

Computer Systems

Muchang Bahng

Spring 2024

Contents

1	Computers	3
2	Representation of Data	5
2.1	Representation of Numbers	5
2.2	Representation of General Sets	8
2.2.1	Prefix-free Encoding	8
2.3	Representing Letters and Text	9
2.3.1	ASCII	9
2.3.2	Extended ASCII	11
2.3.3	ISO-10646, UCS	12
2.3.4	Unicode, UTF-8	12
2.3.5	All Plaintext depends on Encodings	14
2.4	Text Files	14
2.4.1	Plain Text, Rich Text	15
2.5	Binary Files	15
2.5.1	Viewing Binary Files	15
3	Encoding Schemes	16
3.1	Booleans and Characters	16
3.2	Integer Family	16
3.2.1	Arithmetic Operations on Binary Numbers	20
3.3	Float Family	20
4	Memory	21
4.1	Debugging and Object Dumping	24
4.2	Endian Architecture	24
4.3	Type Casting	24
4.4	Pointers	24
4.4.1	Call by Value vs Call by Reference	25
4.4.2	Pointer Errors	25
4.5	Pointer Arithmetic	27
4.6	Global, Stack, and Heap Memory	29
4.7	Dynamic Memory Allocation	32
5	Implementations of Memory Structures in C	32
5.1	Arrays	32
5.2	Strings	32
5.3	Structs	32
5.4	Functions	32
5.5	Classes (for C++)	32
5.6	Input Output	32

6	Central Processing Unit	32
6.1	Circuits	35
6.2	Registers	36
6.2.1	x86 Assembly Registers	37
6.2.2	ARM Assembly Registers	38
6.3	Addressing Modes	38
6.3.1	x86 Assembly Addressing Modes	39
6.3.2	ARM Assembly Addressing Modes	40
6.4	Instructions	40
6.4.1	Moving and Arithmetic	41
6.4.2	Conditionals	41
6.4.3	Control Transfer on Stack	41
6.4.4	Multiple Functions	42
6.4.5	x86-64 Instructions	43
6.4.6	ARM Instructions	73
6.4.7	Buffer Overflows	73
7	Compiling and Linking	73
7.1	Precompiling Stage	74
7.2	Compiling Stage	76
7.3	Objdump	79
7.3.1	ELF and Mach-O Formats	79
7.3.2	Objdump Commands	80
7.4	Assembling Stage and Object Files	84
7.5	Linking Stage and Relocation	86
7.5.1	Relocation	86
7.5.2	Linking with One Object File	88
7.5.3	Global vs External Symbols	88
7.5.4	Linking with Multiple Object Files	90
7.6	Compiler Optimization	94
8	Storage Hierarchy	96
8.1	Expanding on von Neumann Architecture	96
8.2	Disk	98
8.3	Locality	100
8.4	Caches	101
8.4.1	Direct Mapped Cache	104
8.4.2	N way Set-Associative Cache	106
8.4.3	Types of Cache Misses	107
9	Operating Systems	108
9.1	Control Flow	109
9.2	Virtual Memory	115
10	Shared Memory and Concurrency	121
10.1	Process Level Concurrency	123
10.2	Thread Level Concurrency	124
10.3	Atomicity Violation Bugs and Mutex Locks	128
10.4	Deadlock Bugs	135
10.5	Order Violation Bugs and Condition Variables	138

Before we do any coding, we must learn the theory behind how computer systems work, which all starts from memory management and CPU architecture. We will use C with the gcc compiler, along with MIPS and NASM assembler. It is imperative to learn these two since given that you know a high level language pretty well (Python in my case), you want to learn C to appreciate the things Python does for you, and you want to learn Assembly to appreciate the things C does for you.¹

To start off, we want a big overall picture of high a computer works. We introduce this with the simplest model of the computer, the Von Nuemann architecture. It consists of a **central processing unit** (CPU), **memory**, and an **input/output** (I/O) system. We show a diagram of this first for conciseness in Figure 1.

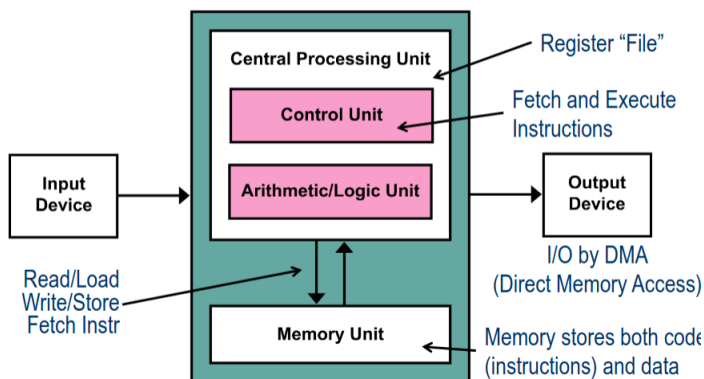


Figure 1: von Neumann Architecture

We will go through these one by one, touching on C and Assembly along the way, but the implementation of these things can differ by the **computer architecture**, so let’s list some of the basic ones.

Definition 0.1 (Computer Architecture)

The **computer architecture** is the design of the computer, which includes the CPU, memory, and I/O system. There are many different architectures, but we will focus on the most common ones.

We first go over some basic theoretical properties of basic data types, focusing on C, and then we cover all the stuff about memory and then all the stuff about the CPU. This is a natural progression since to work with data, you must first know where to store the data and how it is stored (the memory), and then you want to know how data is manipulated (the CPU).

1 Computers

A computer has many parts which are worth remembering.

Definition 1.1 ()

The 5 main parts of a computer is:

1. The **Solid State Drive (SSD)** (or an older version is **HDD**) is where the computer’s long term memory is stored. Most computers usually have 256GB, 512GB, or 1TB of storage.
2. The **Random Access Memory (RAM)** is where the computer’s short term memory is stored, which is usually 4GB, 8GB, 16GB, or 32GB. The RAM determines how well your computer can work with multiple things running at the same time (applications, browser tabs, ...). Accessing this memory is about 20 100 times faster than accessing the SSD memory.
3. The **processor**, or **Central Processing Unit (CPU)**, is the circuitry that executes instruc-

¹<https://www.youtube.com/watch?v=XlvfHORF26M>

tions that make up a computer program. At the hardware level, a CPU is an **integrated circuit**, or a **chip**. This means that a CPU integrates billions of circuits into a chip, where each circuit is really just logic gate (AND, OR, NOT, etc.) made up of a transistor. We can see this in layers:

$$\text{CPU Chip} \implies \text{Circuit} \implies \text{Logic Gate} \implies \text{Transistor}$$

They're effectively minute gates that switch on or off, thereby conveying the ones or zeros that translate into everything you do with the device. One of the most common advancements of CPU technology is in making those transistors smaller and smaller, where the rate of improvement is referred to as *Moore's Law*.

In simplest terms, the CPU takes instructions from a program/application and performs a calculation. It first *fetches* the instruction from RAM, *decodes* what the instruction actually is, and then executes the instruction using relevant parts of the CPU. The CPU performs basic arithmetic, logic, controlling, and input/output (I/O) operations specified by the instructions in the program.

Originally, CPUs had a single processing core. Today's modern CPU consists of multiple cores that allow it to perform multiple instructions at once, effectively cramming several CPUs on a single chip. Almost all CPUs sold today are at least dual-core or quad-core. Additionally, a physical CPU core can perform two lines of execution (threads) at once with a process called *multithreading*. The clock speed should also be noted with CPUs: the gigahertz figure quoted on the CPU. It denotes how many instructions a CPU can handle per second (giga=billions, mega=millions).

4. The **Graphic Processing Unit (GPU)** is a specialized electronic circuit designed to rapidly manipulate and alter memory to accelerate the creation of images intended for output to a display device.
5. The **kernel** is a computer program at the core of a computer's operating system that has complete control over everything in the system. It facilitates interactions between hardware and software components. On most systems, the kernel is one of the first programs loaded on startup. It handles the rest of startup as well as memory, peripherals, and input/output (I/O) requests from software, translating them into data-processing instructions for the central processing unit.

$$\text{CPU, Memory, Devices} \iff \text{Kernel} \iff \text{Applications}$$

The critical code of the kernel is usually loaded into a separate area of memory, which is protected from access by application programs or other, less critical parts of the operating system. The kernel performs its tasks, such as running processes, managing hardware devices such as the hard disk, and handling interrupts, in this protected **kernel space**. In contrast, application programs like browsers, word processors, or audio or video players use a separate area of memory, **user space**.

All components of a computer communicate through a circuit board called the **motherboard**.

Note that all of these parts work in conjunction. For the CPU to function, it still must feed to specialized hardware the numbers they need to function. It needs to tell the graphics card to show an explosion or tell the hard drive to transfer a document to the system's RAM for quicker access.

Definition 1.2 ()

Other parts of the computer include:

1. The **heatsink** is a passive heat exchanger that transfers heat. It is typically a metallic part which can be attached to a device releasing energy in the form of heat, with the aim of dissipating that heat to a surrounding fluid in order to prevent the device overheating.

Certain microcomputers with all these parts (but not packaged) include the **Raspberry Pi** and the **Arduino**.

Definition 1.3 ()

A computer processes bits of information based off of how much electricity is traveling through a wire. Since there can be varying amounts of electricity running through a wire, engineers use the transistor as a "switch" that turns on when the voltage running through the wire is greater than the *threshold voltage*.

For example, a transistor with a threshold voltage of 4.5V will turn on when there is a current of at least 4.5V running through the wire and off otherwise.

Definition 1.4 ()

A **daemon** is a type of program on Unix-like operating systems that runs unobtrusively in the background, rather than under direct control of a user, waiting to be activated.

2 Representation of Data

The inputs of a computational program at its most fundamental level really takes in a **binary string** of 0s and 1s. Note that the choice of 0 and 1 is for convenience, but it must be binary (i.e. Boolean) in some way in order for the model to be physically implemented by transistors. Once information is in digital form, we can *compute* over it and gain insights from data that were not accessible in prior times. In fact, we can represent an unbounded variety of objects using only two symbols 0 and 1.

Therefore, when we say that a program P takes x as an input, we really mean that P takes as input the *representation of x* as a binary string.

Definition 2.1 ()

A *representation scheme* is a way to map an object x to a unique binary string $E(x) \in \{0, 1\}^*$. That is, given a set of objects, E is an injective (not necessarily surjective) map

$$E : X \longrightarrow \{0, 1\}^*$$

2.1 Representation of Numbers

Definition 2.2 (Representation of the Naturals)

A representation for natural numbers (note that in this context, $0 \in \mathbb{N}$) is the (non-surjective) regular binary representation denoted

$$NtS : \mathbb{N} \longrightarrow \{0, 1\}^* \quad (NtS = \text{"Naturals to Strings"})$$

recursively defined as

$$NtS(n) = \begin{cases} 0 & n = 0 \\ 1 & n = 1 \\ NtS(\lceil n/2 \rceil \text{parity}(n)) & n > 1 \end{cases}$$

where given strings $x, y \in \{0, 1\}^*$, xy denotes the concatenation of x and y , and $\text{parity} : \mathbb{N} \longrightarrow \{0, 1\}^*$

is defined

$$\text{parity}(n) = \begin{cases} 0 & n \text{ is even} \\ 1 & n \text{ is odd} \end{cases}$$

Since NtS is injective, its inverse $StN : \text{Im}NtS \subset \{0, 1\}^* \rightarrow \mathbb{N}$ is well-defined.

Definition 2.3 (Representation of the Integers)

To construct a representation scheme for \mathbb{Z} , we can just add one more binary digit to represent the sign of the number. The binary representation $ZtS : \mathbb{Z} \rightarrow \{0, 1\}^*$ is defined

$$ZtS(m) = \begin{cases} 0NtS(m) & m \geq 0 \\ 1NtS(-m) & m < 0 \end{cases}$$

where NtS is defined as before. Again this function must be injective but need not be surjective.

When representing rational numbers, we cannot simply concatenate the numerator and denominator as such

$$a/b \mapsto ZtS(a)ZtS(b)$$

since this map is not surjective (and may overlap with other integers).

Definition 2.4 (Representation of Rationals)

To represent a rational number a/b , we create a separator symbol $|$ and map the rational number as below in the alphabet $\{0, 1, |\}$.

$$q : a/b \mapsto ZtS(a)|ZtS(b)$$

Then, we use a second map that goes through each digit in z and is defined

$$p : \{0, 1, |\} \rightarrow \{00, 11, 01\} \subset \{0, 1\}^2, p(n) = \begin{cases} 00 & n = 0 \\ 11 & n = 1 \\ 01 & n = | \end{cases}$$

Therefore, p maps the length n string $z \in \{0, 1\}^*$ to the length $2n$ string $\omega \in \{0, 1\}^*$. The representation scheme for \mathbb{Q} is simply

$$QtS \equiv p \circ q$$

Example 2.1 ()

Given the rational number $-5/8$,

$$\frac{-5}{8} \mapsto 1101|01000 \mapsto 11110011010011000000$$

This same idea of using separators and compositions of injective functions can be used to represent arbitrary n -tuples of strings (since a finite Cartesian product of countable sets is also countable).

Theorem 2.1 (Representation of Vectors)

All vectors over the field \mathbb{Q} are representable.

Proof.

We can simply create another separator symbol \cdot and have the initial mapping q map to a string over the alphabet $\{0, 1, |, \cdot\}$, which injectively maps to $\{00, 01, 10, 11\}$.

Corollary 2.1 (Representation of Matrices and Tensors)

Matrices (over \mathbb{Q}), which are a collection of vectors are representable in binary. Furthermore, general tensors over the field \mathbb{Q} are representable in binary.

Proof.

Create more separator symbols and map them to a sufficiently large set (which can be extended arbitrarily). For example, to perhaps $\{000, 001, \dots, 111\}$.

Corollary 2.2 (Representation of Graphs)

Directed graphs, which can be represented with their adjacency matrices, can therefore be represented with binary strings.

Theorem 2.2 (Representation of Images)

Every finite-resolution image can be represented as a binary number.

Proof.

Since we can interpret each image as a matrix where each element (a pixel) is a color, and since each color can be represented as a 3-tuple of rational numbers corresponding to the intensities of red, green, and blue (for humans, we can restrict it to three primary colors), all images can eventually be decomposed into binary strings.

Theorem 2.3 (Representation of Reals)

There exists no representation of the reals

$$NtR : \mathbb{R} \longrightarrow \{0, 1\}^*$$

Proof.

By Cantor's theorem, the reals are uncountable. That is, there does not exist a surjective function $NtR : \mathbb{N} \longrightarrow \mathbb{R}$. This implies the nonexistence of an injective inverse; that is, there does not exist an injective function

$$RtS : \mathbb{R} \longrightarrow \{0, 1\}^*$$

However, since \mathbb{Q} is dense in \mathbb{R} , we can approximate every real number x by a rational number a/b to arbitrary accuracy. There are multiple ways to construct these approximations (decimal approximation up to k th digit, finite continued fractions, truncated infinite series, etc.), but computers use the *floating-point approximation*.

Definition 2.5 (Floating-Point Representation)

The **floating-point representation scheme** of a real number $x \in \mathbb{R}$ is its approximation as a number of the form

$$\sigma b \cdot 2^e$$

where $\sigma \in \{0, 1\}$ determines the sign of the representation of x , e is a (potentially negative) integer, and b is a rational number between 1 and 2 expressed as a binary fraction

$$1.b_0b_1b_2\dots b_k = 1 + \frac{b_1}{2} + \frac{b_2}{4} + \dots + \frac{b_k}{2^k}, \quad b_i \in \{0, 1\}$$

where the number k is fixed (determined by the desired accuracy; greater k implies more digits and better accuracy). The $\sigma b \cdot 2^e$ closest to x is the *floating-point representation*, or *approximation*, of x . We can think of σ determining the sign, e the order of magnitude (in base 2) of x , and b the value of the number scaled down to a value in $[1, 2)$, called the *mantissa*.

2.2 Representation of General Sets

Let there exist some set \mathcal{O} consisting of objects. Then, a representation scheme for representing objects in \mathcal{O} consists of an *encoding* function that maps an object in \mathcal{O} to a string, and a *decoding* function that decodes a string back to an object in \mathcal{O} .

Definition 2.6 ()

Let \mathcal{O} be any set. A *representation scheme* for \mathcal{O} is a pair of functions E, D where

$$E : \mathcal{O} \longrightarrow \{0, 1\}^*$$

is an injective function, and the induced mapping D is restriction of the inverse of E to the image of E .

$$D : \text{Im}(E) \subset \{0, 1\}^* \longrightarrow \mathcal{O}$$

This means that $(D \circ E)(o) = o$ for all $o \in \mathcal{O}$. E is known as the *encoding function* and D is known as the *decoding function*.

2.2.1 Prefix-free Encoding

Definition 2.7 (Prefix)

For two strings y, y' , y is a prefix of y' if y' "starts" with y . That is, y is a **prefix** of y' if $|y| \leq |y'|$ and for every $i < |y|$, $y'_i = y_i$.

With this, we can define the concept of prefix free encoding.

Definition 2.8 ()

Let \mathcal{O} be a nonempty set and $E : \mathcal{O} \longrightarrow \{0, 1\}^*$ be a function. E is **prefix-free** if $E(o)$ is nonempty for every $o \in \mathcal{O}$ and there does not exist a distinct pair of objects $o, o' \in \mathcal{O}$ such that $E(o)$ is a prefix of $E(o')$.

Being prefix-free is a nice property that we would like an encoding to have. Informally, this means that no string x representing an object o is an initial substring of string y representing a different object o' . This means that we can simply represent a *list* of objects simply by concatenating the representations of all the list members and still get a valid, injective representation. We formalize this below.

Theorem 2.4 ()

Suppose that $E : \mathcal{O} \rightarrow \{0, 1\}^*$ is prefix free. Then the following map

$$\bar{E} : \mathcal{O}^* \rightarrow \{0, 1\}^*$$

over all finite length tuples of elements in \mathcal{O} is injective, where for every $o_0, o_1, \dots, o_{k-1} \in \mathcal{O}^*$, we define \bar{E} to be the simple concatenation of the separate encodings of o_i :

$$\bar{E}(o_0, \dots, o_{k-1}) \equiv E(o_0)E(o_1)\dots E(o_{k-1})$$

Even if the representation E of objects in \mathcal{O} is prefix free, this does not imply that our representation \bar{E} of lists of such objects will be prefix free as well. In fact, it won't be, since for example, given three objects o, o', o'' , the representation of the list (o, o') will be a prefix of the representation of the list (o, o', o'') .

However, it turns out that in fact we can transform *every* representation into prefix free form, and so will be able to use that transformation if needed to represent lists of lists, lists of lists of lists, and so on.

Some natural representations are prefix free. For example, every *fixed output length* representation (i.e. an injective function $E : \mathcal{O} \rightarrow \{0, 1\}^n$) is automatically prefix-free, since a string x can only be a prefix of an equal length x' if x and x' are identical. Moreover, the approach that was used for representing rational numbers can be used to show the following lemma.

Lemma 2.1 ()

Let $E : \mathcal{O} \rightarrow \{0, 1\}^*$ be a one-to-one function. Then there is a one-to-one prefix-free encoding \bar{E} such that

$$|\bar{E}(o)| \leq 2|E(o)| + 2$$

for every $o \in \mathcal{O}$.

Proof.

The general idea is the use the map $0 \mapsto 00, 1 \mapsto 11$ to "double" every bit in the string x and then mark the end of the string by concatenating to it the pair 01. If we encode a string x in this way, it ensures that the encoding of x is never a prefix of the encoding of a distinct string x' . (Note that this is not the only or even the best way to transform an arbitrary representation into prefix-free form.)

2.3 Representing Letters and Text

We can represent a letter or symbol by a string, and then if this representation is prefix free, we can represent a sequence of symbols by merely concatenating the representation of each symbol. Here are a few examples.

2.3.1 ASCII

The **ASCII** (also called US-ASCII) code, which stands for American Standard Code for Information Interchange is a 7 bit character code where every single bit represents a unique character. ASCII codes represent text in computers, telecommunications equipment, and other devices. Most modern character-encoding schemes are based on ASCII, although they support many additional characters.

The first 32 characters are called the *control characters*: codes originally intended not to represent printable information, but rather to control devices (such as printers) that make use of ASCII, or to provide meta-information about data streams. For example, character 10 (decimal) represents the "line feed" function (which causes a printer to advance its paper) and character 8 represents "backspace." Except for the control characters that prescribe elementary line-oriented formatting, ASCII does not define any mechanism for describing the structure or appearance of text within a document.

Dec	Oct	Hex	Bin	Symbol	Description
0	000	00	0000000	NULL	Null char
1	001	01	0000001	SOH	Start of Heading
2	002	02	0000010	STX	Start of Text
3	003	03	0000011	ETX	End of Text
4	004	04	0000100	EOT	End of Transmission
5	005	05	0000101	ENQ	Enquiry
6	006	06	0000110	ACK	Acknowledgement
7	007	07	0000111	BEL	Bell
8	010	08	0001000	BS	Back Space
9	011	09	0001001	HT	Horizontal Tab
10	012	0A	0001010	LF	Line Feed
11	013	0B	0001011	VT	Vertical Tab
12	014	0C	0001100	FF	Form Feed
13	015	0D	0001101	CR	Carriage Return
14	016	0E	0001110	SO	Shift Out/X-On
15	017	0F	0001111	SI	Shift In/X-Off
16	020	10	0010000	DLE	Data Line Escape
17	021	11	0010001	DC1	Device Control 1
18	022	12	0010010	DC2	Device Control 2
19	023	13	0010011	DC3	Device Control 3
20	024	14	0010100	DC4	Device Control 4
21	025	15	0010101	NAK	Negative Acknowledgement
22	026	16	0010110	SYN	Synchronous Idle
23	027	17	0010111	ETB	End of Transmit Block
24	030	18	0011000	CAN	Cancel
25	031	19	0011001	EM	End of Medium
26	032	1A	0011010	SUB	Substitute
27	033	1B	0011011	ESC	Escape
28	034	1C	0011100	FS	File Separator
29	035	1D	0011101	GS	Group Separator
30	036	1E	0011110	RS	Record Separator
31	037	1F	0011111	US	Unit Separator

The rest of the characters are the ASCII printable characters.

Dec	Oct	Hex	Bin	Sym	Description	Dec	Oct	Hex	Bin	Sym	Description
32	040	20	0100000		Space	80	120	50	1010000	P	Uppercase P
33	041	21	0100001	!	Exclamation	81	121	51	1010001	Q	Uppercase Q
34	042	22	0100010	"	Double quotes	82	122	52	1010010	R	Uppercase R
35	043	23	0100011	#	Number	83	123	53	1010011	S	Uppercase S
36	044	24	0100100	\$	Dollar	84	124	54	1010100	T	Uppercase T
37	045	25	0100101	%	Per cent sign	85	125	55	1010101	U	Uppercase U
38	046	26	0100110	&	Ampersand	86	126	56	1010110	V	Uppercase V
39	047	27	0100111	'	Single quote	87	127	57	1010111	W	Uppercase W
40	050	28	0101000	(Open paren.	88	130	58	1011000	X	Uppercase X
41	051	29	0101001)	Closed paren.	89	131	59	1011001	Y	Uppercase Y
42	052	2A	0101010	*	Asterisk	90	132	5A	1011010	Z	Uppercase Z
43	053	2B	0101011	+	Plus	91	133	5B	1011011	[Opening bracket
44	054	2C	0101100	,	Comma	92	134	5C	1011100	\	Backslash
45	055	2D	0101101	-	Hyphen	93	135	5D	1011101]	Closing bracket
46	056	2E	0101110	.	Period	94	136	5E	1011110	^	Caret
47	057	2F	0101111	/	Slash	95	137	5F	1011111	_	Underscore
48	060	30	0110000	0	Zero	96	140	60	1100000	`	Grave accent
49	061	31	0110001	1	One	97	141	61	1100001	a	Lowercase a
50	062	32	0110010	2	Two	98	142	62	1100010	b	Lowercase b
51	063	33	0110011	3	Three	99	143	63	1100011	c	Lowercase c
52	064	34	0110100	4	Four	100	144	64	1100100	d	Lowercase d
53	065	35	0110101	5	Five	101	145	65	1100101	e	Lowercase e
54	066	36	0110110	6	Six	102	146	66	1100110	f	Lowercase f
55	067	37	0110111	7	Seven	103	147	67	1100111	g	Lowercase g
56	070	38	0111000	8	Eight	104	150	68	1101000	h	Lowercase h
57	071	39	0111001	9	Nine	105	151	69	1101001	i	Lowercase i
58	072	3A	0111010	:	Colon	106	152	6A	1101010	j	Lowercase j
59	073	3B	0111011	;	Semicolon	107	153	6B	1101011	k	Lowercase k
60	074	3C	0111100	<	Less than	108	154	6C	1101100	l	Lowercase l
61	075	3D	0111101	=	Equals	109	155	6D	1101101	m	Lowercase m
62	076	3E	0111110	>	Greater than	110	156	6E	1101110	n	Lowercase n
63	077	3F	0111111	?	Question mark	111	157	6F	1101111	o	Lowercase o
64	100	40	1000000	@	At symbol	112	160	70	1110000	p	Lowercase p
65	101	41	1000001	A	Uppercase A	113	161	71	1110001	q	Lowercase q
66	102	42	1000010	B	Uppercase B	114	162	72	1110010	r	Lowercase r
67	103	43	1000011	C	Uppercase C	115	163	73	1110011	s	Lowercase s
68	104	44	1000100	D	Uppercase D	116	164	74	1110100	t	Lowercase t
69	105	45	1000101	E	Uppercase E	117	165	75	1110101	u	Lowercase u
70	106	46	1000110	F	Uppercase F	118	166	76	1110110	v	Lowercase v
71	107	47	1000111	G	Uppercase G	119	167	77	1110111	w	Lowercase w
72	110	48	1001000	H	Uppercase H	120	170	78	1111000	x	Lowercase x
73	111	49	1001001	I	Uppercase I	121	171	79	1111001	y	Lowercase y
74	112	4A	1001010	J	Uppercase J	122	172	7A	1111010	z	Lowercase z
75	113	4B	1001011	K	Uppercase K	123	173	7B	1111011	{	Opening brace
76	114	4C	1001100	L	Uppercase L	124	174	7C	1111100		Vertical bar
77	115	4D	1001101	M	Uppercase M	125	175	7D	1111101	}	Closing brace
78	116	4E	1001110	N	Uppercase N	126	176	7E	1111110	~	Tilde
79	117	4F	1001111	O	Uppercase O	127	177	7F	1111111		Delete

2.3.2 Extended ASCII

The **Extended ASCII** (EASCII or high ASCII) character encodings are 8-bit or larger encodings that include the standard 7-bit ASCII characters, plus additional characters. Note that this does not mean that the standard ASCII coding has been updated to include more than 128 characters nor does it mean that there is an universal extension to the original ASCII coding. In fact, there are several (over 100) extended ASCII encodings.

With the creation of the 7-bit ASCII format, increased need for more letters and symbols (such as characters in other languages or more punctuation/mathematical symbols). With better computers and software, it became obvious that they could handle text that uses 256-character sets at almost no additional cost in programming or storage. The 8-bit format would allow ASCII to be used unchanged and provide 128 more characters.

But even 256 characters is still not enough to cover all purposes, all languages, or even all European languages,

so the emergence of *many* ASCII-derived 8-bit character sets was inevitable. Translating between these sets (*transcoding*) is complex, especially if a character is not in both sets and was often not done, producing **mojibake** (semi-readable text resulting from text being decoded using an unintended character encoding. The result is a systematic replacement of symbols with completely unrelated ones, often from a different writing system). ASCII can also be used to create graphics, commonly called **ASCII art**.

But ASCII isn't enough. We have lots of languages with lots of characters that computers should ideally display. Unicode assigns each character a unique number, or code point. Computers deal with such numbers as bytes: 8-bit computers would treat an 8-bit byte as the largest numerical unit easily represented on the hardware, 16-bit computers would expand that to 2 bytes, and so forth. Old character encodings like ASCII are from the (pre-) 8-bit era, and try to cram the dominant language in computing at the time, i.e. English, into numbers ranging from 0 to 127 (7 bits). When ASCII got extended by an 8th bit for other non-English languages, the additional 128 numbers/code points made available by this expansion would be mapped to different characters depending on the language being displayed. The **ISO-8859** standards are the most common forms of this mapping:

1. **ISO-8859-1**
2. **ISO-8859-15**, also called **ISO-Latin-1**

But that's not enough when you want to represent characters from more than one language, so cramming all available characters into a single byte just won't work. The following shows ways to do this (that is compatible with ASCII).

2.3.3 ISO-10646, UCS

We can simply expand the value range by adding more bits. The UCS-2 uses 2 bytes (or 16 bits) and UCS-4 uses 4 bytes (32 bits). However, these codings suffer from inherently the same problem as ASCII and ISO-8859 standards, as their value range is still limited, even if the limit is vastly higher. Note that these encode from the ISO-10646, which defines several character encoding forms for the Universal Coded Character Set.

1. UCS-2 can store $2^{16} = 65,536$ characters.
2. UCS-4 can store $2^{32} = 4,294,967,296$ characters.

Notice that UCS encoding has a fixed number of bytes per character, which means that UCS-2 stores each character in 2 bytes, and UCS-4 stores each character in 4 bytes. This is different from **UTF-8** encoding.

ISO 10646 and Unicode have an identical repertoire and numbers—the same characters with the same numbers exist on both standards, although Unicode releases new versions and adds new characters more often. Unicode has rules and specifications outside the scope of ISO 10646. ISO 10646 is a simple character map, an extension of previous standards like ISO 8859. In contrast, Unicode adds rules for collation, normalization of forms, and the bidirectional algorithm for right-to-left scripts such as Arabic and Hebrew. For interoperability between platforms, especially if bidirectional scripts are used, it is not enough to support ISO 10646; Unicode must be implemented.

2.3.4 Unicode, UTF-8

Unicode is the universal character encoding, maintained by Unicode Consortium, and it covers the characters for all the writing systems of the world, modern and ancient. It also includes technical symbols, punctuation, and many other characters used in writing text. As of Unicode Version 13.0, the Unicode standard contains 143,859 characters, stored in the format $U+****$, where $****$ is a number in hexadecimal notation. Notice that these ones are not fixed in the number of bits; that is,

$$U+27BD \text{ and } U+1F886$$

are perfectly viable representations of characters in Unicode. Even though only 143,859 characters are in use, Unicode currently allows for 1,114,112 ($16^5 + 16^4$) code values, and assigns codes covering nearly all

modern text writing systems, as well as many historical ones and for many non-linguistic characters such as printer’s dingbats, mathematical symbols, etc.

*Note that Unicode, along with ISO-10646, is a standard that assigns a name and a value (**Character Code** or **Code-Point**) to each character in its repertoire.* However, the Unicode format must be encoded in a binary format for the computer to understand. When you save a document, the text editor has to explicitly set its encoding to be UTF-8 (or whatever other format) the user wants it to be. Also, when a text editor program reads a file, it needs to select a text encoding scheme to decode it correctly. Even further, when you are typing and entering a letter, the text editor needs to know what scheme you use so that it will save it correctly. Therefore, *UTF-8 encoding is a way to represent these characters digitally in computer memory.* The way that **UTF-8** encodes characters is with the following format:

	1st Byte	2nd Byte	3rd Byte	4th Byte	Number of Free Bits
1	0 xxxxxxx				7
2	110 xxxxx	10xxxxxx			(5+6)=11
3	1110 xxxx	10xxxxxx	10xxxxxx		(4+6+6)=16
4	11110 xxx	10xxxxxx	10xxxxxx	10xxxxxx	(3+6+6+6)=21

From this, we can see that UTF-8 uses a variable number of bytes per character. All UTF encodings work in roughly the same manner: you choose a unit size, which for UTF-8 is 8 bits, for UTF-16 is 16 bits, and for UTF-32 is 32 bits. The standard then defines a few of these bits as *flags* (e.g. the 0, 110, 1110, 11110, ...). If they’re set, then the next unit in a sequence of units is considered part of the same character. If they’re not set, this unit represents one character fully. Thus, the most common (English) characters only occupy one byte in UTF-8 (two in UTF-16, 4 in UTF-32), but other language characters can occupy more bytes. We can see that UTF-8 can encode up to (and slightly more than) $2^{21} = 2,097,152$ characters. UTF-8 is by far the most common encoding for the World Wide Web, accounting for 96.0% of all web pages, and up to 100% for some languages, as of 2021.

For example, let’s take a random character, say with the Unicode value to be U+6C49. Then, we convert this to binary to get

01101100 01001001

But we can’t just store this because this isn’t a prefix-free notation. This is when UTF-8 is needed. Using the chart above, we need to prefix our character with some headers/flags. The binary Unicode value of the character is 16 bits long, so we can store it in 3 bytes (in the format of the third row) as it provides enough space. The headers are not bolded, while the binary values added are.

11100110 10110001 10001001

We can take another example of a character with the Unicode value U+1F886. Converting to binary gets

0001 1111 1000 1000 0110

There are 20 bits, so we will need to store it in 4 bytes (in the format of fourth row) as it provides enough space (21). We convert the 20-bit-long binary Unicode value to a 21-bit-long value (so that it is compatible with the 21 free bits) to get

0 0001 1111 1000 1000 0110

Encoding it in UTF-8 in 4 bytes gives

11110000 10011111 10100010 10000110

There is no need to go beyond 4 bytes since every Unicode value will have at most 5 hexadecimal digits (since $16^5 = 1,048,576$, which is far more than the number of characters there are). There is also another, obsolete, encoding used called the **UTF-7**.

Both the UCS and UTF standards encode the code points as defined in Unicode. In theory, those encodings could be used to encode any number (within the range the encoding supports) - but of course these encodings

were made to encode Unicode code points. Windows handles so-called "Unicode" strings as UTF-16 strings, while most UNIXes default to UTF-8 these days. Communications protocols such as HTTP tend to work best with UTF-8, as the unit size in UTF-8 is the same as in ASCII, and most such protocols were designed in the ASCII era. On the other hand, UTF-16 gives the best average space/processing performance when representing all living languages.

While UTF-7, 8, 16, and 32 all have the nice property of being able to store *any* code point correctly, there are hundreds of encodings that can only store a set amount of characters. If there's no equivalent for the Unicode code point you're trying to represent in the encoding you're trying to represent it in, you usually get a little question mark: ? For example, trying to store Russian or Hebrew letters in these encodings results in a bunch of question marks.

2.3.5 All Plaintext depends on Encodings

Note that **it does not make sense to have a string without knowing what encoding it uses**. We can't just assume that every plaintext is in ASCII, since there are hundreds of extended ASCII encodings. If you have a string, in memory, in a file, or in an email message, you have to know what encoding it is in or you cannot interpret it or display it to users correctly.

For example, when you are sending an email, Gmail is the only client that automatically converts your text to UTF-8, regardless of what you set in the header. The browser also uses a certain encoding, which can be accessed (and changed) under the "view" tab.

2.4 Text Files

The ASCII character set is the most common compatible subset of character sets for English-language text files, and is generally assumed to be the default file format in many situations.

In the Mac, checking the character encoding of a text file can be done with the command

```
1 >>>file -I filename.txt
2 filename.txt: text/plain; charset=us-ascii
```

ASCII covers American English, but for the British Pound sign, the Euro sign, or characters used outside English, a richer character set must be used. In many systems, this is chosen based on the default setting on the computer it is read on. Prior to UTF-8, this was traditionally single-byte encodings (such as ISO-8859-1 through ISO-8859-16) for European languages and wide character encodings for Asian languages. However, most computers use UTF-8 as the natural extension. We can check this firsthand by inputting a non-ASCII character in filename.txt, which would result in

```
1 >>>file -I filename.txt
2 filename.txt: text/plain; charset=utf-8
```

Because encodings necessarily have only a limited repertoire of characters, often very small, many are only usable to represent text in a limited subset of human languages. Unicode is an attempt to create a common standard for representing all known languages, and most known character sets are subsets of the very large Unicode character set. Although there are multiple character encodings available for Unicode, the most common is UTF-8, which has the advantage of being backwards-compatible with ASCII; that is, every ASCII text file is also a UTF-8 text file with identical meaning. UTF-8 also has the advantage that it is easily auto-detectable. Thus, a common operating mode of UTF-8 capable software, when opening files of unknown encoding, is to try UTF-8 first and fall back to a locale dependent legacy encoding when it definitely isn't UTF-8.

Because of their simplicity, text files are commonly used for storage of information. When data corruption occurs in a text file, it is often easier to recover and continue processing the remaining contents. A disadvantage of text files is that they usually have a low entropy, meaning that the information occupies more

storage than is strictly necessary. A simple text file may need no additional metadata (other than knowledge of its character set) to assist the reader in interpretation. A text file may contain no data at all, which is a case of zero-byte file.

2.4.1 Plain Text, Rich Text

Plain text is a loose term for data (e.g. file contents) that represent only characters of readable material but not its graphical representation nor other objects (floating-point numbers, images, etc.). It may also include a limited number of "whitespace" characters that affect simple arrangement of text, such as spaces or line breaks. A plain text file cannot have bold text, fonts, larger font sizes, or any other special text formatting.

Plain text is different from **formatted text** or **rich text**, where style information is included and structured text, where structural parts of the document such as paragraphs, sections, and the like are identified.

Most systems associate plain text file with the file .txt. To view a plaintext file, a text editor must be used.

2.5 Binary Files

Another type of data is a **binary file**, which is a computer file that is not a text file (it is often used as a term meaning *non-text file*). Many binary file formats contain parts that can be interpreted as text; for example, some computer document files containing formatted text, such as older Microsoft Word document files, contain the text of the document but also contain formatting information in binary form.

Binary files are usually thought of as being a sequence of bytes, which means the binary digits (bits) are grouped in eights. Binary files typically contain bytes that are intended to be interpreted as something other than text characters. Compiled computer programs are typical examples; indeed, compiled applications are sometimes referred to, particularly by programmers, as binaries. But binary files can also mean that they contain images, sounds, compressed versions of other files, etc. – in short, any type of file content whatsoever.

2.5.1 Viewing Binary Files

A hex editor or viewer may be used to view file data as a sequence of hexadecimal (or decimal, binary or ASCII character) values for corresponding bytes of a binary file.

If a binary file is opened in a text editor, each group of eight bits will typically be translated as a single character, and the user will see a (probably unintelligible) display of textual characters. If the file is opened in some other application, that application will have its own use for each byte: maybe the application will treat each byte as a number and output a stream of numbers between 0 and 255—or maybe interpret the numbers in the bytes as colors and display the corresponding picture. Other type of viewers (called 'word extractors') simply replace the unprintable characters with spaces revealing only the human-readable text. This type of view is useful for a quick inspection of a binary file in order to find passwords in games, find hidden text in non-text files and recover corrupted documents. It can even be used to inspect suspicious files (software) for unwanted effects. For example, the user would see any URL/email to which the suspected software may attempt to connect in order to upload unapproved data (to steal). If the file is itself treated as an executable and run, then the operating system will attempt to interpret the file as a series of instructions in its machine language

Standards are very important to binary files. For example, a binary file interpreted by the ASCII character set will result in text being displayed. A custom application can interpret the file differently: a byte may be a sound, or a pixel, or even an entire word. Binary itself is meaningless, until such time as an executed algorithm defines what should be done with each bit, byte, word or block. Thus, just examining the binary and attempting to match it against known formats can lead to the wrong conclusion as to what it actually represents. This fact can be used in *steganography*, where an algorithm interprets a binary data file differently to reveal hidden content. Without the algorithm, it is impossible to tell that hidden content exists.

3 Encoding Schemes

In order to get into memory, it is helpful to know the theory behind how primitive types are stored in memory.

Definition 3.1 (Collections of Bits)

There are many words that are used to talk about values of different data types:

1. A **bit** (b) is either 0 or 1.
2. A **Hex** (x) is a collection of 4 bits, with a total of $2^4 = 16$ possible values, and this is used since it is easy to read for humans.
3. A **Byte** (B) is a collection of 8 bits or 2 hex, with a total of $2^8 = 256$ possible values, and most computers will work with Bytes as the smallest unit of memory.

Definition 3.2 (Collections of Bytes)

Sometimes, we want to talk about slightly larger collections, so we group them by how many bytes they have. However, note that these may not always be the stated size, depending on what architecture or language you are using. This is more of a general term, and they may have different names in different languages. If there is a difference, we will state it explicitly.

1. A **word** (w) is 2 Bytes.
2. A **long** (l) is 4 Bytes.
3. A **quad** (q) is 8 Bytes.

Try to know which letter corresponds to which structure, since that will be useful in both C and Assembly.

3.1 Booleans and Characters

Definition 3.3 (Booleans in C)

The most basic type is the boolean, which is simply a bit. In C, it is represented as `bool`, and it is either `true` (1) or `false` (0).

We can manually check the size of the boolean type in C with the following code.

<pre> 1 #include<stdio.h> 2 #include<stdbool.h> 3 4 int main() { 5 printf("%lu\n", sizeof(bool)); 6 return 0; 7 }</pre>	<pre> 1 1 2 . 3 . 4 . 5 . 6 . 7 .</pre>
---	---

Figure 2: We can verify the size of various primitive data types in C with the `sizeof` operator.

3.2 Integer Family

The most primitive things that we can store are integers. Let us talk about how we represent some of the simplest primitive types in C: unsigned short, unsigned int, unsigned long, unsigned long long.

Definition 3.4 (Unsigned Integer Types in C)

In C, there are several integer types. We use this hierarchical method to give flexibility to the programmer on the size of the integer and whether it is signed or not.

1. An **unsigned short** is 2 bytes long and can be represented as a 4-digit hex or 16 bits, with values in $[0 : 65,535]$. Therefore, say that we have
2. An **unsigned int** is 4 bytes long and can be represented as an 8-digit hex or 32 bits, with values in $[0 : 4,294,967,295]$.
3. An **unsigned long** is 8 bytes and can be represented as an 16-digit hex or 64 bits, but they are only guaranteed to be stored in 32 bits in other systems.
4. An **unsigned long long** is 8 bytes and can be represented as an 16-digit hex or 64 bits, and they are guaranteed to be stored in 64 bits in other systems.

Theorem 3.1 (Bit Representation of Unsigned Integers in C)

To encode a signed integer in bits, we simply take the binary expansion of it.

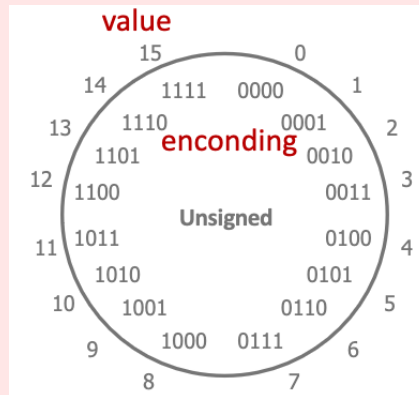


Figure 3: Unsigned encoding of 4-bit integers in C.

Example 3.1 (Bit Representation of Unsigned Integers in C)

We can see for ourselves how these numbers are represented in bits. Printing the values out in binary requires to make new functions, but we can easily convert from hex to binary.

<pre> 1 int main() { 2 3 unsigned short x = 13; 4 unsigned int y = 256; 5 6 printf("%x\n", x); 7 printf("%x\n", y); 8 9 return 0; 10 }</pre>	<pre> 1 d 2 100 3 . 4 . 5 . 6 . 7 . 8 . 9 . 10 .</pre>
--	---

So far, the process of converting unsigned numbers to bits seemed simple. Now let's introduce signed integers.

Definition 3.5 (Signed Integer Types in C)

In C, there are several signed integer types. We use this hierarchical method to give flexibility to the programmer on the size of the integer and whether it is signed or not.

1. A **signed short** is 2 bytes long and can be represented as a 4-digit hex or 16 bits, with values in $[-32,768 : 32,767]$.
2. A **signed int** is 4 bytes long and can be represented as an 8-digit hex or 32 bits, with values in $[-2,147,483,648 : 2,147,483,647]$.
3. A **signed long** is 8 bytes and can be represented as an 16-digit hex or 64 bits, but they are only guaranteed to be stored in 32 bits in other systems.
4. A **signed long long** is 8 bytes and can be represented as an 16-digit hex or 64 bits, and they are guaranteed to be stored in 64 bits in other systems.

To store signed integers, it is intuitive to simply take the first (left-most) bit and have that be the sign. Therefore, we lose one significant figure but gain information about the sign. However, this has some problems: first, there are two representations of zeros: -0 and $+0$. Second, the continuity from -1 to 0 is not natural. It is best explained through an example, which doesn't lose much insight into the general case.

Example 3.2 (Problems with the Signed Magnitude)

Say that you want to develop the signed magnitude representation for 4-bit integers in C. Then, you can imagine the following diagram to represent the numbers.

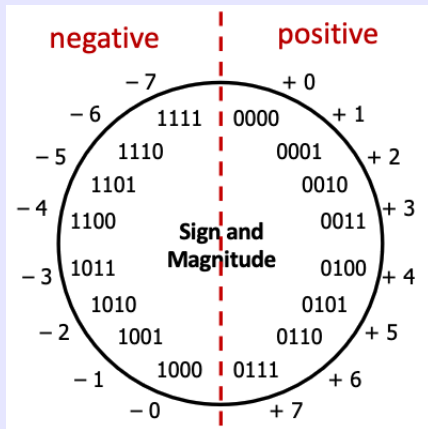


Figure 4: Signed magnitude encoding of 4-bit integers in C.

You can see that there are some problems:

1. There are two representations for 0, which is 0000 and 1000.
2. -1 (1001) plus 1 becomes -2 (1010).
3. The lowest number -7 (1111) plus 1 goes to 0 (0000) when it should go to -6 (1100).
4. The highest number 7 (0111) plus 1 goes to 0 (1000).

An alternative way is to use the two's complement representation, which solves both problems and makes it more natural.

Theorem 3.2 (Bit Representation of Signed Integers in C)

The **two's complement** representation is a way to represent signed integers in binary. It is defined as follows. Given that you want to store a decimal number p in n bits,

1. If p is positive, then take the binary expansion of that number, which should be at most $n - 1$

bits (no overflow), pad it with 0s on the left.

- If p is negative, then you can do two things: First, take the binary expansion of the positive number, flip all the bits, and add 1. Or second, represent $p = q - 2^n$, take the binary representation of q in $n - 1$ bits, and add a 1 to the left.

If you have a binary number $b = b_n b_{n-1} \dots b_1$ then to convert it to a decimal number, you simply calculate

$$q = -b_n 2^{n-1} + b_{n-1} 2^{n-2} + \dots + b_1 \tag{1}$$

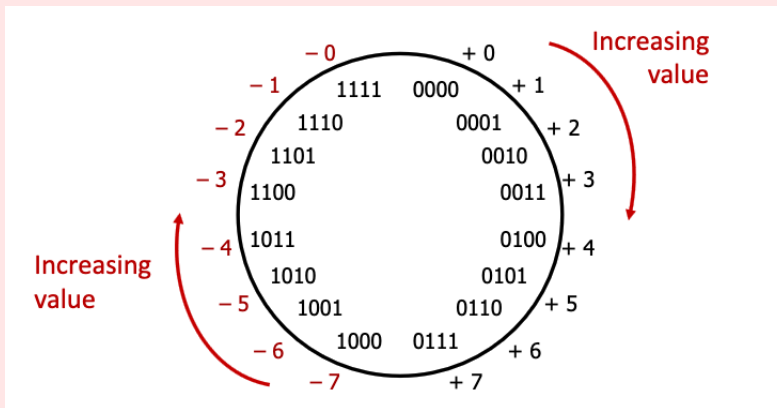


Figure 5: Two's complement encoding of 4-bit integers in C.

Example 3.3 (Bit Representation of Signed Integers in C)

We can see for ourselves how these numbers are represented in bits.

<pre> 1 int main() { 2 3 short short_pos = 13; 4 short short_neg = -25; 5 int int_pos = 256; 6 int int_neg = -512; 7 8 printf("%x\n", short_pos); 9 printf("%x\n", short_neg); 10 printf("%x\n", int_pos); 11 printf("%x\n", int_neg); 12 13 return 0; 14 }</pre>	<pre> 1 d 2 ffe7 3 100 4 ffffffe00 5 . 6 . 7 . 8 . 9 . 10 . 11 . 12 . 13 . 14 .</pre>
--	--

<pre> 1 #include<stdio.h> 2 #include<stdbool.h> 3 4 int main() { 5 printf("%lu\n", sizeof(bool)); 6 printf("%lu\n", sizeof(short)); 7 printf("%lu\n", sizeof(int)); 8 printf("%lu\n", sizeof(long)); 9 printf("%lu\n", sizeof(long long)); 10 return 0; 11 }</pre>	<pre> 1 1 2 2 3 4 4 8 5 8 6 . 7 . 8 . 9 . 10 . 11 .</pre>
--	--

Figure 6: Size of various integer types in C with the `sizeof`.

3.2.1 Arithmetic Operations on Binary Numbers

Theorem 3.3 (Inversion of Binary Numbers)

Given a binary number p , to compute $-p$, simply invert the bits and add 1.

Theorem 3.4 (Addition and Subtraction of Binary Numbers)

Given two binary numbers p and q .

1. To compute $p + q$, simply add the numbers together as you would in base 10, but carry over when the sum is greater than 1.
2. To compute $p - q$, you can invert q to $-q$ and compute $p + (-q)$.

3.3 Float Family

Definition 3.6 (Floating Point Types in C)

In C, there are several floating point types. We use this hierarchical method to give flexibility to the programmer on the size of the integer and whether it is signed or not.

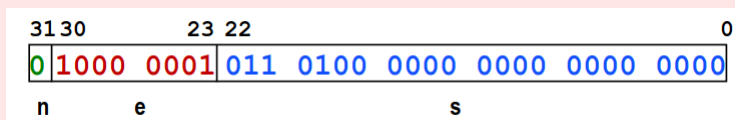
1. A **float** is 4 bytes long and can be represented as an 8-digit hex or 32 bits, with values in $[1.2 \times 10^{-38} : 3.4 \times 10^{38}]$.
2. A **double** is 8 bytes long and can be represented as an 16-digit hex or 64 bits, with values in $[2.3 \times 10^{-308} : 1.7 \times 10^{308}]$.
3. A **long double** is 8 bytes and can be represented as an 16-digit hex or 64 bits, but they are only guaranteed to be stored in 80 bits in other systems.

Theorem 3.5 (Bit Representation of Floating Point Types in C)

Floats are actually like signed magnitude. We have

$$(-1)^n \times 2^{e-127} \times 1.s \tag{2}$$

where



Doubles encode 64 bits, so not we have exponent having 11 bits (so bias is not 1023) and 52 bits for mantissa.

4 Memory

Definition 4.1 (Memory)

The **memory** is where the computer stores data and instructions, which can be thought of as a giant array of memory addresses, with each containing a byte. This data consists of graphical things or even instructions to manipulate other data. It can be visualized as a long array of boxes that each have an **address** (where it is located) and **contents** (what is stored in it).

Memory simply works as a bunch of bits in your computer with each bit having some memory address, which is also a bit. For example, the memory address `0b0010` (2) may have the bit value of `0b1` (1) stored in it.

Addresses	Values
<code>0b0010</code>	1
<code>0b0011</code>	1
<code>0b0100</code>	0
<code>0b0101</code>	1
<code>0b0110</code>	0
<code>0b0111</code>	0
<code>0b1000</code>	0
<code>0b1001</code>	1
<code>0b1010</code>	1

Figure 7: Visualization of memory as a long array of boxes of bits.

However, computers do not need this fine grained level of control on the memory, and they really work at the Byte level rather than the bit level. Therefore, we can visualize the memory as a long array of boxes indexed by *Bytes*, with each value being a byte as well. In short, the memory is **byte-addressable**. In certain architectures, some systems are **word-addressable**, meaning that the memory is addressed by words, which are 4 bytes.^a

Byte Address	Values	Values	Word Address
0x120	10010010 = 0x92	0x92006FB0	0x48
0x121	00000000 = 0x00		
0x122	01101111 = 0x6F		
0x123	10110000 = 0xB0		
0x124	10010110 = 0x96	0x96971199	0x49
0x125	10010111 = 0x97		
0x126	00010001 = 0x11		
0x127	10011001 = 0x99		
0x128	11111110 = 0xFE	0xFE....	0x4A

Figure 8: Visualization of memory as a long array of boxes of bytes. Every address is a byte and its corresponding value at that address is also a byte, though we represent it as a 2-digit hex.

^aNote that in here the size of a word is 2 bytes rather than 4 as stated above. This is just how it is defined in some x86 architectures.

In the examples above, I listed the memory addresses as a 3 hex character (1.5 bytes) for brevity. In reality, the number of bytes that a memory address takes is much longer.

Definition 4.2 (32 and 64 Bit Machines)

There are two types of machines that tend to format these boxes very differently: 32-bit and 64-bit machines.

1. 32 bit machines store addresses in 32 bits, so they can have 2^{32} addresses, which is about 4 GB of memory.
2. 64 bit machines store addresses in 64 bits, so they can have 2^{64} addresses, which is about 16 EB of memory. This does not mean that the actual RAM is 16 EB, but it means that the machine can *handle* that much memory.

```

1  ...
2  0x00007FFF7FBFF860 --> 0b0000000000000000000000001111111111
3                          111101111111101111111111100001100000
4  0x00007FFF7FBFF861 --> 0b0000000000000000000000001111111111
5                          111101111111101111111111100001100001
6  0x00007FFF7FBFF862 --> 0b0000000000000000000000001111111111
7                          111101111111101111111111100001100010
8  0x00007FFF7FBFF863 --> 0b0000000000000000000000001111111111
9                          111101111111101111111111100001100011
10 0x00007FFF7FBFF864 --> 0b0000000000000000000000001111111111
11                          111101111111101111111111100001100100
12  ...
    
```

The numbers typically mean the size of the type that the machine works best with, so all memory addresses will be 32 or 64 bits wide. Most machines are 64-bits, and so everything in this notes will assume that we are working with a 64 bit machine. As we will later see, this is why pointers are 8

bytes long, i.e. 64 bits. This is because the memory addresses are 64 bits long, though all of them are not used.

With this structure in mind and knowing the size of some primitive types, we can now focus on how declaring them works in the backend.

Definition 4.3 (Declaration, Initialization)

Assigning a value to a variable is a two step process, which is often not distinguished in high level languages like Python.

1. You must first **initialize** the variable by setting aside the correct number of bytes in memory.
2. You must then **assign** that variable to be some actual value.

The two step process is often called declaration.

This is the reason why C is statically, or strongly, typed. In order to set aside some memory for a variable, you must know how big that variable will be, which you know by its type. This makes sense. We can first demonstrate how to both initialize and declare a variable.

<pre> 1 int main() { 2 // declaring 3 int x = 4; 4 printf("%p\n", &x); 5 6 // initializing and assigning 7 int y; 8 printf("%p\n", &y); 9 y = 3; 10 printf("%p\n", &y); 11 12 return 0; 13 }</pre>	<pre> 1 0x16d37ee68 2 0x16d37ee64 3 0x16d37ee64 4 . 5 . 6 . 7 . 8 . 9 . 10 . 11 . 12 . 13 .</pre>
--	--

Figure 9: How to declare variables in C. As you can see, by initializing y, the memory address is already assigned and it doesn't change when you assign it. The address is only shown to be 9 hex digits long, but it is actually 16 hex digits long and simply 0 padded on the left.

One question that may come to mind is, what is the value of the variable if you just initialize it? After all the value at that address that is initialized must be either 0s or 1s. Let's find out.

<pre> 1 int main() { 2 int y; 3 printf("%d\n", y); 4 y = 3; 5 printf("%d\n", y); 6 7 return 0; 8 }</pre>	<pre> 1 6298576 2 3 3 . 4 . 5 . 6 . 7 . 8 .</pre>
--	---

Figure 10: The value of an uninitialized variable is some random number.

It may be interesting to see how this random uninitialized value is generated. It is simply the value that was stored in that memory address before, and it is not cleared when you initialize it, so you should not use this

as a uniform random number generator.

4.1 Debugging and Object Dumping

Talk about gdb, lldb, objdump, etc. These are debugging tools that allow you to parse your code line by line. However, to actually see the C code, you must compile it with the debugging flag. This adds a little bit of overhead memory to the binary, but not a lot.

4.2 Endian Architecture

It is intuitive to think that given some multi-byte object like an `int` (4 bytes), the beginning of the `int` would be the lowest address and the end of the `int` would be the highest address, like how consecutive integers are stored in an array. However, this is not always the case (almost always not the case since most computers are little-endian).

Definition 4.4 (Endian Architecture)

Depending on the machine architecture, computers may store these types slightly differently in their *byte* order. Say that we have an integer of value `0xA1B2C3D4` (4 bytes). Then,

1. A **big-endian architecture** (e.g. SPARC, z/Architecture) will store it so that the least significant byte has the highest address.
2. A **little-endian architecture** (e.g. x86, x86-64, RISC-V) will store it so that the least significant byte has the lowest address.
3. A **bi-endian architecture** (e.g. ARM, PowerPC) can specify the endianness as big or little.

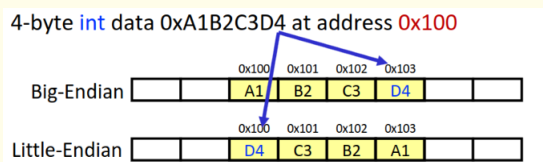


Figure 11: The big vs little endian architectures.

We can simply print out the hex values of primitive types to see how they are stored in memory, but it does not provide the level of details that we want on which bytes are stored where. At this point, we must use certain **debuggers** to directly look at the memory. For x86 architectures, we can use `gdb` and for ARM architectures, we can use `lldb`. At this point, we need to understand assembly to look through debuggers, so we will provide the example here.

Example 4.1 (Endianness of C Int in x86-64)

To do.

Example 4.2 (Endianness of C Int in ARM64)

To do.

4.3 Type Casting

4.4 Pointers

We have learned how to declare/initialize a variable, which frees up some space in the memory and possibly assigns a value to it. One great trait of C is that we can also store the memory address of a variable in another variable called a pointer. You access both the memory and the value at that memory with this pointer variable.

Definition 4.5 (Pointer Variable)

A **pointer** variable/type is a variable that stores the memory address of another variable.

1. You can declare a pointer in the same way that you declare a variable, but you must add a asterisk in front of the variable name.
2. The size of this variable is the size of the memory address, which is 8 bytes in a 64-bit architecture.
3. To get the value of the variable that the pointer points to, called **dereferencing**, you simply put a asterisk in front of the pointer. This is similar to how you put a ampersand in front of a variable to get its memory address.

<pre> 1 int main() { 2 // declare an integer 3 int x = 4; 4 printf("x = %d\n", x); 5 printf("&x = %p\n", &x); 6 7 // declare pointer 8 int *p = &x; 9 printf("p = %p\n", p); 10 printf("*p = %d\n", *p); 11 12 // initialize pointer 13 int *q; 14 q = &x; 15 printf("q = %p\n", q); 16 printf("*q = %d\n", *q); 17 return 0; 18 }</pre>	<pre> 1 x = 4 2 &x = 0x16d49ae68 3 p = 0x16d49ae68 4 *p = 4 5 q = 0x16d49ae68 6 *q = 4 7 . 8 . 9 . 10 . 11 . 12 . 13 . 14 . 15 . 16 . 17 . 18 .</pre>
--	--

Figure 12

Since the size of addresses are predetermined by the architecture, it may not seem like we need to know the underlying data type of what it points to, so why do we need to write strongly type the underlying data type? Remember that to do pointer arithmetic, you need to know how large the underlying data type is so that you can know how many bytes to move when traversing down an array.

One of the reasons why pointers are so valuable is that they allow you to pass by reference, which is a way to change the value of a variable in a function.

4.4.1 Call by Value vs Call by Reference

Definition 4.6 (Call by Value)

Definition 4.7 (Call by Reference)

4.4.2 Pointer Errors

Just like for regular variables, you may be curious on the value of an unassigned pointer. Let's take a look.

Example 4.3 (Uninitialized Pointers)

<pre> 1 int main() { 2 int x = 4; 3 int *p; 4 printf("p = %p\n", p); 5 printf("*p = %x\n", *p); 6 7 return 0; 8 }</pre>	<pre> 1 p = 0x10249ff20 2 *p = d100c3ff 3 . 4 . 5 . 6 . 7 . 8 .</pre>
--	---

Figure 13: The value of an uninitialized pointer is some random address and at a random address it would be some random byte.

This is clearly not good, especially since the program compiles correctly and runs without any errors. This kind of pointer that hasn't been initialized is called a wild pointer.

Definition 4.8 (Wild Pointer)

A **wild pointer** is a pointer that has not been initialized to a known value.

To fix this, we must always initialize a pointer to a known value. This may come at a disadvantage, since now we can't reap the benefits of initializing first and assigning later. A nice compromise is to initialize the pointer to a null pointer.

Definition 4.9 (Null Pointer)

A **null pointer** is a pointer that has been initialized to a known value, which is the address 0x0. You can set the type of the pointer and then initialize it to NULL.

<pre> 1 int main() { 2 int *p = NULL; 3 printf("p = %p\n", p); 4 5 // the code below gives seg fault 6 /* printf("*p = %d\n", *p); */ 7 8 int x = 4; 9 p = &x; 10 printf("p = %p\n", p); 11 printf("*p = %d\n", *p); 12 return 0; 13 }</pre>	<pre> 1 p = 0x0 2 p = 0x16da72e5c 3 *p = 4 4 . 5 . 6 . 7 . 8 . 9 . 10 . 11 . 12 . 13 .</pre>
---	---

Figure 14: Initializing a null pointer. It is a good practice to initialize a pointer to a null value.

Therefore, the null pointer allows you to set the type of the underlying data type, but the actual address will be 0x0. You cannot dereference a null pointer, and doing so will give you a segmentation fault. There may be times when you do not even know the data type of the pointer, and for this you can use the void pointer, which now doesn't know the type of the variable that it points to but it does allocate address.

Definition 4.10 (Void Pointer)

A **void pointer** is a pointer that does not know the type of the variable that it points to. We can initialize it by simply setting the underlying type to be void. This initializes the address, which should always be 8 bytes, but trying to access the value of the variable is not possible.

<pre> 1 int main() { 2 void *p; 3 printf("p = %p\n", p); 4 int x = 4; 5 p = &x; 6 printf("%d", *((int*)p)); 7 return 0; 8 }</pre>	<pre> 1 p = 0x102553f54 2 4 3 . 4 . 5 . 6 . 7 . 8 .</pre>
---	---

Figure 15: Initialize a void pointer and then use typecasting to access the value of the variable that it points to.

4.5 Pointer Arithmetic

With pointers out of the way, we can talk about how arrays are stored in memory.

Definition 4.11 (Array)

A C array is a collection of elements of the same type, which are stored in contiguous memory locations. You can initialize and declare arrays in many ways, and access their elements with the index, e.g. `arr[i]`.

1. You declare an array of some constant number of elements n with the elements themselves.

```
1  int arr[5] = {1, 2, 3, 4, 5};
```

2. You declare an array with out its size n and simply assign them. Then n is automatically determined.

```
1  int arr[] = {1, 2, 3, 4, 5};
```

3. You initialize an array of some constant size c , and then you assign each element of the array.

```

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

Unfortunately, C does not provide a built-in way to get the size of the array (like `len` in Python), so we must keep track of the size of the array ourselves. Furthermore, the address of the array is the address of where it begins at, i.e. the address of the first element.

You can literally see that the elements of the array are contiguous in memory by iterating through each element and printing out its address.

<pre> 1 int main(void) { 2 // initialize array 3 int arr[5]; 4 for (int val = 1; val < 6; val++) { 5 arr[val-1] = val * val; 6 } 7 8 int* p = &arr[0]; 9 for (int i = 0; i < 5; i++) { 10 printf("Value at position %d : %d\n", i, 11 arr[i]); 12 printf("Address at position %d : %p\n", 13 i, p + i); 14 } 15 return 0; 16 } </pre>	<pre> 1 Value at position 0 : 1 2 Address at position 0 : 0x7ffd8636b0d0 3 Value at position 1 : 4 4 Address at position 1 : 0x7ffd8636b0d4 5 Value at position 2 : 9 6 Address at position 2 : 0x7ffd8636b0d8 7 Value at position 3 : 16 8 Address at position 3 : 0x7ffd8636b0dc 9 Value at position 4 : 25 10 Address at position 4 : 0x7ffd8636b0e0 11 . 12 . 13 . 14 . 15 . 16 . 17 . </pre>
---	--

Figure 16: Ints are 4 bytes long, so the address of the next element is 4 bytes away from the previous element, making this a contiguous array.

The most familiar implementation of an array is a string in C.

Definition 4.12 (String)

A string is an array of characters, which is terminated by a null character `\0`. You can initialize them in two ways:

1. You can declare a string with the characters themselves, which you must make sure to end with the null character.

```
1 char str[6] = {'H', 'e', 'l', 'l', 'o', '\0'};
```

2. You can declare them with double quotes, which automatically adds the null character.

```
1 char str[] = "Hello";
```

Note that for whatever string we initialize, the size of the array is the number of characters plus 1.

To access elements of an array, you simply use the index of the element, e.g. `arr[i]`, but in the backend, this is implemented with *pointer arithmetic*.

Definition 4.13 (Pointer Arithmetic)

Pointer arithmetic is the arithmetic of pointers, which is done by adding or subtracting an integer to a pointer.

1. If you add an integer n to a pointer p , e.g. `p + n`, then the new pointer will point to the n th element after the current element, with the next element being `sizeof(type)` bytes away from the previous element.
2. If you subtract an integer n from a pointer, then the pointer will point to the n th element before the current element.

This is why you can access the elements of an array with the index, since the index is simply the number of elements away from the first element.

Example 4.4 (Pointer Arithmetic with Arrays of Ints and Chars)

Ints have a size of 4 bytes and chars 1 byte. You can see that using pointer arithmetic, the addresses of the elements of ints increment by 4 and those of the char array increment by 1.

<pre> 1 int main() { 2 int integers[3] = {1, 2, 3}; 3 char characters[3] = {'a', 'b', 'c'}; 4 int *p = &integers[0]; 5 char *q = &characters[0]; 6 7 printf("Array of Integers\n"); 8 for (int i = 0; i < 3; i++) { 9 printf("%p\n", integers+i); } 10 11 printf("Array of Characters\n"); 12 for (int i = 0; i < 3; i++) { 13 printf("%p\n", characters+i); } 14 return 0; 15 }</pre>	<pre> 1 Array of Integers 2 0x16d39ee58 3 0x16d39ee5c 4 0x16d39ee60 5 . 6 Array of Characters 7 0x16d39ee50 8 0x16d39ee51 9 0x16d39ee52 10 . 11 . 12 . 13 . 14 . 15 .</pre>
--	--

Therefore, we can think of accessing the elements of an array as simply pointer arithmetic.

Theorem 4.1 (Bracket Notation is Pointer Arithmetic)

The bracket notation is simply pointer arithmetic in the backend.

<pre> 1 int main() { 2 int arr[3] = {1, 2, 3}; 3 int *p = &arr[0]; 4 5 for (int i = 0; i < 3; i++) { 6 printf("%d\n", arr[i]); 7 printf("%d\n", *(p+i)); 8 } 9 return 0; 10 }</pre>	<pre> 1 1 2 1 3 2 4 2 5 3 6 3 7 . 8 . 9 . 10 .</pre>
--	---

Figure 17: Accessing the elements of the list using both ways is indeed the same.

4.6 Global, Stack, and Heap Memory

Everything in a program is stored in memory, variables, functions, and even the code itself. However, we will find out that they are stored in different parts of the memory. When a program runs, its application memory consists of four parts, as visualized in the Figure 18.

1. The **code** is where the code text is stored.
2. The **global memory** is where all the global variables are stored.
3. The **stack** is where all of the functions and local variables are stored.
4. The **heap** is variable and can expand to as much as the RAM on the current system. We can specifically store whatever variables we want in the heap.

We provide a visual of these four parts first, and we will go into them later.

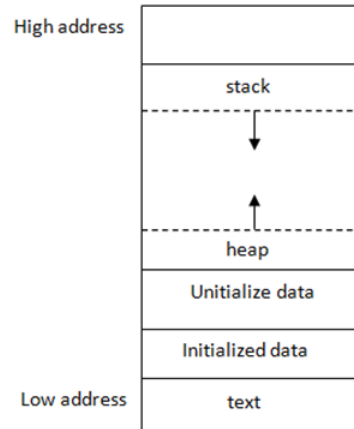


Figure 18: The four parts of memory in a C program.

Definition 4.14 (Code Memory)

This is where the code text is stored. It is read-only and is not modifiable.

In high level languages, we always talk about local and global scope. That is, variables defined within functions have a local scope in the sense that anything we modify in the local scope does not affect the global scope. We can now understand what this actually means by examining the backend. The global scope variables are stored in the global memory, and all local variables (and functions) are stored in the stack.

Definition 4.15 (Global Memory)

This is where all the global variables are stored.

Definition 4.16 (Stack Memory)

This is where all of the functions and local variables are stored. As we will see later, the compiler will always run the main function, which must exist in your file. By the main function is a function itself, and therefore it has its own local scope.

Then, when you initialize any functions or local variables within those functions (which will be the majority of your code), all these will be stored in the stack, which is an literally an implementation of the stack data structure. It is LIFO, and the first thing that goes in is the `main` function and its local variables, which is referred to as the **stack frame**. You can't free memory in the stack unless its in the top of the stack.

To see what happens in the stack, we can go through an example.

Example 4.5 (Going through the Stack)

Say that you have the following code:

```

1  int total;
2  int Square(int x) {
3      return x*x;
4  }
5  int SquareOfSum(int x, int y) {
6      int z = Square(x + y);
7      return z;

```

```
8  }
9  int main() {
10     int a = 4, b = 8;
11     total = SquareOfSum(a, b);
12     printf("output = %d", total);
13     return 0;
14 }
```

The memory allocation of this program will run as such:

1. The `total` variable is initialized and is put into global memory.
2. `main` is called. It is put into the stack.
3. The local variables `a=4` and `b=8` are initialized and are put into the stack.
4. The `SquareOfSum` function is called and put into the stack.
5. The input local variables `x=4`, `y=8`, `z` are initialized and put into the stack.
6. `x + y=12` is computed and put into the stack.
7. The `Square` function is called and put into the stack.
8. The `x=12` local variable of `Square` is initialized and put into the stack.
9. The CPU computes `x*x=144` and returns the output. The `Square` function is removed from the stack.
10. We assign `z=144` and `SquareOfSum` returns it. Now `SquareOfSum` is removed from the stack.
11. `total=144` is assigned in the global memory still.
12. The `printf` function is called and put into the stack.
13. The `printf` function prints the output and is removed from the stack.
14. The `main` function returns 0 and is removed from the stack, ending our application.

One limitation of the stack is that its total available memory is fixed from the start, ranging from 1MB to 8MB, and so you can't initialize arrays of billions of integers in the stack. It will cause a memory overflow. In fact, the memory of the stack, along with the global and text memory, are assigned at compile time, making it a **static memory**.

Since the stack is really just a very small portion of available memory, the heap comes into rescue, which is the pool of memory available to you in RAM.

Definition 4.17 (Heap Memory)

The **heap memory** (nothing to do with the heap data structure) is a variable length (meaning it can grow at runtime) and **dynamically allocated** (meaning that we can assign memory addresses during runtime) memory that is limited to your computer's hardware. Unlike simply initializing variables to allocate memory as in the stack, we must use the **malloc** and **free** functions in C, and **new** and **delete** operations in C++.

Definition 4.18 (malloc)

Definition 4.19 (free)

The stack can store pointer variables that point to the memory address in the heap. So the only way to access variables in the heap is through pointer reference, and the stack provides you that window to access that big pool of heap memory.

One warning: if you allocate another address, the previous address does not get deallocated off the memory.

Definition 4.20 (Memory Leak)

On the other hand, if you free an address but have a pointer still pointing to that address, this is also a problem called the dangling pointer.

Definition 4.21 (Dangling Pointer)

At this point, we might be wondering why we need both a stack and a heap. Well the benefits of heaps are clearer since you can dynamically allocate memory, and you don't have the LIFO paradigm that is blocking you from deallocating memory that has been allocated in the beginning of your program. A problem with just having heap is that stacks can be orders of magnitude times faster when allocating/deallocating from it than the heap, and the sequence of function calls is naturally represented as a stack.

4.7 Dynamic Memory Allocation

Let's talk about how `malloc` and `free` are implemented in C. If you make a for loop and simply print all the addresses that you allocate to. You will find that they can be quite random. After a program makes some calls to `malloc` and `free`, the heap memory can become fragmented, meaning that there are chunks of free heap space interspersed with chunks of allocated heap space. The heap memory manager typically keeps lists of different ranges of sizes of heap space to enable fast searching for a free extent of a particular size. In addition, it implements one or more policies for choosing among multiple free extents that could be used to satisfy a request.

The `free` function may seem odd in that it only expects to receive the address of the heap space to free without needing the size of the heap space to free at that address. That's because `malloc` not only allocates the requested memory bytes, but it also allocates a few additional bytes right before the allocated chunk to store a header structure. The header stores metadata about the allocated chunk of heap space, such as the size. As a result, a call to `free` only needs to pass the address of heap memory to free. The implementation of `free` can get the size of the memory to free from the header information that is in memory right before the address passed to `free`.

5 Implementations of Memory Structures in C

5.1 Arrays

5.2 Strings

5.3 Structs

5.4 Functions

5.5 Classes (for C++)

5.6 Input Output

We have `standard in`, `standard out`, and `standard error`.

6 Central Processing Unit

Now let's talk about how functions work on a deeper level. When we write a command, like `int x = 4`, we are manually looking for an address (in the stack, global, or heap) and rewriting the bits that are at that address. Functions are just an automated way to do this, and all these modifications and computations are done by the CPU.

Definition 6.1 (Central Processing Unit)

The CPU is responsible for taking instructions (data) from memory and executing them.

1. The CPU is composed of **registers** (different from the cache), which are small, fast storage locations. These registers can either be **general purpose** (can be used with most instructions) or **special purpose** (can be accessed through special instructions, or have special meanings/uses, or are simply faster when used in a specific way).
2. The CPU also has an **arithmetic unit** and **logic unit**, which is responsible for performing arithmetic and logical operations.
3. The CPU also has a **control unit**, which is responsible for fetching instructions from memory through the **databus**, which is literally a wire connecting the CPU and RAM, and executing them.

It executes instructions from memory one at a time and executes them, known as the **fetch-execute cycle**. It consists of 4 main operations.

1. **Fetch:** The **program counter**, which holds the memory address of the next instruction to be executed, tells the control unit to fetch the instruction from memory through the databus.
2. **Decode:** The fetched data is passed to the **instruction decoder**, which figures out what the instruction is and what it does and stores them in the registers.
3. **Execute:** The arithmetic and logic unit then carries out these operations.
4. **Store:** Then it puts the results back on the databus, and stores them back into memory.

The CPU's **clock cycle** is the time it takes for the CPU to execute one instruction. More specifically, the clock cycle refers to a single oscillation of the clock signal that synchronizes the operations of the processor and the memory (e.g. fetch, decode, execute, store), and decent computers have clock cycles of at least 2.60GHz (2.6 billion clock cycles per second).

Therefore, in order to actually do computations on the data stored in the memory, the CPU must first fetch the data, perform the computations, and then store the results back into memory. This can be done in two ways.

1. **Load and Store Operations:** CPUs use load instructions to move data from memory to registers (where operations can be performed more quickly) and store instructions to move the modified data back into memory.
2. If the data is too big to fit into the registers, the CPU will use the **cache** to store the data, and in worse cases, the actual memory itself. Compilers optimize code by maximizing the use of registers for operations to minimize slow memory access. This is why you often see assembly code doing a lot in registers.

To clarify, let us compare registers and memory. Memory is addressed by an unsigned integer while registers have names like `%rsi`. Memory is much bigger at several GB, while the total register space is much smaller at around 128 bytes (may differ depending on the architecture). The memory is much slower than registers, which is usually on a sub-nanosecond timescale. The memory is dynamic and can grow as needed while the registers are static and cannot grow.

The specific structure/architecture of the CPU is determined by the instruction set architecture (ISA), which can be thought of as a subset of the general computer architecture.

Definition 6.2 (Instruction Set Architecture)

The **ISA** or just **architecture** of a CPU is a high level description of what it can do. Some differences are listed here:

1. What instructions it can execute.
2. The instruction length and decoding, along with its complexity.
3. The performance vs power efficiency.

ISAs can be classified into two types.

1. The **complex instruction set computer** (CISC) is characterized by a large set of complex instructions, which can execute a variety of low-level operations. This approach aims to reduce the number of instructions per program, attempting to achieve higher efficiency by performing more operations with fewer instructions.
2. The **reduced instruction set computer** (RISC) emphasizes simplicity and efficiency with a smaller number of instructions that are generally simpler and more uniform in size and format. This approach facilitates faster instruction execution and easier pipelining, with the philosophy that simpler instructions can provide greater performance when optimized.

Just like how memory addressing is different between 32 and 64 bit machines, CPUs also use these schemes. While 32-bit processors have 2^{32} possible addresses in their cache, it turns out that 64-bit processors have a 48-address space. This is because CPU manufacturers took a shortcut. They use an instruction set which allows a full 64-bit address space, but current CPUs just only use the last 48-bits. The alternative was wasting transistors on handling a bigger address space which wasn't going to be needed for many years (since 48-bits is about 256TB). Just a bit of history for you. Finally, just to briefly mention, the input/output device, as the name suggests, processes inputs and displays outputs, which is how you can see what the program does.

Example 6.1 (x86 Architecture)

The x86 architecture is a CISC architecture, which is the most common architecture for personal computers. Here are important properties:

1. It is a complex instruction set computer (CISC) architecture, which means that it has a large set of complex instructions^a.
2. Byte-addressing is enabled and words are stored in little-endian format.
3. In the x86_64 architecture, registers are 8 bytes long (and 4 bytes in x86_32) and there are 16 total general purpose registers, for a total of only 128 bytes (very small compared to many GB of memory). Other special purpose registers are also documented in the wikipedia page, but it is not fully documented.

^ahttps://en.wikipedia.org/wiki/X86_instruction_listings

Example 6.2 (ARM Architecture)

Mainly in phones, tablets, laptops.

Example 6.3 (MIPS Architecture)

MIPS is a RISC architecture, which is used in embedded systems such as digital home and networking equipment.

Definition 6.3 (Input/Output Device)

The input device can read/load/write/store data from the outside world. The output device, which has **direct memory address**, can display data to the outside world.

One final note to mention, there are many assembly languages out there and various syntaxes.

Example 6.4 (Assembly Syntax)

The two most popular syntaxes are AT&T and Intel.

1. **Intel Syntax:** Specifies memory operands without any special prefixes. Square brackets [] are used to denote memory addresses. For example, `mov eax, [ebx]` means move the contents of the

memory location pointed to by ebx into eax.

2. **AT&T Syntax:** Memory operands are denoted with parentheses () and include the % prefix for registers. An instruction moving data from a memory location into a register might look like `movl (%ebx), %eax`, with additional prefixes for immediate values and segment overrides.

Example 6.5 (Assembly Languages)

The various assembly languages are as follows:

1. **x86 Assembly :** The assembly language for Intel and AMD processors using the x86 architecture. Both AT&T and Intel syntax are available. Tools or environments often allow switching between the two, with AT&T being the default in GNU tools like GDB.
2. **ARM Assembly :** The assembly language for ARM processors. Has its own unique syntax, not categorized as AT&T or Intel. ARM syntax is closely tied to its instruction set architecture and is distinct from the x86 conventions.
3. **MIPS Assembly :** The assembly language for MIPS processors. MIPS uses its own assembly language syntax, which is neither AT&T nor Intel. MIPS syntax is designed around the MIPS instruction set architecture.
4. **PowerPC Assembly :** The assembly language for PowerPC processors. PowerPC has its own syntax style, tailored to its architecture and instruction set, distinct from the AT&T and Intel syntax models.
5. **6502 Assembly :** Used in many early microcomputers and gaming consoles. Utilizes a syntax unique to the 6502 processor, not following AT&T or Intel conventions.
6. **AVR Assembly :** The assembly language for Atmel's AVR microcontrollers. AVR assembly follows its own syntax style, designed specifically for AVR microcontrollers and not based on AT&T or Intel syntax.
7. **Z80 Assembly :** Associated with the Z80 microprocessor, used in numerous computing devices in the late 20th century. Z80 assembly language has its own syntax that does not adhere to AT&T or Intel syntax guidelines.

The most common one is the x86_64, which is the one that we will be focusing on, with the AT&T syntax.

6.1 Circuits

Let's go over some common logic gates since this is at the basis of how to construct arithmetic operations.

Definition 6.4 (AND, NOT, OR)

Definition 6.5 (XOR, NAND, NOR)

Definition 6.6 (NAND)

Talk about how to construct arithmetic operations with these gates such as adding two integers or multiplying them, and not just that, but other operations that we may need in a programming language.

Theorem 6.1 (Implementation of Moving Data in Circuits)

Theorem 6.2 (Implementation of Addition, Subtraction in Circuits)**Theorem 6.3 (Implementation of Multiplication in Circuits)****Theorem 6.4 (Implementation of Bitwise Operations in Circuits)****Theorem 6.5 (Implementation of Bitshift Operations)**

We also want some sort of conditionals. This then can be used to implement loops by checking some conditional.

Theorem 6.6 (Implementation of Conditionals in Circuits)

As a bonus, we talk about the difference between volatile and non-volatile memory. We already learned that RAM is volatile, and this is simple to implement in a circuit since we can manually set all the bits to 0 or just deplete all power. If this is the case, then how does non-volatile memory like SSDs maintain their state?

Theorem 6.7 (Implementation of Volatile Memory)**Theorem 6.8 (Implementation of Non-Volatile Memory)**

6.2 Registers

To understand anything that the CPU does, we must understand assembly language. In here, everything is done within registers, and we can see how the CPU fetches, decodes, and executes instructions. So what exactly are these registers?

Definition 6.7 (Register)

A register is a small, fast storage location within the CPU. It is used to store data that is being used immediately, and is the only place where the CPU can perform operations, which is why it must move data from memory to registers before it can perform operations on it. Everything in a register is in binary, at most 8 bytes, or 64 bits.

There are very specific types of registers that you should know. All of these registers are implemented for all assembly languages and are integral to the workflow of the CPU.

1. **parameter registers** which store the parameters of a function.
2. **Return registers** which store return values of functions.
3. **stack pointers** which point to the top of the stack (at the top of the current stack frame).
4. **frame pointers** which point to the base of the current stack frame.
5. **instruction pointers** which point to the next instruction to be executed.

6.2.1 x86 Assembly Registers

The specific type of registers that are available to a CPU depends on the computer architecture, or more specifically, the ISA, but here is a list of common ones for the x86-64. We have `%rax`, `%rbx`, `%rcx`, `%rdx`, `%rsi`, `%rdi`, `%rbp`, `%rsp`, `%r8`, `%r9`, `%r10`, `%r11`, `%r12`, `%r13`, `%r14`, `%r15`. Therefore, the x86-64 Intel CPU has a total of 16 registers for storing 64 bit data. However, it is important to know which registers are used for what.

Definition 6.8 (Parameter Registers)

Compilers typically store the first six parameters of a function in registers

$$\%rdi, \%rsi, \%rdx, \%rcx, \%r8, \%r9, \quad (3)$$

respectively.

Definition 6.9 (Return Register)

The return value of a function is stored in the

$$\%rax \quad (4)$$

register.

Definition 6.10 (Stack and Frame Pointers)

The `%rsp` register is the **stack pointer**, which points to the top of the stack. The `%rbp` register is the **frame pointer**, or **base pointer**, which points to the base of the current stack frame. In a typical function prologue, `%rbp` is set to the current stack pointer (`%rsp`) value, and then `%rsp` is adjusted to allocate space for the local variables of the function. This establishes a fixed point of reference (`%rbp`) for accessing those variables and parameters, even as the stack pointer (`%rsp`) moves.

Definition 6.11 (Instruction Pointer)

The `%rip` register is the **instruction pointer**, which points to the next instruction to be executed. Unlike all the registers that we have shown so far, programs cannot write directly to `%rip`.

Definition 6.12 (Notation for Accessing Lower Bytes of Registers)

Sometimes, we need a more fine grained control of these registers, and x86-64 provides a way to access the lower bits of the 64 bit registers. We can visualize them with the diagram below.

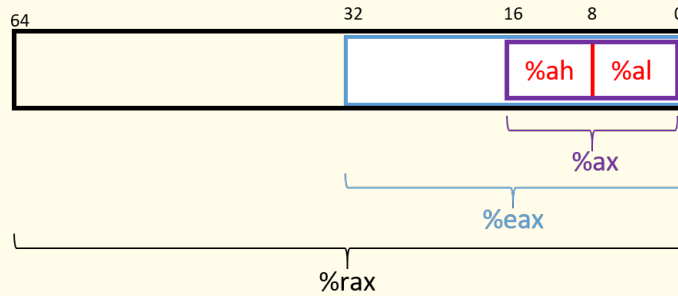


Figure 19: The names that refer to subsets of register `%rax`.

A complete list is shown below.

64-bit Register	32-bit Register	Lower 16 Bits	Lower 8 Bits
<code>%rax</code>	<code>%eax</code>	<code>%ax</code>	<code>%al</code>
<code>%rbx</code>	<code>%ebx</code>	<code>%bx</code>	<code>%bl</code>
<code>%rcx</code>	<code>%ecx</code>	<code>%cx</code>	<code>%cl</code>
<code>%rdx</code>	<code>%edx</code>	<code>%dx</code>	<code>%dl</code>
<code>%rdi</code>	<code>%edi</code>	<code>%di</code>	<code>%dil</code>
<code>%rsi</code>	<code>%esi</code>	<code>%si</code>	<code>%sil</code>
<code>%rsp</code>	<code>%esp</code>	<code>%sp</code>	<code>%spl</code>
<code>%rbp</code>	<code>%ebp</code>	<code>%bp</code>	<code>%bpl</code>
<code>%r8</code>	<code>%r8d</code>	<code>%r8w</code>	<code>%r8b</code>
<code>%r9</code>	<code>%r9d</code>	<code>%r9w</code>	<code>%r9b</code>
<code>%r10</code>	<code>%r10d</code>	<code>%r10w</code>	<code>%r10b</code>
<code>%r11</code>	<code>%r11d</code>	<code>%r11w</code>	<code>%r11b</code>
<code>%r12</code>	<code>%r12d</code>	<code>%r12w</code>	<code>%r12b</code>
<code>%r13</code>	<code>%r13d</code>	<code>%r13w</code>	<code>%r13b</code>
<code>%r14</code>	<code>%r14d</code>	<code>%r14w</code>	<code>%r14b</code>
<code>%r15</code>	<code>%r15d</code>	<code>%r15w</code>	<code>%r15b</code>

Table 1: Register mapping in x86-64 architecture

6.2.2 ARM Assembly Registers

6.3 Addressing Modes

Registers being 8 bytes mean that we can store memory addresses, and if we can store memory addresses, we can access memory, i.e. the values at those memory addresses. There are 4 ways to do this, called **addressing modes**: immediate, normal, displacement, and indexed. When we parse an instruction, its operands are either

1. Constant (literal) values
2. Registers
3. Memory forms

Definition 6.13 (Immediate Addressing)

Immediate addressing is when the operand is a constant value, used with a \$ sign.

$$\text{\$val} \tag{5}$$

Definition 6.14 (Normal Addressing)

Normal addressing is when the operand is a register, used with a % sign and the following syntax. The parentheses are used to dereference the memory address like dereferencing a pointer in C.

$$(\text{R}) = \text{Mem}[\text{Reg}[\text{R}]] \tag{6}$$

where R is the register name, Reg[R] is the value in the register, and Mem[Reg[R]] is the value in the memory address pointed to by the register.

Definition 6.15 (Displacement Addressing)

When we have a memory address stored in a register, we can add an offset to it to access a different memory address.

$$D(\text{R}) = \text{Mem}[\text{Reg}[\text{R}] + D] \tag{7}$$

where R is the register name and D is a constant displacement that specifies offset.

Definition 6.16 (Indexed Addressing)

Indexed addressing gives us more flexibility, allowing us to multiply the value in the register by a constant and add it to the value in another register. The general formula is shown as the top, but there are special cases:

$$\begin{aligned} D(\text{Rb}, \text{Ri}, \text{S}) &= \text{Mem}[\text{Reg}[\text{Rb}] + \text{S} * \text{Reg}[\text{Ri}] + D] \\ D(\text{Rb}, \text{Ri}) &= \text{Mem}[\text{Reg}[\text{Rb}] + \text{Reg}[\text{Ri}] + D] \\ (\text{Rb}, \text{Ri}, \text{S}) &= \text{Mem}[\text{Reg}[\text{Rb}] + \text{S} * \text{Reg}[\text{Ri}]] \\ (\text{Rb}, \text{Ri}) &= \text{Mem}[\text{Reg}[\text{Rb}] + \text{Reg}[\text{Ri}]] \\ (, \text{Ri}, \text{S}) &= \text{Mem}[\text{S} * \text{Reg}[\text{Ri}]] \end{aligned}$$

where D is a constant displacement of 1, 2, or 4 bytes, Rb is the base register (can be any of 8 integer registers), Ri is the index register (can be any register except `rsp`), and S is the scale factor (1, 2, 4, or 8).

6.3.1 x86 Assembly Addressing Modes

Example 6.6 (Immediate Addressing)

```
1 movq $0x4, %rax
```

Example 6.7 (Normal Addressing)

The following example shows the source operand being a memory address, with normal addressing, and the destination operand being a register.

```
1 movq (%rax), %rbx
```

Example 6.8 (Displacement Addressing)

The following example shows the source operand being a memory address and the destination operand being a register. They are both addressed normally.

```
1 movq 8(%rdi), %rdx
```

Example 6.9 (Indexed Addressing)

The following shows the source operand being a memory address and the destination operand being a register. Say that `%rdx = 0xf000` and `%rcx = 0x0100`. Then

$$0x80(,%rdx,2) = \text{Mem}[2*0xF000 + 0x80] = \text{Mem}[0x1E080] \quad (8)$$

We see that

```
1 movq 0x100(%rdi, %rsi, 8), %rdx
```

6.3.2 ARM Assembly Addressing Modes**6.4 Instructions**

Now that we've gotten a sense of what these registers are and some commonalities between them, let's do some operations on them with instructions.

Definition 6.17 (Instruction)

An instruction is a single line of assembly code. It consists of some instruction followed by its (one or more) operands. The instruction is a mnemonic for a machine language operation (e.g. `mov`, `add`, `sub`, `jmp`, etc.). The **size specifier** can be appended to this instruction mnemonic to specify the size of the operands.

1. **b** (byte) for 1 byte
2. **w** (word) for 2 bytes
3. **l** (long) for 4 bytes
4. **q** (quad word) for 8 bytes

Note that due to backwards compatibility, word means 2 bytes in instruction names. Furthermore, the maximum size is 8 bytes since that is the size of each register in `x86_64`. An operand can be of 3 types, determined by their **mode of access**:

1. **Immediate addressing** is denoted with a `$` sign, e.g. a constant integer data `$1`.
2. **Register addressing** is denoted with a `%` sign with the following register name, e.g. `%rax`.
3. **Memory addressing** is denoted with the hexadecimal address in memory, e.g. `0x034AB`.

Like higher level programming languages, we can perform operations, do comparisons, and jump to different parts of the code. Instructions can be generally categorized into three types:

1. **Data Movement:** These instructions move data between memory and registers or between the registry and registry. Memory to memory transfer cannot be done with a single instruction.

```

1 %reg = Mem[address]    # load data from memory into register
2 Mem[address] = %reg    # store register data into memory

```

2. **Arithmetic Operation:** Perform arithmetic operation on register or memory data.

```

1 %reg = %reg + Mem[address]    # add memory data to register
2 %reg = %reg - Mem[address]    # subtract memory data from register
3 %reg = %reg * Mem[address]    # multiply memory data to register
4 %reg = %reg / Mem[address]    # divide memory data from register

```

3. **Control Flow:** What instruction to execute next both unconditional and conditional (if statements) ones. With if statements, loops can then be defined.

```

1 jmp label    # jump to label
2 je label     # jump to label if equal
3 jne label    # jump to label if not equal
4 jg label     # jump to label if greater
5 jl label     # jump to label if less
6 call label   # call a function
7 ret         # return from a function

```

Now unlike compiled languages, which are translated into machine code by a compiler, assembly code is translated into machine code through a two-step process. First, we **assemble** the assembly code into an **object file** by an **assembler**, and then we **link** the object file into an executable by a **linker**. Some common assemblers are **NASM** (Netwide Assembler) and **GAS/AS** (GNU Assembler), and common linkers are **ld** (GNU Linker) and **lld** (LLVM Linker), both installable with **sudo pacman -S nasm ld**.

6.4.1 Moving and Arithmetic

Again, it is more important to have a general feel of what instructions every assembly language should and get the ideas down rather than the syntax. We list them here, beginning with simply moving.

Definition 6.18 (Moving)

Next we want to have some sort of arithmetic to do calculations and to compare values.

Definition 6.19 (Arithmetic Operations)

6.4.2 Conditionals

Definition 6.20 (Conditionals)

6.4.3 Control Transfer on Stack

These are really the three basic functions needed to do anything in assembly, but let's talk about an important implementation called the **control transfer**. Say that you want to compute a function.

1. Then we must retrieve the data from the memory.

2. We must load it into our registers in the CPU and perform some computation.
3. Then we must store the data back into memory.

Let's begin with a refresher on how the call stack is managed. Recall that `%rsp` is the stack pointer and always points to the top of the stack. The register `%rbp` represents the base pointer (also known as the frame pointer) and points to the base of the current stack frame. The stack frame (also known as the activation frame or the activation record) refers to the portion of the stack allocated to a single function call. The currently executing function is always at the top of the stack, and its stack frame is referred to as the active frame. The active frame is bounded by the stack pointer (at the top of stack) and the frame pointer (at the bottom of the frame). The activation record typically holds local variables for a function.

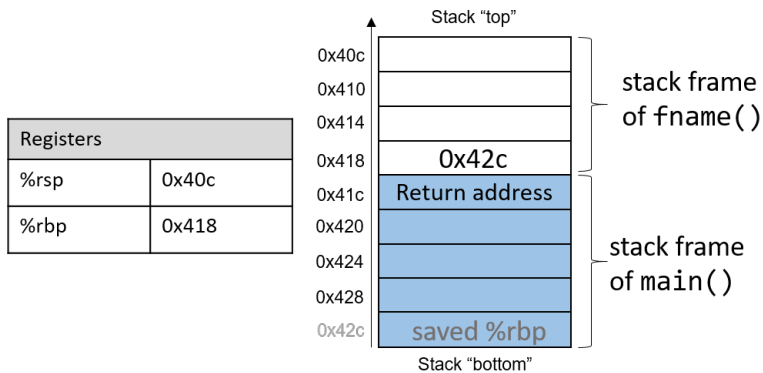


Figure 20: The current active frame belongs to the callee function (`fname`). The memory between the stack pointer and the frame pointer is used for local variables. The stack pointer moves as local values are pushed and popped from the stack. In contrast, the frame pointer remains relatively constant, pointing to the beginning (the bottom) of the current stack frame. As a result, compilers like GCC commonly reference values on the stack relative to the frame pointer. In Figure 1, the active frame is bounded below by the base pointer of `fname`, which is stack address `0x418`. The value stored at address `0x418` is the "saved" `%rbp` value (`0x42c`), which itself is an address that indicates the bottom of the activation frame for the main function. The top of the activation frame of `main` is bounded by the return address, which indicates where in the main function program execution resumes once the callee function `fname` finishes executing.

Once we have done this we are really done. Formally, this is called Turing complete (?).

Definition 6.21 (Control Transfers)

We list some.

1. Push
2. Pop
3. Call to call a function
4. Return to return from a function
5. Continue
6. Get out of stack with leave.

Example 6.10 (Control Transfer Example)

We show this with a minimal example with psuedocode.

6.4.4 Multiple Functions

Now what happens if there are multiple functions calling each other? Take a look at the following example with two functions.

Example 6.11 (Multiple Functions Example)

There is a bit of a concern here from the previous example. The main function had two functions that returned two values. As the subfunction stack frame is removed from the stack, the return value is stored in the `%rax` register. If another function is called right after, then the return value of the second function will overwrite that of the previous one. This was not a problem in the previous example since the return value of the `assign` function was not used. However, if it was, then the return value of the `adder` function would have overwritten it. This is known as register saving.

1. For **caller-saved registers**, the caller function is responsible for saving the value of the register before calling a function and restoring it after the function returns. The caller should save values in its stack frame before calling the callee function, e.g. by pushing all the return values of each callee in the caller stack frame. Then it will restore values after the call.

Therefore, if we have a set of registers $\{\%reg\}$, the caller must take everything and push them in the caller stack frame. Then it will restore them after the call.

2. For **callee-saved registers**, it is the callee's responsibility to save any data in these registers before using the registers.

Therefore, if we have a set of registers $\{\%reg\}$, then inside the callee stack frame, the callee must take everything and push them in the callee stack frame. Once it computes the final return value, then it will restore all the saved register values from the callee stack frame back into the registers for the caller to use.

Ideally, we want *one* calling convention to simply separate implementation details between caller and callee. In general, however, neither is best. If the caller isn't using a register, then caller-save is better, and if callee doesn't need a register, then callee-save is better. If we do need to save, then callee save generally makes smaller programs, so we compromise and use a combination of both caller-save and callee-save.

6.4.5 x86-64 Instructions

Let's talk about moving instructions first.

Definition 6.22 (mov)

Let's talk about the `mov` instruction which copies data from the source to the destination (the data in the source still remains!) and has the syntax

$$\text{mov_src, dest} \quad (9)$$

1. The source can be a register (`%rsi`), a value (`$0x4`), or a memory address (`0x4`).
2. The destination can be a register or a memory address.
3. The `_` is defined to be one of the size operands, which determine how big the data is. For example, we can call `movq` to move 8 bytes of data (which turns out to be the maximum size of a register).

A good diagram to see is the following:

	Source	Dest	Src, Dest	C Analog
movq	Imm	Reg	movq \$0x4, %rax	var_a = 0x4;
		Mem	movq \$-147, (%rax)	*p_a = -147;
	Reg	Reg	movq %rax, %rdx	var_d = var_a;
		Mem	movq %rax, (%rdx)	*p_d = var_a;
	Mem	Reg	movq (%rax), %rdx	var_d = *p_a;

Even with just the mov instruction, we can look at a practical implementation of a C program in Assembly.

Example 6.12 (Swap Function)

Let us take a look at a function that swaps two integers. Let's see what they do.

1. In C, we dereference both `xp` and `yp` (note that they are pointers to longs, so they store 8 bytes), and assign these two values to two temporary variables. Then, we assign the value of `yp` to `xp` and the value of `xp` to `yp`.
2. In Assembly, we first take the registers `%rdi` and `%rsi`, which are the 1st and 2nd arguments of the function, dereference them with the parantheses, and store them in the temporary registers `%rax` and `%rdx`. Then, we store the value of `%rdx` into the memory address of `%rdi` and the value of `%rax` into the memory address of `%rsi`. Note that the input values (the actual of)

```

1 void swap(long *xp, long *yp) {
2     long t0 = *xp;
3     long t1 = *yp;
4     *xp = t1;
5     *yp = t0;
6 }
    
```

```

1 swap:
2     movq (%rdi), %rax
3     movq (%rsi), %rdx
4     movq %rdx, (%rdi)
5     movq %rax, (%rsi)
6     ret
    
```

Definition 6.23 (movz and movs)

The `movz` and `movs` instructions are used to move data from the source to the destination, but with zero and sign extension, respectively. It is used to copy from a smaller source value to a larger destination, with the syntax

```

movz__ src, dest
movs__ src, dest
    
```

where the first `_` is the size of the source and the second `_` is the size of the destination.

1. The source can be from a memory or register.
2. The destination must be a register.

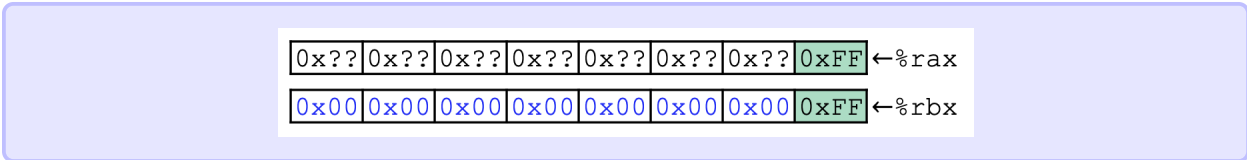
Example 6.13 (Simple example with movz)

Take a look at the code below.

```

1 movz bq %al, %rbx
    
```

The `%al` represents the last byte of the `%rax` register. It is 1 byte long. The `%rbx` register is 8 bytes long, so we can fill in the rest of the 7 bytes with zeros.

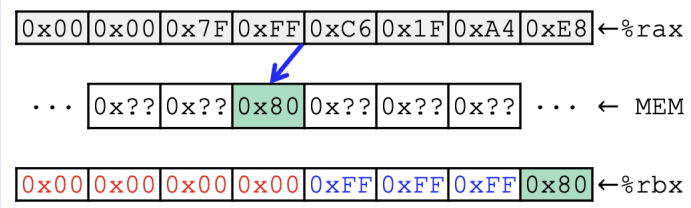


Example 6.14 (Harder example with movs)

Take a look at the code below.

```
1  movsbl (%rax), %ebx
```

You want to move the value at the memory address in %rax into %ebx. Since the source size is set to 1 byte, you take that byte, say it is 0x80, from the memory, and then sign extend (by a size of 4 bytes!) into %ebx. Note that therefore, the first four bytes of %rbx will not be affected since it's not a part of %ebx. An exception to this is that in x86-64, any instruction that generates a 32-bit long word value for a register also sets the high-order 32 bits of the register to 0, so this ends up clearing the first 4 bytes to 0.



Now we can talk about control transfer. Say that you have the following C and Assembly code.

<pre>1 int add(int x) { 2 return x + 2; 3 } 4 5 int main() { 6 int a = 2; 7 int b = add(a); 8 return 0; 9 }</pre>	<pre>1 add: 2 movq %rdi, %rax 3 addq \$2, %rax 4 ret 5 main: 6 movq \$3, %rdi 7 call add 8 movq \$0, %rax 9 ret</pre>
---	--

Figure 21: A simple function.

If you go through the instructions, you see that in main, you first move \$3 into the %rdi register. Then, you call the add function, and within it you also have the %rdi register. This is a conflict in the register, and we don't want to simply overwrite the value of %rdi in the main function. Simply putting it to another register isn't a great idea since we can't always guarantee that it will be free. Therefore, we must use the memory itself.

Recall the stack, which we can think of as a giant array in which data gets pushed and popped in a last-in-first-out manner. The stack is used to store data and return addresses, and is used to manage function calls. Visually, we want to think of the elements getting pushed in from the bottom (upside down) towards lower memory addresses.

Definition 6.24 (Stack Pointer)

Note that every time we want to push or pop something from the stack, we must know *where* to push or pop it. This is where the **stack pointer** comes in. It is a special register that always points to the top of the stack, and is used to keep track of the stack.

Definition 6.25 (Push and Pop)

The **push** and **pop** instructions are used to push and pop data onto and off the stack, respectively.

```
push_ src          rsp = rsp - 8; Mem[rsp] = src
pop_ dest          dest = Mem[rsp]; rsp = rsp + 8
```

1. When we push the source, we fetch the value at the source and store it at the memory address pointed to by the stack pointer `%rsp`. Then, we decrement `%rsp` by 8.
2. When we pop from the stack, we fetch the value at the memory address pointed to by the stack pointer `%rsp` and store it in the destination. Then, we increment `%rsp` by 8.

Note that no matter what the size of the operand, we always subtract 8 from the stack pointer. This is because the stack grows downwards, and we want to make sure that the next element is pushed into the next available space.

Note that the register `%rsp` is the stack pointer, which points to the top of the stack. The stack is used to store data and return addresses, and is used to manage function calls.

Definition 6.26 (Push and Pop)

The **push** and **pop** instructions are used to push and pop data onto and off the stack, respectively.

```
push_ src          rsp = rsp - 8; Mem[rsp] = src
pop_ dest          dest = Mem[rsp]; rsp = rsp + 8
```

The `_` is a size operand, which determines how big the data is.

Definition 6.27 (Call and Ret)

The **call** instruction pushes the return address onto the stack and jumps to the function. The **ret** instruction pops the return address from the stack and jumps to it.

We also talked about how there is instruction code that is even below the stack that is stored. This is where all the machine code/assembly is stored, and we want to find out where we are currently at in this code. This is done with the program counter.

Definition 6.28 (Program Counter, Instruction Pointer)

The **program counter**, or **instruction pointer**, is a special register **rip** that points to the current instruction in the program. It is used to keep track of the next instruction to be executed.

Let's go through one long example to see in detail how this is calculated.

Example 6.15 (Evaluating a Function)

Say that we have the following C code.

```

1  int adder2(int a) {
2      return a + 2;
3  }
4
5  int main() {
6      int x = 40;
7      x = adder2(x);
8      printf("x is: %d\n", x);
9      return 0;
10 }

```

When we compile this program, we can view its full assembly code by calling `objdump -d a.out`. The output is quite long, so we will focus on the instruction for the `adder2` function.

```

1  000000000400526 <adder2>:
2  400526:    55                push  %rbp
3  400527:    48 89 e5          mov   %rsp,%rbp
4  40052a:    89 7d fc          mov   %edi,-0x4(%rbp)
5  40052d:    8b 45 fc          mov   -0x4(%rbp),%eax
6  400530:    83 c0 02          add  $0x2,%eax
7  400533:    5d                pop   %rbp
8  400534:    c3                retq

```

Figure 22: The output of `objdump` for the `adder2` function. The leftmost column represents the addresses (in hex) of where the actual instructions lie. The second column represents the machine code that is being executed. The third column represents the assembly code.

Note some things. Since `adder2` is taking in an integer input value, we want to load it into the lower 32 bits (4 bytes) of the `%rdi` register, which is the first parameter. So we use `%edi`. Likewise for the return value, we want to output an `int` so we use `%eax` rather than `%rax`. Let's go through some of the steps.

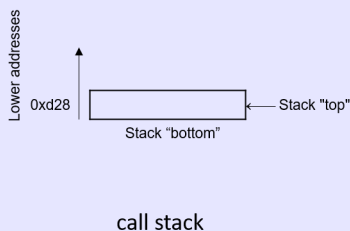
1. By the time we get into calling `adder2`, we can take a look at the relevant registers.

```

0x526 push %rbp
0x527 mov  %rsp, %rbp
0x52a mov  %edi, -0x4(%rbp)
0x52d mov  -0x4(%rbp), %eax
0x530 add  $0x2, %eax
0x533 pop  %rbp
0x534 retq

```

Registers	
<code>%eax</code>	0x123
<code>%edi</code>	0x28
<code>%rsp</code>	0xd28
<code>%rbp</code>	0xd40
<code>%rip</code>	0x526



- (a) First, the `%eax` is filled with garbage, which are leftovers from previous programs that haven't been overwritten yet.
- (b) Second, the `%edi=0x28` since we have set `x=40` in `main`, before calling `adder2`, so it lingers on.
- (c) `%rsp=0xd28` since that is where the top of the stack is.
- (d) `%rbp=0xd40`
- (e) `%rip=0x526` since that is where we are currently at in our instruction (we are about to do

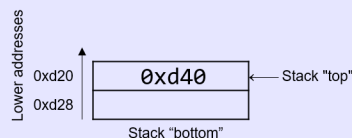
it, but haven't done it yet).

- When we execute the first line of code, we simply push the value at `%rbp` into the stack. The top of the stack gets decremented by 8 and the value at `%rbp` is stored there. This means that the top of the stack is at `%rsp=0xd20` and the next instruction will be at `%rip=0x527`.

```

➔ 0x526 push %rbp
   0x527 mov %rsp, %rbp
   0x52a mov %edi, -0x4(%rbp)
   0x52d mov -0x4(%rbp), %eax
   0x530 add $0x2, %eax
   0x533 pop %rbp
   0x534 retq
    
```

Registers	
<code>%eax</code>	0x123
<code>%edi</code>	0x28
<code>%rsp</code>	0xd20
<code>%rbp</code>	0xd40
<code>%rip</code>	0x527



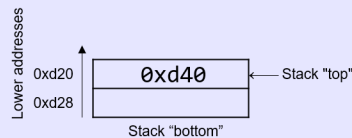
call stack

- The reason we have pushed `%rbp` onto the stack is that we want to save it before it gets overwritten by this next execution. We basically move the value of `%rsp` into `%rbp`, and the `%rip` advances to the next instruction. `%rip` moves to the next instruction.

```

➔ 0x526 push %rbp
   0x527 mov %rsp, %rbp
   0x52a mov %edi, -0x4(%rbp)
   0x52d mov -0x4(%rbp), %eax
   0x530 add $0x2, %eax
   0x533 pop %rbp
   0x534 retq
    
```

Registers	
<code>%eax</code>	0x123
<code>%edi</code>	0x28
<code>%rsp</code>	0xd20
<code>%rbp</code>	0xd20
<code>%rip</code>	0x52a



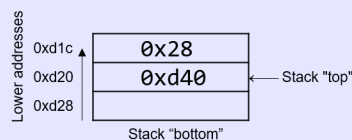
call stack

- Now we want to take our first argument `%edi` and store it in memory. Note that since this is 4 bytes, we can move this value into memory that is 4 bytes below the stack (`-0x4(%rbp)`). Note that the storing the value of `%edi` into memory doesn't affect the stack pointer `%rsp`. As far as the program is concerned, the top of this stack is still address `0xd20`.


```

0x526 push %rbp
0x527 mov %rsp, %rbp
→ 0x52a mov %edi, -0x4(%rbp)
0x52d mov -0x4(%rbp), %eax
0x530 add $0x2, %eax
0x533 pop %rbp
0x534 retq
    
```

Registers	
%eax	0x123
%edi	0x28
%rsp	0xd20
%rbp	0xd20
%rip	0x52d



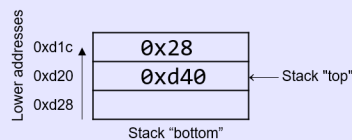
call stack

5. The next instruction simply goes into memory 4 bytes below the stack pointer, takes the value there, and stores it into %eax. This is the value of %edi that we just stored. This may seem redundant since we are making a round trip to memory and back to ultimately move the value of %edi into %eax, but compilers are not smart and just follow these instructions.

```

0x526 push %rbp
0x527 mov %rsp, %rbp
→ 0x52a mov %edi, -0x4(%rbp)
0x52d mov -0x4(%rbp), %eax
0x530 add $0x2, %eax
0x533 pop %rbp
0x534 retq
    
```

Registers	
%eax	0x28
%edi	0x28
%rsp	0xd20
%rbp	0xd20
%rip	0x530



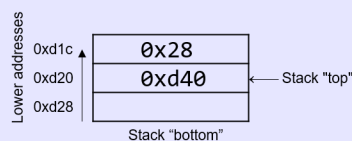
call stack

6. Finally, we add the value \$0x2 to %eax and store it back into %eax.

```

0x526 push %rbp
0x527 mov %rsp, %rbp
0x52a mov %edi, -0x4(%rbp)
0x52d mov -0x4(%rbp), %eax
→ 0x530 add $0x2, %eax
0x533 pop %rbp
0x534 retq
    
```

Registers	
%eax	0x2A
%edi	0x28
%rsp	0xd20
%rbp	0xd20
%rip	0x533



call stack

7. Finally, we pop the value at the top of the stack and store it into %rbp. Note that this is *not* the value 0x28. It is simply the value that is stored at %rsp=0xd20, which is (%rsp)=0xd40.

```

0x526 push %rbp
0x527 mov  %rsp, %rbp
0x52a mov  %edi, -0x4(%rbp)
0x52d mov  -0x4(%rbp), %eax
0x530 add  $0x2, %eax
→ 0x533 pop  %rbp
0x534 retq
    
```

Registers	
%eax	0x2A
%edi	0x28
%rsp	0xd28
%rbp	0xd40
%rip	0x534

8. Finally, we return the value with `retq`.

Note that the final values in the registers `%rsp` and `%rip` are `0xd28` and `0x534`, respectively, which are the same values as when the function started executing! This is normal and expected behavior with the call stack, which just stores temporary variable sand data of each function as it executes a program. Once a function completes executing, the stack returns to the state it was in prior to the function call. Therefore, it is common to see the following two instructions at the beginning of a function:

```

1 push %rbp
2 mov %rsp, %rbp
    
```

and the following two at the end of a function

```

1 pop %rbp
2 retq
    
```

Now arithmetic operations are quite simple.

Definition 6.29 (Add, Subtract, Multiply)

The **add** and **sub** instructions are used to add and subtract data from the destination.

```

add_ src, dest           dest = dest + src
sub_ src, dest           dest = dest - src
    
```

The **imul** instruction is used to multiply data between the source and destination and store it in the destination.

```

imul_ src, dest         dest = dest * src
    
```

Again the `_` is a size operand, which determines how big the data is.

Definition 6.30 (Increment, Decrement)

The **inc** and **dec** instructions are used to increment and decrement the value in the destination.

```

inc_ dest               dest = dest + 1
dec_ dest               dest = dest - 1
    
```

Definition 6.31 (Negative)

The **neg** instruction is used to negate the value in the destination.

$$\text{neg_ dest} \qquad \qquad \text{dest} = -\text{dest}$$

Example 6.16 (Basic Arithmetic Function)

The following represents the same program in C and in assembly. Let's go through each one:

1. In C, we first initialize **a = 4**, then **b = 8**, add them together to get **c**, and then return **c**.
2. In Assembly, we move the value 4 to the **%rax** register, then move the value 8 to the **%rbx** register, add the two values together to store it into **%rax**, and then return the value in the **%rax** register.

<pre> 1 int main() { 2 int a = 4, b = 8; 3 int c = a + b; 4 return c; 5 }</pre>	<pre> 1 main: 2 movq \$4, %rax 3 movq \$8, %rbx 4 addq %rbx, %rax 5 ret</pre>
--	--

It is slightly different in Assembly since rather than storing 4 in some intermediate register, we immediately store it in the return register. In a way it is more optimized, and this is what the compiler does for you so that as few registers are used.

A shorthand way to do this is with **lea**, which stands for load effective address.

Definition 6.32 (Load Effective Address)

The **lea** instruction is used to load the effective address of the source into the destination. For now, we will focus on the arithmetic operations that it can do

$$\begin{aligned} \text{lea_ (src1, src2), dest} & \qquad \qquad \text{dest} = \text{src1} + \text{src2} \\ \text{lea_ (src1, src2, scale), dest} & \qquad \text{dest} = \text{src1} + \text{src2} * \text{scale} \\ \text{lea_ const(src1, src2), dest} & \qquad \text{dest} = \text{src1} + \text{src2} + \text{const} \\ \text{lea_ const(src1, src2, scale), dest} & \qquad \text{dest} = \text{src1} + \text{src2} * \text{scale} + \text{const} \end{aligned}$$

This is useful for doing arithmetic operations on the address of a variable.

Definition 6.33 (Bitwise)

The **and**, **or**, **xor**, and **not** instructions are used to perform bitwise operations on the source and destination.

$$\begin{aligned} \text{and src, dest} & \qquad \qquad \text{dest} = \text{dest} \& \text{ src} \\ \text{or src, dest} & \qquad \qquad \text{dest} = \text{dest} | \text{ src} \\ \text{xor src, dest} & \qquad \qquad \text{dest} = \text{dest} \hat{\& } \text{ src} \\ \text{neg dest} & \qquad \qquad \text{dest} = -\text{dest} \\ \text{not dest} & \qquad \qquad \text{dest} = \sim \text{dest} \end{aligned}$$

Definition 6.34 (Arithmetic and Logical Bit Shift)

The `sal` arithmetic instruction is used to shift the bits of the destination to the left by the number of bits specified in the source. The `shr` instruction is used to shift the bits of the destination to the right by the number of bits specified in the source.

```
sal src, dest           dest = dest << src
shr src, dest           dest = dest >> src
```

The `sar` instruction is used to shift the bits of the destination to the right by the number of bits specified in the source, and fill the leftmost bits with the sign bit. The `shl` instruction is used to shift the bits of the destination to the left by the number of bits specified in the source, and fill the rightmost bits with zeros.

```
sar src, dest           dest = dest >> src
shl src, dest           dest = dest << src
```

Example 6.17 (Harder Arithmetic Example)

The following two codes are equivalent.

<pre>1 long arith(long x, long y, long z) { 2 long t1 = x + y; 3 long t2 = z + t1; 4 long t3 = x + 4; 5 long t4 = y * 48; 6 long t5 = t3 + t4; 7 long rval = t2 * t5; 8 return rval; 9 } 10 . 11 . 12 . 13 . 14 .</pre>	<pre>1 arith: 2 # rax/t1 = x + y 3 leaq (%rdi, %rsi), %rax 4 # rax/t2 = z + t1 5 addq %rdx, %rax 6 #rdx = 3 * y 7 leaq (%rsi, %rsi, 2), %rdx 8 #rdx/t4 = (3*y) * 16 9 salq \$4, %rdx 10 #rcx/t5 = x + t4 + 4 11 leaq 4(%rdi, %rdi), %rcx 12 # rax/rval = t5 * t2 13 imulq %rcx, %rax 14 ret</pre>
---	--

The final thing in our list is condition codes.

Sometimes, we want to move (really copy) some value to another register if some condition is met. This is where we use conditional moves. These conditions are met by the flags register, which is a special register that stores the status of the last operation. It is the value of these flags that determine whether all future conditional statements are met in assembly.

Definition 6.35 (Condition Code Flags)

The flags register in the x86 CPU keeps 4 *condition code* flag bits internally. Think of these as status flags that are *implicitly* set by the most recent arithmetic operation (think of it as side effects). Note that condition codes are NOT set by `lea` or `mov` instructions!

1. **Zero Flag:** if the last operation resulted in a zero value.
2. **Sign Flag:** if the last operation resulted in a negative value (i.e. the most significant bit is 1).
3. **Overflow Flag:** if the last operation resulted in a signed overflow.
4. **Carry Flag:** if the last operation resulted in a carry out of the most significant bit, i.e. an unsigned overflow.

Every operation may or may not changes these flags to test for zero or nonzero, positive or negative,

or overflow conditions, and combinations of these flags express the full range of conditions and cases, e.g. for signed and unsigned values.

Example 6.18 (Zero Flag)

If the code below was just run, then ZF would be set to 1.

```
1 movq $2, %rax
2 subq $2, %rax
```

Example 6.19 (Sign Flag)

If the code below was just run, then SF would be set to 1.

```
1 movq $2, %rax
2 subq $4, %rax
```

Example 6.20 (Overflow Flag)

If either code below was just run, then OF would be set to 1.

```
1 movq $0x7fffffffffffffff, %rax
2 addq $1, %rax
```

```
1 movq 0x8000000000000000, %rax
2 addq 0xfffffffffffffff, %rax
```

This is because in the left in signed arithmetic, we have a positive + positive = negative (result is 0x8000000000000000), which is a signed overflow. Furthermore, in the right we have negative + negative = positive (result is 0x7fffffffffffffff).

Example 6.21 (Carry Flag)

If the code below was just run, then CF would be set to 1.

```
1 movq $0xfffffffffffffff, %rax
2 addq $1, %rax
```

This is because the result is 0x0, which is a carry out of the most significant bit and an unsigned overflow.

It would be tedious to always set these flags manually, so there are two methods that can be used to *explicitly* set these flags.

Definition 6.36 (Compare)

The **cmp** instruction is used to perform a subtraction between the source and destination, and set the flags accordingly, but it does not store the result.

$$\text{cmp_src, dest} \qquad \text{dest - src}$$

The following flags are set if the conditions are met:

1. **ZF** = 1 if `dest == src`
2. **SF** = 1 if `dest < src` (MSB is 1)
3. **OF** = 1 if signed overflow

- CF = 1** if unsigned overflow

Definition 6.37 (Test)

The **test** instruction is used to perform a bitwise AND operation between the source and destination, and set the flags accordingly.

```
test_ src, dest          dest & src
```

The following flags are set if the conditions are met. Note that you can't have carry out (CF) or overflow (OF) if these flags are set.

- ZF = 1** if `dest & src == 0`
- SF = 1** if `dest & src < 0` (MSB is 1)

Example 6.22 (Compare)

Assuming that `%a1 = 0x80` and `%b1 = 0x81`, which flags are set when we execute `cmpb %a1, %b1`? Well we must first compute

$$\%b1 - \%a1 = 0x81 - 0x80 = 0x81 + \sim 0x80 + 1 = 0x81 + 0x7F + 1 = 0x101 = 0x01 \quad (10)$$

- CF=1 since the result is greater than 0xFF (i.e. larger than byte)
- ZF=0 since the result is not 0
- SF=0 since the MSB is 0, i.e. there is unsigned overflow
- OF=0 since there is no signed overflow

For conditional moves and jumps later shown, it basically uses these explicit sets and always compares them to 0. We will see what this means later.

Finally, we can actually set a byte in a register to 1 or 0 based on the value of a flag.

Definition 6.38 (Set)

We can then talk about conditional moves and jumps.

Definition 6.39 (Equality with 0)

The **test** instruction is used to perform a bitwise AND operation between the source and destination, and set the flags accordingly.

```
test_ src, dest          dest & src
```

The **sete** instruction is used to set the destination to 1 if the zero flag is set, and 0 otherwise.

```
sete_ dest          dest = (ZF == 1) ? 1 : 0
```

The **cmovne** instruction is used to move the source to the destination if the zero flag is not set.

```
cmovne_ src, dest          dest = (ZF == 0) ? src : dest
```

Definition 6.40 (Jump)

There are several jump instructions, but essentially they are used to jump to another part of the code. We can use the following mnemonic to jump to a label.

Letter	Word
j	jump
n	not
e	equal
s	signed
g	greater (signed interpretation)
l	less (signed interpretation)
a	above (unsigned interpretation)
b	below (unsigned interpretation)

Table 2: Letter to Word Mapping

Figure 23: Mnemonic for Jump Instructions

For completeness, we include all the jump instructions.

Signed Comparison	Unsigned Comparison	Description
je (jz)		jump if equal (==) or jump if zero
jne (jnz)		jump if not equal (!=)
js		jump if negative
jns		jump if non-negative
jg (jnl)	ja (jnbe)	jump if greater (>)
jge (jnl)	jae (jnb)	jump if greater than or equal (>=)
jl (jnge)	jb (jnae)	jump if less (<)
jle (jng)	jbe (jna)	jump if less than or equal (<=)

Table 3: Comparison Instructions in Assembly

Figure 24: All jump instructions

Definition 6.41 (int)

The `int` instruction is used to generate a software interrupt. It is often used to invoke a system call.

Definition 6.42 (ret)

The `ret` instruction is used to return from a function. It returns the value in the `%rax` register.

Now we can have a basic idea of how if statements can be used as a sequence of conditionals and jump operators. Let's first look at the `goto` version of C.

Definition 6.43 (Goto Syntax)

The goto version processes instructions sequentially as long as there is no jump. This is useful because compilers translating code into assembly designate a jump when a condition is true. Contrast this behavior with the structure of an if statement, where a "jump" (to the else) occurs when conditions are not true. The goto form captures this difference in logic.

<pre> 1 int getSmallest(int x, int y) { 2 int smallest; 3 if (x > y) { //if (conditional) 4 smallest = y; //then statement 5 } 6 else { 7 smallest = x; //else statement 8 } 9 return smallest; 10 } 11 . 12 . 13 . 14 . 15 . </pre>	<pre> 1 int getSmallest(int x, int y) { 2 int smallest; 3 4 if (x <= y) { //if (!conditional) 5 goto else_statement; 6 } 7 smallest = y; //then statement 8 goto done; 9 10 else_statement: 11 smallest = x; //else statement 12 13 done: 14 return smallest; 15 } </pre>
--	--

Figure 25: C vs GoTo code of the same function. While GoTo code allows us to view C more like assembly, it is generally not readable and is not considered best practice.

Now let's see how if statements are implemented by taking a look at this function straight up in assembly.

<pre> 1 int getSmallest(int x, int y) { 2 int smallest; 3 if (x > y) { //if (conditional) 4 smallest = y; //then statement 5 } 6 else { 7 smallest = x; //else statement 8 } 9 return smallest; 10 } 11 . </pre>	<pre> 1 Dump of assembler code for function getSmallest: 2 0x40059a <+4>: mov %edi,-0x14(%rbp) 3 0x40059d <+7>: mov %esi,-0x18(%rbp) 4 0x4005a0 <+10>: mov -0x14(%rbp),%eax 5 0x4005a3 <+13>: cmp -0x18(%rbp),%eax 6 0x4005a6 <+16>: jle 0x4005b0 <getSmallest+26> 7 0x4005a8 <+18>: mov -0x18(%rbp),%eax 8 0x4005ae <+24>: jmp 0x4005b9 <getSmallest+35> 9 0x4005b0 <+26>: mov -0x14(%rbp),%eax 10 0x4005b9 <+35>: pop %rbp 11 0x4005ba <+36>: retq </pre>
--	---

Figure 26: Assembly code of a simple if statement

Again, note that since we are working with int types, the respective parameter registers are %edi and %esi, the respective lower 32-bits of the registers %rdi and %rsi. Let's walk through this again.

1. The first mov instruction copies the value located in register %edi (the first parameter, x) and places it at memory location %rbp-0x14 on the call stack. The instruction pointer (%rip) is set to the address of the next instruction, or 0x40059d.
2. The second mov instruction copies the value located in register %esi (the second parameter, y) and places it at memory location %rbp-0x18 on the call stack. The instruction pointer (%rip) updates to point to the address of the next instruction, or 0x4005a0.

3. The third mov instruction copies x to register %eax. Register %rip updates to point to the address of the next instruction in sequence.
4. The cmp instruction compares the value at location %rbp-0x18 (the second parameter, y) to x and sets appropriate condition code flag registers. Register %rip advances to the address of the next instruction, or 0x4005a6.
5. The jle instruction at address 0x4005a6 indicates that if x is less than or equal to y, the next instruction that should execute should be at location <getSmallest+26> and that %rip should be set to address 0x4005b0. Otherwise, %rip is set to the next instruction in sequence, or 0x4005a8.

With the cmov instruction, this can be a lot shorter. With the gcc compiler with level 1 optimizations turned on, we can see that a lot of redundancies are turned off.

```

1 <getSmallest>:
2 0x400546 <+0>: cmp    %esi,%edi    #compare x and y
3 0x400548 <+2>: mov    %esi,%eax    #copy y to %eax
4 0x40054a <+4>: cmovle %edi,%eax    #if (x<=y) copy x to %eax
5 0x40054d <+7>: retq   #return %eax
    
```

Figure 27: Compiled with gcc -O1 -o getSmallest getSmallest.c

Like if statements, loops in assembly can be implementing using jump functions that revisit some instruction address based on the result on an evaluated condition. Let's take a look at a basic loop function.

<pre> 1 int sumUp(int n) { 2 int total = 0; 3 int i = 1; 4 5 while (i <= n) { 6 total += i; 7 i++; 8 } 9 return total; 10 } 11 . 12 . 13 . 14 . 15 . 16 . </pre>	<pre> 1 Dump of assembler code for function sumUp: 2 0x400526 <+0>: push %rbp 3 0x400527 <+1>: mov %rsp,%rbp 4 0x40052a <+4>: mov %edi,-0x14(%rbp) 5 0x40052d <+7>: mov \$0x0,-0x8(%rbp) 6 0x400534 <+14>: mov \$0x1,-0x4(%rbp) 7 0x40053b <+21>: jmp 0x400547 <sumUp+33> 8 0x40053d <+23>: mov -0x4(%rbp),%eax 9 0x400540 <+26>: add %eax,-0x8(%rbp) 10 0x400543 <+29>: add \$0x1,-0x4(%rbp) 11 0x400547 <+33>: mov -0x4(%rbp),%eax 12 0x40054a <+36>: cmp -0x14(%rbp),%eax 13 0x40054d <+39>: jle 0x40053d <sumUp+23> 14 0x40054f <+41>: mov -0x8(%rbp),%eax 15 0x400552 <+44>: pop %rbp 16 0x400553 <+45>: retq </pre>
---	---

Figure 28: Simple loop function in C and assembly.

Finally, we want to let the reader know the convention of callee and caller saved registers. The compiler tries to pick these registers, and by convention in x86, we have the following.

<code>%rax</code>	Return value - Caller saved	<code>%r8</code>	Argument #5 - Caller saved
<code>%rbx</code>	Callee saved	<code>%r9</code>	Argument #6 - Caller saved
<code>%rcx</code>	Argument #4 - Caller saved	<code>%r10</code>	Caller saved
<code>%rdx</code>	Argument #3 - Caller saved	<code>%r11</code>	Caller Saved
<code>%rsi</code>	Argument #2 - Caller saved	<code>%r12</code>	Callee saved
<code>%rdi</code>	Argument #1 - Caller saved	<code>%r13</code>	Callee saved
<code>%rsp</code>	Stack pointer	<code>%r14</code>	Callee saved
<code>%rbp</code>	Callee saved	<code>%r15</code>	Callee saved

Figure 29: Caller save and callee save registers.

So far, we’ve traced through simple functions in assembly. In this section, we discuss the interaction between multiple functions in assembly in the context of a larger program. We also introduce some new instructions involved with function management.

Definition 6.44 (Leave)

The **leave** instruction is used to deallocate the current stack frame. For example, the `leaveq` instruction is a shorthand that the compiler uses to restore the stack and frame pointers as it prepares to leave a function. When the callee function finishes execution, `leaveq` ensures that the frame pointer is restored to its previous value. It is equivalent to the following two instructions:

```
leaveq                                movq %rbp, %rsp
                                       popq %rbp
```

Definition 6.45 (Call and Return)

The **call** instruction is used to call a function and the **ret** to return from a function. The `callq` and `retq` instructions play a prominent role in the process where one function calls another. Both instructions modify the instruction pointer (register `%rip`).

1. When the caller function executes the `callq` instruction, the current value of `%rip` is saved on the stack to represent the return address, or the program address at which the caller resumes executing once the callee function finishes. The `callq` instruction also replaces the value of `%rip` with the address of the callee function.

```
callq addr <fname>                   push %rip
                                       mov addr, %rip
```

2. The `retq` instruction restores the value of `%rip` to the value saved on the stack, ensuring that the program resumes execution at the program address specified in the caller function. Any value returned by the callee is stored in `%rax` or one of its component registers (e.g., `%eax`). The `retq` instruction is usually the last instruction that executes in any function.

```
retq                                   pop %rip
```

Let's work through an example to solidify our knowledge.

Example 6.23 (Calling Functions in Assembly)

Let's take the following code and trace through main.

1	<code>#include <stdio.h></code>	1	000000000400526 <assign>:	
2		2	400526: 55	push %rbp
3	<code>int assign(void) {</code>	3	400527: 48 89 e5	mov %rsp,%rbp
4	<code>int y = 40;</code>	4	40052a: c7 45 fc 28 00 00 00	movl \$0x28,-0x4(%rbp)
5	<code>return y;</code>	5	400531: 8b 45 fc	mov -0x4(%rbp),%eax
6	<code>}</code>	6	400534: 5d	pop %rbp
7		7	400535: c3	retq
8	<code>int adder(void) {</code>	8		
9	<code>int a;</code>	9	000000000400536 <adder>:	
10	<code>return a + 2;</code>	10	400536: 55	push %rbp
11	<code>}</code>	11	400537: 48 89 e5	mov %rsp,%rbp
12		12	40053a: 8b 45 fc	mov -0x4(%rbp),%eax
13	<code>int main(void) {</code>	13	40053d: 83 c0 02	add \$0x2,%eax
14	<code>int x;</code>	14	400540: 5d	pop %rbp
15	<code>assign();</code>	15	400541: c3	retq
16	<code>x = adder();</code>	16		
17	<code>printf("x is:</code>	17	000000000400542 <main>:	
18	<code>%d\n", x);</code>	18	400542: 55	push %rbp
19	<code>return 0;</code>	19	400543: 48 89 e5	mov %rsp,%rbp
20	<code>}</code>	20	400546: 48 83 ec 10	sub \$0x10,%rsp
21	<code>.</code>	21	40054a: e8 e3 ff ff ff	callq 400526 <assign>
22	<code>.</code>	22	40054f: e8 d2 ff ff ff	callq 400536 <adder>
23	<code>.</code>	23	400554: 89 45 fc	mov %eax,-0x4(%rbp)
24	<code>.</code>	24	400557: 8b 45 fc	mov -0x4(%rbp),%eax
25	<code>.</code>	25	40055a: 89 c6	mov %eax,%esi
26	<code>.</code>	26	40055c: bf 04 06 40 00	mov \$0x400604,%edi
27	<code>.</code>	27	400561: b8 00 00 00 00	mov \$0x0,%eax
28	<code>.</code>	28	400566: e8 95 fe ff ff	callq 400400
29	<code>.</code>	29	<printf@plt>	
30	<code>.</code>	30	40056b: b8 00 00 00 00	mov \$0x0,%eax
31	<code>.</code>	31	400570: c9	leaveq
			400571: c3	retq

Figure 30: C code and its assembly equivalent. Main function calls two other functions.

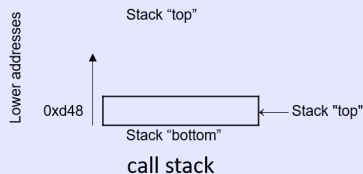
Let's trace through what happens here in detail. This will be long.

1. `%rbp` is the base pointer that is initialized to something. Before we even begin main, say that we have the following initializations, where `%eax`, `%edi` is garbage. `%rsp` denotes where on the stack we are right before calling to main, `%rbp` is the base pointer to the current program, and `%rip` should be the address of the first instruction in main. Again since we work with integers we use the lower 32-bits of the registers. `%rip` now points to the next instruction.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
    
```

Registers	
%eax	650
%edi	1
%rsp	0xd48
%rbp	0x830
%rip	0x542



```

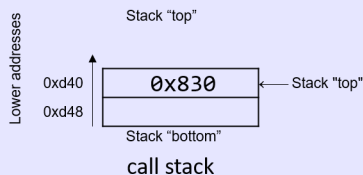
Terminal:
$ ./prog
    
```

2. Now we start the main function. By calling main, the base pointer `%rbp` of the stack outside of the main frame will be overwritten by the base of the main stack frame, so we must save it for when main is done. Therefore, we push it onto the stack where `%rsp` is pointing. `%rip` now points to the next instruction.

```

0x542 <main>:
→ 0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
    
```

Registers	
%eax	650
%edi	1
%rsp	0xd40
%rbp	0x830
%rip	0x543



```

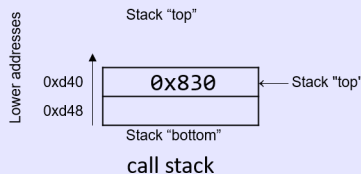
Terminal:
$ ./prog
    
```

3. Then we actually change the location of the base pointer to the top of the stack, which now includes the first instruction in main.

```

0x542 <main>:
0x542 push  %rbp
➔ 0x543 mov  %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov  %eax, -0x4(%rbp)
0x557 mov  -0x4(%rbp), %eax
0x55a mov  %eax, %esi
    
```

Registers	
%eax	650
%edi	1
%rsp	0xd40
%rbp	0xd40
%rip	0x546



```

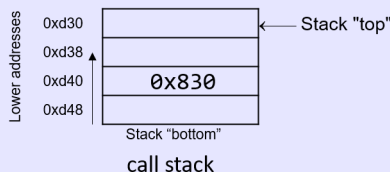
Terminal:
$ ./prog
    
```

4. Now we manually change the stack pointer and have it grow by two bytes (0x10). Therefore, %rsp is decremented by 0x10 and %rip points to the next instruction at 0x54a.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov  %rsp, %rbp
➔ 0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov  %eax, -0x4(%rbp)
0x557 mov  -0x4(%rbp), %eax
0x55a mov  %eax, %esi
    
```

Registers	
%eax	650
%edi	1
%rsp	0xd30
%rbp	0xd40
%rip	0x54a



```

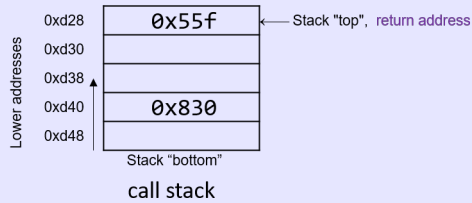
Terminal:
$ ./prog
    
```

5. Now the next instruction pointed at by %rip is the callq instruction, which tells to go to the address of the assign function. We by default first update %rip to point to the next instruction at 0x55f. However, this should not be the actual next instruction that we execute since we are calling another function. Rather, we want to update %rip to address 0x526 where assign is located at, but after completion we also want to know that we want to execute the instruction after it at address 0x55f. Therefore, we should save address 0x55f onto the stack and then update %rip to point to 0x526. This is what we refer to as a **return address**.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
➔ 0x54a callq 0x526 <assign>
➡ 0x55f callq 0x536 <adder>
0x554 mov  %eax, -0x4(%rbp)
0x557 mov  -0x4(%rbp), %eax
0x55a mov  %eax, %esi
    
```

Registers	
%eax	0x0
%edi	1
%rsp	0xd28
%rbp	0xd40
%rip	0x526



```

Terminal:
$ ./prog
    
```

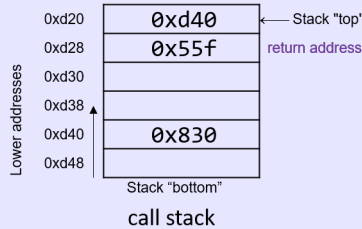
Equivalent to:
 push %rip
 mov 0x526, %rip

6. %rip is incremented to the next address. We step into the `assign` function, which is now a new stack frame, so the first thing we do is save the base pointer of the main stack frame onto the stack since we must immediately update it with the base pointer of the `assign` stack frame, which is where %rsp is pointing to.

```

0x526 <assign>:
➔ 0x526 push  %rbp
0x527 mov   %rsp, %rbp
0x52a mov   $0x28, -0x4(%rbp)
0x531 mov  -0x4(%rbp), %eax
0x534 pop  %rbp
0x535 retq
0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov  %eax, -0x4(%rbp)
0x557 mov  -0x4(%rbp), %eax
0x55a mov  %eax, %esi
    
```

Registers	
%eax	0x0
%edi	1
%rsp	0xd20
%rbp	0xd40
%rip	0x527



```

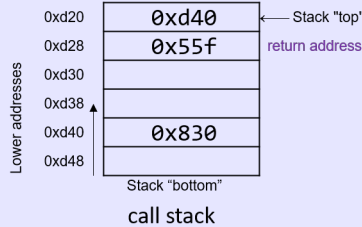
Terminal:
$ ./prog
    
```

7. %rip is incremented to the next address. We then update the base pointer to the top of the stack.

```

0x526 <assign>:
0x526 push  %rbp
➔ 0x527 mov  %rsp, %rbp
0x52a mov  $0x28, -0x4(%rbp)
0x531 mov  -0x4(%rbp), %eax
0x534 pop  %rbp
0x535 retq
0x542 <main>:
0x542 push  %rbp
0x543 mov  %rsp, %rbp
0x546 sub  $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov  %eax, -0x4(%rbp)
0x557 mov  -0x4(%rbp), %eax
0x55a mov  %eax, %esi
    
```

Registers	
%eax	0x0
%edi	1
%rsp	0xd20
%rbp	0xd20
%rip	0x52a



```

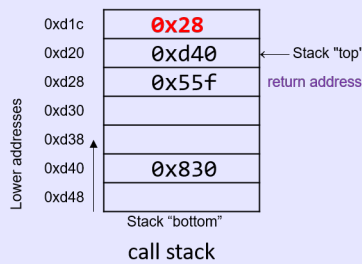
Terminal:
$ ./prog
    
```

8. Now we want to move the number 0x28 (40) into the memory location -0x4(%rbp) of the stack, which is 4 bytes above the frame pointer, which is also the stack pointer. It is common that the frame pointer is used to reference locations on the stack. Note that this does not update the stack pointer.

```

0x526 <assign>:
0x526 push  %rbp
➔ 0x527 mov  %rsp, %rbp
0x52a mov  $0x28, -0x4(%rbp)
0x531 mov  -0x4(%rbp), %eax
0x534 pop  %rbp
0x535 retq
0x542 <main>:
0x542 push  %rbp
0x543 mov  %rsp, %rbp
0x546 sub  $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov  %eax, -0x4(%rbp)
0x557 mov  -0x4(%rbp), %eax
0x55a mov  %eax, %esi
    
```

Registers	
%eax	0x0
%edi	1
%rsp	0xd20
%rbp	0xd20
%rip	0x531



```

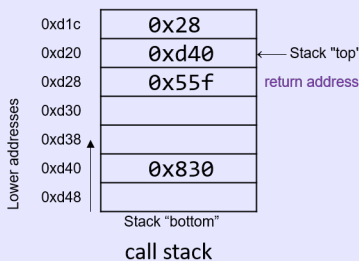
Terminal:
$ ./prog
    
```

9. Now we take the same address where we stored 0x28 to and move it into %eax, effectively loading 40 onto the return value.

```

0x526 <assign>:
0x526 push %rbp
0x527 mov %rsp, %rbp
0x52a mov $0x28, -0x4(%rbp)
➔ 0x531 mov -0x4(%rbp), %eax
0x534 pop %rbp
0x535 retq
0x542 <main>:
0x542 push %rbp
0x543 mov %rsp, %rbp
0x546 sub $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov %eax, -0x4(%rbp)
0x557 mov -0x4(%rbp), %eax
0x55a mov %eax, %esi
    
```

Registers	
%eax	0x28
%edi	1
%rsp	0xd20
%rbp	0xd20
%rip	0x534



```

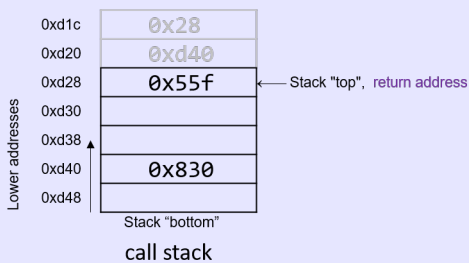
Terminal:
$ ./prog
    
```

10. We see that we will return this value soon, but before we do, we want to make sure that when the assign stack frame gets deleted (not really, but overwritten), we want to restore the base pointer of the main stack frame. We have already saved this before at `%rsp`, which hasn't changed since we only worked with displacements from the base pointer. We retrieve the main stack pointer data and load it back into `%rbp`. Note that this increments `%rsp` by 8 bytes, shrinking the stack, and we are technically out of the assign stack frame.

```

0x526 <assign>:
0x526 push %rbp
0x527 mov %rsp, %rbp
0x52a mov $0x28, -0x4(%rbp)
0x531 mov -0x4(%rbp), %eax
➔ 0x534 pop %rbp
0x535 retq
0x542 <main>:
0x542 push %rbp
0x543 mov %rsp, %rbp
0x546 sub $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov %eax, -0x4(%rbp)
0x557 mov -0x4(%rbp), %eax
0x55a mov %eax, %esi
    
```

Registers	
%eax	0x28
%edi	1
%rsp	0xd28
%rbp	0xd40
%rip	0x535



```

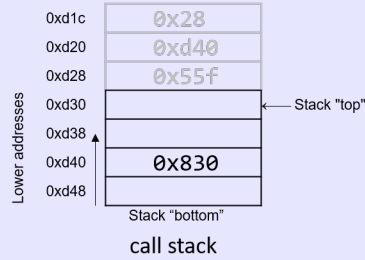
Terminal:
$ ./prog
    
```

11. Note that at this point, since `%rbp` was popped off, the next value that is at the top of the stack

is the address `%rip` that we store earlier, which points to the next execution in `main`. When `retq` executes, this value at the top of the stack is popped into `%rip`, allowing `main` to continue executing within the `main` stack frame. Note that the return value is stored in `%eax`.

```

0x526 <assign>:
0x526 push  %rbp
0x527 mov   %rsp, %rbp
0x52a mov   $0x28, -0x4(%rbp)
0x531 mov   -0x4(%rbp), %eax
0x534 pop  %rbp
➔ 0x535 retq
0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
    
```



```

Terminal:
$ ./prog
    
```

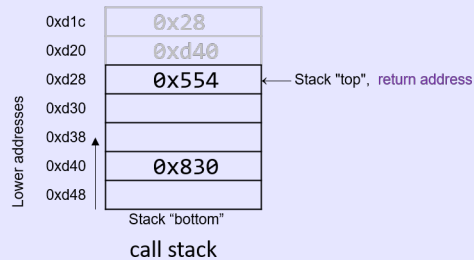
Registers	
<code>%eax</code>	0x28
<code>%edi</code>	1
<code>%rsp</code>	0xd30
<code>%rbp</code>	0xd40
<code>%rip</code>	0x55f

Equivalent to:
`pop %rip`

- Now we execute the next instruction in `%rip` which is a call to the `adder` function. `%rip` is automatically updated to the next address at 0x554, but since this is a `callq` instruction, we first want to store this `%rip` into the stack so we can come back to it, and then update `%rip` to the first instruction in `adder`, which is address 0x536.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
➔ 0x55f callq 0x536 <adder>
➡ 0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
    
```



```

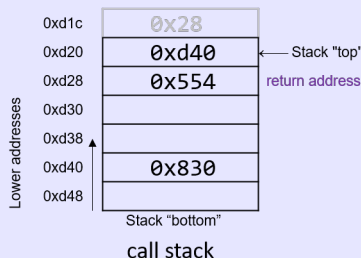
Terminal:
$ ./prog
    
```

Registers	
<code>%eax</code>	0x0
<code>%edi</code>	1
<code>%rsp</code>	0xd28
<code>%rbp</code>	0xd40
<code>%rip</code>	0x536

- Since we are in the `adder` function, this creates a new stack frame and we must update `%rbp`. Again, we don't want to overwrite the base pointer of `main`, so we save it onto the stack by pushing `%rbp`.

```

0x536 <adder>:
→ 0x536 push %rbp
0x537 mov %rsp, %rbp
0x53a mov $-0x4(%rbp), %eax
0x53d add $0x2, %eax
0x540 pop %rbp
0x541 retq
0x542 <main>:
0x542 push %rbp
0x543 mov %rsp, %rbp
0x546 sub $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov %eax, -0x4(%rbp)
0x557 mov -0x4(%rbp), %eax
0x55a mov %eax, %esi
    
```



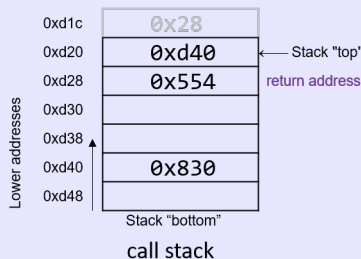
Terminal:
\$./prog

Registers	
%eax	0x0
%edi	1
%rsp	0xd20
%rbp	0xd40
%rip	0x537

14. Then we update %rbp to the current stack pointer.

```

0x536 <adder>:
0x536 push %rbp
→ 0x537 mov %rsp, %rbp
0x53a mov $-0x4(%rbp), %eax
0x53d add $0x2, %eax
0x540 pop %rbp
0x541 retq
0x542 <main>:
0x542 push %rbp
0x543 mov %rsp, %rbp
0x546 sub $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov %eax, -0x4(%rbp)
0x557 mov -0x4(%rbp), %eax
0x55a mov %eax, %esi
    
```



Terminal:
\$./prog

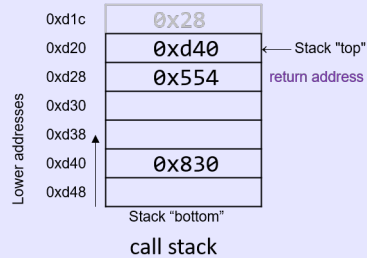
Registers	
%eax	0x0
%edi	1
%rsp	0xd20
%rbp	0xd20
%rip	0x53a

15. This part is a bit tricky. Note that the value of 0x28 still lives at 0xd1c, which is conveniently at address -0x4(%rbp). Therefore, when we call int a; in that corresponding line in adder, we can actually add 2 to it, though it seems like there was no value assigned to it. This is just a trick though. So, we can take these remnant value and store it into %eax.

```

0x536 <adder>:
0x536 push %rbp
0x537 mov %rsp, %rbp
➔ 0x53a mov $-0x4(%rbp), %eax
0x53d add $0x2, %eax
0x540 pop %rbp
0x541 retq
0x542 <main>:
0x542 push %rbp
0x543 mov %rsp, %rbp
0x546 sub $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x55a mov %eax, -0x4(%rbp)
0x557 mov -0x4(%rbp), %eax
0x55a mov %eax, %esi
    
```

Registers	
%eax	0x28
%edi	1
%rsp	0xd20
%rbp	0xd20
%rip	0x53d



Terminal:
\$./prog

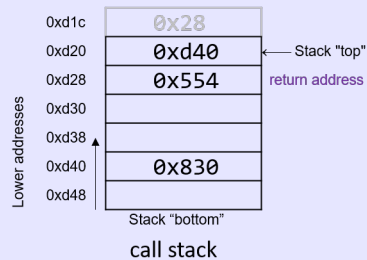
Using an old value on the stack!

16. We then add 2 to it.

```

0x536 <adder>:
0x536 push %rbp
0x537 mov %rsp, %rbp
0x53a mov $-0x4(%rbp), %eax
➔ 0x53d add $0x2, %eax
0x540 pop %rbp
0x541 retq
0x542 <main>:
0x542 push %rbp
0x543 mov %rsp, %rbp
0x546 sub $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x55a mov %eax, -0x4(%rbp)
0x557 mov -0x4(%rbp), %eax
0x55a mov %eax, %esi
    
```

Registers	
%eax	0x2A
%edi	1
%rsp	0xd20
%rbp	0xd20
%rip	0x540



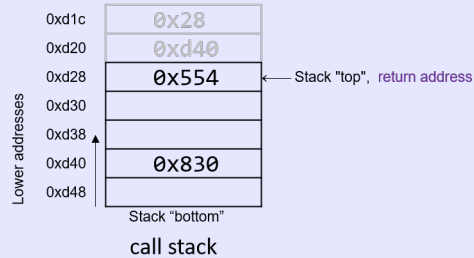
Terminal:
\$./prog

17. Now we are almost done, so we pop the base pointer of the main stack frame, at 0xd40, back into %rbp.

```

0x536 <adder>:
0x536 push %rbp
0x537 mov %rsp, %rbp
0x53a mov $-0x4(%rbp), %eax
0x53d add $0x2, %eax
➔ 0x540 pop %rbp
0x541 retq
0x542 <main>:
0x542 push %rbp
0x543 mov %rsp, %rbp
0x546 sub $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov %eax, -0x4(%rbp)
0x557 mov -0x4(%rbp), %eax
0x55a mov %eax, %esi
    
```

Registers	
%eax	0x2A
%edi	1
%rsp	0xd28
%rbp	0xd40
%rip	0x541



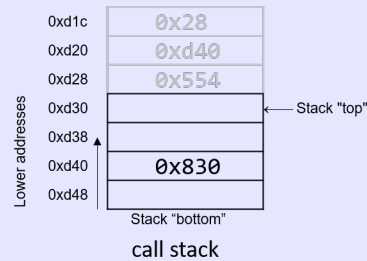
Terminal:
\$./prog

18. We now return the value in %eax and pop the base pointer of the adder stack frame, which simply updates the instruction pointer %rip back to the next instruction in main. This is equivalent to pop %rip, which is equivalent to moving the stack pointer %rsp into %rip and then shrinking the stack by 8 bytes subq \$8, %rsp.

```

0x536 <adder>:
0x536 push %rbp
0x537 mov %rsp, %rbp
0x53a mov $-0x4(%rbp), %eax
0x53d add $0x2, %eax
0x540 pop %rbp
➔ 0x541 retq
0x542 <main>:
0x542 push %rbp
0x543 mov %rsp, %rbp
0x546 sub $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov %eax, -0x4(%rbp)
0x557 mov -0x4(%rbp), %eax
0x55a mov %eax, %esi
    
```

Registers	
%eax	0x2A
%edi	1
%rsp	0xd30
%rbp	0xd40
%rip	0x554



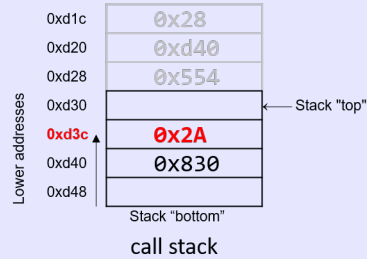
Terminal:
\$./prog

19. Now it is relatively straightforward since we do the rest in main (except for the print statement). The current value in %eax represents the return value of adder. We want to put this in the variable x, which we have already allocated some memory for right above the base pointer in the main stack frame. We move it there. Note that right after, it places this right back into

%eax.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
→ 0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
0x55c mov   $0x400604, %edi
0x561 mov   $0x0, %eax
0x566 callq <printf@plt>
0x56b mov   $0x0, %eax
0x570 leaveq
0x571 retq
    
```



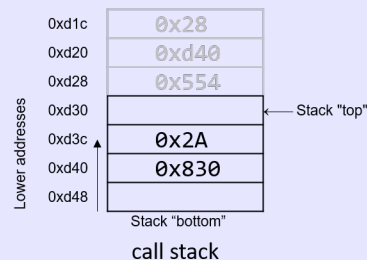
Terminal:
\$./prog

Registers	
%eax	0x2A
%edi	1
%rsp	0xd30
%rbp	0xd40
%rip	0x557

20. the mov instruction at address 0x55a copies the value in %eax (or 0x2A) to register %esi, which is the 32-bit component register associated with %rsi and typically stores the second parameter to a function. We can see why since this will be put into a print statement, which is a function, and x = %esi is the second argument of printf.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
→ 0x55a mov   %eax, %esi
0x55c mov   $0x400604, %edi
0x561 mov   $0x0, %eax
0x566 callq <printf@plt>
0x56b mov   $0x0, %eax
0x570 leaveq
0x571 retq
    
```



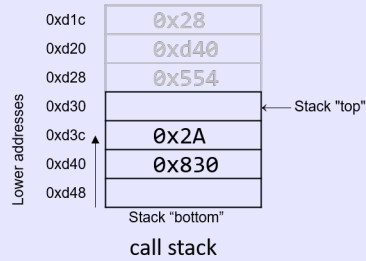
Terminal:
\$./prog

Registers	
%eax	0x2A
%edi	1
%rsp	0xd30
%rbp	0xd40
%rip	0x55c
%esi	0x2A

21. Now we want to retrieve the first argument of the print function. The address at \$0x400604 is some address in the code segment memory that holds the string "x is %d\n".

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
→ 0x55c mov   $0x400604, %edi
0x561 mov   $0x0, %eax
0x566 callq <printf@plt>
0x56b mov   $0x0, %eax
0x570 leaveq
0x571 retq
    
```



Registers	
%eax	0x2A
%edi	0x400604
%rsp	0xd30
%rbp	0xd40
%rip	0x561
%esi	0x2A

Terminal:

```
$ ./prog
```

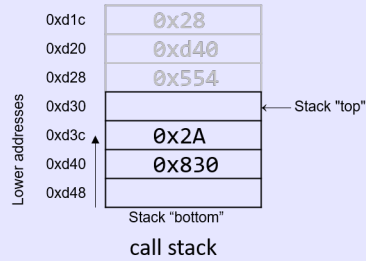
Memory

0x400604	"x is %d\n"
----------	-------------

22. Then we move 0 into the %eax register to clear it.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
→ 0x55c mov   $0x400604, %edi
0x561 mov   $0x0, %eax
0x566 callq <printf@plt>
0x56b mov   $0x0, %eax
0x570 leaveq
0x571 retq
    
```



Registers	
%eax	0x0
%edi	0x400604
%rsp	0xd30
%rbp	0xd40
%rip	0x566
%esi	0x2A

Terminal:

```
$ ./prog
```

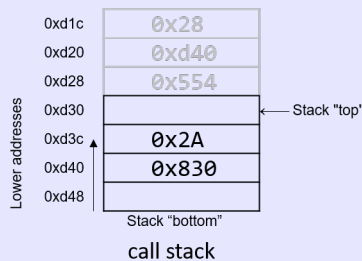
Memory

0x400604	"x is %d\n"
----------	-------------

23. We then call the printf function, which we won't trace through but it outputs to stdout.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
0x55c mov   $0x400604, %edi
0x561 mov   $0x0, %eax
→ 0x566 callq <printf@plt>
0x56b mov   $0x0, %eax
0x570 leaveq
0x571 retq
    
```



Registers	
%eax	0x0
%edi	0x400604
%rsp	0xd30
%rbp	0xd40
%rip	0x56b

Terminal:

```

$ ./prog
x is 42
    
```

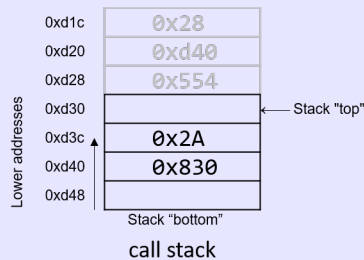
Memory	
0x400604	"x is %d\n"

printf() is called with arguments "x is %d\n" and 42.

24. The print function might have returned something, but we don't care. We want to main function to return 0, so we move 0 into %eax.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
0x55c mov   $0x400604, %edi
0x561 mov   $0x0, %eax
→ 0x566 callq <printf@plt>
0x56b mov   $0x0, %eax
0x570 leaveq
0x571 retq
    
```



Registers	
%eax	0x0
%edi	0x400604
%rsp	0xd30
%rbp	0xd40
%rip	0x570

Terminal:

```

$ ./prog
x is 42
    
```

25. Finally we execute leaveq, which prepares the stack for returning from the function call. It essentially moves the base pointer back to the stack pointer and then pops the base pointer off the stack. The new %rbp is the original base pointer of whatever was outside the main function, 0x830.

```

0x542 <main>:
0x542 push  %rbp
0x543 mov   %rsp, %rbp
0x546 sub   $0x10, %rsp
0x54a callq 0x526 <assign>
0x55f callq 0x536 <adder>
0x554 mov   %eax, -0x4(%rbp)
0x557 mov   -0x4(%rbp), %eax
0x55a mov   %eax, %esi
0x55c mov   $0x400604, %edi
0x561 mov   $0x0, %eax
0x566 callq <printf@plt>
0x56b mov   $0x0, %eax
→ 0x570 leaveq
0x571 retq
    
```

Registers	
%eax	0x0
%edi	0x400604
%rsp	0xd48
%rbp	0x830
%rip	0x571

Terminal:

```

$ ./prog
x is 42
    
```

Equivalent to:

```

mov %rbp, %rsp
pop %rbp
    
```

26. Finally, we execute `retq`, which pops the return address off the stack and puts it into `%rip`.

We have omitted the details of caller and callee saved registers, but they do exist and are important for the general implementations.

For arrays, there's not anything new here. Let's go over some code and follow through it.

```

1 int sumArray(int *array, int length) {
2     int i, total = 0;
3     for (i = 0; i < length; i++) {
4         total += array[i];
5     }
6     return total;
7 }
    
```

This function takes the address of an array and the length of it and sums up all the elements in the array.

```

1 0x400686 <+0>: push %rbp           # save %rbp
2 0x400687 <+1>: mov  %rsp,%rbp      # update %rbp (new stack frame)
3 0x40068a <+4>: mov  %rdi,-0x18(%rbp) # copy array to %rbp-0x18
4 0x40068e <+8>: mov  %esi,-0x1c(%rbp) # copy length to %rbp-0x1c
5 0x400691 <+11>: movl $0x0,-0x4(%rbp) # copy 0 to %rbp-0x4 (total)
6 0x400698 <+18>: movl $0x0,-0x8(%rbp) # copy 0 to %rbp-0x8 (i)
7 0x40069f <+25>: jmp  0x4006be <sumArray+56> # goto <sumArray+56>
8 0x4006a1 <+27>: mov  -0x8(%rbp),%eax  # copy i to %eax
9 0x4006a4 <+30>: cltq                    # convert i to a 64-bit integer
10 0x4006a6 <+32>: lea 0x0(,%rax,4),%rdx  # copy i*4 to %rdx
11 0x4006ae <+40>: mov  -0x18(%rbp),%rax  # copy array to %rax
12 0x4006b2 <+44>: add  %rdx,%rax         # compute array+i*4, store in %rax
13 0x4006b5 <+47>: mov  (%rax),%eax       # copy array[i] to %eax
14 0x4006b7 <+49>: add  %eax,-0x4(%rbp)   # add %eax to total
15 0x4006ba <+52>: addl $0x1,-0x8(%rbp)  # add 1 to i (i+=1)
16 0x4006be <+56>: mov  -0x8(%rbp),%eax  # copy i to %eax
    
```



```

17 0x4006c1 <+59>:  cmp  -0x1c(%rbp),%eax    # compare i to length
18 0x4006c4 <+62>:  jl   0x4006a1 <sumArray+27> # if i<length goto <sumArray+27>
19 0x4006c6 <+64>:  mov  -0x4(%rbp),%eax    # copy total to %eax
20 0x4006c9 <+67>:  pop  %rbp               # prepare to leave the function
21 0x4006ca <+68>:  retq                    # return total

```

6.4.6 ARM Instructions

6.4.7 Buffer Overflows

7 Compiling and Linking

Now let's talk about how this compiling actually happens. *Compiling* is actually an umbrella term that is misused. Turning a C file into an executable file consists of multiple intermediate steps, one of which is actually compiling, but the whole series is sometimes referred to as compiling. A more accurate term would be *building*. Before we get onto it, there are two types of compilers.

Definition 7.1 (GCC, CLang)

The two mainstream compilers used is GCC (with the gdb debugger) and Clang (with lldb). For now, the difference is that

1. gcc is more established.
2. clang is newer and has more features.

A useful flag to know is that we can always specify the name of the (final or intermediary) output file with the `-o` flag.

Definition 7.2 (Complete Build Process)

To actually turn a C file into an executable file, we need to go through a series of steps. We start off with the C code, which are the `.c`, `.cpp`, or `.h` files.

1. **Preprocessing:** The precompiler step expands the *preprocessor directives* (all the `#include` and `#define` statements) and removes comments. This results in a `.i` file. The preprocessor will replace these macros with the actual code. This results in a `.i` file.

```
1 clang/gcc -E main.c -o main.i
```

2. **Compiling:** We take these and generate assembly code. This results in a `.asm` or `.s` file.

```
1 clang/gcc -S main.c -o main.s
```

3. **Assembler:** We take the assembly code and generate machine code in the form of relocatable binary object code (this is machine code, not assembly). This results in a `.o` or `.obj` file.

```
1 clang/gcc -c main.c -o main.o
```

4. **Linking:** We take these object files and link them together to form an executable file. This results in a `.exe` or `.out` file.

The GCC or CLang compiler automates this process for us. For example, `gcc -c hello.c` generates an object file, taking care of the preprocessing, compiling, and assembling code. Then, `gcc hello.o` links the object file to generate an executable file.

There are a lot of questions to be asked here, and we will go through them step by step.

7.1 Precompiling Stage

Just like how Python package managers like conda have specific directories that they find package in, the C library also has a certain directory.

Definition 7.3 (Standard Library Directory)

In Linux systems, there are two main directories you look at:

1. `/usr/include` contains the standard C library headers.
2. `/usr/local/include` contains the headers for libraries that you install yourself.

In Mac Silicon, these directories are a little bit more involved. You must first install the xcode command line developer tools, which will then create these directories.

1. The standard C library headers are in

`/Library/Developer/CommandLineTools/SDKs/MacOSX*.sdk/usr/include.`

In here, we can find all the relevant import files like `stdio.h` and such. When we precompile, the output `.i` file represents a precompiled C file. It still has C code, but it has been optimized to

1. Remove comments.
2. Replace all the `#include` statements with the actual code.
3. Replace all the global variables declared in `#define` with the actual value.

Between x86 and ARM, there are no significant differences in how C files are precompiled.

Example 7.1 ()

Take a look at the following minimal example.

<pre> 1 #include "second.h" 2 #define a 3 3 4 int add(int x, int y) { 5 return x + y; 6 } 7 8 int main() { 9 // test comment 10 int b = 5; 11 int c = add(a, b); 12 int d = subtract(a, b); 13 return 0; 14 }</pre>	<pre> 1 int subtract(int a, int b) { 2 return a - b; 3 } 4 . 5 . 6 . 7 . 8 . 9 . 10 . 11 . 12 . 13 . 14 .</pre>
--	--

Figure 31: I have included a `main.c` file that imports statements from a `second.h` file.

Now, I run `gcc -E main.c -o main.i` to generate the precompiled file, which gives me the following.

```
1 # 1 "main.c"
2 # 1 "<built-in>" 1
3 # 1 "<built-in>" 3
4 # 418 "<built-in>" 3
5 # 1 "<command line>" 1
6 # 1 "<built-in>" 2
7 # 1 "main.c" 2
8 # 1 "./second.h" 1
9 int subtract(int a, int b) {
10     return a - b;
11 }
12 # 2 "main.c" 2
13
14
15 int add(int x, int y) {
16     return x + y;
17 }
18
19 int main() {
20
21     int b = 5;
22     int c = add(3, b);
23     int d = subtract(3, b);
24     return 0;
25 }
```

Figure 32: The precompiled file.

Notice a few things:

1. The header file `second.h` has been replaced with the actual code.
2. The comments have indeed been removed.
3. The global variable `a` has been replaced with the actual value 3.

This leaves us with the question of what all the rest of the lines that start with a `#` are for. They are called *preprocessor directives*.

Definition 7.4 (Preprocessor Directives)

Preprocessor directives are commands that are executed before the actual compilation begins. These directives allow additional actions to be taken on the C source code before it is compiled into object code. Directives are not part of the C language itself, and they are always prefixed with a `#` symbol.

1. `#include` is used to include the contents of a file into the source file. It selects portions of the file to include based on the file name.
2. `#define` is used to define a macro, which is a way to give a name to a constant value or a piece of code.
3. `#ifdef`, `#ifndef`, `#else`, and `#endif` are used for conditional compilation.
4. `#error` is used to generate a compilation error.
5. `#pragma` is used to give the compiler specific instructions.

7.2 Compiling Stage

Once we have precompiled, we can compile the code into assembly code. For the following two examples, we will parse through the general syntax of assembly code. It is quite different between x86 and ARM, so we will use the minimal C code

```

1  int add(int x, int y) {
2      return x + y;
3  }
4
5  int main() {
6      int a = 3;
7      int b = 5;
8      int c = add(a, b);
9      return 0;
10 }

```

for both examples.

Example 7.2 (x86 Compiled Assembly Language)

The assembly code is shown.

```

1  .
2  .file "main.c"
3  .text
4  .globl add
5  .type add, @function
6  add:
7  .LFB0:
8  .cfi_startproc
9  endbr64
10 pushq %rbp
11 .cfi_def_cfa_offset 16
12 .cfi_offset 6, -16
13 movq %rsp, %rbp
14 .cfi_def_cfa_register 6
15 movl %edi, -4(%rbp)
16 movl %esi, -8(%rbp)
17 movl -4(%rbp), %edx
18 movl -8(%rbp), %eax
19 addl %edx, %eax
20 popq %rbp
21 .cfi_def_cfa 7, 8
22 ret
23 .cfi_endproc
24 .LFE0:
25 .size add, .-add
26 .globl main
27 .type main, @function
28 main:
29 .LFB1:
30 .cfi_startproc
31 endbr64
32 pushq %rbp
33 .cfi_def_cfa_offset 16
34 .cfi_offset 6, -16

```

```

35  movq  %rsp, %rbp
36  .cfi_def_cfa_register 6
37  subq  $16, %rsp
38  movl  $3, -12(%rbp)
39  movl  $5, -8(%rbp)
40  movl  -8(%rbp), %edx
41  movl  -12(%rbp), %eax
42  movl  %edx, %esi
43  movl  %eax, %edi
44  call  add
45  movl  %eax, -4(%rbp)
46  movl  $0, %eax
47  leave
48  .cfi_def_cfa 7, 8
49  ret
50  .cfi_endproc
51  .LFE1:
52  .size  main, .-main
53  .ident  "GCC: (Ubuntu 9.4.0-1ubuntu1~20.04.2) 9.4.0"
54  .section .note.GNU-stack,"",@progbits
55  .section .note.gnu.property,"a"
56  .align 8
57  .long  1f - 0f
58  .long  4f - 1f
59  .long  5
60  0:
61  .string  "GNU"
62  1:
63  .align 8
64  .long  0xc0000002
65  .long  3f - 2f
66  2:
67  .long  0x3
68  3:
69  .align 8
70  4:

```

Example 7.3 (ARM Compiled Assembly Language)

The assembly code is shown.

```

1  .
2  .section __TEXT,__text,regular,pure_instructions
3  .build_version macos, 14, 0 sdk_version 14, 4
4  .globl _add ; -- Begin function add
5  .p2align 2
6  _add: ; @add
7  .cfi_startproc
8  ; %bb.0:
9  sub sp, sp, #16
10 .cfi_def_cfa_offset 16
11 str w0, [sp, #12]
12 str w1, [sp, #8]
13 ldr w8, [sp, #12]
14 ldr w9, [sp, #8]

```

```

15  add w0, w8, w9
16  add sp, sp, #16
17  ret
18  .cfi_endproc
19                                     ; -- End function
20  .globl _main                          ; -- Begin function main
21  .p2align 2
22  _main:                                ; @main
23  .cfi_startproc
24  ; %bb.0:
25  sub sp, sp, #48
26  .cfi_def_cfa_offset 48
27  stp x29, x30, [sp, #32]                ; 16-byte Folded Spill
28  add x29, sp, #32
29  .cfi_def_cfa w29, 16
30  .cfi_offset w30, -8
31  .cfi_offset w29, -16
32  mov w8, #0
33  str w8, [sp, #12]                       ; 4-byte Folded Spill
34  stur wzr, [x29, #-4]
35  mov w8, #3
36  stur w8, [x29, #-8]
37  mov w8, #5
38  stur w8, [x29, #-12]
39  ldur w0, [x29, #-8]
40  ldur w1, [x29, #-12]
41  bl _add
42  mov x8, x0
43  ldr w0, [sp, #12]                       ; 4-byte Folded Reload
44  str w8, [sp, #16]
45  ldp x29, x30, [sp, #32]                 ; 16-byte Folded Reload
46  add sp, sp, #48
47  ret
48  .cfi_endproc
49                                     ; -- End function
50  .subsections_via_symbols

```

We can see that in both examples, there are generally two types of codes.

1. The regular CPU operations with registers and memory.
2. Some code starts off with some code that starts with a `..`. Every line that starts with a `.` are called *assembler directives*.

Let's elaborate more on what these directives are.

Definition 7.5 (Assembler Directives)

An **assembler directives** are instructions in assembly language programming that that give commands to the assembler (which then converts this to an object file) about various aspects of the assembly process, but they do not represent actual CPU instructions that execute in the program. Unlike typical assembly language instructions that directly manipulate registers and execute arithmetic or logical operations, directives are used to organize, control, and provide necessary information for the assembly and linking of binary programs. They can manage memory allocation, define symbols, control compilation settings, and much more.

There are general types of directives that are common in both x86 and ARM that we should be aware

about:

1. Section directives.
2. Data allocation directives.
3. Symbol definition directives.
4. Macro and Include directives.
5. Debugging and error handling directives.

Example 7.4 (x86 Assembly Directives)

Let us elaborate on the specific directives in the x86 assembly code, some of which are in the example above.

1. `.file "main.c"` is a directive that tells the assembler that the following code is from the file `main.c`. It is a form of metadata.
2. `.text` is a directive that tells the assembler that the following code is the text section (the text/code portion of memory) of the program. This is where the actual code is stored.
3. `.globl add` is a directive that tells the assembler that the following code is a global function called `add`.
4. `.type add, @function` is a directive that tells the assembler that the following code is a function.

Example 7.5 (ARM Assembly Directives)

You also see that there are symbols that represent memory addresses. Let's elaborate on what symbols mean.

Definition 7.6 (Symbol)

A **symbol** is a name that is used to refer to a memory location. It can be a function name, a global variable, or a local variable.

1. Global symbols are symbols that can be referenced by other object files, e.g. non-static functions and global variables.
2. Local symbols are symbols that are only visible within the object file, e.g. static functions and local variables. The linker won't know about these types.
3. External symbols are referenced by this object file but defined in another object file.

7.3 Objdump

Since we will be using the `objdump` package quite a lot, it is worth mentioning the different commands you will use and store them here as a reference. For first readers, don't expect to know what each of them do, but rather look back at this for a reference.

7.3.1 ELF and Mach-O Formats

Objdump is a command line utility that is used to display information about object files, which are often outputted in a specific format. The two main output file types are called ELF (Executable and Linkable Format) and Mach-O (Mach Object).

Definition 7.7 (ELF)

The **Executable and Linkable Format** (ELF) is a common standard file format for executables, object code, shared libraries, and core dumps. It is analogous to a book, with the following parts:

1. **Header**, which is like the cover of the book. It contains metadata about the file, such as the architecture, the entry point, and the sections.
2. **Sections**, which are like chapters. Each section contains the content for some given purpose or use within the program. e.g. `.binary` is just a block of bytes, `.text` contains the machine code, `.data` contains initialized data, and `.bss` contains uninitialized data.
3. **Symbol Table**, is like a detailed table of contents of all defined symbols such as functions, external (global) variables, local maps, etc.
4. **Relocation records**, which is like the index of the book that lists references to symbols.

The format is generally as such when you run `objdump -d -r hello.o` (d represents disassembly and r represents relocation entries).

```

1  ELF header          # file type
2
3  .text section
4  - code goes here
5
6  .rodata section
7  - read only data
8
9  .data section
10 - initialized global variables
11
12 .bss section
13 - uninitialized global variables
14
15 .symtab section
16 - symbol table (symbol name, type, address)
17
18 .rel.text section
19 - relocation entries for .text section
20 - addresses of instructions that will need to be modified in the executable.
21
22 .rel.data section
23 - relocation info for .data section
24 - addresses of pointer data that will need to be modified in the merged executable.
25
26 .debug section
27 - info for symbolic debugging (gcc -g)

```

Definition 7.8 (Mach-O)

7.3.2 Objdump Commands

Theorem 7.1 (File Headers with Objdump)

Given that you have an object file, the first thing you might want to do is see the file header. You do with this `objdump -f main.o`.

```

1  main.o:      file format elf64-x86-64
2  architecture: i386:x86-64, flags 0x00000011:
3  HAS_RELOC, HAS_SYMS
4  start address 0x0000000000000000

```


Theorem 7.2 (Section with Objdump)

To look at the section headers to get a closer overview, you use `objdump -h main.o`.

```

1 main.o:      file format elf64-x86-64
2
3 Sections:
4 Idx Name          Size      VMA              LMA              File off  Algn
5  0 .text          0000004b  0000000000000000 0000000000000000 00000040 2**0
6                CONTENTS, ALLOC, LOAD, RELOC, READONLY, CODE
7  1 .data          00000000  0000000000000000 0000000000000000 0000008b 2**0
8                CONTENTS, ALLOC, LOAD, DATA
9  2 .bss           00000000  0000000000000000 0000000000000000 0000008b 2**0
10               ALLOC
11  3 .comment       0000002c  0000000000000000 0000000000000000 0000008b 2**0
12               CONTENTS, READONLY
13  4 .note.GNU-stack 00000000  0000000000000000 0000000000000000 000000b7 2**0
14               CONTENTS, READONLY
15  5 .note.gnu.property 00000020  0000000000000000 0000000000000000 000000b8 2**3
16               CONTENTS, ALLOC, LOAD, READONLY, DATA
17  6 .eh_frame      00000058  0000000000000000 0000000000000000 000000d8 2**3
18               CONTENTS, ALLOC, LOAD, RELOC, READONLY, DATA

```

Theorem 7.3 (Disassembly with Objdump)

Now you might actually want to look at the disassembly of the code, which is what we often use it for. To do this, you use `objdump -D main.o` to get the entire output.

1. The leftmost column represents the address of the instruction.
2. The next column represents the machine code of the instruction.
3. The next column represents the assembly code of the instruction.

```

1 main.o:      file format elf64-x86-64
2
3 Disassembly of section .text:
4
5 0000000000000000 <add>:
6   0: f3 0f 1e fa          endbr64
7   ...
8  17: c3                  retq
9
10 0000000000000018 <main>:
11  18: f3 0f 1e fa          endbr64
12  ...
13  4a: c3                  retq
14
15 Disassembly of section .comment:
16
17 0000000000000000 <.comment>:
18   0: 00 47 43             add    %al,0x43(%rdi)
19  ...
20  2a: 30 00                xor    %al,(%rax)
21
22 Disassembly of section .note.gnu.property:
23
24 0000000000000000 <.note.gnu.property>:

```

```

25  0: 04 00          add   $0x0,%al
26  ...
27
28  Disassembly of section .eh_frame:
29
30  0000000000000000 <.eh_frame>:
31  0: 14 00          adc   $0x0,%al
32  ...

```

If you just want to look at the contents of the executable sections, then you can use `objdump -d main.o`.

```

1  main.o:      file format elf64-x86-64
2
3  Disassembly of section .text:
4
5  0000000000000000 <add>:
6  0: f3 0f 1e fa          endbr64
7  4: 55                    push  %rbp
8  5: 48 89 e5             mov   %rsp,%rbp
9  8: 89 7d fc             mov   %edi,-0x4(%rbp)
10 b: 89 75 f8            mov   %esi,-0x8(%rbp)
11 e: 8b 55 fc             mov   -0x4(%rbp),%edx
12 11: 8b 45 f8            mov   -0x8(%rbp),%eax
13 14: 01 d0              add   %edx,%eax
14 16: 5d                  pop   %rbp
15 17: c3                  retq
16
17 0000000000000018 <main>:
18 18: f3 0f 1e fa          endbr64
19 1c: 55                    push  %rbp
20 1d: 48 89 e5             mov   %rsp,%rbp
21 20: 48 83 ec 10          sub   $0x10,%rsp
22 24: c7 45 f4 03 00 00 00 movl  $0x3,-0xc(%rbp)
23 2b: c7 45 f8 05 00 00 00 movl  $0x5,-0x8(%rbp)
24 32: 8b 55 f8            mov   -0x8(%rbp),%edx
25 35: 8b 45 f4            mov   -0xc(%rbp),%eax
26 38: 89 d6              mov   %edx,%esi
27 3a: 89 c7              mov   %eax,%edi
28 3c: e8 00 00 00 00      callq 41 <main+0x29>
29 41: 89 45 fc            mov   %eax,-0x4(%rbp)
30 44: b8 00 00 00 00      mov   $0x0,%eax
31 49: c9                  leaveq
32 4a: c3                  retq

```

If you want to see the source code intermixed with disassembly, then you can use the `-S` flag, but make sure that the object file is a generated with debugging information, i.e. use `gcc -c -g main.c -o main.o`.

```

1  main.o:      file format elf64-x86-64
2
3
4  Disassembly of section .text:
5
6  0000000000000000 <add>:
7  int add(int x, int y) {
8      0: f3 0f 1e fa      endbr64
9      4: 55                  push  %rbp
10     5: 48 89 e5             mov   %rsp,%rbp
11     8: 89 7d fc             mov   %edi,-0x4(%rbp)
12     b: 89 75 f8             mov   %esi,-0x8(%rbp)
13     return x + y;
14     e: 8b 55 fc             mov   -0x4(%rbp),%edx
15    11: 8b 45 f8             mov   -0x8(%rbp),%eax
16    14: 01 d0                add   %edx,%eax
17  }
18    16: 5d                  pop   %rbp
19    17: c3                  retq
20
21  0000000000000018 <main>:
22
23  int main() {
24    18: f3 0f 1e fa      endbr64
25    1c: 55                  push  %rbp
26    1d: 48 89 e5             mov   %rsp,%rbp
27    20: 48 83 ec 10         sub   $0x10,%rsp
28    int a = 3;
29    24: c7 45 f4 03 00 00 00  movl  $0x3,-0xc(%rbp)
30    int b = 5;
31    2b: c7 45 f8 05 00 00 00  movl  $0x5,-0x8(%rbp)
32    int c = add(a, b);
33    32: 8b 55 f8             mov   -0x8(%rbp),%edx
34    35: 8b 45 f4             mov   -0xc(%rbp),%eax
35    38: 89 d6                mov   %edx,%esi
36    3a: 89 c7                mov   %eax,%edi
37    3c: e8 00 00 00 00     callq 41 <main+0x29>
38    41: 89 45 fc             mov   %eax,-0x4(%rbp)
39    return 0;
40    44: b8 00 00 00 00     mov   $0x0,%eax
41  }
42    49: c9                  leaveq
43    4a: c3                  retq

```

Figure 33: Disassembly of the object file back into assembly using `objdump -d -S main.o`.

Note that you can always see this disassembly with debuggers like `gdb` or `lldb`, but `objdump` generally works for all architectures.

Theorem 7.4 (Symbol Table)

If you want to look at all the symbols existing within the object file, you use `objdump -t main.o` (t for table of symbols).

1. The leftmost column represents the address of the symbol.

2. The next column represents the type of the symbol. The `g` and `l` represent global and local symbols, respectively. The `O` and `F` represent object and function symbols, while the `UND` and `ABS` represent undefined and absolute symbols.
3. The next column represents the section that the symbol is in.
4. The next column represents the size of the symbol.
5. The last column represents the name of the symbol.

```

1  main.o:      file format elf64-x86-64
2
3  SYMBOL TABLE:
4  0000000000000000  l    df *ABS*  0000000000000000  main.c
5  0000000000000000  l    d  .text  0000000000000000  .text
6  0000000000000000  l    d  .data  0000000000000000  .data
7  0000000000000000  l    d  .bss  0000000000000000  .bss
8  0000000000000000  l    d  .note.GNU-stack  0000000000000000  .note.GNU-stack
9  0000000000000000  l    d  .note.gnu.property  0000000000000000  .note.gnu.property
10 0000000000000000  l    d  .eh_frame  0000000000000000  .eh_frame
11 0000000000000000  l    d  .comment  0000000000000000  .comment
12 0000000000000000  g    F  .text  0000000000000018  add
13 0000000000000018  g    F  .text  0000000000000033  main
    
```

Theorem 7.5 (Relocation Table)

If you want to look then at the relocation table, then you use `objdump -r main.o`.

1. The leftmost column represents the offset of the relocation (i.e. the location within the section where this relocation needs to be applied).
2. The second column represents the type of relocation.
3. The third column represents the symbol that this relocation references.

```

1  main.o:      file format elf64-x86-64
2
3  RELOCATION RECORDS FOR [.text]:
4  OFFSET          TYPE             VALUE
5  000000000000003d R_X86_64_PLT32  add-0x0000000000000004
6
7
8  RELOCATION RECORDS FOR [.eh_frame]:
9  OFFSET          TYPE             VALUE
10 0000000000000020 R_X86_64_PC32   .text
11 0000000000000040 R_X86_64_PC32   .text+0x0000000000000018
    
```

7.4 Assembling Stage and Object Files

Now, once you have gotten the object file, you cannot simply open it up in a text edit as it is in machine code. To actually interpret anything from it, you must **disassemble** it, meaning that you convert the machine code back into assembly code. The main software that you use to do this is `objdump`. Let's take a look again at the object file.

```

1 Disassembly of section .text:
2
3 0000000000000000 <add>:
4   0: f3 0f 1e fa          endbr64
5   4: 55                    push   %rbp
6   5: 48 89 e5              mov    %rsp,%rbp
7   8: 89 7d fc              mov    %edi,-0x4(%rbp)
8   b: 89 75 f8              mov    %esi,-0x8(%rbp)
9   e: 8b 55 fc              mov    -0x4(%rbp),%edx
10  11: 8b 45 f8              mov    -0x8(%rbp),%eax
11  14: 01 d0                add   %edx,%eax
12  16: 5d                    pop   %rbp
13  17: c3                    retq
14
15 0000000000000018 <main>:
16  18: f3 0f 1e fa          endbr64
17  1c: 55                    push   %rbp
18  1d: 48 89 e5              mov    %rsp,%rbp
19  20: 48 83 ec 10          sub   $0x10,%rsp
20  24: c7 45 f4 03 00 00 00 movl  $0x3,-0xc(%rbp)
21  2b: c7 45 f8 05 00 00 00 movl  $0x5,-0x8(%rbp)
22  32: 8b 55 f8              mov    -0x8(%rbp),%edx
23  35: 8b 45 f4              mov    -0xc(%rbp),%eax
24  38: 89 d6                mov   %edx,%esi
25  3a: 89 c7                mov   %eax,%edi
26  3c: e8 00 00 00 00      callq 41 <main+0x29>
27  41: 89 45 fc              mov   %eax,-0x4(%rbp)
28  44: b8 00 00 00 00      mov   $0x0,%eax
29  49: c9                    leaveq
30  4a: c3                    retq

```

Figure 34: Disassembly of the object file back into assembly using `objdump -d main.o`.

Let's note a couple things.

1. The functions are organized by their starting address followed by their name, e.g.

```

1 0000000000000000 <add>:

```

Within each function, each line of assembly code is shown. To find the total memory the function takes up, you can just take the address of the last line and subtract it from the address of the first line. Or you can literally count the number of bytes in each line (remember 2 hex is 1 byte).

2. The line that calls the `add` function is `0x0 (00 00 00 00)`, with is the *relative target address* intended to be filled in by the linker. The actual assembly line just says that the function continues on to the next line at address `0x41`. This is because the object file is not aware of where it will be loaded into memory, and all lines with this opcode `e8 00 00 00 00` is intended to be filled in by the linker.
3. Look at address `0x3c`. It is calling another function, but the values starting from address `0x3d` is `00 00 00 00`, which is not the actual address of the function but also a dummy address. This is because the object file is not aware of where the function is located in memory.

7.5 Linking Stage and Relocation

7.5.1 Relocation

If the object file is already in machine code, then why do we need a separate linking stage that converts `main.o` into `main` the binary? The reason is stated in the previous section: because the object files uses relative memory addressing and does not know about which memory is accessed in other object files, we need to **relocate** the symbols in the object file to their proper addresses. So how does the linker actually know how to relocate these symbols into their proper addresses? It uses the *relocation table*, which contains information about the addresses that need to be modified in the object file.

```

1  main.o:      file format elf64-x86-64
2
3  RELOCATION RECORDS FOR [.text]:
4  OFFSET          TYPE          VALUE
5  000000000000003d R_X86_64_PLT32  add-0x0000000000000004
6
7
8  RELOCATION RECORDS FOR [.eh_frame]:
9  OFFSET          TYPE          VALUE
10 0000000000000020 R_X86_64_PC32  .text
11 0000000000000040 R_X86_64_PC32  .text+0x0000000000000018

```

Figure 35: Relocation table for `main.o` object file.

Let’s talk about how to actually read this table. We can look at the first entry, which shows an offset of `0x3d`. This represents the offset from the beginning of the `.text` section where the relocation needs to be applied. Looking back at the disassembly file, this address `0x3d` is precisely where there was a dummy address `00 00 00 00`. We want to replace this with the actual address defined in the `VALUE` column, which is `add` (with a slight offset of `0x4`, which is typically used to compensate for the PC-relative addressing mode where the CPU might be adding the length of the instruction to the program counter (PC) before the relocation value is applied). The type of relocation won’t be covered in our scope. Let’s go through each relocation entry:

1. The first entry is for the `add` function. If we look at the disassembly, within the `main` function, the address `0x3d` is where the `add` function is called. The linker will replace the dummy address with the actual address of the `add` function.

```

1  Disassembly of section .text:
2
3  0000000000000000 <add>:
4      0: f3 0f 1e fa      endbr64
5      4: 55                push   %rbp
6      5: 48 89 e5          mov    %rsp,%rbp
7      8: 89 7d fc          mov    %edi,-0x4(%rbp)
8      b: 89 75 f8          mov    %esi,-0x8(%rbp)
9      e: 8b 55 fc          mov    -0x4(%rbp),%edx
10     11: 8b 45 f8          mov    -0x8(%rbp),%eax
11     14: 01 d0            add   %edx,%eax
12     16: 5d                pop    %rbp
13     17: c3                retq
14
15 0000000000000018 <main>:
16     18: f3 0f 1e fa      endbr64
17     1c: 55                push   %rbp
18     1d: 48 89 e5          mov    %rsp,%rbp
19     20: 48 83 ec 10       sub   $0x10,%rsp

```

```

20 24: c7 45 f4 03 00 00 00    movl    $0x3,-0xc(%rbp)
21 2b: c7 45 f8 05 00 00 00    movl    $0x5,-0x8(%rbp)
22 32: 8b 55 f8                    mov     -0x8(%rbp),%edx
23 35: 8b 45 f4                    mov     -0xc(%rbp),%eax
24 38: 89 d6                        mov     %edx,%esi
25 3a: 89 c7                        mov     %eax,%edi
26 3c: e8 00 00 00 00            callq   41 <main+0x29>    <-- here
27 41: 89 45 fc                    mov     %eax,-0x4(%rbp)
28 44: b8 00 00 00 00            mov     $0x0,%eax
29 49: c9                          leaveq
30 4a: c3                          retq

```

2. The second and third entries are for the `.eh_frame` section. We can see that the offset of `0x20` and `0x40` represents the following lines below. They also have dummy addresses that need to be replaced. They are replaced by the address `.text`, which represents the first address in the `.text` section, i.e. the address of the `add` function, and the address `.text+0x18`, which represents the address of the `main` function.

```

1  Disassembly of section .eh_frame:
2
3  0000000000000000 <.eh_frame>:
4    0: 14 00                    adc     $0x0,%al
5    2: 00 00                    add     %al,(%rax)
6    4: 00 00                    add     %al,(%rax)
7    6: 00 00                    add     %al,(%rax)
8    8: 01 7a 52                add     %edi,0x52(%rdx)
9    b: 00 01                    add     %al,(%rcx)
10   d: 78 10                    js      1f <.eh_frame+0x1f>
11   f: 01 1b                    add     %ebx,(%rbx)
12  11: 0c 07                    or      $0x7,%al
13  13: 08 90 01 00 00 1c        or      %dl,0x1c000001(%rax)
14  19: 00 00                    add     %al,(%rax)
15  1b: 00 1c 00                add     %bl,(%rax,%rax,1)
16  1e: 00 00                    add     %al,(%rax)
17  20: 00 00                    add     %al,(%rax)    <-- here for 2nd entry
18  22: 00 00                    add     %al,(%rax)
19  24: 18 00                    sbb    %al,(%rax)
20  26: 00 00                    add     %al,(%rax)
21  28: 00 45 0e                add     %al,0xe(%rbp)
22  2b: 10 86 02 43 0d 06        adc     %al,0x60d4302(%rsi)
23  31: 4f 0c 07                rex.WRXB or $0x7,%al
24  34: 08 00                    or      %al,(%rax)
25  36: 00 00                    add     %al,(%rax)
26  38: 1c 00                    sbb    $0x0,%al
27  3a: 00 00                    add     %al,(%rax)
28  3c: 3c 00                    cmp     $0x0,%al
29  3e: 00 00                    add     %al,(%rax)
30  40: 00 00                    add     %al,(%rax)    <-- here for 3rd entry
31  42: 00 00                    add     %al,(%rax)
32  44: 33 00                    xor     (%rax),%eax

```

Therefore, we can see that the object file generates a “skeleton” code that contains all the instructions, with some dummy addresses that need to be replaced. The relocation table *T* tells us exactly where these dummy addresses are in the code and what they need to be replaced with. Therefore, if we want to call a function `printf` that is in the `text` section at address `0x30`, then we can actually look at the value at `T[30]` to see where the actual address is. At this point, note that we still do not know the actual memory address of `add`. This is determined by the linker.

7.5.2 Linking with One Object File

Now let's see what happens once we link the object file `main.o` into the final executable `main`. If we disassemble it, then we can see a few things:

1. The addresses of all the functions have been changed. `add` starts on address `0x1129` rather than `0x0` and `main` starts on address `0x1141` rather than `0x18`.
2. The dummy address `0x0` of the call to function `add` in `main` have been replaced with the actual addresses `0x1129`.

```

1  0000000000001129 <add>:
2  1129: f3 0f 1e fa          endbr64
3  112d: 55                   push   %rbp
4  112e: 48 89 e5             mov    %rsp,%rbp
5  1131: 89 7d fc             mov    %edi,-0x4(%rbp)
6  1134: 89 75 f8             mov    %esi,-0x8(%rbp)
7  1137: 8b 55 fc             mov    -0x4(%rbp),%edx
8  113a: 8b 45 f8             mov    -0x8(%rbp),%eax
9  113d: 01 d0               add    %edx,%eax
10 113f: 5d                   pop    %rbp
11 1140: c3                   retq
12
13 0000000000001141 <main>:
14 1141: f3 0f 1e fa          endbr64
15 1145: 55                   push   %rbp
16 1146: 48 89 e5             mov    %rsp,%rbp
17 1149: 48 83 ec 10          sub    $0x10,%rsp
18 114d: c7 45 f4 03 00 00 00 movl   $0x3,-0xc(%rbp)
19 1154: c7 45 f8 05 00 00 00 movl   $0x5,-0x8(%rbp)
20 115b: 8b 55 f8             mov    -0x8(%rbp),%edx
21 115e: 8b 45 f4             mov    -0xc(%rbp),%eax
22 1161: 89 d6               mov    %edx,%esi
23 1163: 89 c7               mov    %eax,%edi
24 1165: e8 bf ff ff ff      callq  1129 <add>      <-- replaced with actual address
25 116a: 89 45 fc             mov    %eax,-0x4(%rbp)
26 116d: b8 00 00 00 00      mov    $0x0,%eax
27 1172: c9                   leaveq
28 1173: c3                   retq
29 1174: 66 2e 0f 1f 84 00 00 nopw   %cs:0x0(%rax,%rax,1)
30 117b: 00 00 00
31 117e: 66 90               xchg  %ax,%ax

```

7.5.3 Global vs External Symbols

So far, we have talked about using the `#include` as a precompiling command that says “put all the text from this other file right here.” Take the following code for instance.

<pre> 1 // file1.c 2 #include "sum.h" 3 4 int array[2] = {1, 2}; 5 6 int main() { 7 int val = sum(array, 2); 8 return val; 9 } </pre>	<pre> 1 // sum.h 2 int sum(int *a, int n) { 3 int i, s = 0; 4 for (i = 0; i < n; i++) { 5 s += a[i]; 6 } 7 return s; 8 } 9 . </pre>
---	--

Figure 36: Including a header file in `file1.c` to import functions and variables.

However, there is another way to do this. We can use *external symbols* to access. Rather than simply copying and pasting the code into the file, the `extern` keyword marks that the variable or function exists externally to this source file and does not allocate storage for it.

<pre> 1 // main.c 2 extern int sum(int *array, int n); 3 4 int array[2] = {1, 2}; 5 6 int main(void) { 7 int val = sum(array, 2); 8 return val; 9 } </pre>	<pre> 1 // sum.c 2 int sum(int *array, int n) { 3 int i, s = 0 ; 4 for (int i = 0; i < n; i++) { 5 s += array[i]; 6 } 7 return s; 8 } 9 . </pre>
--	---

Figure 37: Using external symbols to access functions and variables.

One is not a replacement for the other, so what advantage does this have? Well, as we will see, if we have multiple object (source) files, say `A.c`, `B.c`, and `C.c`, that need to reference the same function or variable `var` in `ext.c`, then how would we do this? If we simply put `#include "ext.h"` in all the files, then we would have multiple copies of the same code. This means that for each source there would be its own copy of `var` created and the linker would be unable to resolve this symbol. However, if we put `extern int var;` at the top of each source file, then only one copy of `var` would be created (in `ext.c`), which creates a single instance of `var` for the linker to resolve.²

Therefore, there are three types of symbols (variables, functions, etc.) that we need to consider:

1. **Global symbols** that are defined in the global scope of a C file.
2. **Local symbols** that are defined in the local scope of a C file, e.g. within functions, loops, etc.
3. **External symbols** that are defined in another C file referenced by the `extern` keyword.

Linkers will only know about global and external symbols, and will have no idea that any local symbols exist. With the information of these two types of symbols and the relocation tables of each object file, the linker can then resolve the addresses of all the symbols in the final binary.

The two types of symbols that the linker will know about are the global and external symbols. We can see that external symbols can be problematic if the object files don't know about each other.

²<https://stackoverflow.com/questions/1330114/whats-the-difference-between-using-extern-and-including-header-files>

Example 7.6 (Global and Local Symbols)

Consider the following code where the left file includes the right file.

<pre> 1 // main.c 2 #include "sum.h" 3 4 int array[2] = {1, 2}; 5 6 int main() { 7 int val = sum(array, 2); 8 return val; 9 }</pre>	<pre> 1 // sum.h 2 int sum(int *a, int n) { 3 int i, s = 0; 4 for (i = 0; i < n; i++) { 5 s += a[i]; 6 } 7 return s; 8 } 9 .</pre>
---	---

In the left file,

1. We define the global symbol `main()`.
2. Inside `main`, `val` is a local symbol so the linker knows nothing about it.
3. The `sum` function is an external symbol, and it references a global symbol that's defined in `sum` the right file.
4. The `array` is a global symbol that is defined in the right file.

In the right file, the linker knows nothing of the local symbols `i` or `s`.

7.5.4 Linking with Multiple Object Files

We have seen the case of linking when we simply have one object file. The relocation was simple since the `.text` section is contiguous and so we needed simple translations of addresses to relocate `add` and `main`, along with whatever other sections and files. Now let's consider the case where we have multiple object files.

<pre> 1 // main.c 2 extern int sum(int *array, int n); 3 4 int array[2] = {1, 2}; 5 6 int main(void) { 7 int val = sum(array, 2); 8 return val; 9 }</pre>	<pre> 1 // sum.c 2 int sum(int *array, int n) { 3 int i, s = 0 ; 4 for (int i = 0; i < n; i++) { 5 s += array[i]; 6 } 7 return s; 8 } 9 .</pre>
---	--

Now they have their own object files shown below, where I also put the source code lines to make it easier to parse. Note that again, in `main.o` the call to function `sum` is a dummy address that needs to be replaced. Furthermore, in both `main.o` and `sum.o`, the `.text` section is at address `0x0`, where the addresses of the function `main` and `sum` are, respectively. This causes an overload in the address space.

To demonstrate what happens, we look at how the disassembly, symbol tables, and relocation tables are updated before (with the object files) and after (in the binary) linking.

Example 7.7 (Disassembly of Object Files)

In here, note that both the `array` and `sum` are not initialized and are therefore set to dummy addresses.

```

1 main.o:      file format elf64-x86-64
2 Disassembly of section .text:
3
4 0000000000000000 <main>:
5     extern int sum(int *array, int n);
6
```

```

7  int array[2] = {1, 2};
8
9  int main(void) {
10     0: f3 0f 1e fa      endbr64
11     4: 55              push  %rbp
12     5: 48 89 e5        mov   %rsp,%rbp
13     8: 48 83 ec 10      sub   $0x10,%rsp
14     int val = sum(array, 2);
15     c: be 02 00 00 00  mov   $0x2,%esi
16     11: 48 8d 3d 00 00 00 00  lea   0x0(%rip),%rdi      # 18 <main+0x18>  <-- dummy
        address
17     18: e8 00 00 00 00      callq 1d <main+0x1d>      <-- dummy
        address
18     1d: 89 45 fc          mov   %eax,-0x4(%rbp)
19     return val;
20     20: 8b 45 fc          mov   -0x4(%rbp),%eax
21 }
22     23: c9              leaveq
23     24: c3              retq

```

```

1  sum.o:      file format elf64-x86-64
2  Disassembly of section .text:
3
4  0000000000000000 <sum>:
5  int sum(int *array, int n) {
6     0: f3 0f 1e fa      endbr64
7     4: 55              push  %rbp
8     5: 48 89 e5        mov   %rsp,%rbp
9     8: 48 89 7d e8      mov   %rdi,-0x18(%rbp)
10    c: 89 75 e4        mov   %esi,-0x1c(%rbp)
11    int i, s = 0;
12    f: c7 45 f8 00 00 00 00  movl  $0x0,-0x8(%rbp)
13    for (int i = 0; i < n; i++) {
14    16: c7 45 fc 00 00 00 00  movl  $0x0,-0x4(%rbp)
15    1d: eb 1d          jmp   3c <sum+0x3c>
16    s += array[i];
17    1f: 8b 45 fc          mov   -0x4(%rbp),%eax
18    22: 48 98          cltq
19    24: 48 8d 14 85 00 00 00  lea   0x0(,%rax,4),%rdx
20    2b: 00
21    2c: 48 8b 45 e8      mov   -0x18(%rbp),%rax
22    30: 48 01 d0          add   %rdx,%rax
23    33: 8b 00          mov   (%rax),%eax
24    35: 01 45 f8          add   %eax,-0x8(%rbp)
25    for (int i = 0; i < n; i++) {
26    38: 83 45 fc 01      addl  $0x1,-0x4(%rbp)
27    3c: 8b 45 fc          mov   -0x4(%rbp),%eax
28    3f: 3b 45 e4          cmp   -0x1c(%rbp),%eax
29    42: 7c db          jl   1f <sum+0x1f>
30    }
31    return s;
32    44: 8b 45 f8          mov   -0x8(%rbp),%eax
33    }
34    47: 5d              pop   %rbp
35    48: c3              retq

```

1. In `main.o` at address `0x0`, we have the `main` function and this is because everything is stored relatively to the start of `main`. Once we have linked, `main` shows the absolute addresses of all the instructions.
2. In instruction 11 in `main.o` we can see that `48 8d 3d` is the `lea` instruction, which is the same as that in `main`. However, the address that it was acting on is `0x0` since the array has not been initialized yet. We can see in `main` that the address is now `0x00002ecf`.
3. The comment in `main` indicates that the final relocated address used to access the `array` is `0x4010`. To see relocated addresses in general, just look for the comments and shift them accordingly.

```

1  main:      file format elf64-x86-64
2
3  0000000000001129 <main>:
4      1129:  f3 0f 1e fa          endbr64
5      112d:  55                   push   %rbp
6      112e:  48 89 e5             mov    %rsp,%rbp
7      1131:  48 83 ec 10         sub   $0x10,%rsp
8      1135:  be 02 00 00 00      mov   $0x2,%esi
9      113a:  48 8d 3d cf 2e 00 00 lea   0x2ecf(%rip),%rdi    # 4010 <array>
10     1141:  e8 08 00 00 00     callq 114e <sum>
11     1146:  89 45 fc             mov   %eax,-0x4(%rbp)
12     1149:  8b 45 fc             mov   -0x4(%rbp),%eax
13     114c:  c9                   leaveq
14     114d:  c3                   retq
15
16  000000000000114e <sum>:
17     114e:  f3 0f 1e fa          endbr64
18     1152:  55                   push   %rbp
19     1153:  48 89 e5             mov   %rsp,%rbp
20     1156:  48 89 7d e8         mov   %rdi,-0x18(%rbp)
21     115a:  89 75 e4             mov   %esi,-0x1c(%rbp)
22     ...

```

Example 7.8 (Symbol Tables of Object Files)

Let's take a look at the symbol table of each file as well. Again, all of the addresses of each symbol are `0s` since they are using relative addressing. The `array` and `main` are global symbols since they reside in the global scope, while the `sum` function is an external and undefined symbol.

```

1  main.o:    file format elf64-x86-64
2
3  SYMBOL TABLE:
4  0000000000000000 1  df *ABS*  0000000000000000 main.c
5  0000000000000000 1  d  .text  0000000000000000 .text
6  0000000000000000 1  d  .data  0000000000000000 .data
7  0000000000000000 1  d  .bss  0000000000000000 .bss
8  0000000000000000 1  d  .note.gnu-stack 0000000000000000 .note.gnu-stack
9  0000000000000000 1  d  .note.gnu.property 0000000000000000 .note.gnu.property
10 0000000000000000 1  d  .eh_frame 0000000000000000 .eh_frame
11 0000000000000000 1  d  .comment 0000000000000000 .comment
12 0000000000000000 g  0  .data  0000000000000008 array
13 0000000000000000 g  F  .text  0000000000000025 main
14 0000000000000000  *UND*  0000000000000000 _GLOBAL_OFFSET_TABLE_
15 0000000000000000  *UND*  0000000000000000 sum

```

```

1  sum.o:      file format elf64-x86-64
2
3  SYMBOL TABLE:
4  0000000000000000 1   df *ABS*  0000000000000000 sum.c
5  0000000000000000 1   d  .text  0000000000000000 .text
6  0000000000000000 1   d  .data  0000000000000000 .data
7  0000000000000000 1   d  .bss  0000000000000000 .bss
8  0000000000000000 1   d  .note.gnu-stack 0000000000000000 .note.gnu-stack
9  0000000000000000 1   d  .note.gnu.property 0000000000000000 .note.gnu.property
10 0000000000000000 1   d  .eh_frame 0000000000000000 .eh_frame
11 0000000000000000 1   d  .comment 0000000000000000 .comment
12 0000000000000000 g   F  .text  0000000000000049 sum
    
```

When we have the linked binary, note a few things.

1. In `main.o`, the numbers on the left represents the address of the symbol (all 0s since we haven't linked yet and their final addresses aren't known), while the addresses in `a.out` are all known.
2. In `main.o`, the `sum` function is an external symbol and is undefined. The linker will need to know where this is. In `main`, note that the `sum` function is now a global symbol and is defined, along with the size. We can now see that all the final addresses of each symbol is known, along with their sizes, and the UND marker is now gone as well.
3. Only the size of the global variable is known in `main.o` since we have defined it within the code. However, in `main`, the linker has now assigned an address to it.
4. To see the size in bytes of the array, you can look at the address and how much size it takes up.

```

1  main:      file format elf64-x86-64
2
3  SYMBOL TABLE:
4  ...
5  0000000000004008 g   0  .data  0000000000000000      .hidden __dso_handle
6  00000000000114e g   F  .text  0000000000000049      sum
7  000000000002000 g   0  .rodata 0000000000000004      _IO_stdin_used
8  0000000000011a0 g   F  .text  0000000000000065      __libc_csu_init
9  0000000000004020 g       .bss  0000000000000000      _end
10 000000000001040 g   F  .text  000000000000002f      _start
11 0000000000004018 g       .bss  0000000000000000      __bss_start
12 0000000000001129 g   F  .text  0000000000000025      main
13 0000000000004018 g   0  .data  0000000000000000      .hidden __TMC_END__
14 ...
    
```

Example 7.9 (Relocation Tables)

Ignoring the `.eh_frame`, in `main.o` the relocation table contains entries for `array` and `sum` that must be relocated.

```

1  main.o:      file format elf64-x86-64
2
3  RELOCATION RECORDS FOR [.text]:
4  OFFSET          TYPE          VALUE
5  0000000000000014 R_X86_64_PC32  array-0x0000000000000004
6  0000000000000019 R_X86_64_PLT32  sum-0x0000000000000004
7
8  RELOCATION RECORDS FOR [.eh_frame]:
9  OFFSET          TYPE          VALUE
10 0000000000000020 R_X86_64_PC32  .text
    
```

```

1  sum.o:      file format elf64-x86-64
2
3  RELOCATION RECORDS FOR [.eh_frame]:
4  OFFSET          TYPE          VALUE
5  0000000000000020 R_X86_64_PC32    .text

```

We can see a couple things. Namely, there is nothing to be relocated in `a.out` since everything has been relocated already by the linker. So let's focus on the relocation for `main.o`. In here, we can see that in the `.text` section, there are two things being relocated:

1. The reference to the global variable `array` is being relocated. In this object file, we look at the offset `0x14` from the beginning of the `.text` section, which contains the instruction that needs to access `array`. This relocation record tells the linker to calculate the 32-bit offset from the instruction (at offset `0x14`) to the start of `array`, then adjust it by subtracting 4 bytes.
2. The reference to the `sum` function is being relocated. In this object file, we look at the offset `0x19` from the beginning of the `.text` section, which contains the instruction that needs to access `sum`. This relocation record tells the linker to calculate the 32-bit offset from the instruction (at offset `0x19`) to the start of the `.plt` section, then adjust it by subtracting 4 bytes.

```

1  main:      file format elf64-x86-64

```

7.6 Compiler Optimization

We have learned the complete process of compilers, but compilers can be a little smarter than just translating code line by line. They also come with flags that can optimize the code.

Definition 7.9 (gcc Optimization)

The gcc compiler can optimize the code with the `-O` flag. To run level 1 optimization, we can write

```

1  gcc -O1 -o main main.c

```

The level of optimizations are listed:

1. Level 1 perform basic optimizations to reduce code size and execution time while attempting to keep compile time to a minimum.
2. Level 2 optimizations include most of GCC's implemented optimizations that do not involve a space-performance trade-off.
3. Level 3 performs additional optimizations (such as function inlining) and may cause the program to take significantly longer to compile.

Let's see what common implementation are.

Definition 7.10 (Constant Folding)

Constants in the code are evaluated at compile time to reduce the number of resulting instructions. For example, in the code snippet that follows, macro expansion replaces the statement `int debug = N-5` with `int debug = 5-5`. Constant folding then updates this statement to `int debug = 0`.

```

1  #define N 5
2  int debug = N - 5; //constant folding changes this statement to debug = 0;

```

Definition 7.11 (Constant Propagation)

Constant propagation replaces variables with a constant value if such a value is known at compile time. Consider the following code segment, where the `if (debug)` statement is replaced with `if (0)`.

```

1  int debug = 0;
2
3  int doubleSum(int *array, int length){
4      int i, total = 0;
5      for (i = 0; i < length; i++){
6          total += array[i];
7          if (debug) {
8              printf("array[%d] is: %d\n", i, array[i]);
9          }
10     }
11     return 2 * total;
12 }
```

Definition 7.12 (Dead Code Elimination)

Dead code elimination removes code that is never executed. For example, in the code snippet that follows, the `if (debug)` statement and its body is removed since the value of `debug` is known to be 0.

```

1  int debug = 0;
2
3  int doubleSum(int *array, int length){
4      int i, total = 0;
5      for (i = 0; i < length; i++){
6          total += array[i];
7          if (debug) {                // remove
8              printf("array[%d] is: %d\n", i, array[i]); // remove
9          }                          // remove
10     }
11     return 2 * total;
12 }
```

Definition 7.13 (Simplifying Expressions)

Some instructions are more expensive than others, so things like

1. `2 * total` may be replaced with `total + total` because addition instruction is less expensive than multiplication.
2. `total * 8` may be replaced with `total << 3`
3. `total % 8` may be replaced with `total & 7`

Note that these optimization techniques are in no way a guarantee that the code will run faster since there are many factors and always edge cases (for example, maybe some localities are lost). Furthermore, compiler optimization will never be able to improve runtime complexity (e.g. by replacing bubble sort with quicksort).

8 Storage Hierarchy

8.1 Expanding on von Neumann Architecture

So far, our model of the computer has been a simple von Neumann architecture which consists of a CPU and memory. However, there are many other intricacies that are extremely important in practice, and we'll expand on each one by one.

Definition 8.1 (Computer Architecture)

In our elaborated computer architecture, a computer consists of the components.

1. A **CPU** that consists of an arithmetic logic unit (ALU), registers, and a **bus interface** that controls the input and output.
2. The **IO bridge** that handles communication between everything.
3. The **system bus** that connects the CPU to the IO bridge.
4. The **memory bus** that connects the memory to the IO bridge.
5. The **IO bus** that connects the IO devices and disk to the IO bridge.
6. **IO devices** like mouse, keyboard, and monitor.
7. The **disk controller and disk** that stores data.

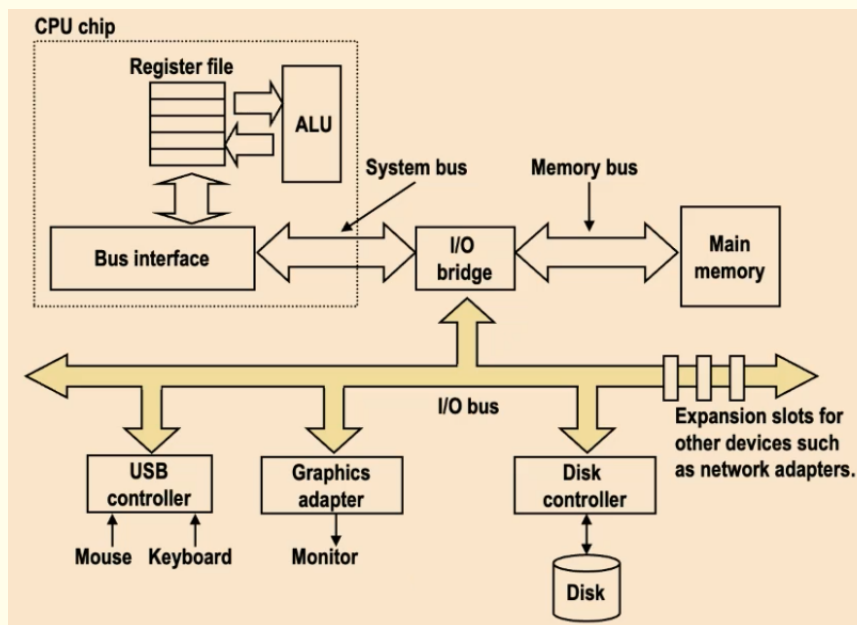


Figure 38: Diagram of the IO bus.

We can see from the diagram above that the CPU can directly access registers (since it's in the CPU itself) and the main memory (since it's connected to the memory bus). However, to access something like the disk, it must go through the disk controller. This gives us our first categorization of memory.

Definition 8.2 (Primary Storage)

Primary storage devices are directly accessible by the CPU and are used to store data that is currently being processed. This includes CPU registers, cache memory, and RAM. In memory, the basic storage unit is normally a **cell** (one bit per cell), which is the physical material that holds information. A **supercell** has address and data widths (number of bits), which is analogous to a lock number and the lock capacity, respectively. It is called random access since it takes approximately the same amount of time to access any cell in memory. There are two primary ways that this is

implemented:

1. **Static RAM (SRAM)** stores data in small electrical circuits (e.g. latches) and is typically the fastest type of memory. However, it is more expensive to build, consumes more power, and occupies more space, limiting the SRAM storage.
2. **Dynamic RAM (DRAM)** stores data using electrical components (e.g. capacitors) that hold an electrical charge. It is called *dynamic* because a DRAM system must frequently refresh the charge of its capacitors to maintain a stored value. It also requires error correction which introduces redundancy.

Device	Capacity	Approx. latency	RAM type
Register	4 - 8 bytes	< 1 ns	SRAM
CPU cache	1 - 32 megabytes	5 ns	SRAM
Main memory	4 - 64 gigabytes	100 ns	DRAM

Table 4: Memory hierarchy characteristics

Definition 8.3 (Secondary Storage)

Secondary storage devices are not directly accessible by the CPU and are used to store data that is not currently being processed. This includes hard drives, SSDs, and magnetic tapes. There are two primary ways:

1. **Spinning disks** store data on a magnetic surface that spins at high speeds.
2. **Solid state drives (SSDs)** store data on flash memory chips.

There are three key components of memory that we should think about:

1. The **capacity**, i.e. amount of data, it can store (how large the water tank is).
2. The **latency**, i.e. amount of time it takes for a device to respond with data after it has been instructed to perform a data retrieval operation (how fast the data flows).
3. The **transfer rate** or **throughput**, i.e. amount of data that can be moved between the device and main memory (how wide the pipe is). Naively, with one channel and sequential transfer the transfer rate is one over the latency.

We must provide a good balance of these three qualities, and also note that there are some physical limitations (i.e. latency cannot be faster than speed of light), and this is more effectively done through a hierarchical memory system.

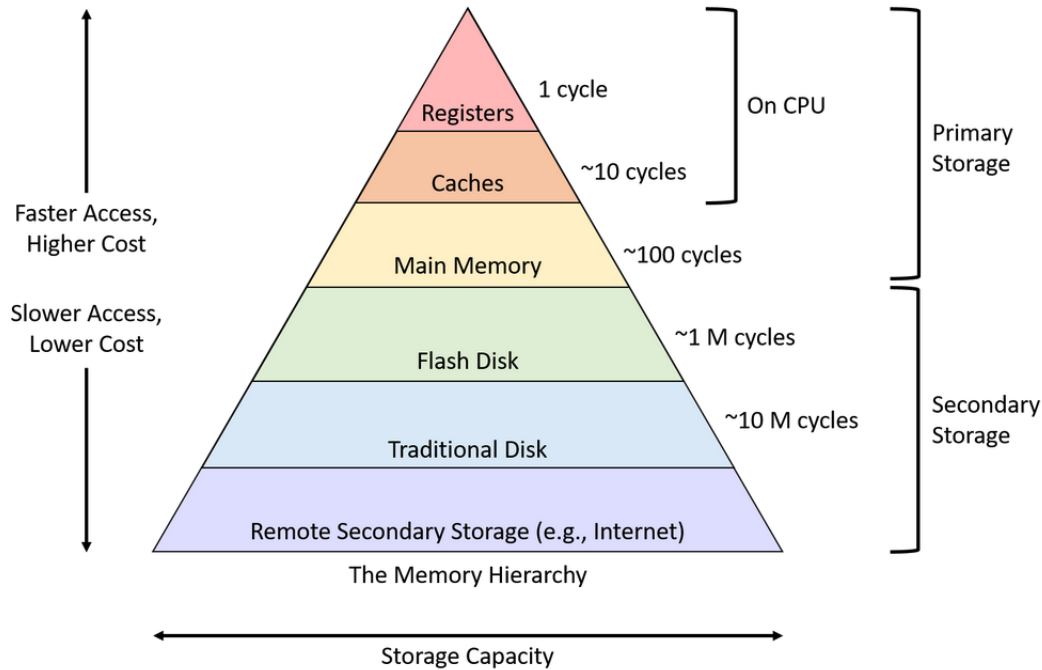


Figure 39: Memory hierarchy.

For example when we want to read from the disk, the CPU must request to the bus interface, which travels through the bus interface, I/O bridge, I/O bus, disk controller, and to the disk itself. Then the data goes back through the disk controller, I/O bus, I/O bridge, through the memory bus, and resides in the main memory. Note that disks are block addressed, so it will transfer the entire block of data into the memory. It must specify a **destination memory address (DMA)**. When the DMA completes, the disk controller notifies the CPU with an *interrupt* (i.e. asserts a special interrupt pin on the CPU), letting it know that the operation has finished. This signal goes through the disk controller to the IO bridge to the CPU. From now on, the CPU knows that there is memory that it can access to run an application loaded in memory.

8.2 Disk

Definition 8.4 (Hard Disk Drives)

Back then, there were **hard disk drives (HDDs)** that literally had a spinning wheel and a needle head that read the data.

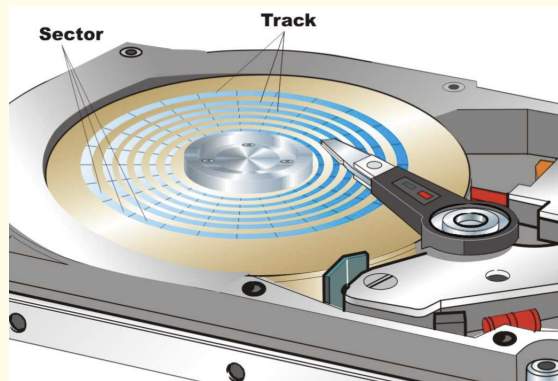


Figure 40: Visual diagram of hard disk drive with its sectors.

1. HDDs are not random access since the data must be sequentially read. This was disadvantageous since the spinning wheel had to spin to the correct location, which took time. The needle also had to move to the correct location, which also took time and therefore read and write speeds were dominated by the time it took to move the needle.
2. The smallest unit of data that can be read is a complete disk sector (not a single byte like RAM).

Definition 8.5 (Solid State Drives)

Now, we have **solid state drives (SSDs)** that store data on flash memory chips. This is advantageous since there are no moving parts, so the latency is much lower and the latency is not dominated by the time it takes to move the needle.

1. SSDs are random access.
2. The smallest unit of data is a **page**, which is usually 4KB and maybe for high scale computers 2-4 MB (but on “Big Data” applications big but computers, it can be up to 1GB).
3. A collection of pages, usually 128 pages, is called a **block**, making is 512KB.

While virtually all RAM and primary storage devices are **byte addressable** (i.e. you can access any byte in memory), secondary storage devices are **block addressable** (i.e. you can only access a block of memory at a time). Therefore, to access a single byte in secondary storage, you must first load the entire block into memory, calculate which byte from that block you want, and then access it. Therefore, you need both the block number x and the offset o to access a byte in secondary storage, which is why it is even slower than accessing RAM.

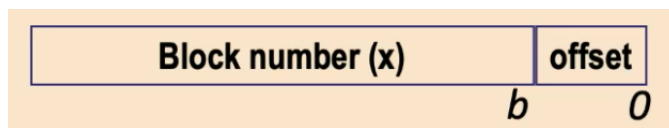


Figure 41: Block offset.

Therefore, you can think of raw data in units of blocks of size 2^b for some b bits.

1. Take the low order b bits of a byte address as an integer, which is the offset of the addressed byte in the block.
2. The rest of the bits are the block number x , which is an unsigned long.
3. You request the block number x , receive the block contents, and then extract the requested byte at

offset in x i.e. calculate `block[x][offset]`.

8.3 Locality

So far, we have abstracted away most of these memory types as a single entity with nearly instantaneous access, but in practice this is not the case. The most simple way is to simply have RAM and our CPU registers, but by introducing more intermediate memory types, we can achieve greater efficiency.

Definition 8.6 (Locality)

Locality is a principle that generally states that a program that accesses a memory location n at time t is likely to access memory location $n + \epsilon$ at time $t + \epsilon$. This principle motivates the design of efficient caches.

1. **Temporal locality** is the idea that if you access a memory location, you are likely to access it again soon.
2. **Spatial locality** is the idea that if you access a memory location, you are likely to access nearby memory locations soon.

This generally means that if you access some sort of memory, the values around that address is also likely to be accessed and therefore it is wise to store it closer to your CPU. In CPUs, both the instructions and the data are stored in the cache, which exploits both kinds of locality (repeated operations for temporal and nearby data for spatial).

Example 8.1 (Locality)

Consider the following code.

```

1  int sum_array(int *array, int len) {
2      int i;
3      int sum = 0;
4
5      for (i = 0; i < len; i++) {
6          sum += array[i];
7      }
8
9      return sum;
10 }
```

1. Temporal Locality

- (a) We cycle through each loop repeatedly with the same add operation, exploiting temporal locality.
- (b) The CPU accesses the same memory (stored in variables `i`, `len`, `sum`, `array`) within each iteration and therefore at similar times.

2. Spatial Locality

- (a) The spatial locality is exploited when the CPU accesses memory locations from each element of the array, which are contiguous in memory.
- (b) Even though the program accesses each array element only once, a modern system loads more than one `int` at a time from memory to the CPU cache. That is, accessing the first array index fills the cache with not only the first integer but also the next few integers after it too. Exactly how many additional integers get moved depends on the cache's **block size**. For example, a cache with a 16 byte block size will store `array[i]` and the elements in `i+1`, `i+2`, `i+3`.

We can see the differences in spatial locality in the following example.

Example 8.2 ()

One may find that simply changing the order of loops can cause a significant speed up in your program. Consider the following code.

<pre> 1 float averageMat_v1(int **mat, int n) { 2 int i, j, total = 0; 3 4 for (i = 0; i < n; i++) { 5 for (j = 0; j < n; j++) { 6 // Note indexing: [i][j] 7 total += mat[i][j]; 8 } 9 } 10 return (float) total / (n * n); 11 }</pre>	<pre> 1 float averageMat_v2(int **mat, int n) { 2 int i, j, total = 0; 3 4 for (j = 0; j < n; j++) { 5 for (i = 0; i < n; i++) { 6 total += mat[i][j]; 7 } 8 } 9 return (float) total / (n * n); 10 } 11 .</pre>
--	--

Figure 42: Two implementations of taking the total sum of all elements in a matrix.

It turns out that the left hand side of the code executes about 5 times faster than the second version. Consider why. When we iterate through the *i* first and then the *j*, we access the values `array[i][j]` and then by spatial locality, the next few values in the array, which are `array[i][j+1]`, ... are stored in the cache.

1. In the left hand side of the code, these next stored values are exactly what is being accessed, and the CPU can access them in the cache rather than having to go into memory.
2. In the right hand side of the code, these next values are *not* being accessed since we want to access `array[i+1][j]`, Unfortunately, this is not stored in the cache and so for every n^2 loops we have to go back to the memory to retrieve it.

8.4 Caches

In theory, a cache should know which subsets of a program’s memory it should hold, when it should copy a subset of a program’s data from main memory to the cache (or vice versa), and how it can determine whether a program’s data is present in the cache. Let’s talk about the third point first. It all starts off with a CPU requesting some memory address, and we want to determine whether it is in the cache or not. To do this, we need to look a little deeper into memory addresses.

Definition 8.7 (Portions of Memory Addresses)

A memory address is a *m*-bit number.^a It is divided up into three portions.

1. The **tag** field with *t* bits at the beginning.
2. The **index** field with *i* bits in the middle.
3. The **offset** field with *o* bits at the end.

The tag plus the index together refers to the **block number**.

Tag	Index	Offset
1010	0000011	00100

Figure 43: Portions of a 16 bit memory address with $t = 4, i = 7, o = 5$.

^a64 in 64-bit machines.

Before we see why we do this, we should also define the portions of a CPU.

Definition 8.8 (CPU Cache)

A **CPU cache** divides its storage space as follows. A cache is essentially an array of sets, where S is the number of sets. Each set is divided into E units called **cache lines/rows**, with each cache line independent of all others and contains two important types of information.

1. The **cache block** stores a subset of program data from main memory, of size 2^o .^a Sometimes, the block is referred to as the cache line. Note that if the cache block size is 2^o bytes, then the block offset field has length $\log_2 2^o = o$.
2. The **metadata** stores the **valid bit** (which tells us if the actual data in memory is valid), and the **tag** of length t (the same as the tag length of the memory address) which tells us the memory address of the data in the cache.

Therefore, the **cache size** is defined to be $C = S \cdot E \cdot B$ (the metadata is not included).

Line	V	Tag	Cache Data Block
0			
1			
2			
→ 3	1	01110010...	
4			
5			
6			

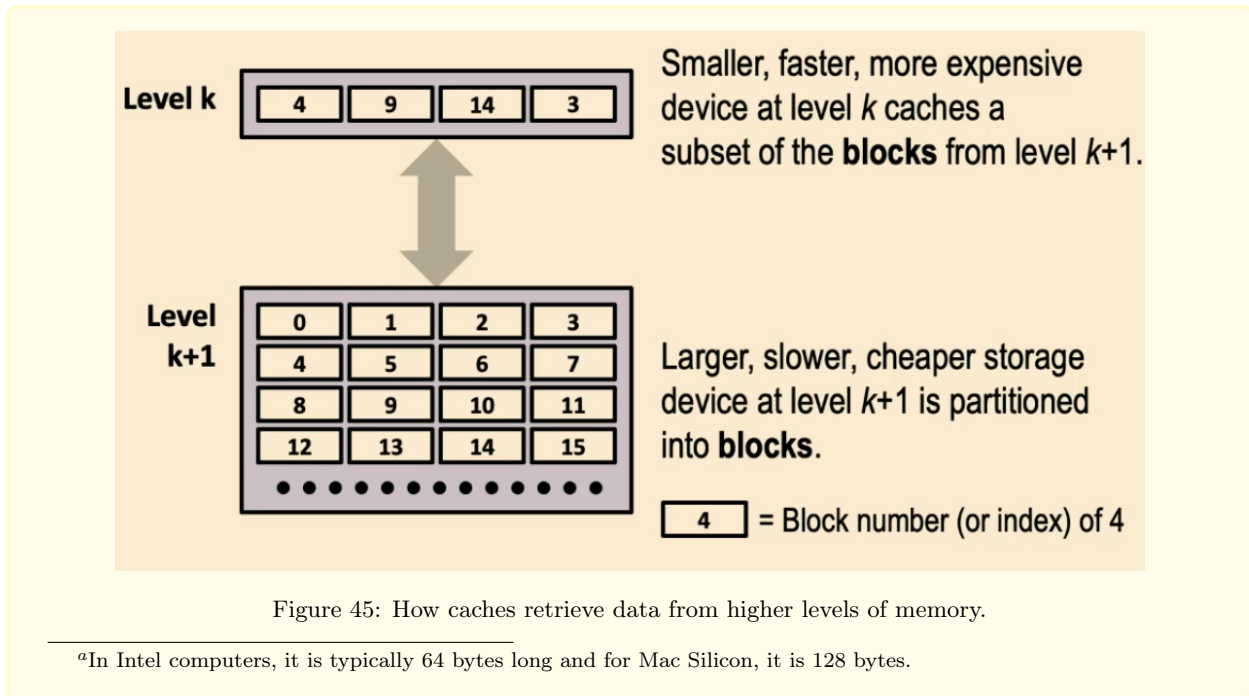
Figure 44: Diagram of a direct-mapped cache.

CPU caches are built-in fast memory (SRAM) that stores stuff. There are two types:

1. **i-cache** stores copies of instructions.
2. **d-cache** stores copies of data from commonly referenced locations.

We saw that caches come in different levels, they all just hold words retrieved from a higher level of memory.

1. CPU registers hold words retrieved from L1 cache.
2. L1 holds cache lines retrieved from L2 cache.
3. L2 cache holds cache lines retrieved from L3 cache or the main memory.
4. Main memory holds disk blocks retrieved from local disks.
5. Local disks hold blocks retrieved from remote disks or network servers.



Example 8.3 (Simple Calculations)

Given a direct-mapped cache specified by a block size of 8 bytes and a cache capacity of 4 KB,

1. the cache can hold 512 blocks.
2. the block offset field is $\log_2 8 = 3$ bits wide.
3. the address $0x1F = 0b00011111$ is in block number 3 since the last three bits are the offset, and whatever is left (passed through the hashamp, which is simply modulo), is the block number.

In **I/O caches**, software keeps copies of cached items in memory, indexed by name via a hash table.

At the lowest level, registers are explicitly program-controlled, but when accessing any sort of higher memory, the CPU doesn't know whether some data is in the cache, memory, or the disk.

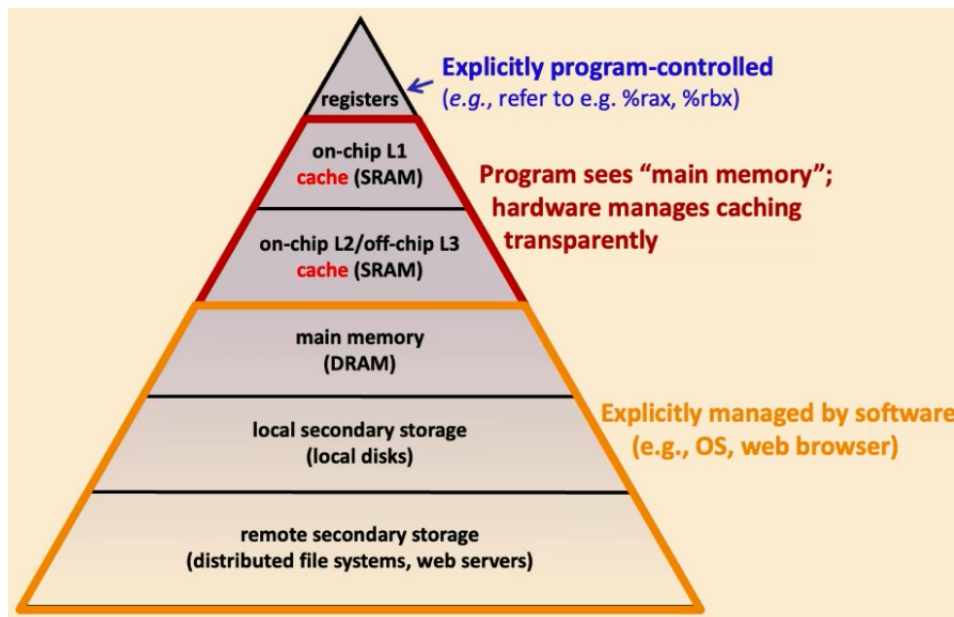


Figure 46

Finally, let's compare software vs hardware caches.

Definition 8.9 (Software Caches)

When implementing caches in software, there are large time differences (DRAM vs disk, local vs remote), and they can be tailored to specific uses cases. They also have flexible and sophisticated approaches with data structures (like trees) and can perform complex computation.

Theoretically, when implementing hash tables, you never actually have to evict something. You can have the values of the table to be a linked list where we add to the head. If there is unlimited chaining, we have a full associative cache, and if we have limited chaining (e.g. 5), it is like a 5-way set associative cache. If it goes out of bound, we can implement LRU by removing the tail of the linked list.

Definition 8.10 (Hardware Caches)

In hardware caches, there are smaller time differences, needs to be as fast as possible, and parallelization is emphasized.

There are slightly different implementations of caching, and for each implementation, we will describe

1. how to load data from memory into the cache,
2. how to retrieve data from the cache,
3. how to write data to the cache.

8.4.1 Direct Mapped Cache

A direct mapped cache is a caching implementation when we assume that $E = 1$, which means that for any given memory address, there is only one possible cache line that can store this data at that memory address. That is, the cache is really just a bunch of sets with one cache line each, and each cache line is completely isolated from the others. Whether we load data from memory into cache or try to retrieve data from the cache, it's really the same process.

Theorem 8.1 (Placement)

To load data from memory into the cache, which happens when there is a **cache miss**, we do the following.

1. The CPU requests a memory address $M = (T, I, O)$.
2. There exists a hashmap H that maps the index I to a cache line.
3. At line $H(I)$, we can get a cache miss and must load from memory into this cache.
4. We wait until the memory has retrieved the data from the portion of the memory. i.e. we wait for the 2^o bytes located at addresses $(T, I, 0 \dots 0)$ to $(T, I, 1 \dots 1)$. Call this data D .
5. The 2^o byte string D is stored in the cache data block at line $M(I)$, ready to be used.

Theorem 8.2 (Lookup)

To see whether a requested memory address is in the cache, we do the following.

1. The CPU requests a memory address $M = (T, I, O)$.
2. There exists a hashmap H that maps the index I to a cache line.
3. At line $H(I)$, check the cache line's valid bit. If it is not valid, then this is a cache miss and we must go to the memory to retrieve the data, leading to the above process.
4. Since there could be multiple I that maps to the same cache line, there will be overlap. But this is where the tag portion comes in. At cache line $H(I)$, the CPU checks the cache tag to see if it matches the memory tag T .
5. If it does, then we have just found a way to identify the first $t + i$ bits of the requested memory address, and we have gotten a cache hit. Now, we know that the cache's data block holds the data that the program is looking for. We use the low-order offset bits of the address to extract the program's desired data from the stored block.

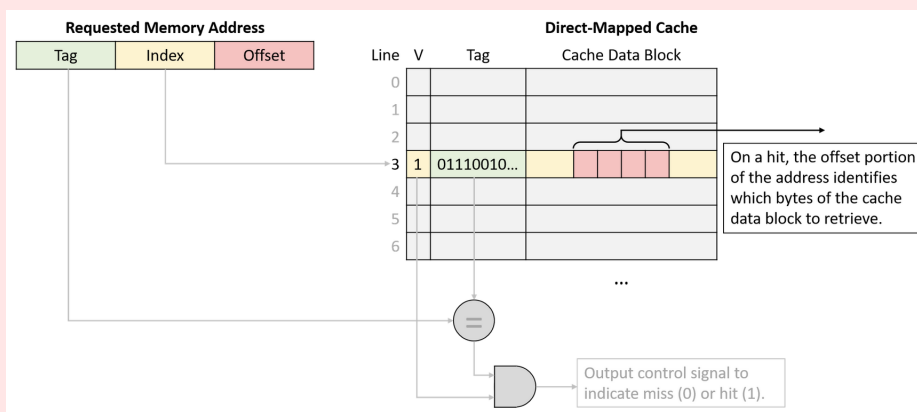


Figure 47: Diagram of a cache request. Note that since the entire data in the memory block stored in the cache, we can take advantage of spatial locality.

So far, we've talked about reading operations, but what about writing to the cache? It is generally implemented in two ways.

Definition 8.11 (Write-Through, Write-Back Cache)

Note that when we write data to cache, it does not need to be immediately written to memory, but rather it can be flushed to memory at a later time. This is efficient since if we have repeated operations on a single memory address, we don't have to go back and forth between the CPU and memory.

1. In a **write-through cache**, a memory write operation modifies the value in the cache and

simultaneously writes the value to the corresponding location in memory. It is always synchronized.

2. In a **write-back cache**, a memory write operation modifies the value stored in the cache's data block, but does *not* update main memory. Instead, the cache sets a **dirty bit** in the metadata to indicate that the cache block has been modified. The modified block is only written back to memory when the block is replaced in the cache.

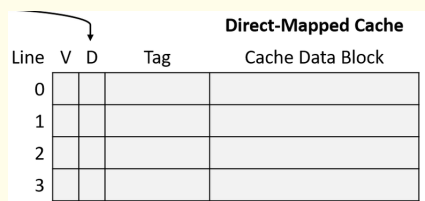


Figure 48: A dirty bit is a one bit flag that indicates whether the data stored in a cache line has been modified. When set, the data in the cache line is out of sync with main memory and must be written back (flushed) back to memory before eviction.

As usual, the difference between the designs reveals a trade-off. Write-through caches are less complex than write-back caches, and they avoid storing extra metadata in the form of a dirty bit for each line. On the other hand, write-back caches reduce the cost of repeated writes to the same location in memory.

Theorem 8.3 (Replacement)

Replacement occurs exactly the same way as if we just did a placement and is trivial. We retrieve the data block from the memory and store it in the cache. Direct-mapping conveniently determines which cache line to evict when loading new data. Given new memory $M = (T, I, O)$, you *must* evict the cache line at $H(I)$.

8.4.2 N way Set-Associative Cache

Note that for both examples, given a fixed hashmap H it is not possible to store data in two memory addresses M_1 and M_2 where both $H(I_1) = H(I_2)$. Therefore, the choice of hashing must be done so that it minimizes the number of collisions. So far, we have only considered memory read operations for which a CPU performs lookups on the cache. Caches must also allow programs to store values. However, there is a better way to do this: just construct it so that each set has more than one cache line, and so data in index portions of different memory addresses can be stored in different cache lines.

In here, we deal with $E \neq 1$, and so there are multiple sets each with multiple lines. This means that the cache is more like a 2D array, and when we want to retrieve an index, we must look through the $H(I)$ th line in *each* set to see if the tag matches.

Theorem 8.4 (Lookup)

To see whether a requested memory address is in the cache, we do the following.

1. The CPU requests a memory address $M = (T, I, O)$.
2. We iterate through each of the S sets in the cache, looking at cache line $M(I)$.
3. For each line, we check if it is valid and if so, whether the line tag matches the memory tag. If we get a hit, then we have found the data in the cache.

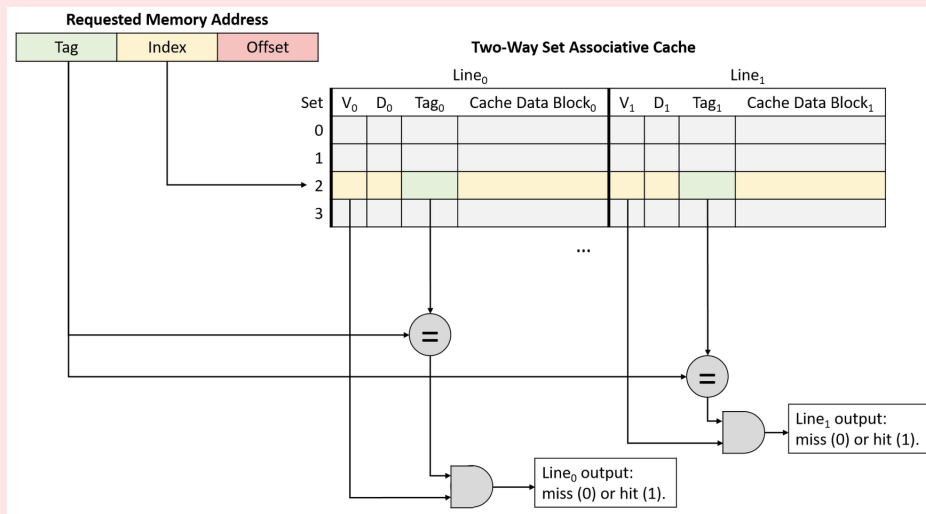


Figure 49: Diagram of a 2 set-associative cache.

If you have a **fully associative cache**, then you have one set with $E = C/B$ lines. Therefore, you can really put any memory address data in any cache line. There is a clear tradeoff here. As we increase N , we can get more flexibility in using all of our cache space, but the time complexity of retrieving and writing data scales linearly. In fact, this linear scan is too slow for a cache, which is why you need to implement some parallel tag search, but this turns out to be quite expensive to build.³

Though we have a more robust implementation with associative mapping, placement and replacement now face the problem of *which* set to place the data in or evict existing data.

Theorem 8.5 (Placement)

To load data from memory into the cache this is trivial since we can just go through the sets, find one where the valid bit is 0, and just place the data there.

In replacement, this is a bit trickier, but using the principle of temporal locality, we can try and replace the least recently used cache. This tries to minimize cache misses, but not slow down the lookup too much.

Theorem 8.6 (Replacement)

To replace data on the cache, we use the **least recently used (LRU)** algorithm. This matches temporal locality, but it also requires some additional state to be kept.

8.4.3 Types of Cache Misses

There are three types of cache misses.

Definition 8.12 (Cold (Compulsory) Miss)

A **cold miss** occurs when the cache is empty and the CPU requests a memory address. This is the first time the CPU is requesting this memory address, and so it must go to the memory to retrieve the data.

³You have to copy the request tag with a circuit and compare it to all the tags in the cache, which turns out to be a much larger circuit.

Definition 8.13 (Capacity Miss)

A **capacity miss** occurs when the cache is full and the CPU requests a memory address that is not in the cache. This is because the cache is full and so the CPU must evict some data to make space for the new data.

Definition 8.14 (Conflict Miss)

A **conflict miss** occurs from premature eviction of a warm block.

Valgrind's cachegrind mode.

9 Operating Systems

Up until now, we've seen the dynamics of how one program works in a computer system. The code, which first resides in the disk, is fetched (through blocks) into memory, and after compiling (precompiling, compiling, assembling, linking), we have a binary. The binary is then loaded into memory in the stack frame, and the CPU executes the instructions. The CPU also has a cache, which stores the most frequently accessed data during the process, taking advantage of locality for efficiency.

Our computer obviously does not just run one program. It runs several, and to run several, we need some control mechanism to manage how these programs interact with the CPU, memory, and disk. For example, one problem is that if we download application A and application B and run their binaries, how do we know whether they share memory addresses and consequently overwrite each other's data?⁴ The operating system takes care of these, which manages *processes* that each have their own *virtual memory space*.

Furthermore, consider some of the components of the computer: the RAM, disk, and IO devices like your keyboard and monitor. For security reasons, it is not wise to let the user applications (e.g. Chrome or Slack) control these devices completely. Their power must be restricted in some way.

1. When you have a Chrome window and resize it, Chrome should not be able to modify the pixels outside that window.
2. When you want to print some statement using `printf`,
3. When you're editing a code file with VSCode, you want to limit the application to save to certain parts of the disk.
4. When you are running Chrome and Slack together, you don't want them to read each other's data directly.

This is also for convenience. Say that if you are creating an application that has the option to save files to disk, you don't want to write the hardware backend to write to the disk. You want to just call a function that writes to the disk, and the OS will take care of the rest.

Definition 9.1 (Operating System)

A common confusion is that people think that the **operating system** describes the computer itself, but it is really just another piece of software. What makes this piece of software so special is that it manages every other software in the computer. It provides generally three services:

1. It **multiplexes** the hardware resources. Since there are many applications/programs with finite CPU resources (number of cores) and shared access to storage devices, the OS schedules some sharing mechanism for execution time on CPU cores and manages access to storage devices.
2. It **abstracts** the hardware platform. Since each CPU core simply executes a sequence of instructions, the OS introduces processes and thread abstractions. Furthermore, it introduces

⁴This is different from linking, where we have relocation tables to ensure that *object files* do not conflict with each other.

filesystems (file/directories) on top of raw storage devices.

3. It **protects** software principals from each other. Since many applications from various users are using the CPU, the OS provides isolation between them. It enforces user access permission (read/write) for files.

The OS is booted by the system firmware (BIOS or UEFI), which lives in ROM (sRAM and therefore non-volatile) and copies the OS from a fixed part of the disk, called the **bootloader**, into the RAM, which itself then loads the OS into memory. Once the OS starts running, it loads the rest of itself from disk, discovers and initializes hardware resources, and initializes its data structures and abstractions to make the system ready for users.

Definition 9.2 (Kernel)

The **kernel** is the actual binary that is loaded into RAM that runs the OS. The kernel code and data resides in a fixed and protected range of addresses, called the **kernel space**, and user programs cannot access kernel space.

9.1 Control Flow

When we worked with jumps (conditional and unconditional), calls, and returns in assembly, all of these operations were with respect to the **program state**, which is the isolated environment that the program is in. One program doesn't have any clue of what is going on anywhere else, such as other programs or input/output signals. This means that given what we have learned,

1. programs cannot to write files to the disk (since that is outside the program).
2. programs cannot be terminated by pressing CTRL + C on the keyboard.
3. programs cannot receive data that arrives from the disk.
4. programs cannot send data to the monitor to display.
5. programs cannot react accordingly when there is an instruction to divide by 0.⁵

To do all these things, we need to have access to the global **system state**, which the OS has access to.

It turns out that it is impossible for jumps and procedure calls to achieve this, and rather the system needs mechanisms for **exceptional control flow** (i.e. control flow that is not within the regular program state), or commonly referred to as **exceptions**. This requires the CPU to enter into a more powerful state than its current place in the program state, called the kernel state. The actual thing that triggers this is called an **interrupt**, which can come from both the hardware and software. In the kernel state, the CPU can access the hardware and perform operations that the program state cannot to handle these exceptions.

Example 9.1 (Interrupts)

Some examples of how the OS can be interrupted is:

1. when one's WiFi card detects a signal.
2. a hard disk drive may interrupt the OS if a read fails due to a bad sector.
3. an application may request a system call to open a file.
4. If you have 10 applications running on 1 CPU core, you may want the CPU core to run to the next application every 10 milliseconds. So, there may be a system call every 10 milliseconds in each program to the OS to switch to the next application.

⁵I guess you can use a conditional jump to check if the divisor is 0 and then jump to a different part of the code.

Definition 9.3 (Execution Modes)

The CPU helps with this by providing two execution modes, which is determined by a special bit in the CPU called the **mode bit**.

1. In **user mode**, the CPU executes only user-level instructions and accesses only the memory locations that the OS makes available to it. It also restricts which hardware components the CPU can directly access.
2. In **kernel mode**, the CPU executes any instructions and accesses any memory location (including those that store OS instructions and data). It can also directly access hardware components and execute special instructions.

Note that the execution mode is *property of the CPU!*

Example 9.2 (Monitor)

A monitor is really just some device that scans a certain portion of memory at a certain frequency that is higher than the human eye can detect. In user mode, if you try to access this memory buffer, you get an exception. No user mode can access this memory buffer.

Example 9.3 (Amazon.com)

When you are on Amazon to search up some product, you want to type in some keyword in the search bar. The web browser, say Chrome, that you are running it on, runs in user mode. When you type in the keyword, Chrome sends a system call to the OS, triggering the kernel mode which retrieves the keys that you pressed, and redirects it to Chrome. The same goes with the location of your mouse. When you move and click on a product, Chrome sends a system call to the OS, which then receives the mouse location and sends it back to Chrome. The application has no way to directly access the hardware.

Now specifically, how does one enter in this kernel mode? We've already hinted at it before, but to elaborate, there are 4 types of exceptions.

Definition 9.4 (Types of Exceptions/Interrupts)

As we have mentioned, we go into kernel mode through exceptional control flows. To go back from the kernel mode to the user mode after the exception handling is done, the kernel must explicitly give back the control to the user program, which is done with a special instruction, which changes the CPU to the user mode again. At this point, it can return back to user mode at the current instruction, next instruction, or abort it.

1. These control flows can either by **synchronous** (caused by an instruction) or **asynchronous** (caused by some other event external to the processor). Asynchronous interrupts are indicated by setting the processor's interrupt pins.
2. Furthermore, **intentional** exceptions transfer control to the OS to perform some function, and **unintentional** exceptions happen when there is a bug.

This gives us 4 categories of exceptions.

1. Intentional synchronous exceptions are **system calls**, aka **traps** (e.g. `printf`, `open`, `close`, `write`, breakpoint traps, special instructions). It returns control to the next instruction.
2. Unintentional synchronous exceptions are **faults** (possibly recoverable) or **aborts** (unrecoverable) (e.g. invalid or protected address or opcode, page fault, overflow, divide by zero). This automatically triggers the kernel mode which then uses an exception handler to kill the process.
3. Intentional asynchronous exceptions are **software interrupts**, which is when software requests an interrupt to be delivered at a later time (e.g. there's some task you want the kernel to do later).
4. Unintentional asynchronous exceptions are **hardware interrupts** caused by an external event

(e.g. IO such as CTRL + C, op completed, timers which may switch to another application every 10ms, power fail, keyboard, mouse click, disk, receiving a network packet). Unlike system calls, which come from executing program instructions, hardware interrupts are delivered to the CPU on an **interrupt bus**.

Once a system call or hardware interrupt is finished, the program continues to resume back in user mode.

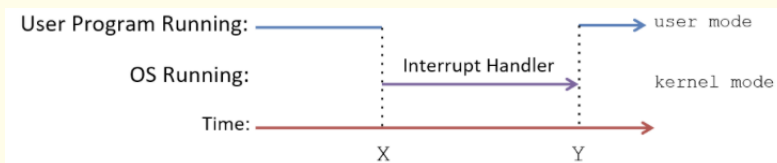


Figure 50: The CPU and interrupts. User code running on the CPU is interrupted (at time X on the time line), and OS interrupt handler code runs. After the OS is done handling the interrupt, user code execution is resumed (at time Y on the time line).

Now the question arises: how does the CPU know where to go when an system call or interrupt occurs? These are done through tables that map some unique ID number to some functionality. These tables are stored in a protected memory space reserved by the kernel.

Definition 9.5 (System Call Table)

This is done through the **system call table**, which is a table of addresses in memory that the CPU can jump to when a system call occurs. Each system call has a unique number k , and the handler function k is called each time system call k occurs.

Example 9.4 (Common System Calls)

Some common system calls, or **syscalls**, are shown below with their unique ID number (in Linux x64).

Table 5: System Call Functions

Number	Name	Description
0	read	Read file
1	write	Write file
2	open	Open file
3	close	Close file
4	stat	Get info about file
57	fork	Create process
59	execve	Execute a program
60	_exit	Terminate process
62	kill	Send signal to process

Example 9.5 (Syscalls of Open)

Look at the following objdump file below. The corresponding C code just calls `open(filename, options)` and the corresponding syscall ID is `0x2`. We are simply loading the syscall ID into the `%eax` register (only needs last 32 bits since the syscall IDs are quite small), which is then executed by the `syscall` instruction to go into the kernel mode.

```

1 0000000000e5d70 <__open>:
2 ...
3 e5d79: b8 02 00 00 00      mov    $0x2,%eax    # 2 is the open syscall number
4 e5d7e: 0f 05              syscall             # return value in %rax
5 e5d80: 48 3d 01 f0 ff ff  cmp    $0xfffffffff001,%rax
6 ...
    
```

A negative number in `%eax` gives an error corresponding to negative `errno`. It is also worth mentioning that `%eax` is used rather than `%rdi` or `%rsi` because we need these two parameter registers as arguments for the `open` function itself.

Note that whether we are in the program stack or the kernel stack, we always have stack pointers and other registers to navigate them. In fact, for every CPU core, it has its own set of registers and its own kernel stack.

Example 9.6 (Syscall of Read)

If we have read syscall, then

1. We use the syscall table to go to the trap handler for the read syscall.
2. The handler identifies the block and allocates a buffer.
3. Then it reads the block from the disk, which may take a while (in CPU time) since it is extremely slow for all IO tasks. The CPU, while waiting, can be put to sleep for other processes to run on the CPU. When the disk is done reading, it (the hardware) can send a hardware interrupt to the CPU, telling it that it is done.
4. Then it copies the block to the user buffer and returns from the syscall back into the user mode in the program state.

Definition 9.6 (Exception Table)

This is done through the **exception table**, which is a table of addresses in memory that the CPU can jump to when an exception occurs. Each type of event has a unique exception number k , and the handler function k is called each time exception k occurs.^a

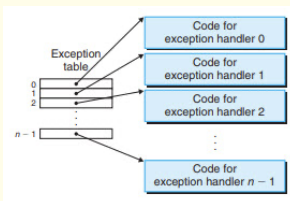


Figure 51: System call table is stored in a protected memory space reserved by the kernel.

^aThis is similar to a hardware implementation of a switch statement in C.

Example 9.7 (Common Exception Numbers)

Some common exception numbers are listed below.

Table 6: Exception Summary

Exception Number	Description	Exception Class
0	Divide Error	Fault
13	General protection fault	Fault
14	Page fault	Fault
18	Machine check	Abort
32-255	OS-defined	Interrupt or trap

From the application’s point of view, even if an interrupt happens, it just thinks it is running line by line.

Definition 9.7 (Process Address Space)

Interrupts can happen at any time, and one way to efficiently support this execution context switch from user mode to kernel mode is to do the following. At boot time, the OS loads its kernel code at a fixed location in RAM. Every time you create a new program state, the OS initializes a CPU register with the starting address of the OS handler function. On an interrupt, the CPU switches to kernel mode and executes OS interrupt handler code instructions that are accessible at the top addresses in every process’s address space. Because every process has the OS mapped to the same location at the top of its address space, the OS interrupt handler code is able to execute quickly in the context of any process that is running on the CPU when an interrupt occurs. This OS code can be accessed only in kernel mode, protecting the OS from user-mode accesses; during regular execution a process runs in user mode and cannot read or write to the OS addresses mapped into the top of its address space.^a

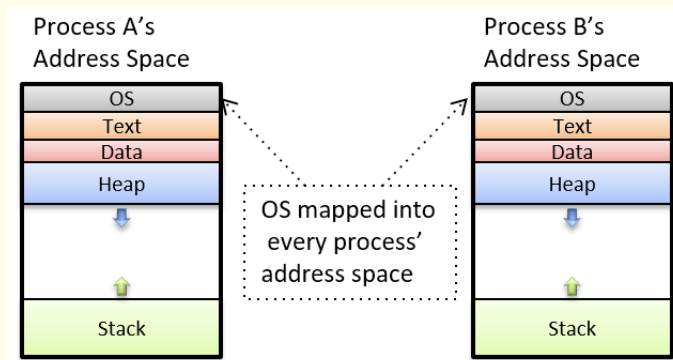


Figure 52: Process address space: the OS kernel is mapped into the top of every process’s address space.

^aHowever, due to security reasons where the user space can read kernel space data, this is obsolete.

In summary, a good visual is that each program runs as independent processes, with its own virtual address space (elaborated next) and the OS mediates access to shared resources.

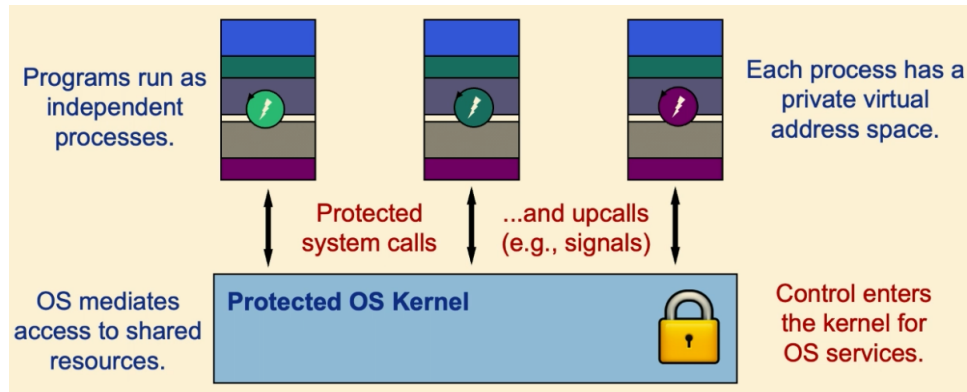


Figure 53: Multiple programs running and controlled by an operating system.

Each process can be in one of three states. It can either be currently running on the state, ready to run, or if there is a long IO operation, it can be blocked, which is then unblocked with a hardware interrupt. Usually anything that involves IO puts the state to blocked (e.g. reading data from disk, the keyboard, or the internet). The pool of processes that are concurrently running is the running and ready states.

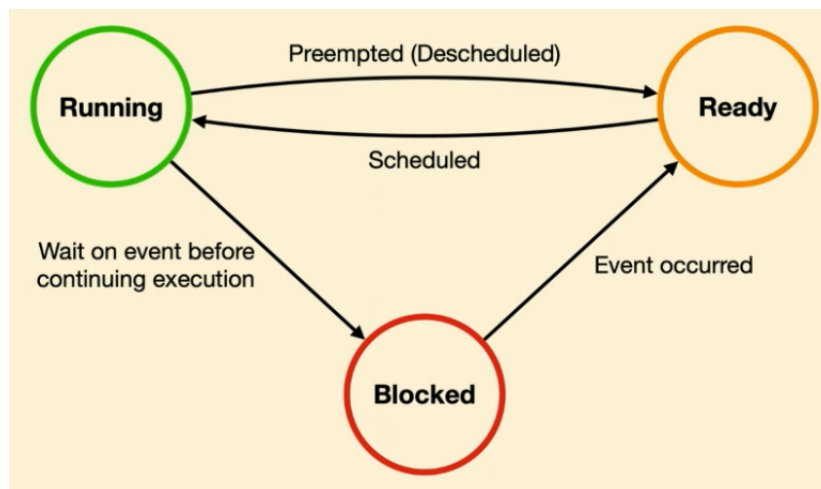


Figure 54: Three states that a single process can be in. The pool of processes that are concurrently running is the running and ready states. The blocked state is waiting to be put back into this pool by a hardware interrupt.

Example 9.8 (Running a Binary)

Therefore, to run a binary file `a.out`,

1. The kernel first loads the binary file from disk into RAM.
2. Then the OS kernel creates a new process with its own virtual memory stack and its global variables, etc.
3. Then the CPU's `%rip` register point to the address of the `main` function.
4. The kernel's virtual memory space is mapped to the top of the process's virtual memory space, where it is not visible to the user mode.

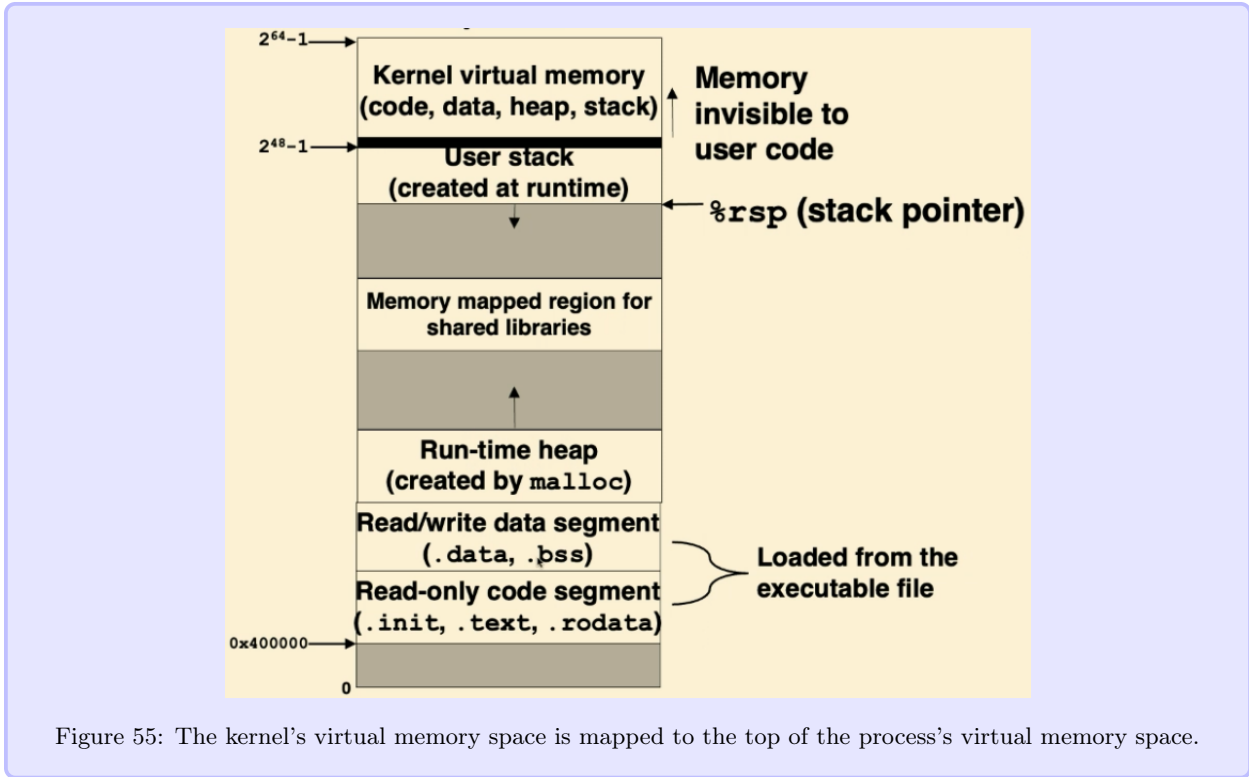


Figure 55: The kernel's virtual memory space is mapped to the top of the process's virtual memory space.

9.2 Virtual Memory

We have mentioned that there is a problem where two different application developers, who have linked their own C files to create binaries, can be installed on one computer and run at the same time. However, the linking has already been finished and the memory addresses of the symbols in each executable are fixed. This can be a problem if there are overlaps in the memory addresses.

Definition 9.8 (Virtual Memory)

The actual main memory of our system is referred to as the **physical memory**. To prevent such overlaps, the kernel and each user process has its own **virtual memory**. That is, there exists a **memory management unit (MMU)** in the CPU that translates virtual addresses to physical addresses through a hashmap.

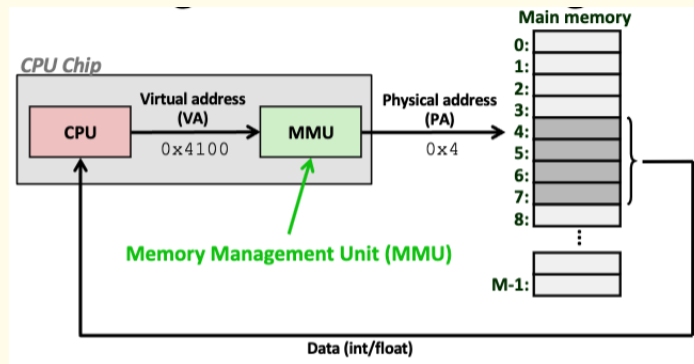


Figure 56: Memory management unit maps each virtual address to a physical address.

This allows the kernel to map the virtual memory of each process to the physical memory.

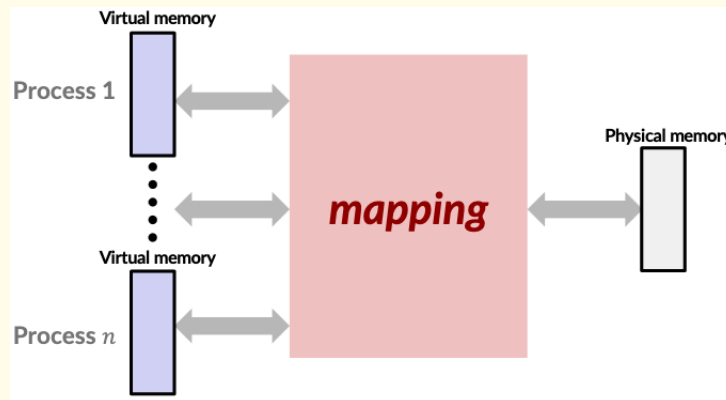


Figure 57: Each process has its own virtual memory space, which is mapped by the MMU to the physical memory space.

Example 9.9 (Virtual and Physical Memory Size)

Given a n -bit machine with 2^m -bytes of memory, $n > m$ and so there are more virtual addresses than physical addresses. If we have a 64-bit machine with 16GB of memory, then there are 2^{64} virtual addresses and $2^3 \cdot 2^{34} = 2^{37}$ bits of physical memory. If there are 8 processes running then there are $8 \cdot 2^{64} = 2^{67}$ bits of virtual memory.

There are many properties of virtual memory that solves a lot of problems and makes things more convenient. The main property is called **indirection** which means that the virtual memory is not the actual physical memory.

1. The first problem is that there are much more virtual addresses than physical addresses. Even storing a table for one process would take up more than all of your RAM. Therefore, for every byte in main memory, there exists one physical address (PA) and zero, one, or more virtual addresses (VA). We will elaborate on the specifics of this implementation later.
2. We also need to have memory management. Every process has its own stack, heap, `.text`, and `.data` sections. We must be able to allocate and deallocate memory and fit this accordingly.
3. We also need to have protection. We need to ensure that one process cannot read or write to another process's memory.
4. While we want isolation, we also want sharing between processes if needed (e.g. signing into Slack using Google on a browser). Furthermore, if there are multiple calls of the `printf` function, we can just have a single copy of the `printf` function in memory rather than having multiple copies for each process. This can be done through the concept of permissions.

Let's talk about how we should actually map these addresses. One property of this mapping is that we want contiguous addresses both in the virtual and the physical level so that we can store arrays, exploit locality, etc. Therefore, we can use larger blocks known as *pages*. Just like how we have divided memory addresses into sections that can be used to map to caches, we can divide the memory addresses into sections that can be used to map to the physical memory. Note that this also takes care of the first problem partially since now we can fit this table in the memory.

Definition 9.9 (Page)

Both in virtual and physical memory, an n -bit address can be divided into a **page number** and an **offset**. The page number is $n - 12$ and the offset is 12 bits. The page number is used to index into a **page table** that maps the page number to a physical address.



Figure 58: A page is a contiguous block of memory addresses.

While the entire page table is stored in memory (at memory stored by a protected CPU register), a portion of the page table is stored in the CPU cache.

1. The virtual page number (VPN) is equivalent to the block number.
2. The page offset is equivalent to the block offset.

Example 9.10 (Page Number)

In a 64-bit machine with 16GB of RAM, you have $2^{64}/2^{12} = 2^{52}$ virtual pages and $2^{37}/2^{12} = 2^{25}$ physical pages.

Therefore, our translation table is really a map from a virtual page number to a physical page number, rather than a virtual address to a physical address. This is created at runtime. Therefore,

1. The virtual page number VP is mapped through some map M to get the physical page number PP .
2. The virtual offset is the same as the physical offset.

Definition 9.10 (Page Table)

The **page table** is a hashmap that maps the virtual page number to the physical page number defined as the mapping

$$H : \underbrace{(VP, m)}_{64 \text{ bits}} \longrightarrow PP \tag{11}$$

Each input-output pair is called a **page table entry (PTE)**, and virtual memory is **fully associative**, meaning that any virtual page can be placed in any physical page, though it requires a large mapping function (the PT), which is different from CPU caches.

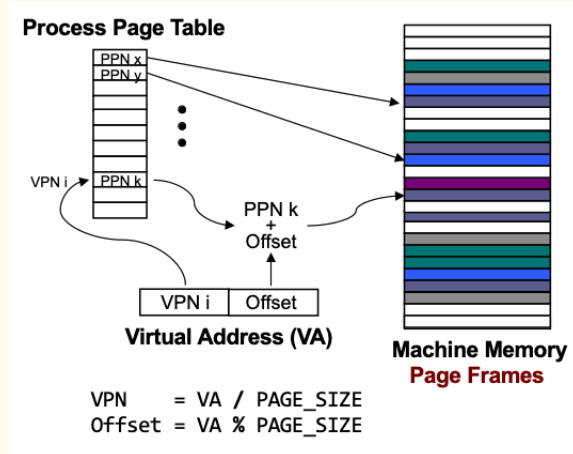


Figure 59: The page table only needs the virtual page number plus the metadata to map to the physical page number. The offset is provided by the virtual memory address itself.

Note that while we want to store the 52-bit VP in the page table, the actual input is still 64-bits, with 12 bits of metadata m . This metadata contains some information about the following

1. A bit that indicates whether the page is a read, write, or executable piece of code (3 bits).
2. A bit that indicates whether the page is valid or not.

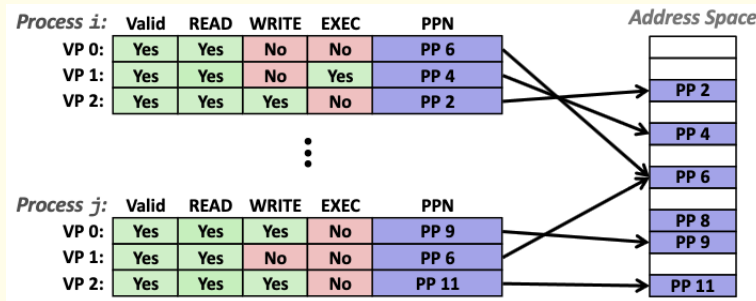


Figure 60: The page table entry contains the physical page number and some metadata.

Therefore, if you malloc, you are really just allocating some virtual memory addresses, which then get mapped to physical memory addresses in one or more pages.

Definition 9.11 (Page Fault)

It is clear that not every virtual page number can be mapped to a physical page number. If it turns out that a **page fault** happens if

1. the virtual page number maps to no physical page (i.e. is not in the page table) in the RAM
2. if some user program tries to access a physical page owned by the kernel
3. if the page number maps to some place in the disk (but it is not in physical RAM)

Page faults can be used in a lot of creative ways, but to reduce the risk of a page fault, e.g. when running out of physical memory, we can move some physical pages into disk and allocate memory by creating a new entry in our page table that maps this application’s virtual page into the now empty physical page.

Note that by this construction, instructions that are contiguous in virtual memory may not be contiguous in physical memory. This may seem like it defeats the purpose of locality, but for most purposes, the 4KB page size will be enough to exploit it. We also see that malloced addresses in the heap (while we have learned

that they were higher on the stack on higher addresses), are not necessarily in higher addresses in physical memory. Therefore, physical memory is scattered, and this is good since you don't need a giant contiguous block of memory to run large programs; you can divide it up into multiple physical pages.

Definition 9.12 (Swap Space)

Sometimes, the memory might not be in physical memory. Since memory is constrained (e.g. only 16GB), if we initialize a large array in the stack or global data, we may run out of memory. Therefore, the OS can flush out some physical pages in memory to disk, which is called **swapping**. The portion of the disk space that can be used in swapping is called the **swap space**.

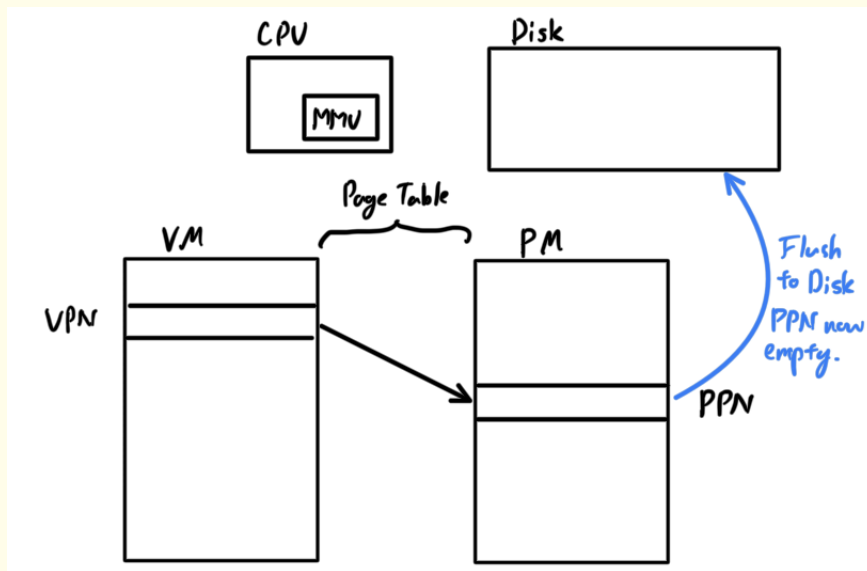


Figure 61: Swapping out physical pages to disk.

This allows us to abstract software into having almost infinite memory. Another important property is that swapping is **write-back** rather than write-through. We really don't want to write to disk every time we modify memory, so some thing may never end up on the disk (e.g. stack for short-lived processes). This is why when we open a file in C or Python, you may have to call `close()` since that will flush the memory to disk.

Example 9.11 (Page Fault)

When we swap out a physical page to disk, the physical page is now empty and accessing the virtual memory at this page table will cause a page fault. Say, when we want to write to a memory address that is swapped into the disk. The following will happen.

1. You execute code normally in user mode.
2. Then you try to write to a memory address that is swapped out, say through a `mov` operation. Say it is the following assembly code.

```

1  80483b7: c7 05 10 9d 04 08 0d          movl $0x0,0x8049d10
    
```

This raises a page fault, an exception, and so the OS goes into kernel mode.

3. The kernel then finds the location of this physical page in the disk. The implementation is OS-specific (e.g. you can store some metadata).
4. Then it must copy the page back from disk into memory, and it may also have to swap out

some other physical page to disk to make space if needed.

5. Then the OS goes back into user mode, which now has access to the relevant memory in disk. Ultimately, the moving operation is called twice. The first time it fails in user mode, and the second time (after the kernel mode, but now back to user mode) it succeeds. Note that this is different from a system call, which returns back to the *next* instruction. This call returns to the current instruction.

Definition 9.13 (Page Sharing)

This also makes protection and sharing to be quite nice. Given two virtual pages VP_1 and VP_2 , owned by two different processes, we can have them share information by mapping to the same physical page PP .

Section	Read	Write	Execute
Stack	1	1	0
Heap	1	1	0
Static Data	1	1	0
Literals/const	1	0	0
Instructions	1	0	1

Table 7: Permissions for different sections of virtual memory.

Example 9.12 (Page Sharing Between Two Applications)

Furthermore, we can set process 1 to have only read permissions and process 2 to have read/write permissions. Therefore, say Google Chrome (process 2) can write your password into some memory, and then Slack (process 1) can read it, copy it into the CPU, and do stuff with it.

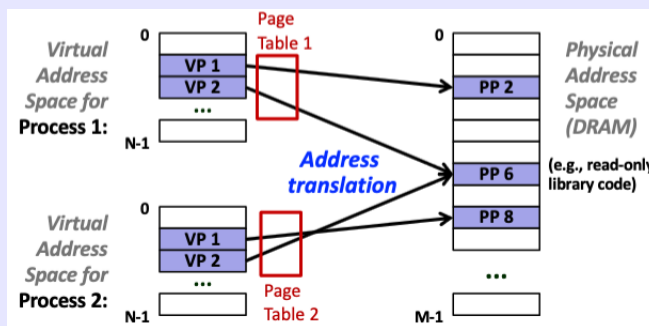


Figure 62: Sharing of data between two processes.

Now that we see how memory is swapped in the backend, we can see why larger memory can sometimes mean faster programs and why thrashing occurs.

Definition 9.14 (Thrashing)

The set of virtual pages that a program is “actively” accessing at any point in time is called its **working set**.

1. If the working set of one process is less than physical memory, then there is good performance for one process.
2. If the working set of all processes is greater than physical memory, then we have **thrashing**, which is a performance meltdown where pages are swapped between memory and disk continuously, and the CPU is always waiting or paging.

Example 9.13 (Computation Exercise)

Suppose that you have 16 KiB pages, 48-bit virtual addresses, and 16 GiB physical memory. How many bits wide are the following fields?

1. Virtual page number : $48 - 14 = 34$ bits.
2. Virtual page offset : 16 KiB is 2^{14} bytes, so we need 14 bits.
3. Physical page number : 16 GiB is 2^{34} bytes, so we need $34 - 14 = 20$ bits.
4. Physical page offset : 16 KiB is 2^{14} bytes, so we need 14 bits.

Furthermore, we have

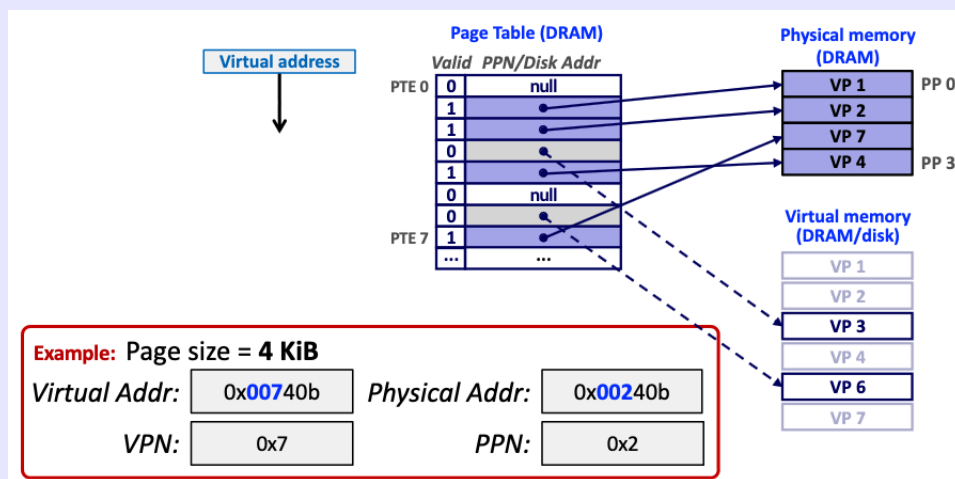


Figure 63: Given the virtual address, we can figure out the physical address, VPN, and PPN easily.

10 Shared Memory and Concurrency

So far, we’ve talked about everything as a sequential process of instructions. In practicality, we have improved this from the memory perspective by implementing caches, virtual memory, and swapping, but in the CPU perspective. In CPUs, we can’t just simply increase the clock frequency indefinitely since there are physical limitations.⁶ The current trend is to increase parallelism to compute faster, which is implemented with cores and threads. Let’s clear some of these definitions up.

Definition 10.1 (Processors, Cores)

In almost every consumer computer, there exists one **processor** (CPU) in it. The CPU can have multiple **cores**. Each core has its own set of registers, L1/L2 cache, and possibly even a shared L3 cache.

Now these cores must run a certain program. Let’s define what this means exactly.

Definition 10.2 (Program, Process)

A **program** can be thought of as a binary executable produced after linking. A **process** is a running instance of some program.

1. It is identified by a **process ID (PID)** number.
2. It is run on a CPU core, with its own registers.
3. It has its own virtual address space, containing the code (instructions), heap, and pagetable

⁶It turns out that power consumption increases faster than clock frequency, so it scales badly.

that maps it to physical memory.

Example 10.1 (Where to look for PIDs)

We can see the PIDs either by using `htop` (for UNIX systems) or by looking at the `/proc` directory in Linux systems. Each directory name represents the PID of the process.

```

1  ubuntu@passionate-blesbok:/proc$ ls
2  1      118   1368  26   44   590  762      cpuinfo    modules
3  10     119   1369  27   448  599  763      crypto     mounts
4  101    12    1370  28   45   6    764      devices    net
5  1012   120   1371  29   46   605  765      diskstats  pagetypeinfo
6  102    1232  138   3    467  608  769      driver     partitions
7  103    1234  139   30   468  611  796      execdomains pressure
8  1031   1261  14    309  47   612  8       fb         sched_debug
9  1037   129   15    31   471  613  801      filesystems schedstat
10 1038   13    16    32   473  629  810      fs         scsi
11 104    132   17    33   474  638  822      interrupts self
12 1043   1342  18    34   475  651  836      iomem     slabinfo
13 105    1353  180   35   48   666  850      ioports   softirqs
14 106    1354  19    356  49   696  872      irq       stat
    
```

(a) You can see the PIDs of the process by looking at the `/proc` directory. This changes quite often as processes are destroyed and created often, so to maybe track this in real time you might want to run `watch -n 0.1 'ls'`.

```

1  [ 0.0%] 4  [ 0.0%]
2  [ 0.0%] 5  [ 0.6%]
3  [ 0.0%] 6  [ 0.0%]
Mem[||||] 217M/11.7G Tasks: 45, 47 thr; 1 running
Swp[ ] 0K/0K Load average: 0.00 0.00 0.00
Uptime: 00:02:35

PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+ Command
736 avahi 20 0 7824 3184 2776 S 0.0 0.0 0:00.02 avahi-daemon: running [passionate-blesbok.local]
769 avahi 20 0 6892 328 0 S 0.0 0.0 0:00.00 avahi-daemon: chroot helper
896 daemon 20 0 3596 1896 1732 S 0.0 0.0 0:00.00 /usr/sbin/atd -f
905 kernoops 20 0 11460 1912 1464 S 0.0 0.0 0:00.00 /usr/sbin/kerneloops --test
913 kernoops 20 0 11460 1896 1448 S 0.0 0.0 0:00.00 /usr/sbin/kerneloops
738 messagebu 20 0 9816 4696 3464 S 0.0 0.0 0:00.15 /usr/bin/dbus-daemon --system --address=systemd: --nofork --nopidfile --systemd-ac
1 root 20 0 164M 10316 7184 S 0.0 0.1 0:00.47 /sbin/init splash
434 root 19 -1 60956 29504 28372 S 0.0 0.2 0:00.12 /lib/systemd/systemd-journald
471 root 20 0 21764 6884 3472 S 0.0 0.1 0:00.22 /lib/systemd/systemd-udev
591 root RT 0 273M 16312 6312 S 0.0 0.1 0:00.00 /sbin/multipathd -d -s
592 root RT 0 273M 16312 6312 S 0.0 0.1 0:00.00 /sbin/multipathd -d -s
593 root RT 0 273M 16312 6312 S 0.0 0.1 0:00.00 /sbin/multipathd -d -s
594 root RT 0 273M 16312 6312 S 0.0 0.1 0:00.01 /sbin/multipathd -d -s
595 root RT 0 273M 16312 6312 S 0.0 0.1 0:00.00 /sbin/multipathd -d -s
596 root RT 0 273M 16312 6312 S 0.0 0.1 0:00.00 /sbin/multipathd -d -s
590 root RT 0 273M 16312 6312 S 0.0 0.1 0:00.04 /sbin/multipathd -d -s
742 root 20 0 232M 6796 5884 S 0.0 0.1 0:00.00 /usr/lib/accounts-service/accounts-daemon
798 root 20 0 232M 6796 5884 S 0.0 0.1 0:00.00 /usr/lib/accounts-service/accounts-daemon
734 root 20 0 232M 6796 5884 S 0.0 0.1 0:00.03 /usr/lib/accounts-service/accounts-daemon
735 root 20 0 5144 456 392 S 0.0 0.0 0:00.00 /usr/sbin/anaconda -d -q -s
802 root 20 0 256M 16844 14432 S 0.0 0.1 0:00.00 /usr/sbin/NetworkManager --no-daemon
804 root 20 0 256M 16844 14432 S 0.0 0.1 0:00.00 /usr/sbin/NetworkManager --no-daemon
739 root 20 0 256M 16844 14432 S 0.0 0.1 0:00.06 /usr/sbin/NetworkManager --no-daemon
760 root 20 0 80948 1584 1380 S 0.0 0.0 0:00.00 /usr/sbin/irqbalance --foreground
744 root 20 0 80948 1584 1380 S 0.0 0.0 0:00.01 /usr/sbin/irqbalance --foreground
746 root 20 0 35696 18904 10312 S 0.0 0.2 0:00.06 /usr/bin/python3 /usr/bin/networkd-dispatcher --run-startup-triggers
754 root 20 0 230M 8628 5988 S 0.0 0.1 0:00.00 /usr/lib/policykit-1/polkitd --no-debug
797 root 20 0 230M 8628 5988 S 0.0 0.1 0:00.02 /usr/lib/policykit-1/polkitd --no-debug
747 root 20 0 230M 8628 5988 S 0.0 0.1 0:00.05 /usr/lib/policykit-1/polkitd --no-debug
918 root 20 0 1583M 28696 17848 S 0.0 0.2 0:00.04 /usr/lib/swapd/swapd
919 root 20 0 1583M 28696 17848 S 0.0 0.2 0:00.01 /usr/lib/swapd/swapd
920 root 20 0 1583M 28696 17848 S 0.0 0.2 0:00.00 /usr/lib/swapd/swapd
921 root 20 0 1583M 28696 17848 S 0.0 0.2 0:00.00 /usr/lib/swapd/swapd
922 root 20 0 1583M 28696 17848 S 0.0 0.2 0:00.00 /usr/lib/swapd/swapd
937 root 20 0 1583M 28696 17848 S 0.0 0.2 0:00.02 /usr/lib/swapd/swapd
944 root 20 0 1583M 28696 17848 S 0.0 0.2 0:00.01 /usr/lib/swapd/swapd
945 root 20 0 1583M 28696 17848 S 0.0 0.2 0:00.02 /usr/lib/swapd/swapd
    
```

(b) You can see the PID number of each process (binary) running on the left column when running `htop` on UNIX systems.

Figure 64: Two different ways to see the PIDs of all current processes.

There is a specific numbering to each process.

1. The process with PID 1 is always the kernel process.
2. The smaller PIDs (perhaps less than 300) are also reserved for the kernel, so don't kill it.

If you go into each process, you can see a few things.

```

1  ubuntu@passionate-blesbok:/proc/750$ sudo ls
2  attr      comm      fd        map_files net        pagemap    sessionid  statm
   uid_map
3  autogroup  coredump_filter fdinfo    maps      ns        personality setgroups
   status    wchan
4  auxv      cpuset    gid_map  mem       numa_maps projid_map smaps      syscall
5  cgroup    cwd       io       mountinfo oom_adj   root       smaps_rollup
   task
6  clear_refs environ    limits   mounts    oom_score sched      stack
   timers
7  cmdline   exe       loginuid mountstats oom_score_adj schedstat  stat
   timerslack_ns
    
```

Figure 65: There are many files in each PID folder that tells you about the process.

1. To get information about the status of this process, you can `cat status`.
2. The virtual address space is stored in `pagemap`. If you're on an 64-bit machine, this file will be extremely big, so just `cat pagemap` won't work. Therefore you should try `cat maps`, which shows you something like the following.

```

1  ubuntu@passionate-blesbok:/proc/750$ sudo cat maps
2  aaaac673c000-aaaac67e3000 r-xp 00000000 08:01 2576
   /usr/sbin/rsyslogd
3  aaaac67f3000-aaaac67f6000 r--p 000a7000 08:01 2576
   /usr/sbin/rsyslogd
4  aaaac67f6000-aaaac67fd000 rw-p 000aa000 08:01 2576
   /usr/sbin/rsyslogd
5  aaaac67fd000-aaaac67fe000 rw-p 00000000 00:00 0
6  aaaadfe94000-aaaadfed7000 rw-p 00000000 00:00 0
   [heap]
    
```

In here, you can see that the lefthand column represents the range of virtual memory address. The next column gives us the permissions (read, write, executable, shared/private).

10.1 Process Level Concurrency

Definition 10.3 (Context Switch)

Let us first start off with a single core system. At this point, everything is sequential, and to run all these processes at once^a we want to use system calls to transition between these processes. This is called a **context switch**.

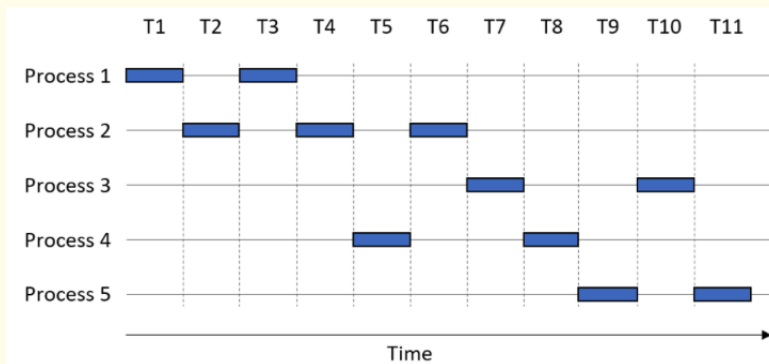


Figure 66: 5 processes may be executed as such on a single core.

Note that due to context switches, the **CPU time**, which is the time it takes to run a process on a CPU, is much shorter than the **wall-clock time**, which is the time a human perceives a process takes to complete.

^anot programs, since there can be multiple instances of one program, like two Chrome instances. In fact, Chrome produces multiple processes to help run each part of the browser, so one program may translate to multiple processes.

Note that context switches are expensive. To do one, you must essentially replace two things.

1. First, you need to clear out all the register values. This can be done by storing them in the current stack at the VAS, which then gets mapped through the page table into the physical address space.
2. Now the register values (like the instruction and stack pointers) are stored safely in the stack in the VAS, the actual page table must be swapped out too since each process must have its own virtual address space.

Since it is quite expensive to context switch all the time, the simplest thing to do is add more cores, which gives us the double benefit of distributing the process workload *and* having to do less context switches. This is called **physical concurrency**, and given the same workload, it speeds up our computation absolutely. However, this can physically take us so far due to the limited number of cores, and we must go further and use **logical concurrency**.

10.2 Thread Level Concurrency

It turns out that it is much more expensive to reload the page table of a new process rather than clearing out the register values. So, perhaps maybe we can try to implement multiple related “processes” that *share* the same VAS, but have their own execution stream (i.e. own stack and registers). This is precisely the concept of a *thread*.

Definition 10.4 (Threads)

Threads are multiple execution streams within a single process. To summarize them, a thread is an execution context within a process that has a...

1. thread ID
2. its own stack frame
3. its own register context^a

This is all that is really needed to execute some computation. Now, given that there are some number of threads in process *K*, they *share* the same virtual address space (VAS), are all under the same PID, share the same code, static data, heap, and file table. The individual stacks living within the VAS are protected from each other to avoid stack overflow.

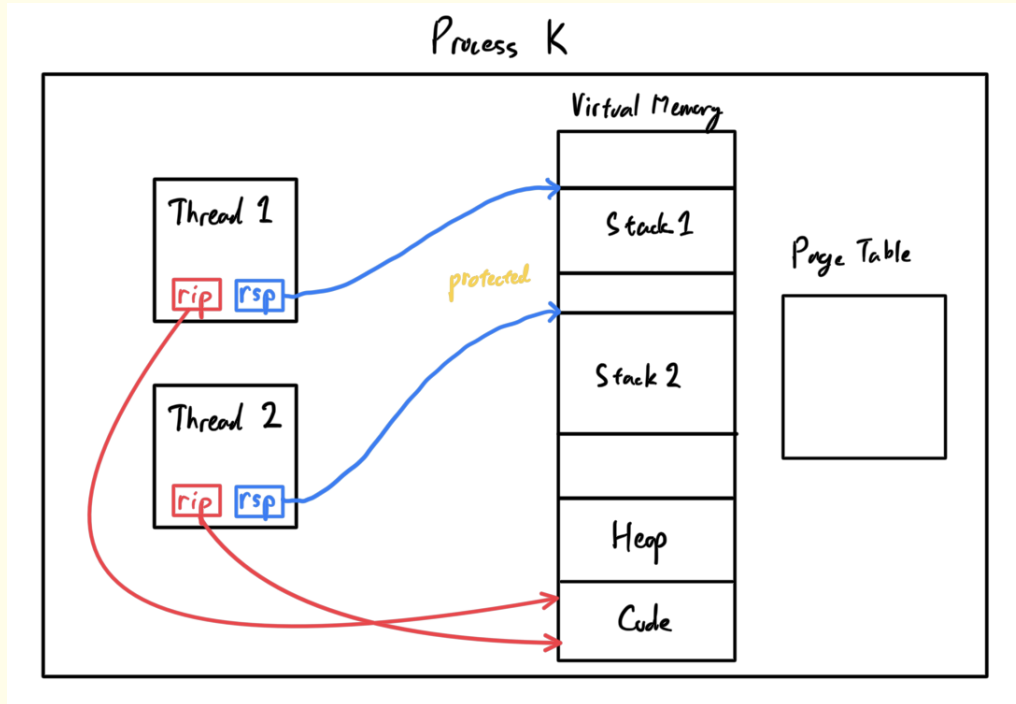


Figure 67: When there are two threads of a single process, the threads share the same virtual memory space. However, they each have their own set of registers. For example, they each have their own instruction pointer that points to the next line of code, along with their own stack pointer. Furthermore, to prevent stack overflow, there are protection mechanisms that prevent one stack from growing past a certain limit into another stack owned by a different thread.

Therefore, we can speed up our program in two ways.

1. If we have one core, we can do context switching faster between each thread (since we only have to load the register values).
2. If we have multiple cores, we can take thread 1 and have it run on one core while taking thread 2 and running it on another core. This is really analogous to having two separate processes on two cores, but these two processes simply share the same VAS, with the same code, data, and heap.

^aThis does not mean that each process has its own physical registers. It has its own *value* that is loaded into the registers. The physical number of registers is determined by the number of CPU cores.

Threads are advantageous for multiple reasons. First, by utilizing multiple cores we can speed up our program to reduce our *CPU time*. However, if we are sharing threads between one core, we're not actually speeding up anything at all but rather reducing our *wall-clock time*. The main speedup that we will feel is that latency heavy tasks will get offloaded to other threads, while more relevant programs can be run on the main thread. This is explained more in the following example.

Example 10.2 (Mobile Application)

If we have a single threaded messaging mobile app, then this is painfully slow since if we want to scroll down our messages while also sending and receiving messages, then we would have to wait for the message to receive from the server, into our disk, and into our memory, before the app responds when scrolling.

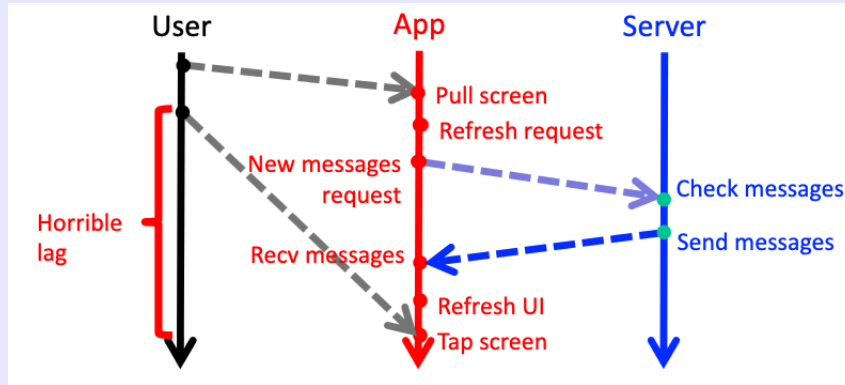


Figure 68: Single threaded app.

However, if we have a multithreaded app with one thread for the app UI and the other one for the server through a background thread, then we can have good UI response time.

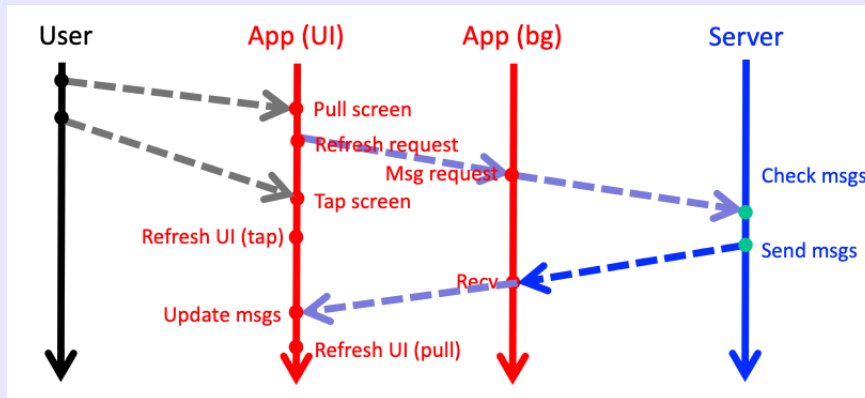


Figure 69: Multithreaded app. Methods on the UI thread must be fast to ensure user satisfaction while anything slow can run on a background thread.

The following law gives us a certain bound on how much parallelization can help us. Note that this does not talk about the responsiveness of an application due to clever thread sharing. It just says given a certain amount of computational task, how much can we reduce the CPU time with parallelization?

Theorem 10.1 (Amdahl's Law)

Say that we have code that runs in 1 second. Given that proportion f of our code can be parallelized, and the speedup for that portion is N , then the new time that our program will take is

$$T_{\text{new}} = (1 - f) + f/N \tag{12}$$

since the sequential part $1 - f$ cannot be sped up, and the remaining parallel part f can be sped up by distributing over N cores. Therefore, defining the speedup as $T_{\text{new}}/T_{\text{old}}$, we get our total parallelized speedup is

$$\frac{1}{(1 - f) + f/N} \tag{13}$$

Note that it is bounded by the sequential portion as $N \rightarrow \infty$.

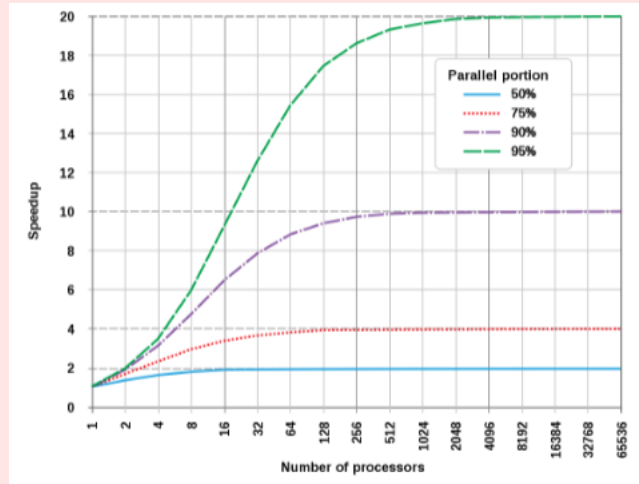


Figure 70: Amdahl’s law for different f ’s over multiple number of processors N .

This is implemented in C with the `pthread.h` library, which is included in the standard library directory and follows the POSIX (Portable Operating System Interface) standard. Essentially, we want to do the following:

1. Define a function that will be called for each thread. It must return a void pointer `void *` and its arguments must also be a void pointer `void *`. Think of this as our new main function for each stack that will be created from each thread.⁷ Since we are only restricted to a function taking in one void pointer argument, it is common to define a new struct like `arg_t` that contains all the parameters you need to run each thread. The void pointer can be typecast into the struct pointer at the beginning of each thread function.
2. We create `pthread_t` objects, which are the thread objects.
3. We call the `pthread_create` function that takes in the pointer of the thread object, some settings, the function to be called, and its arguments. At this point, the operating system will determine how these threads will be run, so you can’t make any sequential assumptions about them.
4. Then we join them using `pthread_join`, which basically waits until all the threads are complete before `main` continues.

We will show two examples that go over this process. But more importantly, the concurrency of these two examples will show unpredictable behavior.

Example 10.3 (Simply Print out Thread Number)

We can make threads to print out the number. But these aren’t really in the same order.

⁷Called a function pointer?

<pre> 1 #include <stdio.h> 2 #include <pthread.h> 3 #include <stdlib.h> 4 5 void* thread(void* args) { 6 printf("Thread %d\n", *(int*)args); 7 return NULL; 8 } 9 10 int main(int argc, char *argv[]) { 11 int size = 10; 12 pthread_t threads[size]; 13 int rc, i; 14 15 // thread creation 16 for (i = 0; i < size; i++) { 17 rc = pthread_create(&threads[i], NULL, thread, &i); 18 } 19 20 // join waits for the threads to finish 21 for (i = 0; i < size; i++) { 22 rc = pthread_join(threads[i], NULL); 23 } 24 return 0; 25 } </pre>	<pre> 1 Thread 1 2 Thread 2 3 Thread 6 4 Thread 3 5 Thread 8 6 Thread 5 7 Thread 2 8 Thread 4 9 Thread 3 10 Thread 3 11 . 12 . 13 . 14 . 15 . 16 . 17 . 18 . 19 . 20 . 21 . 22 . 23 . 24 . 25 . </pre>
---	---

Figure 71: Threads output shows that the order in which the functions are called cannot be predicted.

10.3 Atomicity Violation Bugs and Mutex Locks

We see that there are some parts of the code that are not meant to be parallelized.

Definition 10.5 (Atomicity-Violation Bugs)

This bug happens when the desired **atomicity** (indivisibility) among multiple memory accesses is violated.

Example 10.4 (Atomicity-Violation in SQL)

For example, if we have two threads doing the following (in MySQL):

```

1  Thread1::
2  if (thd->proc_info) {
3      ...
4      fputs(thd->proc_info);
5      ...
6  }
7
8  Thread2::
9  thd->proc_info = NULL;

```

If we pass the if statement but within it, `thd->proc_info` becomes NULL, then this would be very bad. Therefore, we should put locks around.


```
1 pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
2
3 Thread1::
4 pthread_mutex_lock(&lock);
5 if (thd->proc_info) {
6     ...
7     fputs(thd->proc_info);
8     ...
9 }
10 pthread_mutex_unlock(&lock);
11
12 Thread2::
13 pthread_mutex_lock(&lock);
14 thd->proc_info = NULL;
15 pthread_mutex_unlock(&lock);
```

Example 10.5 (Incrementing Shared Counter between Two Threads)

The volatile keyword for counter means that it can be changed by all threads.

<pre> 1 #include <stdio.h> 2 #include <pthread.h> 3 #include <stdlib.h> 4 5 static volatile int counter = 0; 6 7 void* thread(void* args) { 8 printf("%s : Start \n", (char*)args); 9 10 for (int i = 0; i < 10 * 1000 * 1000; i++) { 11 counter += 1; 12 } 13 14 printf("%s : End \n", (char*)args); 15 return NULL; 16 } 17 18 int main(int argc, char *argv[]) { 19 pthread_t thread1, thread2; 20 21 int rc; 22 rc = pthread_create(&thread1, NULL, thread, "A"); 23 rc = pthread_create(&thread2, NULL, thread, "B"); 24 25 rc = pthread_join(thread1, NULL); 26 rc = pthread_join(thread2, NULL); 27 28 printf("Counter : %d\n", counter); 29 return 0; 30 } </pre>	<pre> 1 A : Start 2 B : Start 3 A : End 4 B : End 5 Counter : 10229646 6 . 7 . 8 A : Start 9 B : Start 10 B : End 11 A : End 12 Counter : 14965289 13 . 14 . 15 A : Start 16 B : Start 17 A : End 18 B : End 19 Counter : 10086690 20 . 21 . 22 . 23 . 24 . 25 . 26 . 27 . 28 . 29 . 30 . </pre>
--	---

Figure 72: From looking at the behavior, we can see that the start and end times of thread A and thread B is complete unpredictable. It is only within the sequential nature within each thread (i.e. X : start must come before X : End that is predictable. More so, after all the increments, the total sum collected by counter isn't even 20 million!). The actual speed and ordering of this is determined at runtime.

To really see what's going on here, we must look into the assembly code behind this. If we focus on line 11 when the counter is incremented and look at its assembly, we see

<pre> 1 7a7: 8b 05 0a 0b 20 00 mov 0x200b0a(%rip),%eax # 200c1c <counter> 2 7ad: 83 c0 01 add \$0x1,%eax 3 7b0: 89 05 01 0b 20 00 mov %eax,0x200b01(%rip) # 200c1c <counter> </pre>

What this means is that every thread consists of loading the counter value from the instruction pointer plus an offset into %eax, then adding 1 to it, and then storing it back to the same memory location (the rip changes, so the memory address between the 1st and 3rd line will be slightly different, but they are the same address). If we have threads 1 and 2, we can have the following possible interweavings:

1. (1 loads, 1 adds, 1 stores, 2 loads, 2 adds, 2 stores). This would result in a +2.
2. (1 loads, 1 adds, 2 loads, 2 adds, 2 stores, 1 stores). This would result in a +1 since thread 2 loads before 1 could store the incremented version.
3. (1 loads, 2 loads, 1 adds, 2 adds, 1 stores, 2 stores). This would also result in a +1 for the same reasons as before.

There is a big problem here: it seems that these things overwrite each other.

Definition 10.6 (Data Race)

We have just seen an example of a **data race**, which occurs if two or more threads concurrently accesses the same memory location with at least one write. The section of code where a data race can occur is called the **critical section**.

So how do we address this challenge of concurrency? This is where locks and mutexes come in.

Definition 10.7 (Lock)

A **lock** is a construct to enforce mutual exclusion in conflicting code sections (critical sections). It is implemented as a special data object in memory. We can use the API methods

1. **acquire()** or **lock()** is called when going into a critical section.
2. and **release()** or **unlock()** is called when going out of a critical section.

If the lock is already acquired by a thread and not released yet, then other threads will not be able to acquire the lock and execute the next instructions. There are two ways this is implemented. First, an incoming thread can wait if another thread holds the lock, called a **spinlock**,^a or it can be blocked, called a **mutex** (mutual exclusion).^{b,c}

^aThis is mainly implemented through kernel and used almost exclusively for OS development, not application development.

^bThis is mainly implemented in user space.

^cThis has a FIFO queue.

To implement this in C, there's a few things that we have to do.

1. First, make a `pthread_mutex_t` global variable.
2. Then, make initialize the mutex before you create the threads and destroy the mutex after you join the threads in the `main` function.
3. Finally, put the specific locks and unlocks in the locations of the functions that the thread calls.

There are few strategies to correct the counter example above, which all produce the correct output `Counter: 20000000`.

1. We can put the lock around the entire for loop of the `thread` function. However, this essentially takes us back to the sequential regime.

```
1  static volatile int counter = 0;
2  pthread_mutex_t mutex; // global declaration of mutex
3
4  void* thread(void* args) {
5      pthread_mutex_lock(&mutex); //acquire the mutex lock
6      for (int i = 0; i < 10 * 1000 * 1000; i++) {
7          counter += 1;
8      }
9      pthread_mutex_unlock(&mutex); //release the mutex lock
10     return NULL;
11 }
12
13 int main(int argc, char *argv[]) {
14     pthread_t thread1, thread2;
15     int rc;
16     rc = pthread_mutex_init(&mutex, NULL); //initialize the mutex
17     rc = pthread_create(&thread1, NULL, thread, "A");
18     rc = pthread_create(&thread2, NULL, thread, "B");
19     rc = pthread_join(thread1, NULL);
20     rc = pthread_join(thread2, NULL);
21     pthread_mutex_destroy(&mutex); //destroy (free) the mutex
22     printf("Counter : %d\n", counter);
23     return 0;
24 }
```

Figure 73: Putting the mutex locks around the entire for loop.

2. A better idea to actually implement parallelization is to put the locks around the line that says `counter += 1;`. This is the critical code that loads the counter value from the stack into the register, increments it by 1, and sends it back to the stack. We should isolate this so that no other threads can execute during these 3 assembly lines.

```
1  static volatile int counter = 0;
2  pthread_mutex_t mutex; // global declaration of mutex
3
4  void* thread(void* args) {
5      for (int i = 0; i < 10 * 1000 * 1000; i++) {
6          pthread_mutex_lock(&mutex); //acquire the mutex lock
7          counter += 1;
8          pthread_mutex_unlock(&mutex); //release the mutex lock
9      }
10     return NULL;
11 }
12
13 int main(int argc, char *argv[]) {
14     pthread_t thread1, thread2;
15     int rc;
16     rc = pthread_mutex_init(&mutex, NULL); //initialize the mutex
17     rc = pthread_create(&thread1, NULL, thread, "A");
18     rc = pthread_create(&thread2, NULL, thread, "B");
19     rc = pthread_join(thread1, NULL);
20     rc = pthread_join(thread2, NULL);
21     pthread_mutex_destroy(&mutex); //destroy (free) the mutex
22     printf("Counter : %d\n", counter);
23     return 0;
24 }
```

Figure 74: Now we put the locks within the counter. However, this is much slower since locking and unlocking are relatively expensive. It runs in 0.274s.

3. The first two tries are not ideal, but what we can do is have each thread store its local work all within each of its stack, and then when it communicates with the shared memory on `counter`, this is where the locks should come in place. This has the double benefit of locking/unlocking very few times, along with protecting the critical section of the code.

```

1  static volatile int counter = 0;
2  pthread_mutex_t mutex; // global declaration of mutex
3
4  void* thread(void* args) {
5      int my_counter = 0;
6      for (int i = 0; i < 10 * 1000 * 1000; i++) {
7          my_counter += 1;
8      }
9      pthread_mutex_lock(&mutex); //acquire the mutex lock
10     counter += my_counter;
11     pthread_mutex_unlock(&mutex); //release the mutex lock
12     return NULL;
13 }
14
15 int main(int argc, char *argv[]) {
16     pthread_t thread1, thread2;
17     int rc;
18     rc = pthread_mutex_init(&mutex, NULL); //initialize the mutex
19     rc = pthread_create(&thread1, NULL, thread, "A");
20     rc = pthread_create(&thread2, NULL, thread, "B");
21     rc = pthread_join(thread1, NULL);
22     rc = pthread_join(thread2, NULL);
23     pthread_mutex_destroy(&mutex); //destroy (free) the mutex
24     printf("Counter : %d\n", counter);
25     return 0;
26 }

```

Figure 75: In here, we have each thread increment its own local version of the counter and store it in `my_counter`. Then, we increment the global counter by the total `my_counter`, which will require one mutex lock. It runs in 0.049s.

Example 10.6 (Inserting into Linked List)

Inserting into a linked list is a sequential process, but what if we want to parallelize it with multiple threads? Consider the following code.

```

1  typedef struct __node_t {
2      int key;
3      struct __node_t *next;
4  } node_t;
5
6  typedef struct __list_t {
7      node_t *head;
8  } list_t;
9
10 void List_Init(list_t *L) {
11     L -> head = NULL;
12 }
13
14 void List_Insert(list_t *L, int key) {
15     // insert a new node with the key value at the beginning of the list
16     node_t *new = malloc(sizeof(node_t));
17     assert(new);
18     new -> key = key;
19     new -> next = L -> head;
20     L -> head = new;

```

```

21 }
22
23 int main(void) {
24     list_t L;
25     list_t *Lp = &L;
26     List_Init(Lp);
27     for (int i = 0; i < 10; i++) {
28         List_Insert(Lp, i);
29     }
30     return 0;
31 }

```

Say that we want to parallelize the creation of the length 10 linked list. Simply initializing some threads and calling `List_Insert` 10 times won't properly create this linked list. This is because two threads may overwrite where `L -> head` points to. The most naive thing to do is to put locks around the whole function, but now we are back to the sequential regime! If we think about it, the `malloc` calls, asserting that there is viable memory, and assigning the input key value to `new -> key` does not overwrite anything else. It is only when we assign `new -> next` to the head value of `L` that things may get overwritten, so we put the locks shown below, while slightly modifying the `list` struct to have the `pthread_mutex_t` attribute.

```

1  typedef struct __list_t {
2      node_t *head;
3      pthread_mutex_t lock;
4  } list_t;
5
6  void List_Init(list_t *L) {
7      L -> head = NULL;
8      pthread_mutex_init(&L -> lock, NULL);
9  }
10
11 void List_Insert(list_t *L, int key) {
12     // insert a new node with the key value at the beginning of the list
13     node_t *new = malloc(sizeof(node_t));
14     assert(new);
15     new -> key = key;
16     pthread_mutex_lock(&L->lock)
17     new -> next = L -> head;
18     L -> head = new;
19     pthread_mutex_unlock(&L->lock)
20 }

```

10.4 Deadlock Bugs

Locks are extremely useful to segment out a portion of code that should be uninterrupted. However, there are many consequences of misusing it.

Definition 10.8 (Deadlock Bugs)

A **deadlock bug** happens when you make locks such that the program cannot run anymore. These aren't specific to threads, but also processes as well. They require the four preconditions:

1. Mutual exclusion: you must have a lock to begin with to have a deadlock, so this is trivial.
2. Hold and Wait: Threads must have the ability to hold resources (e.g. the philosophers must hold forks which prevent others from taking them) .

3. No Preemption: Preemption refers to the ability to take the fork out of someone else's hand. So if you have a thread, you can first take lock A, then look at whether lock B is available. If not, then unlock A and try again. However, this can lead to a bunch of threads just simply picking up and putting down forks, causing a *livelock*.
4. Circular Wait: That is, there exists a circular chain of threads such that each thread holds a resource needed by the next thread. A strategy is to define a fixed acquisition order for locks (e.g. lock A always before lock B).

You shouldn't try to hold multiple locks at once, but if you must, you should have a strategy to avoid deadlock. Choosing a lock order is the recommended way.

Example 10.7 (Bank Accounts)

Let's go through an example. Suppose we had the following code below.

```

1  struct account {
2      pthread_mutex_t lock;
3      int balance;
4  };
5
6  struct arg_t {
7      struct account fromAcct;
8      struct account toAcct;
9      int amt;
10 };
11
12 pthread_mutex_t mutex; // global declaration of mutex
13
14 void* Transfer(void* args) {
15
16     struct arg_t* data = (struct arg_t*)args;
17
18     struct account* fromAcct = &(data -> fromAcct);
19     struct account* toAcct = &(data -> toAcct);
20     int amt = data -> amt;
21
22
23     pthread_mutex_lock(&fromAcct->lock);
24     pthread_mutex_lock(&toAcct->lock);
25
26     fromAcct->balance -= amt;
27     toAcct->balance += amt;
28
29     pthread_mutex_unlock(&fromAcct->lock);
30     pthread_mutex_unlock(&toAcct->lock);
31
32     return NULL;
33 }
```

Suppose that Threads 0 and 1 are executing concurrently and represent users A and B, respectively. Now consider the situation in which A and B want to transfer money to each other: A wants to transfer 20 dollars to B, while B wants to transfer 40 to A.

Both threads concurrently execute the `Transfer` function. Thread 0 acquires the lock of `acctA` while Thread 1 acquires the lock of `acctB`. Now consider what happens. To continue executing, Thread 0 needs to acquire the lock on `acctB`, which Thread 1 holds. Likewise, Thread 1 needs to acquire the lock on `acctA` to continue executing, which Thread 0 holds. Since both threads are blocked on each

other, they are in deadlock.

This can be simply fixed by rearranging the locks so that each lock/unlock pair surrounds only the balance update statement associated with it.

```

1 void *Transfer(void *args){
2     //argument passing removed to increase readability
3     //...
4
5     pthread_mutex_lock(&fromAcct->lock);
6     fromAcct->balance -= amt;
7     pthread_mutex_unlock(&fromAcct->lock);
8
9     pthread_mutex_lock(&toAcct->lock);
10    toAcct->balance += amt;
11    pthread_mutex_unlock(&toAcct->lock);
12
13    return NULL;
14 }
```

Example 10.8 (Dining Philosopher's Problem)

There are 5 philosophers at a roundtable with 5 plates of food with 5 forks. Each philosopher needs two forks (both on their left and on their right) to start eating their plate. However, there can be a deadlock if every philosopher takes the fork on their right and is always waiting for the left fork (this happens due to a circular dependency).

1. One way to resolve this issue is to set an ID to every fork and have the philosophers always take the lower ID fork first before trying to take the higher ID. Then this resolves and at least one philosopher can eat.

Example 10.9 (Set-Intersection Problem)

Again, like the bank account section, we have the following code that can create deadlocks if there is a code that attempts to compute $A \cap B$ and $B \cap A$ at similar times.

```

1 set_t *set_intersection(set_t *s1, set_t *s2) {
2     set_t *result = set_create();
3     pthread_mutex_lock(&s1->lock);
4     pthread_mutex_lock(&s2->lock);
5
6     for (int i = 0; i < s1->size; i++) {
7         if (set_contains(s2, s1->data[i])) {
8             set_add(result, s1->data[i]);
9         }
10    }
11
12    pthread_mutex_unlock(&s2->lock);
13    pthread_mutex_unlock(&s1->lock);
14    return result;
15 }
```

Like using the IDs and ordering, we can resolve this by grabbing locks in a defined order, say from high to low. Therefore, we replace the locks as such.

```

1 if (&m1 > &m2) {
```

```

2 pthread_mutex_lock(&m1);
3 pthread_mutex_lock(&m2);
4 } else {
5 pthread_mutex_lock(&m2);
6 pthread_mutex_lock(&m1);
7 }

```

However, this still creates a deadlock if we compute $A \cap A$, so we should place a deadlock.

It turns out that we can make logically equivalent code without locks! By using atomic primitives that have implementations in C, we can create wait-free algorithms. For example, consider the `int CompAndSwap(int* addr, int expected, int new)` function. If `*addr == expected`, then we set `*addr` to `new` and return 1. Otherwise, we do nothing and return 0. Then, the two pieces of code are equivalent, but the advantage of the RHS is that there is never a deadlock.

1. In the left, we pass by pointer to increment a variable.
2. In the right, we first set `old` to the value of `val`, and if the value of `val` is equal to `old` (which may not always be true if some other thread modifies this value), then we set `old = *val` again and increment it by `amt`.

<pre> 1 void add(int *val, int amt) { 2 pthread_mutex_lock(&m); 3 *val += amt; 4 pthread_mutex_unlock(&m); 5 } 6 . </pre>	<pre> 1 void add(int *val, int amt) { 2 do { 3 int old = *val; 4 } 5 while(!CompAndSwap(val, old, old+amt)); 6 } </pre>
---	---

However, the left hand, despite the risk of deadlocks, is still recommended. The mutex API allows for better distribution of CPU.

10.5 Order Violation Bugs and Condition Variables

Definition 10.9 (Order Violation Bugs)

An **order violation bug** happens when the desired order between two memory accesses is flipped, i.e. A should always be executed before B, but the order is not enforced during execution.

1. For example, one thread that acts on a linked list may already assume that the linked list is initialized when it is not. Therefore, we should have these threads wait until initialization.
2. Another example is if there is a set of producers and consumers that manage the flow of items through a buffer. If there is no good, then consumers cannot consume and they must wait. If the supply is full then the producers must wait. This also requires some process of waiting.

We can use a regular conditional to check, but this may not be efficient due to the following example.

Example 10.10 (Condition Checks with Vanilla While Loops)

Say that you want to call the following function on a linked list, but you must wait until the list is not NULL.

```

1 int getItem(list_t *list) {
2 pthread_mutex_lock(&list -> m);
3 // TODO: wait until list is non-empty
4 int item = list -> head -> item;
5 list -> head = list -> head -> next;

```

```

6 pthread_mutex_unlock(&list -> m);
7 return item;
8 }

```

To check whether `list` is not a null pointer, we can simply use a while loop to check.

```

1 int tryGetItem(list_t *list, int *out) {
2     int success = 0;
3     pthread_mutex_lock(&list -> m);
4     if (list -> head) {
5         success = 1;
6         *out = list -> head -> item;
7         list -> head = list -> head -> next;
8     }
9     pthread_mutex_unlock(&list -> m);
10    return success;
11 }
12 .
13 .

```

```

1 int getItem( list_t *list) {
2     pthread_mutex_lock(&list -> m);
3     // TODO: wait until list is non-empty
4     while (list -> head == NULL) {
5         pthread_mutex_unlock(&list -> m);
6         yield(); // optional
7         pthread_mutex_lock(&list -> m);
8     }
9     int item = list -> head -> item;
10    list -> head = list -> head -> next;
11    pthread_mutex_unlock(&list -> m);
12    return item;
13 }

```

Figure 76: On the LHS, this returns 0 if the retrieval is not successful, and so you must wrap this function within some while loop to check if it is actually successful. On the RHS, we can use the `yield()` function which gives control to the OS to schedule another thread.

This constant checking through a while loop leads to a waste of CPU resources. Therefore, we want to put the thread to sleep while there is no element in `list` so other processes can use the CPU core. This is the motivation behind *conditional variables*.

Definition 10.10 (Condition Variables)

Condition Variables force a thread to be blocked until a particular condition is reached. This construct is useful for scenarios in which a condition must be met before the thread does some work.

1. Every CV is bound to exactly one mutex. This is because the state of a condition, even if true on one thread, can be changed immediately by another thread, so some sort of locking is needed.
2. Condition variables have the type `pthread_cond_t`.
3. To initialize a condition variable, use the `pthread_cond_init` function.
4. To destroy it, use `pthread_cond_destroy`.
5. `pthread_cond_wait(&cond, &mutex)` takes the address of a condition variable `cond` and a mutex `mutex` as its arguments. It causes the calling thread to block on the condition variable `cond` until another thread signals it (or "wakes" it up).
6. The `pthread_cond_signal(&cond)` function causes the calling thread to unblock (or signal) another thread that is waiting on the condition variable `cond` (based on scheduling priority). If no threads are currently blocked on the condition, then the function has no effect. Unlike `pthread_cond_wait`, the `pthread_cond_signal` function can be called by a thread regardless of whether or not it owns the mutex in which `pthread_cond_wait` is called.

It is important to know about the states that a thread can be in. It can either be currently running/active, ready to run (perhaps through a syscall), or blocked.



Figure 77: Note that the only way a thread can be blocked is if it is waiting for a condition to happen. If that condition happened and a signal arrives at the thread, then it “wakes up.” Condition variables allow these thread to go in and out of the blocked state.

With this, the general design pattern is as such.

Theorem 10.2 (Condition Variable Design Pattern)

To implement this effectively, we must first identify a state that will be accessed by 2+ threads concurrently and add locks to protect the shared state. If we need to wait on some condition, we use condition variables.

<pre> 1 methodThatWaits() { 2 pthread_mutex_lock(&m); 3 4 // Read/write shared state 5 6 while (!checkSharedState()) { 7 pthread_cond_wait(&cv, &m); 8 } 9 10 // Read/write shared state 11 12 pthread_mutex_unlock(&m); 13 } </pre>	<pre> 1 methodThatSignals() { 2 pthread_mutex_lock(&m); 3 4 // Read/write shared state 5 6 // If checkSharedState() is now true 7 pthread_cond_signal(&cv); 8 9 // Read/write shared state 10 11 pthread_mutex_unlock(&m); 12 } 13 . </pre>
--	---

The implementation is quite complex at first, so let’s go through an example.

Example 10.11 (Soda Machine)

Say that we want to model a soda machine, with consumers taking soda from the machine and producers filling the machine up.

1. We’d like to create some variables that encode the state of the soda machine. This is the **shared state**.

```

1  static volatile int numSodas;
2  #define MaxSodas 100;
                
```

2. We also want to implement *one* lock to protect all shared states, say

```

1  pthread_mutex_t sodaLock;
                
```

This allows us to implement mutual exclusion so that only one thread can manipulate the

machine (state) at a time.

3. The **ordering constraints** are that the consumer must wait if the machine is empty (CV `hasSoda`), and the producer must wait if the machine is full (CV `hasRoom`).
4. The first thing we must do is make sure that the consumer and producer function has a lock and unlock over its body since both functions modify the vending machine.

<pre> 1 consumer() { 2 pthread_mutex_lock(&sodaLock); 3 4 // take a soda from machine 5 6 pthread_mutex_unlock(&sodaLock); 7 }</pre>	<pre> 1 producer() { 2 pthread_mutex_lock(&sodaLock); 3 4 // add a soda to machine 5 6 pthread_mutex_unlock(&sodaLock); 7 }</pre>
--	---

5. Moreover, the consumer and producer's actions of taking or adding soda is dependent on the state already. For the consumer, it should **wait** if the machine is empty for a signal that notifies that it is not empty. Once the consumer receives the signal, it takes a soda and can send a signal that the machine is not full to the producer function.

<pre> 1 consumer() { 2 pthread_mutex_lock(&sodaLock); 3 // wait if empty 4 // take a soda from machine 5 // notify that it is not full 6 pthread_mutex_unlock(&sodaLock); 7 }</pre>	<pre> 1 producer() { 2 pthread_mutex_lock(&sodaLock); 3 // wait if full 4 // add a soda to machine 5 // notify that it is not empty 6 pthread_mutex_unlock(&sodaLock); 7 }</pre>
---	--

6. To put this into code, we finally have

<pre> 1 consumer() { 2 pthread_mutex_lock(&sodaLock); 3 4 while (numSodas == 0) { 5 // while empty 6 wait(sodaLock, hasSoda); 7 } 8 numSodas -= 1; // take soda 9 signal(hasRoom); 10 pthread_mutex_unlock(&sodaLock); 11 }</pre>	<pre> 1 producer() { 2 pthread_mutex_lock(&sodaLock); 3 4 while (numSodas == MaxSodas) { 5 wait(sodaLock, hasRoom); 6 } 7 8 numSodas += 1; // add soda 9 signal(hasSoda); 10 pthread_mutex_unlock(&sodaLock); 11 }</pre>
--	---

Let's go through this. From the consumer function, we have:

1. The consumer function can start by first acquiring a lock on the `sodaLock` mutex using `pthread_mutex_lock()`. This ensures that only one consumer thread can access the shared resources (e.g., `numSodas`) at a time.
2. It then enters a while loop^a that checks if `numSodas` is zero. If `numSodas` is zero, it means there are no sodas available for consumption. In this case, the consumer thread calls the `wait()` function, which unlocks the `sodaLock` mutex and waits for a signal on the `hasSoda` condition variable. This allows other threads to acquire the lock and proceed.
3. When the consumer thread receives a signal indicating that sodas are available (`hasSoda` condition variable), it wakes up and reacquires the lock on `sodaLock`.
4. The consumer thread then decrements `numSodas` by 1, indicating the consumption of a soda.
5. After consuming a soda, the consumer thread signals the `hasRoom` condition variable using the `signal()` function. This notifies any waiting producer threads that there is now room available in the soda buffer.
6. Finally, the consumer thread unlocks the `sodaLock` mutex using `pthread_mutex_unlock()`,

allowing other threads to access the shared resources.

From the producer function, we have:

1. The producer function can start by first acquiring a lock on the `sodaLock` mutex using `pthread_mutex_lock()`. This ensures that only one producer thread can access the shared resources at a time.
2. It then enters a while loop that checks if `numSodas` is equal to `MaxSodas`. If `numSodas` is equal to `MaxSodas`, it means the soda buffer is full, and the producer cannot add more sodas. In this case, the producer thread calls the `wait()` function, which unlocks the `sodaLock` mutex and waits for a signal on the `hasRoom` condition variable. This allows other threads to acquire the lock and proceed.
3. When the producer thread receives a signal indicating that there is room available in the soda buffer (`hasRoom` condition variable), it wakes up and reacquires the lock on `sodaLock`.
4. The producer thread then increments `numSodas` by 1, indicating the production of a new soda.
5. After producing a soda, the producer thread signals the `hasSoda` condition variable using the `signal()` function. This notifies any waiting consumer threads that a soda is now available for consumption.
6. Finally, the producer thread unlocks the `sodaLock` mutex using `pthread_mutex_unlock()`, allowing other threads to access the shared resources.

Note that we must have a while loop since if the producer ended up broadcasting the `hasSoda` condition to say 10 threads when there are 5 sodas, then 5 of those threads will get a soda while 5 may not, and this possibility should be detected by the while loop.

^aWe want this to be a while loop since we want to reevaluate the condition even after reacquiring the lock.