

Introduction to Computer Science

PTI CS 103 Lecture Notes

The Prison Teaching Initiative

Authors

TK.

Acknowledgements

TK.

Contents

1	Introduction	1
1.1	Overview of course	1
1.2	Brief history of computer science	2
1.3	Components of a computer	2
1.4	Types of computers	5
1.5	Why computers are useful	6
2	Hello World	8
2.1	Starting DrJava	8
2.2	Writing Hello World	8
2.3	Saving and Compiling in DrJava	9
2.4	Running Code with DrJava	9
2.5	Making Changes	10
3	Hardware	13
3.1	Input and Output Devices	14
3.2	Memory	15
3.3	Central Processing Unit	16
3.4	Conclusion	17
3.5	Learning Objectives	17
4	Numbering Systems	19
4.1	Decimal Numbers	19
4.2	Radix Decomposition	19
4.3	Binary Numbers	20
4.4	Binary Arithmetic	20
4.5	Binary to Decimal Conversion	21
4.6	Decimal to Binary Conversion	21
4.7	Conclusion	22
4.8	Learning Objectives	22
5	Programming Languages and Linux commands	23
5.1	Most popular programming languages	23
5.2	Overview of a GUI	24
6	Networks I: Protocols	28
6.1	Protocols	28
6.2	Network Layers	30
6.3	Network Security	36

7	Networks II: The Internet	39
7.1	Local Area Networks	39
7.2	Wide Area Networks	41
8	Control Structures	43
8.1	The <code>if</code> statement	43
8.2	The <code>else</code> statement	44
8.3	Nested conditionals	45
8.4	The <code>else if</code> statement	46
8.5	Curly braces	48
8.6	Common mistakes	49
9	Business Software	52
9.1	Introduction	52
9.2	Word Processing	52
9.3	Spreadsheets	53
9.4	Presentations	53
9.5	Project	53
10	Methods	54
10.1	Methods	55
10.2	Abstraction	62
10.3	Modularity	63
11	Arrays	66
11.1	Creating arrays	67
11.2	Indexing	67
11.3	Array length	68
11.4	Default initialization	68
11.5	Bounds checking	68
11.6	Empty arrays	69
11.7	Enhanced for loop	70
11.8	Exchanging and shuffling	70
12	ArrayLists	73
12.1	Creating an <code>ArrayList</code>	73
12.2	<code>add()</code>	73
12.3	<code>get()</code>	74
12.4	<code>contains()</code>	74
12.5	<code>ArrayLists</code> with custom classes	75
12.6	Arrays versus <code>ArrayLists</code>	75
13	Systems Development	77
13.1	Introduction	77
13.2	Code Organization	77
13.3	User-oriented Design	80
13.4	Testing	81
	Appendices	83

Introduction

1.1 Overview of course

Knowing just a little bit of computer science can get you started right away in writing code for real-world applications. In this course, you will learn about the fascinating subject of computer science. You will develop algorithmic thinking skills that will improve day-to-day critical thinking and problem-solving. Perhaps most importantly, you will learn to code, which will not only open up new job opportunities but also make you more effective in most areas of business.

[AFM: Is this still true?] In the first semester, the first two classes of each week focus on computer science theory and special topics. The final day of each week will be a lab day, where we actually start practicing coding skills.

In the second semester, we will focus more on practical coding, with a single day a week for theory and 2 lab periods per week for coding.

Broadly, we will cover the following topics:

- How modern computers work
 - Hardware
 - Software
 - Computer networks and information systems
- Algorithms for quickly solving complex problems
 - Searching
 - Sorting
- Data structures
 - Arrays
 - ArrayLists
- Applications of Computer science
 - Basic coding in Java
 - How to use productivity software

Exercise 1.1.0: Get a piece of graph paper, and build a maze on it. Your maze should have one input, and one output. An example maze is shown on Figure 1.1.

Next, get another student's maze. You are programming a robot that can solve this maze. These robots can take a sequence of commands, where each command is "Turn Left [number]" (TL[n]) and "Move Forwards [number]" (MF[n]). Create a sequence of inputs to the robots that would solve that maze. For example, the following robot completes the provided maze: MF[3], TL[3], MF[6], TL[1], MF[4], TL[3], MF[3], TL[1], MF[4].

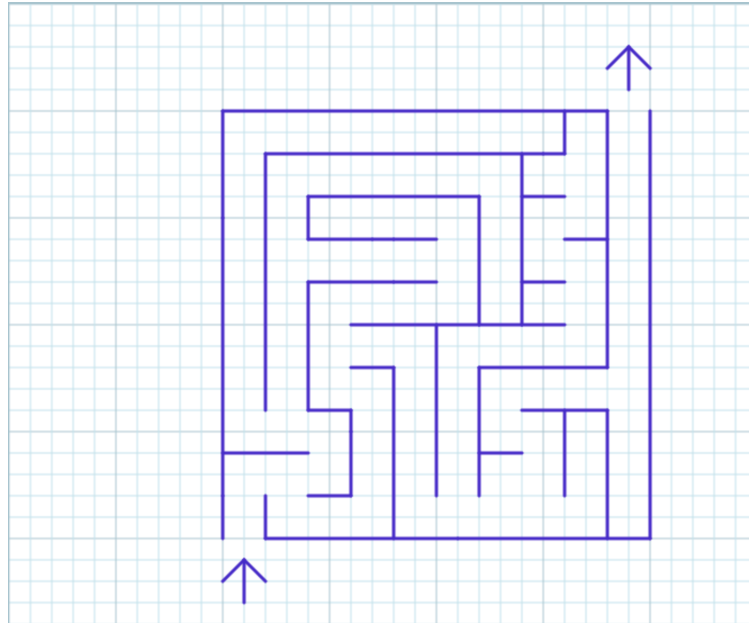


Figure 1.1: An example maze. The input is in the lower left, the output is in the upper right.

What types of commands could have made this easier? Could your robot solve arbitrary mazes? What types of commands would be helpful for making your robot solve any maze?

1.2 Brief history of computer science

Timeline:¹

- Invention of the abacus (2700-2300 BC, Sumerians)
- Design of first modern-style computer (Charles Babbage, 1837)
- Design of first computer algorithm (Ada Lovelace, 1843)
- Invention of first electronic digital computer (Konrad Zuse, 1941)
- Invention of the transistor (Bell labs, 1947)
- Invention of the first computer network (early Internet) (DARPA, 1968)
- Invention of the World Wide Web (Sir Tim Berners-Lee, 1990)

1.3 Components of a computer

There are 4 primary components of a computer:

- Hardware
- Software
- Data
- User

¹https://www.worldsciencefestival.com/infographics/a_history_of_computer_science/

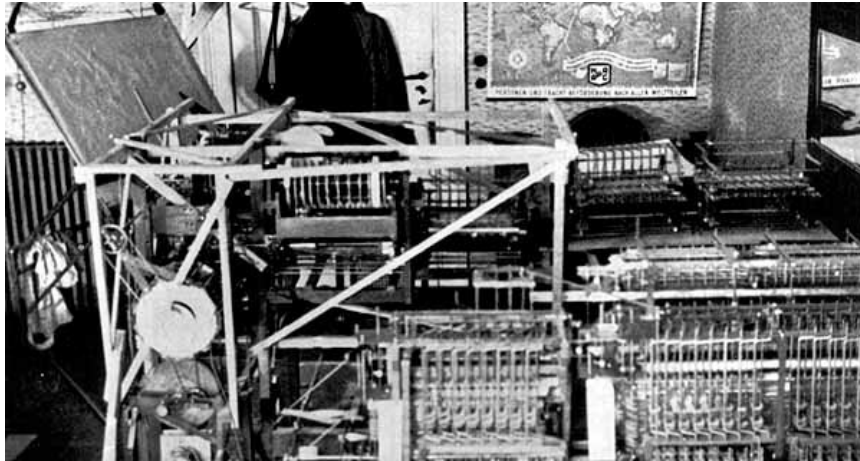


Figure 1.2: Construction of Konrad Zuse's Z1, the first modern computer, in his parents' apartment. <https://history-computer.com/ModernComputer/Relays/Zuse.html>

1.3.1 Hardware

Computer hardware consists of physical, electronic devices. These are the parts you actually can see and touch. Some examples include

- Central processing unit (CPU)
- Monitor
- CD drive
- Keyboard
- Computer data storage
- Graphic card
- Sound card
- Speakers
- Motherboard

We will discuss these components in more detail in lecture 3.

1.3.2 Software

Software (otherwise known as *programs*, *applications*, or *apps*) are organized sets of instructions for controlling the computer.

There are two main classes of software:

- Applications software: programs allowing the human to interact directly with the computer
- Systems software: programs the computer uses to control itself

Some more familiar applications software include

- Microsoft Word: allows the user to edit text files
- Internet Explorer: connects the user to the world wide web
- iTunes: organizes and plays music files

While applications software allows the user to interact with the computer, systems software keeps the computer running. The operating system (OS) is the most common example of systems software, and it schedules tasks and manages storage of data.

We will dive deeper into the details of both applications and systems software in lecture 4.



Figure 1.3: Examples of hardware components of a personal computer.
<https://www5.cob.ilstu.edu/dsmath1/tag/computer-hardware/>

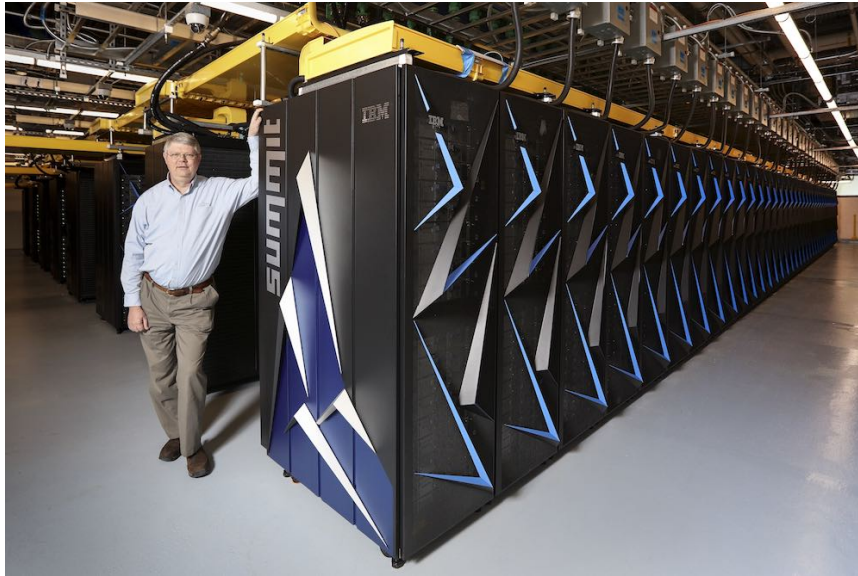


Figure 1.4: Summit, a world-class supercomputing cluster at Oak Ridge National Laboratory in Tennessee. <https://insidehpc.com/2018/11/new-top500-list-lead-doe-supercomputers/>

1.3.3 Data

Data is fundamentally information of any kind. One key benefit of computers is their ability to reliably store massive quantities of data for a long time. Another is the speed with which they can do calculations on data once they receive instructions from a human user.

While humans can understand data with a wide variety of perceptions (taste, smell, hearing, touch, sight), computers read and write everything internally as *bits*, sequences of 0s and 1s.

Computers have software and hardware which allow them to convert their internal 0s and 1s into text, numerals, and images displayed on the monitor; and sounds which can be played through the speaker.

Similarly, computers have hardware and software that convert information from the real world into bits: a microphone converts sound, a camera converts pictures, and a text editor converts character symbols.

1.3.4 Users

Of course, there would be no data and no meaningful calculations without the human user. Computers are ultimately tools for making humans more powerful.

As we will see in the next section, however, different types of computers have different roles for the user.

1.4 Types of computers

1.4.1 Supercomputers

These are the most powerful computers out there. They are used for problems that take along time to calculate. They are rare because of their size and expense, and therefore primarily used by big organizations like universities, corporations, and the government.

The user of a supercomputer typically gives the computer a list of instructions, and allows the supercomputer to run on its own over the course of hours or days to complete its task.

1.4.2 Mainframe computers

Although not as powerful as supercomputers, mainframe computers can handle more data and run much faster than a typical personal computer. Often, they are given instructions only periodically by computer programmers, and then run on their own for months at a time to store and process incoming data. For example, census number-crunching, consumer statistics, and transactions processing all use mainframe computers

1.4.3 Personal computers

These are the familiar computers we use to interact with applications every day. Full-size desktop computers and laptop computers are examples

1.4.4 Embedded computers

In the modern “digital” age, nearly all devices we use have computers embedded within them. From cars to washing machines to watches to heating systems, most everyday appliances have a computer within them that allows them to function.

1.4.5 Mobile computers

In the past 2 decades, mobile devices have exploded onto the scene, and smartphones have essentially become as capable as standalone personal computers for many tasks.

1.5 Why computers are useful

Computers help us in most tasks in the modern age. We can use them, for example, to:

- Write a letter.
- Do our taxes.
- Play video games.
- Watch videos.
- Surf the internet.
- Keep in touch with friends.
- Date.
- Order food.
- Control robots and self-driving cars.

Exercise 1.5.0: What are some other tasks a computer can accomplish?

This is why the job market for computer scientists continues to expand, and why computer skills are more and more necessary even in non-computational jobs.

According to a Stackoverflow survey from 2018,² 9% of professional coders on the online developer community have only been coding for 0-2 years. This demonstrates two things:

1. The job market for people with coding skills is continually expanding
2. It doesn't take much to become a coder

Some examples of careers in computer science include

- IT management / consulting

²<https://insights.stackoverflow.com/survey/2018/>

- Game developer
- Web developer
- UI/UX designer
- Data analyst
- Database manager

Exercise 1.5.0: Think about building a simple video game. What are the hardware, software, data needed for this game? What would be the minimal hardware, and what hardware would you need to add features like voice chat or online play. What types of computers would best run this game? Why?

Do the same for a text editor.

References

Computer Science: An Interdisciplinary Approach, Robert Sedgewick and Kevin Wayne.

University of Wisconsin-Madison CS 202 Lectures, Andrea Arpaci-Dusseau.
<http://pages.cs.wisc.edu/~dusseau/Classes/CS202-F11/>

Hello World

It is now time for us to start practicing coding skills. In this course, we will be focusing on learning how to program in Java, which is just one of many programming languages. While each programming languages has its own quirks and features specific to it, many of the concepts you will learn in this class will exist not only in Java but also in other languages.

We will now write our first program in Java called “Hello World.” We will use DrJava to help us write, compile, and run this program.

2.1 Starting DrJava

TODO: how to start DrJava on laptop that students are given TODO: add screenshot of DrJava when just opened and label parts

On the left-hand side of the window, there is a list of currently-open files. When you first start up DrJava, the only thing there should say “(Untitled)”. The main part of the window is where you will write code. There should be numbers along the left-hand side to make it easy for you to see which line you are writing code on. The bottom part of the window is where we can interact with the program, for example by seeing the output of compiling the program or by seeing the output of running the program.

2.2 Writing Hello World

In the main part of the window, type the following code:

```
1  class HelloWorld {
2
3      public static void main(String[] args) {
4          System.out.println("Hello World!");
5      }
6
7  }
```

We will learn more about what the different parts of this program do later on, but here is a brief explanation for now:

1. In line 1, `class HelloWorld` makes a new *class* called `HelloWorld`. In Java, all code must be enclosed in a class. The class contains all the code between the open curly brace { on line 1 and the close curly brace } on line 7.
2. Line 3 makes a new method called `main`. We will learn later what `public static void` and `String[] args` mean. For now, all you need to know is that every Java program needs a `main` method and starts by executing the code within the `main` method. Our `main` method contains all the code between the open curly brace { on line 3 and the close curly brace } on line 5.

3. Line 4 is a print statement. `System.out.println` will print out whatever is specified between the open parenthesis (and close parenthesis) that follows immediately after. `System.out.println` will begin a new line after printing out whatever is specified in the parentheses. In this case, we put the *string literal* `"Hello World!"` in between these parentheses. String literals refer to text between quotation marks.

2.3 Saving and Compiling in DrJava

In the upper part of the window, you should see a toolbar with many buttons. You can now click the “Save” button. Because the class name is `HelloWorld`, the name of your file will be `HelloWorld.java`. The `.java` file extension is used for *source code*, which in this case is human-readable code written in Java.

Clicking the “Compile” button will first save your file (if it hasn’t already been saved) and then *compile* it. Compiling converts your human-readable source code into code that is easier for the computer to read. In the case of Java, this will generate a `.class` file. The `HelloWorld.class` file contains the compiled code. After clicking “Compile,” the bottom part of the window should say **Compilation completed**. If there are problems with your program, it may not be able to be compiled. In such a case, this bottom part of the window will instead indicate that there was an issue with compilation. We can try seeing how this might look now.

Try deleting the final curly brace `}` from the program and hit “Compile.” The bottom window should now indicate that something has gone wrong. The message **Error: reached end of file while parsing** arises because the end of the file was reached without finding the final curly brace. Add the curly brace back in and hit “Compile” again to recompile the program. Note that in the bottom part of the window, the “Compiler Output” tab is selected. This tab is where you can view any messages from compilation.

Let’s see another example of a compilation error. Try deleting the semicolon `;` from the program and hit “Compile.” The message **Error: ';' expected** arises because the curly brace on line 5 was reached without the semicolon being encountered. In Java, every statement must end in a semicolon. In this case, the print statement on line 4 must end in a semicolon. Add the semicolon back in and recompile the program. Compilation should succeed.

2.4 Running Code with DrJava

Now that the program has been written and compiled, it is now time to run it. We can use DrJava to run a compiled program by clicking the “Run” button in the toolbar. In the bottom part of the window, you should see something like the following:

```
1 > run HelloWorld
2 Hello World!
3 >
```

The `run HelloWorld` part indicates that the program `HelloWorld` should be run, and the next line is the output of the program. Note that in the bottom part of the window, the “Interactions” tab is now selected. This tab allows you to interact with your program. Clicking the “Run” button is equivalent to clicking the “Interactions” tab, typing `run HelloWorld`, and then hitting enter/return key on your keyboard.

If you would like, you can type `run HelloWorld` after the final `>` and type the enter/return key to run your program again.

2.5 Making Changes

Now we’ve successfully written, compiled, and ran our first Java program. As an exercise, let’s make our “Hello World!” more enthusiastic. Change the `Hello World!` in line 4 to `Hello World!!`. You can compile and rerun your program to see the results. In general, you can make changes to your program, recompile it, and run it again to see the results of your changes.

Let’s make some other changes. Try adding `//print hello` to the end of line 4, so that your line 4 now looks like this:

```
1      System.out.println("Hello World!!"); // print hello
```

Try compiling and running the program. The behavior should remain unchanged from the previous time you ran the program. Two forward slashes `//` denote the beginning of an *inline comment* and have the effect of turning the rest of the line into a *comment*. Comments are ignored by the compiler and do not have any effect on the program’s behavior. They are useful for adding explanations or providing information to people that are reading the code. In this case, the comment `print hello` can be viewed as explaining the purpose of the code on that line.

Now try adding `//` to the beginning of line 4, so that it looks as follows:

```
1  //      System.out.println("Hello World!!"); //print hello
```

Try compiling and running the program now. You should find that there is now no output. The two forward slashes that you added turned the whole line into a comment. As far as the compiler is concerned, your program is the same as if the entire line 4 did not exist.

Uncomment line 4 by deleting the forward slashes you added at the beginning of it. Now add the following lines after line 4:

```
1  System.out.println("Hello!"); // print hello again
2  System.out.println("Goodbye World!"); // say goodbye
```

Your code now has *multiple statements*. In particular, you now have three. The first one on line 4 will print `Hello World!!`, the second one will print `Hello!`, and the third will print `Goodbye World!`. Note that each of these statements is followed by a semicolon. Deleting any one of these semicolons will result in a compilation error (give it a try!). Your code should now look as follows:

```
1  class HelloWorld {
2
3      public static void main(String[] args) {
4          System.out.println("Hello World!!"); //print hello
5          System.out.println("Hello!"); // print hello again
6          System.out.println("Goodbye World!"); // say goodbye
7      }
8
9  }
```

Try compiling it and running it. The output should look as follows:

```
1  Hello World!!
2  Hello!
3  Goodbye World!
```

Let’s add some description of what we’re doing by adding a *block comment*. A block comment begins with `/*` and is ended whenever `*/` is encountered. It may (but does

not necessarily need to) span multiple lines. Right before the main method, add a block comment to get the following:

```

1  class HelloWorld {
2
3      /* Main method to print out the following:
4         Hello World!!
5         Hello!
6         Goodbye World!
7     */
8     public static void main(String[] args) {
9         System.out.println("Hello World!!"); //print hello
10        System.out.println("Hello!"); // print hello again
11        System.out.println("Goodbye World!"); // say goodbye
12    }
13
14 }
```

If you compile and run the program now, it should have the same behavior as before, since whatever occurs inside comments does not affect the program's behavior.

Inline comments can be nested inside block comments, so for example, the following program will compile without issue:

```

1  class HelloWorld {
2
3      /* Main method to print out the following:
4         Hello World!!
5         Hello!
6         Goodbye World!
7     */
8     public static void main(String[] args) {
9         /*
10         System.out.println("Hello World!!"); //print hello
11         System.out.println("Hello!"); // print hello again
12         System.out.println("Goodbye World!"); // say goodbye
13         */
14     }
15
16 }
```

Note that since we have commented out all the code in the main method, this program is actually the same as the following program that we considered earlier:

```

1  class HelloWorld {
2
3      public static void main(String[] args) {
4          // System.out.println("Hello World!!"); //print hello
5      }
6
7  }
```

Let's now consider the following program that has two print statements and no comments:

```

1  class HelloWorld {
2
3      public static void main(String[] args) {
4          System.out.println("Hello!");
5          System.out.println("Goodbye!");
6      }
7
8  }
```

Try compiling and running the program. Note the output. What happens when we remove the linebreaks and spaces between each pair of print statements? Try removing them to get the following:

```
1  class HelloWorld {  
2  
3      public static void main(String[] args) {  
4          System.out.println("Hello!"); System.out.println("Goodbye!");  
5      }  
6  
7  }
```

Now compile and run the program. What happens? The behavior should still be such that it prints the same output:

```
1  Hello!  
2  Goodbye!
```

Whitespace such as linebreaks and spaces do not have any meaning in Java. While spaces are necessary to separate things like `class` and `HelloWorld` or `String[]` and `args` so that Java can distinguish between them, spaces are otherwise unimportant. Linebreaks have no significance at all in a Java program except to denote when an inline comment should end.

That being said, whitespace should be used to make your code more legible. Although you can write all your Java programs on a single line, it is considered very bad style! For demonstrative purposes, here the original Hello World Java program all on two lines:

```
1  class HelloWorld{public static void main(String[] args)  
2      {System.out.println("Hello World!");}}
```

Try compiling and running it. It should print `Hello World!` as its output. You can get rid of all the spacing after the close parenthesis `)` to get it all on one line—it just doesn't fit on the page if it's all on one line here!

Hardware

Today's lecture will focus on *computer systems* [AFM: should we change the title to systems, this is funky to say the chapter is hardware, and then say we focus on systems], groupings of devices that work together to perform a common task (or tasks). A user will interact with a computer through a variety of input and output devices (e.g. keyboards, mice, speakers, microphone, and monitors). A user's input will be processed, some computations will be performed, and then the resulting output will be displayed to the user. When most of us think about computers, we often think of a desktop or laptop computer, equipped with a keyboard, mouse, and monitor as seen in Figure 3.1. However, many things we interact with daily are computerized, including cell phones, cars, traffic lights, smart watches, televisions, and manufacturing lines. Today each of these items have sensors to perceive the real world, use an embedded computing device to understand the sensory input, and use a combination of display and mechanical devices to interact with the real world.

Example 3.0.1: For intersections across busy roadways, some traffic lights are computerized to optimize road traffic. These lights will stay green along the busier of the two roads, and use cameras or pressure sensors to detect the presence of cars along the less busy of the two roads. When cars arrive, the lights switch, allowing the cars on the less busy road to cross.

Today we will introduce three fundamental parts of computer systems: input and output devices, memory, and the central processing unit (CPU). These components work together to perform the basic building blocks of input processing, storage, control, and output. Understanding how the three parts work together will allow us to create powerful information processing tools. We will introduce each of these parts in turn. In Figure 3.2, we see how these parts come together to form a computer system.



Figure 3.1: A variety of computer systems: desktops, a laptop, a tablet, and a smart phone. [AFM: Source for pic.]

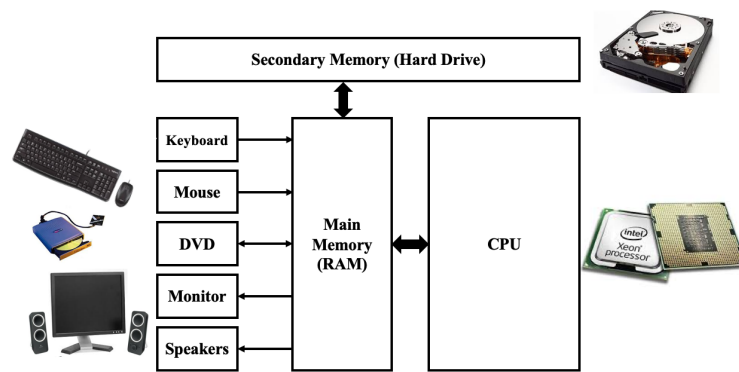


Figure 3.2: Interconnected parts of a computer system (keyboard, mouse, monitor, DVD player / burner, speaker, hard drive, CPU). [\[AFM: source for pic\]](#)

3.1 Input and Output Devices

Input and output (IO) devices allow the computer to interact with users and the world directly. Without these devices, a computer system would be very boring, always performing the same computation each time it's used. Even if it did compute a different value we would be unable to examine the value. A computer needs to be able to accept input and provide output. The first computers would occupy a large room in an office building and connect to a terminal (a keyboard and a text screen) in another room for users to interact with. Thanks to the hard work of electrical engineers, computers can now fit in the palm of your hand while being much more powerful. Likewise, many more types of input and output devices are now available. We still have the keyboard and monitor, and the mouse was invented for interacting with graphical displays. Today's phones are more computer than phone, coming equipped with speakers, microphones, touch screens, cameras, fingerprint scanners, radio transmitters, and much more. Computers even come embedded in other devices like cars, traffic lights, X-ray machines, and thermostats to both control and monitor the devices. As shown in Figure 3.2, these devices connect to the rest of the computer through the computer's memory. This kind of input and output is called memory mapped I/O (input and output). Creators of input (or output) devices are assigned a section of the computer's memory to write (read resp.) data. The computer will then read (write resp.) data to those locations to communicate with the given device. The creators of these devices, agree upon a known format to read and write data.

Example 3.1.1: You can think of communication between devices and computer as similar to leaving messages for a friend in a locker. Only you and your friend have access to this locker, which only holds space for one message. The format you agree upon is which language you'll use to speak (e.g. English) and any special keywords or phrases. You might agree with your friend that if either of you write a message saying "The morning is upon us" that the other will wait until "The night has come" before leaving any new messages.

The format that devices and computers communicate in are generally very simple and structured to permit fast and easily understandable communication for computers and devices.

Example 3.1.2: A monitor is a graphical display for computers. Let's consider a monitor connected to a computer that only displays in black and white images that are 20 x 20 pixels large. The monitor and keyboard agree upon using the following

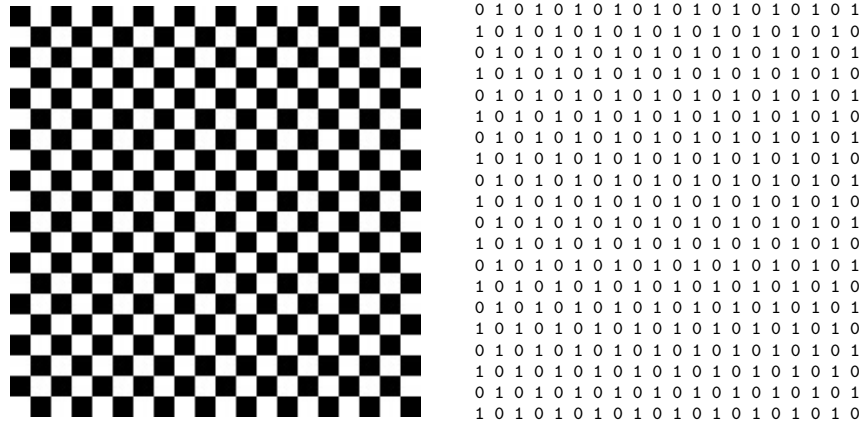


Figure 3.3: An example checkered image and its encoding — newlines and spaces added for readability.

format to communicate. The format is black and white images that are 20 x 20 pixels large. Each pixel's value is represented at 0 for black and 1 for white. Then an image is represented as a $400 = (20 \times 20)$ long sequence of pixel values. The sequence is ordered left to right, top to bottom. Now that both the monitor and computer agree upon the communication format, the computer can write images to the section of memory dedicated to the monitor and the monitor will read the image and display the image on it's screen. Figure 3.3 displays an example image, a 20 x 20 chekerboard with its encoding.

Note: While this is a simplified example, this is similar to how modern graphical displays communicate with computers.

3.2 Memory

With *memory*, computers gain the ability to store and recall data. This is very similar to physical storage of items. Figure 3.4 shows three storage locations — a storage closet, a garage, and a warehouse. Each of the three locations make a tradeoff between convenience of location and storage capacity. The closet can contain a few things and is the same room you need it. The garage can fit even more things and is only a walk outside (or through) your home, and the warehouse can fit practically anything you would want to store but you have to drive to the warehouse to pick-up or store your items. Similarly, a computer's memory makes the same trade-offs.

There are two major types of memory, *Main Memory* (RAM) and *Secondary Memory*, like hard disks, solid-state drves, tape drives, and more. Main memory is *volatile*, meaning that the contents of the memory is not preserved when a computer is turned off and back on. On the other hand Secondary Memory is meant to be *persistent* or *nonvolatile*, and does not go away when the computer is turned off. Main Memory can be thought of as the “scratch paper” the computer uses for computations. Computers will also use Main Memory as a conduit for communicating between the CPU and all other parts of the computer. Main memory is closer to a garage (where you can lose items when you turn off the lights) — there is enough room to fit most items you use regularly and is close enough to not worry about the time it takes to get to the garage.

In most modern computers, programs are treated as data. The individual instructions that combine to form a program are stored in memory just as data is. It is the job of the computer to properly understand if a segment of memory is data or a program. The computer is able to fetch data from Secondary Memory to Main Memory or persist data in Main Memory to Secondary Memory when needed. However, this process of transferring data between Secondary and Main Memory can cost a lot of time relative to keeping data in Main Memory only.



Figure 3.4: Storage closet, garage, and warehouse trading off between capacity and locality. [\[AFM: picture references\]](#)

3.3 Central Processing Unit

The final part of a computer we will introduce today is the *central processing unit* (CPU), also known as the *processor* or *main processor*. The CPU is the physical circuitry of a computer that performs instructions. The CPU has several key components: the control logic, arithmetic and logic unit (ALU), registers, program counter (PC), and clock. These components work together to fetch, decode, and execute all instructions — the building blocks of all programs. Instructions vary between different brands of CPUs, but, in general, they will include arithmetic, control, read (from memory), and write (to memory) functionalities. Example 3.3.1 shows several instructions that together would perform $x = x + y$, given x is stored in memory location 16 and y at memory location 20. These instructions are quite low level, and harder for humans to read than the programs we will write in this course. However, the programs we write will be translated into these instructions to be easily understood and executed by the CPU.

Example 3.3.1:

```
load R1 16    -- Load value at memory location 16 into register 1
load R2 20    -- load value at memory location 20 into register 2
add  R1 R2 R1 -- add the value in register 1 to the value in register 2
              and store in register 1
store R1 16   -- store the value in register 1 to memory location 16
```

A key component of a CPU is its clock. The clock allows the CPU to progress in time, triggering the time to progress from time t to time $t + 1$. A single time segment is referred to a clock cycle. You might have heard about a computer’s CPU speed (e.g. my computer runs at 2.4 GHz = 2.4 billion clock cycles a second). This is determined by how fast the clock transitions from one time step to the next. This clock drives the progress of the CPU. The time an instruction takes is measured in instruction cycles. The rate of the CPU’s clock is determined by the slowest operation of the CPU (e.g. fetch, decode, execute stage).

In addition to clocks, the CPU contains a group of memory locations called registers. A single register is capable of holding a single word of information (the smallest unit of data in a computer). The key benefit of registers is the ability for the CPU to immediately read and write the contents of the CPU. The value of registers can be updated on each clock trigger (i.e. on the change from time t to $t + 1$). Most modern CPU’s will have between 16 and 64 registers that programs may use. For comparison, accessing Main memory can take 10’s or even 100’s of instruction cycles to access while registers are immediately available to the CPU.

The control unit and program counter (PC) will fetch, decode, and output the controls for the execution of each instruction. The program counter holds the location of the next instruction to be executed. The next instruction is then fetched to the CPU. The CPU’s control unit then decodes the fetched instruction and outputs control signals (commands) to main memory and the ALU. The CPU may read data from memory (e.g. store) and then the ALU will then execute the action specified by the control unit

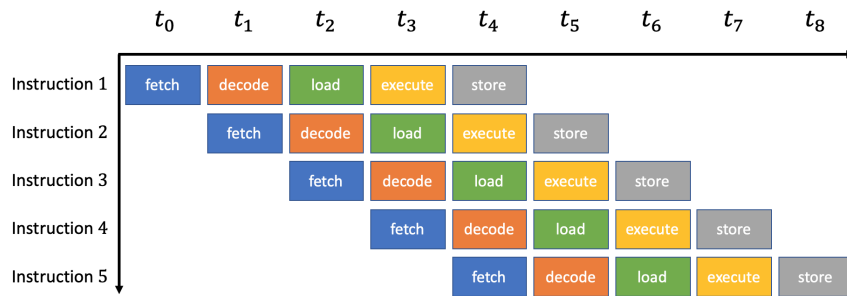


Figure 3.5: Instruction Pipelining. [AFM: citation]

(e.g. add, subtract, multiply, compare to 0, etc.), and then possibly write the output to memory.

In most computers, the CPU and its constituent parts are responsible for all computing needs of the computer. In some select systems, there will be additional hardware to perform specialized operations (e.g. graphics processing units for processing / producing images). It is the CPU's responsibility to control the computer and coordinate with devices to execute programs. As such, the CPU has seen a quick evolution to increase its processing power. Electrical engineers, originally focused on making the CPU smaller and smaller and thus quicker, following Moore's Law : every two years the size of a CPU shrinks in half. Additionally, CPU's were designed with a pipelining architecture (i.e. multiple instructions are executed in quick succession). This is done by noting that each stage of the five stages — fetch, load, decode, execute, and store — can be performed independently. Thus, while one instruction is being executed, the next instruction can be decoded. Figure 3.5 shows how pipelining is performed by executing the different parts of the pipeline in parallel for five consecutive instructions.

Due to the decline of Moore's Law in recent years, many CPU designers focus on increasing processing power by improving parallelism (i.e. being able to execute multiple instructions at a time). This allows instructions that do not depend on each other to be executed at a time. These CPU are referred to as multi-core, as they have multiple cores that each have their own set of registers, control unit, ALU, and PC but share main memory and a single clock. In this class we, will not teach how to effectively harness parallel programming, but note that this is an important progress in how hardware has evolved.

3.4 Conclusion

In this chapter, we covered the three fundamental parts of a computer system: input and output devices, Main and Secondary Memory, and the CPU. We discussed their roles, relationships, and basic capabilities. I hope that this will help you better understand how hardware works at a high level to better improve how to write programs that will eventually run on these computer system. In the next chapter we will begin discussing the concept of Software and how its similarities and differences to hardware and where the boundary between the two lies.

3.5 Learning Objectives

[AFM: weird to have this at the end] After covering this chapter, you should be able to answer the following questions:

1. What three parts comprise a computer system?
2. What are examples of common input and output devices?
3. Name an uncommon example of a computer system and explain how it may work.

4. How does memory mapped input and output work?
5. Name four kinds of memory devices and explain the difference between Main and Secondary Memory.
6. What parts comprise the central processing unit (CPU)?
7. Describe three possible methods to increase the computation power of a CPU.

Numbering Systems

How do we represent numbers? Based on how we represent numbers, certain operations can be easier or harder. Today, we'll consider two representations: decimal and binary numbers. Decimal numbers — decimal notation or base 10 — are what we're all used to seeing in our daily lives. For example we could consider the string of numerals:

19,874

And we know this represents the value nineteen thousand eight hundred seventy four. However, we could consider representing numbers differently. Depending on how we represent numbers specific operations become easier or harder. Let's go back to our example and now we'll multiply by 10.

$$\begin{array}{r} 19,874 \\ \times \quad 10 \\ \hline 198,740 \end{array}$$

We see it is very easy to find the answer — we only need to add a 0 at the end of our number.

Today we'll see how to find the value of numbers represented in decimal and binary notation. And how to convert between the two notations. We'll start with the decimal notation that we see in our everyday lives.

4.1 Decimal Numbers

When we see numbers in our daily lives there's an implicit assumption that these numbers are decimal numbers, whether this is a price, the temperature, or time. We can be explicit about what representation we are using by subscripting our numbers with the base.

$$19,874_{10} \quad \text{or} \quad 12.38_{10}$$

In this class, unless otherwise stated we'll consider all numbers to be represented in base 10.

4.2 Radix Decomposition

How do we actually go from numbers in base 10 to finding their value? We can break down the number by each position. Let's go back to our running example:

$$\begin{aligned}
 19,874_{10} &= 1 \times 10,000 + 9 \times 1,000 + 8 \times 100 + 7 \times 10 + 4 \times 1 \\
 &= 1 \times 10^4 + 9 \times 10^3 + 8 \times 10^2 + 7 \times 10^1 + 4 \times 10^0
 \end{aligned}$$

How do we know the value is nineteen thousand eight hundred seventy four? We can break it down digit by digit and add together the values. We have 4 in the 1s (the right-most position), 7 in the 10's position, 8 in the hundreds position, 9 in the thousands, and 1 in the ten-thousands position. Or in other words, at each successive position we multiply the value of the digit at that position by the next power of 10 to determine its value. This even works for non-whole numbers. Consider the number:

$$\begin{aligned}
 12.38_{10} &= 1 \times 10 + 2 \times 1 + 3/10 + 8/100 \\
 &= 1 \times 10^1 + 2 \times 10^0 + 3 \times 10^{-1} + 8 \times 10^{-2}
 \end{aligned}$$

We find 1 in the tens position, 2 in the ones position, 3 in the tenths position, and 8 in the hundredths position. For non-whole numbers we treat every digit to the right of the decimal (or radix) point exactly the same as we do for whole numbers. Everything to the left of the radix point we successively divide by the next power of ten. Now that we know how to determine the value of a decimal number, how can we do this for other number representations. We'll now transition to considering binary numbers.

4.3 Binary Numbers

What are binary numbers? They're just another way to represent numbers; however, instead of having ten digits (one to nine) we have two bits (zero and one). Last week, we learned about computer hardware. A computer's primary purpose is to compute; so we need to be able to store numbers on a computer's hardware. Without getting too technical, a computer represents numbers in a sequence of transistors, each storing an electrical charge. These charges can be on (high voltage) or off (low voltage) or somewhere inbetween. However, like with a light bulb that may burn brighter or be dimmer based on the incoming charge, we only consider if the state is on or off.

In effect, computers represent numbers in binary. To better understand how computers work, we'll practice with binary numbers. Let's consider a few example numbers.

$$1_2 = 1_{10} \quad 11_2 = 3_{10} \quad 101_2 = 5_{10} \quad 110.01_2 = 6.25_{10}$$

Here we see how we represent 1, 3, 5, and 6.25 in binary notation.

4.4 Binary Arithmetic

Just as we perform arithmetic on decimal numbers, we can perform arithmetic on binary numbers. The process is exactly the same, except we work with bits instead of digits. Let's start by looking at addition and multiplication.

$$\begin{array}{r}
 10110101_2 \\
 + 1010110_2 \\
 \hline
 100001011_2
 \end{array}
 \qquad
 \begin{array}{r}
 110101_2 \\
 \times 1101_2 \\
 \hline
 110101_2 \\
 11010100_2 \\
 + 110101000_2 \\
 \hline
 1010110001_2
 \end{array}$$

Here we see how to add the binary representations of 181 and 86 and get 267 as the result. We also see how to do long form multiplication of 53 and 13 to get 689. You'll notice this is exactly the same as for decimal numbers but we force ourselves to only work with bits, carrying the one to the next place when necessary. Both subtraction and division work similarly.

Now we'll turn our attention to three important arithmetic operations on binary numbers **and**, **or**, and **exclusive or (xor)**. These are generally called bit-wise operations as they apply to each bit place without considering the result of other placements. For these operations, we'll pad the shorter of the two numbers with zeros when needed.

a	b	a and b	a or b	a xor b
0	0	0	0	0
0	1	0	1	1
1	0	0	1	1
1	1	1	1	0

Let's practice a few examples!

101011001_2	101010110_2	1100011_2
and 11101_2	or 111011_2	xor 110110_2
11001_2	10111111_2	1010101_2

We notice that all of the operators perform the operation to each bit position independently of the rest of the positions in the numbers.

4.5 Binary to Decimal Conversion

How do we go from binary numbers to decimal and back? We'll now examine how to take binary numbers and convert them to their decimal notation. This works exactly like the decimal decomposition we learned at the beginning of class. And in fact this works to convert a number represented in any base to decimal.

$$\begin{aligned}
 1010110001_2 &= 2^9 + 2^7 + 2^5 + 2^4 + 2^0 \\
 &= 512 + 128 + 32 + 16 + 1 \\
 &= 689_{10}
 \end{aligned}$$

4.6 Decimal to Binary Conversion

Now, let's look at the opposite conversion, converting decimal numbers to binary. This will use successive division and subtraction as opposed to addition in multiplication. We will successively divide by two, if the remainder is non-zero we'll keep the bit at the current position. We'll continue until we can no longer divide by two. Again, this works for any base, not just base two.

Initial Value	Result	Remainder
267	133	1
133	66	1
66	33	0
33	16	1
16	8	0
8	4	0
4	2	0
2	1	0
1	0	1

$$267_{10} = 100001011_2$$

4.7 Conclusion

Today, we covered two very important number representations. In particular, decimal notation that we're used to seeing in everyday life and binary notation that is useful for understanding how computers manipulate numbers. We covered binary addition, multiplication, and three bitwise arithmetic operations (**and**, **or**, and **xor**). Then we showed how we can convert between binary and decimal notation for numbers. Next lecture, we'll review what we've learned about binary numbers and introduce Hexadecimal numbers, another important number representation in computer science.

4.8 Learning Objectives

After covering this chapter, you should be able to answer the following questions:

1. How many digits are in $142,123,098_{10}$?
2. How many bits are in 100101110010101_2 ?
3. Convert 421_{10} to binary.
4. How many bits do you need to represent $123,042,982_{10}$?
5. How many digits are required to represent 10110110_2 ?

Programming Languages and Linux commands

Today's lecture will focus on programming languages and how to give a computer instructions through a command line. Programming languages allow programmers to write instructions that can be interpreted by a computer to perform operations. Computers can not understand programming languages directly, so they have to be translated to a format that the computer can understand. In Java that transformation is referred to as compilation. When Java code is compiled, a program called a compiler takes the Java code as input and as output produces code that's easier for a computer to execute directly.

For instance, suppose that you have a text file with a large number of itemized expenses and you'd like to calculate the total of the expenses. Writing instructions for this task that the computer's hardware could understand directly could take more time than summing the entire list manually. Using a programming language to specify the instructions allows you to efficiently tell the computer to perform this task.

5.1 Most popular programming languages

What programming languages are most used depends on the application that the programming language is best used for. Java is one of the most popular programming languages. It's used extensively by large companies; it powers many of the services of companies like Oracle and Amazon. It's also one of the primary languages used to develop applications for the Android mobile operating system.

JavaScript is another very popular language. Despite its similarity to Java in name, it is a very different language from Java. It's the primary language used for writing code to interact with browsers in websites. It powers features of websites like generating dynamic content, sending requests to a server, making sure that the information entered in forms is correct, and many more things. It's also increasingly used to write applications. It's a popular choice among newer and smaller companies.

The number and diversity of programming languages can be intimidating, but learning new programming languages becomes easier once you have a solid understanding of fundamental concepts. In this course we use one programming language, but the ideas and concepts you learn can help you learn any programming language you need to. We recommend focusing on one programming language initially to learn fundamental concepts rather than spending time trying to learn multiple languages. We've decided to focus on Java because it's easier to learn, widely used, and illustrates many important concepts in programming. 40% of developers who responded to a large survey

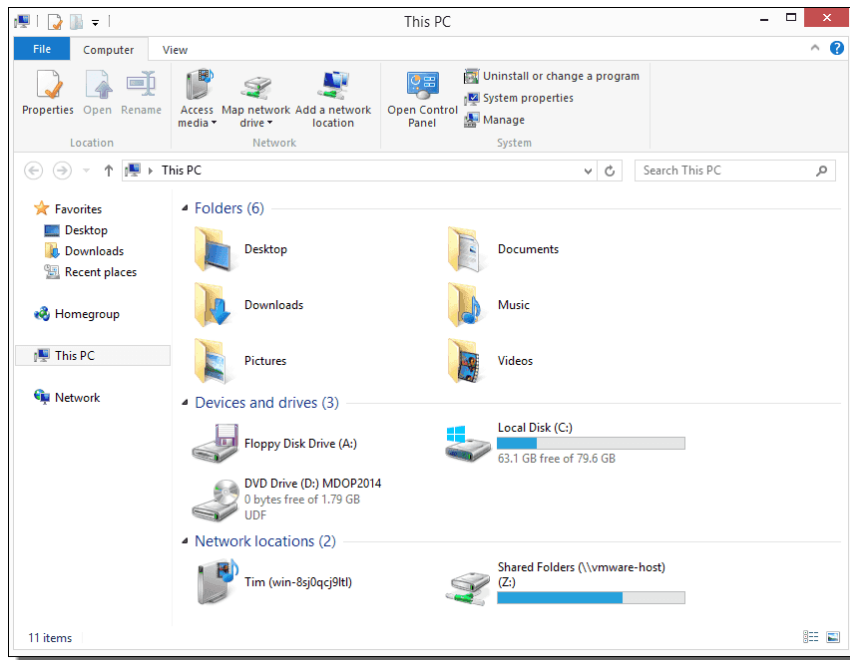


Figure 5.1: Windows file explorer

of developers stated that they know Java!

5.2 Overview of a GUI

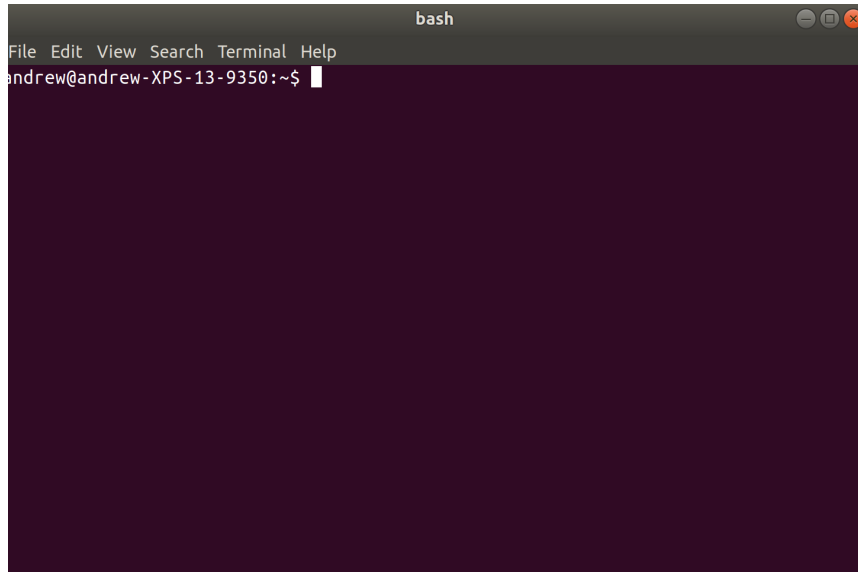
Using a graphical user interface (GUI) or command line are the two primary ways to interact with your computer and complete tasks like edit files, create folders, and run programs. A graphical user interface like a file explorer (see figure 7.1) provides an easy way for users to complete tasks on their computer. A graphical user interface is the primary way most users of a computer interact with the computer. It's popular because users can use their mouse and intuitive commands to complete basic tasks.

The graphical user interface abstracts the details of the commands being issued to the computer away from the user of the program. This abstraction provides simplicity, but it also limits the flexibility of what a user can do with the computer. Furthermore, it means that users can only interact with a system that has a graphical user interface installed. This doesn't include servers, computers without standard operating system, or computers that a user is connected to remotely.

The command line is a program you can access through the GUI of a linux machine. It's also installed on servers without a GUI and can be accessed directly through monitors (this was the only way to interact with computers before the invention of the GUI).

A command line allows the user to issue a broader range of commands and interact with computers without a graphical user interface. It gives users the ability to interface directly with the operating system and issue commands that it can understand. This increases the range of the commands that the user can issue, giving them more flexibility to perform complex and custom tasks.

We will discuss the command line used in the Linux operating system because it is one of the most command and intuitive command lines.

**Figure 5.2:** Linux command line

Command	Description
ls	List all the files in the current directory
touch	Create a new file
cd [directory]	Changes the current directory to the directory given
mkdir	Create a new directory within the current directory
man [command]	Display the manual for a certain command
cat [file]	Display the contents of file in the terminal
mv [file] [location]	Moves files to a new location. By default it moves file to a different name within the same directory.
cp [file] [location]	Creates a new file with identical contents to the file specified in a new location. If a path for the file isn't given, it will copy the file to the same directory.

A folder in Linux is referred to as a directory. When you open a command line, you're in a certain directory within the file system. All commands that you issue that relate to interacting with the file system like creating files, editing files, or renaming files will be issued in the context of this directory. This is like navigating to a particular folder within the GUI file explorer and completing all your operations relative to that folder.

Figure 7.2 shows the command line interface for a Linux computer. The prompt shown is where the user can type commands. By default it prepends the user's name, the name of the computer, and a dollar sign. The white rectangle represents the location of my cursor prompting the user to enter commands.

To make a new directory, the user can use the 'mkdir' command. This is like making a new folder within an existing folder in the window's file explorer.

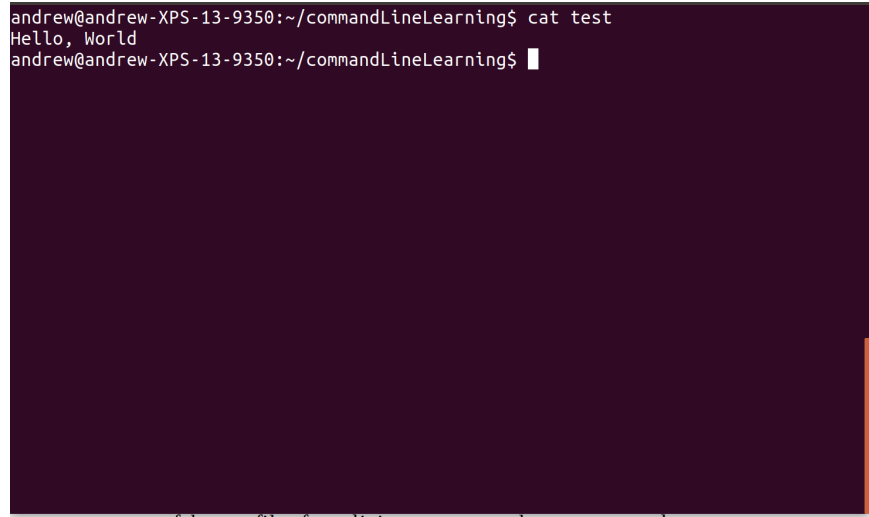
To navigate to that directory the user can use the 'cd' command, short for change directory. Figure 7.3 displays an example of creating a new directory and changing to use that directory. In this example the `~` sign prepending the command represents the user's home directory. The home directory is the outer most folder for that user's account on the server or machine.

```
andrew@andrew-XPS-13-9350:~$ mkdir commandLineLearning
andrew@andrew-XPS-13-9350:~$ cd commandLineLearning/
andrew@andrew-XPS-13-9350:~/commandLineLearning$
```

Figure 5.3: Creating a directory

```
andrew@andrew-XPS-13-9350:~$ mkdir commandLineLearning
andrew@andrew-XPS-13-9350:~$ cd commandLineLearning/
andrew@andrew-XPS-13-9350:~/commandLineLearning$ touch test
andrew@andrew-XPS-13-9350:~/commandLineLearning$ ls
test
andrew@andrew-XPS-13-9350:~/commandLineLearning$ cd ..
andrew@andrew-XPS-13-9350:~$
```

Figure 5.4: Listing files and navigating

A terminal window with a dark purple background. The prompt is 'andrew@andrew-XPS-13-9350:~/commandLineLearning\$'. The command 'cat test' has been entered, and the output 'Hello, World' is displayed on the next line. The prompt is now 'andrew@andrew-XPS-13-9350:~/commandLineLearning\$' with a cursor at the end.

```
andrew@andrew-XPS-13-9350:~/commandLineLearning$ cat test
Hello, World
andrew@andrew-XPS-13-9350:~/commandLineLearning$
```

Figure 5.5: Cat command example

In figure 7.3 I create a `commandLineLearning` directory and change my current directory to that directory. If I want to create a file in that directory, I can use the `'touch'` command to create a new file. If I want to see the files in the directory, I can use the `'ls'` command to list all of the files in the current directory. If I want to move back to the directory I was in before, I can use the `'cd'` command with the parameters `'..'` to move to the directory my current directory is contained within. In Linux, `'..'` refers to the directory that contains your current directory, or the parent directory. These commands are summarized in figure 7.4.

Try these commands out. You can edit the contents of `test` using any file editing software. To view the contents of the `test` file after edit it, you can use the `cat` command.

References

Stack Overflow Developer Survey 2019. (n.d.).

Networks I: Protocols

Computing devices have been in use for thousands of years, at least since the invention of the Sumerian abacus between four and five thousand years ago. Even as computers became more sophisticated, a long process culminating in the development of transistor-based digital computers, they remained specialized machines, largely for use in government or business settings. The explosion in personal computing, which has flooded the planet with billions of computers that pervade almost every aspect of society, coincided roughly with the development of computer networking and the Internet. Communicating via computers – with friends, news outlets, vendors, banks, or even artificial intelligence systems – has revolutionized the way the world works.

In the next two chapters, we'll learn how computer networks work, from the wires in the ground up to the more familiar technologies in a browser or email client. The Internet is one of the largest and most complex machines ever built, so we certainly won't be able to cover every aspect of it, but we'll focus on a few important topics. After reading, you should be able to tell a story about what happens when you open a website or send an email, and understand how to get under the hood and see how different pieces are working together to complete the task.

In this chapter, we'll focus specifically on the *protocols* that allow computers to communicate with one another. We'll cover what a protocol is, and why detailed protocols are important for ensuring reliable networking. In the following chapter, we'll look at how these protocols work together and allow us to build large computer networks, including the global Internet.

In addition to its inherent interest, understanding networks is very practical. Many companies look for Internet Technology or IT specialists who understand how networks work. Taking a test like the Network+ certification or A+ certification allows access to these jobs. These tests are knowledge based, so you don't need practical experience to do well on them. They qualify you for entry level positions that pay around \$50,000 a year. If you're interested in preparing for these exams, there are many books that can help you prepare. For instance, "CompTIA A+ Certification All-in-One Exam Guide, Ninth Edition (Exams 220-901 & 220-902)."

6.1 Protocols

Long before there were computers, there were protocols. A protocol is a set of rules and procedures that specify how two entities should interact. In particular, protocols can apply to people, in the form of social protocols. These can be simple – for example, there is an unspoken protocol regarding shaking someone's hand – or quite complex, as in the protocol for diplomatic exchanges between different countries. In fact, the US Department of State employs a Chief of Protocol, whose office is tasked with ensuring that visiting diplomats and heads of state are appropriately welcomed and not offended by any "breach of protocol."

Rather than delve into complex diplomatic protocol, let's consider the simple example: the "handshake protocol." It's so easy to shake someone's hand that we might not realize how many rules and procedures we follow when we do it. Here are a few possible "breaches of protocol" that could occur when Alice shakes Bob's hand.

- Alice extends a left hand, while Bob extends a right hand.
- After locking hands, Alice refuses to let go for a full 15 seconds.
- After locking hands, Alice lets go immediately without pausing for a shake.
- Alice refuses to make eye contact for the duration of the handshake.
- Alice shakes her hand left to right, or forward and backward.
- Alice makes no effort to grip Bob's hand, leaving the handshake loose.
- Alice uses a "death grip," causing Bob discomfort.
- Alice extends her arm fully and locks her elbow, keeping Bob at as great a distance as possible.
- Alice barely extends her arm, forcing Bob to stand very close.
- Alice sneezes into her hand immediately before initiating the handshake.

Any of these examples would cause some problems; many would lead Bob to wonder if Alice had ever shaken hands before, while a few might leave Bob offended or wondering if Alice is interested in communicating. Clearly it is important that Alice and Bob abide by some rules so that their handshake goes smoothly and gives a friendly message to both parties.

How might we construct such a set of rules, or a protocol? It would be helpful to first divide the above gaffes into a few categories. We could do this chronologically: there are potential problems with initiating the handshake, completing the handshake itself, and terminating the handshake.

Let's first write a protocol for initiating a handshake. The two parties need to come to an appropriate distance, verify that their hands are clean, make eye contact, and both extend the same hand. We can condense this into a step-by-step procedure:

Handshake Initiation

1. Verify that your right hand is acceptably clean, and maintain this status until handshake termination.
2. Look at the other party, and wait for eye contact to be established.
3. After making eye contact, approach the other party and come to within about one arm length apart.
4. Extend your right hand to the halfway point between you and the other party, at about waist height, and orient your hand so that your open palm is facing left.
5. Wait for the other party to extend their hand likewise.

There are a surprising number of steps for something we all do almost unconsciously!

Exercise 6.1.0: Write similarly detailed protocols for Handshake Completion and Handshake Termination. Make sure your protocol eliminates the possibility of any of the gaffes listed above.

Computers have no social intuition whatsoever. This is why protocols are important: a computer needs to be told every single step of an interaction as simple as a handshake in order to execute it properly. Building a computer network requires a number of different protocols to establish how messages are sent on wires, routed to the appropriate destination, and more.

One of the most important parts of a protocol is that everyone follows the same one. Some aspects of protocols are completely arbitrary, and it is important for everyone to follow the same arbitrary convention. For example, handshakes could use left hands instead of right hands. Luckily the whole world shakes right hands, but anyone who's tried driving in both the US and UK knows that not all protocols are globally uniform.

In order to ensure that networking protocols are consistent around the world, a body known as the Internet Engineering Task Force (IETF) carefully specifies all the protocols and makes them publicly available. Publications of the IETF are known as Requests for Comments, or RFCs. The name "Request for Comments" originated in the early days of networking when the documents really were intended for discussion and subsequent revision. Nowadays there is a category called standards-track RFCs which specify the standards of the Internet and which should be followed by anyone building networking technology.

6.2 Network Layers

In the previous section we saw an example of a protocol for initiating a handshake. The protocol was more detailed than you might have expected, but it also didn't include all the details. For example, let's delve into step 2: "Look at the other party, and wait for eye contact to be established." First, what does it mean to "look at the other party"? This requires identifying the location of another human, and sending motor signals to the neck and to the eyeballs in such a way as to orient them in the direction of the other human's eyeballs. If we were to include all these details, our protocol would become extremely complex and hard to follow.

In real life, we avoid this complexity by never even thinking about these details of movement, facial recognition, and the like. Our brains can automatically send the right signals to our muscles, interpret signals from our senses, and so on. We say that the brain is segmented into two layers: conscious and unconscious. The conscious layer sends high-level signals, such as the command "Look at Bob," to the unconscious layer. The unconscious layer is responsible for parsing those commands and executing them by sending the right signals to muscles. This layered architecture allows the protocols to remain simple by referring to the capabilities of lower layers.

Communication networks also use a layered architecture. We will come to the layers that make up the Internet in a moment, but long before the Internet we had the postal system. The postal system is built up from layers, and in many ways the Internet is simply a much faster and more automated version of the same layers. Let's consider the following example of a letter being sent from Alice at Acme, Inc. to Bob at Boxes, Inc.:

1. Alice tells her employee, Jane, to order 50 new boxes. Jane drafts a letter to Bob at Boxes, Inc.: "Please send 50 boxes." She seals it in an envelope, opens her address book and copies down Bob's mailing address, and drops the envelope off at the company mailroom. A worker in the mailroom sees that the letter is addressed to a destination outside Acme, Inc., and so he places the envelope in a postal box outside.
2. A postal worker picks up the envelope from the postal box, and sees that the letter is addressed to a location outside the city. He brings the letter to the local post office, and places it into a bin to be sorted.
3. Another postal worker picks up the envelope from the bin, and puts it on a mail truck destined for Bob's city.
4. The mail truck drives along a highway from Alice's city to Bob's city, and gives the letter to the post office there.

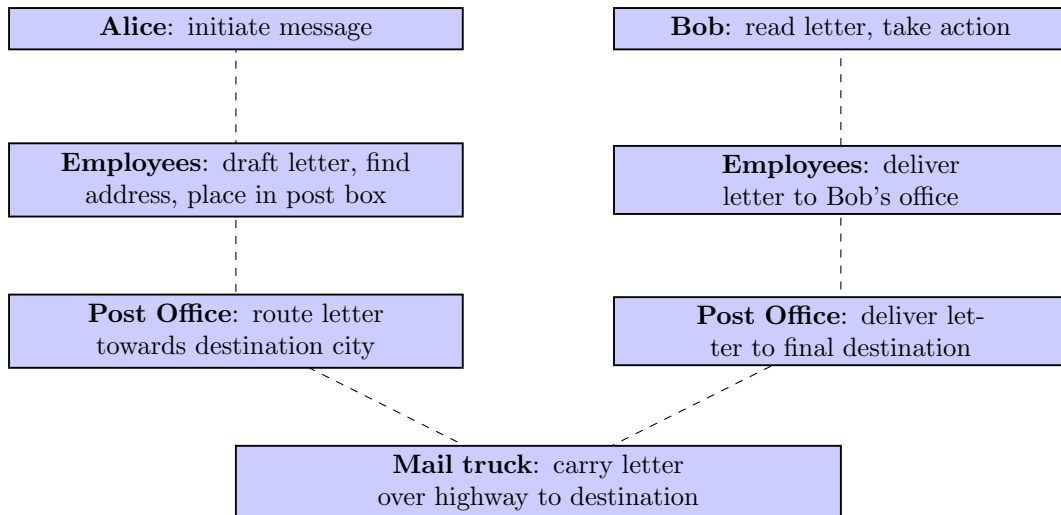


Figure 6.1: Layers of the postal system.

5. A postal worker in Bob's city carries the envelope to Boxes, Inc., and places it in their mailbox.
6. A worker in the mailroom sees the envelope addressed to Bob. She brings it to Bob's office and puts it on his desk where he'll see it.
7. Bob reads the letter, and asks an employee to begin producing 50 new boxes for Acme, Inc.

There's a lot going on here – when you think about it, it's amazing that the letter managed to reach Bob at all, given everything that had to go correctly. The trick here is the layered architecture. Alice didn't need to know what highway the letter would be carried on, and the mail truck driver didn't need to know where to find Bob's address. Everyone in the pipeline has a specialty, and no one needs to understand the full process for it to work well. Figure 6.1 shows the different steps arranged in layers, from Alice and Bob at the top down to the truck driver at the bottom.

These layers all have counterparts in the Internet. At the top, we have the *application layer*: programs which rely on the Internet and know how to initiate messages, much like Alice and Bob. The application layer relies on the *transport layer* to determine an address where the messages should go, just as Alice relied on Jane to find Bob's mailing address. The transport layer relies on the *Internet layer* to use the address to send the message to the right location, just as the Post Office routed Alice's letter to Bob's city. Finally, the transport layer relies on the *link layer* to physically transmit data over wires or wireless protocols, just as the Post Office relies on mail trucks to move envelopes.

This set of four layers is commonly known as TCP/IP, named for the Transmission Control Protocol (TCP) which commonly runs the transport layer, and the Internet Protocol (IP) which commonly runs the Internet layer. The following sections contain more details on each of the layers and the protocols they use. You should always keep in mind that these layers work together just like the layers of the postal system work together.

6.2.1 Application Layer

The application layer is the highest layer of TCP/IP. It consists of various protocols that specify how to send different types of data over the Internet. You might have heard of one of these protocols, HTTP. HTTP stands for HyperText Transfer Protocol. It comes at the beginning of a web address. When you type `http://www.google.com`, you're telling your computer to use HTTP to access `www.google.com`.

There are many other protocols in the application layer. HTTPS is a secure version of HTTP, which uses security protocols we'll discuss in Section 6.3. Another popular

Protocol	Port	Usage
HTTP	80	Web traffic
HTTPS	443	Secure web traffic
FTP	21	File transfer
SSH	22	Remote shell access
IMAP	993	Receiving email
SMTP	587	Sending email

Table 6.1: Commonly used Internet protocols.

protocol is SSH, or Secure Shell, which programmers use to access a Linux terminal on another Internet-connected machine. The file transfer protocol, or FTP, is used for exchanging files between different computers. A variety of protocols, such as POP3, IMAP, and SMTP, are used for exchanging emails.

When the application layer receives data from the transport layer, it needs to know which protocol it was sent with. In order to sort out all the different data, the application layer uses *ports*. All Internet communications happen over some numbered port. For example, HTTP uses port 80. When you load `http://www.google.com`, your computer sends the request to Google over port 80. Google sends its home page back over port 80. Table 6.1 lists the most commonly used protocols and the ports they use.

If you write a program that uses the Internet, you will interact primarily with the application layer. A very common example is writing a program which runs on a web server. A web server is a computer which is connected to the Internet, and runs a program that constantly listens for incoming communications on ports 80 and 443 – the ports used for the World Wide Web. When it receives a request, it will run a program and send the output of that program back to the requester. Every website is hosted by a web server, which knows how to send the data of the website to anyone who requests it.

A simple and famous example is `isitfriday.net`. This website is hosted by a server which listens for requests. When it gets a request for `isitfriday.net`, it will run a program which checks if it is Friday (you’ll learn how to use “if” logic in programs in the next chapter). If it is Friday, it sends back a website that says “Yes.” Otherwise, it sends back a website that says “Not yet.”

6.2.2 Transport Layer

When the application layer is ready to transmit a message, it will pass it down to the transport layer. By far the most common transport layer protocol is TCP, for Transmission Control Protocol. We will focus on TCP in this section. Another common protocol in the transport layer is TLS, for Transport Layer Security, which runs alongside TCP. We will discuss TLS more in Section 6.3.

When TCP is asked to transfer a message, it will prepare to send it piece by piece over the network. The individual pieces are called *packets*. These packets are relatively small, about a kilobyte. For reference, an average-sized image would typically be broken into a few thousand packets before being transferred.

The reason for this is that a major role of TCP is to guarantee reliability. The application layer assumes that its message will be sent in full and without errors, while the Internet layer does not guarantee accurate or successful transmission, so it is the job of the transport layer to ensure robust communication. By breaking a message into small packets, TCP helps to reduce the chance that any individual packet fails to transmit. In addition, if a packet does fail to transmit, it doesn’t take too much time to recover, since TCP can re-send the relatively small packet.

In addition to providing reliability, TCP is responsible for initiating connections between two computers. When TCP receives a request to transfer a message to a destination, it will first complete a “handshake” with that destination. The goal of the handshake is to establish that both parties are connected to the Internet and ready and willing to communicate with one another. In addition, the parties need to synchronize a *sequence number* which specifies the numbers that will label subsequent packets.

We have already thought about handshakes between humans in some detail, so now let's delve into how TCP facilitates handshakes between computers. If humans were communicating using a TCP-style handshake, the conversation might look like this:

1. Alice says to Bob, "Let's start talking. I'll send some packets, starting with this one, #3874."
2. Bob says to Alice, "I hear you Alice, and I'm looking forward to packet #3875. My packets will start with this one, #8452."
3. Alice says to Bob, "Great, I'm looking forward to packet #8453."

This is the *three-way handshake* model. First, Alice notifies Bob of her intent to communicate and synchronizes her sequence number. Then Bob sends an acknowledgement of Alice's sequence number, and also synchronizes his own sequence number. Finally, Alice acknowledges Bob's sequence number, and then they are ready to communicate.

There are two kinds of messages being sent in the handshake, synchronization and acknowledgement. These are typically abbreviated as **SYN** and **ACK**. The TCP handshake protocol works as follows.

TCP Handshake

1. *A* generates a random number x . *A* sends SYN x to *B*.
2. *B* generates a random number y . *B* sends ACK $x + 1$ to *A*, and then sends SYN y to *A*.
3. *A* sends ACK $y + 1$ to *B*.

6.2.3 Internet Layer

Perhaps the most interesting layer of the TCP/IP stack is the Internet layer, which almost always consists of the Internet Protocol (IP). This layer corresponds to the Post Office in our analogy: it is responsible for routing data where it needs to go. Just as the Post Office uses mailing addresses to specify locations, the Internet layer uses IP addresses to specify computers on the network.

An IP address is a sequence of 32 bits, grouped into blocks of 8. When they are written out, the groups of 8 bits are converted to decimal numbers between 0 and 255. The groups are separated by dots, so sometimes an IP address is referred to as a "dotted quad." For example, the IP address 172.217.15.78 points to the server which hosts www.google.com. These addresses are stored in a sort of phone book called the Domain Name System, or DNS. When you type a website into your browser, it will first look up its IP address in the Domain Name System, and then request that IP address over the network.

When the Internet layer receives a packet from the transport layer for transmission, it will look at the IP address of the destination. It knows how to find a server which is closer to that destination. For example, if a server in New York City sees a packet destined for an IP address in Washington, DC., it might decide to forward that packet to a server in Philadelphia, PA.

You can view this process yourself using a Unix utility called **traceroute**. Open a terminal and type **traceroute** followed by any website. You'll see how the Internet Protocol routes your packets through different servers in order to reach their destination. For example, below is the route from a personal computer to youtube.com. The address in the last line, 172.217.12.238, specifies the location of one of YouTube's servers.

```
# traceroute youtube.com
 1  192.168.1.1 (192.168.1.1)  0.448 ms  0.958 ms  0.966 ms
 2  * * *
```

```

3  ae1312-21.ARTNVAFC-MSE01-AA-IE1.verizon-gni.net (100.41.21.152)
   7.779 ms  7.766 ms  7.628 ms
4  0.ae10.GW13.IAD8.ALTER.NET (140.222.225.219)  9.139 ms  9.117 ms
   9.082 ms
5  72.14.218.232 (72.14.218.232)  8.507 ms  7.277 ms  4.378 ms
6  * * *
7  216.239.54.106 (216.239.54.106)  4.779 ms  108.170.232.0
   (108.170.232.0)  12.022 ms  108.170.246.33 (108.170.246.33)
   16.086 ms
8  108.170.246.67 (108.170.246.67)  12.976 ms  108.170.232.19
   (108.170.232.19)  12.769 ms  108.170.246.66 (108.170.246.66)
   10.935 ms
9  iad30s15-in-f14.1e100.net (172.217.12.238)  10.570 ms
   216.239.47.127 (216.239.47.127)  16.793 ms
   iad30s15-in-f14.1e100.net (172.217.12.238)  10.485 ms

```

The IP addresses we’ve described all correspond to version 4 of the Internet Protocol, or IPv4. There are about four billion unique IPv4 addresses, which seemed like plenty at the time of its creation. However, the growth of the Internet and the proliferation of mobile devices has led to the exhaustion of almost all IPv4 addresses. In order to ensure that the Internet can continue to grow, a transition to IPv6, which provides for many more addresses, is slowly taking place.

In the transport layer, we delved into the step-by-step procedure used for handshaking. These sorts of procedures are one important aspect of the protocols that run networks. Another important aspect is the format in which data is sent. The link layer, which we’ll come to in the next section, can only send sequences of bits – zeroes and ones – so it is important for both parties to agree on which bits mean what.

To give an example of this, we will look inside the packets sent by IPv4. You should not try to memorize the following, but instead try to appreciate the importance of having a detailed packet structure.

An IP packet is divided into its *header* and the *data*. The header contains information about how the packet is being transmitted, while the data is the actual contents of the packet. The format of IPv4 headers is shown in Table 6.2. The various segments specify the following information.

- **Version:** the first four bits specify the version of IP being used. For IPv4, these bits are always 0100, binary for 4.
- **IHL:** this is the Internet Header Length. This specifies the number of 32 bit blocks, called “words,” used by the header. In Table 6.2 we show a header with 5 words, so these bits are 0101. This is almost always the case, although it is possible to have a larger IHL value and include additional data in the header.
- **DSCP:** this stands for Differentiated Services Code Point. These bits are used to specify the type of service for the packet. This allows the network to give priority to some packets, such as those used for voice or video calling, which need to arrive as quickly as possible.
- **ECN:** this stands for Explicit Congestion Notification. This allows the network to flag packets which are experiencing congestion along their routes.
- **Total Length:** this specifies the total length of the packet, in bytes including both the header and data. Since it is 16 bits, the maximum value is $2^{16} - 1 = 65535$ bytes, or about 64 KB.
- **Identification:** these bits are used to group packets which have been fragmented due to a maximum transmission size imposed by the link layer, which we’ll come to in the next section.
- **Flags:** these are flags regarding fragmentation, used primarily to specify if this is the last in a series of fragments or if there are more fragments to come.
- **Fragment Offset:** if a packet has been fragmented, these bits specify where the current block of data falls into the sequence of fragments.

Bit	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31																																	
0	Version			IHL			DSCP				ECN		Total Length																																																				
32	Identification										Flags				Fragment Offset																																																		
64	Time to Live					Protocol																Header Checksum																																											
96	Source IP Address																																																																
128	Destination IP Address																																																																

Table 6.2: The IPv4 header format.

- **Time to Live:** this keeps track of how many hops a packet has made through the network. Every time the packet is forwarded along a new link in the network, this value is decreased by one. If it hits zero, the packet is dropped and the sender is notified. This is how the `traceroute` utility works: it sends packets with small Time to Live values, and waits for the notifications to see where the packet was after one hop, two hops, etc.
- **Protocol:** this specifies the transport layer protocol which is making use of IP. For example, for TCP the protocol value is set to 6.
- **Header Checksum:** these bits are used for correcting errors in the header. The idea of a “checksum” is encoded in the word: we could imagine transmitting a stream of bits, followed by a few bits that count the number of ones in the original stream. If a zero was flipped to a one, or a one was flipped to a zero, then the receiver will notice that the checksum does not match and can ask for the packet to be resent. In practice a checksum is a more complicated function of the bits, but the idea is the same.
- **Source IP Address:** this is the address of the device sending the packet. Note that it spans 32 bits, since an IPv4 address consists of four numbers each from 0 – 255, which require 8 bits each.
- **Destination IP Address:** this is the address of the device meant to receive the packet.

6.2.4 Link Layer

Finally, we come to the actual nuts and bolts of the Internet: how are individual pieces of data sent across the network? This part of networking is handled by the link layer.

There are two main link layer technologies you will encounter. The first is Ethernet, which is used for wired Internet connections. An Ethernet cable, also known as an RJ-45 cable, is shown in Figure 6.2. Each end of the cable has an 8-pin connector. Inside the cable, a twisted pair of wires can carry signals in both directions without having them collide with each other.

The other primary link layer technology is 802.11, more commonly known as Wi-Fi. With a wired network, signals can be sorted into separate wires so they don’t overlap, but over a wireless network signals are broadcast and have the potential to interfere with one another. In order to prevent this, Wi-Fi is designed to allow devices to operate at slightly different frequencies which are arranged into channels.

There are two major variants of Wi-Fi: 2.4 gigahertz (GHz) and 5 GHz. These are two different frequency bands. Within each of them there are multiple channels. At 2.4 GHz, which was the first frequency available for Wi-Fi, there are 11 channels. However, the channels overlap with one another, so it is best for devices to use channels 1, 6, and 11, reducing the effective number of channels down to 3. The newer 5 GHz protocol alleviates this problem by making hundreds of channels available. Most new devices can operate at either 2.4 GHz or 5 GHz, but some older devices only support 2.4 GHz.

The details of link layer technologies are complex, but luckily you rarely have to deal with them directly. In most cases, Ethernet or Wi-Fi have been engineered to the point where little effort is required to make them work properly, and you can focus on other layers of the network while trusting that signals are being carried appropriately by the link layer.

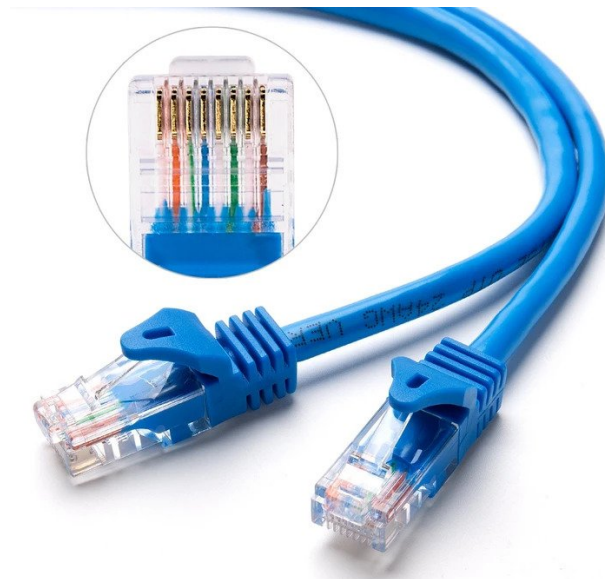


Figure 6.2: An Ethernet cable used to network computers together.

6.3 Network Security

One of the most important issues with computer networks is security. The Internet is frequently used for transmitting confidential messages, facilitating financial transactions, and other applications where security is paramount. For this reason, there are protocols dedicated to ensuring that the Internet can be secure. The most common such protocol is Transport Layer Security (TLS), which is a successor to an older protocol called Secure Sockets Layer (SSL). TLS is implemented at the transport layer of a network. Web traffic which uses TLS/SSL security is sent using the HTTPS protocol at the application layer, rather than the unsecured HTTP protocol.

Another common network security application is a Virtual Private Network, or VPN. A VPN allows computers connected to different LANs to communicate over the Internet as if they belonged to the same LAN. This is often used by companies to allow employees to access private resources even when they are away from the office. Security is very important for a VPN, and it can be provided using TLS or other protocols.

The details of TLS are complicated, and we will not get into them here. Instead, we will cover what exactly security means in a network setting, and look at three goals of any secure network: privacy, authenticity, and reliability. Achieving these goals requires using tools from cryptography, which you may learn about in a future course.

6.3.1 Privacy

When two parties want to securely transmit a message, one of their most straightforward goals is privacy. The two parties do not want a third party to be able to listen in on their communication.

Following a standard convention in cryptography, we will denote the two parties by Alice and Bob, and the adversarial third party by Carol. Alice is sending a message to Bob, and Carol wants to listen in. Since network signals travel over cables, Carol could listen in to the signal by physically tapping into a wire. Even more easily, on a wireless network, all signals are broadcast over the air and can be received by anyone who wants to listen.

In order to protect their messages from being read by Carol, Alice and Bob use an encryption algorithm. Alice and Bob have some shared key, and when Alice wants to send a message, she first uses her key to turn the plain text of the message into a cipher text. If the algorithm is secure, then the cipher text cannot be converted back to the plain text without knowledge of the key. This way, if Carol manages to read the message,

she will only see nonsensical cipher text. Only Bob, who has the key, can produce the plain text.

A key issue here is how Alice and Bob come to have the same key. If Alice tells Bob her key over the network, then Carol could retrieve the key and use it to decrypt all future messages. This is a fundamental problem in cryptography called the key-sharing problem. It is solved using a technique called public key cryptography. In TLS, public key cryptography is implemented using either the RSA algorithm, or a more sophisticated method called elliptic curve cryptography (ECC). These algorithms allow Alice and Bob to generate shared keys which Carol cannot intercept.

6.3.2 Authenticity

Even if Carol cannot read messages sent between Alice and Bob, she can still do damage. One potential attack would be for Carol to impersonate Alice, by sending a message to Bob which says it comes from Alice. This is why TLS must also protect the authenticity of messages: if a message says it comes from Alice, Bob needs some way of verifying that this is actually the case.

The public key cryptography algorithms used for ensuring privacy can also provide authenticity. In these algorithms, every party on the network has a private key which is never shared with anyone. This private key is important for generating the shared keys used for encrypting messages. It can also be used to generate a “digital signature,” which could only be produced by someone in possession of the private key. As long as the private key has never been shared, a digital signature can provide a guarantee of authenticity.

6.3.3 Integrity

Finally, even if Carol cannot read messages or impersonate Alice or Bob, she could try to modify one of their messages. This is called a man-in-the-middle (MITM) attack, where Carol intercepts Alice’s communications, edits them, and then forwards them to Bob with Alice’s signature intact.

To protect against this, TLS uses a message authentication code (MAC) which is sent along with any message. This is a short code sent with the message, which depends on the shared key and on the message. If the algorithms are secure, then it should be nearly impossible to generate a MAC without knowledge of both the shared key and the message. When Bob receives the message along with the MAC, he verifies that the MAC is correct. If it is not, then he knows the message has been altered in transit and discards it.

6.3.4 Security versus Trust

An enormous amount of expertise, testing, and iteration has gone into making TLS secure. However, it is still relatively common to hear about security problems on the Internet. Where do these problems come from?

In this context, it is important to recognize a key difference between security and trust. Network security provided by TLS achieves the three goals outlined above: privacy, authenticity, and integrity. When Alice communicates with Bob, she can be sure that only Bob can read her messages, that messages she receives from Bob really come from Bob, and that the messages have not been altered in transit. However, *TLS does not guarantee that Alice can trust Bob.*

There are many ways in which Bob could be unreliable. For example, Bob could be attempting to execute a phishing attack by impersonating someone else. TLS guarantees that Bob cannot truly impersonate `facebook.com`, but he could own `facebook.com` and, using TLS security, provide a website that looks just like `facebook.com` and asks for your login credentials. If you provide them to `facebook.com`, then Bob could read your password and log into your actual Facebook account.

Even if Bob is not malicious, he could be irresponsible. For example, if Bob runs an online vendor which collects credit card information to make payments, it is crucial that Bob takes appropriate security measures on his local network to protect that credit card

information from becoming known to an attacker. TLS does not guarantee that Bob has taken the necessary precautions. Many of the most substantial security problems arise on the Internet when many customers have trusted a vendor like Bob with their personal details like credit card numbers, and Bob suffers a data breach, putting this data into the wrong hands.

Network security is an evolving field as companies and governments work to protect against these sorts of vulnerabilities. There are many opportunities for employment in this area, as organizations increasingly need to rely on dedicated security professionals in order to keep their data and their clients' data safe.

Networks II: The Internet

In the previous chapter, we delved into Internet protocols and learned how the layered architecture of the Internet enables messages to be carried reliably without any single protocol becoming too complex. Now that we understand how to network computers together, we can discuss how these networks are laid out in practice. The most important distinction is between Local Area Networks, used to network a single home or business together, and Wide Area Networks, used to network larger areas or even the whole world.

7.1 Local Area Networks

The network used to connect computers within a single home, business, or organization is called a local area network, or a LAN. Even though most LANs are connected to the outside Internet, they can also be used to allow devices on the same network to connect with each other without going through the Internet. In this section, we will review some concepts of local area networking and discuss using switches to manage network traffic.

7.1.1 Network Topology

When building a LAN, a network engineer has to connect computers either via Ethernet cables or wireless connections. For simplicity, we will focus on wired connections in this section.

In order to allow for a group of computers to communicate with each other, each computer must be able to reach any other computer via some number of hops in the network. The *topology* of the network is the layout of the connections in the network, and this sets how a pair of computers communicates.

One simple way to allow computers to communicate would be to connect all the computers in a line. Then each computer can communicate with the computers directly adjacent to it. This is called a linear topology, also known as daisy-chaining. Figure 7.1 shows a network with linear topology.

With a linear topology, if a computer wants to send a message to a computer which is not its neighbor, the intermediate computers have to carry the message along. In order to facilitate this, the destination address is included along with a packet. If a computer receives a packet and sees that the destination address is not its own, it will forward the packet along in the direction of the destination address.



Figure 7.1: A network with linear topology.



Figure 7.2: A network with bus topology.

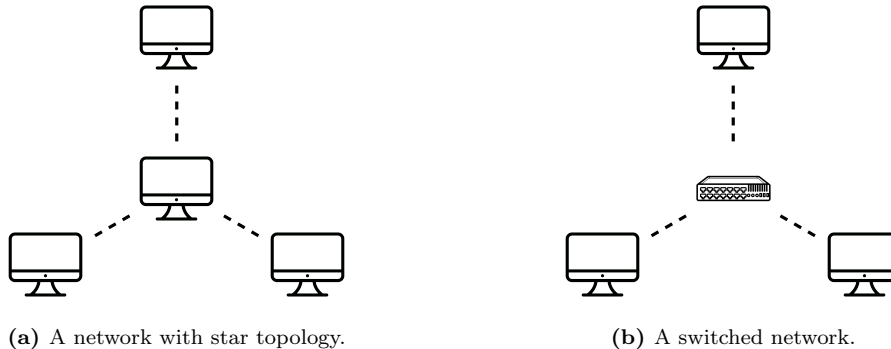


Figure 7.3

One limitation of a linear topology is the difficulty of adding new devices to the network. Every time a new computer is being set up, it would have to be inserted somewhere in the line, causing temporary disruptions in the network. A possible alternative is a bus topology, in which all computers are connected to a central bus line. This layout is shown in Figure 7.2.

With a bus topology, a device can be added to the network simply by linking it to the bus line. When a computer wants to send a message on the network, it specifies the destination address and sends the message to the bus. All other computers receive the message, since they are also on the bus. The computer for which the message was intended can then read the message.

A clear problem with the bus and linear topologies is that several computers on the network are involved whenever two computers have to communicate. In the bus topology in particular, every message is broadcast, so there can be competition for network resources when multiple devices try to send messages at once.

In order to alleviate these difficulties, a network can be arranged in a star topology. In a star topology, shown in Figure 7.3a all computers are connected to one central device. When a computer wants to send a message, it will send it to the central device, which then forwards it on to the destination address. This way every communication happens with no more than two hops across the network, and only the central device is involved in facilitating communications.

In a star topology, the central device places an important and unique role. For this reason, it is very often replaced with a specialized network device called a switch. We will discuss switches in the following section.

7.1.2 Switches

A *switch* is a device dedicated to moving packets through a network. It is designed to play the role of the central device in the star topology in Figure 7.3a. A typical switch may have dozens of Ethernet ports used to connect devices to it. The switch can receive packets from devices through these ports, and knows how to read the destination address of each packet. It forwards every packet along the cable pointing to the destination address, and nowhere else. This way, an Ethernet cable connected to a device is only every carrying packets sent from that device or intended for that device. A switched network is depicted in Figure 7.3b.

It is useful to think about a switch in terms of the network layers we covered in Section 6.2. A switch is typically an Internet layer device. The Internet layer runs on top of the link layer, and so switches have link layer technology as well: these are the

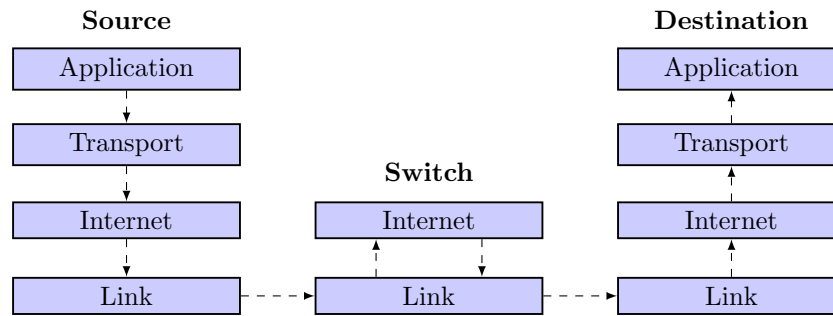


Figure 7.4: The flow of data through network layers in a switched network.

Ethernet ports and cables which allow devices to connect. Once a packet arrives at the switch, it uses the IP protocol in order to forward that packet to the appropriate device on the network. The flow of data from device to device follows the pattern in Figure 7.4.

Imagine a LAN with 50 devices connected to a switch. There are over a thousand different pairs of devices which may need to communicate with another, and the switch enables them to all communicate as if they had a dedicated cable between them, while actually using only 50 cables. This is the essential feature of any network: with relatively few physical connections, they enable an enormous number of potential communication channels. The Internet at large is designed in roughly the same way. We will discuss the details of the organization of the Internet, or more generally any wide area network, in the following section.

7.2 Wide Area Networks

Almost every LAN needs a way of communicating with the outside world. A network at a larger scale than a LAN is called a wide area network, or WAN. By far the most familiar example of a WAN is the Internet, but a WAN need not be this large. For example, a campus or metropolitan area may have its network organized as a WAN.

The link layer technology which carries signals over a WAN is similar to that of a LAN. A WAN typically needs higher capacity Ethernet cables, but the general principles are the same. The primary difference between a LAN and a WAN comes in at the Internet layer, where the IP protocol distinguishes between public and private networks. In this section we will discuss this difference, and the systems and protocols used to connect public networks together.

7.2.1 Gateways and Routing

If you have used a home network, you may be more familiar with routers than switches. This is because a switch is designed to connect all the devices into an area into a local network, but typically we also want to connect a local network to the outside Internet. This is the role of a router. Routers are responsible for connecting different networks together into a larger network.

In a typical setup, several devices on a local network will be connected to a router. The router assigns private IP addresses to each of these devices. Private IP addresses belong to one of three special blocks reserved for this purpose. The most common block is 192.168.*.*, any address ranging from 192.168.0.0 to 192.168.255.255. There are about 65,000 such addresses. A block like this can be written as 192.168.0.0/16, meaning the first 16 bits are fixed and the rest can change. The other private IP address ranges are 172.16.0.0/12, with about 1 million addresses, and 10.0.0.0/8, with about 16 million addresses.

When a device (such as a switch) within the LAN receives a packet, it will look at the destination IP address. If it knows where to route the packet – that is, if the packet is destined for another device on the local network – then it will forward the packet

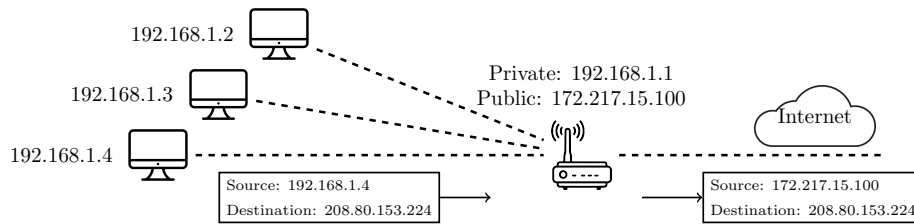


Figure 7.5: The default gateway forwards packets from a LAN to a WAN.

appropriately. If it does not know where to route the packet, then it will forward it on to the device designated as the *default gateway* for the network. The default gateway is typically the router.

When the router receives a packet with a destination address outside the local network, it knows how to forward it appropriately. First, it will change the source address of the packet. The packet originated from a device on the private network, so its source address is a private address, such as 192.168.1.4. The router will replace this source address with its own address on the WAN. This will be a public IP address, such as 172.217.15.100. It then sends the packet out to the public Internet and awaits a reply. When it receives a reply, it will forward the packet on to the correct device on the network. This process is depicted in Figure 7.5.

7.2.2 Autonomous Systems

Once a packet is routed to a wide area network, such as the Internet, it needs to be sent to the network which contains the destination address. Typically, the first destination for a packet is a router belonging to an Internet Service Provider (ISP). For example, in the `traceroute` trace on page 33, the packet ends up at a router owned by Verizon, a popular ISP. This router belongs to an *autonomous system* (AS) operated by Verizon. An autonomous system is simply a connection of routers on the Internet which are under some common administrative control, such as by an ISP.

Autonomous systems are managed globally using autonomous system numbers, or ASNs, which are similar to IP addresses. These numbers are allocated by the Internet Assigned Numbers Authority (IANA) to Regional Internet Registries (RIRs), which then assign the numbers to autonomous systems.

The routers in an autonomous system can forward packets to destinations within the same AS, but if a packet is destined for a network on another AS, it needs to be routed appropriately. This requires different ASs to communicate with one another. This can happen via *peering*, in which two ASs mutually agree to carry traffic to one another free of charge, or via *transit*, in which an AS charges a fee to carry traffic from other ASs.

When two ASs need to communicate with one another, they use the Border Gateway Protocol (BGP). This is the core routing protocol of the Internet, and is responsible for making decisions on how to get packets to a destination in the quickest way possible.

Control Structures

So far, we have seen programs in which each statement is executed in order, one by one. Today we will learn about *conditionals*, which allow us to execute statements depending on certain conditions. This is our first exposure to the idea of *control flow*, which refers to the order (or sequence) in which statements of a program are executed.

In this chapter, we will first learn about `if` statements, which allow us to write programs based on conditions. Then we will learn about `else` statements. We will combine these `if` and `else` statements into more complex nested structures, and then finally learn about `else if` statements. Lastly, this chapter ends with a list a common mistakes to avoid when writing conditionals.

8.1 The `if` statement

Definition 8.1.1. An *if statement* is a *conditional statement* that consists of the reserved keyword `if`, followed by a boolean expression enclosed in parentheses, followed by a statement typically enclosed in curly braces. If the boolean expression (or *condition*) evaluates to true, the statement is executed. Otherwise, it is skipped.

An `if` statement allows us to write programs that decide whether to execute a particular statement, based on a boolean expression which we call the *condition*. Below is an example of a simple `if` statement:

```
1  if (count > 20) {  
2      System.out.println("Count exceeded");  
3  }
```

The condition in this example is `count > 20`. It is a boolean expression that evaluates to either true or false. That is, `count` is either greater than 20 or not. If it is, “Count exceeded” is printed. Otherwise, the `println` statement is skipped.

Example 8.1.2: What does the following piece of code print if `x` is 101? What about if `x` is 200? What about if `x` is 8?

```
1  if (x > 100) {  
2      System.out.println("Big number!");  
3  }  
4  System.out.println("Hi there");  
5  if (x % 2 == 0) { // checks if a number is even*  
6      System.out.println("Even number!");  
7  }
```

* Recall that the modulo operator (%) computes the remainder of some number. By

checking if the remainder of some number when divided by 2 is 0, we are checking if that number is even.

Answer: If `x` is 101, the first boolean expression will evaluate to true, so the program will print “Big number!” Then the program will print “Hi there” regardless of any condition. Finally, since 101 is not even, the program will not print “Even number!”

Using the same reasoning to trace the execution with `x` as 200, the following statements will be printed: “Big number! Hi there! Even number!”

When `x` is 8, the output is “Hi there! Even number!”

Example 8.1.3: How would you fill in the boolean expression below to take the absolute value of an integer `x`? (Hint: to take an absolute value of a negative number, you must negate it).

```
1  if (/* Insert boolean expression here */) {
2      x = -x;
3  }
```

Answer: Insert the boolean expression `x < 0`.

8.2 The else statement

Sometimes we want to do one thing if a condition is true and another thing if that condition is false. We can add an **else** clause to an **if** statement to handle this kind of situation. Below is an example of a simple **if-else** statement:

```
1  if (price > 20) {
2      System.out.println("Too expensive");
3  } else {
4      System.out.println("Affordable");
5  }
```

This example prints either “Too expensive” or “Affordable”, depending on whether `price` is greater than 20. Only one or the other is ever executed, never both. This is because boolean conditions only evaluate to either true or false.

Note that it is not possible to have an **else** statement without having a corresponding **if** statement. Any **else** statement must be attached to a preceding **if** statement.

Example 8.2.1: What does the following piece of code do?

```
1  if (balance <= 0) {
2      System.out.println("Unable to withdraw");
3  } else {
4      System.out.println("Withdraw successful");
5  }
```

Answer: It prints “Unable to withdraw” if balance is 0 or lower. Otherwise, it prints “Withdraw successful”.

Example 8.2.2: Rewrite the code below using an **else** statement.

```
1  if (rating >= 4) {
```



```
2     System.out.println("Would recommend");
3 } if (rating < 4) {
4     System.out.println("Would not recommend")
5 }
```

Answer: Since the second boolean expression is the negation of the first boolean expression, we can replace the second if statement with an else.

```
1 if (rating >= 4) {
2     System.out.println("Would recommend");
3 } else {
4     System.out.println("Would not recommend")
5 }
```

Example 8.2.3: Does the following code compile?

```
1 else {
2     System.out.println("Hello");
3 }
```

Answer: No, since an else statement must have a corresponding if statement.

Example 8.2.4: Does the following code compile?

```
1 if (duration > 60) {
2     System.out.println("Sorry, that's too long");
3 } else {
4     System.out.println("I can fit that in my schedule");
5 } else {
6     System.out.println("Let me know when you are free")
7 }
```

Answer: No, since an if statement may only have at most 1 corresponding else statement.

8.3 Nested conditionals

It is possible to combine if and else statements in more interesting ways, by nesting them. Consider the following example:

```
1 if (x < y) {
2     System.out.println("x is less than y");
3 } else {
4     if (x > y) {
5         System.out.println("x is greater than y");
6     } else {
7         System.out.println("x is equal to y");
8     }
9 }
```

This example correctly prints out the relationship between two integer variables `x` and `y` using nested conditionals. No matter what the values of `x` and `y` are, only a single statement will be printed.

Example 8.3.1: Fill in each blank below with `if` or `else` to describe a person's height.

```

1  /*-----*/ (height < 60) {
2      System.out.println("Relatively short");
3  } /*-----*/ {
4      /*-----*/ (height > 72) {
5          System.out.println("Relatively tall");
6      } /*-----*/ {
7          System.out.println("Pretty average");
8      }
9  }

```

Answer: The four blanks should contain `if`, `else`, `if`, and `else`, in that order.

Example 8.3.2: Given integer variables `a` and `b` and the following piece of code, what value gets stored in variable `c`?

```

1  int c;
2  if (a > b) {
3      c = a;
4  } else {
5      if (b > a) {
6          c = b;
7      } else {
8          c = 0;
9      }
10 }

```

Answer: The larger of the two variables, `a` and `b`, gets stored in `c`. If `a` and `b` are equal, 0 gets stored in `c`.

Example 8.3.3: Imagine that the following program represents the person's reaction to different types of food. Given a pizza (`green` = false, `bitter` = false, `warm` = true, `cheesy` = true), how will this person respond? What about when given a kiwi (`green` = true, `bitter` = false, `warm` = false, `cheesy` = false)? What about a strawberry (`green` = false, `bitter` = false, `warm` = false, `cheesy` = false)?

```

1  if (green || bitter) {
2      System.out.println("No thank you");
3  } else {
4      if (warm && cheesy) {
5          System.out.println("Yum, yes please!");
6      } else {
7          System.out.println("Thank you, I'll take some");
8      }
9  }

```

Answer: The person responds "Yum, yes please!" to pizza, "No thank you" to kiwi, and "Thank you, I'll take some" to strawberry.

8.4 The else if statement

Nested conditionals are useful in many applications, but they can start feeling clunky as the nesting becomes deep. For example, consider the following piece of code:

```

1  if (grade > 90) {
2      System.out.println("You earned an A");
3  } else {
4      if (grade > 80) {
5          System.out.println("You earned a B");
6      } else {
7          if (grade > 70) {
8              System.out.println("You earned a C");
9          } else {
10             if (grade > 60) {
11                 System.out.println("You earned a D");
12             } else {
13                 System.out.println("You earned an F");
14             }
15         }
16     }
17 }

```

To avoid having such deeply nested conditionals, we introduce the `else if` clause. Using `else if` statements, we can rewrite the above program as:

```

1  if (grade > 90) {
2      System.out.println("You earned an A");
3  } else if (grade > 80) {
4      System.out.println("You earned a B");
5  } else if (grade > 70) {
6      System.out.println("You earned a C");
7  } else if (grade > 60) {
8      System.out.println("You earned a D");
9  } else {
10     System.out.println("You earned an F");
11 }

```

You can imagine an `else if` statement as shorthand for an `if` nested inside an `else`. For example, the following snippets of code behave identically (for any boolean expressions A and B):

```

1  // Version 1 (using nested conditionals)
2  if (A) {
3      System.out.println("A");
4  } else {
5      if (B) {
6          System.out.println("B");
7      } else {
8          System.out.println("C");
9      }
10 }
11
12 // Version 2 (equivalent, using else if)
13 if (A) {
14     System.out.println("A");
15 } else if (B) {
16     System.out.println("B");
17 } else {
18     System.out.println("C");
19 }

```

It is very common to see a conditional block of code that consists of 1 `if` statement, 0 or more `else if` statements, and 1 `else` statements. In these cases, recall that only one of these statements will ever execute.

Example 8.4.1: Imagine A and B are boolean expressions. If both of them evaluate to true, what does the following code print out?

```

1  if (A) {
2      System.out.println("apple");
3  } else if (B) {
4      System.out.println("banana");
5  } else {
6      System.out.println("carrot");
7  }
8  System.out.println("dragonfruit");

```

Answer: The code will print “apple” and then “dragonfruit”. Notice that “banana” does not get printed even though B is true since the entire `if`, `else if`, `else` chain will only ever print one of “apple”, “banana”, and “carrot”.

Example 8.4.2: Imagine X, Y, and Z are boolean expressions. If X evaluates to true, but Y and Z evaluate to false, what does the following code print out?

```

1  if (X && Y) {
2      System.out.println("xylophone");
3  } else if (!Z) {
4      System.out.println("zoo");
5  }

```

Answer: This code will print “zoo”. Notice that the boolean expressions used with `if` statements can be complex and contain boolean operators such as `&&` and `||` and `!`, as long as the expression evaluates to a true or false.

8.5 Curly braces

So far, we have been using curly braces to enclose each of our `if`, `else if`, and `else` statements. These curly braces can be omitted if there is only a single statement. For example, in the snippet of code below, the first set of curly braces can be omitted, while the second cannot.

```

1  if (guess == answer) {
2      System.out.println("You guessed correctly!");
3  } else {
4      System.out.println("Sorry, you guessed incorrectly.");
5      System.out.println("The answer was " + answer);
6  }

```

If all the curly braces were left out (as in the code below), the program would first either print “You guessed correctly!” or “Sorry, you guessed incorrectly.”, but then also print “The answer was ...” in all cases, regardless of the condition.

```

1  if (guess == answer)
2      System.out.println("You guessed correctly!");
3  else
4      System.out.println("Sorry, you guessed incorrectly.");
5      System.out.println("The answer was " + answer);
6      // ^^^ Careful! This is not part of the else clause!

```

Here it is important to note that whitespace and indentation are ignored by Java. Indentation has no effect on the behavior of a program. Proper indentation is extremely important for human readability. When used incorrectly, however, misleading indentation can result in unexpected behavior.

Example 8.5.1: Does the following code compile and print the larger of two integers **a** and **b** correctly?

```
1  if (a > b)      System.out.println("a");
2  else {
3  System.out.println("b");
4  }
```

Answer: Yes! Even though the whitespacing and indentation look funny, this piece of code compiles correctly.

8.6 Common mistakes

When writing conditional structures, beware of the following common mistakes:

8.6.1 Forgetting parentheses

In Java, the parentheses surrounding the boolean expression in conditional structures are required. For example, the following code will not compile.

```
1  if count > 10
2      System.out.println("So many!");
```

8.6.2 Accidental semi-colons

Accidentally semi-colons immediately after the parentheses around the boolean expression in an if statement is one of the trickiest mistakes to detect. The following code compiles, but it behaves unexpectedly.

```
1  if (count > 10);
2      System.out.println("So many!");
```

The code above is misleading because it is identical to the following:

```
1  if (count > 10) {
2      ;
3  }
4  System.out.println("So many!");
```

The accidental semi-colon is treated as a single, “do-nothing” statement.

8.6.3 Missing curly braces

As mentioned before, indentation helps make code more readable to humans, but it can also make code more confusing if used incorrectly. It is a common mistake to accidentally forget curly braces.

```
1  if (leaves = 4) {
2      System.out.println("You found a four-leaf clover, how lucky!");
3  } else
4      System.out.println("Sorry, not a four-leaf clover.");
5      System.out.println("Keep looking!");
```

This code above is misleading because it prints “Keep looking!” regardless of whether a four-leaf clover was found. The writer of this program probably meant to add curly braces around the `else` clause.

8.6.4 `else` without an `if`

Although the following example looks fine at first glance, it does not compile. This is because the `else` clause is nested *inside* the `if` statement instead of acting as an *alternative* option to the `if` statement.

```
1  if (chanceOfRain > 50) {
2      System.out.println("Bring an umbrella!");
3      else {
4          System.out.println("Hm, I don't think it will rain today.")
5      }
6  }
```

To fix the code above, move the `else` statement *outside* of the `if` as below:

```
1  if (chanceOfRain > 50) {
2      System.out.println("Bring an umbrella!");
3  } else {
4      System.out.println("Hm, I don't think it will rain today.")
5  }
```

8.6.5 Assignment v.s. equality operator

It is very easy to mix up the assignment operator (the single equals sign `=`) with the equality operator (the double equals sign `==`). The following code is incorrect and results in a compilation error:

```
1  if (temperature = 0) {
2      System.out.println("It's freezing!");
3  }
```

The assignment operator cannot be used here because conditional structures require a boolean expression (i.e. something that evaluates to either true or false).

Exercises

Exercise 8.6.0: What output is produced by the following code fragment?

```
1  int num = 87;
2  int max = 25;
3  if (num >= max*2)
4      System.out.println("apple");
5      System.out.println("orange");
6  System.out.println("pear");
```

TODO: modify this exercise, or make sure we can reference the Java Foundations textbook

Exercise 8.6.0: What output is produced by the following code fragment?

```
1  int limit = 100;
2  int num1 = 15;
3  int num2 = 40;
4  if (limit <= limit) {
5      if (num1 == num2)
6          System.out.println ("lemon");
7      System.out.println ("lime");
8  }
9  System.out.println ("grape");
```

TODO: modify this exercise, or make sure we can reference the Java Foundations textbook

Exercise 8.6.0: Given a day of the week encoded as 0=Sun, 1=Mon, 2=Tue, ...6=Sat, and a boolean indicating if we are on vacation, we need to decide when to set our alarm clock to ring. We want to wake up at “7:00” on weekdays and “10:00” on weekends. However, if we are on vacation, it should be “10:00” on weekdays” and “off” on weekends. Imagine you are given an integer variable `day` and a boolean variable `vacation`. Write a program that prints out the appropriate alarm time.

TODO: modify this exercise, or make sure we can reference codingbat.com

Exercise 8.6.0: Your cell phone rings. Normally you answer, except in the morning you only answer if it is your mom calling. In all cases, if you are asleep, you do not answer. Given three booleans, `isMorning`, `isMom`, and `isAsleep`, write a program that correctly prints out whether or not you will pick up your phone (“Answer” or “Do not answer”).

For example, if `isMorning`, `isMom`, and `isAsleep` are all false, print “Answer”.

TODO: modify this exercise, or make sure we can reference codingbat.com

Exercise 8.6.0: You have a lottery ticket with a 3-digit number. If your ticket has the number 777, you win \$100. If your ticket has a number with the same 3 digits (like 222, for example), then you win \$10. Otherwise, you don’t earn any money. Given three integers, `a`, `b`, and `c`, that represent the digits of your lottery ticket, write a program that prints out how much money you earn.

For example, if your ticket number is 123 (`a = 1`, `b = 2`, `c = 3`), print “\$0”.

TODO: modify this exercise, or make sure we can reference codingbat.com

References

Computer Science: An Interdisciplinary Approach, Robert Sedgewick and Kevin Wayne.

Lewis, John, Peter DePasquale, and Joseph Chase. Java Foundations: Introduction to Program Design and Data Structures. Addison-Wesley Publishing Company, 2010.

Programming exercises from codingbat.com

Business Software

9.1 Introduction

This lecture will go over three primary types of business productivity software most businesses use on a daily basis:

1. Word processing software: for typing and formatting text and images on a page (such as essays and reports)
2. Presentation software: for creating slides with information and graphics (such as for meetings and briefings)
3. Spreadsheet software: for storing data, doing computations, and making graphs

9.2 Word Processing

Word processing software is used for creating, editing, and formatting documents. Many tech companies have their own version of word processing software. A few more familiar examples include

- Microsoft Word
- Apple Pages
- Google Docs

Beyond providing an interface to type the words of your report or essay into, word processors allow a variety of styling and formatting options to customize the document to look the way you like it. For example, all of the word processors above support

- Automatic page numbering
- Adjustable margins
- Varieties of text size, text font, text color
- Addition of images and captions, along with image formatting
- Addition of tables and lists
- Spelling and grammar checking

9.3 Spreadsheets

Spreadsheet software is used for dealing with data. A few examples of spreadsheet software include

- Microsoft Excel
- Apple Numbers
- Google Sheets

Beyond providing a grid in which you can type numbers, spreadsheet software facilitate a variety of operations on the data. styling. A few key features supported by most spreadsheet applications include

- Statistical analysis
- Sorting
- Creating charts and graphs
- Formatting (e.g. adding color, borders, etc.)

9.4 Presentations

Presentation software is used for presenting data and graphics to an audience. Some common examples include

- Microsoft Powerpoint
- Apple Keynote
- Google Slides

All of these software give the user a series of "slides" that can be clicked through during a presentation. It allows the presenter to focus the audience's attention on a single point at a time. The slides are not meant to contain all the information that the presenter wants to get across. Instead, the slides only mention the key takeaways, and the presenter should verbally fill in the details. For that reason, slides often consist of just a few bullet points and an image. Of course, there are a variety of creative ways to present - this is just the most common format!

A few particularly important features offered by almost all presentation software include

- Automatic styling (color schemes, layouts, etc)
- Add images and videos
- Slide transitions
- Animations (e.g. showing one line of text on a slide at a time)

9.5 Project

Methods

Sometimes it is useful to abstract certain computations into methods. In this chapter, we will be looking at how and why we can use methods to achieve more concise and easily-maintainable code.

We will first look at some code that does not use methods and identify some problems with this code. We will then, in Section 10.1, introduce what a method is, including how to define and use them. In Sections 10.2 and 10.3, we will then talk about how we can use methods to improve the problems that we noticed with our code. In particular, in Section 10.2, we talk about how methods allow us to abstract some of our code to get more concise and maintainable code, and in Section 10.3, we talk about how methods allow us to divide our code up into chunks that we can consider in isolation, allowing us to manage and maintain our code more easily.

Let us consider the following problem:

We have a class of nine students, split up into three groups of three. For each group, we want to report the maximum score for each group. We will use the variable names `group1s1`, `group1s2`, and `group1s3` to refer to the three scores achieved by the students in Group 1, and we will use a similar naming convention for the variable names for the scores achieved by students in Group 2 and Group 3.

The following snippet of code computes the maximum score in Group 1, storing the result in `group1Max`.

```
1  int group1Max = group1s1;
2  if (group1s2 > group1Max) {
3      group1Max = group1s2;
4  }
5  if (group1s3 > group1Max) {
6      group1Max = group1s3;
7  }
```

If we now wish to extend this to compute the maximum score for each group, we can copy the code and rename the variables to achieve the following code:

```
1  int group1Max = group1s1;
2  if (group1s2 > group1Max) {
3      group1Max = group1s2;
4  }
5  if (group1s3 > group1Max) {
6      group1Max = group1s3;
7  }
8
9  int group2Max = group2s1;
10 if (group2s2 > group2Max) {
11     group2Max = group2s2;
12 }
13 if (group2s3 > group2Max) {
```

```

14     group2Max = group2s3;
15 }
16
17 int group3Max = group3s1;
18 if (group3s2 > group3Max) {
19     group3Max = group3s2;
20 }
21 if (group3s3 > group3Max) {
22     group3Max = group3s3;
23 }

```

There are a couple of things about the process of getting this resulting code that aren't particularly ideal. The first is that it is very easy to make a mistake when copying the code and renaming the variables. For example, we might have missed renaming `group1Max` to `group3Max` in the computation of the maximum score for Group 3, leading us to compute an incorrect maximum value in some cases.

The second issue is that there is a lot of redundancy. If we look at the code, we can notice that there seem to be a lot of redundant parts across the code for the three different groups. If we had ended up with the wrong computation initially and needed to fix it later, then we would have to make sure we fixed it everywhere, leading to repeated work and the possibility that we didn't fix one of the copies. For example, if we initially started out with `group1s3 < group1Max` instead of `group1s3 > group1Max`, and then copied and renamed variables to compute the maximum scores for the other groups, we would end up with the code below:

```

1  int group1Max = group1s1;
2  if (group1s2 > group1Max) {
3      group1Max = group1s2;
4  }
5  if (group1s3 < group1Max) {
6      group1Max = group1s3;
7  }
8
9  int group2Max = group2s1;
10 if (group2s2 > group1Max) {
11     group2Max = group2s2;
12 }
13 if (group2s3 < group2Max) {
14     group2Max = group2s3;
15 }
16
17 int group3Max = group3s1;
18 if (group3s2 > group3Max) {
19     group3Max = group3s2;
20 }
21 if (group3s3 < group3Max) {
22     group3Max = group3s3;
23 }

```

To fix our code, we would not only have to replace `group1s3 < group1Max` with `group1s3 > group1Max`, but we would also need to replace `group2s3 < group2Max` with `group2s3 > group2Max` and `group3s3 < group3Max` with `group3s3 > group3Max`.

10.1 Methods

In Java, we can use abstractions in our code in the form of *methods*. Methods contain bits of code that may compute things using *arguments* that are passed into the method.

An example of the syntax for a *method definition* is given below:

```

1  public bool negate(bool arg) {
2      return !arg; // method body
3  }

```

We can identify the different components of the method definition:

- The *access modifier* of the method is **public**,
- the **bool** before **negate** gives the *return type* of the method
- **negate** is the *name of the method*,
- **arg** is the *name of the first argument/parameter* of the method,
- the **bool** before **arg** it is the *type* of the **arg** parameter,
- and the code between the curly braces is the method *body*.

Inside the method body, there is code that may use the arguments of the method. In the example above, the body only consists of the **return** statement. The **return** statement is followed by something of the return type of the method. This expression, after being evaluated, gives the *return value* of the method. As soon as a **return** statement is encountered, a method finishes its execution. *Access modifiers* will be described later on when we talk about Classes, but note that if an access modifier is omitted, it the method has the *default* modifier.

Other important terms to note are as follows: The *parameter list* of a method consists of the types of its parameters in the order they are given in the method definition. For the above method, it is simply **bool**. The *method signature* consists of the method name and parameter list. For the above method, it is **negate(bool)**. A method is identified by its method signature.

Example 10.1.1: Consider the following method:

```
1  int squareSum(int arg0, int arg1) {
2      int res = arg0 + arg1;
3      return res * res;
4  }
```

Give the following for the method:

- Access modifier
- Return type
- Arguments
- Method signature

Answer:

- Access modifier: default
- Return type: **int**
- Arguments: **arg0, arg1**
- Method signature: **squareSum(int, int)**

Example 10.1.2: Consider the following method:

```
1  private void mystery (int age, String name) {
2      if (age >= 18) {
3          System.out.println(name);
4      }
5  }
```

Give the following for the method:

- Access modifier
- Return type
- Name
- Parameter list

Answer:

- Access modifier: `private`
- Return type: `void`
- Name: `mystery`
- Parameter list: `int, String`

We may *call* or *invoke* the above `negate` method from elsewhere in the code by, for example, using `negate(true)` if we want to call `negate` with argument `true`. We must specify the name of the method that we wish to call, followed by the arguments that we wish to call it on, separated by commas if there are more than one. The arguments provided in the call must correspond to the types of the arguments specified in the method definition: if the method definition has the type `bool` for its first argument, then the first argument supplied in the call should also be of type `bool`, and if the definition has the type `int` for its second argument, then the first argument in the call should also be of type `int`, and so on. A method invocation evaluates to the return value it computes, where this computation happens with the actual arguments given in the call substituted for the arguments specified in the definition, so in this case, the invocation `id(true)` evaluates to `false`. A method invocation has the same type as its return type, which in this case, is `bool`.

So, for example, we can do something like in following code snippet:

```
1  bool f = negate(true);
2  bool t = negate(negate(true));
```

After executing the above code, the variable `f` contains the value `false` and the variable `t` contains the value `true`.

Note that we can also have methods that do not return anything. Such methods have `void` return types. They may contain a `return` statement with no expression following `return`. Calls to these methods cannot be used inside of other expressions. Two examples of methods with `void` return types are below:

```
1  public void printSum(double a, double b) {
2      System.out.println(a + b);
3  }
```

```
1  public void printBigger(double a, double b) {
2      if (a > b) {
3          System.out.println(a);
4          return;
5      }
6      System.out.println(b);
7  }
```

Note that the latter example only ever prints the value of either `a` or `b` but never both because the `return` statement finishes the execution of the method.

Some methods also do not take any arguments at all. One such example is below:

```

1 public int constantOne() {
2     return 1;
3 }

```

This method can be called using `constantOne()`.

Example 10.1.3: What value does `res` have after the following code snippet is executed?

```

1 int res = constantOne() + constantOne();

```

Answer: The value of `res` is 2.

Note that for every method that does not have a `void` return type, for every possible control-flow path through the method body, there must be a `return` statement that is eventually reached. For example, the following method does not meet this requirement because when `arg` is `false`, the `return` statement is not reached:

```

1 public int twoIfTrue(bool arg) {
2     if (arg) {
3         return 2;
4     }
5 }

```

Example 10.1.4: Consider the following method:

```

1 public void mystery (bool arg0, int arg1) {
2     if (arg0) {
3         return -1 * arg1;
4     }
5     return arg0;
6 }

```

What do each of the following evaluate to?

- `mystery(true, -2)`
- `mystery(true, 0)`
- `mystery(true, 4)`
- `mystery(false, -2)`
- `mystery(false, 0)`
- `mystery(false, 4)`

Answer:

- `mystery(true, -2)` evaluates to 2
- `mystery(true, 0)` evaluates to 0
- `mystery(true, 4)` evaluates to -4
- `mystery(false, -2)` evaluates to -2
- `mystery(false, 0)` evaluates to 0

- `mystery(false, 4)` evaluates to 4

What does the method do?

Answer: It negates its second argument `arg1` whenever its first argument `arg0` is `true` and otherwise just returns its second argument `arg1`.

Example 10.1.5: How many methods can be called by the following method? What do you think each of their return types and parameter lists are?

```

1 public void sayHello(int times, String name, String day) {
2     if (shouldStop(times))
3         return;
4     String hello = getHelloStr(name, day);
5     System.out.println(hello);
6     sayHello(times - 1, name, day);
7 }

```

Answer: Three methods can be called: `shouldStop`, `getHelloStr` and `sayHello`.

The return type of `shouldStop` is `bool`, and it has parameter list `int`. The return type of `getHelloStr` is `String`, and it has parameter list `String, String`. The return type of `sayHello` is `void`, and it has parameter list `int, String, String`.

Example 10.1.6: Write a method with signature `printInOrder(int, int, int)` that takes three `ints` `a`, `b`, and `c` as arguments and prints them out, one per line, in increasing order. E.g. `printInOrder(5, 9, 7)` should output the following:

```

1 5
2 7
3 9

```

Answer: A possible solution is the following:

```

1 public void printInOrder(int a, int b, int c) {
2     int smallest = a;
3     int middle = b;
4     int biggest = c;
5     if (smallest > middle) {
6         int tmp = smallest;
7         smallest = middle;
8         middle = tmp;
9     }
10    if (smallest > biggest) {
11        int tmp = smallest;
12        smallest = biggest;
13        biggest = tmp;
14    }
15    if (middle > biggest) {
16        int tmp = middle;
17        middle = biggest;
18        biggest = tmp;
19    }
20    System.out.println(smallest);
21    System.out.println(middle);
22    System.out.println(biggest);
23 }

```

10.1.1 Overloading

As mentioned above, methods are identified by their signature, not just their names. Because of this, we can actually have methods with the same name (in the same Class) as long as they have different parameter lists.

For example, we can have the following two methods because one has signature `printBigger(double, double)` and the other has signature `printBigger(int, int)`:

```
1 public void printBigger(double a, double b) {
2     if (a > b) {
3         System.out.println(a);
4         return;
5     }
6     System.out.println(b);
7 }
8 public void printBigger(int a, int b) {
9     if (a > b) {
10        System.out.println(a);
11        return;
12    }
13    System.out.println(b);
14 }
```

In this case, we say that the `printBigger` method is *overloaded*.

However, we cannot have the following two methods because they both have the same signature `twoIfTrue(bool)`:

```
1 public int twoIfTrue(bool arg) {
2     if (arg) {
3         return 2;
4     }
5 }
6 public double twoIfTrue(bool arg) {
7     if (arg) {
8         return 2.0;
9     }
10 }
```

Example 10.1.7: Consider the following two methods that we have seen in previous examples:

```
1 private void mystery (int age, String name) {
2     if (age >= 18) {
3         System.out.println(name);
4     }
5 }
6 public void mystery (bool arg0, int arg1) {
7     if (arg0) {
8         return -1 * arg1;
9     }
10    return arg0;
11 }
```

Are we allowed to have both of these methods together (in the same Class)? Why or why not?

Answer: We are allowed to have both of these methods together because they have different signatures. One has signature `mystery(int, String)` and the other has signature `mystery(bool, int)`.

Example 10.1.8: Which of the following are valid overloadings?

A:

```

1 public int absValue(int x) {
2     if (x < 0) {
3         return x * -1;
4     }
5     return x;
6 }
7 public double absValue(int x) {
8     if (x < 0) {
9         return 0.0 - x;
10    }
11    return x - 0.0;
12 }
```

B:

```

1 public int absValue(int x) {
2     if (x < 0) {
3         return x * -1;
4     }
5     return x;
6 }
7 public double absValue(double x) {
8     if (x < 0) {
9         x = 0.0 - x;
10    }
11    return x;
12 }
```

C:

```

1 public int absValue(int x) {
2     if (x < 0) {
3         return x * -1;
4     }
5     return x;
6 }
7 public int absValueX(int x) {
8     if (x < 0) {
9         return x * -1;
10    }
11    return x;
12 }
```

Answer:

- A is not a valid overloading because both methods have the same signature.
- B is a valid overloading because both methods have different signatures.
- C is valid code, but this is not a valid overloading. It is not overloading at all because the methods have different names.

10.1.2 Methods Summary

In summary, every method definition has the following:

- An access modifier (if not explicitly stated, this will be the *default* modifier)
- A return type (possibly `void`)

- A name
- A (possibly empty) list of parameter types and names that can be passed to the method when it is called
- A body
- A return statement reached on every control path in the body, where each `return` is followed by an expression of the method's return type (unless the return type is `void`)

Each method is identified by its signature, so you may have methods with the same name but different parameter lists.

10.2 Abstraction

Let us now return to our problem of finding the maximum scores for each of the three groups. The code for each of the three groups, though similar, have different variable names; however, this is the only way in which they are different. The process of distilling the similarity among different pieces of code is the process of *abstraction*. In this case, we would like to *abstract* away the details of each copy to achieve a piece of code that describes all of their behavior.

After abstracting away the variable names, the code for each group is such that, if we use the appropriate variable names for the students in the group instead of `x`, `y`, and `z`, that the following code snippet would describe all of the three different operations, with `max` storing the desired result:

```

1  int max = x;
2  if (y > max) {
3      max = y;
4  }
5  if (z > max) {
6      max = z;
7  }
```

We can wrap this code snippet up in a method definition that returns the value we care about:

```

1  public int max3(int x, int y, int z) {
2      int max = x;
3      if (y > max) {
4          max = y;
5      }
6      if (z > max) {
7          max = z;
8      }
9      return max;
10 }
```

If we use this method in place of the redundant code, our original code that compute the maximum scores for each of the three groups then becomes as follows:

```

1  int group1Max = max3(group1s1, group1s2, group1s3);
2  int group2Max = max3(group2s1, group2s2, group2s3);
3  int group3Max = max3(group3s1, group3s2, group3s3);
```

This resulting code is much more concise and is an example of *code reuse*, where the same exact code (i.e. the code inside the body of `max3`) is being reused in several places. This code is also perhaps easier to understand if you know that `max3` simply calculates the maximum of its three arguments. In the original code, after going through the code

for one of the groups and realizing that it computes a maximum, you would have to go through the calculations of the other groups to make sure that they also compute a maximum. Here, it is easy to see that the same computation is happening for each group (but with different inputs).

We might further notice that both of the conditional statements in the `max3` method compute very similar things (i.e. the maximum of two numbers), and might perform further abstraction to achieve the following method:

```

1  int max2(int x, int y) {
2      int max = x;
3      if (y > max) {
4          max = y;
5      }
6  }
```

We can then adjust our `max3` method to call this one:

```

1  int max3(int x, int y, int z) {
2      return max2(max2(x, y), z);
3  }
```

Note that we do not need to change our code for computing `group1Max`, `group2Max`, or `group3Max` in order to use the method `max2`. The changes to `max3` are sufficient.

10.3 Modularity

Let us again consider the case in which our calculation of the maximum score is incorrect because we used `<` instead of `>` in the last comparison. Let us assume that we are still using the version of `max3` that does not call `max2`, but we mistakenly have the comparison `z < max` instead of `z > max` in our method:

```

1  public int max3(int x, int y, int z) {
2      int max = x;
3      if (y > max) {
4          max = y;
5      }
6      if (z < max) {
7          max = z;
8      }
9      return max;
10 }
```

In order to fix, this we simply need to change `z < max` to `z > max` once. The code that calculates the values of `group1Max`, `group2Max`, and `group3Max` by calling `max3` is fixed by this one change because for all the groups, we call the same method. Here we have achieved a *separation of concerns* in the two parts of our code:

- We have one part of our code (the `max3` method) that is concerned with finding the maximum of three *arbitrary* numbers, but it is not concerned with *which* numbers specifically for which it is finding the maximum.
- We have another part that is concerned with computing the maximum scores for each group, *given that* we have some other code that can calculate the maximum of three numbers, but it is not concerned with *how* this maximum is found, as long as it is done correctly

This separation of concerns is referred to as *modularity*, and we can regard the code inside of different methods as being in different *modules*. We have already seen that modularity can allow us to do things like fix bugs in one part of our code without having to touch other parts of our code.

Modularity also lets us do other things, like naturally divide up labor. If you wanted to split up the work of writing a program that finds the maximum of each group's three scores with a friend, you might volunteer to write the code that deals with finding the maximum scores for each group provided that your friend writes code that finds the maximum of three numbers contained in a method with signature `max3(int, int, int)` that returns an `int` that is the maximum of its three `int` arguments.

Given that you know your friend will write such a `max3` method, you can, without seeing your friend's code, write the following code (that we have already seen):

```
1 int group1Max = max3(group1s1, group1s2, group1s3);
2 int group2Max = max3(group2s1, group2s2, group2s3);
3 int group3Max = max3(group3s1, group3s2, group3s3);
```

Once your friend finishes writing `max3`, then your code should work as expected.

An advantage of modularity is that it is not necessary to know the details of how the modules that you use are implemented nor used. In the above example, you do not need to know how exactly your friend implements `max3` (maybe it calls `max2` and maybe it does not) and your friend does not need to know how you use `max3` (maybe you also use it to find the maximum of the three maximum scores and maybe you do not).

Example 10.3.1: What would you add to the following code to calculate the maximum of `group1Max`, `group2Max`, and `group3Max` and store it in `allGroupMax`?

```
1 int group1Max = max3(group1s1, group1s2, group1s3);
2 int group2Max = max3(group2s1, group2s2, group2s3);
3 int group3Max = max3(group3s1, group3s2, group3s3);
```

Answer:

```
1 int group1Max = max3(group1s1, group1s2, group1s3);
2 int group2Max = max3(group2s1, group2s2, group2s3);
3 int group3Max = max3(group3s1, group3s2, group3s3);
4 // added code below:
5 int allGroupMax = max3(group1Max, group2Max, group3Max);
```

Example 10.3.2: Consider the case where not only do you want to compute the maximum scores for each group for Group 1, Group 2, and Group 3, but you also want to do the same for Group 4 and Group 5. Unfortunately, while Group 4 and Group 5 also have three scores per group (`group4s1`, `group4s2`, `group4s3`, `group5s1`, `group5s2`, and `group5s3`), these scores are all `doubles` rather than `ints`.

You want write the following code:

```
1 int group1Max = max3(group1s1, group1s2, group1s3);
2 int group2Max = max3(group2s1, group2s2, group2s3);
3 int group3Max = max3(group3s1, group3s2, group3s3);
4 double group4Max = max3(group4s1, group4s2, group4s3);
5 double group5Max = max3(group5s1, group5s2, group5s3);
```

What code do you ask your friend to write so that your code works?

Answer: You can ask your friend to write an overloaded `max3`: one overloading (the one we have seen so far) should return an `int` and have method signature `max3(int, int, int)` and the other should return a `double` and have method signature `max3(double,`

`double, double)`). Both versions of `max3` should return the maximum value of their three arguments.

Example 10.3.3: Give the signatures of methods that you would need to write to make the following code work:

```

1  int computeAverages(int a, int b, int c, int d, int e, int f, int g, int
   h) {
2      int avg0 = sum(a, b) / 2;
3      int avg1 = sum(c, d, e) / 3;
4      int avg2 = sum(f, g, h) / 3;
5      int avgAvg = sum(avg0, avg1, avg2) / 3;
6      return avgAvg;
7  }
```

Answer: `sum(int, int)` and `sum(int, int, int)`

Example 10.3.4: This is the same code as in the previous example, but someone made a mistake and divided by 4 instead of 3 when taking the average of three numbers:

```

1  int computeAverages(int a, int b, int c, int d, int e, int f, int g, int
   h) {
2      int avg0 = sum(a, b) / 2;
3      int avg1 = sum(c, d, e) / 4;
4      int avg2 = sum(f, g, h) / 4;
5      int avgAvg = sum(avg0, avg1, avg2) / 4;
6      return avgAvg;
7  }
```

Improve the original code by using abstraction.

Answer:

```

1  int avg(int a, int b) {
2      return sum(a, b) / 2;
3  }
4  int avg(int a, int b, int c) {
5      return sum(a, b, c) / 3;
6  }
7  int computeAverages(int a, int b, int c, int d, int e, int f, int g, int
   h) {
8      int avg0 = avg(a, b);
9      int avg1 = avg(c, d, e);
10     int avg2 = avg(f, g, h);
11     int avgAvg = sum(avg0, avg1, avg2) / 4;
12     return avgAvg;
13 }
```

Arrays

Consider this snippet of code:

```

1  if      (day == 0) System.out.println("Monday");
2  else if (day == 1) System.out.println("Tuesday");
3  else if (day == 2) System.out.println("Wednesday");
4  else if (day == 3) System.out.println("Thursday");
5  else if (day == 4) System.out.println("Friday");
6  else if (day == 5) System.out.println("Saturday");
7  else if (day == 6) System.out.println("Sunday");

```

What does this code do? It prints the day of the week after conditioning on the value of an integer `day`. But this code is repetitive. It would be useful if we had some way of creating a list of days of the week, and then just specifying which of those days we wanted to print. Something like this:

```

1  System.out.println(DAYS_OF_WEEK[day]);

```

To achieve this in Java, we need arrays.

Definition 11.0.1. An *array* is an ordered and fixed-length list of values that are of the same type. We can access data in an array by *indexing*, which means referring to specific values in the array by number. If an array has n values, then we think of it as being numbered from 0 to $n-1$.

To *loop* or *iterate* over an array means that our program accesses every value in the array, typically in order. For example, if we looped over the array in the diagram, that would mean that we looked at the value at the 0th index, then the value at the 1st index, then the value at the 2nd index, and so on.

When we say that the array is "ordered" is that the relationship between an index and its stored value is unchanged (unless we explicitly modify it). If we loop over an unchanged array multiple times, we will always access the same values.

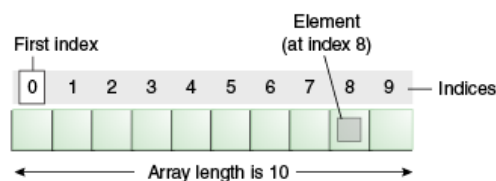


Figure 11.1: Diagram of an array (Credit: <https://www.geeksforgeeks.org/arrays-in-java/>)

Arrays are *fixed-length*, meaning that after we have created an array, we cannot change its length. We will see in the next chapter [TK: confirm] that `ArrayLists` are an array-like data structure that allows for changing lengths.

Finally, all the values in an array must be of the same type. For example, an array can hold all floating point numbers or all characters or all strings. But an array cannot hold values of different types.

11.1 Creating arrays

The syntax for creating an array in Java has three parts:

1. Array type
2. Array name
3. Either: array size or specific values

For example, this code creates an array of size `n = 10` and fills it with all `0.0s`

```
1 double[] arr;           // Declare array
2 arr = new double[n];    // Initialize the array
3 for (int i = 0; i < n; i++) { // Iterate over array
4     arr[i] = 0.0;        // Initialize elements to 0.0
5 }
```

The key steps are: we first declare and initialize the array. We then loop over the array to initialize specific values. We can also initialize the array at compile time, for example

```
1 String[] DAYS_OF_WEEK = {
2     // Indices:
3     // 0      1      2      3      4      5      6
4     "Mon", "Tue", "Wed", "Thu", "Fri", "Sat", "Sun"
5 };
```

Notice the difference in syntax. When creating an empty array, we must specify a size. When initialize an array at compile time with specific values, the size is implicit in the number of values provided.

Finally, in Java, it is acceptable to move the brackets to directly after the type declaration to directly after the name declaration. For example, these two declarations are equivalent:

```
1 int arr[];
2 int[] arr;
```

11.2 Indexing

Consider the array `DAYS_OF_WEEK` from the previous section. We can *index* the array using the following syntax:

```
1 System.out.println(DAYS_OF_WEEK[3]); // Prints "Thu"
```

In Java, array's are said to use *zero-based indexing* because the first element in the array is accessed with the number 0 rather than 1.

Example 11.2.1: What does `System.out.println(DAYS_OF_WEEK[1]);` print?

Example 11.2.2: What does this code do? What number does it print?

```

1  double sum = 0.0;
2  double[] arr = { 1, 2, 2, 3, 4, 7, 9 }
3  for (int i = 0; i < arr.length; i++) {
4      sum += arr[i];
5  }
6  System.out.println(sum / arr.length);

```

11.3 Array length

As mentioned previously, arrays are *fixed-length*. After you have created an array, its length is unchangeable. You can access the length of an array `arr[]` with the code `arr.length`.

Example 11.3.1: What does `System.out.println(DAYS_OF_WEEK.length);` print?

Example 11.3.2: Write a `for` loop to print the days of the week in order (Monday through Sunday) using an array rather than seven `System.out.println` function calls.

11.4 Default initialization

In Java, the default initial values for numeric primitive types is 0 and `false` for the `boolean` type.

Example 11.4.1: Consider this code from earlier:

```

1  double[] arr;
2  arr = new double[n];
3  for (int i = 0; i < n; i++) {
4      arr[i] = 0.0;
5  }

```

Rewrite this code to be a single line.

11.5 Bounds checking

Consider this snippet of code.

Example 11.5.1: Where is the bug?

```

1  int[] arr = new int[100];
2  for (int i = 0; i <= 100; ++i) {
3      System.out.println(arr[i]);

```



```
4 }
```

The issue is that the program attempts to access the value `arr[100]`, while the last element in the array is `arr[99]`.

This kind of bug is called an “off-by-one error” and is so common it has a name. In general, an off-by-one-error is one in which a loop iterates one time too many or too few.

Example 11.5.2: Where is the off-by-one-error?

```
1 int[] arr = new int[100];
2 for (int i = 0; i < array.length; i++) {
3     arr[i] = i;
4 }
5 for (int i = 100; i > 0; --i) {
6     System.out.println(arr[i]);
7 }
```

Example 11.5.3: Fill in the missing code in this `for` loop to print the numbers in reverse order, i.e. 5, 4, 3, 2, 1:

```
1 int[] arr = { 1, 2, 3, 4, 5 };
2 for (???) {
3     System.out.println(arr[i]);
4 }
```

11.6 Empty arrays

This code prints five values, one per line, but we never specified which values. What do you think it prints?

```
1 int[] arr = new int[5];
2 for (int i = 0; i < arr.length; i++) {
3     System.out.println(arr[i]);
4 }
```

In Java, an uninitialized or empty array is given a default value:

- For `int`, `short`, `byte`, or `long`, the default value is 0.
- For `float` or `double`, the default value is 0.0
- For `boolean` values, the default value is `false`.
- For `char`, the default value is the null character `'\0000'`.

Note that an array can be partially initialized.

Example 11.6.1: What does this code print?

```
1 char[] alphabet = new char[26];
2 alphabet[0] = 'a';
3 alphabet[1] = 'b';
```

```

4  for (int i = 0; i < alphabet.length; i++) {
5      System.out.println(alphabet[i]);
6  }

```

11.7 Enhanced for loop

So far, we have seen how to iterate over arrays by indexing each element with a number:

```

1  char[] vowels = {'a', 'e', 'i', 'o', 'u'};
2  for (int i = 0; i < vowels.length; ++ i) {
3      System.out.println(vowels[i]);
4  }

```

We can perform the same iteration without using indices using an “enhanced for loop” or for-each loops:

```

1  char[] vowels = {'a', 'e', 'i', 'o', 'u'};
2  for (char item: vowels) {
3      System.out.println(item);
4  }

```

11.8 Exchanging and shuffling

Two common tasks when manipulating arrays are *exchanging two values* and *shuffling* values. (*Sorting* is more complicated and will be address later.)

To exchange to values, consider the following code:

```

1  double[] arr = { 1.0, 2.0, 3.0, 4.0, 5.0, 6.0 };
2  int i = 1;
3  int j = 4;
4  double tmp = arr[i];
5  arr[i] = arr[j];
6  arr[j] = tmp;

```

Example 11.8.1: What are the six values in the array, in order?

To shuffle the array, consider the following code:

```

1  int n = arr.length;
2  for (int i = 0; i < n; i++) {
3      int r = i + (int) (Math.random() * (n-i));
4      String tmp = arr[r];
5      arr[r] = arr[i];
6      arr[i] = tmp;
7  }

```

Example 11.8.2: What does this code do:

```

1  for (int i = 0; i < n/2; i++) {
2      double tmp = arr[i];
3      arr[i] = arr[n-1-i];

```

```

4     arr[n-i-1] = tmp;
5 }

```

Exercises

Exercise 11.8.2: Write a program that reverses the order of values in an array.

Exercise 11.8.2: What is wrong with this code snippet?

```

1  int[] arr;
2  for (int i = 0; i < 10; i++) {
3      arr[i] = i;
4  }

```

Exercise 11.8.2: Rewrite this snippet using an enhanced **for-each** loop (for now, it is okay to re-define the array):

```

1  char[] vowels = {'a', 'e', 'i', 'o', 'u'};
2  for (int i = array.length; i >= 0; i--) {
3      char letter = vowels[i];
4      System.out.println(letter);
5  }

```

Exercise 11.8.2: Write a program that uses **for** loops to print the following pattern:

```

1  1*****
2
3  12*****
4
5  123*****
6
7  1234*****
8
9  12345*****
10
11 123456***
12
13 1234567**
14
15 12345678*
16
17 123456789

```

Exercise 11.8.2: Write a program `HowMany.java` that takes an arbitrary number of command line arguments and prints how many there are.

References

Computer Science: An Interdisciplinary Approach, Robert Sedgewick and Kevin Wayne.

12

ArrayLists

A *collection* is a group of objects. Today, we'll be looking at a very useful collection, the **ArrayList**. A *list* is an ordered collection, and an **ArrayList** is one type of list. We will see at the end of this chapter how an **ArrayList** is different from an **array**.

12.1 Creating an ArrayList

Let's create a class **NameTracker** and follow along in it. Before we can use an **ArrayList**, we have to import it:

```
1  import java.util.ArrayList;
2
3  public class NameTracker {
4
5      public static void main(String args[]) {
6      }
7  }
```

Next, we call the constructor; but we have to declare the type of object the **ArrayList** is going to hold. This is how you create a new **ArrayList** holding **String** objects.

```
1  ArrayList<String> names = new ArrayList<String>();
```

Notice the word “String” in angle brackets. This is the Java syntax for constructing an **ArrayList** of **String** objects.

12.2 add()

We can add a new **String** to **names** using the **add()** method.

```
1  names.add("Ana");
```

Example 12.2.1: Exercise: Write a program that asks the user for some names and then stores them in an **ArrayList**. Here is an example program:

```
Please give me some names:
Sam
Alecia
Trey
```

```
Enrique
Dave

Your name(s) are saved!
```

We can see how many objects are in our ArrayList using the `size()` method.

```
1 System.out.println(names.size()); // 5
```

Example 12.2.2: Modify your program to notify the user how many words they have added.

```
Please give me some names:
Mary
Judah

Your 2 name(s) are saved!
```

12.3 `get()`

Remember how the `String.charAt()` method returns the `char` at a particular index? We can do the same with names. Just call `get()`:

```
1 names.add("Noah");
2 names.add("Jeremiah");
3 names.add("Ezekiel");
4 System.out.println(names.get(2)); // 'Ezekiel'
```

Example 12.3.1: Update your program to repeat the names back to the user in reverse order. Your solution should use a for loop and the `size()` method. For example:

```
Please give me some names:
Ying
Jordan

Your 2 name(s) are saved! They are:
Jordan
Ying
```

12.4 `contains()`

Finally, we can ask our names ArrayList whether or not it has a particular string.

```
1 names.add("Veer");
2 System.out.println(names.contains("Veer")); // true
```

Example 12.4.1: Update your program to check if a name was input by the user. For example:

```
Please give me some names:
Ying
Jordan

Search for a name:
Ying

Yes!
```

12.5 ArrayLists with custom classes

An `ArrayList` can hold any type of object! For example, here is a constructor for an `ArrayList` holding an instance of a `Person` class:

```
1 ArrayList<Person> people = new ArrayList<Person>();
```

where `Person` is defined as

```
1 public class Person {
2
3     String name;
4     int age;
5
6     public Person(String name, int age) {
7         this.name = name;
8         this.age = age;
9     }
10
11     public String getName() {
12         return this.name;
13     }
14
15     public int getAge() {
16         return this.age;
17     }
18 }
```

Example 12.5.1: Modify our program to save the user's input names as `Person` instances. Rather than storing `String` objects in the `ArrayList`, store `Person` objects by constructing them with the input name. You'll need to use the `Person` constructor to get a `Person` instance!

12.6 Arrays versus ArrayLists

The most important difference between an array and an `ArrayLists` is that an `ArrayLists` is dynamically resized. Notice that when we created an `ArrayLists`, we did not specify a size. But when we create an array, we must:

```
1 int[] numbers1 = new int[10];
2 ArrayList<Integer> numbers2 = new ArrayList<Integer>();
```

In the first line of code, we create an array of length 10, and we cannot add more than 10 elements to `numbers1`. In the second line of code, we create an `ArrayList`, but we do not need to specify a size.

The second distinction between arrays and `ArrayLists` is syntax. In Java, nearly everything is a class. An `ArrayList` is a class with useful methods such as `get()` and `contains()` and is therefore more “Java-like”. Arrays are not quite primitives such as `ints` and not quite classes. Most likely, their syntax was borrowed from the C or C++ programming languages.

Exercises

Exercise 12.6.0: Write a class `BlueBook` that tells the user the price of their car, depending on the make, model, and year. You should use `Car.java` and the stencil file provided, `BlueBook.java`. Your program depends on what cars your `BlueBook` supports, but here is an example program:

```
What is your car's make?  
Toyota  
What is your Toyota's model?  
Corolla  
What is your Toyota Corolla's year?  
1999  
  
Your 1999 Toyota Corolla is worth $2000.
```

Exercise 12.6.0: Notify the user if the car is not in your `BlueBook`.

Exercise 12.6.0: Clean up `main` by putting your code for creating the `ArrayList` in a separate method. What type should the method return?

Exercise 12.6.0: If the car is not in the `BlueBook`, ask the user to input the relevant data, construct a new `Car` instance, add it to your `ArrayList`.

References

1. https://github.com/accesscode-2-1/unit-0/blob/master/lessons/week-3/2015-03-24_arraylists.md

Systems Development

13.1 Introduction

Creating reliable software requires strong coding skills, so you can get the computer to do what you want. It additionally requires good “systems development” skills so that (1) your code is organized / quickly updateable by yourself and others and (2) designing the product your users want, rather than what you *think* they want. Systems development is a useful concept not just for computer science projects, but for any type of product development.

For code organization, we use *Application Programming Interfaces* (APIs), which are a set of functions and procedures to allow different portions of code in a big project to communicate with each other. For user-guided product design, we use *iterative design*, which is the process of making many imperfect iterations of a project to prototype to users before trying to make the final product.

13.2 Code Organization

13.2.1 Modularity

Definition 13.2.1. *Modularity* means that code can be separated into many smaller components that are relatively independent.

Modularity has two primary benefits:

- **Teamwork:** It allows different people to work on different portions of the code without interfering with one another. This is how software companies like Google, Amazon, and Microsoft can have tens of thousands of teammembers all working on the same project.
- **Updateability:** It ensures that if one portion of the code needs to be updated or fixed, the rest of the code will not be impacted.

To illustrate this, we will consider modularity’s usefulness in cars. Cars are incredibly complicated, but they can be understood by considering them as a collection of interacting components. Each component has its own functionality and relationship to the rest of the car.

A few components of a car include:

- **Tires:** a round object that can attach to an axle
- **Engine:** a device that can spin an axle
- **Axle:** a rod that can be connected in the center to an engine, and at the edges to tires

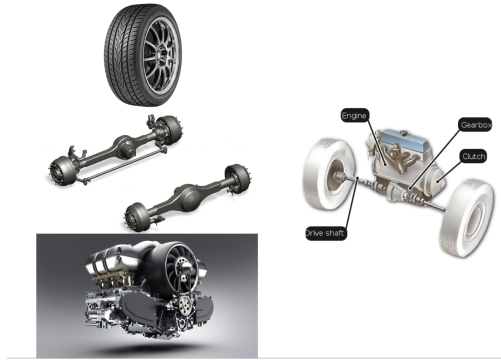


Figure 13.1: Tires, engine, and axle can be combined to make the core axle-tire-engine unit of a car.



Figure 13.2: A variety of tires are available, and each must simply (1) be built to roll and (2) be connectable to an axle

There are different ways to design each of these. They can each be made from a number of materials, and there are different sizes and varieties of each (e.g. 4 vs 6 cylinder engine, snow vs all-season tires, etc.). Furthermore, any *combination* of these different varieties can be put together: we can change the type of tires we have without thinking about what type of axle or motor is in the car.

We will now look at the benefits of this modularity in car design:

- **Teamwork:** Car companies can divide the work amongst multiple teams, and those teams won't interfere with each other as they work. One team can spend months developing the axle: finding the right alloy, shape, and weight to handle the load from the car. Another team can spend months developing the engine: trying out different arrangements of the cylinders or different numbers of cylinders. What's more, car manufacturers almost always outsource the production of their tires to tire manufacturers: they leave the decisions on the type of rubber and treads to a team totally outside of their company. The modularity of the car allows these teams to work totally separately, with the only communication between them the above list: the tire must roll and attach to an axle; the engine must spin the axle; and the axle must connect to the engine and the tires. Once the teams agree on the way they will bolt the systems together, they no longer have to communicate.
- **Updateability:** Modularity allows car enthusiasts to upgrade and customize their vehicles without having to redesign the whole thing. If a vehicle owner wants to replace their V6 with a V8, they can do that as long as the connections between the motor and the rest of the vehicle are the same (i.e. the connection to the axle, to the fuel source, etc). If a vehicle owner wants to replace their all season tires for snow tires in the winter, they can do that as long as the connection with the rest of the vehicle is the same (i.e. the connection to the axle is the same).

13.2.2 APIs

Definition 13.2.2. An *Application Programming Interface*, or API, is a contract between components in a system, expressing what each component can expect from the others.

Building on the previous section's example of a vehicle, we can consider the API between the engine manufacturer and the vehicle manufacturer.

The engine manufacturer can expect the vehicle to:

- Have an axle.
- Have a fuel source.
- Have a cooling source.
- Have a control system.

The vehicle manufacturer the engine to:

- Attach to an axle.
- Spin at a given revolution speed.

As we saw in the previous section, this type of contract allows division of labor between different teams working on a project to ensure that they can work independently: as the engine manufacturer,

Another API exists between a driver and the manufacturer: the driver can expect the car to have a:

- Device to turn on/off the car.
- Device to steer the car.
- Device to accelerate the car.
- Device to decelerate the car.

The car maker can expect the driver to have

- Arms and hands that can turn and push.
- Feet and legs that can press.

Example 13.2.3: Make the API between an ATM and a user.

From the perspective of a user, an ATM's purpose is to take in a credit card and a specified dollar amount from a user, and output the requested amount of money. The user can expect an ATM to:

- Read their credit card.
- Specify the card's PIN number.
- Specify a requested dollar amount.
- Output money and notify their bank of the transaction.

An ATM can expect the user to:

- Have a credit card.
- Type with their fingers.
- Read and respond to prompts on the screen.

Note that a problem with the API translates to a problem with the end-product. This API doesn't guarantee that a blind person will be able to use the ATM, for example (they won't be able to read).

Example 13.2.4: Make the API between an ATM and a bank.

From the perspective of a bank, an ATM's purpose is to send in credit card info and the requested dollar amount, and fulfill a user request for money only if the bank authenticates the request.

The bank can expect an ATM to

- Accurately send the bank credit card info and a PIN.
- Accurately send the bank the requested dollar amount.
- Complete the user's transaction only if the bank sends back a confirmation.

The ATM can expect the bank to

- Receive credit card info and a PIN.
- Send back a confirmation or denial.

13.3 User-oriented Design

13.3.1 Waterfall vs Iterative Design

Example 13.3.1: build the highest tower possible in 6 minutes using the following ingredients:^a

- 10 sticks of dry spaghetti
- One foot of string
- One foot of tape

^aAdapted from <https://tinkerlab.com/spaghetti-tower-marshmallow-challenge/>

This activity generally demonstrates that iteration trumps pre-planning. It's faster to just trying out imperfect designs than to try to wait for a perfect idea. With a 6 minute time limit, iteration tends to work out better than pre-planning.

Definition 13.3.2. *Waterfall design* is a development process in which each stage of development is finished before the next is started.

[AFM: I'm unsure about putting this here, shouldn't it just go in a bibliography?]The components of waterfall design, adapted from <https://airbrake.io/blog/sdlc/waterfall-model>, are the following:

1. Requirements: define what the application should do (essentially, write the API between a user and your product).
2. Design: decide what the product will look like based on the requirements, and how it will be implemented.
3. Implementation: build the product based on the design.
4. Test: ensure the product works as expected.
5. Deployment: release the product to the users.

In waterfall design, mistakes early in the process can kill you. You might spend a lot of time and money going through the full waterfall and developing a final product and then realize that the requirements were wrong. Going back to a previous example, suppose you had designed an ATM thinking that users would be able to read and later found out that blind users must also be able to access the ATM. You would have to completely redesign the system, perhaps having to put speakers into the ATM so that the prompts can be read aloud to the user.

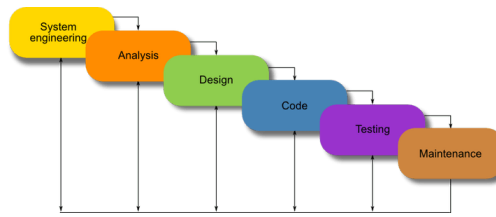


Figure 13.3: Waterfall design. <https://airbrake.io/blog/sdlc/waterfall-model>

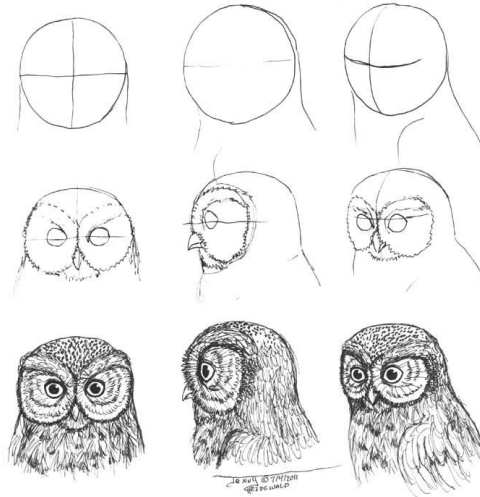


Figure 13.4: Iterative drawing of an owl. <https://www.pinterest.com/pin/469289223655955022/>

Definition 13.3.3. In *iterative design*, many iterations of a project are built and floated to users with the expectation that the requirements and design may need to be adjusted in further iterations.

Iterative design entails going through the same 4 steps of specifying requirements, drafting a design, implementing, then deploying. However, less time and money is invested into trying to make the first iteration perfect. The first iteration of a project might look horrible, but users will be able to tell you the fundamental flaws (such as missing requirements) before you start implementing a perfect product for the wrong problem.

Iterative design is also used in drawing. As shown in Figure 13.4, professional artists start with a rough sketch of an object before starting to fill in details. The first iteration (top) doesn't look very good, but you might find out that you're missing basic requirements: you might be missing certain body parts, or decide you'd like to add a background or another object. By the second iteration, things look a little better, but you still might find more fundamental errors along the way. By the time you're ready to fill in all of the details in the bottom step, you're confident you're drawing the figure you want to draw.

13.4 Testing

Code rarely works the way we would like it to the first time around.

Definition 13.4.1. A *software bug*, is an error in a computer program which causes an incorrect or unexpected result, or to behave in unintended way (Wikipedia).

Medium estimated that in 2016 approximately a trillion dollars was lost to the US economy due to software bugs. As one tangible example of the cost of bugs, in 1962

NASA's 18 million dollar Mariner 1 had to be self-destructed mid-flight because of a missing hyphen in the control code.¹

For this reason, developers are always expected to carefully test the code they write to ensure there are no bugs. In addition, even relatively small software firms often have an entire team whose sole job is to test code written by the developers. Their job is to use the product in a variety of potentially unexpected ways. In essence, they try their best to break the code. The product is not ready for the market until the testing team is unable to break it, and all of its behaviors are as expected by the API for the product.

¹<https://medium.com/@ryancohane/financial-cost-of-software-bugs-51b4d193f107>

Appendices

Possible appendices topics

- Comparing floating point numbers.
- Type conversion.
- Useful math functions, e.g. `min()`, `log`.