only
- n. A process changes from running to ready state on a timer interrupt or any I/O wait

Your answer is partially correct.

You have selected too many options.

The correct answers are: Typically, it's represented as a number in the PCB, A process in ready state is ready to be scheduled, Processes in the ready queue are in the ready state, A process that is running is not on the ready queue, A running process may terminate, or go to wait or become ready again, A process changes from running to ready state on a timer interrupt, Changing from running state to waiting state results in "giving up the CPU", A process can self-terminate only when it's running, A process waiting for I/O completion is typically woken up by the particular interrupt handler code

**Question 2**

Incorrect

Mark 0.00 out of 1.00

For each line of code mentioned on the left side, select the location of sp/esp that is in use

`jmp *%eax`  
in entry.S

0x7c00 to 0x10000



`ljmp $(SEG_KCODE<<3), $start32`  
in bootasm.S

0x10000 to 0x7c00



`call bootmain`  
in bootasm.S

0x7c00 to 0x10000



`cli`  
in bootasm.S

0x7c00 to 0



`readseg((uchar*)elf, 4096, 0);`  
in bootmain.c

The 4KB area in kernel image, loaded in memory, named as 'stack'



Your answer is incorrect.

The correct answer is: `jmp *%eax`

`in entry.S` → The 4KB area in kernel image, loaded in memory, named as 'stack', `ljmp $(SEG_KCODE<<3), $start32`

`in bootasm.S` → Immateriel as the stack is not used here, `call bootmain`

`in bootasm.S` → 0x7c00 to 0, `cli`

`in bootasm.S` → Immateriel as the stack is not used here, `readseg((uchar*)elf, 4096, 0);`

`in bootmain.c` → 0x7c00 to 0

**Question 3**

Correct

Mark 0.25 out of 0.25

Order the following events in boot process (from 1 onwards)

Boot loader	2	✓
Shell	6	✓
BIOS	1	✓
OS	3	✓
Init	4	✓
Login interface	5	✓

Your answer is correct.

The correct answer is: Boot loader → 2, Shell → 6, BIOS → 1, OS → 3, Init → 4, Login interface → 5

**Question 4**

Partially correct

Mark 0.30 out of 0.50

Consider the following command and its output:

```
$ ls -lht xv6.img kernel
-rw-rw-r-- 1 abhijit abhijit 4.9M Feb 15 11:09 xv6.img
-rwxrwxr-x 1 abhijit abhijit 209K Feb 15 11:09 kernel*
```

Following code in bootmain()

```
readseg((uchar*)elf, 4096, 0);
```

and following selected lines from Makefile

```
xv6.img: bootblock kernel
dd if=/dev/zero of=xv6.img count=10000
dd if=bootblock of=xv6.img conv=notrunc
dd if=kernel of=xv6.img seek=1 conv=notrunc
```

```
kernel: $(OBJS) entry.o entryother initcode kernel.ld
$(LD) $(LDFLAGS) -T kernel.ld -o kernel entry.o $(OBJS) -b binary initcode entryother
$(OBJDUMP) -S kernel > kernel.asm
$(OBJDUMP) -t kernel | sed '1,/SYMBOL TABLE/d; s/ .* / /; /^$$/d' > kernel.sym
```

Also read the code of bootmain() in xv6 kernel.

Select the options that describe the meaning of these lines and their correlation.

- a. Although the size of the kernel file is 209 Kb, only 4Kb out of it is the actual kernel code and remaining part is all zeroes.
- b. The kernel is compiled by linking multiple .o files created from .c files; and the entry.o, initcode, entryother files ✓
- c. The kernel.ld file contains instructions to the linker to link the kernel properly ✓
- d. The bootmain() code does not read the kernel completely in memory
- e. readseg() reads first 4k bytes of kernel in memory
- f. Although the size of the xv6.img file is ~5MB, only some part out of it is the bootloader+kernel code and remaining part is all zeroes.
- g. The kernel.asm file is the final kernel file
- h. The kernel disk image is ~5MB, the kernel within it is 209 kb, but bootmain() initially reads only first 4kb, and the later part is not read as it is user programs.
- i. The kernel disk image is ~5MB, the kernel within it is 209 kb, but bootmain() initially reads only first 4kb, and the later part is read using program headers in bootmain(). ✓

Your answer is partially correct.

You have correctly selected 3.

The correct answers are: The kernel disk image is ~5MB, the kernel within it is 209 kb, but bootmain() initially reads only first 4kb, and the later part is read using program headers in bootmain(), readseg() reads first 4k bytes of kernel in memory, The kernel is compiled by linking multiple .o files created from .c files; and the entry.o, initcode, entryother files, The kernel.ld file contains instructions to the linker to link the kernel properly, Although the size of the xv6.img file is ~5MB, only some part out of it is the bootloader+kernel code and remaining part is all zeroes.

**Question 5**

Partially correct

Mark 0.50 out of 1.00

```
int f() {  
    int count;  
    for (count = 0; count < 2; count++) {  
        if (fork() == 0)  
            printf("Operating-System\n");  
    }  
    printf("TYCOMP\n");  
}
```

The number of times "Operating-System" is printed, is:

Answer:

The correct answer is: 7.00

**Question 6**

Partially correct

Mark 0.40 out of 0.50

Select Yes/True if the mentioned element must be a part of PCB

Select No/False otherwise.

**Yes****No**

<input checked="" type="radio"/>	<input checked="" type="radio"/>	PID	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Process context	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	List of opened files	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Process state	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Parent's PID	✗
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Pointer to IDT	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Function pointers to all system calls	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Memory management information about that process	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Pointer to the parent process	✗
<input checked="" type="radio"/>	<input checked="" type="radio"/>	EIP at the time of context switch	✓

PID: Yes

Process context: Yes

List of opened files: Yes

Process state: Yes

Parent's PID: No

Pointer to IDT: No

Function pointers to all system calls: No

Memory management information about that process: Yes

Pointer to the parent process: Yes

EIP at the time of context switch: Yes

**Question 7**

Incorrect

Mark 0.00 out of 1.00

Select all the correct statements about code of bootmain() in xv6

```

void
bootmain(void)
{
    struct elfhdr *elf;
    struct proghdr *ph, *eph;
    void (*entry)(void);
    uchar* pa;

    elf = (struct elfhdr*)0x10000; // scratch space

    // Read 1st page off disk
    readseg((uchar*)elf, 4096, 0);

    // Is this an ELF executable?
    if(elf->magic != ELF_MAGIC)
        return; // let bootasm.S handle error

    // Load each program segment (ignores ph flags).
    ph = (struct proghdr*)((uchar*)elf + elf->phoff);
    eph = ph + elf->phnum;
    for(; ph < eph; ph++){
        pa = (uchar*)ph->paddr;
        readseg(pa, ph->filesz, ph->off);
        if(ph->memsz > ph->filesz)
            stosb(pa + ph->filesz, 0, ph->memsz - ph->filesz);
    }

    // Call the entry point from the ELF header.
    // Does not return!
    entry = (void(*)(void))(elf->entry);
    entry();
}

```

Also, inspect the relevant parts of the xv6 code. binary files, etc and run commands as you deem fit to answer this question.

- a. The kernel file gets loaded at the Physical address 0x10000 +0x80000000 in memory. ✗
- b. The elf->entry is set by the linker in the kernel file and it's 0x80000000 ✗
- c. The kernel ELF file contains actual physical address where particular sections of 'kernel' file should be loaded ✓
- d. The kernel file in memory is not necessarily a continuously filled in chunk, it may have holes in it. ✓
- e. The kernel file has only two program headers ✓
- f. The elf->entry is set by the linker in the kernel file and it's 0x80000000 ✗
- g. The readseg finally invokes the disk I/O code using assembly instructions ✓
- h. The elf->entry is set by the linker in the kernel file and it's 8010000c ✓
- i. The kernel file gets loaded at the Physical address 0x10000 in memory. ✓
- j. The condition if(ph->memsz > ph->filesz) is never true. ✗
- k. The stosb() is used here, to fill in some space in memory with zeroes ✓

Your answer is incorrect.

The correct answers are: The kernel file gets loaded at the Physical address 0x10000 in memory., The kernel file in memory is not necessarily a continuously filled in chunk, it may have holes in it., The elf->entry is set by the linker in the kernel file and it's 8010000c, The readseg finally invokes the disk I/O code using assembly instructions, The stosb() is used here, to fill in some space in memory with zeroes, The kernel ELF file contains actual physical address where particular sections of 'kernel' file should be loaded, The kernel file has only two program headers

**Question 8**

Partially correct

Mark 0.13 out of 0.25

Which of the following are NOT a part of job of a typical compiler?

- a. Check the program for logical errors ✓
- b. Convert high level language code to machine code
- c. Process the # directives in a C program
- d. Invoke the linker to link the function calls with their code, extern globals with their declaration
- e. Check the program for syntactical errors
- f. Suggest alternative pieces of code that can be written

Your answer is partially correct.

You have correctly selected 1.

The correct answers are: Check the program for logical errors, Suggest alternative pieces of code that can be written

**Question 9**

Correct

Mark 0.25 out of 0.25

Rank the following storage systems from slowest (first) to fastest(last)

Cache	6	✓
Hard Disk	3	✓
RAM	5	✓
Optical Disks	2	✓
Non volatile memory	4	✓
Registers	7	✓
Magnetic Tapes	1	✓

Your answer is correct.

The correct answer is: Cache → 6, Hard Disk → 3, RAM → 5, Optical Disks → 2, Non volatile memory → 4, Registers → 7, Magnetic Tapes → 1

**Question 10**

Partially correct

Mark 0.21 out of 0.50

Which of the following parts of a C program do not have any corresponding machine code ?

- a. local variable declaration
- b. global variables
- c. function calls ✗
- d. #directives ✓
- e. expressions
- f. pointer dereference
- g. typedefs ✓

Your answer is partially correct.

You have correctly selected 2.

The correct answers are: #directives, typedefs, global variables

**Question 11**

Correct

Mark 0.25 out of 0.25

Match a system call with it's description

pipe	create an unnamed FIFO storage with 2 ends - one for reading and another for writing	✓
dup	create a copy of the specified file descriptor into smallest available file descriptor	✓
dup2	create a copy of the specified file descriptor into another specified file descriptor	✓
exec	execute a binary file overlaying the image of current process	✓
fork	create an identical child process	✓

Your answer is correct.

The correct answer is: pipe → create an unnamed FIFO storage with 2 ends - one for reading and another for writing, dup → create a copy of the specified file descriptor into smallest available file descriptor, dup2 → create a copy of the specified file descriptor into another specified file descriptor, exec → execute a binary file overlaying the image of current process, fork → create an identical child process

**Question 12**

Correct

Mark 0.25 out of 0.25

Match the register with the segment used with it.

eip	cs	✓
edi	es	✓
esi	ds	✓
ebp	ss	✓
esp	ss	✓

Your answer is correct.

The correct answer is: eip → cs, edi → es, esi → ds, ebp → ss, esp → ss

**Question 13**

Correct

Mark 0.25 out of 0.25

What's the trapframe in xv6?

- a. A frame of memory that contains all the trap handler code
- b. The sequence of values, including saved registers, constructed on the stack when an interrupt occurs, built by hardware only
- c. The IDT table
- d. A frame of memory that contains all the trap handler code's function pointers
- e. A frame of memory that contains all the trap handler's addresses
- f. The sequence of values, including saved registers, constructed on the stack when an interrupt occurs, built by hardware + code in trapasm.S ✓
- g. The sequence of values, including saved registers, constructed on the stack when an interrupt occurs, built by code in trapasm.S only

Your answer is correct.

The correct answer is: The sequence of values, including saved registers, constructed on the stack when an interrupt occurs, built by hardware + code in trapasm.S

**Question 14**

Incorrect

Mark 0.00 out of 0.50

Select all the correct statements about linking and loading.

Select one or more:

- a. Continuous memory management schemes can support dynamic linking and dynamic loading. ✗
- b. Loader is last stage of the linker program ✗
- c. Continuous memory management schemes can support static linking and dynamic loading. (may be inefficiently) ✓
- d. Dynamic linking and loading is not possible without demand paging or demand segmentation. ✓
- e. Dynamic linking essentially results in relocatable code. ✓
- f. Continuous memory management schemes can support static linking and static loading. (may be inefficiently) ✓
- g. Loader is part of the operating system ✓
- h. Static linking leads to non-relocatable code ✗
- i. Dynamic linking is possible with continuous memory management, but variable sized partitions only. ✗

Your answer is incorrect.

The correct answers are: Continuous memory management schemes can support static linking and static loading. (may be inefficiently), Continuous memory management schemes can support static linking and dynamic loading. (may be inefficiently), Dynamic linking essentially results in relocatable code., Loader is part of the operating system, Dynamic linking and loading is not possible without demand paging or demand segmentation.

**Question 15**

Incorrect

Mark 0.00 out of 0.25

In bootasm.S, on the line

```
ljmp    $(SEG_KCODE<<3), $start32
```

The SEG\_KCODE << 3, that is shifting of 1 by 3 bits is done because

- a. The value 8 is stored in code segment
- b. The code segment is 16 bit and only upper 13 bits are used for segment number
- c. The code segment is 16 bit and only lower 13 bits are used for segment number ✗
- d. While indexing the GDT using CS, the value in CS is always divided by 8
- e. The ljmp instruction does a divide by 8 on the first argument

Your answer is incorrect.

The correct answer is: The code segment is 16 bit and only upper 13 bits are used for segment number

**Question 16**

Partially correct

Mark 0.07 out of 0.50

Order the events that occur on a timer interrupt:

Change to kernel stack

1	✗
---	---

Jump to a code pointed by IDT

2	✗
---	---

Jump to scheduler code

5	✗
---	---

Set the context of the new process

4	✗
---	---

Save the context of the currently running process

3	✓
---	---

Execute the code of the new process

6	✗
---	---

Select another process for execution

7	✗
---	---

Your answer is partially correct.

You have correctly selected 1.

The correct answer is: Change to kernel stack → 2, Jump to a code pointed by IDT → 1, Jump to scheduler code → 4, Set the context of the new process → 6, Save the context of the currently running process → 3, Execute the code of the new process → 7, Select another process for execution → 5

**Question 17**

Incorrect

Mark 0.00 out of 1.00

Consider the two programs given below to implement the command (ignore the fact that error checks are not done on return values of functions)

```
$ ls . /tmp/asdfksdf >/tmp/ddd 2>&1
```

**Program 1**

```
int main(int argc, char *argv[]) {
    int fd, n, i;
    char buf[128];

    fd = open("/tmp/ddd", O_WRONLY | O_CREAT, S_IRUSR | S_IWUSR);
    close(1);
    dup(fd);
    close(2);
    dup(fd);
    execl("/bin/ls", "/bin/ls", ".", "/tmp/asldjfaldfs", NULL);
}
```

**Program 2**

```
int main(int argc, char *argv[]) {
    int fd, n, i;
    char buf[128];

    close(1);
    fd = open("/tmp/ddd", O_WRONLY | O_CREAT, S_IRUSR | S_IWUSR);
    close(2);
    fd = open("/tmp/ddd", O_WRONLY | O_CREAT, S_IRUSR | S_IWUSR);
    execl("/bin/ls", "/bin/ls", ".", "/tmp/asldjfaldfs", NULL);
}
```

Select all the correct statements about the programs

Select one or more:

- a. Both programs are correct ✗
- b. Program 2 makes sure that there is one file offset used for '2' and '1' ✗
- c. Only Program 2 is correct ✗
- d. Program 2 does 1>&2 ✗
- e. Program 2 ensures 2>&1 and does not ensure >/tmp/ddd ✗
- f. Program 1 makes sure that there is one file offset used for '2' and '1' ✓
- g. Program 1 is correct for >/tmp/ddd but not for 2>&1 ✗
- h. Program 1 does 1>&2 ✗
- i. Both program 1 and 2 are incorrect ✗
- j. Program 2 is correct for >/tmp/ddd but not for 2>&1 ✗
- k. Only Program 1 is correct ✓
- l. Program 1 ensures 2>&1 and does not ensure >/tmp/ddd ✗

Your answer is incorrect.

The correct answers are: Only Program 1 is correct, Program 1 makes sure that there is one file offset used for '2' and '1'

**Question 18**

Correct

Mark 0.25 out of 0.25

Select the option which best describes what the CPU does during its powered ON lifetime

- a. Ask the user what is to be done, and execute that task
- b. Ask the OS what is to be done, and execute that task
- c. Fetch instructions specified by location given by PC, Decode and Execute it, during execution increment PC or change PC as per the instruction itself, Ask the User or the OS what is to be done next, repeat
- d. Fetch instructions specified by location given by PC, Decode and Execute it, during execution increment PC or change PC as per ✓ the instruction itself, repeat
- e. Fetch instruction specified by OS, Decode and execute it, repeat
- f. Fetch instructions specified by location given by PC, Decode and Execute it, during execution increment PC or change PC as per the instruction itself, Ask OS what is to be done next, repeat

The correct answer is: Fetch instructions specified by location given by PC, Decode and Execute it, during execution increment PC or change PC as per the instruction itself, repeat

**Question 19**

Partially correct

Mark 0.86 out of 1.00

Consider the following code and MAP the file to which each fd points at the end of the code.

```
int main(int argc, char *argv[]) {
    int fd1, fd2 = 1, fd3 = 1, fd4 = 1;

    fd1 = open("/tmp/1", O_WRONLY | O_CREAT, S_IRUSR|S_IWUSR);
    fd2 = open("/tmp/2", O_RDONLY);
    fd3 = open("/tmp/3", O_WRONLY | O_CREAT, S_IRUSR|S_IWUSR);
    close(0);
    close(1);
    dup(fd2);
    dup(fd3);
    close(fd3);
    dup2(fd2, fd4);
    printf("%d %d %d %d\n", fd1, fd2, fd3, fd4);
    return 0;
}
```

1	closed	✗
fd4	/tmp/2	✓
fd2	/tmp/2	✓
fd1	/tmp/1	✓
2	stderr	✓
0	/tmp/2	✓
fd3	closed	✓

Your answer is partially correct.

You have correctly selected 6.

The correct answer is: 1 → /tmp/3, fd4 → /tmp/2, fd2 → /tmp/2, fd1 → /tmp/1, 2 → stderr, 0 → /tmp/2, fd3 → closed

**Question 20**

Incorrect

Mark 0.00 out of 2.00

Following code claims to implement the command

```
/bin/ls -l | /usr/bin/head -3 | /usr/bin/tail -1
```

Fill in the blanks to make the code work.

Note: Do not include space in writing any option. x[1][2] should be written without any space, and so is the case with [1] or [2]. Pay attention to exact syntax and do not write any extra character like ';' or = etc.

```
int main(int argc, char *argv[]) {
```

```
    int pid1, pid2;
```

```
    int pfd[
```

```
    x ] [2];
```

```
    pipe(
```

```
    x );
```

```
    pid1 =
```

```
    x ;
```

```
    if(pid1 != 0) {
```

```
        close(pfd[0]
```

```
    x );
```

```
        close(
```

```
    x );
```

```
        dup(
```

```
    x );
```

```
        execl("/bin/ls", "/bin/ls", "
```

```
    x ", NULL);
```

```
    }
```

```
    pipe(
```

```
    x );
```

```
    x = fork();
```

```
    if(pid2 == 0) {
```

```
        close(
```

```
    x ;
```

```
        close(0);
```

```
        dup(
```

```
    x );
```

```
        close(pfd[1]
```

```
✗ );
close(
  
✗ );
dup(
  
✗ );
execl("/usr/bin/head", "/usr/bin/head", "  
  
✗ ", NULL);
} else {
close(pfd
  
✗ );
close(
  
✗ );
dup(
  
✗ );
close(pfd
  
✗ );
execl("/usr/bin/tail", "/usr/bin/tail", "  
  
✗ ", NULL);
}  
}
```

**Question 21**

Partially correct

Mark 0.11 out of 1.00

Select all the correct statements about calling convention on x86 32-bit.

- a. Return address is one location above the ebp ✓
- b. Parameters may be passed in registers or on stack ✓
- c. Space for local variables is allocated by subtracting the stack pointer inside the code of the called function ✓
- d. The ebp pointers saved on the stack constitute a chain of activation records ✓
- e. The two lines in the beginning of each function, "push %ebp; mov %esp, %ebp", create space for local variables ✗
- f. Parameters may be passed in registers or on stack ✓
- g. The return value is either stored on the stack or returned in the eax register ✗
- h. Parameters are pushed on the stack in left-right order
- i. during execution of a function, ebp is pointing to the old ebp
- j. Space for local variables is allocated by subtracting the stack pointer inside the code of the caller function ✗
- k. Compiler may allocate more memory on stack than needed ✓

Your answer is partially correct.

You have selected too many options.

The correct answers are: Compiler may allocate more memory on stack than needed, Parameters may be passed in registers or on stack, Return address is one location above the ebp, during execution of a function, ebp is pointing to the old ebp, Space for local variables is allocated by subtracting the stack pointer inside the code of the called function, The ebp pointers saved on the stack constitute a chain of activation records

**Question 22**

Correct

Mark 1.00 out of 1.00

Match the program with its output (ignore newlines in the output. Just focus on the count of the number of 'hi')

main() { int i = fork(); if(i == 0) execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); }

hi ✓

main() { fork(); execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); }

hi hi ✓

main() { int i = NULL; fork(); printf("hi\n"); }

hi hi ✓

main() { execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); }

hi ✓

Your answer is correct.

The correct answer is: main() { int i = fork(); if(i == 0) execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); } → hi, main() { fork(); execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); } → hi hi, main() { int i = NULL; fork(); printf("hi\n"); } → hi hi, main() { execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); } → hi

**Question 23**

Incorrect

Mark 0.00 out of 0.50

Some part of the bootloader of xv6 is written in assembly while some part is written in C. Why is that so?

Select all the appropriate choices

- a. The code in assembly is required for transition to protected mode, from real mode; but calling convention was applicable all the time ✗
- b. The setting up of the most essential memory management infrastructure needs assembly code ✓
- c. The code for reading ELF file can not be written in assembly ✗
- d. The code in assembly is required for transition to protected mode, from real mode; after that calling convention applies, hence code can be written in C ✓

Your answer is incorrect.

The correct answers are: The code in assembly is required for transition to protected mode, from real mode; after that calling convention applies, hence code can be written in C, The setting up of the most essential memory management infrastructure needs assembly code

**Question 24**

Incorrect

Mark 0.00 out of 0.50

xv6.img: bootblock kernel

```
dd if=/dev/zero of=xv6.img count=10000
dd if=bootblock of=xv6.img conv=notrunc
dd if=kernel of=xv6.img seek=1 conv=notrunc
```

Consider above lines from the Makefile. Which of the following is incorrect?

- a. The size of the kernel file is nearly 5 MB ✓
- b. The kernel is located at block-1 of the xv6.img ✗
- c. The xv6.img is of the size 10,000 blocks of 512 bytes each and occupies 10,000 blocks on the disk. ✗
- d. The size of xv6.img is exactly = (size of bootblock) + (size of kernel) ✗
- e. The bootblock is located on block-0 of the xv6.img ✗
- f. The xv6.img is of the size 10,000 blocks of 512 bytes each and occupies upto 10,000 blocks on the disk. ✓
- g. The bootblock may be 512 bytes or less (looking at the Makefile instruction) ✗
- h. The xv6.img is the virtual disk that is created by combining the bootblock and the kernel file. ✗
- i. The size of the xv6.img is nearly 5 MB ✗
- j. xv6.img is the virtual processor used by the qemu emulator ✓
- k. Blocks in xv6.img after kernel may be all zeroes. ✗

Your answer is incorrect.

The correct answers are: xv6.img is the virtual processor used by the qemu emulator, The xv6.img is of the size 10,000 blocks of 512 bytes each and occupies upto 10,000 blocks on the disk., The size of the kernel file is nearly 5 MB, The size of xv6.img is exactly = (size of bootblock) + (size of kernel)

**Question 25**

Incorrect

Mark 0.00 out of 1.00

Select the sequence of events that are NOT possible, assuming a non-interruptible kernel code

Select one or more:

a. P1 running

P1 makes system call

timer interrupt

Scheduler

P2 running

timer interrupt

Scheuler

P1 running

P1's system call return

b. P1 running

P1 makes sytem call and blocks

Scheduler

P2 running

P2 makes sytem call and blocks

Scheduler

P1 running again



c. P1 running

P1 makes system call

system call returns

P1 running

timer interrupt

Scheduler running

P2 running

d. P1 running

P1 makes sytem call and blocks

Scheduler

P2 running

P2 makes sytem call and blocks

Scheduler

P3 running

Hardware interrupt

Interrupt unblocks P1

Interrupt returns

P3 running

Timer interrupt

Scheduler

P1 running



e.

P1 running

P1 makes sytem call

Scheduler

P2 running

P2 makes sytem call and blocks

Scheduler

P1 running again

f. P1 running

keyboard hardware interrupt

keyboard interrupt handler running

interrupt handler returns

P1 running

P1 makes sytem call

system call returns



P1 running  
timer interrupt  
scheduler  
P2 running

Your answer is incorrect.

The correct answers are: P1 running

P1 makes system call and blocks

Scheduler

P2 running

P2 makes system call and blocks

Scheduler

P1 running again, P1 running

P1 makes system call

timer interrupt

Scheduler

P2 running

timer interrupt

Scheuler

P1 running

P1's system call return,

P1 running

P1 makes system call

Scheduler

P2 running

P2 makes system call and blocks

Scheduler

P1 running again

**Question 26**

Correct

Mark 0.25 out of 0.25

Which of the following are the files related to bootloader in xv6?

- a. bootasm.s and entry.S
- b. bootasm.S and bootmain.c ✓
- c. bootasm.S, bootmain.c and bootblock.c
- d. bootmain.c and bootblock.S

Your answer is correct.

The correct answer is: bootasm.S and bootmain.c

**Question 27**

Correct

Mark 0.25 out of 0.25

Match the following parts of a C program to the layout of the process in memory

Instructions	Text section	✓
Local Variables	Stack Section	✓
Dynamically allocated memory	Heap Section	✓
Global and static data	Data section	✓

Your answer is correct.

The correct answer is:

Instructions → Text section, Local Variables → Stack Section,  
Dynamically allocated memory → Heap Section,  
Global and static data → Data section

**Question 28**

Incorrect

Mark 0.00 out of 0.50

What will this program do?

```
int main() {  
    fork();  
    execl("/bin/ls", "/bin/ls", NULL);  
    printf("hello");  
}
```

- a. one process will run ls, another will print hello
- b. run ls once ✗
- c. run ls twice
- d. run ls twice and print hello twice
- e. run ls twice and print hello twice, but output will appear in some random order

Your answer is incorrect.

The correct answer is: run ls twice

**Question 29**

Correct

Mark 0.25 out of 0.25

What is the OS Kernel?

- a. The code that controls hardware, abstracts access to hardware resources using system calls, creates an environment for processes to be created and run ✓ correct
- b. The set of tools like compiler, linker, loader, terminal, shell, etc.
- c. Only the system programs like compiler, linker, loader, etc.
- d. Everything that I see on my screen

The correct answer is: The code that controls hardware, abstracts access to hardware resources using system calls, creates an environment for processes to be created and run

**Question 30**

Correct

Mark 0.50 out of 0.50

Which of the following is/are not saved during context switch?

- a. Program Counter
- b. General Purpose Registers
- c. Bus ✓
- d. Stack Pointer
- e. MMU related registers/information
- f. Cache ✓
- g. TLB ✓

Your answer is correct.

The correct answers are: TLB, Cache, Bus

**Question 31**

Partially correct

Mark 0.10 out of 0.25

Select the order in which the various stages of a compiler execute.

Linking	3	
Syntactical Analysis	2	
Pre-processing	1	
Intermediate code generation	does not exist	
Loading	4	

Your answer is partially correct.

You have correctly selected 2.

The correct answer is: Linking → 4, Syntactical Analysis → 2, Pre-processing → 1, Intermediate code generation → 3, Loading → does not exist

**Question 32**

Partially correct

Mark 0.08 out of 0.50

Order the sequence of events, in scheduling process P1 after process P0

context of P0 is saved in P0's PCB	2	
context of P1 is loaded from P1's PCB	3	
Process P1 is running	5	
timer interrupt occurs	6	
Process P0 is running	1	
Control is passed to P1	4	

Your answer is partially correct.

You have correctly selected 1.

The correct answer is: context of P0 is saved in P0's PCB → 3, context of P1 is loaded from P1's PCB → 4, Process P1 is running → 6, timer interrupt occurs → 2, Process P0 is running → 1, Control is passed to P1 → 5

**Question 33**

Not answered

Marked out of 1.00

Select the correct statements about interrupt handling in xv6 code

- a. On any interrupt/syscall/exception the control first jumps in vectors.S
- b. The trapframe pointer in struct proc, points to a location on user stack
- c. Each entry in IDT essentially gives the values of CS and EIP to be used in handling that interrupt
- d. xv6 uses the 64th entry in IDT for system calls
- e. The CS and EIP are changed only after pushing user code's SS,ESP on stack
- f. The trapframe pointer in struct proc, points to a location on kernel stack
- g. The function trap() is called only in case of hardware interrupt
- h. The CS and EIP are changed only immediately on a hardware interrupt
- i. All the 256 entries in the IDT are filled
  
- j. On any interrupt/syscall/exception the control first jumps in trapasm.S
- k. The function trap() is called irrespective of hardware interrupt/system-call/exception
- l. xv6 uses the 0x64th entry in IDT for system calls
- m. Before going to alltraps, the kernel stack contains upto 5 entries.

Your answer is incorrect.

The correct answers are: All the 256 entries in the IDT are filled, Each entry in IDT essentially gives the values of CS and EIP to be used in handling that interrupt, xv6 uses the 64th entry in IDT for system calls, On any interrupt/syscall/exception the control first jumps in trapasm.S, Before going to alltraps, the kernel stack contains upto 5 entries., The trapframe pointer in struct proc, points to a location on kernel stack, The function trap() is called irrespective of hardware interrupt/system-call/exception, The CS and EIP are changed only after pushing user code's SS,ESP on stack

[◀ \(Assignment\) Change free list management in xv6](#)

Jump to...

**Started on** Tuesday, 22 March 2022, 1:59:34 PM

**State** Finished

**Completed on** Tuesday, 22 March 2022, 4:14:53 PM

**Time taken** 2 hours 15 mins

**Grade** 24.57 out of 40.00 (61%)

**Question 1**

Partially correct

Mark 0.10 out of 1.00

Select all the correct statements w.r.t user and kernel threads

Select one or more:

a. many-one model can be implemented even if there are no kernel threads ✓

b. many-one model gives no speedup on multicore processors

c. all three models, that is many-one, one-one, many-many , require a user level thread library ✓

d. A process blocks in many-one model even if a single thread makes a blocking system call

e. A process may not block in many-one model, if a thread makes a blocking system call ✗

f. one-one model increases kernel's scheduling load ✓

g. one-one model can be implemented even if there are no kernel threads

Your answer is partially correct.

You have correctly selected 3.

The correct answers are: many-one model can be implemented even if there are no kernel threads, all three models, that is many-one, one-one, many-many , require a user level thread library, one-one model increases kernel's scheduling load, many-one model gives no speedup on multicore processors, A process blocks in many-one model even if a single thread makes a blocking system call

**Question 2**

Partially correct

Mark 0.50 out of 1.00

Select all correct statements w.r.t. Major and Minor page faults on Linux

- a. Thrashing is possible only due to major page faults
- b. Minor page fault may occur because the page was freed, but still tagged and available in the free page list
- c. Minor page fault may occur because of a page fault during fork(), on code of an already running process ✓
- d. Major page faults are likely to occur in more numbers at the beginning of the process ✓
- e. Minor page fault may occur because the page was a shared memory page ✓
- f. Minor page faults are an improvement of the page buffering techniques

The correct answers are: Minor page fault may occur because the page was a shared memory page, Minor page fault may occur because of a page fault during fork(), on code of an already running process, Minor page fault may occur because the page was freed, but still tagged and available in the free page list, Major page faults are likely to occur in more numbers at the beginning of the process, Thrashing is possible only due to major page faults, Minor page faults are an improvement of the page buffering techniques

**Question 3**

Correct

Mark 2.00 out of 2.00

W.r.t. Memory management in xv6,

xv6 uses physical memory upto 224 MB only  
Mark statements True or False

True	False	
<input checked="" type="radio"/>	<input type="radio"/> ✗	The stack allocated in entry.S is used as stack for scheduler's context for first processor
<input type="radio"/> ✗	<input checked="" type="radio"/>	The switchkvm() call in scheduler() is invoked after control comes to it from swtch() scheduler(), thus demanding execution in new process's context
<input checked="" type="radio"/>	<input type="radio"/> ✗	The switchkvm() call in scheduler() changes CR3 to use page directory kpgdir
<input checked="" type="radio"/>	<input type="radio"/> ✗	The kernel code and data take up less than 2 MB space
<input checked="" type="radio"/>	<input type="radio"/> ✗	The switchkvm() call in scheduler() is invoked after control comes to it from sched(), thus demanding execution in kernel's context
<input checked="" type="radio"/>	<input type="radio"/> ✗	xv6 uses physical memory upto 224 MB only
<input checked="" type="radio"/>	<input type="radio"/> ✗	The free page-frame are created out of nearly 222 MB
<input type="radio"/> ✗	<input checked="" type="radio"/>	The switchkvm() call in scheduler() changes CR3 to use page directory of new process
<input checked="" type="radio"/>	<input type="radio"/> ✗	PHYSTOP can be increased to some extent, simply by editing memlayout.h
<input checked="" type="radio"/>	<input type="radio"/> ✗	The process's address space gets mapped on frames, obtained from ~2MB:224MB range
<input type="radio"/> ✗	<input checked="" type="radio"/>	The kernel's page table given by kpgdir variable is used as stack for scheduler's context

The stack allocated in entry.S is used as stack for scheduler's context for first processor: True

The switchkvm() call in scheduler() is invoked after control comes to it from swtch() scheduler(), thus demanding execution in new process's context: False

The switchkvm() call in scheduler() changes CR3 to use page directory kpgdir: True

The kernel code and data take up less than 2 MB space: True

The switchkvm() call in scheduler() is invoked after control comes to it from sched(), thus demanding execution in kernel's context: True

xv6 uses physical memory upto 224 MB only: True

The free page-frame are created out of nearly 222 MB: True

The switchkvm() call in scheduler() changes CR3 to use page directory of new process: False

PHYSTOP can be increased to some extent, simply by editing memlayout.h: True

The process's address space gets mapped on frames, obtained from ~2MB:224MB range: True

The kernel's page table given by kpgdir variable is used as stack for scheduler's context: False

**Question 4**

Correct

Mark 1.00 out of 1.00

Given that a kernel has 1000 KB of total memory, and holes of sizes (in that order) 300 KB, 200 KB, 100 KB, 250 KB. For each of the requests on the left side, match it with the chunk chosen using the specified algorithm.

Consider each request as first request.

50 KB, worst fit	300 KB	✓
100 KB, worst fit	300 KB	✓
200 KB, first fit	300 KB	✓
150 KB, best fit	200 KB	✓
220 KB, best fit	250 KB	✓
150 KB, first fit	300 KB	✓

The correct answer is: 50 KB, worst fit → 300 KB, 100 KB, worst fit → 300 KB, 200 KB, first fit → 300 KB, 150 KB, best fit → 200 KB, 220 KB, best fit → 250 KB, 150 KB, first fit → 300 KB

**Question 5**

Partially correct

Mark 0.60 out of 1.00

Choice of the global or local replacement strategy is a subjective choice for kernel programmers. There are advantages and disadvantages on either side. Out of the following statements, that advocate either global or local replacement strategy, select those statements that have a logically **CONSISTENT** argument. (That is any statement that is logically correct about either global or local replacement)

**Consistent**    **Inconsistent**

<input checked="" type="radio"/>	<input type="radio"/> X	Local replacement can lead to under-utilisation of memory, because a process may not use all the pages allocated to it all the time.	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> X	Global replacement may give highly variable per process completion time because number of page faults become un-predictable.	X
<input checked="" type="radio"/>	<input type="radio"/> X	Global replacement can be preferred when greater throughput (number of processes completing per unit time) is a concern, because each process tries to complete at the expense of others, thus leading to overall more processes completing (unless thrashing occurs).	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> X	Local replacement results in more predictable per-process completion time because number of page faults can be better predicted.	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> X	Local replacement can be preferred when avoiding thrashing is a major concern because with local replacement and minimum number of frames allocated, a process is always able to progress and cascading inter-process page faults are avoided.	X

Local replacement can lead to under-utilisation of memory, because a process may not use all the pages allocated to it all the time.: Consistent

Global replacement may give highly variable per process completion time because number of page faults become un-predictable.: Consistent

Global replacement can be preferred when greater throughput (number of processes completing per unit time) is a concern, because each process tries to complete at the expense of others, thus leading to overall more processes completing (unless thrashing occurs).: Consistent

Local replacement results in more predictable per-process completion time because number of page faults can be better predicted.: Consistent

Local replacement can be preferred when avoiding thrashing is a major concern because with local replacement and minimum number of frames allocated, a process is always able to progress and cascading inter-process page faults are avoided.: Consistent

**Question 6**

Correct

Mark 1.00 out of 1.00

Map the functionality/use with function/variable in xv6 code.

Setup kernel part of a page table, and switch to that page table

kvmalloc()

Setup kernel part of a page table, mapping kernel code, data, read-only data, I/O space, devices

setupkvm()

return a free page, if available; 0, otherwise

kalloc()

Create page table entries for a given range of virtual and physical addresses; including page directory entries if needed

mappages()

Return address of page table entry in a given page directory, for a given virtual address; creates page table if necessary

walkpgdir()

Array listing the kernel memory mappings, to be used by setupkvm()

kmap[]

Your answer is correct.

The correct answer is: Setup kernel part of a page table, and switch to that page table → kvmalloc(), Setup kernel part of a page table, mapping kernel code, data, read-only data, I/O space, devices → setupkvm(), return a free page, if available; 0, otherwise → kalloc(), Create page table entries for a given range of virtual and physical addresses; including page directory entries if needed → mappages(), Return address of page table entry in a given page directory, for a given virtual address; creates page table if necessary → walkpgdir(), Array listing the kernel memory mappings, to be used by setupkvm() → kmap[]

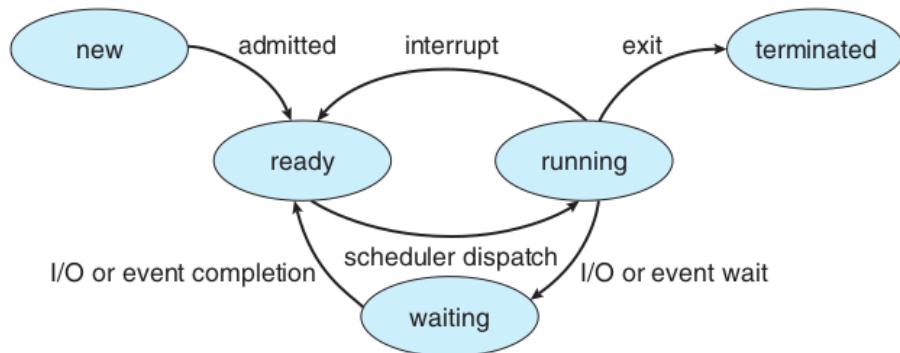
## Question 7

Partially correct

Mark 0.20 out of 1.00

Mark statements True/False w.r.t. change of states of a process. Note that a statement is true only if the claim and argument both are true.

Reference: The process state diagram (and your understanding of how kernel code works). Note - the diagram does not show zombie state!



**Figure 3.2** Diagram of process state.

True	False	
<input checked="" type="radio"/>	<input type="radio"/>	Every forked process has to go through ZOMBIE state, at least for a small duration.
<input type="radio"/>	<input checked="" type="radio"/>	A process only in RUNNING state can become TERMINATED because scheduler moves it to ZOMBIE state first
<input checked="" type="radio"/>	<input type="radio"/>	A process in WAITING state can not become RUNNING because the event it's waiting for has not occurred and it has not been moved to ready queue yet
<input checked="" type="radio"/>	<input type="radio"/>	Only a process in READY state is considered by scheduler
<input type="radio"/>	<input checked="" type="radio"/>	A process in READY state can not go to WAITING state because the resource on which it will WAIT will not be in use when process is in READY state.

Every forked process has to go through ZOMBIE state, at least for a small duration.: True

A process only in RUNNING state can become TERMINATED because scheduler moves it to ZOMBIE state first: False

A process in WAITING state can not become RUNNING because the event it's waiting for has not occurred and it has not been moved to ready queue yet: True

Only a process in READY state is considered by scheduler: True

A process in READY state can not go to WAITING state because the resource on which it will WAIT will not be in use when process is in READY state.: False

**Question 8**

Correct

Mark 2.00 out of 2.00

Consider the reference string

6 4 2 0 1 2 6 9 2 0 5

If the number of page frames is 3, then total number of page faults (including initial), using FIFO replacement is:

Answer: 

#6# 6,4# 6,4,2 #0,4,2# 0,1,2 #0,1,6 #9,1,6# 9,2,6# 9,2,0 #5,2,0

The correct answer is: 10

**Question 9**

Incorrect

Mark 0.00 out of 1.00

Select all the correct statements about linking and loading.

Select one or more:

- a. Continuous memory management schemes can support static linking and dynamic loading. (may be inefficiently)
- b. Dynamic linking essentially results in relocatable code.
- c. Continuous memory management schemes can support static linking and static loading. (may be inefficiently)
- d. Loader is last stage of the linker program
- e. Dynamic linking and loading is not possible without demand paging or demand segmentation.
- f. Static linking leads to non-relocatable code
- g. Continuous memory management schemes can support dynamic linking and dynamic loading.
- h. Dynamic linking is possible with continuous memory management, but variable sized partitions only.
- i. Loader is part of the operating system

Your answer is incorrect.

The correct answers are: Continuous memory management schemes can support static linking and static loading. (may be inefficiently), Continuous memory management schemes can support static linking and dynamic loading. (may be inefficiently), Dynamic linking essentially results in relocatable code., Loader is part of the operating system, Dynamic linking and loading is not possible without demand paging or demand segmentation.

**Question 10**

Correct

Mark 1.00 out of 1.00

Consider a computer system with a 32-bit logical address and 4- KB page size. The system supports up to 512 MB of physical memory. How many entries are there in each of the following?

Write answer as a decimal number.

A conventional, single-level page table:

1048576



An inverted page table:

131072

**Question 11**

Incorrect

Mark 0.00 out of 1.00

W.r.t. xv6 code, match the state of a process with a code that sets the state

RUNNING	Choose...
ZOMBIE	Choose...
EMBRYO	Choose...
SLEEPING	Choose...
UNUSED	Choose...
RUNNABLE	scheduler() <span style="color: red;">✖</span>

The correct answer is: RUNNING → scheduler(), ZOMBIE → exit(), called by process itself, EMBRYO → fork()->allocproc() before setting up the UVM, SLEEPING → sleep(), called by any process blocking itself, UNUSED → wait(), called by parent process, RUNNABLE → wakeup(), called by an interrupt handler

**Question 12**

Not answered

Marked out of 1.00

Select the correct statements about interrupt handling in xv6 code

- a. The trapframe pointer in struct proc, points to a location on user stack
- b. The CS and EIP are changed only immediately on a hardware interrupt
- c. The CS and EIP are changed only after pushing user code's SS,ESP on stack
- d. The trapframe pointer in struct proc, points to a location on kernel stack
- e. On any interrupt/syscall/exception the control first jumps in vectors.S
- f. The function trap() is called irrespective of hardware interrupt/system-call/exception
- g. All the 256 entries in the IDT are filled
- h. The function trap() is called only in case of hardware interrupt
- i. xv6 uses the 64th entry in IDT for system calls
- j. xv6 uses the 0x64th entry in IDT for system calls
- k. Before going to alptraps, the kernel stack contains upto 5 entries.
- l. Each entry in IDT essentially gives the values of CS and EIP to be used in handling that interrupt
- m. On any interrupt/syscall/exception the control first jumps in trapasm.S

Your answer is incorrect.

The correct answers are: All the 256 entries in the IDT are filled, Each entry in IDT essentially gives the values of CS and EIP to be used in handling that interrupt, xv6 uses the 64th entry in IDT for system calls, On any interrupt/syscall/exception the control first jumps in vectors.S, Before going to alptraps, the kernel stack contains upto 5 entries., The trapframe pointer in struct proc, points to a location on kernel stack, The function trap() is called irrespective of hardware interrupt/system-call/exception, The CS and EIP are changed only after pushing user code's SS,ESP on stack

**Question 13**

Correct

Mark 1.00 out of 1.00

The complete range of virtual addresses (after main() in main.c is over), from which the free pages used by kalloc() and kfree() is derived, are:

- a. end, (4MB + PHYSTOP)
- b. P2V(end), PHYSTOP
- c. end, P2V(4MB + PHYSTOP)
- d. end, PHYSTOP
- e. end, 4MB
- f. end, P2V(PHYSTOP) ✓
- g. P2V(end), P2V(PHYSTOP)

Your answer is correct.

The correct answer is: end, P2V(PHYSTOP)

**Question 14**

Correct

Mark 1.00 out of 1.00

Select all the correct statements about MMU and its functionality (on a non-demand paged system)

Select one or more:

- a. Illegal memory access is detected in hardware by MMU and a trap is raised ✓
- b. Illegal memory access is detected by operating system
- c. MMU is a separate chip outside the processor
- d. The operating system interacts with MMU for every single address translation
- e. Logical to physical address translations in MMU are done in hardware, automatically ✓
- f. Logical to physical address translations in MMU are done with specific machine instructions
- g. MMU is inside the processor ✓
- h. The Operating system sets up relevant CPU registers to enable proper MMU translations ✓

Your answer is correct.

The correct answers are: MMU is inside the processor, Logical to physical address translations in MMU are done in hardware, automatically, The Operating system sets up relevant CPU registers to enable proper MMU translations, Illegal memory access is detected in hardware by MMU and a trap is raised

**Question 15**

Incorrect

Mark 0.00 out of 2.00

Order the following events, in the creation of init() process in xv6:

1. ✗ initcode is selected by scheduler for execution
2. ✗ kernel stack is allocated for initcode process
3. ✗ values are set in the trapframe of initcode
4. ✗ sys\_exec runs
5. ✗ initcode process is set to be runnable
6. ✗ code is set to start in forkret() when process gets scheduled
7. ✗ Arguments are setup on process stack for /init
8. ✗ trapframe and context pointers are set to proper location
9. ✗ trap() runs
10. ✗ userinit() is called
11. ✗ the header of "/init" ELF file is ready by kernel
12. ✗ empty struct proc is obtained for initcode
13. ✗ Stack is allocated for "/init" process
14. ✗ function pointer from syscalls[] array is invoked
15. ✗ memory mappings are created for "/init" process
16. ✗ page table mappings of 'initcode' are replaced by mappings of 'init'
17. ✗ initcode calls exec system call
18. ✗ initcode process runs
19. ✗ name of process "/init" is copied in struct proc

Your answer is incorrect.

Grading type: Relative to the next item (including last)

Grade details: 0 / 19 = 0%

Here are the scores for each item in this response:

1. 0 / 1 = 0%
2. 0 / 1 = 0%
3. 0 / 1 = 0%
4. 0 / 1 = 0%
5. 0 / 1 = 0%
6. 0 / 1 = 0%
7. 0 / 1 = 0%
8. 0 / 1 = 0%
9. 0 / 1 = 0%
10. 0 / 1 = 0%
11. 0 / 1 = 0%
12. 0 / 1 = 0%
13. 0 / 1 = 0%
14. 0 / 1 = 0%
15. 0 / 1 = 0%

- 16. 0 / 1 = 0%
- 17. 0 / 1 = 0%
- 18. 0 / 1 = 0%
- 19. 0 / 1 = 0%

The correct order for these items is as follows:

1. userinit() is called
2. empty struct proc is obtained for initcode
3. kernel stack is allocated for initcode process
4. trapframe and context pointers are set to proper location
5. code is set to start in forkret() when process gets scheduled
6. kernel memory mappings are created for initcode
7. values are set in the trapframe of initcode
8. initcode process is set to be runnable
9. initcode is selected by scheduler for execution
10. initcode process runs
11. initcode calls exec system call
12. trap() runs
13. function pointer from syscalls[] array is invoked
14. sys\_exec runs
15. the header of "/init" ELF file is ready by kernel
16. memory mappings are created for "/init" process
17. Stack is allocated for "/init" process
18. Arguments on setup on process stack for /init
19. name of process "/init" is copied in struct proc
20. page table mappings of 'initcode' are replaced by makpings of 'init'

**Question 16**

Partially correct

Mark 0.56 out of 1.00

Mark the statements as True or False, w.r.t. mmap()

True	False	
<input checked="" type="radio"/>	<input checked="" type="radio"/>	mmap() can be implemented on both demand paged and non-demand paged systems.
<input checked="" type="radio"/>	<input checked="" type="radio"/>	MAP_FIXED guarantees that the mapping is always done at the specified address
<input checked="" type="radio"/>	<input checked="" type="radio"/>	MAP_PRIVATE leads to a mapping that is copy-on-write
<input checked="" type="radio"/>	<input checked="" type="radio"/>	on failure mmap() returns (void *)-1
<input checked="" type="radio"/>	<input checked="" type="radio"/>	MAP_SHARED leads to a mapping that is copy-on-write
<input checked="" type="radio"/>	<input checked="" type="radio"/>	mmap() results in changes to buffer-cache of the kernel.
<input checked="" type="radio"/>	<input checked="" type="radio"/>	on failure mmap() returns NULL
<input checked="" type="radio"/>	<input checked="" type="radio"/>	mmap() results in changes to page table of a process.
<input checked="" type="radio"/>	<input checked="" type="radio"/>	mmap() is a system call

mmap() can be implemented on both demand paged and non-demand paged systems.: True

MAP\_FIXED guarantees that the mapping is always done at the specified address: False

MAP\_PRIVATE leads to a mapping that is copy-on-write: True

on failure mmap() returns (void \*)-1: True

MAP\_SHARED leads to a mapping that is copy-on-write: False

mmap() results in changes to buffer-cache of the kernel.: False

on failure mmap() returns NULL: False

mmap() results in changes to page table of a process.: True

mmap() is a system call: True

**Question 17**

Incorrect

Mark 0.00 out of 1.00

If one thread opens a file with read privileges then

Select one:

- a. other threads in the same process can also read from that file
- b. none of these
- c. any other thread cannot read from that file
- d. other threads in the another process can also read from that file

Your answer is incorrect.

The correct answer is: other threads in the same process can also read from that file

**Question 18**

Partially correct

Mark 0.60 out of 1.00

Mark the statements about named and un-named pipes as True or False

True	False	
<input checked="" type="radio"/>	<input type="radio"/> 	Named pipe exists as a file
<input checked="" type="radio"/>	<input type="radio"/> 	Un-named pipes are inherited by a child process from parent.
<input type="radio"/> 	<input checked="" type="radio"/>	The buffers for named-pipe are in process-memory while the buffers for the un-named pipe are in kernel memory.
<input checked="" type="radio"/>	<input type="radio"/> 	Both types of pipes are an extension of the idea of "message passing".
<input type="radio"/> 	<input checked="" type="radio"/>	A named pipe has a name decided by the kernel.
<input checked="" type="radio"/>	<input type="radio"/> 	Un-named pipes can be used for communication between only "related" processes, if the common ancestor created it.
<input checked="" type="radio"/>	<input type="radio"/> 	Both types of pipes provide FIFO communication.
<input type="radio"/> 	<input checked="" type="radio"/>	Named pipes can be used for communication between only "related" processes.
<input checked="" type="radio"/>	<input type="radio"/> 	Named pipes can exist beyond the life-time of processes using them.
<input type="radio"/> 	<input checked="" type="radio"/>	The pipe() system call can be used to create either a named or un-named pipe.

Named pipe exists as a file.: True

Un-named pipes are inherited by a child process from parent.: True

The buffers for named-pipe are in process-memory while the buffers for the un-named pipe are in kernel memory.: False

Both types of pipes are an extension of the idea of "message passing": True

A named pipe has a name decided by the kernel.: False

Un-named pipes can be used for communication between only "related" processes, if the common ancestor created it.: True

Both types of pipes provide FIFO communication.: True

Named pipes can be used for communication between only "related" processes.: False

Named pipes can exist beyond the life-time of processes using them.: True

The pipe() system call can be used to create either a named or un-named pipe.: False

**Question 19**

Partially correct

Mark 0.67 out of 1.00

Select the most common causes of use of IPC by processes

- a. More modular code
- b. Breaking up a large task into small tasks and speeding up computation, on multiple core machines ✓
- c. More security checks
- d. Sharing of information of common interest ✓
- e. Get the kernel performance statistics

The correct answers are: Sharing of information of common interest, Breaking up a large task into small tasks and speeding up computation, on multiple core machines, More modular code

**Question 20**

Correct

Mark 1.00 out of 1.00

For each function/code-point, select the status of segmentation setup in xv6

after seginit() in main()	gdt setup with 5 entries (0 to 4) on one processor	✓
bootmain()	gdt setup with 3 entries, at start32 symbol of bootasm.S	✓
after startothers() in main()	gdt setup with 5 entries (0 to 4) on all processors	✓
entry.S	gdt setup with 3 entries, at start32 symbol of bootasm.S	✓
kvmalloc() in main()	gdt setup with 3 entries, at start32 symbol of bootasm.S	✓
bootasm.S	gdt setup with 3 entries, at start32 symbol of bootasm.S	✓

Your answer is correct.

The correct answer is: after seginit() in main() → gdt setup with 5 entries (0 to 4) on one processor, bootmain() → gdt setup with 3 entries, at start32 symbol of bootasm.S, after startothers() in main() → gdt setup with 5 entries (0 to 4) on all processors, entry.S → gdt setup with 3 entries, at start32 symbol of bootasm.S, kvmalloc() in main() → gdt setup with 3 entries, at start32 symbol of bootasm.S, bootasm.S → gdt setup with 3 entries, at start32 symbol of bootasm.S

**Question 21**

Partially correct

Mark 0.50 out of 1.00

Mark whether the given sequence of events is possible or not-possible. Also, select the reason for your answer.

For each sequence it's a not-possible sequence if some important event is not mentioned in the sequence.

Assume that the kernel code is non-interruptible and uniprocessor system.

Process P1, user code executing

Timer interrupt

Context changes to kernel context

Generic interrupt handler runs

Generic interrupt handler calls Scheduler

Scheduler selects P2 for execution

After scheduler, Process P2 user code executing

This sequence of events is:  not-possible ✓

Because

Generic interrupt handler can not call scheduler ✗

**Question 22**

Partially correct

Mark 0.63 out of 1.00

Mark the statements as True or False, w.r.t. passing of arguments to system calls in xv6 code.

True	False	
<input checked="" type="radio"/> ✗	<input type="radio"/> ✗	Integer arguments are stored in eax, ebx, ecx, etc. registers
<input checked="" type="radio"/> ✗	<input checked="" type="radio"/> ✗	String arguments are NOT copied in kernel memory, but just pointed to by a kernel memory pointer
<input checked="" type="radio"/> ✗	<input checked="" type="radio"/> ✗	The functions like argint(), argstr() make the system call arguments available in the kernel.
<input checked="" type="radio"/> ✗	<input type="radio"/> ✗	String arguments are first copied to trapframe and then from trapframe to kernel's other variables.
<input checked="" type="radio"/> ✗	<input checked="" type="radio"/> ✗	The arguments to system call originally reside on process stack.
<input checked="" type="radio"/> ✗	<input type="radio"/> ✗	The arguments to system call are copied to kernel stack in trapasm.S
<input checked="" type="radio"/> ✗	<input checked="" type="radio"/> ✗	Integer arguments are copied from user memory to kernel memory using argint()
<input checked="" type="radio"/> ✗	<input checked="" type="radio"/> ✗	The arguments are accessed in the kernel code using esp on the trapframe.

Integer arguments are stored in eax, ebx, ecx, etc. registers: False

String arguments are NOT copied in kernel memory, but just pointed to by a kernel memory pointer: True

The functions like argint(), argstr() make the system call arguments available in the kernel.: True

String arguments are first copied to trapframe and then from trapframe to kernel's other variables.: False

The arguments to system call originally reside on process stack.: True

The arguments to system call are copied to kernel stack in trapasm.S: False

Integer arguments are copied from user memory to kernel memory using argint(): True

The arguments are accessed in the kernel code using esp on the trapframe.: True

**Question 23**

Not answered

Marked out of 1.00

Given below is a sequence of reference bits on pages before the second chance algorithm runs. Before the algorithm runs, the counter is at the page marked (x). Write the sequence of reference bits after the second chance algorithm has executed once. In the answer write PRECISELY one space BETWEEN each number and do not mention (x).

0 0 1(x) 1 0 1 1

Answer:



The correct answer is: 0 0 0 0 0 1 1

**Question 24**

Correct

Mark 2.00 out of 2.00

For the reference string

3 4 3 5 2

using FIFO replacement policy for pages,

consider the number of page faults for 2, 3 and 4 page frames.

Select the correct statement.

Select one:

- a. Exhibit Balady's anomaly between 3 and 4 frames
- b. Do not exhibit Balady's anomaly
- c. Exhibit Balady's anomaly between 2 and 3 frames



Your answer is correct.

The correct answer is: Do not exhibit Balady's anomaly

**Question 25**

Correct

Mark 1.00 out of 1.00

For the reference string

3 4 3 5 2

using LRU replacement policy for pages,

consider the number of page faults for 2, 3 and 4 page frames.

Select the most correct statement.

Select one:

- a. LRU will never exhibit Balady's anomaly ✓
- b. Exhibit Balady's anomaly between 2 and 3 frames
- c. This example does not exhibit Balady's anomaly
- d. Exhibit Balady's anomaly between 3 and 4 frames

Your answer is correct.

The correct answer is: LRU will never exhibit Balady's anomaly

**Question 26**

Partially correct

Mark 0.55 out of 1.00

Select all the correct statements about process states.

Note that in this question you lose marks for every incorrect choice that you make, proportional to actual number of incorrect choices.

- a. Process state is implemented as a string
- b. Process state is stored in the PCB ✓
- c. A process becomes ZOMBIE when another process bites into its memory
- d. Process state is stored in the processor ✗
- e. The scheduler can change state of a process from RUNNABLE to RUNNING and vice-versa
- f. The scheduler can change state of a process from RUNNABLE to RUNNING ✓
- g. A process becomes ZOMBIE when it calls exit() ✓
- h. Process state is changed only by interrupt handlers
- i. Process state can be implemented as just a number

Your answer is partially correct.

You have correctly selected 3.

The correct answers are: Process state is stored in the PCB, Process state can be implemented as just a number, The scheduler can change state of a process from RUNNABLE to RUNNING, A process becomes ZOMBIE when it calls exit()

**Question 27**

Partially correct

Mark 0.38 out of 1.00

Consider a demand-paging system with the following time-measured utilizations:

CPU utilization : 20%

Paging disk: 97.7%

Other I/O devices: 5%

For each of the following, indicate whether it will (or is likely to) improve CPU utilization (even if by a small amount). Explain your answers.

a. Install a faster CPU : Yes ✗

b. Install a bigger paging disk. : Yes ✗

c. Increase the degree of multiprogramming. : Yes ✗

d. Decrease the degree of multiprogramming. : Yes ✓

e. Install more main memory.: Yes ✓

f. Install a faster hard disk or multiple controllers with multiple hard disks. : Yes ✓

g. Add prepaging to the page-fetch algorithms. :

May be ✗

h. Increase the page size. : May be ✗

**Question 28**

Incorrect

Mark 0.00 out of 1.00

Suppose a kernel uses a buddy allocator. The smallest chunk that can be allocated is of size 32 bytes. One bit is used to track each such chunk, where 1 means allocated and 0 means free. The chunk looks like this as of now:

10011010

Now, there is a request for a chunk of 50 bytes.

After this allocation, the bitmap, indicating the status of the buddy allocator will be

Answer: 10110010

✗

The correct answer is: 11111010

**Question 29**

Partially correct

Mark 0.75 out of 1.00

Select the correct points of comparison between POSIX and System V shared memory.

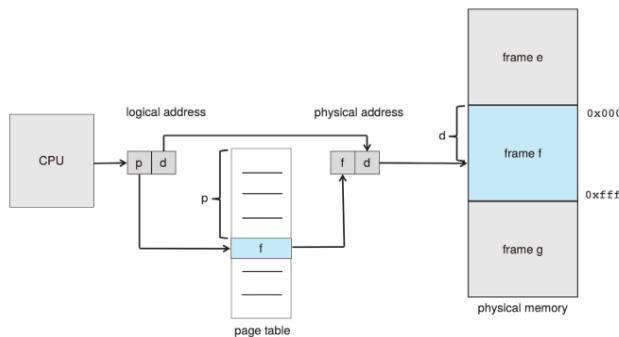
- a. POSIX shared memory is newer than System V shared memory ✓
- b. POSIX shared memory is "thread safe", System V is not ✓
- c. System V is more prevalent than POSIX even today ✓
- d. POSIX allows giving name to shared memory, System V does not

The correct answers are: POSIX shared memory is newer than System V shared memory, POSIX shared memory is "thread safe", System V is not, POSIX allows giving name to shared memory, System V does not, System V is more prevalent than POSIX even today

**Question 30**

Partially correct

Mark 0.67 out of 1.00

**Figure 9.8** Paging hardware.

Mark the statements as True or False, w.r.t. the above diagram (note that the diagram does not cover all details of what actually happens!)

True	False	
<input checked="" type="radio"/>	<input type="radio"/>	The combining of f and d is done by MMU
<input type="radio"/>	<input checked="" type="radio"/>	There are total 3 memory references in this diagram
<input checked="" type="radio"/>	<input type="radio"/>	The split of logical address into p and d is done by MMU
<input checked="" type="radio"/>	<input type="radio"/>	The page table is in physical memory and must be continuous
<input type="radio"/>	<input checked="" type="radio"/>	Using the offset d in the physical page-frame is done by MMU
<input checked="" type="radio"/>	<input type="radio"/>	The logical address issued by CPU is the same one generated by compiler

The combining of f and d is done by MMU: True

There are total 3 memory references in this diagram: False

The split of logical address into p and d is done by MMU: True

The page table is in physical memory and must be continuous: True

Using the offset d in the physical page-frame is done by MMU: False

The logical address issued by CPU is the same one generated by compiler: True

**Question 31**

Partially correct

Mark 0.50 out of 1.00

Select all the correct statements about signals

Select one or more:

- a. SIGKILL definitely kills a process because its code runs in kernel mode of CPU
- b. Signals are delivered to a process by another process ✗
- c. The signal handler code runs in kernel mode of CPU
- d. SIGKILL definitely kills a process because it can't be caught or ignored, and its default action terminates the process ✓
- e. The signal handler code runs in user mode of CPU ✓
- f. A signal handler can be invoked asynchronously or synchronously depending on signal type ✓
- g. Signal handlers once replaced can't be restored
- h. Signals are delivered to a process by kernel

Your answer is partially correct.

You have correctly selected 3.

The correct answers are: Signals are delivered to a process by kernel, A signal handler can be invoked asynchronously or synchronously depending on signal type, The signal handler code runs in user mode of CPU, SIGKILL definitely kills a process because it can't be caught or ignored, and its default action terminates the process

**Question 32**

Correct

Mark 1.00 out of 1.00

The data structure used in kalloc() and kfree() in xv6 is

- a. Singly linked circular list
- b. Singly linked NULL terminated list ✓
- c. Double linked NULL terminated list
- d. Doubly linked circular list

Your answer is correct.

The correct answer is: Singly linked NULL terminated list

**Question 33**

Partially correct

Mark 1.78 out of 2.00

Match the description of a memory management function with the name of the function that provides it, in xv6

Load contents from ELF into existing pages	loaduvm()	✓
Mark the page as in-accessible	clearpteu()	✓
setup the kernel part in the page table	setupkvm()	✓
Switch to kernel page table	switchkvm()	✓
Create a copy of the page table of a process	copyuvm()	✓
Copy the code pages of a process	No such function	✓
Setup and load the user page table for initcode process	inituvm()	✓
Switch to user page table	switchuvm()	✓
Load contents from ELF into pages after allocating the pages first	inituvm()	✗

The correct answer is: Load contents from ELF into existing pages → loaduvm(), Mark the page as in-accessible → clearpteu(), setup the kernel part in the page table → setupkvm(), Switch to kernel page table → switchkvm(), Create a copy of the page table of a process → copyuvm(), Copy the code pages of a process → No such function, Setup and load the user page table for initcode process → inituvm(), Switch to user page table → switchuvm(), Load contents from ELF into pages after allocating the pages first → No such function

**Question 34**

Partially correct

Mark 0.60 out of 1.00

Mark the statements as True or False, w.r.t. thrashing

True	False	
<input checked="" type="radio"/> ✘	<input type="radio"/> ✓	Thrashing occurs because some process is doing lot of disk I/O.
<input checked="" type="radio"/> ✓	<input checked="" type="radio"/> ✘	Processes keep changing their locality of reference, and a high rate of page faults occur when they are changing the locality.
<input checked="" type="radio"/> ✘	<input checked="" type="radio"/> ✓	mmap() solves the problem of thrashing.
<input checked="" type="radio"/> ✓	<input checked="" type="radio"/> ✘	The working set model is an attempt at approximating the locality of a process.
<input checked="" type="radio"/> ✓	<input checked="" type="radio"/> ✘	Thrashing is particular to demand paging systems, and does not apply to pure paging systems.
<input checked="" type="radio"/> ✘	<input type="radio"/> ✓	Processes keep changing their locality of reference, and least number of page faults occur when they are changing the locality.
<input checked="" type="radio"/> ✘	<input checked="" type="radio"/> ✓	Thrashing can occur even if entire memory is not in use.
<input checked="" type="radio"/> ✓	<input checked="" type="radio"/> ✘	During thrashing the CPU is under-utilised as most time is spent in I/O
<input checked="" type="radio"/> ✓	<input checked="" type="radio"/> ✘	Thrashing can be limited if local replacement is used.
<input checked="" type="radio"/> ✓	<input checked="" type="radio"/> ✘	Thrashing occurs when the total size of all processes's locality exceeds total memory size.

Thrashing occurs because some process is doing lot of disk I/O.: False

Processes keep changing their locality of reference, and a high rate of page faults occur when they are changing the locality.: True

mmap() solves the problem of thrashing.: False

The working set model is an attempt at approximating the locality of a process.: True

Thrashing is particular to demand paging systems, and does not apply to pure paging systems.: True

Processes keep changing their locality of reference, and least number of page faults occur when they are changing the locality.: False

Thrashing can occur even if entire memory is not in use.: False

During thrashing the CPU is under-utilised as most time is spent in I/O: True

Thrashing can be limited if local replacement is used.: True

Thrashing occurs when the total size of all processes's locality exceeds total memory size.: True

**Question 35**

Correct

Mark 1.00 out of 1.00

After virtual memory is implemented

(select T/F for each of the following) One Program's size can be larger than physical memory size

True	False	
<input checked="" type="radio"/>	<input type="radio"/> ✗	Cumulative size of all programs can be larger than physical memory size
<input checked="" type="radio"/>	<input type="radio"/> ✗	Code need not be completely in memory
<input checked="" type="radio"/>	<input type="radio"/> ✗	One Program's size can be larger than physical memory size
<input type="radio"/> ✗	<input checked="" type="radio"/>	Virtual addresses become available to executing process
<input type="radio"/> ✗	<input checked="" type="radio"/>	Virtual access to memory is granted to all processes
<input checked="" type="radio"/>	<input type="radio"/> ✗	Relatively less I/O may be possible during process execution
<input checked="" type="radio"/>	<input type="radio"/> ✗	Logical address space could be larger than physical address space

Cumulative size of all programs can be larger than physical memory size: True

Code need not be completely in memory: True

One Program's size can be larger than physical memory size: True

Virtual addresses become available to executing process: False

Virtual access to memory is granted to all processes: False

Relatively less I/O may be possible during process execution: True

Logical address space could be larger than physical address space: True

◀ (Optional Assignment) lseek system call in xv6

Jump to...

Feedback on Quiz-2 ►

**Started on** Thursday, 18 March 2021, 2:46 PM

**State** Finished

**Completed on** Thursday, 18 March 2021, 3:50 PM

**Time taken** 1 hour 4 mins

**Grade** 10.36 out of 20.00 (52%)

**Question 1**

Partially correct

Mark 0.57 out of 1.00

Mark True, the actions done as part of code of swtch() in swtch.S, in xv6

**True**

**False**

<input checked="" type="radio"/>	<input checked="" type="radio"/>	Restore new callee saved registers from kernel stack of new context	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Save old callee saved registers on kernel stack of old context	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Save old callee saved registers on user stack of old context	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Switch from old process context to new process context	✗
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Switch from one stack (old) to another(new)	✗
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Restore new callee saved registers from user stack of new context	✓
<input checked="" type="radio"/>	<input checked="" type="radio"/>	Jump to code in new context	✗

Restore new callee saved registers from kernel stack of new context: True

Save old callee saved registers on kernel stack of old context: True

Save old callee saved registers on user stack of old context: False

Switch from old process context to new process context: False

Switch from one stack (old) to another(new): True

Restore new callee saved registers from user stack of new context: False

Jump to code in new context: False

**Question 2**

Partially correct

Mark 0.17 out of 0.50

For each function/code-point, select the status of segmentation setup in xv6

bootmain()	gdt setup with 3 entries, right from first line of code of bootloader	✗
kvmalloc() in main()	gdt setup with 5 entries (0 to 4) on one processor	✗
after startothers() in main()	gdt setup with 5 entries (0 to 4) on all processors	✓
after seginit() in main()	gdt setup with 5 entries (0 to 4) on all processors	✗
bootasm.S	gdt setup with 3 entries, right from first line of code of bootloader	✗
entry.S	gdt setup with 3 entries, at start32 symbol of bootasm.S	✓

Your answer is partially correct.

You have correctly selected 2.

The correct answer is: bootmain() → gdt setup with 3 entries, at start32 symbol of bootasm.S, kvmalloc() in main() → gdt setup with 3 entries, at start32 symbol of bootasm.S, after startothers() in main() → gdt setup with 5 entries (0 to 4) on all processors, after seginit() in main() → gdt setup with 5 entries (0 to 4) on one processor, bootasm.S → gdt setup with 3 entries, at start32 symbol of bootasm.S, entry.S → gdt setup with 3 entries, at start32 symbol of bootasm.S

**Question 3**

Partially correct

Mark 0.38 out of 1.00

Compare paging with demand paging and select the correct statements.

Select one or more:

- a. The meaning of valid-invalid bit in page table is different in paging and demand-paging. ✓
- b. Demand paging requires additional hardware support, compared to paging. ✓
- c. Paging requires some hardware support in CPU
- d. With paging, it's possible to have user programs bigger than physical memory. ✗
- e. Both demand paging and paging support shared memory pages. ✓
- f. Demand paging always increases effective memory access time.
- g. With demand paging, it's possible to have user programs bigger than physical memory. ✓
- h. Calculations of number of bits for page number and offset are same in paging and demand paging. ✓
- i. TLB hit ration has zero impact in effective memory access time in demand paging.
- j. Paging requires NO hardware support in CPU

Your answer is partially correct.

You have correctly selected 5.

The correct answers are: Demand paging requires additional hardware support, compared to paging., Both demand paging and paging support shared memory pages., With demand paging, it's possible to have user programs bigger than physical memory., Demand paging always increases effective memory access time., Paging requires some hardware support in CPU, Calculations of number of bits for page number and offset are same in paging and demand paging., The meaning of valid-invalid bit in page table is different in paging and demand-paging.

**Question 4**

Partially correct

Mark 0.44 out of 0.50

Suppose a processor supports base(relocation register) + limit scheme of MMU.

Assuming this, mark the statements as True/False

True	False	
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The OS may terminate the process while handling the interrupt of memory violation
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The hardware detects any memory access beyond the limit value and raises an interrupt
<input type="radio"/> <input checked="" type="checkbox"/>	<input type="radio"/>	The hardware may terminate the process while handling the interrupt of memory violation
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The OS sets up the relocation and limit registers when the process is scheduled
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The compiler generates machine code assuming continuous memory address space for process, and calculating appropriate sizes for code, and data;
<input type="radio"/> <input checked="" type="checkbox"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The process sets up its own relocation and limit registers when the process is scheduled
<input type="radio"/> <input checked="" type="checkbox"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The OS detects any memory access beyond the limit value and raises an interrupt
<input type="radio"/> <input checked="" type="checkbox"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The compiler generates machine code assuming appropriately sized segments for code, data and stack.

The OS may terminate the process while handling the interrupt of memory violation: True

The hardware detects any memory access beyond the limit value and raises an interrupt: True

The hardware may terminate the process while handling the interrupt of memory violation: False

The OS sets up the relocation and limit registers when the process is scheduled: True

The compiler generates machine code assuming continuous memory address space for process, and calculating appropriate sizes for code, and data;: True

The process sets up its own relocation and limit registers when the process is scheduled: False

The OS detects any memory access beyond the limit value and raises an interrupt: False

The compiler generates machine code assuming appropriately sized segments for code, data and stack.: False

**Question 5**

Correct

Mark 0.50 out of 0.50

Consider the following list of free chunks, in continuous memory management:

10k, 25k, 12k, 7k, 9k, 13k

Suppose there is a request for chunk of size 9k, then the free chunk selected under each of the following schemes will be

Best fit:

9k



First fit:

10k



Worst fit:

25k

**Question 6**

Partially correct

Mark 0.50 out of 1.00

Select all the correct statements about MMU and its functionality

Select one or more:

- a. MMU is a separate chip outside the processor
- b. MMU is inside the processor ✓
- c. Logical to physical address translations in MMU are done with specific machine instructions
- d. The operating system interacts with MMU for every single address translation ✗
- e. Illegal memory access is detected in hardware by MMU and a trap is raised ✓
- f. The Operating system sets up relevant CPU registers to enable proper MMU translations
- g. Logical to physical address translations in MMU are done in hardware, automatically ✓
- h. Illegal memory access is detected by operating system

Your answer is partially correct.

You have correctly selected 3.

The correct answers are: MMU is inside the processor, Logical to physical address translations in MMU are done in hardware, automatically, The Operating system sets up relevant CPU registers to enable proper MMU translations, Illegal memory access is detected in hardware by MMU and a trap is raised

**Question 7**

Incorrect

Mark 0.00 out of 0.50

Assuming a 8- KB page size, what is the page numbers for the address 874815 reference in decimal :  
 (give answer also in decimal)

Answer: 2186



The correct answer is: 107

**Question 8**

Incorrect

Mark 0.00 out of 0.25

Select the compiler's view of the process's address space, for each of the following MMU schemes:  
 (Assume that each scheme,e.g. paging/segmentation/etc is effectively utilised)

Segmentation, then paging	Many continuous chunks each of page size	
Relocation + Limit	Many continuous chunks of same size	
Segmentation	one continuous chunk	
Paging	many continuous chunks of variable size	

Your answer is incorrect.

The correct answer is: Segmentation, then paging → many continuous chunks of variable size, Relocation + Limit → one continuous chunk, Segmentation → many continuous chunks of variable size, Paging → one continuous chunk

**Question 9**

Incorrect

Mark 0.00 out of 0.50

Suppose the memory access time is 180ns and TLB hit ratio is 0.3, then effective memory access time is (in nanoseconds);

Answer: 192



The correct answer is: 306.00

**Question 10**

Correct

Mark 0.50 out of 0.50

In xv6, The struct context is given as

```
struct context {
    uint edi;
    uint esi;
    uint ebx;
    uint ebp;
    uint eip;
};
```

Select all the reasons that explain why only these 5 registers are included in the struct context.

- a. The segment registers are same across all contexts, hence they need not be saved ✓
- b. esp is not saved in context, because context{} is on stack and it's address is always argument to swtch() ✓
- c. xv6 tries to minimize the size of context to save memory space
- d. esp is not saved in context, because it's not part of the context
- e. eax, ecx, edx are caller save, hence no need to save ✓

Your answer is correct.

The correct answers are: The segment registers are same across all contexts, hence they need not be saved, eax, ecx, edx are caller save, hence no need to save, esp is not saved in context, because context{} is on stack and it's address is always argument to swtch()

**Question 11**

Partially correct

Mark 0.83 out of 1.50

Arrange the following events in order, in page fault handling:

Disk interrupt wakes up the process

7	✓
---	---

The reference bit is found to be invalid by MMU

1	✓
---	---

OS makes available an empty frame

6	✗
---	---

Restart the instruction that caused the page fault

9	✓
---	---

A hardware interrupt is issued

3	✗
---	---

OS schedules a disk read for the page (from backing store)

5	✓
---	---

Process is kept in wait state

4	✗
---	---

Page tables are updated for the process

8	✓
---	---

Operating system decides that the page was not in memory

2	✗
---	---

Your answer is partially correct.

You have correctly selected 5.

The correct answer is: Disk interrupt wakes up the process → 7, The reference bit is found to be invalid by MMU → 1, OS makes available an empty frame → 4, Restart the instruction that caused the page fault → 9, A hardware interrupt is issued → 2, OS schedules a disk read for the page (from backing store) → 5, Process is kept in wait state → 6, Page tables are updated for the process → 8, Operating system decides that the page was not in memory → 3

**Question 12**

Incorrect

Mark 0.00 out of 0.50

Suppose a kernel uses a buddy allocator. The smallest chunk that can be allocated is of size 32 bytes. One bit is used to track each such chunk, where 1 means allocated and 0 means free. The chunk looks like this as of now:

00001010

Now, there is a request for a chunk of 70 bytes.

After this allocation, the bitmap, indicating the status of the buddy allocator will be

Answer: 11101010



The correct answer is: 11111010

**Question 13**

Incorrect

Mark 0.00 out of 0.25

The complete range of virtual addresses (after main() in main.c is over), from which the free pages used by kalloc() and kfree() is derived, are:

- a. end, 4MB
- b. P2V(end), P2V(PHYSTOP)
- c. end, P2V(4MB + PHYSTOP)
- d. P2V(end), PHYSTOP ✗
- e. end, (4MB + PHYSTOP)
- f. end, PHYSTOP
- g. end, P2V(PHYSTOP)

Your answer is incorrect.

The correct answer is: end, P2V(PHYSTOP)

**Question 14**

Partially correct

Mark 0.33 out of 0.50

Match the pair

Hashed page table	Linear search on collision done by OS (e.g. SPARC Solaris) typically	✓
Inverted Page table	Linear/Parallel search using frame number in page table	✗
Hierarchical Paging	More memory access time per hierarchy	✓

Your answer is partially correct.

You have correctly selected 2.

The correct answer is: Hashed page table → Linear search on collision done by OS (e.g. SPARC Solaris) typically, Inverted Page table → Linear/Parallel search using page number in page table, Hierarchical Paging → More memory access time per hierarchy

**Question 15**

Partially correct

Mark 0.29 out of 0.50

After virtual memory is implemented

(select T/F for each of the following) One Program's size can be larger than physical memory size

True	False	
<input checked="" type="radio"/>	<input type="radio"/> ✗	Code need not be completely in memory
<input checked="" type="radio"/>	<input type="radio"/> ✗	Cumulative size of all programs can be larger than physical memory size
<input type="radio"/> ✗	<input checked="" type="radio"/>	Virtual access to memory is granted
<input checked="" type="radio"/>	<input type="radio"/> ✗	Logical address space could be larger than physical address space
<input type="radio"/> ✗	<input checked="" type="radio"/>	Virtual addresses are available
<input checked="" type="radio"/>	<input checked="" type="radio"/> ✗	Relatively less I/O may be possible during process execution
<input checked="" type="radio"/>	<input type="radio"/> ✗	One Program's size can be larger than physical memory size

Code need not be completely in memory: True

Cumulative size of all programs can be larger than physical memory size: True

Virtual access to memory is granted: False

Logical address space could be larger than physical address space: True

Virtual addresses are available: False

Relatively less I/O may be possible during process execution: True

One Program's size can be larger than physical memory size: True

**Question 16**

Partially correct

Mark 0.64 out of 1.00

W.r.t. Memory management in xv6,

xv6 uses physical memory upto 224 MB only  
Mark statements True or False**True      False**

<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The switchkvm() call in scheduler() is invoked after control comes to it from sched(), thus demanding execution in kernel's context	✓
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The stack allocated in entry.S is used as stack for scheduler's context for first processor	✓
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The switchkvm() call in scheduler() changes CR3 to use page directory kpgdir	✓
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The free page-frame are created out of nearly 222 MB	✗
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The kernel code and data take up less than 2 MB space	✓
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The switchkvm() call in scheduler() changes CR3 to use page directory of new process	✗
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The switchkvm() call in scheduler() is invoked after control comes to it from swtch() scheduler(), thus demanding execution in new process's context	✓
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	PHYSTOP can be increased to some extent, simply by editing memlayout.h	✓
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	xv6 uses physical memory upto 224 MB only	✗
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The process's address space gets mapped on frames, obtained from ~2MB:224MB range	✓
<input checked="" type="radio"/>	<input type="radio"/> <input checked="" type="checkbox"/>	The kernel's page table given by kpgdir variable is used as stack for scheduler's context	✗

The switchkvm() call in scheduler() is invoked after control comes to it from sched(), thus demanding execution in kernel's context: True

The stack allocated in entry.S is used as stack for scheduler's context for first processor: True

The switchkvm() call in scheduler() changes CR3 to use page directory kpgdir: True

The free page-frame are created out of nearly 222 MB: True

The kernel code and data take up less than 2 MB space: True

The switchkvm() call in scheduler() changes CR3 to use page directory of new process: False

The switchkvm() call in scheduler() is invoked after control comes to it from swtch() scheduler(), thus demanding execution in new process's context: False

PHYSTOP can be increased to some extent, simply by editing memlayout.h: True

xv6 uses physical memory upto 224 MB only: True

The process's address space gets mapped on frames, obtained from ~2MB:224MB range: True

The kernel's page table given by kpgdir variable is used as stack for scheduler's context: False

**Question 17**

Incorrect

Mark 0.00 out of 1.50

Consider the reference string

6 4 2 0 1 2 6 9 2 0 5

If the number of page frames is 3, then total number of page faults (including initial), using LRU replacement is:

Answer:  ✖

#6# 6,4# 6,4,2 # 0,4,2#0,1,2#6,1,2#6,9,2#0,9,2#0,5,2

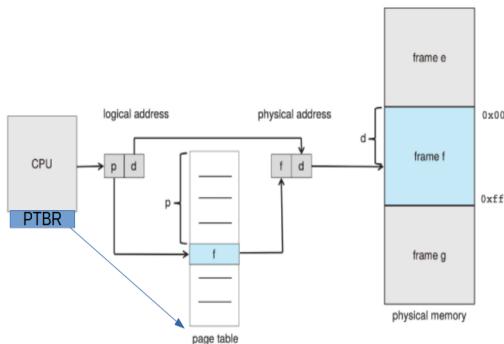
The correct answer is: 9

**Question 18**

Partially correct

Mark 0.31 out of 0.50

Consider the image given below, which explains how paging works.



**Figure 9.8** Paging hardware.

Mention whether each statement is True or False, with respect to this image.

True	False	
<input checked="" type="radio"/>	<input type="radio"/>	The PTBR is present in the CPU as a register
<input type="radio"/>	<input checked="" type="radio"/>	The page table is indexed using frame number
<input checked="" type="radio"/>	<input type="radio"/>	The page table is indexed using page number
<input type="radio"/>	<input checked="" type="radio"/>	The locating of the page table using PTBR also involves paging translation
<input type="radio"/>	<input checked="" type="radio"/>	Size of page table is always determined by the size of RAM
<input checked="" type="radio"/>	<input type="radio"/>	The page table is itself present in Physical memory
<input checked="" type="radio"/>	<input type="radio"/>	Maximum Size of page table is determined by number of bits used for page number
<input checked="" type="radio"/>	<input type="radio"/>	The physical address may not be of the same size (in bits) as the logical address

The PTBR is present in the CPU as a register: True

The page table is indexed using frame number: False

The page table is indexed using page number: True

The locating of the page table using PTBR also involves paging translation: False

Size of page table is always determined by the size of RAM: False

The page table is itself present in Physical memory: True

Maximum Size of page table is determined by number of bits used for page number: True

The physical address may not be of the same size (in bits) as the logical address: True

**Question 19**

Correct

Mark 2.00 out of 2.00

Given below is shared memory code with two processes sharing a memory segment.

The first process sends a user input string to second process. The second capitalizes the string. Then the first process prints the capitalized version.

Fill in the blanks to complete the code.

**// First process**

```
#define SHMSZ 27

int main()
{
    char c;
    int shmid;
    key_t key;
    char *shm, *s, string[128];
    key = 5679;
    if ((shmid =
        shmget
        ✓ (key, SHMSZ, IPC_CREAT | 0666)) < 0) {
        perror("shmget");
        exit(1);
    }
    if ((shm =
        shmat
        ✓ (shmid, NULL, 0)) == (char *) -1) {
        perror("shmat");
        exit(1);
    }
    s = shm;
    *s = '$';
    scanf("%s", string);
    strcpy(s + 1, string);
    *s =
        @
        ✓ ';' //note the quotes
    while(*s != '
        $
        ')
        sleep(1);
        printf("%s\n", s + 1);
        exit(0);
}
```

**//Second process**

```
#define SHMSZ 27

int main()
{
    int shmid;
    key_t key;
    char *shm, *s;
    int i;
    char string[128];
    key =
        5679
```

```

✓ ;
if ((shmid = shmget(key, SHMSZ, 0666)) < 0) {
    perror("shmget");
    exit(1);
}
if ((shm = shmat(shmid, NULL, 0)) == (char *) -1) {
    perror("shmat");
    exit(1);
}
s =

✓ ;
while(*s != '@')
    sleep(1);
for(i = 0; i < strlen(s + 1); i++)
    s[i + 1] = toupper(s[i + 1]);
*s = '$';
exit(0);
}

```

**Question 20**

Partially correct

Mark 0.25 out of 0.50

Map the functionality/use with function/variable in xv6 code.

return a free page, if available; 0, otherwise

Create page table entries for a given range of virtual and physical addresses; including page directory entries if needed

Array listing the kernel memory mappings, to be used by setupkvm()

Setup kernel part of a page table, mapping kernel code, data, read-only data, I/O space, devices

Return address of page table entry in a given page directory, for a given virtual address; creates page table if necessary

Setup kernel part of a page table, and switch to that page table

 kinit1()

 mappages()

 kmap[]

 kvmalloc()

 walkpgdir()

 setupkvm()

Your answer is partially correct.

You have correctly selected 3.

The correct answer is: return a free page, if available; 0, otherwise → kalloc(), Create page table entries for a given range of virtual and physical addresses; including page directory entries if needed → mappages(), Array listing the kernel memory mappings, to be used by setupkvm() → kmap[], Setup kernel part of a page table, mapping kernel code, data, read-only data, I/O space, devices → setupkvm(), Return address of page table entry in a given page directory, for a given virtual address; creates page table if necessary → walkpgdir(), Setup kernel part of a page table, and switch to that page table → kvmalloc()

**Question 21**

Partially correct

Mark 1.53 out of 2.50

Order events in xv6 timer interrupt code

(Transition from process P1 to P2's code.)

P2 is selected and marked RUNNING

12 ✓

Change of stack from user stack to kernel stack of P1

3 ✓

Timer interrupt occurs

2 ✓

alltraps() will call iret

17 ✗

change to context of P2, P2's kernel stack in use now

13 ✓

P2's trap() will return to alltraps

16 ✗

jump in vector.S

4 ✓

P2 will return from sched() in yield()

14 ✗

yield() is called

8 ✓

trap() is called

7 ✓

Process P2 is executing

18 ✗

P1 is marked as RUNNABLE

9 ✓

P2's yield() will return in trap()

15 ✗

Process P1 is executing

1 ✓

sched() is called,

11 ✗

change to context of the scheduler, scheduler's stack in use now

10 ✗

jump to alltraps

5 ✓

Trapframe is built on kernel stack of P1

6 ✓

Your answer is partially correct.

You have correctly selected 11.

The correct answer is: P2 is selected and marked RUNNING → 12, Change of stack from user stack to kernel stack of P1 → 3, Timer interrupt occurs → 2, alltraps() will call iret → 18, change to context of P2, P2's kernel stack in use now → 13, P2's trap() will return to alltraps → 17, jump in vector.S → 4, P2 will return from sched() in yield() → 15, yield() is called → 8, trap() is called → 7, Process P2 is executing → 14, P1 is marked as RUNNABLE → 9, P2's yield() will return in trap() → 16, Process P1 is executing → 1, sched() is called, → 10, change to context of the scheduler, scheduler's stack in use now → 11, jump to alltraps → 5, Trapframe is built on kernel stack of P1 → 6

**Question 22**

Incorrect

Mark 0.00 out of 1.00

Given that the memory access time is 200 ns, probability of a page fault is 0.7 and page fault handling time is 8 ms,  
The effective memory access time in nanoseconds is:

Answer:  ✖

The correct answer is: 5600060.00

**Question 23**

Correct

Mark 0.25 out of 0.25

Select the state that is not possible after the given state, for a process:

- New:  Running ✓
- Ready :  Waiting ✓
- Running:  None of these ✓
- Waiting:  Running ✓

**Question 24**

Partially correct

Mark 0.63 out of 1.00

Select the correct statements about sched() and scheduler() in xv6 code

- a. scheduler() switches to the selected process's context ✓
- b. When either sched() or scheduler() is called, it does not return immediately to caller ✓
- c. After call to swtch() in sched(), the control moves to code in scheduler()
- d. Each call to sched() or scheduler() involves change of one stack inside swtch() ✓
- e. After call to swtch() in scheduler(), the control moves to code in sched()
- f. When either sched() or scheduler() is called, it results in a context switch ✓
- g. sched() switches to the scheduler's context ✓
- h. sched() and scheduler() are co-routines

Your answer is partially correct.

You have correctly selected 5.

The correct answers are: sched() and scheduler() are co-routines, When either sched() or scheduler() is called, it does not return immediately to caller, When either sched() or scheduler() is called, it results in a context switch, sched() switches to the scheduler's context, scheduler() switches to the selected process's context, After call to swtch() in scheduler(), the control moves to code in sched(), After call to swtch() in sched(), the control moves to code in scheduler(), Each call to sched() or scheduler() involves change of one stack inside swtch()

**Question 25**

Correct

Mark 0.25 out of 0.25

The data structure used in kalloc() and kfree() in xv6 is

- a. Doubly linked circular list
- b. Singly linked circular list
- c. Double linked NULL terminated list
- d. Singly linked NULL terminated list



Your answer is correct.

The correct answer is: Singly linked NULL terminated list

[◀ \(Assignment\) lseek system call in xv6](#)

Jump to...

Dashboard / My courses / Computer Engineering & IT / CEIT-Even-sem-20-21 / QS-Even-sem-2020-21 / 16 May - 22 May / End Sem Exam OS-2021

Started on Saturday, 22 May 2021, 8:00 AM

State Finished

Completed on Saturday, 22 May 2021, 9:30 AM

Time taken 1 hour 30 mins

Grade 26.12 out of 40.00 (65%)

Question 1

Incorrect

Mark 0.00 out of 1.00

A 4 GB disk with 1 KB of block size would require these many number of **blocks** for its free block bitmap:

Answer: 4096 ✖

The correct answer is: 512

Question 2

Correct

Mark 1.00 out of 1.00

Given that the memory access time is 110 ns, probability of a page fault is 0.5 and page fault handling time is 12 ms,

The effective memory access time in nanoseconds is:

Answer: 6000165 ✓

The correct answer is: 6000055.00

Question 3

Incorrect

Mark 0.00 out of 1.00

The maximum size of a file in number of blocks of BSIZE in xv6 code is

(write a number only)

Answer: 268 ✖

The correct answer is: 138

Question 4

Incorrect

Mark 0.00 out of 1.00

Calculate the average waiting time using

Round Robin scheduling with time quantum of 5 time units  
for the following workload

assuming that they arrive in the order written below.

Process Burst Time

P1	5
P2	7
P3	6
P4	2

Write only a number in the answer upto two decimal points.

Answer: 40.75 ✖

The correct answer is: 10.25



**Question 5**

Correct

Mark 1.00 out of 1.00

For the reference string

4 2 5 1 0 1 2 5 4 1 2

the number of page faults, including initial ones,  
with FIFO replacement and 2 frames are :

Answer: 10 ✓

4 -

4 2

5 2

5 1

0 1

-

2 1

2 5

4 5

4 1

2 1

The correct answer is: 10

**Question 6**

Correct

Mark 1.00 out of 1.00

Assuming a 16- KB page size, what is the page number for the address 428517 reference in decimal :

(give answer also in decimal)

Answer: 27 ✓

The correct answer is: 26



**Question 7**

Correct

Mark 1.00 out of 1.00

In the code below assume that each function can be executed concurrently by many threads/processes.  
Ignore syntactical issues, and focus on the semantics.

This program is an example of

```
spinlock a, b; // assume initialized
thread1() {
    spinlock(b);
    //some code;
    spinlock(a);
    //some code;
    spinunlock(b);
    spinunlock(a);
}
thread2() {
    spinlock(a);
    //some code;
    spinlock(b);
    //some code;
    spinunlock(b);
    spinunlock(a);
}
```

- a. Deadlock ✓
- b. Self Deadlock
- c. None of these
- d. Deadlock or livelock depending on actual race
- e. Livelock

Your answer is correct.

The correct answer is: Deadlock



**Question 8**

Partially correct

Mark 1.33 out of 2.00

Match the snippets of xv6 code with the core functionality they achieve, or problems they avoid.  
"..." means some code.

```
static inline uint
xchg(volatile uint *addr, uint newval)
{
    uint result;

    // The + in "+m" denotes a read-modify-write operand.
    asm volatile("lock; xchgl %0, %1" :
        "+m" ("addr"), "=a" (result) :
        "1" (newval) :
        "cc");
    return result;
}
```

Atomic compare and swap instruction (to be expanded inline into code)



```
void
sleep(void *chan, struct spinlock *lk)
{
    ...
    if(lk != &ptable.lock){
        acquire(&ptable.lock);
        release(lk);
    }
}
```

If you don't do this, a process may be running on two processors parallelly



```
void
acquire(struct spinlock *lk)
{
    ...
    __sync_synchronize();
}
```

Tell compiler not to reorder memory access beyond this line



Your answer is partially correct.

You have correctly selected 2.

The correct answer is: static inline uint  
xchg(volatile uint \*addr, uint newval)

```
{
    uint result;
```

// The + in "+m" denotes a read-modify-write operand.

```
asm volatile("lock; xchgl %0, %1" :
    "+m" ("addr"), "=a" (result) :
    "1" (newval) :
    "cc");
return result;
} → Atomic compare and swap instruction (to be expanded inline into code), void
sleep(void *chan, struct spinlock *lk)
{
    ...
    if(lk != &ptable.lock){
        acquire(&ptable.lock);
        release(lk);
    } → Avoid a self-deadlock, void
    acquire(struct spinlock *lk)
{
    ...
    __sync_synchronize(); → Tell compiler not to reorder memory access beyond this line
}
```

**Question 9**

Correct

Mark 1.00 out of 1.00

Predict the output of the program given here.

Assume that all the path names for the programs are correct. For example "/usr/bin/echo" will actually run echo command.

Assume that there is no mixing of print output on screen if two of them run concurrently.

In the answer replace a new line by a single space.

For example::

good

output

should be written as good output

--

```
main() {  
    int i;  
    i = fork();  
    if(i == 0)  
        execl("/usr/bin/echo", "/usr/bin/echo", "hi", 0);  
    else  
        wait(0);  
    fork();  
    execl("/usr/bin/echo", "/usr/bin/echo", "one", 0);  
}
```

Answer: hi one one 

The correct answer is: hi one one

**Question 10**

Partially correct

Mark 1.67 out of 2.00

Select all the blocks that may need to be written back to disk (if updated, of-course), as "Yes", when an operation of deleting a file is carried out on ext2 file system.

An option has to be correct entirely to be marked "Yes"

Superblock

Yes 

One or multiple data blocks of the parent directory

No 

One or more data bitmap blocks for the parent directory

No 

Block bitmap(s) for all the blocks of the file

No 

Possibly one block bitmap corresponding to the parent directory

Yes 

Data blocks of the file

No 

Your answer is partially correct.

only one data block of parent directory. multiple blocks not possible. an entry is always contained within one single block

You have correctly selected 5.

The correct answer is: Superblock → Yes, One or multiple data blocks of the parent directory → No, One or more data bitmap blocks for the parent directory → No, Block bitmap(s) for all the blocks of the file → Yes, Possibly one block bitmap corresponding to the parent directory → Yes, Data blocks of the file → No

**Question 11**

Correct

Mark 1.00 out of 1.00

Select all the correct statements about bootloader.

Every wrong selection will deduct marks proportional to  $1/n$  where n is total wrong choices in the question.

You will get minimum a zero.

- a. Modern Bootloaders often allow configuring the way an OS boots
- b. Bootloaders allow selection of OS to boot from
- c. Bootloader must be one sector in length
- d. The bootloader loads the BIOS
- e. LILO is a bootloader



Your answer is correct.

The correct answers are: LILO is a bootloader, Modern Bootloaders often allow configuring the way an OS boots, Bootloaders allow selection of OS to boot from



## Question 12

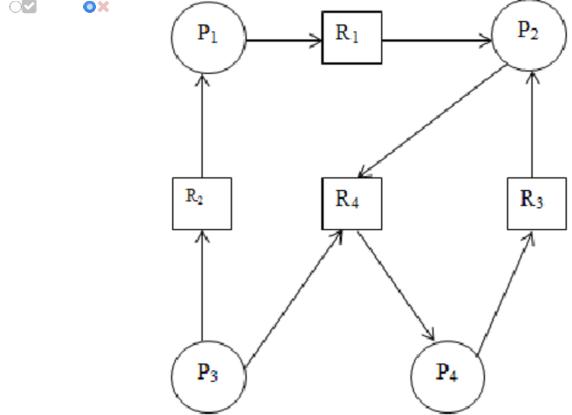
Incorrect

Mark 0.00 out of 1.00

For each of the resource allocation diagram shown,  
infer whether the graph contains at least one deadlock or not.

Yes

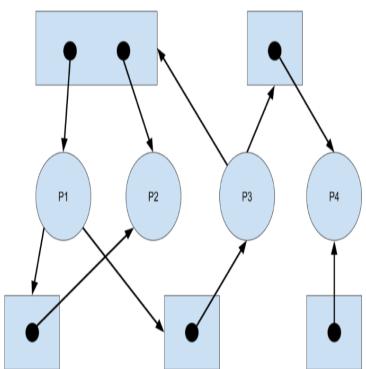
No



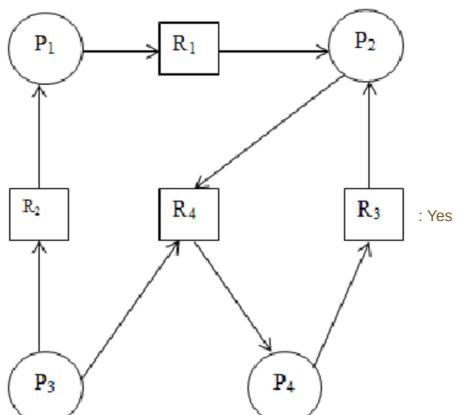
✗

✗

✓

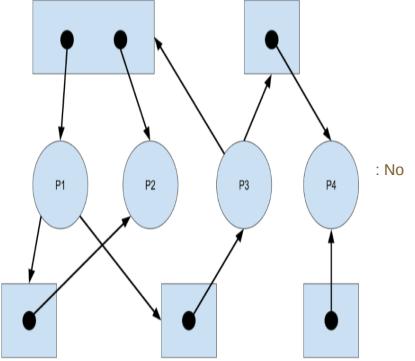


✗



: Yes



**Question 13**

Partially correct

Mark 0.71 out of 1.00

Mark the statements about device drivers by marking as True or False.

True	False	
<input checked="" type="radio"/>	<input type="radio"/> ✘	It's possible that a particular hardware has multiple device drivers available for it.
<input checked="" type="radio"/>	<input type="radio"/> ✘	xv6 has device drivers for IDE disk and console.
<input checked="" type="radio"/>	<input type="radio"/> ✘	A disk driver converts OS's logical view of disk into physical locations on disk.
<input checked="" type="radio"/>	<input type="radio"/> ✘	A device driver code is specific to a hardware device
<input checked="" type="radio"/>	<input type="radio"/> ✘	All devices of the same type (e.g. 2 hard disks) can typically use the same device driver
<input checked="" type="radio"/>	<input type="radio"/> ✘	Writing a device driver mandatorily demands reading the technical documentation about the hardware.
<input type="radio"/> ✘	<input checked="" type="radio"/>	Device driver is an intermediary between the end-user and OS

It's possible that a particular hardware has multiple device drivers available for it.: True

xv6 has device drivers for IDE disk and console.: True

A disk driver converts OS's logical view of disk into physical locations on disk.: True

A device driver code is specific to a hardware device: True

All devices of the same type (e.g. 2 hard disks) can typically use the same device driver: True

Writing a device driver mandatorily demands reading the technical documentation about the hardware.: True

Device driver is an intermediary between the end-user and OS: False

**Question 14**

Partially correct

Mark 0.33 out of 1.00

Consider this program.

Some statements are identified using the // comment at the end.

Assume that `=` is an atomic operation.

```
#include <stdio.h>
#include <pthread.h>
long c = 0, c1 = 0, c2 = 0, run = 1;
void *thread1(void *arg) {
    while(run == 1) { //E
        c = 10; //A
        c1 = c2 + 5; //B
    }
}
void *thread2(void *arg) {
    while(run == 1) { //F
        c = 20; //C
        c2 = c1 + 3; //D
    }
}
int main() {
    pthread_t th1, th2;
    pthread_create(&th1, NULL, thread1, NULL);
    pthread_create(&th2, NULL, thread2, NULL);
    sleep(2);
    run = 0;
    printf(stdout, "c = %ld c1+c2 = %ld c1 = %ld c2 = %ld \n", c, c1+c2, c1, c2);
    fflush(stdout);
}
```

Which statements are part of the critical Section?

Yes	No	
<input checked="" type="radio"/> <input checked="" type="checkbox"/>	<input type="radio"/> <input checked="" type="checkbox"/>	F
<input checked="" type="radio"/> <input checked="" type="checkbox"/>	<input checked="" type="radio"/> <input type="checkbox"/>	D
<input checked="" type="radio"/> <input checked="" type="checkbox"/>	<input type="radio"/> <input checked="" type="checkbox"/>	C
<input checked="" type="radio"/> <input checked="" type="checkbox"/>	<input type="radio"/> <input checked="" type="checkbox"/>	A
<input checked="" type="radio"/> <input checked="" type="checkbox"/>	<input checked="" type="radio"/> <input type="checkbox"/>	B
<input checked="" type="radio"/> <input checked="" type="checkbox"/>	<input type="radio"/> <input checked="" type="checkbox"/>	E

F: No

D: Yes

C: No

A: No

B: Yes

E: No

**Question 15**

Partially correct

Mark 1.43 out of 2.00

Mark statements as T/F

All statements are in the context of preventing deadlocks.

**True****False**

<input checked="" type="radio"/>	<input type="radio"/>	A process holding one resources and waiting for just one more resource can also be involved in a deadlock.	✓
<input type="radio"/>	<input checked="" type="radio"/>	If a resource allocation graph contains a cycle then there is a guarantee of a deadlock	✗
<input type="radio"/>	<input checked="" type="radio"/>	The lock ordering to be followed to avoid circular wait is a code in OS that checks for compliance with decided order	✗
<input checked="" type="radio"/>	<input type="radio"/>	Circular wait is avoided by enforcing a lock ordering	✓
<input checked="" type="radio"/>	<input type="radio"/>	Hold and wait means a thread/process holding some locks and waiting for acquiring some.	✓
<input checked="" type="radio"/>	<input type="radio"/>	Deadlock is possible if all the conditions are met at the same time: Mutual exclusion, hold and wait, no pre-emption, circular wait.	✓
<input checked="" type="radio"/>	<input type="radio"/>	Mutual exclusion is a necessary condition for deadlock because it brings in locks on which deadlock happens	✓

A process holding one resources and waiting for just one more resource can also be involved in a deadlock.: True

If a resource allocation graph contains a cycle then there is a guarantee of a deadlock: False

The lock ordering to be followed to avoid circular wait is a code in OS that checks for compliance with decided order: False

Circular wait is avoided by enforcing a lock ordering: True

Hold and wait means a thread/process holding some locks and waiting for acquiring some.: True

Deadlock is possible if all the conditions are met at the same time: Mutual exclusion, hold and wait, no pre-emption, circular wait.: True

Mutual exclusion is a necessary condition for deadlock because it brings in locks on which deadlock happens: True

**Question 16**

Correct

Mark 1.00 out of 1.00

Match the left side use(or non-use) of a synchronization primitive with the best option on the right side.

This is the smallest primitive made available in software, using the hardware provided atomic instructions

 spinlock ✓

This tool is useful for event-wait scenarios

 semaphore ✓

This tool is more useful on multiprocessor systems

 spinlock ✓

This tool is quite attractive in solving the main bounded buffer problem

 semaphore ✓

This tool is very useful for waiting for 'something'

 condition variables ✓

Your answer is correct.

The correct answer is: This is the smallest primitive made available in software, using the hardware provided atomic instructions → spinlock, This tool is useful for event-wait scenarios → semaphore, This tool is more useful on multiprocessor systems → spinlock, This tool is quite attractive in solving the main bounded buffer problem → semaphore, This tool is very useful for waiting for 'something' → condition variables

**Question 17**

Correct

Mark 1.00 out of 1.00

The permissions -rwx--x--x on a file mean

- a. The file can be read only by the owner
- b. 'cat' on the file by owner will not work
- c. 'cat' on the file by any user will work
- d. 'rm' on the file by any user will work
- e. The file can be executed by anyone
- f. The file can be written only by the owner



Your answer is correct.

The correct answers are: The file can be executed by anyone, The file can be read only by the owner, The file can be written only by the owner, 'rm' on the file by any user will work

**Question 18**

Incorrect

Mark 0.00 out of 1.00

Note: for this question you get full marks if you select all and only correct options, you get ZERO if at least one option is wrong or not selected.

Select all the correct statements about log structured file systems.

- a. a transaction is said to be committed when all operations are written to file system
- b. log may be kept on same block device or another block device
- c. file system recovery may end up losing data
- d. even if file systems followed immediate writes (i.e. non-delayed writes), it could still require recovery
- e. file system recovery recovers all the lost data



Your answer is incorrect.

The correct answers are: file system recovery may end up losing data, log may be kept on same block device or another block device, even if file systems followed immediate writes (i.e. non-delayed writes), it could still require recovery

## Question 19

Incorrect

Mark 0.00 out of 1.00

Consider the structure of directory entry in ext2, as shown in this diagram.

	inode	rec_len	file_type	name_len	name
0	21	12	1	2	.
12	22	12	2	2	.
24	53	16	5	2	h o m e
40	67	28	3	2	u s r
52	0	16	7	1	o l d f i l e
68	34	12	4	2	s b i n

Select the correct statements about the directory entry in ext2 file system.

The correct formula for rec\_len is (when entries are continuously stored)

- a.  $\text{rec\_len} = \text{sizeof(inode entry)} + \text{sizeof(name len entry)} + \text{sizeof(file type entry)} + (\text{strlen(name)} + (-1) * (\text{strlen(name)} \% 4))$
- b.  $\text{rec\_len} = \text{sizeof(inode entry)} + \text{sizeof(name len entry)} + \text{sizeof(file type entry)} + (\text{strlen(name)} + (\text{strlen(name)} - 4) \% 4)$
- c.  $\text{rec\_len} = \text{sizeof(inode entry)} + \text{sizeof(name len entry)} + \text{sizeof(file type entry)} + (\text{strlen(name)} + 4 - (\text{strlen(name)} \% 4))$  ✗
- d.  $\text{rec\_len} = \text{sizeof(inode entry)} + \text{sizeof(name len entry)} + \text{sizeof(file type entry)} + (\text{strlen(name)} + (-1) * (\text{strlen(name)} - 4))$
- e.  $\text{rec\_len} = \text{sizeof(inode entry)} + \text{sizeof(name len entry)} + \text{sizeof(file type entry)} + (\text{strlen(name)} \% 4)$
- f.  $\text{rec\_len} = \text{sizeof(inode entry)} + \text{sizeof(name len entry)} + \text{sizeof(file type entry)} + \text{strlen(name)}$

Your answer is incorrect.

The correct answer is:  $\text{rec\_len} = \text{sizeof(inode entry)} + \text{sizeof(name len entry)} + \text{sizeof(file type entry)} + (\text{strlen(name)} + (-1) * (\text{strlen(name)} - 4))$

## Question 20

Partially correct

Mark 0.50 out of 1.00

Mark whether the given sequence of events is possible or not-possible. Also, select the reason for your answer.

For each sequence it's a not-possible sequence if some important event is not mentioned in the sequence.

Assume that the kernel code is non-interruptible and uniprocessor system.

Process P1 executing a system call  
 Timer interrupt  
 Generic interrupt handler runs  
 Scheduler runs  
 Scheduler selects P2 for execution  
 P2 returns from timer interrupt handler  
 Process p2, user code executing

This sequence of events is:  ✓

Because

✗

**Question 21**

Incorrect

Mark 0.00 out of 1.00

The given semaphore implementation faces which problem?

Assume any suitable code for signal()

Note: blocks means waits in a wait queue.

```
struct semaphore {  
    int val;  
    spinlock lk;  
};  
sem_init(semaphore *s, int initval) {  
    s->val = initval;  
    s->sl = 0;  
}  
wait(semaphore *s) {  
    spinlock(&(s->sl));  
    while(s->val <=0)  
        ;  
    (s->val)--;  
    spinunlock(&(s->sl));  
}
```

- a. blocks holding a spinlock
- b. deadlock
- c. too much spinning, bounded wait not guaranteed
- d. not holding lock after unblock



Your answer is incorrect.

The correct answer is: deadlock



## Question 22

Partially correct

Mark 0.80 out of 1.00

Mark statements True/False w.r.t. change of states of a process.

Reference: The process state diagram (and your understanding of how kernel code works)

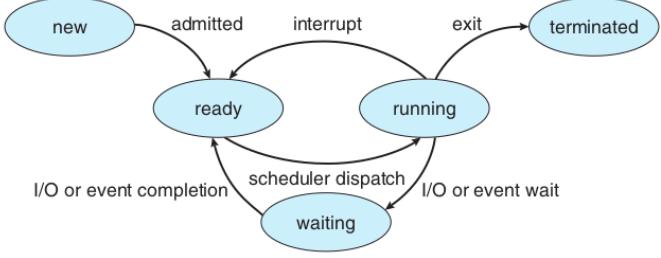


Figure 3.2 Diagram of process state.

## True

## False

<input type="radio"/> ✗	<input checked="" type="checkbox"/> ✗	A process in RUNNING state only can become TERMINATED because scheduler moves it to ZOMBIE state	✓
<input checked="" type="checkbox"/> ✗	<input type="radio"/> ✗	A process in READY state can not go to WAITING state because the resource on which it will WAIT will not be in use when process is in READY state.	✗
<input checked="" type="checkbox"/> ✗	<input type="radio"/> ✗	A process in WAITING state can not become RUNNING because the event it's waiting for has not occurred	✓
<input checked="" type="checkbox"/> ✗	<input type="radio"/> ✗	Every process has to go through ZOMBIE state, at least for a small duration.	✓
<input checked="" type="checkbox"/> ✗	<input type="radio"/> ✗	Only a process in READY state is considered by scheduler	✓

A process in RUNNING state only can become TERMINATED because scheduler moves it to ZOMBIE state: False

A process in READY state can not go to WAITING state because the resource on which it will WAIT will not be in use when process is in READY state.: False

A process in WAITING state can not become RUNNING because the event it's waiting for has not occurred: True

Every process has to go through ZOMBIE state, at least for a small duration.: True

Only a process in READY state is considered by scheduler: True

**Question 23**

Correct

Mark 1.00 out of 1.00

Select T/F for statements about Volume Managers.

Do pay attention to the use of the words physical partition and physical volume.

**True      False**

<input checked="" type="radio"/>	<input type="radio"/> ✗	The volume manager can create further internal sub-divisions of a physical partition for efficiency or features.	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> ✗	A logical volume can be extended in size but upto the size of volume group	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> ✗	A logical volume may span across multiple physical volumes	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> ✗	The volume manager stores additional metadata on the physical disk partitions	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> ✗	A physical partition should be initialized as a physical volume, before it can be used by volume manager.	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> ✗	A volume group consists of multiple physical volumes	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	<input type="radio"/> ✗	A logical volume may span across multiple physical partitions	<input checked="" type="checkbox"/> since a physical volume is made up of physical partitions, and a volume can span across multiple PVs, it can also span across multiple PP

The volume manager can create further internal sub-divisions of a physical partition for efficiency or features.: True

A logical volume can be extended in size but upto the size of volume group: True

A logical volume may span across multiple physical volumes: True

The volume manager stores additional metadata on the physical disk partitions: True

A physical partition should be initialized as a physical volume, before it can be used by volume manager.: True

A volume group consists of multiple physical volumes: True

A logical volume may span across multiple physical partitions: True

**Question 24**

Correct

Mark 1.00 out of 1.00

Map the block allocation scheme with the problem it suffers from

(Match pairs 1-1, match a scheme with the problem that it suffers from relatively the most, compared to others)

Continuous allocation	need for compaction	<input checked="" type="checkbox"/>
Linked allocation	Too many seeks	<input checked="" type="checkbox"/>
Indexed Allocation	Overhead of reading metadata blocks	<input checked="" type="checkbox"/>

Your answer is correct.

The correct answer is: Continuous allocation → need for compaction, Linked allocation → Too many seeks, Indexed Allocation → Overhead of reading metadata blocks

**Question 25**

Correct

Mark 1.00 out of 1.00

This one is not a system call:

- a. open
- b. read
- c. write
- d. scheduler



Your answer is correct.

The correct answer is: scheduler



**Question 26**

Correct

Mark 1.00 out of 1.00

Match the pairs.

This question is based on your general knowledge about operating systems/related concepts and their features.

Java threads	monitors,re-entrant locks, semaphores	✓
Linux threads	atomic-instructions, spinlocks, etc.	✓
POSIX threads	semaphore, mutex, condition variables	✓

Your answer is correct.

The correct answer is: Java threads → monitors,re-entrant locks, semaphores, Linux threads → atomic-instructions, spinlocks, etc., POSIX threads → semaphore, mutex, condition variables

**Question 27**

Correct

Mark 1.00 out of 1.00

Consider the following list of free chunks, in continuous memory management:

7k, 15k, 21k, 14k, 19k, 6k

Suppose there is a request for chunk of size 5k, then the free chunk selected under each of the following schemes will be

Best fit:	6k	✓
First fit:	7k	✓
Worst fit:	21k	✓

**Question 28**

Correct

Mark 1.00 out of 1.00

This one is not a scheduling algorithm

- a. Round Robin
- b. SJF
- c. Mergesort
- d. FCFS



Your answer is correct.

The correct answer is: Mergesort

## Question 29

Correct

Mark 1.00 out of 1.00

Mark whether the concept is related to scheduling or not.

Yes	No	
<input checked="" type="radio"/>	<input type="radio"/>	timer interrupt
<input checked="" type="radio"/>	<input type="radio"/>	context-switch
<input checked="" type="radio"/>	<input type="radio"/>	ready-queue
<input type="radio"/>	<input checked="" type="radio"/>	file-table
<input checked="" type="radio"/>	<input type="radio"/>	runnable process

timer interrupt: Yes

context-switch: Yes

ready-queue: Yes

file-table: No

runnable process: Yes



**Question 30**

Partially correct

Mark 1.00 out of 2.00

Map ext2 data structure features with their purpose

**Many copies of Superblock** Choose...**Free blocks count in superblock and group descriptor**

Redundancy to ensure the most crucial data structure is not lost

**Used directories count in group descriptor**

is redundant and helps do calculations of directory entries faster

**Combining file type and access rights in one variable**

saves 1 byte of space

**rec\_len field in directory entry**

Try to keep all the data of a directory and its file close together in a group

**File Name is padded**

aligns all memory accesses on word boundary, improving performance

**Inode bitmap is one block**

limits total number of files that can belong to a group

**Block bitmap is one block**

Limits the size of a block group, thus improvising on purpose of a group

**Mount count in superblock**

to enforce file check after certain amount of mounts at boot time

**Inode table location in Group Descriptor**

is redundant and helps do calculations of directory entries faster

**Inode table**

All inodes are kept together so that one disk read leads to reading many inodes together, effectively doing a buffering of subsequent inode reads, and to save space on disk

**A group**

Redundancy to ensure the most crucial data structure is not lost



Your answer is partially correct.

You have correctly selected 6.

The correct answer is: **Many copies of Superblock** → Redundancy to ensure the most crucial data structure is not lost, **Free blocks count in superblock and group descriptor** → Redundancy to help fsck restore consistency, **Used directories count in group descriptor** → attempt is made to evenly spread the first-level directories, this count is used there, **Combining file type and access rights in one variable** → saves 1 byte of space, **rec\_len field in directory entry** → allows holes and linking of entries in directory, File Name is padded → aligns all memory accesses on word boundary, improving performance, **Inode bitmap is one block** → limits total number of files that can belong to a group, **Block bitmap is one block** → Limits the size of a block group, thus improvising on purpose of a group, **Mount count in superblock** → to enforce file check after certain amount of mounts at boot time, **Inode table location in Group Descriptor** → Obvious, as it's per group and not per file-system, **Inode table** → All inodes are kept together so that one disk read leads to reading many inodes together, effectively doing a buffering of subsequent inode reads, and to save space on disk, **A group** → Try to keep all the data of a directory and its file close together in a group

**Question 31**

Partially correct

Mark 1.85 out of 2.00

Mark True/False

Statements about scheduling and scheduling algorithms

**True****False**

<input checked="" type="radio"/>	<input type="radio"/> ✗	The nice() system call is used to set priorities for processes	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	Aging is used to ensure that low-priority processes do not starve in priority scheduling.	✓
<input type="radio"/>	<input checked="" type="radio"/> ✗	In non-pre-emptive priority scheduling, the highest priority process is scheduled and runs until it gives up CPU.	✗
<input checked="" type="radio"/>	<input type="radio"/> ✗	xv6 code does not care about Processor Affinity	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	In pre-emptive priority scheduling, priority is implemented by assigning more time quantum to higher priority process.	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	A scheduling algorithm is non-preemptive if it does context switch only if a process voluntarily relinquishes CPU or it terminates.	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	Processor Affinity refers to memory accesses of a process being stored on cache of that processor	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	Response time will be quite poor on non-interruptible kernels	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	Shortest Remaining Time First algorithm is nothing but pre-emptive Shortest Job First algorithm	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	On Linuxes the CPU utilisation is measured as the time spent in scheduling the idle thread	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	Generally the voluntary context switches are much more than non-voluntary context switches on a Linux system.	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	Pre-emptive scheduling leads to many race conditions in kernel code.	✓
<input checked="" type="radio"/>	<input type="radio"/> ✗	Statistical observations tell us that most processes have large number of small CPU bursts and relatively smaller numbers of large CPU bursts.	✓

The nice() system call is used to set priorities for processes.: True

Aging is used to ensure that low-priority processes do not starve in priority scheduling.: True

In non-pre-emptive priority scheduling, the highest priority process is scheduled and runs until it gives up CPU.: True

xv6 code does not care about Processor Affinity: True

In pre-emptive priority scheduling, priority is implemented by assigning more time quantum to higher priority process.: True

A scheduling algorithm is non-preemptive if it does context switch only if a process voluntarily relinquishes CPU or it terminates.: True

Processor Affinity refers to memory accesses of a process being stored on cache of that processor: True

Response time will be quite poor on non-interruptible kernels: True

Shortest Remaining Time First algorithm is nothing but pre-emptive Shortest Job First algorithm: True

On Linuxes the CPU utilisation is measured as the time spent in scheduling the idle thread: True

Generally the voluntary context switches are much more than non-voluntary context switches on a Linux system.: True

Pre-emptive scheduling leads to many race conditions in kernel code.: True

Statistical observations tell us that most processes have large number of small CPU bursts and relatively smaller numbers of large CPU bursts.: True

**Question 32**

Partially correct

Mark 1.17 out of 2.00

The unix file semantics demand that changes to any open file are visible immediately to any other processes accessing that file at that point in time.

Select the data-structure/programmatic features that ensure the implementation of unix semantics. (Assume there is no mmap())

Yes	No	
<input type="radio"/> <input checked="" type="checkbox"/>	All processes accessing the same file share the file descriptor among themselves	✓
<input type="radio"/> <input checked="" type="checkbox"/>	The pointer entry in the file descriptor array entry points to the data of the file directly	✓
<input checked="" type="checkbox"/> <input type="radio"/>	There is only one global file structure per on-disk file.	✗
<input type="radio"/> <input checked="" type="checkbox"/>	All file accesses are made using only global variables	✓
<input checked="" type="checkbox"/> <input type="radio"/>	The 'file offset' is shared among all the processes that access the file.	✗
<input type="radio"/> <input checked="" type="checkbox"/>	No synchronization is implemented so that changes are made available immediately.	✓
<input checked="" type="checkbox"/> <input type="radio"/>	A single spinlock is to be used to protect the unique global 'file structure' representing the file, thus synchronizing access, and making other processes wait for earlier process to finish writing so that writes get visible immediately.	✗
<input checked="" type="checkbox"/> <input type="radio"/>	There is only one in-memory copy of the on disk file's contents in kernel memory/buffers	✓
<input checked="" type="checkbox"/> <input type="radio"/>	The file descriptors in every PCB are pointers to the same global file structure.	✗
<input type="radio"/> <input checked="" type="checkbox"/>	The file descriptor array is external to PCB and all processes that share a file, have pointers to same file-descriptors' array	✓
<input checked="" type="checkbox"/> <input type="radio"/>	All file structures representing any open file, give access to the same in-memory copy of the file's contents	✓
<input checked="" type="checkbox"/> <input type="radio"/>	The 'file offset' index is stored outside the file-structure to which file-descriptor array points	✗

All processes accessing the same file share the file descriptor among themselves: No

The pointer entry in the file descriptor array entry points to the data of the file directly: No

There is only one global file structure per on-disk file.: No

All file accesses are made using only global variables: No

The 'file offset' is shared among all the processes that access the file.: No

No synchronization is implemented so that changes are made available immediately.: No

A single spinlock is to be used to protect the unique global 'file structure' representing the file, thus synchronizing access, and making other processes wait for earlier process to finish writing so that writes get visible immediately.: No

There is only one in-memory copy of the on disk file's contents in kernel memory/buffers: Yes

The file descriptors in every PCB are pointers to the same global file structure.: No

The file descriptor array is external to PCB and all processes that share a file, have pointers to same file-descriptors' array: No

All file structures representing any open file, give access to the same in-memory copy of the file's contents: Yes

The 'file offset' index is stored outside the file-structure to which file-descriptor array points: No

**Question 33**

Partially correct

Mark 0.33 out of 2.00

Map the function in xv6's file system code, to its perceived logical layer.

namei	inode	✗
filestat()	Choose...	
dirlookup	directory	✓
ialloc	file descriptor	✗
stati	Choose...	
ideintr	buffer cache	✗
bread	Choose...	
balloc	file descriptor	✗
sys_chdir()	system call	✓
skipelem	system call	✗
commit	system call	✗
bmap	system call	✗

Your answer is partially correct.

You have correctly selected 2.

The correct answer is: namei → pathname lookup, filestat() → file descriptor, dirlookup → directory, ialloc → inode, stati → inode, ideintr → disk driver, bread → buffer cache, balloc → block allocation on disk, sys\_chdir() → system call, skipelem → pathname lookup, commit → logging, bmap → inode

[◀ Course Exit Feedback](#)[Jump to...](#)[xv6-public-master ►](#)

**Started on** Monday, 24 January 2022, 7:07:42 PM

**State** Finished

**Completed on** Monday, 24 January 2022, 8:08:11 PM

**Time taken** 1 hour

**Grade** 8.90 out of 20.00 (45%)

#### Question 1

Complete

Mark 0.80 out of 1.00

Match the register with the segment used with it.

ebp	ss
eip	cs
edi	ds
esp	ss
esi	ds

The correct answer is: ebp → ss, eip → cs, edi → es, esp → ss, esi → ds

#### Question 2

Complete

Mark 1.00 out of 1.00

```
int value = 5;
int main()
{
    pid_t pid;
    pid = fork();
    if (pid == 0) { /* child process */
        value += 15;
        return 0;
    }
    else if (pid > 0) { /* parent process */
        wait(NULL);
        printf("%d", value); /* LINE A */
    }
    return 0;
}
```

What's the value printed here at LINE A?

Answer:

The correct answer is: 5

**Question 3**

Complete

Mark 0.50 out of 0.50

Is the command "cat README > done &" possible on xv6? (Note the & in the end)

- a. no
- b. yes

The correct answer is: yes

**Question 4**

Complete

Mark 0.00 out of 2.00

xv6.img: bootblock kernel

```
dd if=/dev/zero of=xv6.img count=10000
dd if=bootblock of=xv6.img conv=notrunc
dd if=kernel of=xv6.img seek=1 conv=notrunc
```

Consider above lines from the Makefile. Which of the following is incorrect?

- a. The xv6.img is the virtual disk that is created by combining the bootblock and the kernel file.
- b. The xv6.img is of the size 10,000 blocks of 512 bytes each and occupies upto 10,000 blocks on the disk.
- c. The size of xv6.img is exactly = (size of bootblock) + (size of kernel)
- d. xv6.img is the virtual processor used by the qemu emulator
- e. The size of the kernel file is nearly 5 MB
- f. Blocks in xv6.img after kernel may be all zeroes.
- g. The xv6.img is of the size 10,000 blocks of 512 bytes each and occupies 10,000 blocks on the disk.
- h. The kernel is located at block-1 of the xv6.img
- i. The bootblock is located on block-0 of the xv6.img
- j. The bootblock may be 512 bytes or less (looking at the Makefile instruction)
- k. The size of the xv6.img is nearly 5 MB

The correct answers are: xv6.img is the virtual processor used by the qemu emulator, The xv6.img is of the size 10,000 blocks of 512 bytes each and occupies upto 10,000 blocks on the disk., The size of the kernel file is nearly 5 MB, The size of xv6.img is exactly = (size of bootblock) + (size of kernel)

**Question 5**

Complete

Mark 0.43 out of 1.00

Rank the following storage systems from slowest (first) to fastest(last)

You can drag and drop the items below/above each other.

Registers
Cache
Main memory
Nonvolatile memory
Magnetic tapes
Optical disk
Hard-disk drives

The correct order for these items is as follows:

1. Magnetic tapes
2. Optical disk
3. Hard-disk drives
4. Nonvolatile memory
5. Main memory
6. Cache
7. Registers

**Question 6**

Complete

Mark 1.00 out of 1.00

How does the distinction between kernel mode and user mode function as a rudimentary form of protection (security) ?

Select one:

- a. It disallows hardware interrupts when a process is running
- b. It prohibits one process from accessing other process's memory
- c. It prohibits a user mode process from running privileged instructions
- d. It prohibits invocation of kernel code completely, if a user program is running

The correct answer is: It prohibits a user mode process from running privileged instructions

**Question 7**

Complete

Mark 0.00 out of 2.00

Which of the following are NOT a part of job of a typical compiler?

- a. Suggest alternative pieces of code that can be written
- b. Check the program for syntactical errors
- c. Check the program for logical errors
- d. Convert high level language code to machine code
- e. Process the # directives in a C program
- f. Invoke the linker to link the function calls with their code, extern globals with their declaration

The correct answers are: Check the program for logical errors, Suggest alternative pieces of code that can be written

**Question 8**

Complete

Mark 0.00 out of 2.00

Match the program with its output (ignore newlines in the output. Just focus on the count of the number of 'hi')

main() { execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); }

hi hi

main() { int i = fork(); if(i == 0) execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); }

hi hi hi

main() { int i = NULL; fork(); printf("hi\n"); }

No output

main() { fork(); execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); }

hi

The correct answer is: main() { execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); } → hi, main() { int i = fork(); if(i == 0) execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); } → hi, main() { int i = NULL; fork(); printf("hi\n"); } → hi hi, main() { fork(); execl("/usr/bin/echo", "/usr/bin/echo", "hi\n", NULL); } → hi hi

**Question 9**

Complete

Mark 0.83 out of 2.00

Select all statements that correctly explain the use/purpose of system calls.

Select one or more:

- a. Provide an environment for process creation
- b. Switch from user mode to kernel mode
- c. Handle ALL types of interrupts
- d. Handle exceptions like division by zero
- e. Run each instruction of an application program
- f. Allow I/O device access to user processes
- g. Provide services for accessing files

The correct answers are: Switch from user mode to kernel mode, Provide services for accessing files, Allow I/O device access to user processes, Provide an environment for process creation

**Question 10**

Complete

Mark 0.50 out of 0.50

Compare multiprogramming with multitasking

- a. A multiprogramming system is not necessarily multitasking
- b. A multitasking system is not necessarily multiprogramming

The correct answer is: A multiprogramming system is not necessarily multitasking

**Question 11**

Complete

Mark 0.60 out of 1.00

Select all the correct statements about two modes of CPU operation

Select one or more:

- a. The two modes are essential for a multitasking system
- b. Some instructions are allowed to run only in user mode, while all instructions can run in kernel mode
- c. The software interrupt instructions change the mode from user mode to kernel mode and jumps to predefined location simultaneously
- d. The two modes are essential for a multiprogramming system
- e. There is an instruction like 'iret' to return from kernel mode to user mode

The correct answers are: The two modes are essential for a multiprogramming system, The two modes are essential for a multitasking system, There is an instruction like 'iret' to return from kernel mode to user mode, The software interrupt instructions change the mode from user mode to kernel mode and jumps to predefined location simultaneously, Some instructions are allowed to run only in user mode, while all instructions can run in kernel mode

**Question 12**

Complete

Mark 0.50 out of 0.50

Is the terminal a part of the kernel on GNU/Linux systems?

- a. no
- b. yes

The correct answer is: no

**Question 13**

Complete

Mark 1.00 out of 1.00

Why should a program exist in memory before it starts executing ?

- a. Because the hard disk is a slow medium
- b. Because the processor can run instructions and access data only from memory
- c. Because the variables of the program are stored in memory
- d. Because the memory is volatile

The correct answer is: Because the processor can run instructions and access data only from memory

**Question 14**

Complete

Mark 1.33 out of 2.00

Which of the following instructions should be privileged?

Select one or more:

- a. Read the clock.
- b. Access memory management unit of the processor
- c. Access I/O device.
- d. Turn off interrupts.
- e. Set value of a memory location
- f. Set value of timer.
- g. Access a general purpose register
- h. Modify entries in device-status table
- i. Switch from user to kernel mode.

The correct answers are: Set value of timer., Access memory management unit of the processor, Turn off interrupts., Modify entries in device-status table, Access I/O device., Switch from user to kernel mode.

**Question 15**

Complete

Mark 0.07 out of 2.00

Select all the correct statements about calling convention on x86 32-bit.

- a. The ebp pointers saved on the stack constitute a chain of activation records
- b. The return value is either stored on the stack or returned in the eax register
- c. The two lines in the beginning of each function, "push %ebp; mov %esp, %ebp", create space for local variables
- d. Parameters may be passed in registers or on stack
- e. Return address is one location above the ebp
- f. Parameters may be passed in registers or on stack
- g. Compiler may allocate more memory on stack than needed
- h. Space for local variables is allocated by subtracting the stack pointer inside the code of the called function
- i. Parameters are pushed on the stack in left-right order
- j. during execution of a function, ebp is pointing to the old ebp
- k. Space for local variables is allocated by subtracting the stack pointer inside the code of the caller function

The correct answers are: Compiler may allocate more memory on stack than needed, Parameters may be passed in registers or on stack, Parameters may be passed in registers or on stack, Return address is one location above the ebp, during execution of a function, ebp is pointing to the old ebp, Space for local variables is allocated by subtracting the stack pointer inside the code of the called function, The ebp pointers saved on the stack constitute a chain of activation records

**Question 16**

Complete

Mark 0.33 out of 0.50

Order the following events in boot process (from 1 onwards)

Shell	3
BIOS	1
Init	4
OS	6
Login interface	5
Boot loader	2

The correct answer is: Shell → 6, BIOS → 1, Init → 4, OS → 3, Login interface → 5, Boot loader → 2

[◀ \(Task\) Compulsory xv6 task](#)

Jump to...

[\(Optional Assignment\) Shell Programming\(Conformance tests\) ▶](#)

**Started on** Wednesday, 9 February 2022, 7:00:12 PM

**State** Finished

**Completed on** Wednesday, 9 February 2022, 7:46:38 PM

**Time taken** 46 mins 26 secs

**Grade** 3.00 out of 11.00 (27%)

**Question 1**

Complete

Mark 0.00 out of 1.00

The number of GDT entries setup during boot process of xv6 is

- a. 2
- b. 3
- c. 0
- d. 256
- e. 4
- f. 255

The correct answer is: 3

**Question 2**

Complete

Mark 0.00 out of 1.00

x86 provides which of the following type of memory management options?

- a. segmentation and one level paging
- b. segmentation or one or two level paging
- c. segmentation and two level paging
- d. segmentation or paging
- e. segmentation and one or two level paging
- f. segmentation only

The correct answer is: segmentation and one or two level paging

**Question 3**

Complete

Mark 0.00 out of 1.00

which of the following is not a difference between real mode and protected mode

- a. in real mode general purpose registers are 16 bit, in protected mode they are 32 bit
- b. in real mode the addressable memory is more than in protected mode
- c. in real mode the addressable memory is less than in protected mode
- d. in real mode the segment is multiplied by 16, in protected mode segment is used as index in GDT
- e. processor starts in real mode

The correct answer is: in real mode the addressable memory is more than in protected mode

**Question 4**

Complete

Mark 0.00 out of 1.00

The kernel ELF file contains how many Program headers?

- a. 4
- b. 2
- c. 3
- d. 9
- e. 10

The correct answer is: 3

**Question 5**

Not answered

Marked out of 0.50

code line, MMU setting: Match the line of xv6 code with the MMU setup employed

Answer:

The correct answer is: inb \$0x64,%al

**Question 6**

Complete

Mark 1.00 out of 1.00

The kernel is loaded at Physical Address

- a. 0x000100000
- b. 0x00010000
- c. 0x80100000
- d. 0x800000000

The correct answer is: 0x000100000

**Question 7**

Complete

Mark 0.00 out of 1.00

Why is the code of entry() in Assembly and not in C?

- a. Because the symbol entry() is inside the ELF file
- b. Because the kernel code must begin in assembly
- c. Because it needs to setup paging
- d. There is no particular reason, it could also be in C

The correct answer is: Because it needs to setup paging

**Question 8**

Complete

Mark 1.00 out of 1.00

The ljmp instruction in general does

- a. change the CS and EIP to 32 bit mode, and jumps to next line of code
- b. change the CS and EIP to 32 bit mode, and jumps to new value of EIP
- c. change the CS and EIP to 32 bit mode
- d. change the CS and EIP to 32 bit mode, and jumps to kernel code

The correct answer is: change the CS and EIP to 32 bit mode, and jumps to new value of EIP

**Question 9**

Complete

Mark 0.00 out of 1.00

The variable \$stack in entry.S is

- a. located at the value given by %esp as setup by bootmain()
- b. located at 0x7c00
- c. a memory region allocated as a part of entry.S
- d. located at less than 0x7c00
- e. located at 0

The correct answer is: a memory region allocated as a part of entry.S

**Question 10**

Not answered

Marked out of 0.50

Match the pairs of which action is taken by whom

Answer:

The correct answer is: kernel

**Question 11**

Complete

Mark 0.00 out of 1.00

ELF Magic number is

- a. 0xFFFFFFFF
- b. 0
- c. 0xELFELFELF
- d. 0xELF
- e. 0x0x464CELF
- f. 0x464C457FL
- g. 0x464C457FU

The correct answer is: 0x464C457FU

**Question 12**

Complete

Mark 1.00 out of 1.00

The right side of line of code "entry = (void(\*)(void))(elf->entry)" means

- a. Convert the "entry" in ELF structure into void
- b. Get the "entry" in ELF structure and convert it into a function pointer accepting no arguments and returning nothing
- c. Get the "entry" in ELF structure and convert it into a function void pointer
- d. Get the "entry" in ELF structure and convert it into a void pointer

The correct answer is: Get the "entry" in ELF structure and convert it into a function pointer accepting no arguments and returning nothing

[◀ Homework questions: Basics of MM, xv6 booting](#)

Jump to...

[\(Code\) Files, redirection, dup, \(IPC\)pipe ▶](#)

Started on Monday, 21 February 2022, 7:00:28 PM

State Finished

Completed on Monday, 21 February 2022, 7:49:24 PM

Time taken 48 mins 56 secs

Grade 8.30 out of 10.00 (83%)

Question 1

Complete

Mark 0.80 out of 1.00

Match the elements of C program to their place in memory

Local Static variables	Data
Global variables	Data
Code of main()	Code
Function code	Code
Arguments	Stack
Mallocoed Memory	Heap
#include files	No Memory needed
#define MACROS	No memory needed
Local Variables	Stack
Global Static variables	Data

The correct answer is: Local Static variables → Data, Global variables → Data, Code of main() → Code, Function code → Code, Arguments → Stack, Mallocoed Memory → Heap, #include files → No memory needed, #define MACROS → No Memory needed, Local Variables → Stack, Global Static variables → Data

**Question 2**

Complete

Mark 1.00 out of 1.00

What will be the output of this program

```
int main() {
    int fd;
    printf("%d ", open("/etc/passwd", O_RDONLY));
    close(1);
    fd = printf("%d ", open("/etc/passwd", O_RDONLY));
    close(fd);
    fd = printf("%d ", open("/etc/passwd", O_RDONLY));
}
```

- a. 3 1 2
- b. 3 3 3
- c. 3 1 1
- d. 1 1 1
- e. 3 4 5
- f. 2 2 2

The correct answer is: 3 1 1

**Question 3**

Complete

Mark 1.00 out of 1.00

Arrange in correct order, the files involved in execution of system call

vectors.S	2
trap.c	4
usys.S	1
trapasm.S	3

The correct answer is: vectors.S → 2, trap.c → 4, usys.S → 1, trapasm.S → 3

**Question 4**

Complete

Mark 1.00 out of 1.00

The "push 0" in vectors.S is

- a. A placeholder to match the size of struct trapframe
- b. Place for the error number value
- c. To indicate that it's a system call and not a hardware interrupt
- d. To be filled in as the return value of the system call

The correct answer is: Place for the error number value

**Question 5**

Complete

Mark 0.00 out of 1.00

A process blocks itself means

- a. The kernel code of an interrupt handler, moves the process to a waiting queue and calls scheduler
- b. The application code calls the scheduler
- c. The kernel code of system call, called by the process, moves the process to a waiting queue and calls scheduler
- d. The kernel code of system call calls scheduler

The correct answer is: The kernel code of system call, called by the process, moves the process to a waiting queue and calls scheduler

**Question 6**

Complete

Mark 1.00 out of 1.00

Select the odd one out

- a. Kernel stack of new process to Process stack of new process
- b. Process stack of running process to kernel stack of running process
- c. Kernel stack of running process to kernel stack of scheduler
- d. Kernel stack of scheduler to kernel stack of new process
- e. Kernel stack of new process to kernel stack of scheduler

The correct answer is: Kernel stack of new process to kernel stack of scheduler

**Question 7**

Complete

Mark 0.50 out of 0.50

Match the File descriptors to their meaning

0 Standard Input

2 Standard error

1 Standard output

The correct answer is: 0 → Standard Input, 2 → Standard error, 1 → Standard output

**Question 8**

Complete

Mark 1.00 out of 1.00

Which of the following is not a task of the code of swtch() function

- a. Switch stacks
- b. Change the kernel stack location
- c. Load the new context
- d. Jump to next context EIP
- e. Save the return value of the old context code
- f. Save the old context

The correct answers are: Save the return value of the old context code, Change the kernel stack location

**Question 9**

Complete

Mark 0.50 out of 0.50

Match the names of PCB structures with kernel

linux struct task\_struct

xv6 struct proc

The correct answer is: linux → struct task\_struct, xv6 → struct proc

**Question 10**

Complete

Mark 0.50 out of 0.50

Match the MACRO with its meaning

KERNBASE	2 GB
KERNLINK	2.224 GB
PHYSTOP	224 MB

The correct answer is: KERNBASE → 2 GB, KERNLINK → 2.224 GB, PHYSTOP → 224 MB

**Question 11**

Complete

Mark 1.00 out of 1.00

The trapframe, in xv6, is built by the

- a. hardware, trapasm.S
- b. hardware, vectors.S
- c. hardware, vectors.S, trapasm.S, trap()
- d. hardware, vectors.S, trapasm.S
- e. vectors.S, trapasm.S

The correct answer is: hardware, vectors.S, trapasm.S

**Question 12**

Complete

Mark 0.00 out of 0.50

Which of the following state transitions are not possible?

- a. Running -> Waiting
- b. Ready -> Waiting
- c. Waiting -> Terminated
- d. Ready -> Terminated

The correct answers are: Ready -&gt; Terminated, Waiting -&gt; Terminated, Ready -&gt; Waiting

[◀ Description of some possible course mini projects](#)

Jump to...

[\(Code\) mmap related programs ▶](#)

**Started on** Saturday, 26 February 2022, 5:18:30 PM

**State** Finished

**Completed on** Saturday, 26 February 2022, 6:30:44 PM

**Time taken** 1 hour 12 mins

**Grade** 8.55 out of 15.00 (57%)

**Question 1**

Complete

Mark 0.50 out of 0.50

Map the technique with its feature/problem

static linking	large executable file
dynamic loading	allocate memory only if needed
dynamic linking	small executable file
static loading	wastage of physical memory

The correct answer is: static linking → large executable file, dynamic loading → allocate memory only if needed, dynamic linking → small executable file, static loading → wastage of physical memory

**Question 2**

Complete

Mark 0.00 out of 1.00

Calculate the EAT in NANO-seconds (upto 2 decimal points) w.r.t. a page fault, given

Memory access time = 299 ns

Average page fault service time = 6 ms

Page fault rate = 0.8

Answer: 4.80

The correct answer is: 4800059.80

**Question 3**

Complete

Mark 1.00 out of 1.00

Given six memory partitions of 300 KB , 600 KB , 350 KB , 200 KB , 750 KB , and 125 KB (in order), how would the first-fit, best-fit, and worst-fit algorithms place processes of size 115 KB and 500 KB (in order)?

best fit 115 KB	125 KB
best fit 500 KB	600 KB
worst fit 500 KB	635 KB
worst fit 115 KB	750 KB
first fit 500 KB	600 KB
first fit 115 KB	300 KB

The correct answer is: best fit 115 KB → 125 KB, best fit 500 KB → 600 KB, worst fit 500 KB → 635 KB, worst fit 115 KB → 750 KB, first fit 500 KB → 600 KB, first fit 115 KB → 300 KB

**Question 4**

Complete

Mark 0.29 out of 1.00

Compare paging with demand paging and select the correct statements.

Select one or more:

- a. TLB hit ration has zero impact in effective memory access time in demand paging.
- b. Both demand paging and paging support shared memory pages.
- c. Calculations of number of bits for page number and offset are same in paging and demand paging.
- d. Demand paging requires additional hardware support, compared to paging.
- e. The meaning of valid-invalid bit in page table is different in paging and demand-paging.
- f. Paging requires NO hardware support in CPU
- g. With paging, it's possible to have user programs bigger than physical memory.
- h. With demand paging, it's possible to have user programs bigger than physical memory.
- i. Paging requires some hardware support in CPU
- j. Demand paging always increases effective memory access time.

The correct answers are: Demand paging requires additional hardware support, compared to paging., Both demand paging and paging support shared memory pages., With demand paging, it's possible to have user programs bigger than physical memory., Demand paging always increases effective memory access time., Paging requires some hardware support in CPU, Calculations of number of bits for page number and offset are same in paging and demand paging., The meaning of valid-invalid bit in page table is different in paging and demand-paging.

**Question 5**

Complete

Mark 0.36 out of 0.50

Map the parts of a C code to the memory regions they are related to

local variables	stack
global un-initialized variables	bss
static variables	data
global initialized variables	data
function arguments	stack
malloced memory	stack
functions	stack

The correct answer is: local variables → stack, global un-initialized variables → bss, static variables → data, global initialized variables → data, function arguments → stack, malloced memory → heap, functions → code

**Question 6**

Complete

Mark 0.75 out of 1.00

which of the following, do you think, are valid concerns for making the kernel pageable?

- a. The kernel must have some dedicated frames for it's own work
- b. No data structure of kernel should be pageable
- c. The kernel's own page tables should not be pageable
- d. The disk driver and disk interrupt handler should not be pageable
- e. No part of kernel code should be pageable.
- f. The page fault handler should not be pageable

The correct answers are: The kernel's own page tables should not be pageable, The page fault handler should not be pageable, The kernel must have some dedicated frames for it's own work, The disk driver and disk interrupt handler should not be pageable

## Question 7

Complete

Mark 0.75 out of 1.00

W.r.t the figure given below, mark the given statements as True or False.

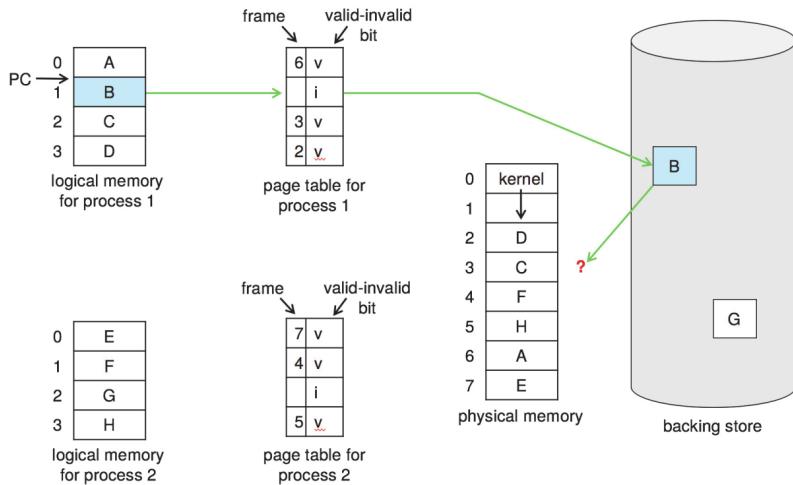


Figure 10.9 Need for page replacement.

True      False

- Kernel occupies two page frames
- Handling this scenario demands two disk I/Os
- Local replacement means chose any of the frames 2, 3, 6
- Local replacement means chose any of the frame from 2 to 7
- Global replacement means chose any of the frame from 0 to 7
- Global replacement means chose any of the frame from 2 to 7
- The kernel's pages can not used for replacement if kernel is not pageable.
- Page 1 of process 1 needs a replacement

Kernel occupies two page frames: True

Handling this scenario demands two disk I/Os: True

Local replacement means chose any of the frames 2, 3, 6: True

Local replacement means chose any of the frame from 2 to 7: False

Global replacement means chose any of the frame from 0 to 7: False

Global replacement means chose any of the frame from 2 to 7: True

The kernel's pages can not used for replacement if kernel is not pageable.: True

Page 1 of process 1 needs a replacement: True

**Question 8**

Complete

Mark 0.67 out of 1.00

Shared memory is possible with which of the following memory management schemes ?

Select one or more:

- a. paging
- b. segmentation
- c. continuous memory management
- d. demand paging

The correct answers are: paging, segmentation, demand paging

**Question 9**

Complete

Mark 0.60 out of 1.00

Given below is the "maps" file for a particular instance of "vim.basic" process.

Mark the given statements as True or False, w.r.t. the contents of the map file.

55a43501b000-55a435049000 r--p 00000000 103:05 917529	/usr/bin/vim.basic
55a435049000-55a435248000 r-xp 0002e000 103:05 917529	/usr/bin/vim.basic
55a435248000-55a4352b6000 r--p 0022d000 103:05 917529	/usr/bin/vim.basic
55a4352b7000-55a4352c5000 r--p 0029b000 103:05 917529	/usr/bin/vim.basic
55a4352c5000-55a4352e2000 rw-p 002a9000 103:05 917529	/usr/bin/vim.basic
55a4352e2000-55a4352f0000 rw-p 00000000 00:00 0	[heap]
55a436bc9000-55a436e5b000 rw-p 00000000 00:00 0	/usr/lib/x86_64-linux-
7f275b0a3000-7f275b0a6000 r--p 00000000 103:05 917901	/usr/lib/x86_64-linux-
gnu/libnss_files-2.31.so	gnu/libnss_files-2.31.so
7f275b0a6000-7f275b0ad000 r-xp 00003000 103:05 917901	/usr/lib/x86_64-linux-
gnu/libnss_files-2.31.so	gnu/libnss_files-2.31.so
7f275b0ad000-7f275b0af000 r--p 0000a000 103:05 917901	/usr/lib/x86_64-linux-
gnu/libnss_files-2.31.so	gnu/libnss_files-2.31.so
7f275b0af000-7f275b0b0000 r--p 0000b000 103:05 917901	/usr/lib/x86_64-linux-
gnu/libnss_files-2.31.so	gnu/libnss_files-2.31.so
7f275b0b0000-7f275b0b1000 rw-p 0000c000 103:05 917901	/usr/lib/x86_64-linux-
gnu/libnss_files-2.31.so	gnu/libnss_files-2.31.so
7f275b0b1000-7f275b0b7000 rw-p 00000000 00:00 0	
7f275b0b7000-7f275b8f5000 r--p 00000000 103:05 925247	/usr/lib/locale/locale-archive
7f275b8f5000-7f275b8fa000 rw-p 00000000 00:00 0	
7f275b8fa000-7f275b8fc000 r--p 00000000 103:05 924216	/usr/lib/x86_64-linux-
gnu/libogg.so.0.8.4	gnu/libogg.so.0.8.4
7f275b8fc000-7f275b901000 r-xp 00002000 103:05 924216	/usr/lib/x86_64-linux-
gnu/libogg.so.0.8.4	gnu/libogg.so.0.8.4
7f275b901000-7f275b904000 r--p 00007000 103:05 924216	/usr/lib/x86_64-linux-
gnu/libogg.so.0.8.4	gnu/libogg.so.0.8.4
7f275b904000-7f275b905000 ---p 0000a000 103:05 924216	
gnu/libogg.so.0.8.4	gnu/libogg.so.0.8.4
7f275b905000-7f275b906000 r--p 0000a000 103:05 924216	
gnu/libogg.so.0.8.4	gnu/libogg.so.0.8.4
7f275b906000-7f275b907000 rw-p 0000b000 103:05 924216	
gnu/libogg.so.0.8.4	gnu/libogg.so.0.8.4
7f275b907000-7f275b90a000 r--p 00000000 103:05 924627	
gnu/libvorbis.so.0.4.8	gnu/libvorbis.so.0.4.8
7f275b90a000-7f275b921000 r-xp 00003000 103:05 924627	
gnu/libvorbis.so.0.4.8	gnu/libvorbis.so.0.4.8
7f275b921000-7f275b932000 r--p 0001a000 103:05 924627	
gnu/libvorbis.so.0.4.8	gnu/libvorbis.so.0.4.8
7f275b932000-7f275b933000 ---p 0002b000 103:05 924627	
gnu/libvorbis.so.0.4.8	gnu/libvorbis.so.0.4.8
7f275b933000-7f275b934000 r--p 0002b000 103:05 924627	
gnu/libvorbis.so.0.4.8	gnu/libvorbis.so.0.4.8
7f275b934000-7f275b935000 rw-p 0002c000 103:05 924627	
gnu/libvorbis.so.0.4.8	gnu/libvorbis.so.0.4.8
7f275b935000-7f275b937000 rw-p 00000000 00:00 0	
7f275b937000-7f275b938000 r--p 00000000 103:05 917914	/usr/lib/x86_64-linux-
gnu/libutil-2.31.so	gnu/libutil-2.31.so
7f275b938000-7f275b939000 r-xp 00001000 103:05 917914	/usr/lib/x86_64-linux-
gnu/libutil-2.31.so	gnu/libutil-2.31.so
7f275b939000-7f275b93a000 r--p 00002000 103:05 917914	/usr/lib/x86_64-linux-
gnu/libutil-2.31.so	gnu/libutil-2.31.so
7f275b93a000-7f275b93b000 r--p 00002000 103:05 917914	/usr/lib/x86_64-linux-
gnu/libutil-2.31.so	gnu/libutil-2.31.so
7f275b93b000-7f275b93c000 rw-p 00003000 103:05 917914	/usr/lib/x86_64-linux-

```

gnu/libutil-2.31.so
7f275b93c000-7f275b93e000 r--p 00000000 103:05 915906      /usr/lib/x86_64-linux-
gnu/libz.so.1.2.11
7f275b93e000-7f275b94f000 r-xp 00002000 103:05 915906      /usr/lib/x86_64-linux-
gnu/libz.so.1.2.11
7f275b94f000-7f275b955000 r--p 00013000 103:05 915906      /usr/lib/x86_64-linux-
gnu/libz.so.1.2.11
7f275b955000-7f275b956000 ---p 00019000 103:05 915906      /usr/lib/x86_64-linux-
gnu/libz.so.1.2.11
7f275b956000-7f275b957000 r--p 00019000 103:05 915906      /usr/lib/x86_64-linux-
gnu/libz.so.1.2.11
7f275b957000-7f275b958000 rw-p 0001a000 103:05 915906      /usr/lib/x86_64-linux-
gnu/libz.so.1.2.11
7f275b958000-7f275b95c000 r--p 00000000 103:05 923645      /usr/lib/x86_64-linux-
gnu/libexpat.so.1.6.11
7f275b95c000-7f275b978000 r-xp 00004000 103:05 923645      /usr/lib/x86_64-linux-
gnu/libexpat.so.1.6.11
7f275b978000-7f275b982000 r--p 00020000 103:05 923645      /usr/lib/x86_64-linux-
gnu/libexpat.so.1.6.11
7f275b982000-7f275b983000 ---p 0002a000 103:05 923645      /usr/lib/x86_64-linux-
gnu/libexpat.so.1.6.11
7f275b983000-7f275b985000 r--p 0002a000 103:05 923645      /usr/lib/x86_64-linux-
gnu/libexpat.so.1.6.11
7f275b985000-7f275b986000 rw-p 0002c000 103:05 923645      /usr/lib/x86_64-linux-
gnu/libexpat.so.1.6.11
7f275b986000-7f275b988000 r--p 00000000 103:05 924057      /usr/lib/x86_64-linux-
gnu/libltdl.so.7.3.1
7f275b988000-7f275b98d000 r-xp 00002000 103:05 924057      /usr/lib/x86_64-linux-
gnu/libltdl.so.7.3.1
7f275b98d000-7f275b98f000 r--p 00007000 103:05 924057      /usr/lib/x86_64-linux-
gnu/libltdl.so.7.3.1
7f275b98f000-7f275b990000 r--p 00008000 103:05 924057      /usr/lib/x86_64-linux-
gnu/libltdl.so.7.3.1
7f275b990000-7f275b991000 rw-p 00009000 103:05 924057      /usr/lib/x86_64-linux-
gnu/libltdl.so.7.3.1
7f275b991000-7f275b995000 r--p 00000000 103:05 921934      /usr/lib/x86_64-linux-
gnu/libtdb.so.1.4.3
7f275b995000-7f275b9a3000 r-xp 00004000 103:05 921934      /usr/lib/x86_64-linux-
gnu/libtdb.so.1.4.3
7f275b9a3000-7f275b9a9000 r--p 00012000 103:05 921934      /usr/lib/x86_64-linux-
gnu/libtdb.so.1.4.3
7f275b9a9000-7f275b9aa000 r--p 00017000 103:05 921934      /usr/lib/x86_64-linux-
gnu/libtdb.so.1.4.3
7f275b9aa000-7f275b9ab000 rw-p 00018000 103:05 921934      /usr/lib/x86_64-linux-
gnu/libtdb.so.1.4.3
7f275b9ab000-7f275b9ad000 rw-p 00000000 00:00 0
7f275b9ad000-7f275b9af000 r--p 00000000 103:05 924631      /usr/lib/x86_64-linux-
gnu/libvorbisfile.so.3.3.7
7f275b9af000-7f275b9b4000 r-xp 00002000 103:05 924631      /usr/lib/x86_64-linux-
gnu/libvorbisfile.so.3.3.7
7f275b9b4000-7f275b9b5000 r--p 00007000 103:05 924631      /usr/lib/x86_64-linux-
gnu/libvorbisfile.so.3.3.7
7f275b9b5000-7f275b9b6000 ---p 00008000 103:05 924631      /usr/lib/x86_64-linux-
gnu/libvorbisfile.so.3.3.7
7f275b9b6000-7f275b9b7000 r--p 00008000 103:05 924631      /usr/lib/x86_64-linux-
gnu/libvorbisfile.so.3.3.7
7f275b9b7000-7f275b9b8000 rw-p 00009000 103:05 924631      /usr/lib/x86_64-linux-
gnu/libvorbisfile.so.3.3.7
7f275b9b8000-7f275b9ba000 r--p 00000000 103:05 924277      /usr/lib/x86_64-linux-
gnu/libpcre2-8.so.0.9.0
7f275b9ba000-7f275ba1e000 r-xp 00002000 103:05 924277      /usr/lib/x86_64-linux-
gnu/libpcre2-8.so.0.9.0

```

7f275ba1e000-7f275ba46000 r--p 00066000 103:05 924277	/usr/lib/x86_64-linux-
gnu/libpcre2-8.so.0.9.0	
7f275ba46000-7f275ba47000 r--p 0008d000 103:05 924277	/usr/lib/x86_64-linux-
gnu/libpcre2-8.so.0.9.0	
7f275ba47000-7f275ba48000 rw-p 0008e000 103:05 924277	/usr/lib/x86_64-linux-
gnu/libpcre2-8.so.0.9.0	
7f275ba48000-7f275ba6d000 r--p 00000000 103:05 917893	/usr/lib/x86_64-linux-
gnu/libc-2.31.so	
7f275ba6d000-7f275bbe5000 r-xp 00025000 103:05 917893	/usr/lib/x86_64-linux-
gnu/libc-2.31.so	
7f275bbe5000-7f275bc2f000 r--p 0019d000 103:05 917893	/usr/lib/x86_64-linux-
gnu/libc-2.31.so	
7f275bc2f000-7f275bc30000 ---p 001e7000 103:05 917893	/usr/lib/x86_64-linux-
gnu/libc-2.31.so	
7f275bc30000-7f275bc33000 r--p 001e7000 103:05 917893	/usr/lib/x86_64-linux-
gnu/libc-2.31.so	
7f275bc33000-7f275bc36000 rw-p 001ea000 103:05 917893	/usr/lib/x86_64-linux-
gnu/libc-2.31.so	
7f275bc36000-7f275bc3a000 rw-p 00000000 00:00 0	
7f275bc3a000-7f275bc41000 r--p 00000000 103:05 917906	/usr/lib/x86_64-linux-
gnu/libpthread-2.31.so	
7f275bc41000-7f275bc52000 r-xp 00007000 103:05 917906	/usr/lib/x86_64-linux-
gnu/libpthread-2.31.so	
7f275bc52000-7f275bc57000 r--p 00018000 103:05 917906	/usr/lib/x86_64-linux-
gnu/libpthread-2.31.so	
7f275bc57000-7f275bc58000 r--p 0001c000 103:05 917906	/usr/lib/x86_64-linux-
gnu/libpthread-2.31.so	
7f275bc58000-7f275bc59000 rw-p 0001d000 103:05 917906	/usr/lib/x86_64-linux-
gnu/libpthread-2.31.so	
7f275bc59000-7f275bc5d000 rw-p 00000000 00:00 0	
7f275bc5d000-7f275bcce000 r--p 00000000 103:05 917016	/usr/lib/x86_64-linux-
gnu/libpython3.8.so.1.0	
7f275bcce000-7f275bf29000 r-xp 00071000 103:05 917016	/usr/lib/x86_64-linux-
gnu/libpython3.8.so.1.0	
7f275bf29000-7f275c142000 r--p 002cc000 103:05 917016	/usr/lib/x86_64-linux-
gnu/libpython3.8.so.1.0	
7f275c142000-7f275c143000 ---p 004e5000 103:05 917016	/usr/lib/x86_64-linux-
gnu/libpython3.8.so.1.0	
7f275c143000-7f275c149000 r--p 004e5000 103:05 917016	/usr/lib/x86_64-linux-
gnu/libpython3.8.so.1.0	
7f275c149000-7f275c190000 rw-p 004eb000 103:05 917016	/usr/lib/x86_64-linux-
gnu/libpython3.8.so.1.0	
7f275c190000-7f275c1b3000 rw-p 00000000 00:00 0	
7f275c1b3000-7f275c1b4000 r--p 00000000 103:05 917894	/usr/lib/x86_64-linux-
gnu/libdl-2.31.so	
7f275c1b4000-7f275c1b6000 r-xp 00001000 103:05 917894	/usr/lib/x86_64-linux-
gnu/libdl-2.31.so	
7f275c1b6000-7f275c1b7000 r--p 00003000 103:05 917894	/usr/lib/x86_64-linux-
gnu/libdl-2.31.so	
7f275c1b7000-7f275c1b8000 r--p 00003000 103:05 917894	/usr/lib/x86_64-linux-
gnu/libdl-2.31.so	
7f275c1b8000-7f275c1b9000 rw-p 00004000 103:05 917894	/usr/lib/x86_64-linux-
gnu/libdl-2.31.so	
7f275c1b9000-7f275c1bb000 rw-p 00000000 00:00 0	
7f275c1bb000-7f275c1c0000 r-xp 00000000 103:05 923815	/usr/lib/x86_64-linux-
gnu/libgpm.so.2	
7f275c1c0000-7f275c3bf000 ---p 00005000 103:05 923815	/usr/lib/x86_64-linux-
gnu/libgpm.so.2	
7f275c3bf000-7f275c3c0000 r--p 00004000 103:05 923815	/usr/lib/x86_64-linux-
gnu/libgpm.so.2	
7f275c3c0000-7f275c3c1000 rw-p 00005000 103:05 923815	/usr/lib/x86_64-linux-
gnu/libgpm.so.2	

```

7f275c3c1000-7f275c3c3000 r--p 00000000 103:05 923315          /usr/lib/x86_64-linux-
gnu/libacl.so.1.1.2253
7f275c3c3000-7f275c3c8000 r-xp 00002000 103:05 923315          /usr/lib/x86_64-linux-
gnu/libacl.so.1.1.2253
7f275c3c8000-7f275c3ca000 r--p 00007000 103:05 923315          /usr/lib/x86_64-linux-
gnu/libacl.so.1.1.2253
7f275c3ca000-7f275c3cb000 r--p 00008000 103:05 923315          /usr/lib/x86_64-linux-
gnu/libacl.so.1.1.2253
7f275c3cb000-7f275c3cc000 rw-p 00009000 103:05 923315          /usr/lib/x86_64-linux-
gnu/libacl.so.1.1.2253
7f275c3cc000-7f275c3cf000 r--p 00000000 103:05 923446          /usr/lib/x86_64-linux-
gnu/libcanberra.so.0.2.5
7f275c3cf000-7f275c3d9000 r-xp 00003000 103:05 923446          /usr/lib/x86_64-linux-
gnu/libcanberra.so.0.2.5
7f275c3d9000-7f275c3dd000 r--p 0000d000 103:05 923446          /usr/lib/x86_64-linux-
gnu/libcanberra.so.0.2.5
7f275c3dd000-7f275c3de000 r--p 00010000 103:05 923446          /usr/lib/x86_64-linux-
gnu/libcanberra.so.0.2.5
7f275c3de000-7f275c3df000 rw-p 00011000 103:05 923446          /usr/lib/x86_64-linux-
gnu/libcanberra.so.0.2.5
7f275c3df000-7f275c3e5000 r--p 00000000 103:05 924431          /usr/lib/x86_64-linux-
gnu/libselinux.so.1
7f275c3e5000-7f275c3fe000 r-xp 00006000 103:05 924431          /usr/lib/x86_64-linux-
gnu/libselinux.so.1
7f275c3fe000-7f275c405000 r--p 0001f000 103:05 924431          /usr/lib/x86_64-linux-
gnu/libselinux.so.1
7f275c405000-7f275c406000 ---p 00026000 103:05 924431          /usr/lib/x86_64-linux-
gnu/libselinux.so.1
7f275c406000-7f275c407000 r--p 00026000 103:05 924431          /usr/lib/x86_64-linux-
gnu/libselinux.so.1
7f275c407000-7f275c408000 rw-p 00027000 103:05 924431          /usr/lib/x86_64-linux-
gnu/libselinux.so.1
7f275c408000-7f275c40a000 rw-p 00000000 00:00 0
7f275c40a000-7f275c418000 r--p 00000000 103:05 924540          /usr/lib/x86_64-linux-
gnu/libtinfo.so.6.2
7f275c418000-7f275c427000 r-xp 0000e000 103:05 924540          /usr/lib/x86_64-linux-
gnu/libtinfo.so.6.2
7f275c427000-7f275c435000 r--p 0001d000 103:05 924540          /usr/lib/x86_64-linux-
gnu/libtinfo.so.6.2
7f275c435000-7f275c439000 r--p 0002a000 103:05 924540          /usr/lib/x86_64-linux-
gnu/libtinfo.so.6.2
7f275c439000-7f275c43a000 rw-p 0002e000 103:05 924540          /usr/lib/x86_64-linux-
gnu/libtinfo.so.6.2
7f275c43a000-7f275c449000 r--p 00000000 103:05 917895          /usr/lib/x86_64-linux-
gnu/libm-2.31.so
7f275c449000-7f275c4f0000 r-xp 0000f000 103:05 917895          /usr/lib/x86_64-linux-
gnu/libm-2.31.so
7f275c4f0000-7f275c587000 r--p 000b6000 103:05 917895          /usr/lib/x86_64-linux-
gnu/libm-2.31.so
7f275c587000-7f275c588000 r--p 0014c000 103:05 917895          /usr/lib/x86_64-linux-
gnu/libm-2.31.so
7f275c588000-7f275c589000 rw-p 0014d000 103:05 917895          /usr/lib/x86_64-linux-
gnu/libm-2.31.so
7f275c589000-7f275c58b000 rw-p 00000000 00:00 0
7f275c5ae000-7f275c5af000 r--p 00000000 103:05 917889          /usr/lib/x86_64-linux-gnu/ld-
2.31.so
7f275c5af000-7f275c5d2000 r-xp 00001000 103:05 917889          /usr/lib/x86_64-linux-gnu/ld-
2.31.so
7f275c5d2000-7f275c5da000 r--p 00024000 103:05 917889          /usr/lib/x86_64-linux-gnu/ld-
2.31.so
7f275c5db000-7f275c5dc000 r--p 0002c000 103:05 917889          /usr/lib/x86_64-linux-gnu/ld-
2.31.so

```

```
7f275c5dc000-7f275c5dd000 rw-p 0002d000 103:05 917889          /usr/lib/x86_64-linux-gnu/ld-
2.31.so
7f275c5dd000-7f275c5de000 rw-p 00000000 00:00 0
7ffd22d2f000-7ffd22d50000 rw-p 00000000 00:00 0          [stack]
7ffd22db0000-7ffd22db4000 r--p 00000000 00:00 0          [vvar]
7ffd22db4000-7ffd22db6000 r-xp 00000000 00:00 0          [vdso]
ffffffff600000-ffffffff601000 --xp 00000000 00:00 0          [vsyscall]
```

**True      False**

- The size of the heap is one page
- This is a virtual memory map (not physical memory map)
- vim.basic uses the math library
- The 5th entry 55a4352c5000-55a4352e2000 may correspond to "data" of the vim.basic
- The size of the stack is one page

The size of the heap is one page: False

This is a virtual memory map (not physical memory map): True

vim.basic uses the math library: True

The 5th entry 55a4352c5000-55a4352e2000 may correspond to "data" of the vim.basic: True

The size of the stack is one page: False

**Question 10**

Complete

Mark 0.00 out of 1.00

Select all the correct statements, w.r.t. Copy on Write

- a. Fork() used COW technique to improve performance of new process creation.
- b. Vfork() assumes that there will be no write, but rather exec()
- c. COW helps us save memory
- d. If either parent or child modifies a COW-page, then a copy of the page is made and page table entry is updated
- e. use of COW during fork() is useless if child called exit()
- f. use of COW during fork() is useless if exec() is called by the child

The correct answers are: Fork() used COW technique to improve performance of new process creation., If either parent or child modifies a COW-page, then a copy of the page is made and page table entry is updated, COW helps us save memory, Vfork() assumes that there will be no write, but rather exec()

**Question 11**

Complete

Mark 1.00 out of 1.00

Assuming a 8- KB page size, what is the page numbers for the address 1093943 reference in decimal :  
 (give answer also in decimal)

Answer:

The correct answer is: 134

**Question 12**

Complete

Mark 0.00 out of 1.00

Order the following events, related to page fault handling, in correct order

1. Page fault handler detects that it's a page fault and not illegal memory access
2. Disk interrupt handler runs
3. MMU detects that a page table entry is marked "invalid"
4. Page faulted process is moved to ready-queue
5. Empty frame is found
6. Page fault interrupt is generated
7. Other processes scheduled by scheduler
8. Page table of page faulted process is updated
9. Disk Interrupt occurs
10. Disk read is issued
11. Page fault handler in kernel starts executing
12. Page faulting process is made to wait in a queue

The correct order for these items is as follows:

1. MMU detects that a page table entry is marked "invalid"
2. Page fault interrupt is generated
3. Page fault handler in kernel starts executing
4. Page fault handler detects that it's a page fault and not illegal memory access
5. Empty frame is found
6. Disk read is issued
7. Page faulting process is made to wait in a queue
8. Other processes scheduled by scheduler
9. Disk Interrupt occurs
10. Disk interrupt handler runs
11. Page table of page faulted process is updated
12. Page faulted process is moved to ready-queue

**Question 13**

Complete

Mark 1.00 out of 1.00

Page sizes are a power of 2 because

Select one:

- a. Certain bits are reserved for offset in logical address. Hence page size =  $2^{(\text{no.of offset bits})}$
- b. Power of 2 calculations are highly efficient
- c. Certain bits are reserved for offset in logical address. Hence page size =  $2^{(32 - \text{no.of offset bits})}$
- d. operating system calculations happen using power of 2
- e. MMU only understands numbers that are power of 2

The correct answer is: Certain bits are reserved for offset in logical address. Hence page size =  $2^{(\text{no.of offset bits})}$

**Question 14**

Complete

Mark 1.00 out of 1.00

Given below is the output of the command "ps -eo min\_flt,maj\_flt,cmd" on a Linux Desktop system. Select the statements that are consistent with the output

```
626729 482768 /usr/lib/firefox/firefox -contentproc -parentBuildID 20220202182137 -prefsLen 9256 -prefMapSize 264738 -appDir /usr/lib/firefox/browser 6094 true rdd
2167 687 /usr/sbin/apache2 -k start
1265185 222 /usr/bin/gnome-shell
102648 111 /usr/sbin/mysqld
9813 0 bash
15497 370 /usr/bin/gedit --gapplication-service
```

- a. All of the processes here exhibit some good locality of reference
- b. Firefox has likely been running for a large amount of time
- c. The bash shell is mostly busy doing work within a particular locality
- d. Apache web-server has not been doing much work

The correct answers are: Firefox has likely been running for a large amount of time, Apache web-server has not been doing much work, The bash shell is mostly busy doing work within a particular locality, All of the processes here exhibit some good locality of reference

**Question 15**

Complete

Mark 0.14 out of 1.00

Suppose two processes share a library between them. The library consists of 5 pages, and these 5 pages are mapped to frames 9, 15, 23, 4, 7 respectively. Process P1 has got 6 pages, first 3 of which consist of process's own code/data and 3 correspond to library's pages 0, 2, 4. Process P2 has got 7 pages, first 3 of which consist of process's own code/data and remaining 4 correspond to library's pages 0, 1, 3, 4. Fill in the blanks for page table entries of P1 and P2.

Page table of P1, Page 5	<input type="text" value="23"/>
Page table of P1, Page 3	<input type="text" value="9"/>
Page table of P2, Page 0	<input type="text" value="6"/>
Page table of P2, Page 3	<input type="text" value="7"/>
Page table of P2, Page 4	<input type="text" value="9"/>
Page table of P2, Page 1	<input type="text" value="5"/>
Page table of P1, Page 4	<input type="text" value="9"/>

The correct answer is: Page table of P1, Page 5 → 7, Page table of P1, Page 3 → 9, Page table of P2, Page 0 → 9, Page table of P2, Page 3 → 4, Page table of P2, Page 4 → 7, Page table of P2, Page 1 → 15, Page table of P1, Page 4 → 23

**Question 16**

Complete

Mark 0.50 out of 1.00

For the reference string

3 4 3 5 2

the number of page faults (including initial ones) using

FIFO replacement and 2 page frames is :

FIFO replacement and 3 page frames is :

◀ (Code) mmap related programs

Jump to...

Points from Mid-term feedback ►

**Started on** Monday, 7 March 2022, 7:00:12 PM

**State** Finished

**Completed on** Monday, 7 March 2022, 8:00:04 PM

**Time taken** 59 mins 52 secs

**Grade** 9.78 out of 15.00 (65%)

**Question 1**

Complete

Mark 1.00 out of 1.00

Why is there a call to kinit2? Why is it not merged with knit1?

- a. call to seginit() makes it possible to actually use PHYSTOP in argument to kinit2()
- b. When kinit1() is called there is a need for few page frames, but later kinit2() is called to serve need of more page frames
- c. Because there is a limit on the values that the arguments to kinit1() can take.
- d. kinit2 refers to virtual addresses beyond 4MB, which are not mapped before kalloc() is called

The correct answer is: kinit2 refers to virtual addresses beyond 4MB, which are not mapped before kalloc() is called

**Question 2**

Complete

Mark 1.50 out of 1.50

Arrange the following in the correct order of execution (w.r.t. 'init')

initcode() returns in forkret()

6

'initcode' process is marked RUNNABLE

3

'initcode' struct proc is created

2

userinit() is called

1

mpmain() calls scheduler()

4

initcode() calls exec("/init", ...)

8

initcode() returns from trapret()

7

scheduler() schedules initcode() process

5

The correct answer is: initcode() returns in forkret() → 6, 'initcode' process is marked RUNNABLE → 3, 'initcode' struct proc is created → 2, userinit() is called → 1, mpmain() calls scheduler() → 4, initcode() calls exec("/init", ...) → 8, initcode() returns from trapret() → 7, scheduler() schedules initcode() process → 5

**Question 3**

Complete

Mark 0.00 out of 2.00

exec() does this: curproc->tf->eip = elf.entry, but userinit() does this: p->tf->eip = 0; Select all the statements from below, that collectively explain this

- a. the 'entry' in initcode is anyways 0
- b. exec() loads from ELF file and the address of first instruction to be executed is given by 'entry'
- c. the initcode is created using objcopy, which discards all relocation information and symbols (like entry)
- d. elf.entry is anyways 0, so both statements mean the same
- e. In userinit() the function inituvm() has mapped the code of 'initcode' to be starting at virtual address 0
- f. the code of 'initcode' is loaded at physical address 0

The correct answers are: exec() loads from ELF file and the address of first instruction to be executed is given by 'entry', In userinit() the function inituvm() has mapped the code of 'initcode' to be starting at virtual address 0, the initcode is created using objcopy, which discards all relocation information and symbols (like entry)

**Question 4**

Complete

Mark 1.00 out of 1.00

What does userinit() do ?

- a. initializes the users
- b. sets up the 'initcode' process to start execution in forkret()
- c. sets up the 'initcode' process to start execution in forkret()
- d. sets up the 'initcode' process to start execution in trapret()
- e. sets up the 'init' process to start execution in forkret()
- f. initializes the process 'init' and starts executing it

The correct answer is: sets up the 'initcode' process to start execution in forkret()

**Question 5**

Complete

Mark 0.67 out of 1.00

Which of the following is done by mappages()?

- a. allocate page directory if required
- b. allocate page frame if required
- c. allocate page table if required
- d. create page table mappings for the range given by "va" and "va + size"
- e. create page table mappings to the range given by "pa" and "pa + size"

The correct answers are: create page table mappings for the range given by "va" and "va + size", allocate page table if required, create page table mappings to the range given by "pa" and "pa + size"

**Question 6**

Complete

Mark 0.00 out of 1.00

Select the statement that most correctly describes what setupkvm() does

- a. creates a 2-level page table setup with virtual->physical mappings specified in the kmap[] global array
- b. creates a 2-level page table setup with virtual->physical mappings specified in the kmap[] global array and makes kpgdir point to it
- c. creates a 1-level page table for the use by the kernel, as specified in kmap[] global array
- d. creates a 2-level page table for the use of the kernel, as specified in gdtdesc

The correct answer is: creates a 2-level page table setup with virtual->physical mappings specified in the kmap[] global array

**Question 7**

Complete

Mark 0.00 out of 1.00

The approximate number of page frames created by kinit1 is

- a. 4
- b. 16
- c. 4000
- d. 3000
- e. 2000
- f. 1000
- g. 10

The correct answer is: 3000

**Question 8**

Complete

Mark 1.20 out of 1.50

Which of the following is DONE by allocproc() ?

- a. setup kernel memory mappings for the process
- b. Select an UNUSED struct proc for use
- c. ensure that the process starts in trapret()
- d. allocate kernel stack for the process
- e. allocate PID to the process
- f. ensure that the process starts in forkret()
- g. setup the trapframe and context pointers appropriately
- h. setup the contents of the trapframe of the process properly

The correct answers are: Select an UNUSED struct proc for use, allocate PID to the process, allocate kernel stack for the process, setup the trapframe and context pointers appropriately, ensure that the process starts in forkret()

**Question 9**

Complete

Mark 1.00 out of 1.00

Map the virtual address to physical address in xv6

KERNLINK	0x100000
0xFE000000	0xFE000000
80108000	0x108000
KERNBASE	0

The correct answer is: KERNLINK → 0x100000, 0xFE000000 → 0xFE000000, 80108000 → 0x108000, KERNBASE → 0

**Question 10**

Complete

Mark 0.42 out of 1.00

Select all the correct statements about initcode

- a. code of initcode is loaded at virtual address 0
- b. The data and stack of initcode is mapped to one single page in userinit()
- c. code of 'initcode' is loaded along with the kernel during booting
- d. the size of 'initcode' is 2c
- e. code of initcode is loaded in memory by the kernel during userinit()
- f. initcode is the 'init' process
- g. initcode essentially calls exec("/init",...)

The correct answers are: code of 'initcode' is loaded along with the kernel during booting, the size of 'initcode' is 2c, The data and stack of initcode is mapped to one single page in userinit(), initcode essentially calls exec("/init",...)

**Question 11**

Complete

Mark 1.00 out of 1.00

The variable 'end' used as argument to kinit1 has the value

- a. 801154a8
- b. 81000000
- c. 80110000
- d. 80102da0
- e. 8010a48c
- f. 80000000

The correct answer is: 801154a8

**Question 12**

Complete

Mark 1.00 out of 1.00

What does seginit() do?

- a. Nothing significant, just repetition of earlier GDT setup but with 2-level paging setup done
- b. Nothing significant, just repetition of earlier GDT setup but with free frames list created now
- c. Adds two additional entries to GDT corresponding to Code and Data segments, but to be used in privilege level 0
- d. Adds two additional entries to GDT corresponding to Code and Data segments, but to be used in privilege level 3
- e. Nothing significant, just repetition of earlier GDT setup but with kernel page table allocated now

The correct answer is: Adds two additional entries to GDT corresponding to Code and Data segments, but to be used in privilege level 3

**Question 13**

Complete

Mark 1.00 out of 1.00

Does exec() code around clearptau() lead to wastage of one page frame?

- a. no
- b. yes

The correct answer is: yes

[◀ Questions for test on kalloc/kfree/kvmalloc, etc.](#)

Jump to...

[\(Optional Assignment\) Slab allocator in xv6 ►](#)