

# Sysprog Summary

November 29, 2020

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
1.1	Computer arithmetic . . . . .	6
<b>2</b>	<b>Introduction to C</b>	<b>7</b>
2.1	Workflow . . . . .	7
2.1.1	Questions . . . . .	7
2.2	control flow in C . . . . .	8
2.3	Functions . . . . .	8
2.4	Basic Types in C . . . . .	8
2.4.1	Declarations . . . . .	8
2.4.2	Integers and floats . . . . .	9
2.4.3	C99 extended integer types . . . . .	9
2.4.4	Booleans . . . . .	9
2.4.5	void . . . . .	10
2.4.6	Operators . . . . .	10
2.4.7	Casting . . . . .	10
2.5	Arrays in C . . . . .	10
2.5.1	Strings . . . . .	10
<b>3</b>	<b>Representing Integers in C</b>	<b>12</b>
3.1	Encodings and operators . . . . .	12
3.1.1	Representing Integers . . . . .	12
3.1.2	Bit-level operations in C . . . . .	12
3.2	Integer Ranges . . . . .	13
3.2.1	Encoding Integers . . . . .	13
3.2.2	Numeric ranges . . . . .	13
3.2.3	Signed vs Unsigned in C . . . . .	13
3.2.4	Sign extension . . . . .	14
3.3	Integer addition and subtraction in C . . . . .	14
3.3.1	Unsigned addition . . . . .	14
3.3.2	Two's complement addition . . . . .	14
3.3.3	Integer multiplication in C . . . . .	15
3.4	Integer multiplication and division using shifts . . . . .	15
3.4.1	Power of 2 multiply with shift . . . . .	15
3.4.2	Unsigned power of 2 divide with shift . . . . .	15
3.4.3	Signed power of 2 divide with shift . . . . .	16
<b>4</b>	<b>Pointers</b>	<b>17</b>
4.1	Loading . . . . .	17
4.2	Stack . . . . .	17
4.2.1	Stack Frames . . . . .	17
4.2.2	Call Chain example: . . . . .	18
4.3	Pointers . . . . .	19
4.3.1	Addresses and & . . . . .	19
4.3.2	Pointers . . . . .	19
4.3.3	Box and arrow diagrams . . . . .	19
4.3.4	Address Space Layout Randomization . . . . .	19
4.3.5	NULL . . . . .	19
4.4	Pointer arithmetic . . . . .	20
4.5	Arrays and Pointers . . . . .	20
4.6	Passing by Reference . . . . .	20
4.6.1	Pass by value . . . . .	20
4.7	Pass by Reference . . . . .	21
4.8	C Pointer Declarations . . . . .	21
<b>5</b>	<b>Dynamic memory allocation</b>	<b>22</b>
5.1	The C memory API . . . . .	23
5.1.1	Deallocation . . . . .	23
5.1.2	Resize . . . . .	23
5.1.3	size_t . . . . .	24

5.1.4	Summary . . . . .	24
5.2	Managing the heap . . . . .	24
5.2.1	Memory Corruption . . . . .	24
5.2.2	Memory Leaks . . . . .	25
5.3	Structures and unions . . . . .	25
5.3.1	Structured data . . . . .	25
5.3.2	Unions . . . . .	26
5.4	Type definitions . . . . .	26
5.4.1	typedef . . . . .	26
5.4.2	Struct tags and typedefs . . . . .	27
5.4.3	C namespaces . . . . .	27
5.5	Dynamic Data Structures . . . . .	27
5.5.1	Singly-Linked List . . . . .	27
<b>6</b>	<b>Wrapping up C</b>	<b>29</b>
6.1	The C Preprocessor . . . . .	29
6.1.1	Macro Definitions . . . . .	29
6.2	Modularity . . . . .	30
6.3	Function Pointers . . . . .	31
6.4	Assertions . . . . .	32
6.4.1	Assertions . . . . .	32
6.5	goto . . . . .	32
6.5.1	When to use goto: . . . . .	32
6.6	setjmp() and longjmp() . . . . .	33
6.7	Coroutines . . . . .	33
<b>7</b>	<b>Implementing dynamic memory allocation</b>	<b>35</b>
7.1	Dynamic Memory allocation . . . . .	35
7.1.1	Explicit allocation . . . . .	36
7.2	The Explicit memory allocation Problem . . . . .	36
7.2.1	Performance goal: Throuput . . . . .	37
7.2.2	Performance goal: Peak memory utilization . . . . .	37
7.2.3	Fragmentation . . . . .	37
7.2.4	How to free blocks . . . . .	38
7.2.5	Keeping track of free blocks . . . . .	38
7.3	Coalescing . . . . .	39
7.4	Explicit free lists . . . . .	39
7.4.1	concept . . . . .	39
7.4.2	Allocation . . . . .	40
7.4.3	freeing . . . . .	40
7.4.4	summary . . . . .	40
7.5	Segregated free lists . . . . .	40
<b>8</b>	<b>Basic x86 architecture</b>	<b>41</b>
8.1	What is an instruction set architecture . . . . .	41
8.1.1	CISC: Complex Instruction Set . . . . .	41
8.1.2	RISC: Reduced Instruction Set . . . . .	41
8.2	Basics of machine code . . . . .	42
8.2.1	Compiling into assembly . . . . .	42
8.2.2	Assembly data types . . . . .	42
8.2.3	Assembly code operations . . . . .	42
8.2.4	Object Code . . . . .	42
8.2.5	Machine instruction example . . . . .	43
8.3	x86 architecture . . . . .	43
8.3.1	8086 Registers . . . . .	43
8.3.2	80386 Registers . . . . .	43
8.3.3	x86-64 integer registers . . . . .	44
8.3.4	Moving data . . . . .	44
8.3.5	Simple memory addressing modes . . . . .	44
8.3.6	Address Computation . . . . .	45
8.4	x86 integer arithmetic . . . . .	45
8.4.1	ordinary arithmetic operations . . . . .	45
8.4.2	Single operand instructions . . . . .	45
8.4.3	Using leal for arithmetic expressions . . . . .	46
8.5	condition codes . . . . .	46
8.5.1	Compare Instruction . . . . .	46
8.5.2	Test Instruction . . . . .	46
8.5.3	SetX Instructions . . . . .	46
8.5.4	Jump Instructions . . . . .	47
<b>9</b>	<b>Compiling C control flow</b>	<b>48</b>
9.1	if-then-else statements . . . . .	48
9.1.1	Conditional move . . . . .	48
9.2	do-while loops . . . . .	48

9.3	while loops . . . . .	49
9.4	for loops . . . . .	49
9.4.1	Summary of loop variants . . . . .	49
9.5	compact switch statements . . . . .	50
9.6	sparse switch statements . . . . .	50
9.7	Procedure call and return . . . . .	51
9.8	Calling conventions . . . . .	51
9.8.1	Register saving conventions . . . . .	52
9.8.2	Procedure Calls: Locals in the red zone . . . . .	53
9.8.3	Procedure Calls: Long swap . . . . .	53
9.8.4	Procedure Calls: non-leaf w/o stack frame . . . . .	53
9.8.5	Procedure Calls: call using a jump . . . . .	53
<b>10</b>	<b>Compiling C data Structures</b>	<b>54</b>
10.1	One-dimensional arrays . . . . .	54
10.1.1	Array allocation . . . . .	54
10.1.2	Array access . . . . .	54
10.1.3	Array Loop example . . . . .	55
10.2	Nested Array . . . . .	55
10.2.1	Nested array row access . . . . .	55
10.2.2	Nested array row access code . . . . .	55
10.2.3	Nested array element access code . . . . .	56
10.3	Multi-level arrays . . . . .	56
10.3.1	Element access in multi-level array . . . . .	56
10.3.2	Referencing examples . . . . .	56
10.4	Structures . . . . .	57
10.4.1	Accessing structure members . . . . .	57
10.5	Alignment . . . . .	57
10.5.1	Alignment with structures . . . . .	57
10.6	Arrays of structures . . . . .	58
10.7	Unions . . . . .	58
<b>11</b>	<b>Linking</b>	<b>59</b>
11.1	Object files . . . . .	59
11.2	Linker symbols . . . . .	60
11.2.1	Resolving symbols . . . . .	60
11.2.2	Relocating code and data . . . . .	60
11.2.3	Strong and weak symbols . . . . .	60
11.3	Static libraries . . . . .	61
11.3.1	Loading executable object files . . . . .	62
11.4	Shared libraries . . . . .	62
11.4.1	Dynamic linking at load time . . . . .	62
11.4.2	Dynamic linking at runtime . . . . .	62
<b>12</b>	<b>Code Vulnerabilities</b>	<b>63</b>
12.1	Worms and Viruses . . . . .	63
12.2	Stack Overflow bug . . . . .	63
12.3	Stopping overrun bugs . . . . .	64
12.4	XDR vulnerability . . . . .	64
12.5	CTF . . . . .	64
<b>13</b>	<b>Floating Point</b>	<b>65</b>
13.1	Representing Floating Point . . . . .	65
13.1.1	Fractional binary numbers . . . . .	65
13.1.2	Floating point representation . . . . .	65
13.2	Types of IEEE floating point numbers . . . . .	66
13.2.1	Precisions . . . . .	66
13.2.2	Floating point in C . . . . .	66
13.2.3	Normalized Values . . . . .	66
13.2.4	Denormalized values . . . . .	66
13.2.5	Special values . . . . .	67
13.2.6	Properties of encoding . . . . .	67
13.3	Floating point ranges . . . . .	67
13.4	Floating-point rounding and arithmetic . . . . .	67
13.4.1	Rounding . . . . .	67
13.4.2	Rounding binary numbers . . . . .	68
13.4.3	Creating a floating point number . . . . .	68
13.5	Floating point addition and multiplication . . . . .	68
13.5.1	Floating point multiplication . . . . .	68
13.5.2	Floating point addition . . . . .	68
13.6	SSE floating point . . . . .	69
13.6.1	SSE3 instruction names . . . . .	70
13.6.2	SSE3 Examples . . . . .	70
13.6.3	Constants . . . . .	70

13.6.4	Vector Instructions . . . . .	71
<b>14</b>	<b>Optimizing Compilers</b>	<b>72</b>
14.1	Optimizing Compilers . . . . .	72
14.2	Code motion and precomputation . . . . .	73
14.3	Strength reduction . . . . .	73
14.4	Common subexpressions . . . . .	73
14.5	Optimization blocker: Procedure calls . . . . .	73
14.6	Optimization blocker: memory aliasing . . . . .	73
14.7	Blocking and unrolling . . . . .	74
<b>15</b>	<b>Architecture and Optimization</b>	<b>75</b>
15.1	modern processor design . . . . .	76
15.1.1	Sequential processor stages . . . . .	76
15.1.2	Pipelined hardware . . . . .	76
15.1.3	Performance . . . . .	76
15.1.4	Superscalar processor . . . . .	77
15.2	Superscalar processor performance . . . . .	77
15.2.1	Intel Haswell CPU schematics . . . . .	77
15.2.2	Latency vs Throughput . . . . .	77
15.2.3	Data Hazards . . . . .	77
15.2.4	Register Renaming . . . . .	78
15.2.5	Instruction Execution . . . . .	78
15.2.6	Meaning for Program performance . . . . .	78
15.3	Reassociation . . . . .	78
15.3.1	Separate Accumulators . . . . .	79
15.4	Combining multiple accumulators and unrolling . . . . .	79
15.4.1	AVX2 SIMD operations . . . . .	80
<b>16</b>	<b>Caches</b>	<b>81</b>
16.1	Intro . . . . .	81
16.1.1	Processor-memory bottleneck . . . . .	81
16.1.2	General Cache concepts . . . . .	81
16.1.3	Two-level cache performance . . . . .	81
16.1.4	Types of cache misses . . . . .	82
16.2	Cache organization . . . . .	82
16.3	Cache reads . . . . .	82
16.3.1	Direct mapped cache . . . . .	82
16.3.2	2-way set-associative cache . . . . .	83
16.4	The memory hierarchy . . . . .	83
16.5	Cache writes . . . . .	83
16.5.1	Write-hit . . . . .	83
16.5.2	Write-miss . . . . .	84
16.5.3	Other hardware cache features . . . . .	84
16.6	Cache optimizations . . . . .	84
16.6.1	Locality . . . . .	84
16.6.2	Locality with matrix multiplication . . . . .	84
16.7	Blocking . . . . .	85
<b>17</b>	<b>Exceptions</b>	<b>86</b>
17.1	Intro . . . . .	86
17.1.1	Control Flow . . . . .	86
17.1.2	Altering the control flow . . . . .	86
17.1.3	Exceptional control flow . . . . .	86
17.1.4	Exceptions . . . . .	86
17.2	Exception vectors and kernel mode . . . . .	87
17.2.1	Exception vectors . . . . .	87
17.2.2	x86 exception vectors . . . . .	87
17.2.3	Kernel . . . . .	87
17.3	Synchronous Exceptions . . . . .	88
17.3.1	Fault Examples . . . . .	88
17.4	Asynchronous Exceptions (Interrupts) . . . . .	88
17.4.1	x86 interrupts . . . . .	88
17.5	Interrupt controllers . . . . .	89
17.5.1	Programmable Interrupt Controllers (PIC) . . . . .	89
17.5.2	Modern PICs . . . . .	89
<b>18</b>	<b>Virtual Memory</b>	<b>90</b>
18.1	Address Translation: . . . . .	90
18.1.1	Address Spaces: . . . . .	90
18.1.2	Address Translation with page table . . . . .	90
18.2	Uses of virtual memory . . . . .	91
18.2.1	Benefits . . . . .	91
18.2.2	VM as a tool for caching . . . . .	91

18.2.3 DRAM cache organization . . . . .	91
18.2.4 Memory Management . . . . .	91
18.2.5 Protection and Sharing . . . . .	92
18.3 The address translation process . . . . .	92
18.3.1 Page Hit: . . . . .	92
18.3.2 Page Fault: . . . . .	93
18.4 Translation lookaside buffer . . . . .	93
18.4.1 TLB hit . . . . .	93
18.4.2 TLB Miss . . . . .	93
18.5 Simple Memory system Example . . . . .	94
18.5.1 Page Table . . . . .	94
18.5.2 TLB . . . . .	94
18.5.3 Cache . . . . .	94
18.5.4 Address translation examples . . . . .	95
18.6 Multi-level page tables . . . . .	95
18.6.1 Terminology . . . . .	95
18.6.2 Linear Page Table size . . . . .	95
18.6.3 2-level page table hierarchy . . . . .	95
18.7 Case study of the core i7 processor . . . . .	96
18.7.1 Core i7 TLB . . . . .	96
18.7.2 Translating with the i7 TLB . . . . .	97
18.7.3 x86-64 paging . . . . .	97
18.8 Core i7 cache . . . . .	97
18.8.1 cache access . . . . .	97
18.8.2 Cache size . . . . .	98
18.9 Caches revisited . . . . .	98
18.9.1 Virtually-indexed, virtually tagged . . . . .	98
18.9.2 Virtually-indexed, physically tagged . . . . .	99
18.9.3 Physically-indexed, physically tagged . . . . .	99
18.9.4 Write Buffers . . . . .	99

# Chapter 1

## Introduction

### 1.1 Computer arithmetic

Arithmetic operation have important mathematical properties but we cannot assume all usual mathematical properties, due to finiteness of representations e.g  $x^2 \geq 0$  is valid for floats but not for ints and associativity is valid for ints but not for floats.

C and C++ dont provide any memory protections:

- Out of bounds array references
- Invalid pointer values
- Abuses malloc/free

Besides asymptotic complexity constant factors matter to hence even the exact operation count does not predict performance, Optimization happens at multiple levels e.g algorithms,data representations, procedures and loops.

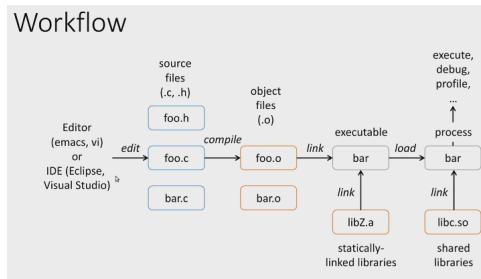
## Chapter 2

# Introduction to C

**#include** Takes a file and includes it into the compilation of the current file

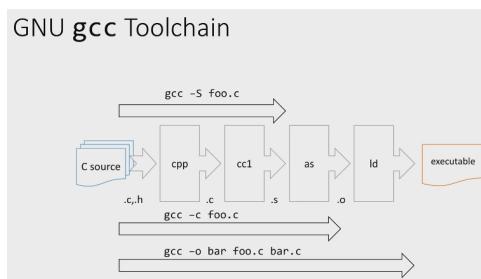
**return 0** In C there are no exceptions, hence to indicate that everything went well we return 0 otherwise we return a non-zero value

### 2.1 Workflow



The steps:

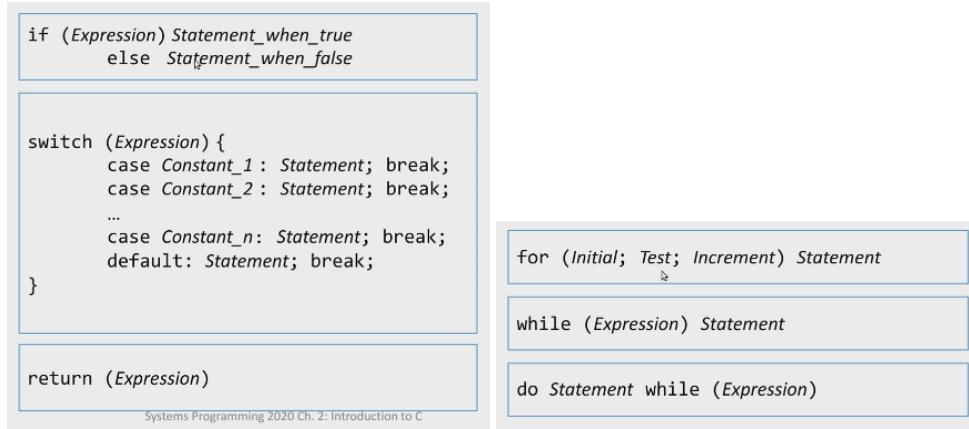
- Edit a file
- Get different kinds of c files
- c files get compiled into object code, which is machine code which doesn't connect to any part of the program i.e each of the files get compiled in isolation
- the linker converts the object files into executables (i.e code which can run)
- loading takes the executable and turns it into a running process in memory and in doing so pulls in the last bits of code needed to run.



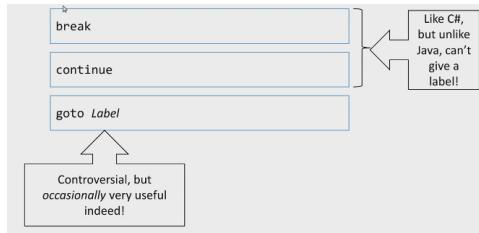
#### 2.1.1 Questions

- What is the difference between an object file and an executable? The object contains machine code, typically it only contains part of the program, hence it contains a lot of placeholders (it knows what it needs to do but not how)

## 2.2 control flow in C



For the switch statement break is needed to leave, otherwise it will continue with the next case, we define a default case which is executed if it doesn't match any cases.



continue stops the current iteration of the loop and starts with the next iteration. "goto" is a unconditional branch.

## 2.3 Functions

Each function similar to java has:

- Name
- Return type
- Argument types
- Body

Every C program needs to have a function called main()

```

/* program to print arguments from command line */
#include <stdio.h>

int main(int argc, char *argv[])
{
    int i;
    printf("argc = %d\n", argc);
    for (i=0; i<argc; ++i) {
        printf("argv[%d]: %s\n", i, argv[i]);
    }
    return 0;
}

```

char \*argv[] can be seen as a list of strings. "int argc" indicates how many strings are in the list  
printf printf takes multiple arguments, the first one is a string each % is replaced with the arguments that follow

## 2.4 Basic Types in C

### 2.4.1 Declarations

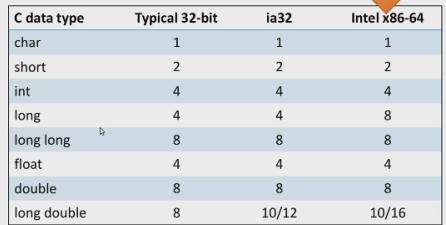
Are like java e.g int *my\_int*, double *some\_float* = 0.123 etc. when we declare a variable inside the block, the scope is just the block after the block, the variable disappears, unless we declare it as static. A static variable persists between calls i.e the next time you call the function the variable will keep its value.

Outside a block, the scope is the entire **program** i.e not just this file but all files that have been compiled in the program. (Called Global Variables). We can limit the scope of a global variable to just the file its in by using static.

## 2.4.2 Integers and floats

The size of an integer in c i.e the max value it can have depends on the machine you are running on.

• Types and sizes:



C data type	Typical 32-bit	ia32	Intel x86-64
char	1	1	1
short	2	2	2
int	4	4	4
long	4	4	8
long long	8	8	8
float	4	4	4
double	8	8	8
long double	8	10/12	10/16

• Integers are **signed** by default  
 • use signed or unsigned to clarify

Hence depending on the machine and using large integers the code could have different performance. Integers in C are either signed or unsigned. They are signed by default.

Rules for arithmetic on integers and floats:

- Implicit conversion between integer types, i.e adding a larger int and a smaller int results in the sum with bit size of the larger int.
- Implicit conversion between floating types
- Explicit conversion between any 2 types can be done with casting

## 2.4.3 C99 extended integer types

```
#include <stdint.h>

int8_t          a;
int16_t         b;
int32_t         c;
int64_t         d;

uint8_t          x;
uint16_t        y;
uint32_t        z;
uint64_t        w;
```

Signed integers,  
precise size in bits

Unsigned integers,

## 2.4.4 Booleans

Booleans in C are also integers. 0 is interpreted as false and anything that's non-zero is interpreted as true. The NOT operator "!" turns zero into a non-zero and vice-versa.

C99 supports a bool type, it is completely optional and still an integer. You can use it by:

```
#include <stdbool.h>
```

In most programming languages we distinguish between statements (things that do stuff, while, if, assignment, etc) and expressions (set of terms which evaluate to a value). In C all statements are also expressions.

```
int rc;
if (rc = call_some_fn()) {
    fprintf(stderr, "Failed with return code %d\n", rc);
    exit(1);
} // Carry on: call succeeded.
```

```
FILE *f;
if (!f = fopen("myfilename", "r")) {
    fprintf(stderr, "Failed with return code %d\n", errno);
    exit(1);
} // Carry on: call succeeded.
```

the exit method ends the program regardless, the argument decides how the program is stopped.

## 2.4.5 void

void is type that has no value, and used for functions which dont return a particular value (procedures). Also used for untyped pointers (pointers that point purely to an address i.e just a pointer we dont know what's there)

## 2.4.6 Operators

Operator	Associativity
0 [] > .	Left-to-right
! ~ ++ -- * - * & (type) sizeof	Right-to-left
* / %	Left-to-right
+ -	Left-to-right
<<>>	Left-to-right
<<>>=	Left-to-right
== !=	Left-to-right
&	Left-to-right
^	Left-to-right
	Left-to-right
&&	Left-to-right
	Left-to-right
? :	Right-to-left
= += -= *= /= %= &= ^=  = <<>>=	Right-to-left
,	Left-to-right

Decreasing precedence

Left to right associativity:  $A+B+C \rightarrow (A+B)+C$

Right to left associativity:  $A+=B+=C \rightarrow A+=(B+=C)$

Post-increment  $i++$ : Take the current value of  $i$ , add one to it and return the value it used to be. Pre-increment  $++i$ : First add one to  $i$  and then return the value of  $i$ . Post/pre-increment works for pointers too.

## 2.4.7 Casting

An operator which takes a value of a certain type and returns a value of a completely different type. Bit-representation usually does not change.

## 2.5 Arrays in C

An array is a finite vector of variables, all the same type. Indexing starts from 0. In an  $n$  element array the last element is  $a[n-1]$ . In most languages bound checking is done and you get an error when you go out of bounds. In C the compiler **does not** check the array bounds!

Multi-dimensional arrays are stored in a contiguous array in row-major.

**Array initializers:** Arrays can be initialized when they are defined:

```

/* a[0] = 3,
   a[1] = 7,
   a[2] = 9 */
int a[3] = {3, 7, 9};

/* list[0]=0.0, ...,
   list[99]=0.0 */
float list[100] = {};

int a[3][3] = {
    { 1, 2, 3},
    { 4, 5, 6},
    { 7, 8, 9},
};
```

### 2.5.1 Strings

C has no real string type, instead a string is an array of char's terminated with null '\0'

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

... is the same as: `char str[6] = "hello";`

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

int main(int argc, char *argv[])
{
    char name1[12], name2[12];
    char mixed[25], title[20];

    strcpy(name1, "Rosalinda");
    strcpy(name2, "Zeke");
    strcpy(title, "This is the title.");
    ...
}

/* returns 1 if name1 > name2 */
if (strcmp(name1, name2) > 0) {
    strcpy(mixed, name1);
} else {
    strcpy(mixed, name2);
}
strcpy(mixed, name1);
printf("The biggest name alphabetically
is %s\n", mixed);

printf("%s\n", title);
printf("Name 1 is %s\n", name1);
printf("Name 2 is %s\n", name2);
...
}

string.c
```

This is the title.

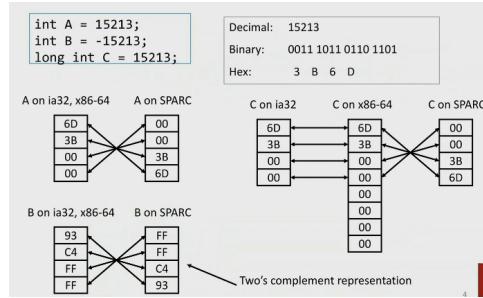
Name 1 is Rosalinda  
Name 2 is Zeke  
The biggest name  
alphabetically is Zeke  
Both names are Rosalinda  
Zeke

# Chapter 3

## Representing Integers in C

### 3.1 Encodings and operators

#### 3.1.1 Representing Integers



The negative in binary is the bits inverted + 1.

#### 3.1.2 Bit-level operations in C

<ul style="list-style-type: none"> <li>Operations &amp;,  , ~, ^ available in C           <ul style="list-style-type: none"> <li>Apply to any “integral” data type               <ul style="list-style-type: none"> <li>long, int, short, char, unsigned</li> </ul> </li> <li>View arguments as bit vectors</li> <li>Arguments applied bit-wise</li> </ul> </li> </ul>	<b>&amp;,  , ~, ^</b>	<ul style="list-style-type: none"> <li>Width w bit vector represents subsets of {0, ..., w-1}</li> <li><math>a_j = 1</math> if <math>j \in A</math>:</li> </ul> <table border="0"> <tr> <td><b>01101001</b></td> <td>{0, 3, 5, 6 }</td> </tr> <tr> <td>76543210</td> <td></td> </tr> </table>	<b>01101001</b>	{0, 3, 5, 6 }	76543210																					
<b>01101001</b>	{0, 3, 5, 6 }																									
76543210																										
<ul style="list-style-type: none"> <li>Examples (char data type):           <ul style="list-style-type: none"> <li><math>\sim 0x41 \rightarrow 0xBE</math> <math>-01000001_2 \rightarrow 10111110_2</math></li> <li><math>\sim 0x00 \rightarrow 0xFF</math> <math>-00000000_2 \rightarrow 11111111_2</math></li> <li><math>0x69 \&amp; 0x55 \rightarrow 0x41</math> <math>01101001_2 \&amp; 01010101_2 \rightarrow 01000001_2</math></li> <li><math>0x69 \mid 0x55 \rightarrow 0x7D</math> <math>01101001_2 \mid 01010101_2 \rightarrow 01111011_2</math></li> </ul> </li> </ul>	<b>01101001</b> 76543210	<ul style="list-style-type: none"> <li>Width w bit vector represents subsets of {0, ..., w-1}</li> <li><math>a_j = 1</math> if <math>j \in A</math>:</li> </ul> <table border="0"> <tr> <td><b>01101001</b></td> <td>{0, 3, 5, 6 }</td> </tr> <tr> <td>76543210</td> <td></td> </tr> </table> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Operator</th> <th>Operation</th> <th>Result</th> <th>Meaning</th> </tr> </thead> <tbody> <tr> <td>&amp;</td> <td>Intersection</td> <td><b>01000001</b></td> <td>{0, 6 }</td> </tr> <tr> <td> </td> <td>Union</td> <td><b>01111101</b></td> <td>{2, 3, 4, 5, 6 }</td> </tr> <tr> <td>^</td> <td>Symmetric difference</td> <td><b>00111100</b></td> <td>{2, 3, 4, 5 }</td> </tr> <tr> <td>~</td> <td>Complement</td> <td><b>10101010</b></td> <td>{1, 3, 5, 7 }</td> </tr> </tbody> </table>	<b>01101001</b>	{0, 3, 5, 6 }	76543210		Operator	Operation	Result	Meaning	&	Intersection	<b>01000001</b>	{0, 6 }		Union	<b>01111101</b>	{2, 3, 4, 5, 6 }	^	Symmetric difference	<b>00111100</b>	{2, 3, 4, 5 }	~	Complement	<b>10101010</b>	{1, 3, 5, 7 }
<b>01101001</b>	{0, 3, 5, 6 }																									
76543210																										
Operator	Operation	Result	Meaning																							
&	Intersection	<b>01000001</b>	{0, 6 }																							
	Union	<b>01111101</b>	{2, 3, 4, 5, 6 }																							
^	Symmetric difference	<b>00111100</b>	{2, 3, 4, 5 }																							
~	Complement	<b>10101010</b>	{1, 3, 5, 7 }																							

the contrast logic operations in C:

<ul style="list-style-type: none"> <li>• <b>&amp;&amp;,    , !</b> <ul style="list-style-type: none"> <li>View 0 as “False”</li> <li>Anything nonzero as “True”</li> <li>Always return 0 or 1</li> <li>Early termination</li> </ul> </li> </ul>	<b>&amp;&amp;,    , !</b>
<ul style="list-style-type: none"> <li>Examples (char data type)           <ul style="list-style-type: none"> <li><math>!0x41 \rightarrow 0x00</math></li> <li><math>!0x00 \rightarrow 0x01</math></li> <li><math>!!0x41 \rightarrow 0x01</math></li> <li><math>0x69 \&amp;&amp; 0x55 \rightarrow 0x01</math></li> <li><math>0x69     0x55 \rightarrow 0x01</math></li> </ul> </li> </ul>	

Shift Operations

<ul style="list-style-type: none"> <li>Left shift: <math>x \ll y</math> <ul style="list-style-type: none"> <li>Shift bit-vector <math>x</math> left <math>y</math> positions           <ul style="list-style-type: none"> <li>Throw away extra bits on left</li> <li>Fill with 0's on right</li> </ul> </li> </ul> </li> <li>Right shift: <math>x \gg y</math> <ul style="list-style-type: none"> <li>Shift bit-vector <math>x</math> right <math>y</math> positions           <ul style="list-style-type: none"> <li>Throw away extra bits on right</li> <li>Logical shift               <ul style="list-style-type: none"> <li>Fill with 0's on left</li> </ul> </li> <li>Arithmetic shift               <ul style="list-style-type: none"> <li>Replicate most significant bit on right</li> </ul> </li> </ul> </li> <li><b>Undefined behavior</b> <ul style="list-style-type: none"> <li>Shift amount <math>&lt; 0</math> or <math>\geq</math> word size</li> </ul> </li> </ul> </li></ul>	<table border="1"> <tr> <td>Argument x</td> <td>01100010</td> </tr> <tr> <td><math>\ll 3</math></td> <td>00010<b>0000</b></td> </tr> <tr> <td>Log. <math>\gg 2</math></td> <td><b>000</b>11000</td> </tr> <tr> <td>Arith. <math>\gg 2</math></td> <td><b>000</b>11000</td> </tr> </table> <table border="1"> <tr> <td>Argument x</td> <td>10100010</td> </tr> <tr> <td><math>\ll 3</math></td> <td>00010<b>0000</b></td> </tr> <tr> <td>Log. <math>\gg 2</math></td> <td><b>001</b>01000</td> </tr> <tr> <td>Arith. <math>\gg 2</math></td> <td><b>11</b>101000</td> </tr> </table>	Argument x	01100010	$\ll 3$	00010 <b>0000</b>	Log. $\gg 2$	<b>000</b> 11000	Arith. $\gg 2$	<b>000</b> 11000	Argument x	10100010	$\ll 3$	00010 <b>0000</b>	Log. $\gg 2$	<b>001</b> 01000	Arith. $\gg 2$	<b>11</b> 101000
Argument x	01100010																
$\ll 3$	00010 <b>0000</b>																
Log. $\gg 2$	<b>000</b> 11000																
Arith. $\gg 2$	<b>000</b> 11000																
Argument x	10100010																
$\ll 3$	00010 <b>0000</b>																
Log. $\gg 2$	<b>001</b> 01000																
Arith. $\gg 2$	<b>11</b> 101000																

Java writes this ">>>".

In C we do not have different operators for logical and arithmetic shift, if we have a signed int then we do an arithmetic shift, and with unsigned int logical shift

## 3.2 Integer Ranges

### 3.2.1 Encoding Integers

Unsigned		Two's complement													
$B2U(X) = \sum_{i=0}^{w-1} x_i \cdot 2^i$		$B2T(X) = -x_{w-1} \cdot 2^{w-1} + \sum_{i=0}^{w-2} x_i \cdot 2^i$													
short int x = 15213; short int y = -15213;		Sign Bit													
• A C short is 2 bytes long:															
<table border="1"> <thead> <tr> <th></th> <th>Decimal</th> <th>Hex</th> <th>Binary</th> </tr> </thead> <tbody> <tr> <td>x</td> <td>15213</td> <td>3B 6D</td> <td>00111011 01101101</td> </tr> <tr> <td>y</td> <td>-15213</td> <td>C4 93</td> <td>11000100 10010011</td> </tr> </tbody> </table>					Decimal	Hex	Binary	x	15213	3B 6D	00111011 01101101	y	-15213	C4 93	11000100 10010011
	Decimal	Hex	Binary												
x	15213	3B 6D	00111011 01101101												
y	-15213	C4 93	11000100 10010011												
• Sign bit <ul style="list-style-type: none"> <li>For 2's complement, most significant bit = 1 indicates negative</li> </ul>															

### 3.2.2 Numeric ranges

<ul style="list-style-type: none"> <li>Unsigned values           <ul style="list-style-type: none"> <li>UMin = 0               <ul style="list-style-type: none"> <li>000...0</li> </ul> </li> <li>UMax = <math>2^w - 1</math> <ul style="list-style-type: none"> <li>111...1</li> </ul> </li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>Two's complement values           <ul style="list-style-type: none"> <li>TMin = <math>-2^{w-1}</math> <ul style="list-style-type: none"> <li>100...0</li> </ul> </li> <li>TMax = <math>2^{w-1} - 1</math> <ul style="list-style-type: none"> <li>011...1</li> </ul> </li> </ul> </li> </ul>	<table border="1"> <thead> <tr> <th></th> <th>8</th> <th>16</th> <th>32</th> <th>64</th> </tr> </thead> <tbody> <tr> <td>UMax</td> <td>255</td> <td>65,535</td> <td>4,294,967,295</td> <td>18,446,744,073,709,551,615</td> </tr> <tr> <td>TMax</td> <td>127</td> <td>32,767</td> <td>2,147,483,647</td> <td>9,223,372,036,854,775,807</td> </tr> <tr> <td>TMin</td> <td>-128</td> <td>-32,768</td> <td>-2,147,483,648</td> <td>-9,223,372,036,854,775,808</td> </tr> </tbody> </table>		8	16	32	64	UMax	255	65,535	4,294,967,295	18,446,744,073,709,551,615	TMax	127	32,767	2,147,483,647	9,223,372,036,854,775,807	TMin	-128	-32,768	-2,147,483,648	-9,223,372,036,854,775,808	<ul style="list-style-type: none"> <li>Observations:           <ul style="list-style-type: none"> <li><math> TMin  = TMax + 1</math> <ul style="list-style-type: none"> <li>Asymmetric range</li> </ul> </li> <li>UMax = <math>2 * TMax + 1</math></li> </ul> </li> <li>C Programming           <ul style="list-style-type: none"> <li>#include &lt;limits.h&gt;</li> <li>Declares constants, e.g.,               <ul style="list-style-type: none"> <li>ULONG_MAX</li> <li>LONG_MAX</li> <li>LONG_MIN</li> </ul> </li> <li>Values platform specific</li> </ul> </li> </ul>
	8	16	32	64																			
UMax	255	65,535	4,294,967,295	18,446,744,073,709,551,615																			
TMax	127	32,767	2,147,483,647	9,223,372,036,854,775,807																			
TMin	-128	-32,768	-2,147,483,648	-9,223,372,036,854,775,808																			

### 3.2.3 Signed vs Unsigned in C

By default constants are considered to be signed Integers, they are unsigned if they have "U" as a suffix. The bit composition in explicit casting remains the same hence a large unsigned value could be a negative signed value and vice versa. Implicit casting can also occur e.g during assignments or function calls.

When mixing unsigned and signed values in a single expression, signed values are implicitly cast to unsigned (including comparison operators!)

Constant 1	Constant 2	Relation	Evaluation
0	0U	==	Unsigned
-1	0	<	Signed
-1	0U	>	Unsigned
2147483647	-2147483647-1	>	Signed
2147483647U	-2147483647-1	<	Unsigned
-1	-2	>	Signed
(unsigned)-1	-2	>	Unsigned
2147483647	2147483648U	<	Unsigned
2147483647	(int) 2147483648U	>	signed

### 3.2.4 Sign extension

- Task:
  - Given  $w$ -bit signed integer  $x$
  - Convert it to  $w+k$ -bit integer with same value
- Rule:
  - Make  $k$  copies of sign bit:
  - $X' = X_{w-1}, \dots, X_0, X_{w-1}, X_{w-2}, \dots, X_0$

Converting from smaller to larger integer data type  
C automatically performs sign extension for signed values

	Decimal	Hex	Binary
x	15213	3B 6D	00111011 01101101
ix	15213	00 00 3B 6D	00000000 00000000 00111011 01101101
y	-15213	C4 93	11000100 10010011
iy	-15213	FF FF C4 93	11111111 11111111 11000100 10010011

## 3.3 Integer addition and subtraction in C

### 3.3.1 Unsigned addition

Operands:  $w$  bits       $u$        $\boxed{\dots}$   
                                    $+ v$        $\boxed{\dots}$   
                                    $\underline{u + v}$        $\boxed{\dots}$   
                                   Discard carry:  $w$  bits       $UAdd_w(u, v)$        $\boxed{\dots}$

- Standard addition function
  - Ignores carry output
- Implements modular arithmetic

$$s = UAdd_w(u, v) = u + v \bmod 2^w$$

$$UAdd_w(u, v) = \begin{cases} u + v, & u + v < 2^w \\ u + v - 2^w, & u + v \geq 2^w \end{cases}$$

The sum wraps around if the true sum  $\geq 2^w$ . The properties are:

- Modular addition forms an Abelian group
  - Closed** under addition  
 $0 \leq UAdd_w(u, v) \leq 2^w - 1$
  - Commutative**  
 $UAdd_w(u, v) = UAdd_w(v, u)$
  - Associative**  
 $UAdd_w(t, UAdd_w(u, v)) = UAdd_w(UAdd_w(t, u), v)$
  - 0** is additive **identity**  
 $UAdd_w(u, 0) = u$
  - Every element has additive **inverse**  
Let:  $UComp_w(u) = 2^w - u$   
Then:  $UAdd_w(u, UComp_w(u)) = 0$

this forms a different abelian group than the normal addition

### 3.3.2 Two's complement addition

Operands:  $w$  bits       $u$        $\boxed{\dots}$   
                                    $+ v$        $\boxed{\dots}$   
                                    $\underline{u + v}$        $\boxed{\dots}$   
                                   Discard carry:  $w$  bits       $TAdd_w(u, v)$        $\boxed{\dots}$

- TAdd and UAdd have identical bit-level behavior
  - Signed vs. unsigned addition in C:

```
int s, t, u, v;
s = (int)((unsigned)u + (unsigned)v);
t = u + v;
```

will give:  $s == t$

Hence the only difference between Uadd and Tadd is how the result is interpreted. However the overflow is different:

- Functionality
  - True sum requires  $w+1$  bits
  - Drop off MSB
  - Treat remaining bits as 2's comp. integer

True Sum

$$TAdd_w(u, v) = \begin{cases} u + v + 2^w & u + v < TMin_w \\ u + v & TMin_w \leq u + v \leq TMax_w \\ u + v - 2^w & TMax_w < u + v \end{cases}$$

(Neg. overflow)  
(Pos. overflow)

The mathemimatical properties of Tadd:

- Group isomorphic to unsigneds with UAdd

- Since both have identical bit patterns

$$TAdd_w(u, v) = U2T \left( UAdd_w(T2U(u), T2U(v)) \right)$$

- 2's complement under TAdd forms a group

- Closed, commutative, associative,
  - 0 is additive identity
  - Every element has additive inverse

$$TComp_w(u) = \begin{cases} -u & u \neq TMin_w \\ TMin_w & u = TMin_w \end{cases}$$

### 3.3.3 Integer multiplication in C

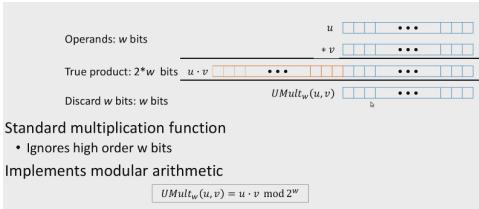
Multiplication is worse than addition, because the number of extra bits which could be needed is much larger

- Computing exact product of w-bit numbers x, y
    - Either signed or unsigned
  - Ranges
    - Unsigned (up to  $2^w$  bits):  

$$0 \leq x * y \leq (2^w - 1)2 = 2^{2w} - 2^{w+1} + 1$$
    - Two's complement min (up to  $2^{w-1}$  bits):  

$$x * y \geq (-2^{w-1}) * (2^{w-1} - 1) = -2^{2w-2} + 2^{w-1}$$
    - Two's complement max (up to  $2^w$  bits, but only for  $(TM_{in_w})^2$ ):  

$$x * y \leq (-2^{w-1})^2 = 2^{w-2}$$
  - Maintaining exact results
    - Would need to keep expanding word size with each product computed
    - Done in software by "arbitrary precision" arithmetic packages



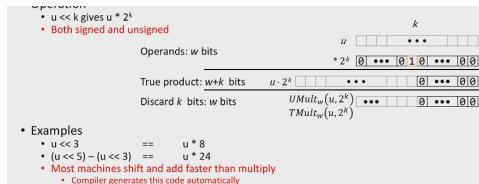
Unsigned multiplication with addition forms a commutative ring. Signed multiplication is isomorphic algebra to unsigned multiplication and addition. Difference to normal mathematics is that there is no guarantee that the addition of two positive numbers signed numbers is also positive.

### 3.4 Integer multiplication and division using shifts

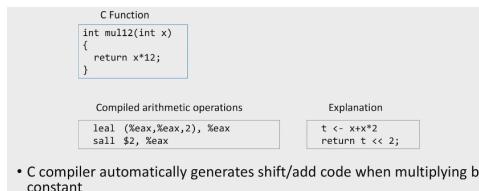
### 3.4.1 Power of 2 multiply with shift

For both signed and unsigned numbers we have:

$u << k$  gives  $u \cdot 2^k$



Shifting and addition/subtraction are quite simple operations. On the other hand multiplication of random numbers are more expensive. But usually not necessary to write the multiplication explicitly because the compiler writes the multiplications as shifts and additions/subtractions:



### 3.4.2 Unsigned power of 2 divide with shift

Same as multiplication but right logical shift where the bits after the binary are truncated. Hence

$u \gg k$  gives  $\lfloor \frac{u}{2k} \rfloor$

The compiler will rewrite division into shifts and additions as well.

### 3.4.3 Signed power of 2 divide with shift



The rounding we want is toward 0. The correct power of 2 division is given by:

• Quotient of negative number by power of 2

- We want  $\lceil s/2^k \rceil$  (round toward 0)
- We compute it as  $\lceil (s + 2^k - 1)/2^k \rceil$
- In C:  $(s + (1< $k$ )-1)>>k$
- Biases the dividend toward 0

• Case 1: No rounding

Dividend: 
$$\begin{array}{r} s \\ \hline +2^k - 1 \\ \hline \end{array}$$

Blasing has no effect

Divisor: 
$$\begin{array}{r} s \\ \hline /2^k \\ s/2^k \end{array}$$

• Case 2: Rounding:

Dividend: 
$$\begin{array}{r} s \\ \hline +2^k - 1 \\ \hline \end{array}$$

Blasing adds 1 to final result

Divisor: 
$$\begin{array}{r} s \\ \hline /2^k \\ s/2^k \end{array}$$

When doing a signed division the compiler will do a check if the number is less than zero, if so then it will add the bias and proceed.

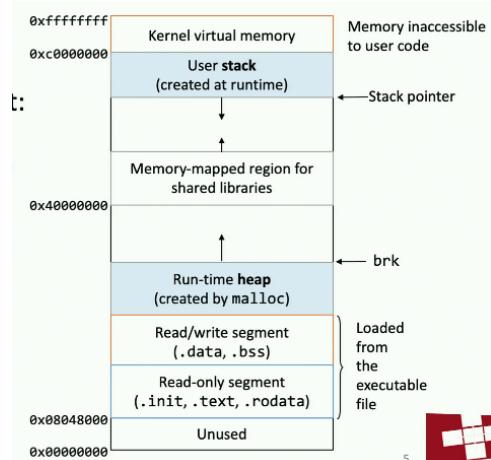
# Chapter 4

## Pointers

Pointers are variables which store a memory address. The Operating system gives each process an address space. The address space contains the virtual memory which is only visible by the process. Memory is byte addressable.

### 4.1 Loading

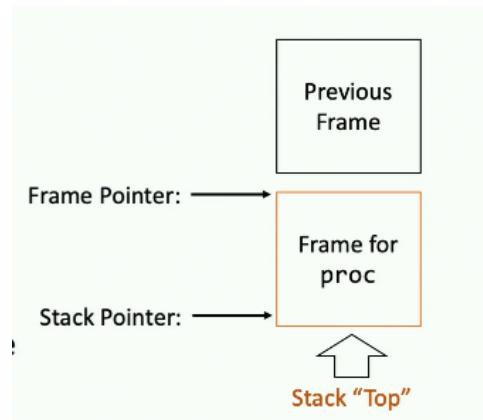
The file is C program is compiled, and we get an executable file, which is what we are actually trying to run. The OS needs to load this file, hence it inspects what's inside, and will load different parts of the program into different parts of the virtual address space.



### 4.2 Stack

The Stack is created at runtime, everytime a procedure is called the stack grows, when a procedure completes the stack shrinks, the stack grows downwards. We keep track where the stack is in memory with the stack pointer.

#### 4.2.1 Stack Frames

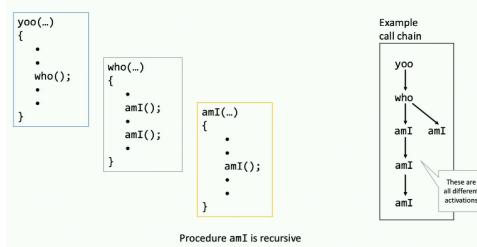


The unit of allocation in a stack is a stack frame. For each instance of a procedure we get a stack frame. The following information is in a stack frame:

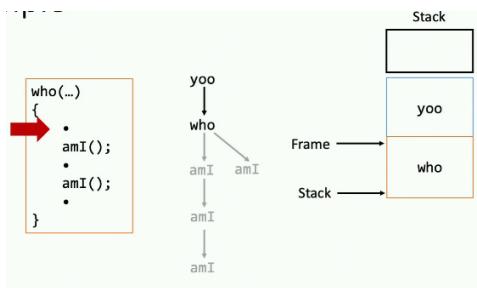
- Local variables
- Return information
- Temporary space

#### 4.2.2 Call Chain example:

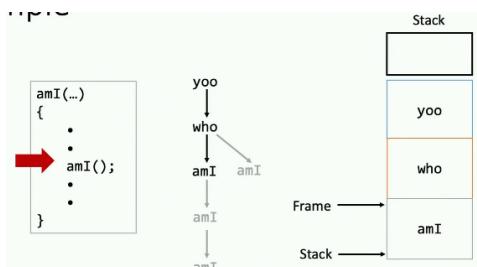
Given the following call chain:



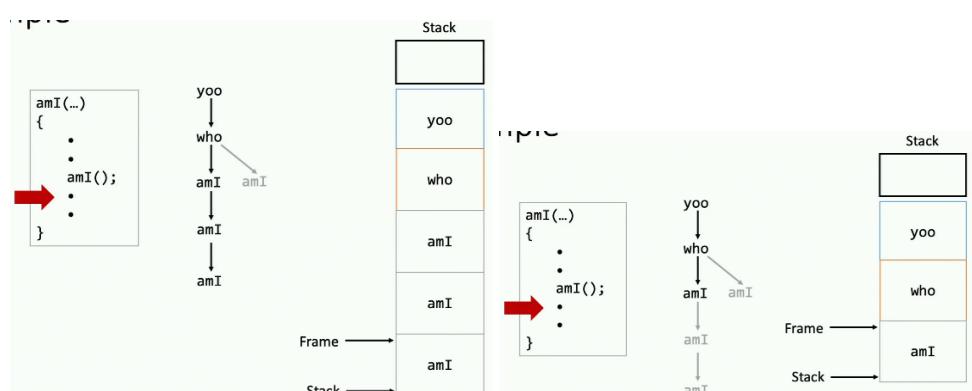
We start with an empty stack, yoo is called first, hence we have a stack frame for yoo. yoo then calls who so we add a stack frame for who on top of the stack. In general whichever procedure you're currently in, that stack frame will be on top of the stack.



The Frame pointer keeps track of the beginning of the Frame, the stack pointer the end of the frame. When we enter the procedure amI, we create a new frame and hence adjust the two pointers.



For each recursive call of amI we add a stack frame, and when they return we pop the stack.



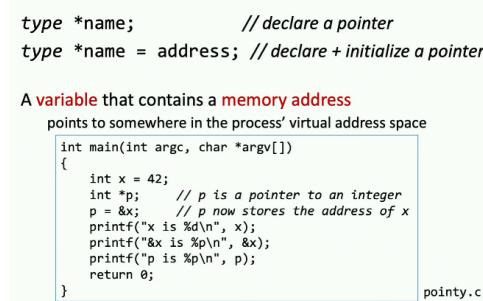
## 4.3 Pointers

### 4.3.1 Addresses and &

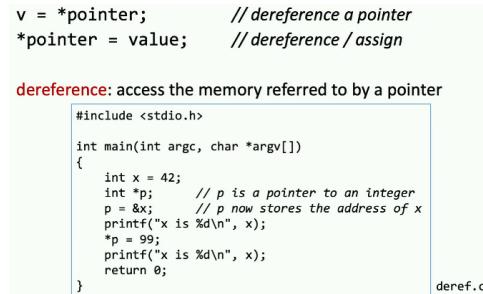
`&` is the reference operator. Given a variable `x`, `&x` produces the virtual address where the value of `x` is stored. We can use `%p` in the `printf()` method to print the address of `x` out.

### 4.3.2 Pointers

Pointers are variables which store memory addresses.



Pointers can be dereferenced which means that we get the value stored at the memory address of the pointer.

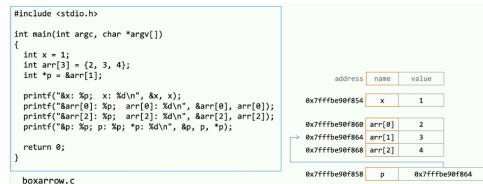


A pointer to a pointer (**double pointer**) can be initialized as follows  $\ast\ast dp = \&p$  where p is a pointer.



### 4.3.3 Box and arrow diagrams

A way to visualize pointers in a program.



#### 4.3.4 Address Space Layout Randomization

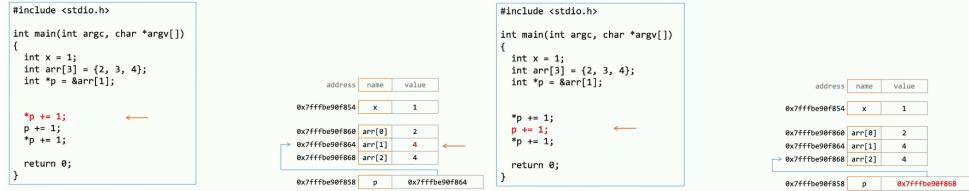
The Base of the stack will be started at a random address by the OS, hence everytime we run a programm the address of the variables will have a different address. This is done as a security measure.

### 4.3.5 NULL

A memory location guaranteed to be invalid. Null is of type "void \*". Any attempt to dereference Null results in a segmentation fault. This is useful to dynamically allocate memory.

## 4.4 Pointer arithmetic

When incrementing the pointer by one the address it points to after depends on the type, in the example below the pointer is an integer pointer, hence we increment by the size of an integer

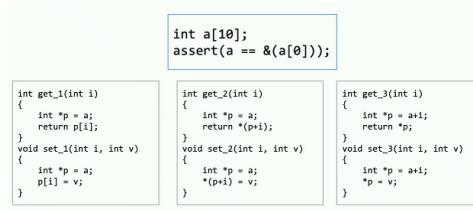


The amount of memory a value takes up depends on the machine and compiler. We can use the methods:

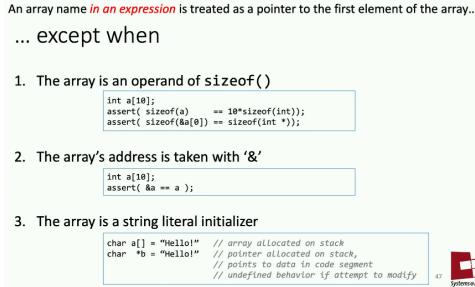
`sizeof(type)` or `sizeof(value)`

## 4.5 Arrays and Pointers

An array name in an expression is treated as a pointer to the first element of the array

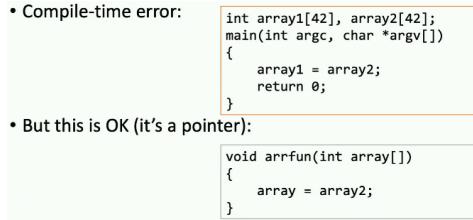


Arrays are not pointers. An array is a collection of homogeneous data elements stored at contiguous memory addresses. A pointer is a variable that stores a memory address



In the compiler  $A[i]$  is always rewritten as  $*(A + i)$

An array name as a function parameter is a pointers. Arrays cannot be renamed:

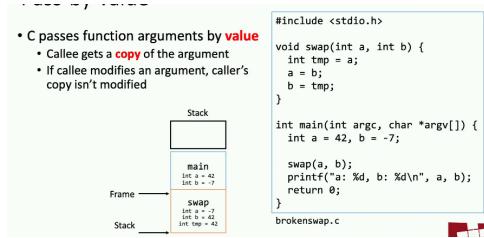


The top one doesn't work, because arrays are not pointers. The bottom one works because in the method array is a parameter hence it's a pointer to the first element and the assignment will work. array2 will be treated as the address of the first element of the array.

## 4.6 Passing by Reference

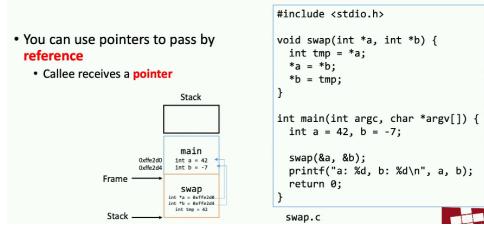
### 4.6.1 Pass by value

By default in C variables are passed by value. That means a copy of the arguments are made.



When we return swap the frame is popped from the stack, but the values in main have not been modified.

## 4.7 Pass by Reference



## 4.8 C Pointer Declarations

### C Pointer Declarations

int *p	p is a pointer to int
int *p[13]	p is an array[13] of pointer to int
int *(p[13])	p is an array[13] of pointer to int
int **p	p is a pointer to a pointer to an int
int (*p)[13]	p is a pointer to an array[13] of int
int *f()	f is a function returning a pointer to int
int (*f)()	f is a pointer to a function returning int
int (*(*f())[13])()	f is a function returning ptr to an array[13] of pointers to functions returning int
int (*(*x[3]))()[5]	x is an array[3] of pointers to functions returning pointers to array[5] of ints

Systems Programming 2020 Ch. 4: Pointers

## Chapter 5

# Dynamic memory allocation

**static memory allocation:** Memory is allocated when the program is loaded and deallocated when the program exits.  
(global variables)

```
// a global variable
int counter = 0;

int main(int argc, char *argv[])
{
    counter++;
    return 0;
}
```

**automatic memory allocation:** allocated when a function is called and deallocated when function returns. We write C code, this gets compiled into an executable. When the compiler sees procedure calls in the code, in assembly it creates instructions which include push and pop instructions which modify the stack pointer which create and deallocate stack frames.

```
int foo(int a)
{
    int x = a + 1;
    return x
}

int main(int argc, char *argv[])
{
    int y = foo(10);
    return 0;
}
```

We want memory that is more dynamic e.g persists across multiple function calls, memory too big to fit on the stack, memory allocation for functions with an unknown return size etc.  
The stack size is usually 8MB.

**Dynamically allocated memory:** The program explicitly requests a new block of memory. C requires manual memory management in contrast to other languages which have automatic memory management (garbage collector). In manual memory management the code explicitly deallocates memory. This gives us much more control over memory but requires much more care.

## 5.1 The C memory API

### malloc()

```
// declared in stdlib.h
typedef unsigned long size_t;
void *malloc(size_t sz);
```

- Allocates a block of memory of the given size
  - returns a pointer to the first byte of that memory
  - returns NULL if the memory cannot be allocated
- you should **assume** the memory initially contains garbage
- typically use sizeof() to calculate the size you need

```
long *arr = (long *)malloc(10*sizeof(long));
if (arr == NULL) {
    return ERRCODE;
}
arr[0] = 5L; // etc.
```

The return of malloc is a pointer to the beginning of the memory allocation. If you request more memory than is available malloc will return null. malloc returns a void pointer so that we can cast it to the appropriate type later. malloc has no guarantee to how the memory is initialized.

The method calloc() zeros out the requested memory

### calloc()

```
// declared in stdlib.h
void *calloc(size_t nm, size_t sz);
```

- Allocates a block of memory size ( $nm \times sz$ )
  - Like malloc(), a pointer to the first byte of that memory
  - returns NULL if the memory cannot be allocated
- zeros** the memory - unlike malloc()
- Slightly slower, but less error-prone, more readable

```
long *arr = (long *)calloc(10, sizeof(long));
if (arr == NULL) {
    return ERRCODE;
}
arr[0] = 5L; // etc.
```

#### 5.1.1 Deallocation

Deallocation is done with the free method. For each malloc/calloc there must be an appropriate free call

```
// declared in stdlib.h
void free(void *);
```

- Releases memory at the pointer
  - Must point to the first byte of malloc'ed memory
  - After freeing, memory might be reallocated by some future malloc()/calloc()
  - Good practice to NULL the pointer after freeing.

```
long *arr = (long *)calloc(10, sizeof(long));
if (arr == NULL) {
    return ERRCODE;
}
// Do something ...
free(arr);
arr = NULL;
```

Its good practice to null the pointer after freeing it, because after you free it, its not pointing to anything meaningful anymore. Dereferencing it will then have a clear error.

#### 5.1.2 Resize

### realloc(): changing the size of the block

```
// declared in stdlib.h
void *realloc(void *ptr, size_t size);
```

- Allocations have a **fixed size**.
- Like free(), realloc() must point to first byte of a malloc'ed block
  - Might (almost certainly will!) **copy** the data to a new location
- ⇒ **Always** use the new address returned!
- Returns NULL if failed.

```
long *arr;
if (!(arr = (long *)malloc(10*sizeof(long)))) {
    return ERRCODE;
}
// ... do something ...
if (!(arr = (long *)realloc(arr, 20*sizeof(long)))) {
    return ERRCODE;
}
```

realloc does not guarantee how the memory region is initialized.

### 5.1.3 size\_t

- **size\_t**:
  - An **unsigned** integer of some size
  - Return type of `sizeof()`
  - Large enough to hold the size of the largest possible array in memory  
⇒ can be used to store a pointer
- **ptrdiff\_t**:
  - An unsigned integer of some size
  - Result of subtracting two pointers
  - Used for array loops, size calculations, etc.
  - Ends up identical to `size_t`

### 5.1.4 Summary

```
/* This program simply reads integers
   into a dynamic array until eof.
   The array is expanded as needed */

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

// Initial array size
#define INIT_SIZE 8

int main(int argc, char *argv[])
{
    int num; // Num. of integers
    int *arr; // Array of ints.
    size_t sz; // Array size
    int m, in; // Index & input num

    // Allocate the initial space.
    sz = INIT_SIZE;
    arr = (int *)calloc(sz, sizeof(int));
    if (arr == NULL) {
        fprintf(stderr, "calloc failed.\n");
        exit(1);
    }

    // Read in the numbers.
    num = 0;
    while (scanf("%d", &in) == 1) {
        // See if there's room.
        if (num == sz) {
            // There's not. Get more.
            sz *= 2;
            arr = (int *)realloc(arr, sz * sizeof(int));
            if (arr == NULL) {
                fprintf(stderr, "realloc failed.\n");
                exit(1);
            }
        }

        // Store the number.
        arr[num++] = in;
    }

    // Print out the numbers
    for (m = 0; m < num; m++)
        printf("%d\n", arr[m]);
    free(arr); // Always good to free
    return 0;
}
```

```
// Read in the numbers.

num = 0;
while (scanf("%d", &in) == 1) {
    // See if there's room.
    if (num == sz) {
        // There's not. Get more.
        sz *= 2;
        arr = (int *)realloc(arr, sz * sizeof(int));
        if (arr == NULL) {
            fprintf(stderr, "realloc failed.\n");
            exit(1);
        }
    }

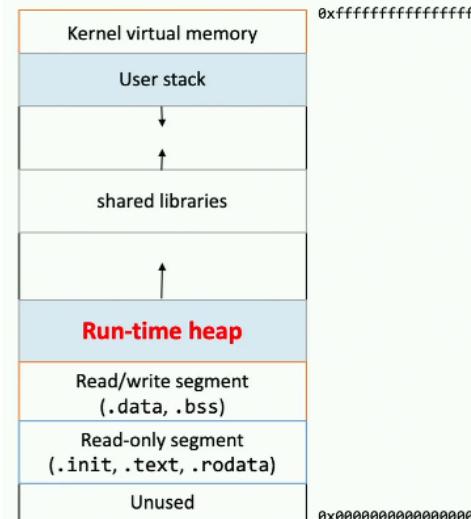
    // Store the number.
    arr[num++] = in;
}

// Print out the numbers
for (m = 0; m < num; m++)
    printf("%d\n", arr[m]);
free(arr); // Always good to free
return 0;
}
```

the method `scanf` is used to get user input. It blocks the program until the input is given.

## 5.2 Managing the heap

The heap ("free store") is a large pool of unused memory, used for dynamically allocated data structures. `malloc()` allocates chunks of memory in the heap, and maintains bookkeeping data in the heap to track allocated blocks. `free` returns the allocated chunks.



### 5.2.1 Memory Corruption

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

int main(int argc, char **argv) {
    int a[2];
    int *b = malloc(2 * sizeof(int)), *c;

    a[2] = 5; // assign past the end of an array
    a[0] += 2; // assume malloc zeroes out memory
    c = b+3; // mess up your pointer arithmetic
    free(&(a[0])); // pass pointer to free() that wasn't malloc'ed
    free(b);
    free(b); // double-free the same block
    b[0] = 5; // use a free()'d pointer

    // any more!
    return 0;
}
```

## 5.2.2 Memory Leaks

A memory leak happens when code fails to deallocate memory that will no longer be used.

```
// assume we have functions FileLen,
// ReadFileIntoBuffer, and NumWordsInString

int NumWordsInFile(char *filename)
{
    char *filebuf = (char *)malloc(FileLen(filename)+1);
    if (filebuf == NULL) {
        return -1;
    }

    ReadFileIntoBuffer(filename, filebuf);

    // Leak! we never free(filebuf);
    return NumWordsInString(filebuf);
}
```

**Memory footprint:** The amount of memory your program uses implications of a leak: For short-lived programs this might be ok, but for long-lived programs this is usually bad and might:

- slow down over time due to VM thrashing
- use up all available memory and crash
- starve other programs out of memory

## 5.3 Structures and unions

### 5.3.1 Structured data

Allows us to have a collection of data of the same type in memory. A bit like a class but there are no methods or constructors.

The diagram illustrates the definition of a struct `Point` and its use in a `main` function. On the left, a box contains the struct definition:

```
struct typename {
    type name;
    type name;
    ...
};
```

Below it, another box shows the actual struct definition:

```
// New structured data type called "struct Point"
struct Point {
    int x;
    int y;
};

struct Point origin = { 0, 0 };
```

To the right, a larger box contains the `main` function code:

```
struct Point {
    int x, y;
};

int main(int argc, char **argv) {
    struct Point p1 = {0, 0}; // p1 on the stack
    struct Point *p1_ptr = &p1;

    p1.x = 1;
    p1_ptr->y = 2;
    return 0;
}
```

A callout box on the right points to the `p1_ptr` assignment in the `main` function, stating: `p->x` is nicer than `(*p).x`.

you can assign the value of a struct from a struct of the same type. This copies the entire contents (unlike for arrays)

```
#include <stdio.h>

struct Point {
    int x, y;
};

int main(int argc, char **argv) {
    struct Point p1 = {0, 2};
    struct Point p2 = {4, 6};

    printf("p1: %d,%d  p2: %d,%d\n", p1.x, p1.y, p2.x, p2.y);
    p2 = p1;
    printf("p1: %d,%d  p2: %d,%d\n", p1.x, p1.y, p2.x, p2.y);

    return 0;
}
```

Even if the struct has an array in it the contents will be copied over by the assignment operator. This is because the copy operations copies the region of memory for that struct and doesn't care about the types.

Structs can be passed as arguments, by default they are passed by value, to pass by reference we pass a pointer to the struct. We can also return structs (right picture)

```

#include <stdio.h>

struct point {
    int x, y;
};

void DoubleXBroken(struct point p) {
    p.x *= 2;
}

void DoubleXWorks(struct point *p) {
    p->x *= 2;
}

int main(int argc, char *argv) {
    struct point a = {1,1};

    DoubleXBroken(a);
    printf("(%d,%d)\n", a.x, a.y);

    DoubleXWorks(&a);
    printf("(%d,%d)\n", a.x, a.y);

    return 0;
}

```

```

// a complex number is a + bi
struct complex {
    double real; // real component
    double imag; // imaginary component
};

struct complex AddComplex(struct complex x, y)
{
    struct complex retval;

    retval.real = x.real + y.real;
    retval.imag = x.imag + y.imag;
    return retval;
}

struct complex MultiplyComplex(struct complex x, y)
{
    struct complex retval;

    retval.real = (x.real * y.real) - (x.imag * y.imag);
    retval.imag = (x.imag * y.real) - (x.real * y.imag);
    return retval;
}

```

### 5.3.2 Unions

Like a struct but only allocates enough memory to hold the largest member.

```

union u {
    int ival;
    float fval;
    char *sval;
};
union u my_uval;

```

Only use a union if you're interested in only one of the members.

## 5.4 Type definitions

### 5.4.1 typedef

Introduces a new type definition i.e a new name for a type, hence it can be used wherever the original type could be.

```

typedef unsigned uint32_t;
uint32_t ui;
...
typedef int **myptr;
int *p;
myptr mp = &p;
...
typedef struct skbuf skbuf_t;
skbuf_t *sptr;

```

This helps us avoiding complex declarations e.g:

- Use `typedef` to build up the declaration
- Instead of `int (*(*x[3]))()[5]` :

```
typedef int fiveints[5];
typedef fiveints* p5i;
typedef p5i (*f_of_p5is)();
f_of_p5is x[3];
```

x is an array of 3 elements, each of which is a pointer to a function returning an array of 5 ints

### 5.4.2 Struct tags and `typedef`s

Struct tags are the names that we give the structs e.g point in struct `point{...}`.

Struct tags and <code>typedef</code> s	<code>• You can:</code>
	<pre>typedef struct list_el el_t;</pre>
Confusing.	<code>• Or even:</code>
Better to stay with always using tags.	<pre>typedef struct list_el {     unsigned long val;     struct list_el *next; } el_t;  struct list_el my_list; el_t my_other_list;</pre> <code>• Or even:</code> <pre>typedef struct list_el {     unsigned long val;     struct list_el *next; } list_el;  struct list_el my_list; list_el my_other_list;</pre>

### 5.4.3 C namespaces

There are 4 namespaces in C:

- Label names (`goto...`)
- Tags (one namespace for all struct, union,enum's)
- Member names (one namespace for each struct,union,enum)
- Everything else (including `typedef`)

```
// structure tag
struct id {
    int id; // structure member
};

// different structure
struct id2 {
    char id; // structure member
};

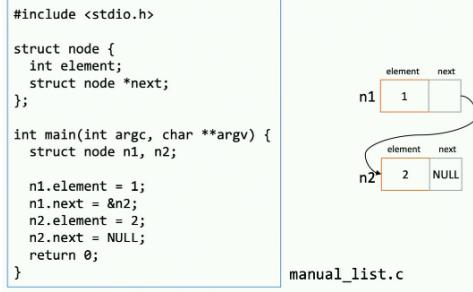
// ordinary identifier
void id() {
    goto id;

    id: // label
}
```

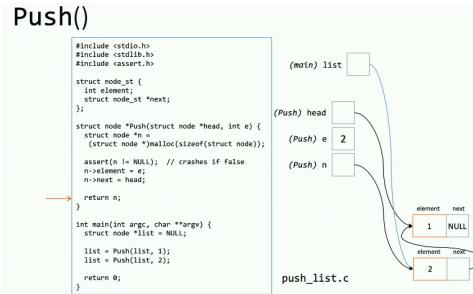
## 5.5 Dynamic Data Structures

### 5.5.1 Singly-Linked List

A list consisting of nodes. Each node in the list contains some element as its payload and a pointer to the next node in the linked list. The last node in the list contains `NULL` instead of the next pointer. We can implement this using structs. We represent each node as a struct.



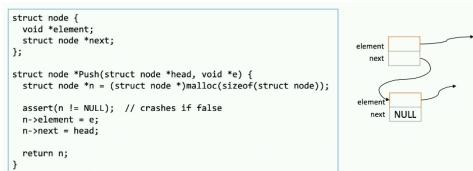
We create a method push() which dynamically adds nodes to the list.



We notice that the above implementation of push() has a leak as we have not deallocated the malloc(). We can check if the program has a leak in the terminal with (assuming our file is called *push\_list.c*)

```
$ gcc -o push_list -g -Wall push_list.c
$ valgrind -leak-check=full ./push_list.c
```

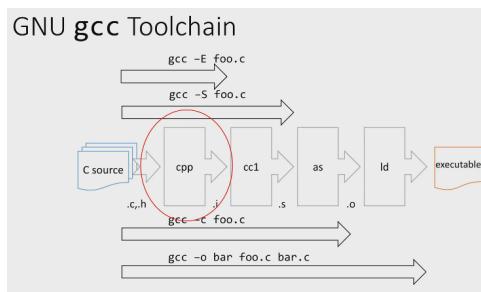
We can create a generic linked list by using "void \*\*" as a type.



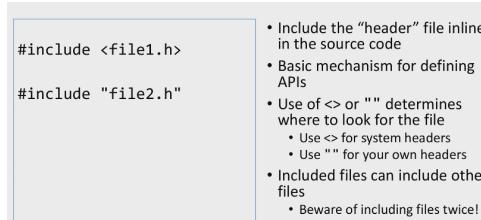
# Chapter 6

## Wrapping up C

### 6.1 The C Preprocessor

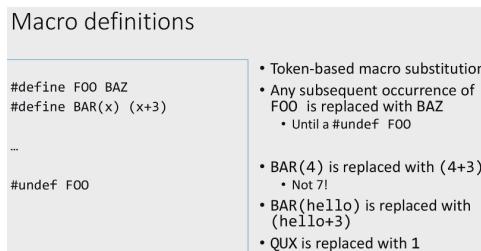


The cpp is not part of the C language but rather an extra language that sits on top.

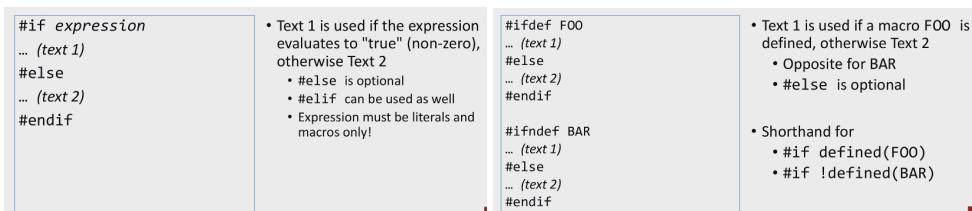


The cpp will go look for the files in the include statement. Hence using <, > or "", will tell the cpp where to start looking. Since included files can include other files we could get a cyclic dependance.

#### 6.1.1 Macro Definitions



C macros purely work on text. The default definition of any macro is 1. Preprocessor conditionals:



In C we can use a backslash to continue the code on the next line. This is used for readability

- Macro expanding to a compound statement:

```
#define SKIP_SPACES(p, limit) \
{ char *lim = (limit); \
  while (p < lim) { \
    if (*p++ != ' ') { \
      p--; break; }}} \
```

### Solution: swallow the semicolon

- Make macro expand to single statement:

```
#define SKIP_SPACES(p, limit) \
do { char *lim = (limit); \
  while (p < lim) { \
    if (*p++ != ' ') { \
      p--; break; }}} \
```

- Problem: a semicolon is a null statement:

```
if (*p != 0)
  SKIP_SPACES (p, my_limit);
else ...
```

Invalid C  
(expand it  
and see!)

- Previous example expands to:

```
if (*p != 0)
  do { .. } while(0);
else ...
```

See: <http://gcc.gnu.org/onlinedocs/cpp/Swallowing-the-Semicolon.html>

## 6.2 Modularity

In C there is a difference between Declarations and Definitions:

- A **declaration** says something exists, somewhere:

```
char *strncpy(char *dest, const char *src, size_t n);
```

A "prototype"

- A **definition** says what it is:

```
char *strncpy(char *dest, const char *src, size_t n)
{
  size_t i;
  for (i = 0; i < n && src[i] != '\0'; i++) {
    dest[i] = src[i];
  }
  for ( ; i < n; i++) {
    dest[i] = '\0';
  }
  return dest;
}
```

**Compilation Units:** What gets fed into the C compiler (cc1) i.e compilation unit is what comes out of the cpp meaning the c file after all the substitutions have been made.

Declarations can be annotated with:

- **static:** definition (also static) is in this compilation unit, and cant be seen outside it
- **extern:** definition is somewhere else, either in this compilation unit or another.

### Global variables are also declared

Declarations:

```
extern const char *banner; // Defined in some other file
static int priv_count; // Defined in this file
```

Elsewhere in this compilation unit:

```
static int priv_count = 0; // Only in scope in this unit
```

In some other compilation unit:

```
const char *banner = "Welcome to Barreelfish";
```

A Module is a self contained piece of a larger program. It consists of:

- Externally visible:
  - functions to be invoked
  - typedefs and perhaps global variables
  - cpp macros
- Internal functions types, global variables that clients shouldnt look at

Modularity in C is done using header files:

## C header files

- Specify interfaces with header files
- Module **foo** has interface in **foo.h**
  - Clients of foo **#include "foo.h"**
  - foo.h contains **no definitions**
    - only external declarations
- Implementation typically in **foo.c**
  - Also includes foo.h
  - Contains **no external declarations**
    - only definitions and internal declarations

There is cpp idiom for header files to avoid a cyclic dependency:

```
"file.h":  
#ifndef __FILE_H__  
#define __FILE_H__  
  
... (declarations, macros)  
  
#endif // __FILE_H__
```

- cpp boilerplate ensures file contents only appear once
- Never #include a .c file!

the **#ifndef** ensures that the body will only be used once in the file.

Example: Linked List:

```
#include "ll.h"  
#include <stdio.h>  
#include <stdlib.h>  
#include <assert.h>  
  
struct node *Push(struct node *head, void *element) {  
    ...  
}  
ll.c  
struct node {  
    void *element;  
    struct node *next;  
};  
  
extern struct node *Push(struct node *head, void *element);  
ll.h
```

```
#include "ll.h"  
  
int main(int argc, char **argv) {  
    struct node *list = NULL;  
    char *hi = "hello";  
    char *bye = "goodbye";  
  
    list = Push(list, (void *) hi);  
    list = Push(list, (void *) bye);  
    return 0;  
}  
example_ll_customer.c
```

## 6.3 Function Pointers

A pointer with the address of some code which corresponds to that function. Calling a function pointer, calls whatever function that pointer is pointing to.

```
int (*func)(int *, char);
```

- “func is a pointer to a function which takes two arguments, a pointer to int and a char, and returns an int”
- Rare in OO languages,
  - but c.f. C# “delegates”; Eiffel “agents”
- As with all types, can be used with **typedef**
- Basis for lots of techniques in Systems code

From the Linux kernel...

```
struct file_operations {  
    ...  
    loff_t (*lseek)(struct file *, loff_t, int);  
    ssize_t (*read)(struct file *, char __user *, size_t, loff_t *);  
    ssize_t (*write)(struct file *, const char __user *, size_t, loff_t *);  
    ...  
};
```

- Record (struct) type of function pointers
  - Sometimes called a **vtable**
  - Each member is a different function type
  - Implements a “method” (or “feature”) of some “object”
    - In this case, a struct **file**
  - Provides polymorphism (e.g. files and sockets)
- Not really objects: no protection or hiding
  - But does provide the reuse benefits

## 6.4 Assertions

### 6.4.1 Assertions

```
assert( <scalar expression> );
```

- At run time, evaluate the expression.
- If true, do nothing
- If false:
  - Print “file.c:line: func: Assertion ‘expr’ failed.”
  - Abort (dump core)

Where true is a non-zero value.

Example:

```
#include <assert.h>
// It's a bug to call this with null a or b
void array_copy(int a[], int b[], size_t count)
{
    int i;
    assert(a != NULL);
    assert(b != NULL);
    for(i=0; i<count; i++) {
        a[i] = b[i];
    }
}

int main(int argc, char *argv[])
{
    // This is not going to go well...
    array_copy(NULL, NULL, 0);
    return 0;
}
```

\$ ./a.out  
a.out: assert\_text.c:8: array\_copy: Assertion `a != ((void \*)0)' failed.  
Aborted  
\$

Assert is a controlled crash, hence it is of no use to the user. Assertions are for programmers to find bugs, not for programs to detect errors.

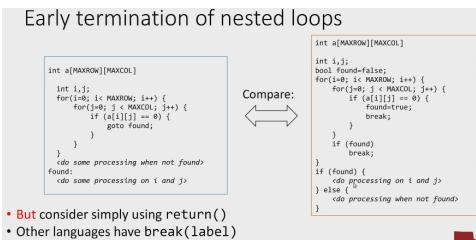
## 6.5 goto

goto jumps to a specific label in a program.

### 6.5.1 When to use goto:

DON'T, however...it can be useful in:

- Early termination of nested loops

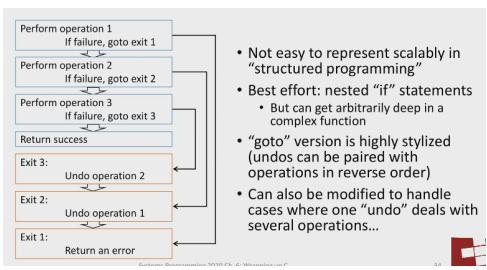


- Cleanup conditions

## Cleanup conditions

Goto is used for *recovery code*

- General idea:
  - Code performs a sequence of operations
  - Any one can fail
  - If one fails, all previous operations must be *undone*
- Canonical example:
  - Malloc'ing a sequence of buffers for data
  - If one fails, must free all previous generated buffers
- Code is often (not always) auto-generated



## 6.6 setjmp() and longjmp()

setjmp()	longjmp()
<pre>#include &lt;setjmp.h&gt; int setjmp(jmp_buf env);</pre> <ul style="list-style-type: none"> <li>• <code>setjmp(env):</code> <ul style="list-style-type: none"> <li>• Saves the current stack state / environment in <code>env</code></li> <li>• Returns 0.</li> </ul> </li> </ul>	<pre>#include &lt;setjmp.h&gt; void longjmp(jmp_buf env, int val);</pre> <ul style="list-style-type: none"> <li>• <code>longjmp(env, val):</code> <ul style="list-style-type: none"> <li>• Causes another return to the point saved by <code>env</code></li> <li>• This new return returns <code>val</code> (or 1 if <code>val</code> is 0)</li> <li>• This can only be done once for each <code>setjmp()</code></li> <li>• It is invalid if the function containing the <code>setjmp</code> returns</li> </ul> </li> </ul>

Example:

**Toy example**

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

static jmp_buf buf;

void second(void) {
    printf("second\n");
    longjmp(buf, 1);
}

void first(void) {
    second();
    printf("first\n");
}

int main() {
    if (!setjmp(buf)) {
        first(); // when executed, setjmp returns 0
    } else if ( // when longjmp jumps back, setjmp returns 1
        printf("main\n");
    }
    return 0;
}
```

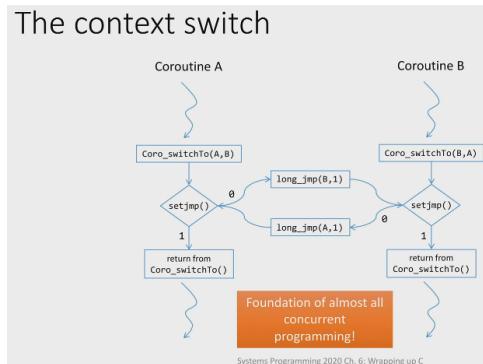
This can be used for Coroutines

## 6.7 Coroutines

A way for functions to call each other without doing subroutine calls. This allows us to have 2 functions running at the same time.

Creation/deletion	Switching
<pre>struct Coro *Coro_new(void) {     struct Coro *self = (struct Coro *)calloc(1, sizeof(struct Coro));     self-&gt;stack = (void *)calloc(1, 65536 + 16);     return self; }  void Coro_free(struct Coro *self) {     free(self-&gt;stack);     free(self); }</pre>	<p style="text-align: center;"><b>Switching</b></p> <div style="border: 1px solid black; padding: 10px;"> <pre>void Coro_switchTo(struct Coro *self,                    struct Coro *next) {     if (setjmp(self-&gt;env) == 0) {         longjmp(next-&gt;env, 1);     } }</pre> <ul style="list-style-type: none"> <li>• First return from <code>setjmp():</code> ⇒ <code>longjmp</code> to the next coroutine</li> <li>• Second return: ⇒ continue where we left off</li> </ul> </div>

This allows us to implement a socalled Context switch:



But how do we start?

## Initialization (the hard bit)

```
struct Coro_Cl {  
    void *context;  
    Coro_func_t *func;  
};
```

```
void Coro_Start_(struct Coro_Cl *block)  
{  
    (block->func)(block->context);  
    printf("Scheduler error: returned from coro start!\n");  
    exit(-1);  
}
```

This is a **closure**: a function bundled with a set of arguments.

## Initialization (the hard bit)

```
void Coro_start(struct Coro *self, struct Coro *other,  
                void *context, Coro_func_t *callback)  
{  
    struct Coro_Cl sblock;  
    struct Coro_Cl *block = &sblock;  
    block->context = context;  
    block->func = callback;  
  
    setjmp(self->env);  
    self->env[8]..._jmpbuf[6] = ((unsigned long)(self->stack));  
    self->env[8]..._jmpbuf[7] = ((long)Coro_Start_);  
    Coro_switchTo(self, other);  
}
```

Machine-dependent hack:  
6 => %rsp  
7 => %rip

Now: replace emit() and getchar() in original code with Coro\_switchTo() and you're done.

## Chapter 7

# Implementing dynamic memory allocation

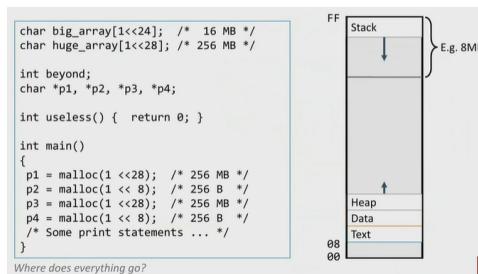
In order to run a java program we need a lot of stuff which actually isn't written in Java. The Java runtime contains memory allocator, memory manager, garbage collector and a low level I/O routine much of which is written in C/C++. The Java compiler generates byte code. The Java trace compiler will translate that into efficient machine code this is the JAVA VM. In C this doesn't exist, there is a C library, but it is not needed to run a C program. Hence this is why C is good for writing Operating Systems.

### 7.1 Dynamic Memory allocation

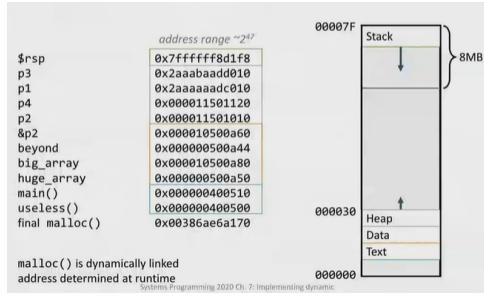
The Memory allocator gives out memory to applications in blocks. It sits in between the application itself which calls malloc and allocates memory out of some area of the virtual address space which has been allocated to it in turn by the OS, when the process starts it gets allocated to a certain amount of memory (the heap). malloc allocates memory out of the heap in arbitrary sized units and hands it to the application and frees it after use. Recall the malloc package:

```
#include <stdlib.h>
void *malloc(size_t size)
Successful:
    Returns a pointer to a memory block of at least size bytes
    (typically aligned to 8- or 16-byte boundary)
    If size == 0, returns NULL
Unsuccessful: returns NULL (0) and sets errno
void free(void *p)
    Returns the block pointed at by p to pool of available memory
    p must come from a previous call to malloc() or realloc()
void *realloc(void *p, size_t size)
    Changes size of block p and returns pointer to new block
    Contents of new block unchanged up to min of old and new size
    Old block has been free()'d (logically, if new != old)
```

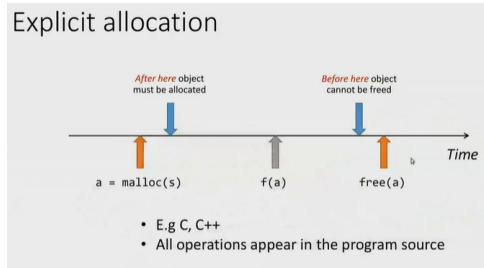
Memory allocation example:



The bottom is empty because we want to have the Null value. Heap is the part that malloc uses, the heap grows upwards when more and more memory is used. The Stack starts at the top of the memory and grows down. The Data segment is data the compiler knows exists. Text contains the machine code

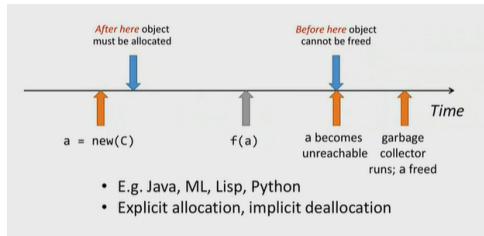


### 7.1.1 Explicit allocation



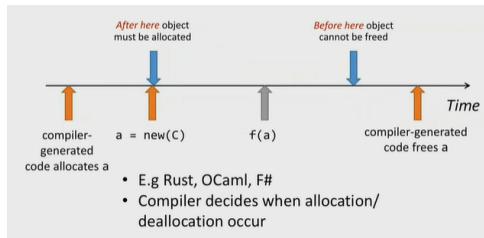
The blue arrow on the right is the earliest point in which the memory can be freed i.e it will not be used any more in the program, the orange arrow on the right is when the program calls free. The blue arrow on the left is the latest possible point in which the memory must be malloced.

In contrast garbage collection has the following timeline



Garbage collection makes it easier to code but decreases performance.

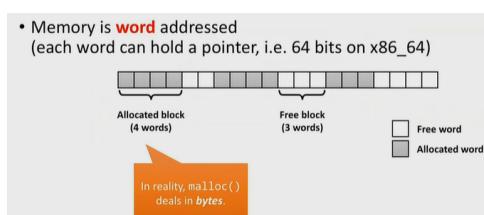
There is a third approach is Compiler support which is something inbetween the above approaches. The compiler looks at the program and figures out that an object is needed, because of the function call. The compiler can prove when the object is no longer needed.



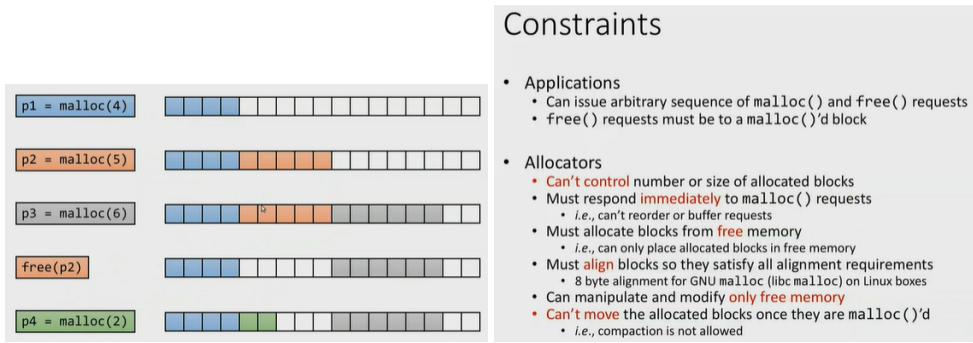
Compiler support cant be used for arbitrary programs and hence have a garbage collection which kicks in

## 7.2 The Explicit memory allocation Problem

### Assumptions



### Allocation Example:



### 7.2.1 Performance goal: Throuput

- Given some sequence of malloc and free requests:
  - $R_0, R_1, \dots, R_k, \dots, R_{n-1}$
- Goals: maximize throughput and peak memory utilization
  - These goals are often conflicting
- Throughput:
  - Number of completed requests per unit time
  - Example:
    - 5,000 malloc() calls and 5,000 free() calls in 10 seconds
    - Throughput is 1,000 operations/second
  - How to do malloc() and free() in O(1)? What's the problem?

### 7.2.2 Performance goal: Peak memory utilization

- Given some sequence of malloc and free requests:
  - $R_0, R_1, \dots, R_k, \dots, R_{n-1}$
- Def: Aggregate payload  $P_k$** 
  - malloc( $p$ ) results in a block with a **payload** of  $p$  bytes
  - After request  $R_i$  has completed, the **aggregate payload**  $P_k$  is the sum of currently allocated payloads
    - all malloc()'d stuff minus all free()'d stuff
- Def: Current heap size =  $H_k$** 
  - Assume  $H_k$  is monotonically nondecreasing
    - reminder: it grows when allocator uses sbrk()
- Def: Peak memory utilization after  $k$  requests**
  - $U_k = (\max_{i \leq k} P_i) / H_k$

### 7.2.3 Fragmentation

The heap gets into a state where it becomes hard to efficiently use it. We distinguish between:

- internal fragmentation
- external fragmentation

#### Internal fragmentation:

- For a given block, internal fragmentation occurs if payload < block size
 

The diagram shows a large rectangular block labeled 'block' at the top. Inside, there is a central area labeled 'payload' with arrows pointing to it from both sides. On either side of the payload are smaller gray rectangular areas labeled 'internal fragmentation' with arrows pointing to them from the center.
- Caused by
  - overhead of maintaining heap data structures
  - padding for alignment purposes
  - explicit policy decisions (e.g., to return a big block to satisfy a small request)
- Depends only on the pattern of previous requests
  - thus, easy to measure

#### External fragmentation:

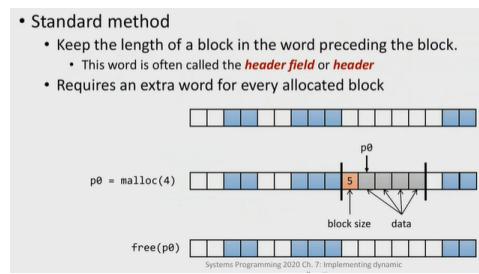
- Occurs when there is enough aggregate heap memory, but no single free block is large enough
 

The diagram shows a sequence of memory allocation requests and their effects on a memory heap. It consists of two columns: requests on the left and state transitions on the right.

  - p1 = malloc(4)**: Initial state with 4 slots.
  - p2 = malloc(5)**: After p1, 4 slots are filled with blue, followed by 5 slots filled with orange.
  - p3 = malloc(6)**: After p2, 4 blue, 5 orange, and 6 slots filled with grey.
  - free(p2)**: After p3, the 5 orange slots are freed, leaving 4 blue, 6 grey, and 1 white slot.
  - p4 = malloc(6)**: After free(p2), the 1 white slot and 6 grey slots are filled with green, creating a large free block of 7 slots.

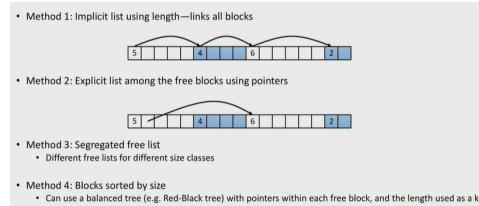
*Oops! (what would happen now?)*
- Depends on the pattern of future requests
  - thus, difficult to measure

## 7.2.4 How to free blocks

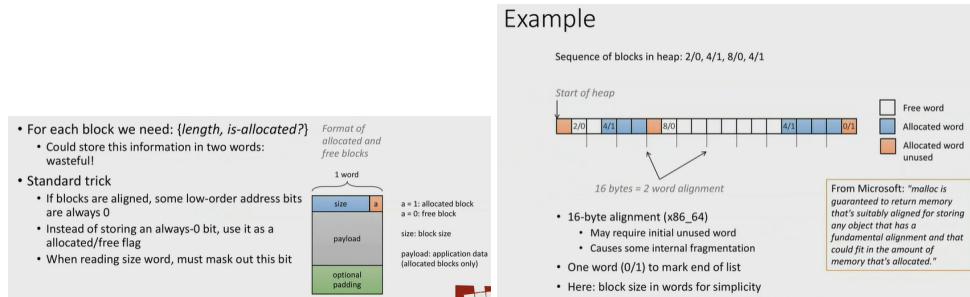


this method introduces internal fragmentation

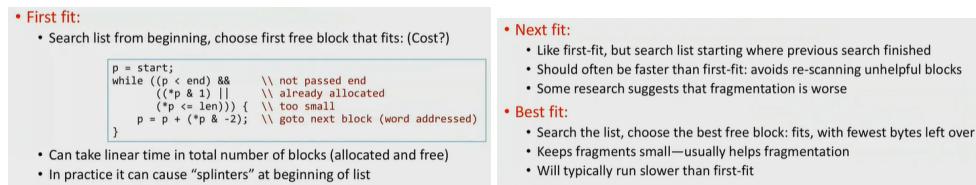
## 7.2.5 Keeping track of free blocks



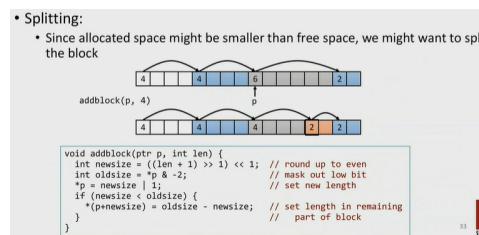
### • Implicit List



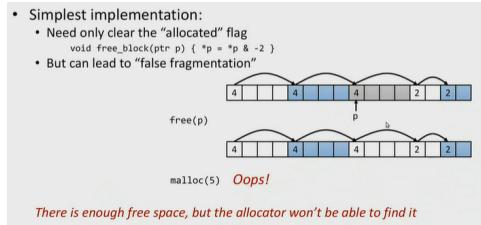
It is the payload which is aligned. 16 byte alignment makes it more efficient for the hardware to access values which are 16 bytes in size. A 16 byte aligned value is guaranteed not to span 2 memory pages and not to span 2 cache lines. The start of the heap is a word that never gets returned to the user because when we allocate the very first block of the heap, the first word is used to hold the size and the flag, which means what we're returning is a pointer to the second word. finding a free block



## Allocating a free block



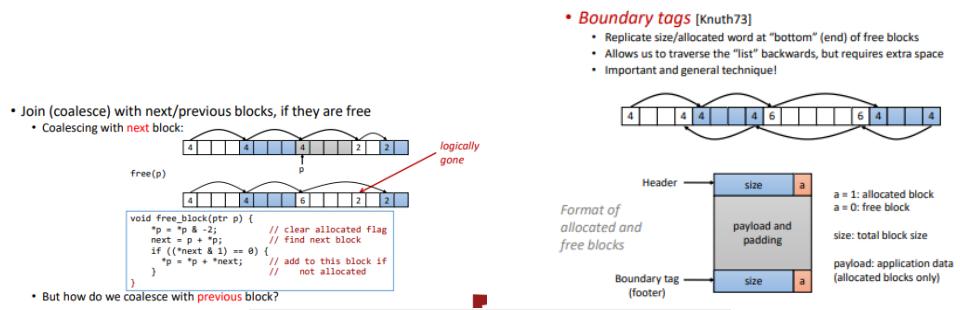
## Freeing a block



## 7.3 Coalescing

Transforming smaller blocks into fewer larger blocks to reduce external fragmentation

### Implicit list: bidirectional coalescing



- Boundary tags** [Knuth73]
    - Replicate size/allocated word at "bottom" (end) of free blocks
    - Allows us to traverse the "list" backwards, but requires extra space
    - Important and general technique!
- Header →

size	a
payload and padding	

Format of allocated and free blocks

Boundary tag (footer) →

size	a
------	---

a = 1: allocated block  
a = 0: free block

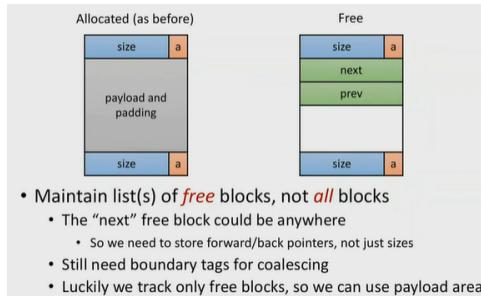
size: total block size  
payload: application data (allocated blocks only)

### Key allocator policies

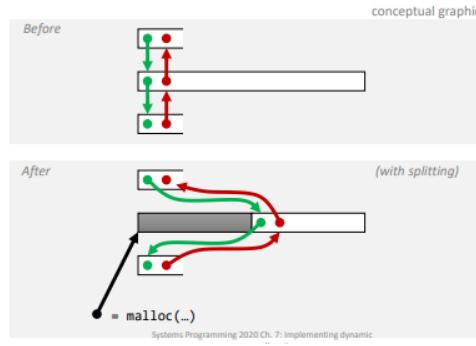
- Placement policy:
  - First-fit, next-fit, best-fit, etc.
  - Trades off lower throughput for less fragmentation
  - Interesting observation:* segregated free lists (see later) approximate a best fit placement policy without having to search entire free list
- Splitting policy:
  - When do we go ahead and split free blocks?
  - How much internal fragmentation are we willing to tolerate?
- Coalescing policy:
  - Immediate coalescing:* coalesce each time free() is called
  - Deferred coalescing:* try to improve performance of free() by deferring coalescing until needed. Examples:
    - Coalesce as you scan the free list for malloc()
    - Coalesce when the amount of external fragmentation reaches some threshold

## 7.4 Explicit free lists

### 7.4.1 concept



### 7.4.2 Allocation



### 7.4.3 freeing

- **Insertion policy:** Where in the free list do you put a newly freed block?
  - LIFO (last-in-first-out) policy
    - Insert freed block at the beginning of the free list
    - **Pro:** simple and constant time
    - **Con:** studies suggest fragmentation is worse than address ordered
  - Address-ordered policy
    - Insert freed blocks so that free list blocks are always in address order:  $addr(prev) < addr(curr) < addr(next)$
    - **Con:** requires search
    - **Pro:** studies suggest fragmentation is lower than LIFO

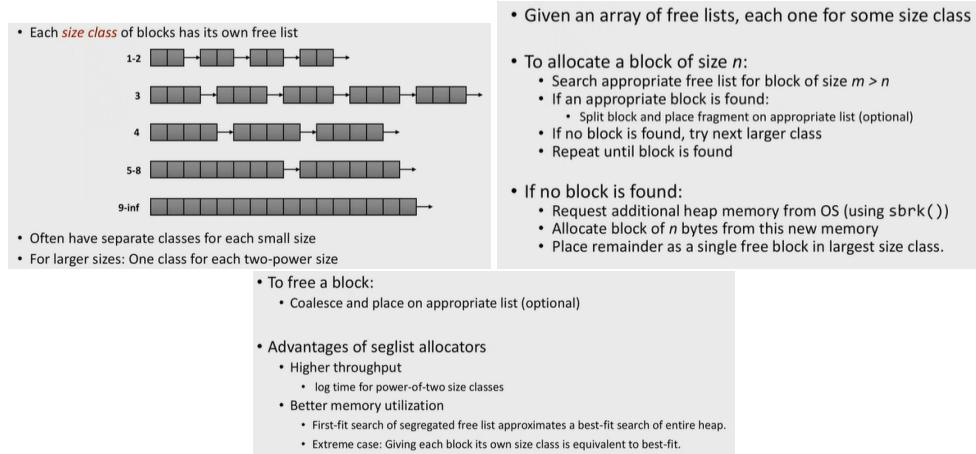


### 7.4.4 summary

Allocation is linear time in number of free blocks which is much faster than implicit lists when most of the memory is full. Its slightly more complicated to allocate and free since we need to splice blocks in and out of the list. Some extra space needed for the pointers.

## 7.5 Segregated free lists

Builds on the explicit free list. The idea is that we have a list for each size class of blocks:

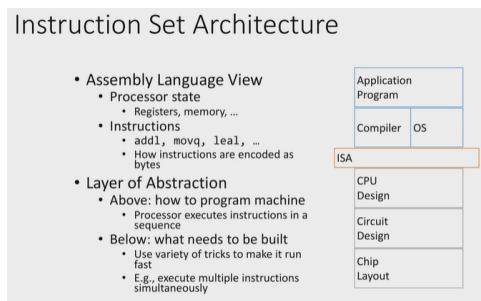


# Chapter 8

## Basic x86 architecture

### 8.1 What is an instruction set architecture

**Architecture:(also instruction set architecture:ISA)** The parts of a processor design that one needs to understand to write assembly code (e.g instruction set specification, registers)



**Microarchitecture:** Implementation of the architecture (cache sizes and core frequency)

#### 8.1.1 CISC: Complex Instruction Set

- Dominant style through mid-80's
- Stack-oriented** instruction set
  - Use stack to pass arguments, save program counter
  - Explicit push and pop instructions
- Arithmetic instructions can access memory
  - addl %eax, 12(%rbx,%rcx,4)
    - requires memory read and write
    - Complex address calculation
- Condition codes**
  - Set as side effect of arithmetic and logical instructions
- Philosophy
  - Add instructions to perform "typical" programming tasks

#### 8.1.2 RISC: Reduced Instruction Set

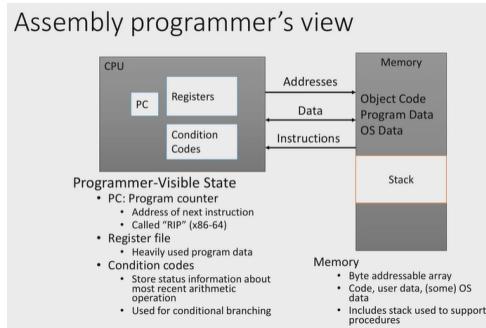
Idea was to have fewer easier instructions and give more work to the compiler. This lead to more space available for cache as well as better compiler optimization.

- Internal project at IBM
  - Popularized by Hennessy (Stanford) and Patterson (Berkeley)
- Fewer, **simpler** instructions
  - Might need **more** to get given task done
  - Can execute them with **small** and **fast** hardware
- Register-oriented** instruction set
  - Many more (typically 32) registers
  - Use for arguments, return pointer, temporaries
- Load-Store architecture
  - Only load and store instructions can access memory
- No **condition codes**
  - Test instructions return 0/1 in register

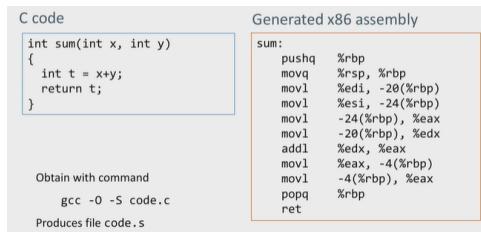
## 8.2 Basics of machine code

There are two common ways to write *x86* Assembler:

- *AT&T* syntax (What we will use)
- Intel syntax: Generally used for Windows machines



### 8.2.1 Compiling into assembly



### 8.2.2 Assembly data types

Integer data of 1,2,4 or 8 bytes which hold data values or addresses (untyped pointer)

Floating point data of 4,8 or 10bytes

No aggregate types(arrays,structures,...) only coniguously allocated bytes in memory

### 8.2.3 Assembly code operations

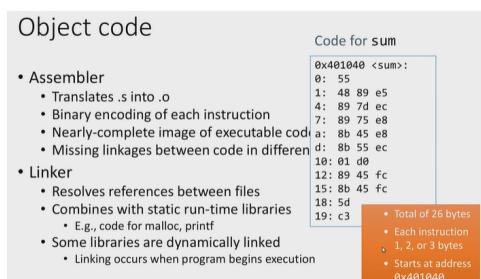
We can perform arithmetic function on register or memory data.

Transfer data between memory and register i.e

- Load data from memory into register
- Store register data into memory

Transfer control i.e Unconditional jumps to/from procedures and conditional branches

### 8.2.4 Object Code



### 8.2.5 Machine instruction example

#### Machine instruction example

```
int t = x+y;
addl $(%rbp),%eax
```

Similar to expression:  
 $x += y$   
More precisely:  
 $\text{int eax;}$   
 $\text{int *rbp;}$   
 $\text{eax += rbp[2]}$

0x401046: 03 45 08

- C Code
  - Add two signed integers
- Assembly
  - Add 2-byte integers
    - “Long” words in GCC parlance
    - Same instruction whether signed or unsigned
  - Operands
    - x: Register %eax
    - y: Memory M[%rbp+8]
    - t: Register %eax
  - Return function value in %eax
- Object Code
  - 3-byte instruction
  - Stored at address 0x401046

#### Disassembling object code

##### Disassembled

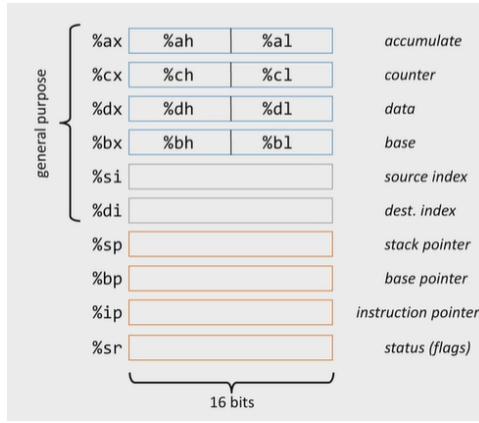
```
0000000000000000 <sum>:
0: 55 push %rbp
1: 48 89 e5 mov %rsp,%rbp
4: 89 7d ec mov %edi,-0x14(%rbp)
7: 89 75 e8 mov %esi,-0x18(%rbp)
a: 8b 45 e8 mov -0x18(%rbp),%eax
d: 8b 55 ec mov -0x14(%rbp),%edx
10: 01 d0 add %edx,%eax
12: 89 45 fc mov %eax,-0x4(%rbp)
15: 8b 45 fc mov -0x4(%rbp),%eax
18: 5d pop %rbp
19: c3 retq
```

##### • Disassembler

- objdump -d p
- Useful tool for examining object code
- Analyzes bit pattern of series of instructions
- Produces approximate rendition of assembly code
- Can be run on either a.out (complete executable) or .o file

## 8.3 x86 architecture

### 8.3.1 8086 Registers

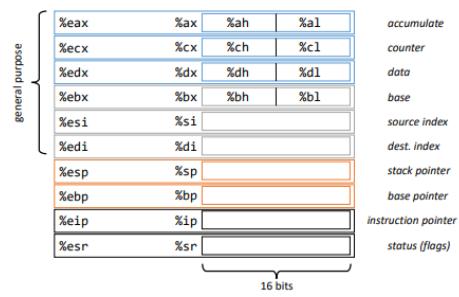


The source/destination index registers were limited by the instruction set to only hold addresses.

The stack pointer was a register only used for the stack.

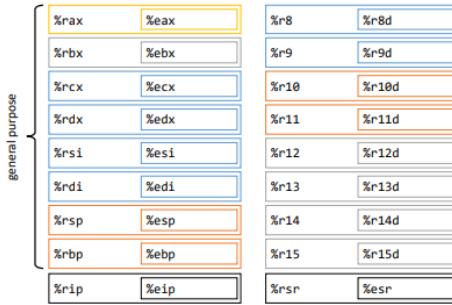
The only general purpose were the top 4 registers which are divided into lo and hi (remnants of the previous 8-bit design)

### 8.3.2 80386 Registers



The 8086 was extended by 16 bits, they kept the old names to refer to the halves and added e to the names to indicate the extended register.

### 8.3.3 x86-64 integer registers



Extension of the 386 by 32 bits with prefix now r. They added 8 more registers. all registers are now general purpose.

### 8.3.4 Moving data

- **movx Source, Dest**
  - x in {b, w, l,q}
  
- **movq Source, Dest:**  
Move 8-byte “quad word”
- **movl Source, Dest:**  
Move 4-byte “long word”
- **movw Source, Dest:**  
Move 2-byte “word”
- **movb Source, Dest:**  
Move 1-byte “byte”

movx Source, Dest:
• Operand Types
• <b>Immediate:</b> Constant integer data <ul style="list-style-type: none"> <li>• Example: \$0x400, \$-533</li> <li>• Like C constant, but prefixed with '\$'</li> <li>• Encoded with 1, 2, 4, 8 bytes</li> </ul>
• <b>Register:</b> One of 16 integer registers <ul style="list-style-type: none"> <li>• Example: %eax, %r14d</li> <li>• Note some (e.g. %rsp, %rbp) reserved for special use</li> <li>• Others have special uses for particular instructions</li> </ul>
• <b>Memory:</b> 1,2,4, or 8 consecutive bytes of memory at address given by register <ul style="list-style-type: none"> <li>• Simplest example: (%rax)</li> <li>• Various other “address modes”</li> </ul>

#### Possible combinations

movl operand combinations				
	Source	Dest	Src,Dest	C Analog
movl	Imm	Reg	movl \$0x4,%eax	temp = 0x4;
	Mem	Reg	movl \$-147,(%rax)	*p = -147;
	Reg	Reg	movl %eax,%edx	temp2 = temp1;
	Reg	Mem	movl %eax,(%rdx)	*p = temp;
	Mem	Reg	movl (%rax),%edx	temp = *p;

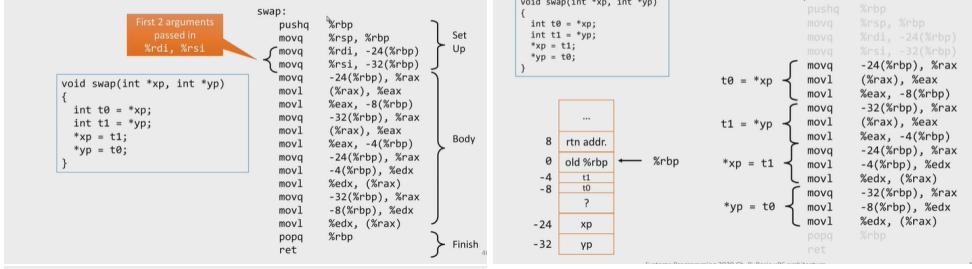
*Cannot do memory-memory transfer with a single instruction*

### 8.3.5 Simple memory addressing modes

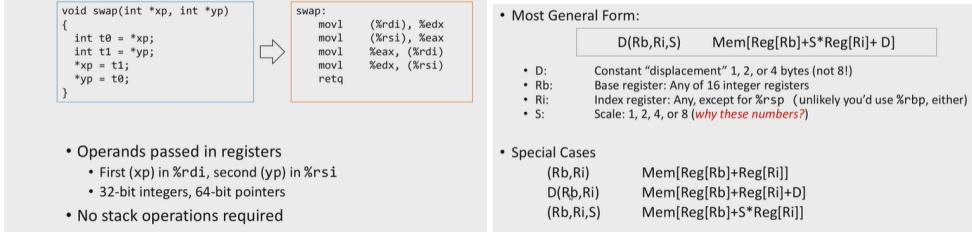
• Normal (R) Mem[Reg[R]]
• Register R specifies memory address
movq (%rcx),%rax
On x86_64, can also be %rip
• Displacement D(R) Mem[Reg[R]+D]
• Register R specifies start of memory region
• Constant displacement D specifies offset
movl 8(%ebp),%edx

Since we can use the instruction pointer as a source, code can be executed anywhere in memory. Displacement is a register offset. Example:

## Using simple addressing modes



With the optimizer on...



The most general form allows us to perform tricks e.g array indexing, accessing fields in a struct, accessing fields in an array of structs. S is one of those sizes depending on the data type being used.

### 8.3.6 Address Computation

Address computation is done with the lea function:

*leaSrc, Dest*

where Src is the address mode expression and Dest gets set to the address denoted by the expression. With this function we can compute addresses without a memory reference or compute arithmetic expressions of the form  $x + k \cdot y, k \in \{1, 2, 4, 8\}$

## 8.4 x86 integer arithmetic

### 8.4.1 ordinary arithmetic operations

- Two-operand instructions (longword variants):

Mnemonic	Format	Computation
addl	Src,Dest	Dest $\leftarrow$ Dest + Src
subl	Src,Dest	Dest $\leftarrow$ Dest - Src
imull	Src,Dest	Dest $\leftarrow$ Dest * Src
sall	Src,Dest	Dest $\leftarrow$ Dest << Src
salr	Src,Dest	Dest $\leftarrow$ Dest >> Src
shrl	Src,Dest	Dest $\leftarrow$ Dest >> Src
xorl	Src,Dest	Dest $\leftarrow$ Dest ^ Src
andl	Src,Dest	Dest $\leftarrow$ Dest & Src
orl	Src,Dest	Dest $\leftarrow$ Dest   Src

- No distinction between signed and unsigned int (why?)

There is no distinction between signed and unsigned values. The processor doesn't have much of a concept of the difference, it's just a register, just a 64-bit value. Hence none of these operations imply signed or unsigned computation because their results are the same.

### 8.4.2 Single operand instructions

- One operand instructions

Mnemonic	Format	Computation
incl	Dest	Dest $\leftarrow$ Dest + 1
decl	Dest	Dest $\leftarrow$ Dest - 1
negl	Dest	Dest $\leftarrow$ -Dest
notl	Dest	Dest $\leftarrow$ ~Dest

- See book for more instructions

### 8.4.3 Using leal for arithmetic expressions

<pre> int arith {     int x, int y, int z     {         int t1 = x*y;         int t2 = z+1;         int t3 = x+4;         int t4 = t1 * 48;         int t5 = t3 + t4;         int rval = t2 * t5;         return rval;     } } </pre>	<pre>         leal    (%rdi,%rsi), %eax    # eax = x + y         addl    %rsi, %eax          # edx = x + eax         leal    (%rsi,%rsi,2), %edx # edx = y * 3         sall    \$4, %edx            # edx *= 16         leal    4(%rdi,%rdx), %ecx # ecx = x + 4 + edx         imull   %ecx, %eax         # eax *= ecx         ret </pre>
---	---

The following happens in the lea code lines:

1. rdi and rsi are the first arguments which are x and y respectively. lea takes rdi as the base and rsi as the index i.e it adds the two together. The result is put into eax.
2. takes y, adds it to y multiplied by 2 and stores it in edx, hence we have a multiplication by 3. There is no instruction on a typical RISC machine which multiplies a number by 3. We multiply by 3 because later we need a multiplication by 48 which is  $3 \times 16$ . A multiplication by 16 can easily be computed a left shift.
3. we add rdi to rdx and add 4. Hence we are doing 2 instructions in one with lea

## 8.5 condition codes

Extra bits in the processor state. Can be thought of as 1-bit registers. They are contained in the status register.

- Single bit registers
  - CF Carry Flag (for unsigned) SF Sign Flag (for signed)
  - ZF Zero Flag OF Overflow Flag (for signed)
- Implicitly set (think of it as *side effect*) by arithmetic operations
  - Example: `addl/addq Src,Dest  $\leftrightarrow$  t = a+b`
  - **CF set** if carry out from most significant bit (unsigned overflow)
  - **ZF set** if t == 0
  - **SF set** if t < 0 (as signed)
  - **OF set** if two's complement (signed) overflow:  
 $(a>0 \&& b<0 \&& t<0) \mid\mid (a<0 \&& b>0 \&& t>0)$
- **Not** set by lea instruction
- Full documentation link on course website

### 8.5.1 Compare Instruction

- Explicit Setting by Compare Instruction
 

```

        cmp1/cmpq Src2,Src1
        cmp1 b,a like computing a - b without setting destination
        CF set if carry out from most significant bit
        (used for unsigned comparisons)
        ZF set if a == b
        SF set if (a-b) < 0 (as signed)
        OF set if two's complement (signed) overflow:
        (a>0 && b<0 && (a-b)<0) || (a<0 && b>0 && (a-b)>0)
      
```

### 8.5.2 Test Instruction

- Explicit Setting by Test instruction
 

```

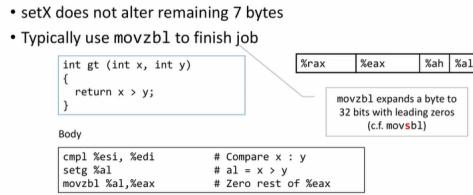
        test1/testq Src2,Src1b
        test1 b,a like computing a & b without setting destination
        • Sets condition codes based on value of Src1 & Src2
        • Useful to have one of the operands be a mask
        ZF set when a & b == 0
        SF set when a & b < 0
      
```

### 8.5.3 SetX Instructions

- SetX Instructions
  - Set single byte based on combinations of condition codes

SetX	Condition	Description
sete	ZF	Equal / Zero
setne	$\sim$ ZF	Not Equal / Not Zero
sets	SF	Negative
setsns	$\sim$ SF	Nonnegative
setg	$\sim$ (SF $\wedge$ OF) $\wedge$ $\sim$ ZF	Greater (Signed)
setge	$\sim$ (SF $\wedge$ OF)	Greater or Equal (Signed)
setl	(SF $\wedge$ OF)	Less (Signed)
setle	(SF $\wedge$ OF) $\mid$ ZF	Less or Equal (Signed)
seta	$\sim$ CF $\wedge$ ZF	Above (unsigned)
setb	CF	Below (unsigned)

The result of the computations get set into the low byte of some other register. If setting the value on a lower order of bits of a register than the top bits are set to 0 e.g if we set eax to a value, then rax gets the same value and we zero out the top bits., Example:



#### 8.5.4 Jump Instructions

##### jX Instructions:

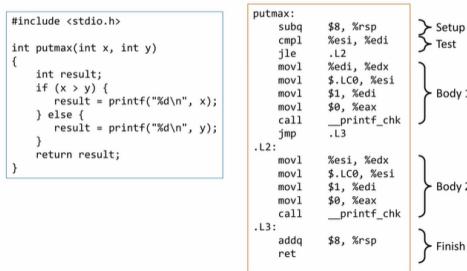
Jump to different part of code depending on condition codes

jX	Condition	Description
jmp	1	Unconditional
je	ZF = 1	Equal / Zero
jne	-ZF	Not Equal / Not Zero
js	SF	Negative
jns	-SF	Non-negative
jg	~(SF^OF) & ~ZF	Greater (Signed)
jge	~(SF^OF)	Greater or Equal (Signed)
jl	(SF^OF)	Less (Signed)
jle	(SF^OF)   ZF	Less or Equal (Signed)
ja	~CF & ~ZF	Above (unsigned)
jb	CF	Below (unsigned)

# Chapter 9

## Compiling C control flow

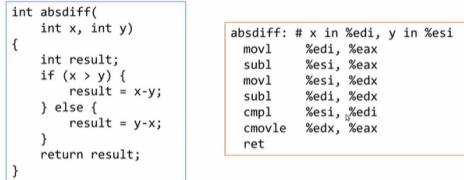
### 9.1 if-then-else statements



Body 1 is the case where  $x > y$ , Body 2 is  $x \leq y$ .

#### 9.1.1 Conditional move

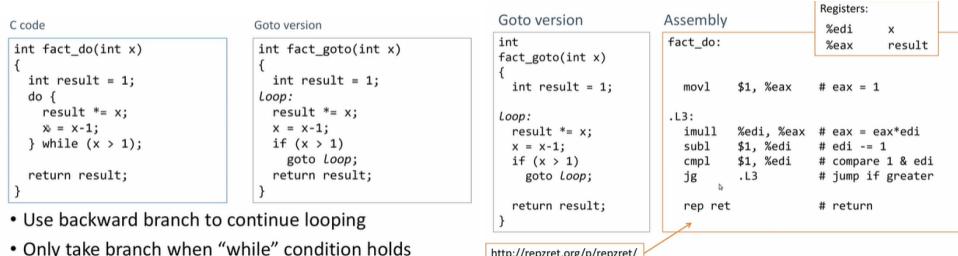
In some cases the compiler can avoid using branch statements and use the so called conditional move instruction:



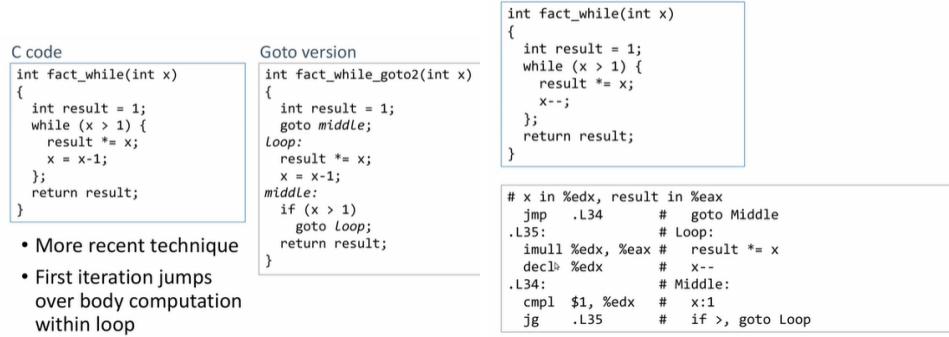
- Conditional move instruction
  - `cmoveC src, dest`
  - Move value from `src` to `dest` if condition `C` holds
  - More efficient than conditional branching (simple control flow)
  - But overhead: both branches are evaluated

Conditional move is more efficient than branching.

### 9.2 do-while loops

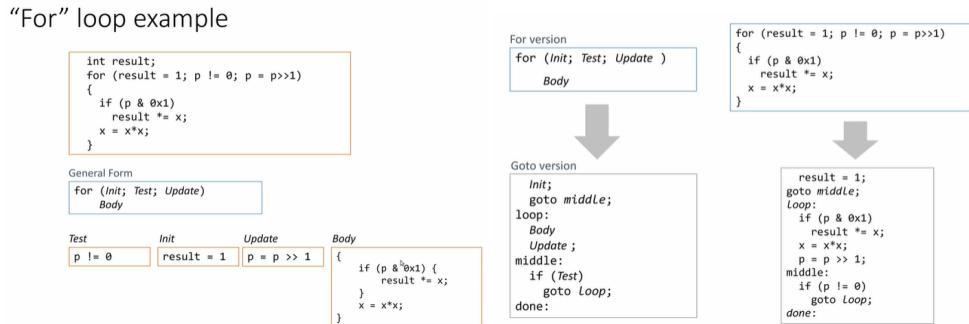


### 9.3 while loops



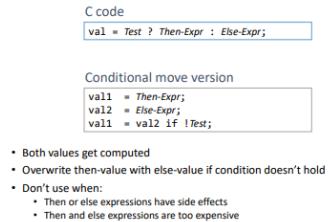
### 9.4 for loops

A for loop is a while loop with an initialisation instruction.



#### 9.4.1 Summary of loop variants

##### General form with conditional move

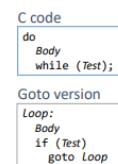


##### General "do-while" translation

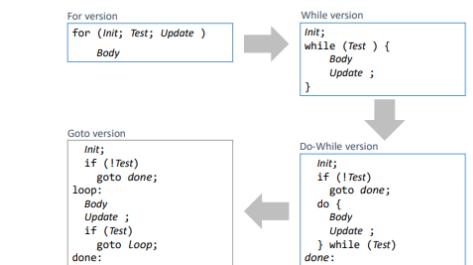
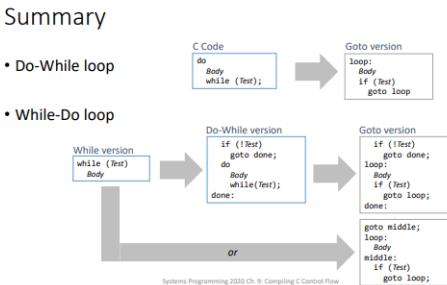
- **Body:**

  - { Statement<sub>1</sub>; Statement<sub>2</sub>; ... Statement<sub>n</sub>; }

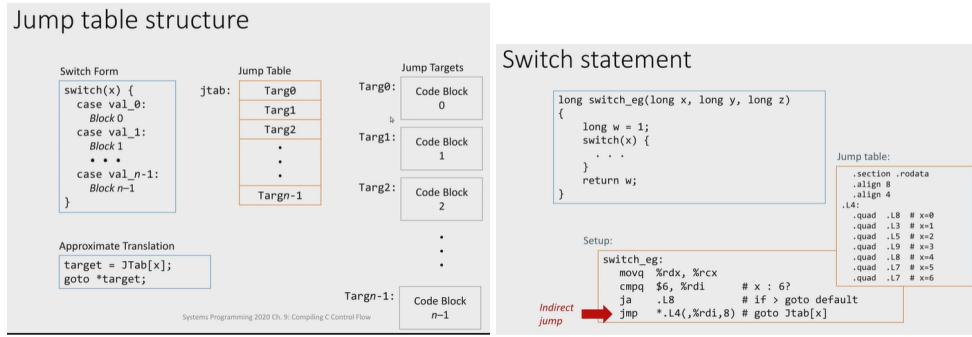
- **Test returns integer**
  - = 0 interpreted as false
  - ≠ 0 interpreted as true



##### "For" → "While" → "Do-While"



## 9.5 compact switch statements

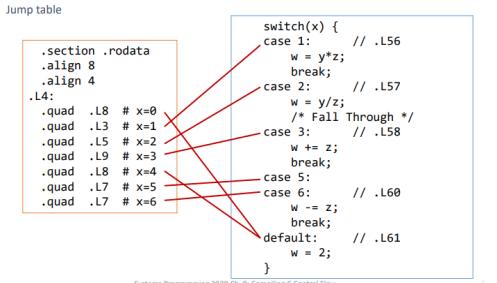


When we execute the case statement, we look at the value to be tested. We get the address to jump to from the jump table. The assembly code is compact because of the indirect jump, which jumps to an address whose value depends on x (rdi).

### Assembly setup explanation

- Table Structure**
  - Each target requires 8 bytes
  - Base address at .L4
- Jumping**
  - Direct:** jmp .L8
  - Jump target is denoted by label .L8
- Indirect:** jmp \*.L4(%rdi,8)
  - Must scale by factor of 8 (labels are 64-bit = 8 Bytes on x86\_64)
  - Fetch target from effective Address .L61 + rdi\*8
    - Only for  $0 \leq x \leq 6$

The order of the jump table might not be in the same order as the cases because there might be values for x which are not checked and hence must be directed to the default case.



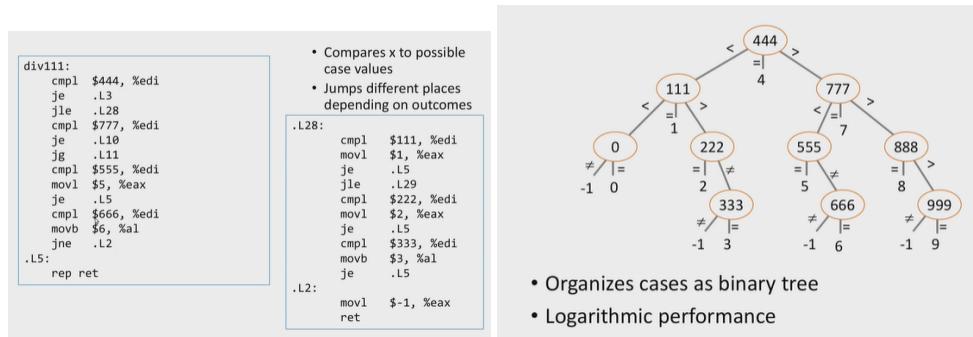
## 9.6 sparse switch statements

The case values have to be constants. If the values are far apart from each other, it's not practical to use a jump table e.g.

```

/* Return x/111 if x is multiple
   && <= 999. -1 otherwise */
int div111(int x)
{
    switch(x) {
    case 0: return 0;
    case 111: return 1;
    case 222: return 2;
    case 333: return 3;
    case 444: return 4;
    case 555: return 5;
    case 666: return 6;
    case 777: return 7;
    case 888: return 8;
    case 999: return 9;
    default: return -1;
}
  
```

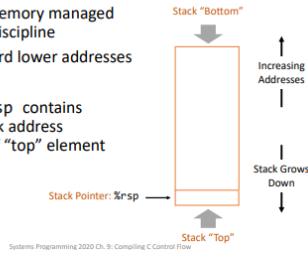
An obvious solution would be to implement this with if else statements, but this grows linear with the cases. When compiled we get the following code:



## 9.7 Procedure call and return

### x86\_64 stack

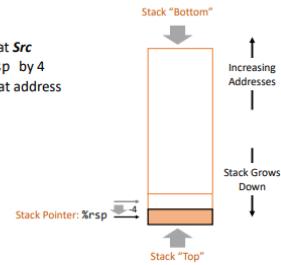
- Region of memory managed with stack discipline
- Grows toward lower addresses
- Register `%rsp` contains lowest stack address = address of "top" element



### x86\_64 stack: push

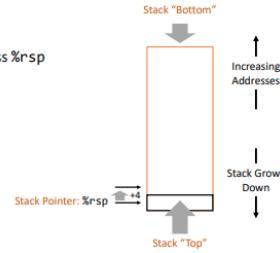
### x86\_64 stack: push

- `pushl Src`
- Fetch operand at `Src`
- Decrement `%rsp` by 4
- Write operand at address given by `%rsp`



### x86\_64 stack: pop

- `popl Dest`
- Read operand at address `%rsp`
- Increment `%rsp` by 4
- Write operand to `Dest`



### Procedure control flow

- Use stack to support procedure call and return
- **Procedure call:** `call Label`

- Push return address on stack
- Jump to `Label`

- **Return address:**

- Address of instruction beyond `call`

- Example from disassembly

```
884854e: e8 3d 06 00 00 call 8048b90 <main> in>
8848553: 50 pushl %eax
```

- Return address = `0x8048553`

- **Procedure return:** `ret`

- Pop address from stack
- Jump to address

## 9.8 Calling conventions

How do you construct a call/return using machine code.

### Recursive factorial

```

int rfact(int x)
{
    int rval;
    if (x <= 1)
        return 1;
    rval = rfact(x-1);
    return rval * x;
}
  
```

```

rfact:
    pusha %rbx
    movl %edi, %ebx
    movl $1, %eax
    cmpl $1,%edi
    jle .L2
    leal -1(%rdi), %edi
    call rfact
    imull %ebx, %eax
.L2:
    popq %rbx
    ret
  
```

- Registers

- `%rax/%eax` used without first saving

- `%rbx/%ebx` used, but saved at beginning & restore at end

### 9.8.1 Register saving conventions

#### Register saving conventions

- When procedure `yoo` calls who:
  - `yoo` is the **caller**
  - who is the **callee**
- Can a register be used for temporary storage?

```
yoo:
    ...
    movl $15213, %edx
    call who
    addl %edx, %eax
    ...
    ret

who:
    ...
    movl 8(%rsp), %edx
    addl $91125, %edx
    ...
    ret
```

- Contents of register `%edx` overwritten by `who`

#### x86-64 integer registers

%rax	Return value, # varargs
%rbx	Callee saved; base ptr
%rcx	Argument #4
%rdx	Argument #3 (& 2 <sup>nd</sup> return)
%rsi	Argument #2
%rdi	Argument #1
%rsp	Stack pointer
%rbp	Callee saved; frame ptr
%r8	Argument #5
%r9	Argument #6
%r10	Static chain ptr
%r11	Used for linking
%r12	Callee saved
%r13	Callee saved
%r14	Callee saved
%r15	Callee saved

Argument Registers are caller save registers.

- Arguments passed to functions via registers
  - If more than 6 integral parameters, then pass rest on stack
  - These registers can be used as caller-saved as well
- All references to stack frame via stack pointer
  - Eliminates 32-bit need to update `%ebp/%rbp`
- Other registers
  - 6+1 callee saved
  - 2 or 3 have special uses

Example:

#### x86-64 stack frame example

```
long sum = 0;
/* Swap a[i] & a[i+1] */
void swap_ele_su(long a[], int i)
{
    swap(a[i], a[i+1]);
    sum += a[i];
}
```

- Keeps values of `a` and `i` in callee save registers
- Must set up stack frame to save these registers

```
swap_ele_su:
    movq %rbx, -16(%rsp)          # Save %rbx
    movslq %esi,%rbx              # Extend & save i
    movq %r12, -8(%rsp)           # Save %r12
    movq %rdi, %r12               # Save a
    leaq (%rdi,%rbx,8), %rdi    # &a[i]
    subq $16, %rsp                # Allocate stack frame
    leaq 8(%rdi), %rsi            # &a[i+1]
    call swap()                  # swap()
    movq (%r12,%rbx,8), %rax    # a[i]
    addq %rax, sum(%rip)          # sum += a[i]
    movq (%rsp), %rbx             # Restore %rbx
    movq 8(%rsp), %r12            # Restore %r12
    addq $16, %rsp                # Deallocate stack frame
    ret
```

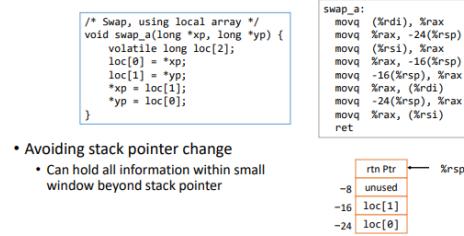


#### Interesting features of the stack frame

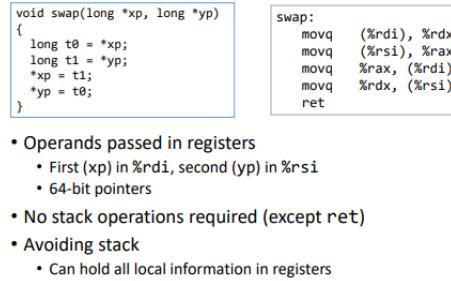
- Allocate entire frame at once
  - All stack accesses can be relative to `%rsp`
  - Do by decrementing stack pointer
  - Can delay allocation, since safe to temporarily use **red zone** - see next slide
- Simple deallocation
  - Increment stack pointer
  - No base/frame pointer needed

We can keep the values of `a` and `i` in registers because the `swap` function has to preserve callee save registers, hence we don't need to write them to memory.

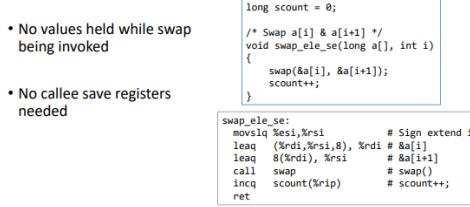
### 9.8.2 Procedure Calls: Locals in the red zone



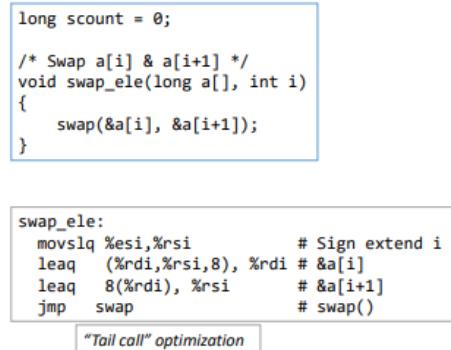
### 9.8.3 Procedure Calls: Long swap



### 9.8.4 Procedure Calls: non-leaf w/o stack frame



### 9.8.5 Procedure Calls: call using a jump



# Chapter 10

## Compiling C data Structures

### 10.1 One-dimensional arrays

#### Basic data types

##### • Integral

- Stored & operated on in general (integer) registers
- Signed vs. unsigned depends on instructions used

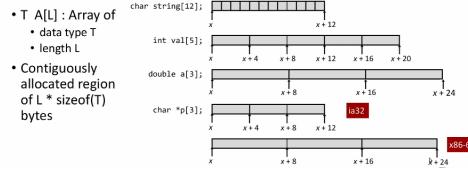
Intel	GAS	Bytes	C
byte	b	1	[unsigned] char
word	w	2	[unsigned] short
double word	l	4	[unsigned] int
quad word	q	8	[unsigned] long int (x86-64)

##### • Floating Point

- Stored & operated on in floating point registers

Intel	GAS	Bytes	C
Single	s	4	float
Double	l	8	double
Extended	t	10/12/16	long double

#### 10.1.1 Array allocation



When allocating multiple arrays its not guaranteed that they will be contiguous in memory.

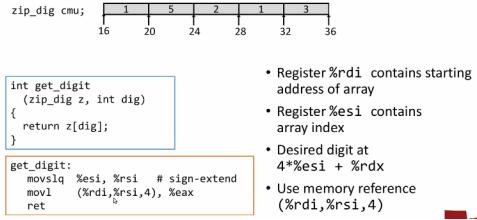
#### 10.1.2 Array access

- Array of data type  $T$  and length  $L$   $T A[L];$
- Identifier A can be used as a pointer to array element 0: Type  $T^*$

int val[5];

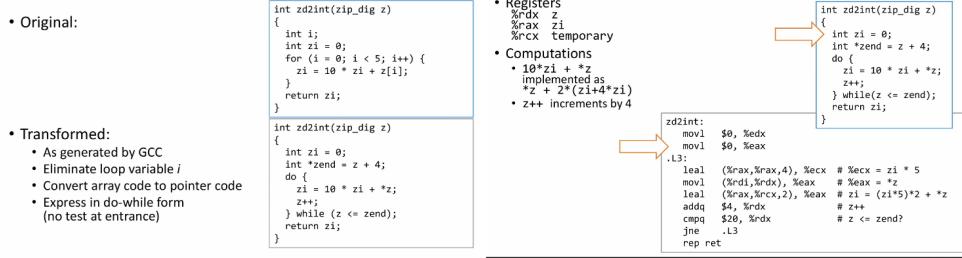
Reference	Type	Value
val[4]	int	3
val	int *	x
val + 1	int *	x + 4
&val[2]	int *	x + 8
val[5]	int	??
*(val+1)	int	5
val + i	int *	x + 4 I

Example:

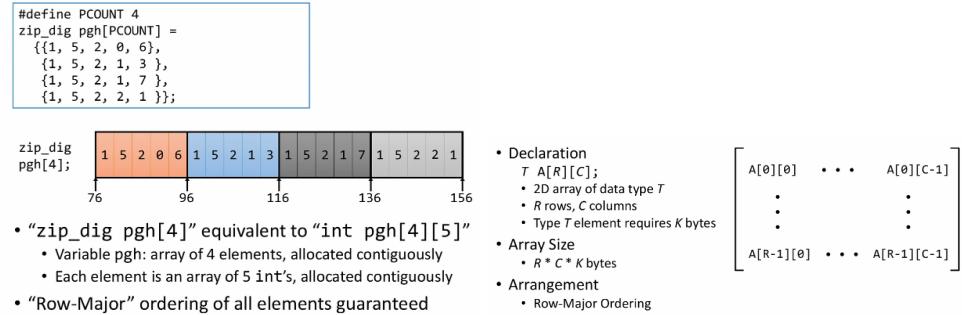


In C there is no bound checking, Out of range behaviour is implementation dependant and there is no guaranteed relative allocation of different arrays.

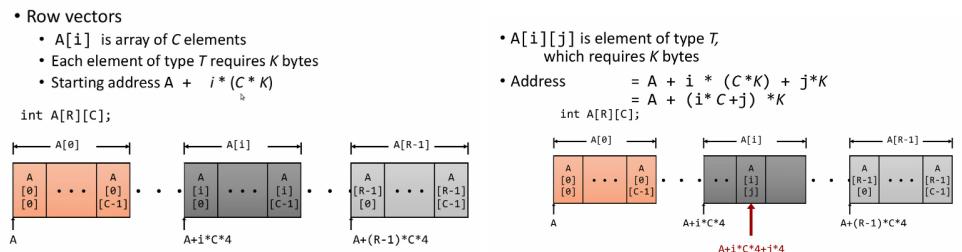
### 10.1.3 Array Loop example



## 10.2 Nested Array



### 10.2.1 Nested array row access



### 10.2.2 Nested array row access code

```

int *get_pgh_zip(int index)
{
    return pgh[index];
}

#define PCOUNT 4
zip_dig pgh[PCOUNT] =
{{1, 5, 2, 0, 6},
 {1, 5, 2, 1, 3},
 {1, 5, 2, 1, 7},
 {1, 5, 2, 2, 1}};

```

- What data type is `pgh[index]`?
- What is its starting address?

```

get_pgh_zip:
    movslq %edi, %rdi
    leaq (%rdi,%rdi,4), %rax # 5 * index
    leaq pgh(%rax,4), %rax # pgh + (20 * index)
    ret

```

### 10.2.3 Nested array element access code

#### Strange referencing examples

• Array elements

- pgf[index][dig] is int
- Address: pgf + 20\*index + 4\*dig

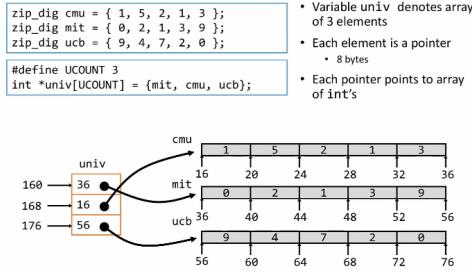
```
int get_pgh_digit
    (int index, int dig)
{
    return pgf[index][dig];
}
```

zip_dig	1 5 2 0 6	1 5 2 1 3	1 5 2 1 7	1 5 2 2 1	
pgf[4];	76	96	116	136	156

• Code does not do any bounds checking

• Ordering of elements within array guaranteed

## 10.3 Multi-level arrays



### 10.3.1 Element access in multi-level array

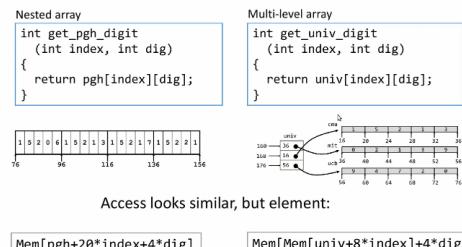
• Computation (x86\_64)

- Element access Mem[Mem[univ+8\*index]+4\*dig]
- Must do two memory reads
- First get pointer to row array
- Then access element within array

```
int get_univ_digit
    (int index, int dig)
{
    return univ[index][dig];
}
```

```
get_univ_digit:
    movslq %rsi, %rdi
    movslq %rsi, %rcx
    movl   %rcx(%rdi,%rdi,8), %eax # Mem[univ+8*index]
    movl   (%eax,%rsi,4), %eax # Mem[...+4*dig]
    ret
```

In comparison to nested arrays we need two memory accesses to access an element in a multi-level array. The first one is to access the first pointer table, the second is to access the element of the array the pointer is pointing to. For each level, we need a memory access more. For nested array we always need one memory access.



Access looks similar, but element:

Mem[pgh+20\*index+4\*dig]      Mem[Mem[univ+8\*index]+4\*dig]

### 10.3.2 Referencing examples

cmu	1 5 2 1 3
mit	0 2 1 3 9
ucb	9 4 7 2 0

• Code does not do any bounds checking

• Ordering of elements in different arrays not guaranteed

## 10.4 Structures

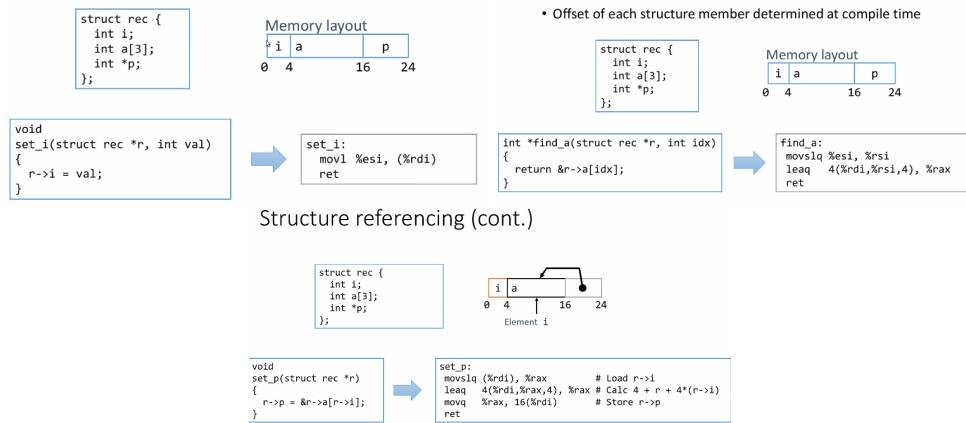
### Structures: concept



Contiguously-allocated region of memory  
Refer to members within structure by names  
Members may be of different types

The members appear in the same order in memory as declared, but there could be padding inbetween

#### 10.4.1 Accessing structure members



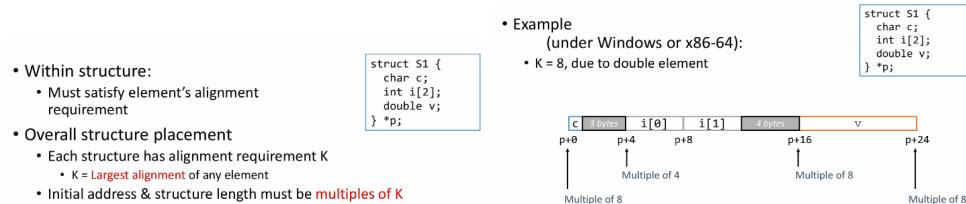
## 10.5 Alignment

If we have a primitive data type which requires K bytes, the address where that data is stored must be a multiple of K. The compiler inserts gaps in structures to ensure correct alignment of fields if necessary.

### Specific cases of alignment (x86-64)

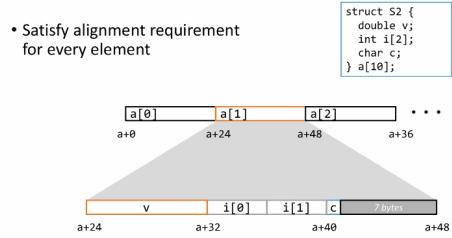
- 1 byte: char, ...  
• no restrictions on address
- 2 bytes: short, ...  
• lowest 1 bit of address must be 0<sub>2</sub>
- 4 bytes: int, float, ...  
• lowest 2 bits of address must be 00<sub>2</sub>
- 8 bytes: double, char \*, ...  
• Windows & Linux:  
• lowest 3 bits of address must be 000<sub>2</sub>
- 16 bytes: long double  
• Linux:  
• lowest 3 bits of address must be 000<sub>2</sub>  
• i.e., treated the same as a 8-byte primitive data type

#### 10.5.1 Alignment with structures



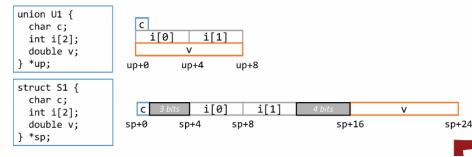
In some cases reordering the fields from largest to smallest size can save memory.

## 10.6 Arrays of structures



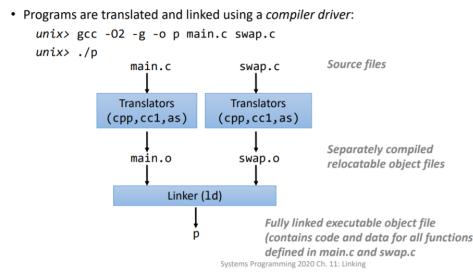
## 10.7 Unions

- Allocate according to largest element
- Can only use one field at a time



# Chapter 11

## Linking



Linkers allow modularity. A program can be written as a collection of smaller source files, rather than one big one. Linkers offer more efficiency as well, as we can compile the files separately. If we change a source file we don't need to recompile all the source files, but only the changed one and then it is relinked. Besides being faster it also saves memory as common functions can be aggregated into a single file. Linkers work as follows:

- Step 1: Symbol resolution
  - Programs define and reference **symbols** (variables and functions):
    - `void swap() { ... } /* define symbol swap */`
    - `swap(); /* reference symbol swap */`
    - `int *xp = &x; /* define xp, reference x */`
  - Symbol definitions are stored (by compiler) in **symbol table**.
    - Symbol table is an array of structs
    - Each entry includes name, type, size, and location of symbol.
  - Linker associates each symbol **reference** with exactly one symbol **definition**.
- Step 2: Relocation
  - **Merges** separate code and data sections into single sections
  - Relocates symbols from their **relative** locations in the `.o` files to their final **absolute** memory locations in the executable.
  - Updates all references to these symbols to reflect their new positions.

the symbol table is located in the object file.

### 11.1 Object files

There are 3 types of object files:

- **Relocatable object file (.o file)**
  - Contains code and data in a form that can be combined with other relocatable object files to form executable object file.
  - Each .o file is produced from **exactly one source (.c file)**
- **Executable object file**
  - Contains code and data in a form that can be copied directly into memory and then executed.
- **Shared object file (.so file)**
  - Special type of relocatable object file that can be loaded into memory and linked **dynamically**, at either load time or run-time.
  - Called **Dynamic Link Libraries (DLLs)** by Windows

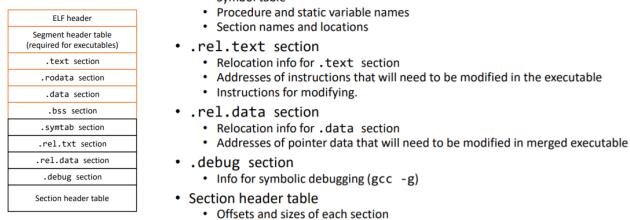
### Executable and Linkable Format (ELF)

- Standard binary format for object files
- Originally proposed by AT&T System V Unix
  - Later adopted by BSD Unix variants and Linux
- One unified format for
  - Relocatable object files (.o),
  - Executable object files
  - Shared object files (.so)
- Generic name: ELF binaries

**Sections:** refers to bits in a file **Segments:** refers to bits in memory

## ELF object file format

- Elf header
  - Word size, byte ordering, file type (.o, exec., .so), machine type, etc.
- Segment header table
  - Page size, virtual addresses memory segments (sections), segment sizes.
- .text section
- .rodata section
- .data section
- .bss section
  - Read only data: jump tables, ...
  - Initialized global variables
  - Uninitialized global variables
    - “Block Sane by Symbol”
    - “Block Sane by Address”
    - Has section header but occupies no space

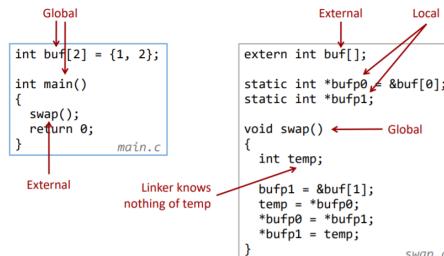


## 11.2 Linker symbols

Linker symbols can fall into one of the three categories:

- Global symbols:** Symbols defined by a module m that can be referenced by other modules e.g non-static C functions and non-static global variables
- External symbols:** Global symbols that are referenced by module m but defined by some other module. (if they are not defined it results in a linker error)
- Local symbols:** Symbols that are defined and referenced exclusively by module m e.g C functions and variables defined with the static attribute. **Local linker symbols are not local program variables**

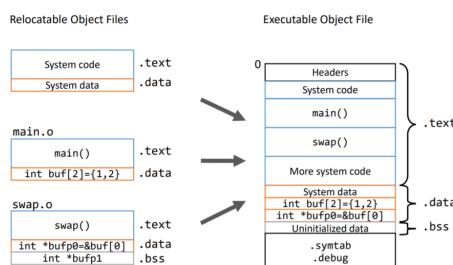
### 11.2.1 Resolving symbols



declaring extern is good practice but not necessary.

“int temp” is not a local symbol, it is a local variable which sits on the stack i.e its purely a feature of the code. There will be some information in the debug section of the object file regarding the value of temp. The linker does not need to deal with temp, because the compiler has already generated all the code to deal with it.

### 11.2.2 Relocating code and data



Each corresponding section of the individual files are merged into a single file. All of the references between the files are matched up i.e there are no more unresolved symbols. All symbols are given an address

### 11.2.3 Strong and weak symbols

**Strong:** Procedures and initialized globals **Weak:** uninitialized globals

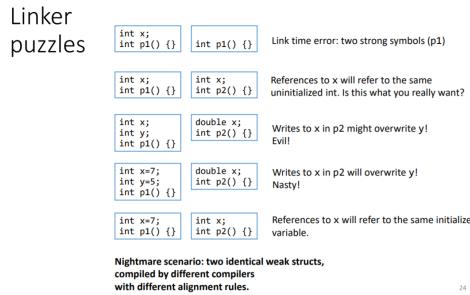


If we compile the left one by itself we get foo in the .data file. If we compile the right one by itself we get foo in the .bss file. However if we compile both together the linker thinks both foo's are the same thing

## The linker's symbol rules

1. Multiple strong symbols are **not allowed**
  - Each item can be defined only once
  - Otherwise: Linker error
2. Given a strong symbol and multiple weak symbol, **choose the strong symbol**
  - References to the weak symbol resolve to the strong symbol
3. If there are multiple weak symbols, **pick an arbitrary one**
  - Can override this with `gcc -fno-common`

Hence if we declare foo in the right as extern it means its a reference to a strong symbol somewhere else and the above rules dont reply. Hence if we use extern we will never use weak symbols.



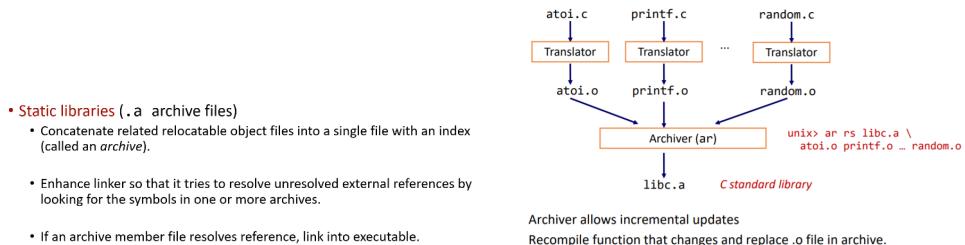
Hence avoid Global variables if possible. Otherwise use static, initialize defined global variables or use "extern" when using external global variables.

## 11.3 Static libraries

### Packaging commonly-used functions

- How to package functions commonly used by programmers?
  - Math, I/O, memory management, string manipulation, etc.
- Awkward, given the linker framework so far:
  - **Option 1:** Put all functions into a single source file
    - Programmers link big object file into their programs
    - Space and time inefficient
  - **Option 2:** Put each function in a separate source file
    - Programmers explicitly link appropriate binaries into their programs
    - More efficient, but burdensome on the programmer

**Archive:** A collection of documents **Library:** A collection of object files The solution to this is static libraries.



Commonly used libraries are:

- libc.a (the C standard library)
- libm.a (the C math library)

The standard C library is always included in the path when we search for something.  
We can include a library to the path with the gcc as follows:

`-l <name_of_library>`

e.g to include the math library we would write `-lm`. In this case the math library is already in a directory the linker already knows about. To add another directory to the search path we write:

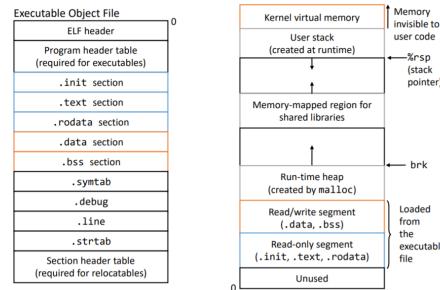
`-L <filepath>`

a “.” for the file path means the current directory

- Linker's algorithm for resolving external references:
  - Scan .o files and .a files in the command line order.
  - During the scan, keep a list of the current unresolved references.
  - As each new .o or .a file, obj, is encountered, try to resolve each unresolved reference in the list against the symbols defined in obj.
  - If any entries in the unresolved list at end of scan, then error.
- Problem:
  - Command line order matters!
  - Moral: put libraries at the end of the command line.

```
unix> gcc -L. -ltest.o -lme
unix> gcc -L. -lme libtest.o
libtest.o: In function 'main':
libtest.o(.text+0x4): undefined reference to `libFun'
```

### 11.3.1 Loading executable object files

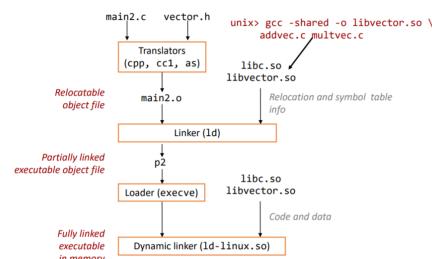


The standard C library knows where the end of the bss file is and hence knows where to start allocating memory for the heap. The loader takes the Executable Object File and distributes it in memory as seen on the right.

## 11.4 Shared libraries

- Static libraries have the following disadvantages:
  - Duplication in the stored executables (every function needs the standard libc)
  - Duplication in the running executables
  - Minor bug fixes of system libraries require each application to explicitly relink
- Solution: shared libraries
  - Object files that contain code and data that are loaded and linked into an application *dynamically*, at either *load-time* or *run-time*
  - Also called: dynamic link libraries, DLLs, .so files
- Dynamic linking can occur when executable is first loaded and run (load-time linking).
  - Common case for Linux, handled automatically by the dynamic linker (ld-linux.so).
  - Standard C Library (libc.so) usually dynamically linked.
- Dynamic linking can also occur after program has begun (run-time linking).
  - In Unix, this is done by calls to the `dlopen()` interface.
    - High-performance web servers.
    - Runtime library interpositioning
- Shared library routines can be shared by multiple processes.
  - More on this when we learn about virtual memory

### 11.4.1 Dynamic linking at load time



### 11.4.2 Dynamic linking at runtime

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

int x[2] = {1, 2};
int y[2] = {3, 4};
int z[2];

int main(int argc, char *argv[])
{
    void *handle;
    void (*addvec)(int *, int *, int *, int);
    char *error;

    /* dynamically load the shared lib that contains addvec() */
    handle = dlopen("./libvector.so", RTLD_LAZY);
    if (!handle) {
        fprintf(stderr, "%s\n", dlerror());
        exit(1);
    }
    ...
}

/* get a pointer to the addvec() function we just loaded */
addvec = dlsym(handle, "addvec");
if ((error = dlerror()) != NULL) {
    fprintf(stderr, "%s\n", error);
    exit(1);
}

/* Now we can call addvec() it just like any other function */
addvec(x, y, z, 2);
printf("z = [%d %d]\n", z[0], z[1]);

/* unload the shared library */
if (dlclose(handle) < 0) {
    fprintf(stderr, "%s\n", dlerror());
    exit(1);
}
return 0;
```

# Chapter 12

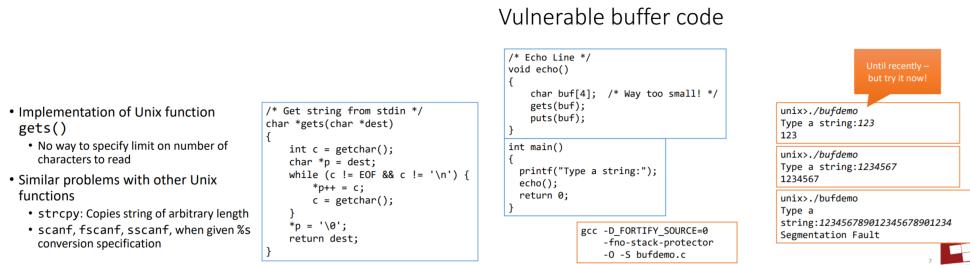
## Code Vulnerabilities

### 12.1 Worms and Viruses

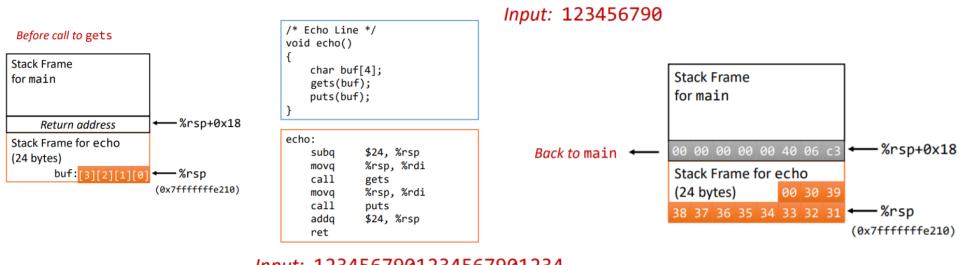
**Worm:** A program that can run by itself, it can propagate a full working version of itself to other computers

**Virus:** Code that can add itself to other programs, it cannot run independently

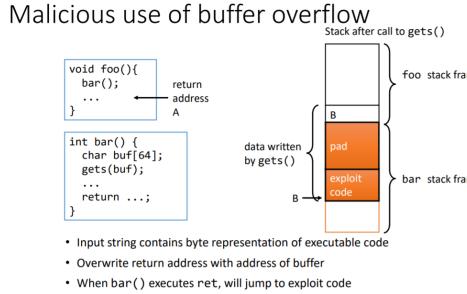
### 12.2 Stack Overflow bug



As soon as the string gets large enough, we get a segmentation fault. The program doesn't have anything obvious that crashes itself, it doesn't corrupt/manipulate its own memory. What happens is that the input can be used to subvert a poorly written program e.g `gets` in order to corrupt memory and cause a segmentation fault. These days modern compilers have protection against this particular attack.



As soon as we fill up the stack frame and add more input, the return address of the buffer becomes random. So rather than causing the program to crash we can make it jump and execute some other code. If we know the address that the stack pointer points to, we can cause the return address to jump into the code we have just added to the stack frame. It was possible to find out the stack pointer before address randomization, because you could test it out on your own machine first, and hence all machines of the same type would have the same stack pointer.



### Exploits based on buffer overflows

- *Buffer overflow bugs allow remote machines to execute arbitrary code on victim machines*
- Internet worm (one vector)
  - Early versions of the finger server (fingerd) used gets() to read the argument sent by the client:
    - finger droh@cs.cmu.edu
  - Worm attacked fingerd server by sending bad argument:
    - finger "exploit-code padding new-return-address"
    - exploit code: executed a root shell on the victim machine with a direct TCP connection to the attacker.

## 12.3 Stopping overrun bugs

```
/* Echo Line */
void echo()
{
    char buf[4]; /* Way too small! */
    fgets(buf, 4, stdin);
    puts(buf);
}
```

- Use library routines that limit string lengths
  - fgets instead of gets
  - strncpy instead of strcpy
  - Don't use scanf with %s conversion specification
    - Use fgets to read the string
    - Or use %ns - where n is a suitable integer

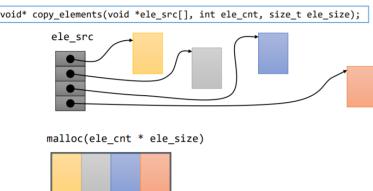
### System-level protections

- Compiler-inserted checks on functions
  - Compiler now understands library calls...
- Randomized stack offsets
  - At start of program, allocate random amount of space on stack
  - Makes it difficult to predict beginning of inserted code
- Nonexecutable code segments
  - In older x86, can mark region of memory as either "read-only" or "writable"
    - Can execute anything readable
  - Add explicit "execute" permission to hardware

## 12.4 XDR vulnerability

### SUN XDR library

- Widely used for transferring data between machines (e.g. for NFS)



### XDR code

```
void* copy_elements(void *ele_src[], int ele_cnt, size_t ele_size) {
    /*
     * Allocate buffer for ele_cnt objects, each of ele_size bytes
     * and copy from locations designated by ele_src
     */
    void *result = malloc(ele_cnt * ele_size);
    if (result == NULL)
        /* malloc failed */
        return NULL;
    void *next = result;
    int i;
    for (i = 0; i < ele_cnt; i++) {
        /* Copy object i to destination */
        memcpy(next, ele_src[i], ele_size);
        /* Move pointer to next memory region */
        next += ele_size;
    }
    return result;
}
```

### XDR vulnerability

```
malloc(ele_cnt * ele_size)
```

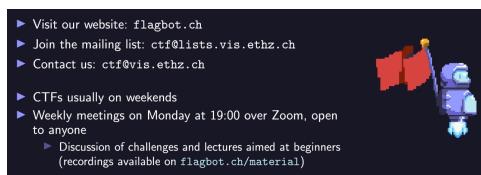
- What if (on a 32-bit machine):

- ele\_cnt =  $2^{20} + 1$
- ele\_size =  $4096 = 2^{12}$
- Allocation = ??

- How can I make this function secure?

Allocating a large amount of memory would most likely fail with malloc, but since this is a large unsigned value (signed \* unsigned = unsigned) it wraps around in this case  $2^{20} + 1 = 1$  which on a 32-bit machine would allocate 32-bits. So malloc would return true and we would end up overwriting an enormous amount of data.

## 12.5 CTF



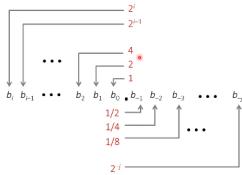
## Chapter 13

# Floating Point

### 13.1 Representing Floating Point

### 13.1.1 Fractional binary numbers

- Representation
    - Bits to right of “binary point” represent fractional powers of 2
    - Represents rational number:



- Observations
    - Divide by 2 by shifting right
    - Multiply by 2 by shifting left
    - Numbers of form  $0.111111\dots$ , are just below 1.0
      - $1/2 + 1/4 + 1/8 + \dots + 1/2^i + \dots \rightarrow 1.0$
      - Use notation  $1.0 - \epsilon$

Value	Representation
$5\frac{3}{4}$	$101.11_2$
$2\frac{7}{8}$	$10.111_2$
$6\frac{3}{64}$	$0.111111_2$

This notation has its limitation, we can only exactly represent numbers of the form  $\frac{x}{2^k}$ . Rational numbers have repeating bit patterns.

Value	Representation
1/3	$0.0101010101[01]..._2$
1/5	$0.001100110011[0011]..._2$
1/10	$0.0001100110011[0011]..._2$

You can't represent 0.110

### 13.1.2 Floating point representation

## Floating point representation (recap from Digital Circuits)

- Numerical form:  
$$(-1)^s \cdot M \cdot 2^E$$
    - Sign bit s determines whether number is negative or positive
    - Significand M normally a fractional value in range [1.0, 2.0).
    - Exponent E weights value by power of two
  - Encoding
    - MSB (Most Significant Bit) s is sign bit s
    - exp field encodes E (but is not equal to E)
    - frac field encodes M (but is not equal to M)

## 13.2 Types of IEEE floating point numbers

### 13.2.1 Precisions

IEEE 754 standards (and others...)

Precision	Significant bits	Exponent bits	Total
Half	11	5	16
Single	24	8	32
Double	53	11	64
Quadruple	113	15	128
Octuple	237	19	256
Google bfloat16	7	8	16
Nvidia TensorFloat	10	8	19
AMD fp24	17	7	24

### 13.2.2 Floating point in C

- C99 guarantees two levels:
  - `float` - single precision
  - `double` - double precision
- `long double` can be double, extended, or quadruple precision
- Conversions/casting
  - Casting between `int`, `float`, and `double` changes bit representation
  - `double/float → int`
    - Truncates fractional part (like rounding toward zero)
    - Not defined when out of range or NaN. Generally sets to TMin
  - `int → double`
    - Exact conversion, as long as `int` has ≤ 53 bit word size
    - `int → float`
      - Will round according to rounding mode

### 13.2.3 Normalized Values

- Condition:  $\exp \neq 000\ldots 0$  and  $\exp \neq 111\ldots 1$
- Exponent coded as **biased value**:  $E = \text{Exp} - \text{Bias}$ 
  - $\text{Exp}$ : unsigned value exp
  - $\text{Bias} = 2^{e-1} - 1$ , where  $e$  is number of exponent bits
    - Single precision: 127 ( $\text{Exp}: 1\ldots 254$ ,  $E: -126\ldots 127$ )
    - Double precision: 1023 ( $\text{Exp}: 1\ldots 2046$ ,  $E: -1022\ldots 1023$ )
- Significand coded with implied **leading 1**:  $M = 1.\text{xxx...x}_2$ 
  - $\text{xxx...x}$ : bits of frac
  - Minimum when  $000\ldots 0$  ( $M = 1.0$ )
  - Maximum when  $111\ldots 1$  ( $M = 2.0 - \epsilon$ )
  - Get extra leading bit for “free”
- Value:  $\text{float } F = 15213.0;$   
 $15213_{10} = 11101101101101_2 = 1.1101101101101_2 \times 2^{13}$
- Significand
  - $M = 1.\underline{\text{1101101101101}}.$
  - $\text{frac} = \underline{\text{1101101101101}}0000000000_2$
- Exponent
  - $E = 13$
  - $\text{Bias} = 127$
  - $\text{Exp} = 140 = 10001100_2$
- Result: 
$$\begin{array}{c|c|c} 0 & 10001100 & 11011011011010000000000 \\ \hline s & \text{exp} & \text{frac} \end{array}$$

### 13.2.4 Denormalized values

- Condition:  $\exp = 000\ldots 0$
- Exponent value:  $E = -\text{Bias} + 1$  (instead of  $E = 0 - \text{Bias}$ )
- Significand coded with implied **leading 0**:  $M = 0.\text{xxx...x}_2$ 
  - $\text{xxx...x}$ : bits of frac
- Cases
  - $\exp = 000\ldots 0$ ,  $\text{frac} = 000\ldots 0$ 
    - Represents value 0
    - Note distinct values: +0 and -0 (why?)
  - $\exp = 000\ldots 0$ ,  $\text{frac} \neq 000\ldots 0$ 
    - Numbers very close to 0.0
    - Lose precision as get smaller
    - Equispaced

There are 2 distinct values for 0 because of the sign bit

Description	exp	frac	Numerical value	{single,double}
Zero	00..00	00..00	0.0	
Smallest pos. denorm.	00..00	00..01	$2^{-[2^{13}+1]} \times 2^{-[126,1022]}$	
Largest denormalized	00..00	11..11	$(1.0 - \epsilon) \times 2^{-[126,1022]}$	
Smallest pos. normalized	00..01	00..00	$2^{-[126,1022]}$	
One	01..11	00..00	1.0	
Largest normalized	11..10	11..11	$(2.0 - \epsilon) \times 2^{[127,1023]}$	

### 13.2.5 Special values

#### Special values

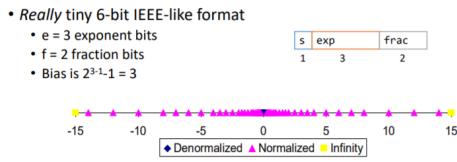
- Condition:  $\text{exp} = 111\dots1$
- Case:  $\text{exp} = 111\dots1, \text{frac} = 000\dots0$ 
  - Represents value  $\infty$  (infinity)
  - Operation that overflows
  - Both positive and negative
  - E.g.  $1.0/0.0 = -1.0/0.0 = +\infty, 1.0/-0.0 = -\infty$
- Case:  $\text{exp} = 111\dots1, \text{frac} \neq 000\dots0$ 
  - Not-a-Number (NaN)
  - Represents case when no numeric value can be determined
  - E.g.,  $\sqrt{-1}, \infty - \infty, \infty * 0$

### 13.2.6 Properties of encoding

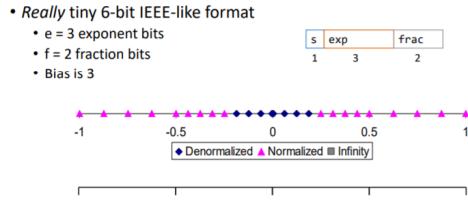
- FP zero same as integer zero
  - All bits = 0 (for +0)
- Can (almost) use unsigned integer comparison
  - Must first compare sign bits
  - Must consider  $-0 = 0$
  - NaNs problematic
    - Will be greater than any other values
    - What should comparison yield?
  - Otherwise OK
    - Denorm vs. normalized
    - Normalized vs. infinity

## 13.3 Floating point ranges

#### Distribution of values



#### Distribution of values (close-up view)



Denormalized numbers have equal spacing around zero because they have the same E value.

## 13.4 Floating-point rounding and arithmetic

### 13.4.1 Rounding

- Rounding modes (illustrate with CHF rounding)

	CHF 1.40	CHF 1.60	CHF 1.50	CHF 2.50	CHF -1.50
Towards zero	1	1	1	2	-1
Round down ( $-\infty$ )	1	1	1	2	-2
Round up ( $+\infty$ )	2	2	2	3	-1
Nearest Even (default)	1	2	2	2	-2

Rounding to nearest even avoids biased rounding. All other rounding types are statistically biased i.e a sum of set of positive numbers will consistently be over- or under- estimated.

- Applying to other decimal places / bit positions
  - When exactly halfway between two possible values
    - Round so that least significant digit is even
  - E.g., round to nearest hundredth

Value	Result	Description
1.2349999	1.23	(less than half way)
1.2350001	1.24	(greater than half way)
1.2350000	1.24	(half-way – round up)
1.2450000	1.24	(half way – round down)

### 13.4.2 Rounding binary numbers

- Binary fractional numbers
  - “Even” when least significant bit is 0
  - “Half way” when bits to right of rounding position =  $100\dots_2$
- Examples
  - Round to nearest  $1/4$  (2 bits right of binary point)

Value	Binary	Rounded	Action	Result
$2\frac{3}{16}$	$10.00011_2$	$10.00_2$	$< \frac{1}{2} : \text{down}$	2
$2\frac{7}{16}$	$10.00110_2$	$10.01_2$	$> \frac{1}{2} : \text{up}$	$2\frac{1}{4}$
$2\frac{1}{8}$	$10.11100_2$	$11.00_2$	$= \frac{1}{2} : \text{up}$	3
$2\frac{1}{8}$	$10.10100_2$	$10.10_2$	$= \frac{1}{2} : \text{down}$	$2\frac{1}{8}$

### 13.4.3 Creating a floating point number

#### Normalize

- Steps
  - Normalize to have leading 1
  - Round to fit within fraction
  - Postnormalize to deal with effects of rounding
- Case study
  - Convert 8-bit unsigned numbers to tiny floating point format

Value	Binary
128	10000000
15	00001101
17	00010001
19	00010011
138	10001010
63	00111111

#### Requirement

Value	Binary	Fraction	Exponent
128	10000000	1.0000000	7
15	00001101	1.1910000	3
17	00010001	1.0001000	4
19	00010011	1.0011000	4
138	10001010	1.0001010	7
63	00111111	1.1111110	5

#### Postnormalize

#### Rounding

1. BBGRXXX
- 
- Round up conditions
    - Round = 1, Sticky = 1  $\Rightarrow > 0.5$
    - Guard = 1, Round = 1, Sticky = 0  $\Rightarrow$  Round to even
  - Round down conditions
    - Round = 0, Sticky = 1  $\Rightarrow < 0.5$
    - Guard = 0, Round = 0, Sticky = 0  $\Rightarrow$  Round to odd

Value	Fraction	GRS	Incr?	Rounded
128	1.0000000	000	N	1.000
15	1.1010000	100	N	1.101
17	1.0001000	010	N	1.000
19	1.0011000	110	Y	1.010
138	1.0001010	011	Y	1.001
63	1.1111110	111	Y	10.000

Value	Rounded	Exp	Adjusted	Result
128	1.000	7		128
15	1.101	3		15
17	1.000	4		16
19	1.010	4		20
138	1.001	7		134
63	10.000	5	1.000/6	64

## 13.5 Floating point addition and multiplication

### 13.5.1 Floating point multiplication

#### Mathematical properties of floating point multiplication

- Exact Result:  $(-1)^s M \cdot 2^E$ 
  - Sign  $s$ :  $S_1 \wedge S_2$
  - Significand  $M$ :  $M_1 \cdot M_2$
  - Exponent  $E$ :  $E_1 + E_2$
- Fixing
  - If  $M \geq 2$ , shift  $M$  right, increment  $E$
  - If  $M$  out of range, overflow
  - Round  $M$  to fit  $\text{frac}$  precision
- Implementation
  - Biggest chore is multiplying significands

Yes

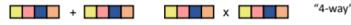
No

Yes

## 13.6 SSE floating point

### Vector instructions: SSE family

- SIMD (single-instruction, multiple data) vector instructions
  - New data types, registers, operations
  - Parallel operation on small (length 2-8) vectors of integers or floats
  - Example:

 "4-way"

- Floating point vector instructions
  - Available with Intel's SSE (streaming SIMD extensions) family
  - SSE starting with Pentium III: 4-way single precision
  - SSE2 starting with Pentium 4: 2-way double precision
  - All x86-64 have SSE3 (superset of SSE2, SSE)

Here we focus on SSE3:

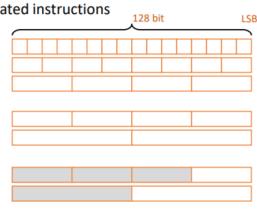
#### SSE3 registers

- All caller saved
- %xmm0 for floating point return value

128 bit = 2 doubles = 4 singles	
%xmm0	Argument #1
%xmm1	Argument #2
%xmm2	Argument #3
%xmm3	Argument #4
%xmm4	Argument #5
%xmm5	Argument #6
%xmm6	Argument #7
%xmm7	Argument #8
%xmm8	
%xmm9	
%xmm10	
%xmm11	
%xmm12	
%xmm13	
%xmm14	
%xmm15	

#### SSE3 registers

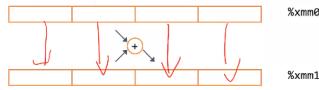
- Different data types and associated instructions
  - Integer vectors:
    - 16-way byte
    - 8-way 2 bytes
    - 4-way 4 bytes
- Floating point vectors:
  - 4-way single
  - 2-way double
- Floating point scalars:
  - single
  - double



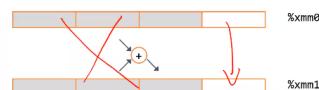
Example:

#### SSE3 instructions: examples

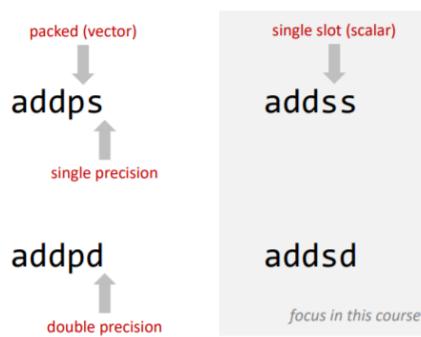
- Single precision **4-way vector add**: addps %xmm0 %xmm1



- Single precision **scalar add**: addss %xmm0 %xmm1



### 13.6.1 SSE3 instruction names



### SSE3 basic instructions

	Single	Double	Effect
Moves:	movss	movsd	D $\leftarrow$ S
Arithmetic:			
	addss	addsd	D $\leftarrow$ D + S
	subss	subsd	D $\leftarrow$ D - S
	mulss	mulsd	D $\leftarrow$ D $\times$ S
	divss	divsd	D $\leftarrow$ D / S
	maxss	maxsd	D $\leftarrow$ max(D,S)
	minss	minsd	D $\leftarrow$ min(D,S)
	sqrtss	sqrtsd	D $\leftarrow$ sqrt(S)

### SSE3 conversion instructions

- Conversions
  - Same operand forms as moves

Instruction	Description
cvtss2sd	single $\rightarrow$ double
cvtsd2ss	double $\rightarrow$ single
cvtz2ss	int $\rightarrow$ single
cvtz2sd	int $\rightarrow$ double
cvtz2ssq	quad int $\rightarrow$ single
cvtz2sdq	quad int $\rightarrow$ double
cvtss2si	single $\rightarrow$ int (truncation)
cvttsd2si	double $\rightarrow$ int (truncation)
cvtss2siq	single $\rightarrow$ quad int (truncation)
cvttsd2siq	double $\rightarrow$ quad int (truncation)

### 13.6.2 SSE3 Examples

#### x86-64 FP code example

##### x86-64 FP code example

- Compute inner product of two vectors
  - Single precision arithmetic
  - Uses SSE3 instructions

```
float ipf (float x[], float y[], int n)
{
    int i;
    float result = 0.0;
    for (i = 0; i < n; i++) {
        result += x[i]*y[i];
    }
    return result;
}
```

```
ipf:
    xorps %xmm1, %xmm1          # result = 0.0
    xorl %eax, %eax
    jmp .L8
.L10:
    movaps (%rdi,%rax,4), %xmm0
    movaps (%rsi,%rax,4), %xmm1
    addps %xmm0, %xmm1
    incl %eax
    movss %xmm1, %eax
    addl %eax, %eax
    .L8:
    cmpl %eax, %eax
    jne .L10
    movaps %xmm1, %xmm0
    ret
```

```
double funct(double a, float x, double b, int i)
{
    return a*x - b/i;
}
```

```
a %xmm0 double
x %xmm1 float
b %xmm2 double
i %edi int

funct:
    cvtss2sd %xmm1, %xmm1      # %xmm1 = (double) x
    mulsd %xmm0, %xmm1          # %xmm1 = a*x
    cvtsi2sd %edi, %xmm0       # %xmm0 = (double) i
    divsd %xmm0, %xmm2          # %xmm2 = b/i
    movsd %xmm1, %xmm0          # %xmm0 = a*x
    subsd %xmm2, %xmm0          # return result
    ret
```

### 13.6.3 Constants

```
double cel2fahr(double temp)
{
    return 1.8 * temp + 32.0;
}

# Constant declarations
.LC2:
.long 3435973837      # Low order four bytes of 1.8
.long 1073532108      # High order four bytes of 1.8
.LC4:
.long 0                 # Low order four bytes of 32.0
.long 1077936128      # High order four bytes of 32.0

# Code
cel2fahr:
    mulsd .LC2(%rip), %xmm0 # Multiply by 1.8
    addsd .LC4(%rip), %xmm0 # Add 32.0
    ret
```

- Here: Constants in decimal format
  - compiler decision
  - hex more readable

- Previous slide: Claim

```
.LC4:
.long 0                  # Low order four bytes of 32.0
.long 1077936128        # High order four bytes of 32.0
```

- Convert to hex format:

```
.LC4:
.long 0x0                # Low order four bytes of 32.0
.long 0x40400000          # High order four bytes of 32.0
```

- Convert to double:

- Remember: e = 11 exponent bits, bias =  $2^{e-1}-1 = 1023$

The top bit is the sign bit, the next 11 bits are the exp field, hence the first 12 bits are the exponent which correspond to 3 hex digits. i.e 0x404 = 1028 is the exponent, since the bias is 1023 we have  $1028-1023 = 5$ , so the actual exponent value is 5. The rest is the significand which is all zeroes. There is an implied leading one. Hence the value is  $1 \cdot 2^5$

#### 13.6.4 Vector Instructions

- Recently, `gcc` can autovectorize to some extent
  - e.g. `-ftree-vectorize`
  - No speed-up guaranteed
  - Very limited
  - `icc` is usually much better
- For highest performance vectorize yourself using *intrinsics*
  - Intrinsics = C interface to vector instructions
- Most recently:
  - Intel AVX512: 8-way double, 16-way single

## Chapter 14

# Optimizing Compilers

### 14.1 Optimizing Compilers

- Use optimization flags,  
default can be no optimization (-O0)!
- Good choices for gcc:  
-O2, -O3, -march=xxx, -m64
- Try different flags and maybe different  
compilers
  - icc is often faster than gcc

Default is O0 because it offers faster compilation. O1,O2,O3 are more optimization respectively but higher number does not guarantee that there is better performance. The compiler only does optimization which preserves the functionality of the code Example:

#### Example

```
double a[4][4];
double b[4][4];
double c[4][4]; // set to zero

/* Multiply 4 x 4 matrices a and b */
void mmn(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < 4; i++) {
        for (j = 0; j < 4; j++)
            for (k = 0; k < 4; k++)
                c[i*4+j] += a[i*4+k]*b[k*4 + j];
    }
}
```

• -O3 -m64 -march=haswell -fno-tree-vectorize: **8 cycles** (yes, 8!)

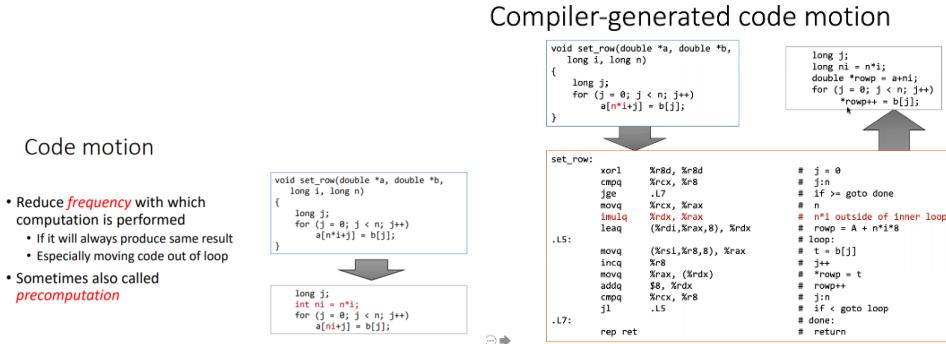
#### Strengths and Weaknesses

- Compilers are **good** at: mapping program to machine
  - register allocation
  - code selection and ordering (scheduling)
  - dead code elimination
  - eliminating minor inefficiencies
- Compilers are **not good** at: improving asymptotic efficiency
  - up to programmer to select best overall algorithm
  - big-O savings are (often) more important than constant factors
    - but constant factors also matter
- Compilers are **not good** at: overcoming “optimization blockers”
  - potential memory aliasing
  - potential procedure side-effects

#### Limitations of optimizing compilers

- **If in doubt, the compiler is conservative**
- Operate under fundamental constraints
  - Must not change program behavior under any possible condition
  - Often prevents it from making optimizations when would only affect behavior under pathological conditions.
- Behavior that may be obvious to the programmer can be obfuscated by languages and coding styles
  - e.g., data ranges may be more limited than variable types suggest
- Most analysis is performed only within procedures
  - Whole-program analysis is too expensive in most cases
- Most analysis is based only on static information
  - Compiler has difficulty anticipating run-time inputs

## 14.2 Code motion and precomputation



## 14.3 Strength reduction

Example: sequence of products



## Strength reduction

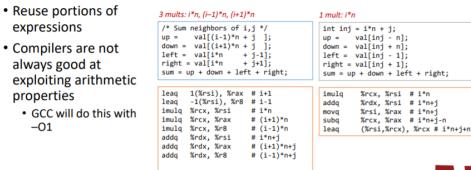
- Replace costly operation with simpler one
  - Usually more specialized ("less strong")
  - Prior example: Shift/add instead of multiply or divide

$$16*x \rightarrow x \ll 4$$

- Usefulness is machine-dependent
  - Depends on cost of multiply or divide instruction
  - On Pentium IV, integer multiply is 10 cycles

## 14.4 Common subexpressions

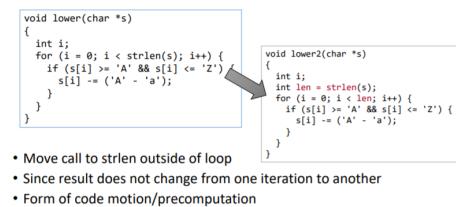
Share common subexpressions



- Compilers can eliminate some common subexpressions
  - But not all! (though they are improving)
  - Some are blocked (see later)
  - Remember: some operations **not associative!**
- If in doubt:
  - Check the resulting assembly code
  - Help the compiler out yourself

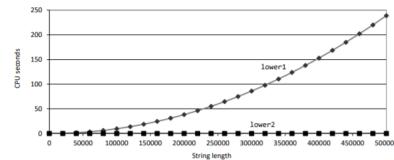
## 14.5 Optimization blocker: Procedure calls

Improving performance



## Lower Case Conversion Performance

- Time doubles when double string length
- Linear performance of lower2



the compiler treats procedures as black boxes and does not replace it with its result, even if its called multiple times in a loop because it might not return the same result given the same inputs, as it could be dependant on global states. Hence in the above example we have to optimize the code because the compiler wont do it.

## 14.6 Optimization blocker: memory aliasing

When two different memory references specify a single location

- Easy to have happen in C
  - Since allowed to do address arithmetic
  - Direct access to storage structures
- Get in habit of introducing local variables
  - Accumulating within loops
  - Your way of telling compiler not to check for aliasing

Example:

### Possible aliasing

- Memory accessed  
⇒ compiler assumes possible side effects
- ```
/* Sum rows of n x n matrix a and store in vector b */
void sum_rows1(double *a, double *b, long n) {
    long i, j;
    for (i = 0; i < n; i++) {
        for (j = 0; j < n; j++) {
            b[j] = 0;
            for (j = 0; j < n; j++)
                b[j] += a[i*n + j];
        }
    }
}
```

Value of B:  
 init: [4, 8, 16]  
 i = 0: [3, 8, 16]  
 i = 1: [3, 22, 16]  
 i = 2: [3, 22, 224]

In memory, we have...

|                 |      |     |
|-----------------|------|-----|
| 0xfffffbef0f868 | A[0] | 0   |
| 0xfffffbef0f870 | A[1] | 1   |
| 0xfffffbef0f878 | A[2] | 2   |
| 0xfffffbef0f880 | A[3] | 3   |
| 0xfffffbef0f888 | A[4] | 22  |
| 0xfffffbef0f890 | B[0] | 0   |
| 0xfffffbef0f898 | B[1] | 22  |
| 0xfffffbef0f8a0 | B[2] | 224 |

### Unaliased version when aliasing happens

- Aliasing still creates interference

- Result different from before

```
/* Sum rows of n X n matrix a
   and store in vector b */
void sum_rows2(double *a, double *b, long n) {
    long i, j, k;
    for (i = 0; i < n; i++) {
        double val = 0;
        for (j = 0; j < n; j++)
            val += a[i*n + j];
        b[i] = val;
    }
}
```

Value of B:  
 init: [4, 8, 16]  
 i = 0: [3, 8, 16]  
 i = 1: [3, 27, 16]  
 i = 2: [3, 27, 224]

What happens in the possible aliasing scenario is that B[0] is overwritten by the first three elements of A becoming 3, then as we sum up the second row of A to calculate B[1], we set B[1] = A[4] to 0 and add A[3]=3 + A[4] (which is now 3) + A[5]=16 hence we get 22. The compiler doesn't optimize the code because it assumes this is the behaviour we want. We can remove aliasing by using temporary variables in this case copying the array elements into "val". The unaliased version still creates some interference because we are still overwriting some memory.

## 14.7 Blocking and unrolling

### More difficult example

- Matrix multiplication:  
 $C = A \cdot B + C$
- Which array elements are reused?
- All of them! But how to take advantage?

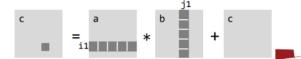
```
c = (double *) calloc(sizeof(double), n*n);
/* Multiply n x n matrices a and b */
void mm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i++) {
        for (j = 0; j < n; j++) {
            for (k = 0; k < n; k++)
                c[i*n+j] += a[i*n+k]*b[k*n+j];
        }
    }
}
```



### Step 1: Blocking (here: 2 x 2)

- Blocking
  - also called tiling
  - = partial unrolling + loop exchange
  - Assumes associativity
    - ⇒ compiler will never do it

```
c = (double *) calloc(sizeof(double), n*n);
/* Multiply n x n matrices a and b */
void mm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i+=2)
        for (j = 0; j < n; j+=2)
            for (k = 0; k < n; k+=2)
                for (i1 = i; i1 < i+2; i1++)
                    for (j1 = j; j1 < j+2; j1++)
                        for (k1 = k; k1 < k+2; k1++)
                            c[i1*n+j1] += a[i1*n+k1]*b[k1*n+j1];
}
```



### Step 2: Unrolling inner loops

- Every array element a[...], b[...], c[...] used twice
- Now scalar replacement can be applied

```
c = (double *) calloc(sizeof(double), n*n);
/* Multiply n x n matrices a and b */
void mm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i+=2)
        for (j = 0; j < n; j+=2)
            for (k = 0; k < n; k+=2)
                <body>
}
<body>
c[i*n+j] = a[i*n+k]*b[k*n+j] + a[i*n+k+1]*b[(k+1)*n+j]
c[i*(n+1)+j] = a[i*(n+1)*k]*b[k*n+j] + a[i*(n+1)*k+1]*b[(k+1)*n+j]
c[i*(n+1)+(j+1)] = a[i*(n+1)*k]*b[(k+1)*n+(j+1)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+1)]
c[i*(n+1)*n+(j+1)] = a[i*(n+1)*k]*b[(k+1)*n+(j+1)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+1)]
c[i*(n+1)*n+(j+2)] = a[i*(n+1)*k]*b[(k+1)*n+(j+2)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+2)]
c[i*(n+1)*n+(j+3)] = a[i*(n+1)*k]*b[(k+1)*n+(j+3)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+3)]
c[i*(n+1)*n+(j+4)] = a[i*(n+1)*k]*b[(k+1)*n+(j+4)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+4)]
c[i*(n+1)*n+(j+5)] = a[i*(n+1)*k]*b[(k+1)*n+(j+5)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+5)]
c[i*(n+1)*n+(j+6)] = a[i*(n+1)*k]*b[(k+1)*n+(j+6)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+6)]
c[i*(n+1)*n+(j+7)] = a[i*(n+1)*k]*b[(k+1)*n+(j+7)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+7)]
c[i*(n+1)*n+(j+8)] = a[i*(n+1)*k]*b[(k+1)*n+(j+8)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+8)]
c[i*(n+1)*n+(j+9)] = a[i*(n+1)*k]*b[(k+1)*n+(j+9)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+9)]
c[i*(n+1)*n+(j+10)] = a[i*(n+1)*k]*b[(k+1)*n+(j+10)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+10)]
c[i*(n+1)*n+(j+11)] = a[i*(n+1)*k]*b[(k+1)*n+(j+11)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+11)]
c[i*(n+1)*n+(j+12)] = a[i*(n+1)*k]*b[(k+1)*n+(j+12)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+12)]
c[i*(n+1)*n+(j+13)] = a[i*(n+1)*k]*b[(k+1)*n+(j+13)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+13)]
c[i*(n+1)*n+(j+14)] = a[i*(n+1)*k]*b[(k+1)*n+(j+14)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+14)]
c[i*(n+1)*n+(j+15)] = a[i*(n+1)*k]*b[(k+1)*n+(j+15)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+15)]
c[i*(n+1)*n+(j+16)] = a[i*(n+1)*k]*b[(k+1)*n+(j+16)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+16)]
c[i*(n+1)*n+(j+17)] = a[i*(n+1)*k]*b[(k+1)*n+(j+17)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+17)]
c[i*(n+1)*n+(j+18)] = a[i*(n+1)*k]*b[(k+1)*n+(j+18)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+18)]
c[i*(n+1)*n+(j+19)] = a[i*(n+1)*k]*b[(k+1)*n+(j+19)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+19)]
c[i*(n+1)*n+(j+20)] = a[i*(n+1)*k]*b[(k+1)*n+(j+20)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+20)]
c[i*(n+1)*n+(j+21)] = a[i*(n+1)*k]*b[(k+1)*n+(j+21)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+21)]
c[i*(n+1)*n+(j+22)] = a[i*(n+1)*k]*b[(k+1)*n+(j+22)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+22)]
c[i*(n+1)*n+(j+23)] = a[i*(n+1)*k]*b[(k+1)*n+(j+23)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+23)]
c[i*(n+1)*n+(j+24)] = a[i*(n+1)*k]*b[(k+1)*n+(j+24)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+24)]
c[i*(n+1)*n+(j+25)] = a[i*(n+1)*k]*b[(k+1)*n+(j+25)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+25)]
c[i*(n+1)*n+(j+26)] = a[i*(n+1)*k]*b[(k+1)*n+(j+26)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+26)]
c[i*(n+1)*n+(j+27)] = a[i*(n+1)*k]*b[(k+1)*n+(j+27)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+27)]
c[i*(n+1)*n+(j+28)] = a[i*(n+1)*k]*b[(k+1)*n+(j+28)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+28)]
c[i*(n+1)*n+(j+29)] = a[i*(n+1)*k]*b[(k+1)*n+(j+29)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+29)]
c[i*(n+1)*n+(j+30)] = a[i*(n+1)*k]*b[(k+1)*n+(j+30)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+30)]
c[i*(n+1)*n+(j+31)] = a[i*(n+1)*k]*b[(k+1)*n+(j+31)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+31)]
c[i*(n+1)*n+(j+32)] = a[i*(n+1)*k]*b[(k+1)*n+(j+32)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+32)]
c[i*(n+1)*n+(j+33)] = a[i*(n+1)*k]*b[(k+1)*n+(j+33)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+33)]
c[i*(n+1)*n+(j+34)] = a[i*(n+1)*k]*b[(k+1)*n+(j+34)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+34)]
c[i*(n+1)*n+(j+35)] = a[i*(n+1)*k]*b[(k+1)*n+(j+35)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+35)]
c[i*(n+1)*n+(j+36)] = a[i*(n+1)*k]*b[(k+1)*n+(j+36)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+36)]
c[i*(n+1)*n+(j+37)] = a[i*(n+1)*k]*b[(k+1)*n+(j+37)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+37)]
c[i*(n+1)*n+(j+38)] = a[i*(n+1)*k]*b[(k+1)*n+(j+38)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+38)]
c[i*(n+1)*n+(j+39)] = a[i*(n+1)*k]*b[(k+1)*n+(j+39)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+39)]
c[i*(n+1)*n+(j+40)] = a[i*(n+1)*k]*b[(k+1)*n+(j+40)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+40)]
c[i*(n+1)*n+(j+41)] = a[i*(n+1)*k]*b[(k+1)*n+(j+41)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+41)]
c[i*(n+1)*n+(j+42)] = a[i*(n+1)*k]*b[(k+1)*n+(j+42)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+42)]
c[i*(n+1)*n+(j+43)] = a[i*(n+1)*k]*b[(k+1)*n+(j+43)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+43)]
c[i*(n+1)*n+(j+44)] = a[i*(n+1)*k]*b[(k+1)*n+(j+44)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+44)]
c[i*(n+1)*n+(j+45)] = a[i*(n+1)*k]*b[(k+1)*n+(j+45)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+45)]
c[i*(n+1)*n+(j+46)] = a[i*(n+1)*k]*b[(k+1)*n+(j+46)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+46)]
c[i*(n+1)*n+(j+47)] = a[i*(n+1)*k]*b[(k+1)*n+(j+47)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+47)]
c[i*(n+1)*n+(j+48)] = a[i*(n+1)*k]*b[(k+1)*n+(j+48)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+48)]
c[i*(n+1)*n+(j+49)] = a[i*(n+1)*k]*b[(k+1)*n+(j+49)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+49)]
c[i*(n+1)*n+(j+50)] = a[i*(n+1)*k]*b[(k+1)*n+(j+50)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+50)]
c[i*(n+1)*n+(j+51)] = a[i*(n+1)*k]*b[(k+1)*n+(j+51)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+51)]
c[i*(n+1)*n+(j+52)] = a[i*(n+1)*k]*b[(k+1)*n+(j+52)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+52)]
c[i*(n+1)*n+(j+53)] = a[i*(n+1)*k]*b[(k+1)*n+(j+53)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+53)]
c[i*(n+1)*n+(j+54)] = a[i*(n+1)*k]*b[(k+1)*n+(j+54)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+54)]
c[i*(n+1)*n+(j+55)] = a[i*(n+1)*k]*b[(k+1)*n+(j+55)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+55)]
c[i*(n+1)*n+(j+56)] = a[i*(n+1)*k]*b[(k+1)*n+(j+56)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+56)]
c[i*(n+1)*n+(j+57)] = a[i*(n+1)*k]*b[(k+1)*n+(j+57)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+57)]
c[i*(n+1)*n+(j+58)] = a[i*(n+1)*k]*b[(k+1)*n+(j+58)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+58)]
c[i*(n+1)*n+(j+59)] = a[i*(n+1)*k]*b[(k+1)*n+(j+59)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+59)]
c[i*(n+1)*n+(j+60)] = a[i*(n+1)*k]*b[(k+1)*n+(j+60)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+60)]
c[i*(n+1)*n+(j+61)] = a[i*(n+1)*k]*b[(k+1)*n+(j+61)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+61)]
c[i*(n+1)*n+(j+62)] = a[i*(n+1)*k]*b[(k+1)*n+(j+62)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+62)]
c[i*(n+1)*n+(j+63)] = a[i*(n+1)*k]*b[(k+1)*n+(j+63)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+63)]
c[i*(n+1)*n+(j+64)] = a[i*(n+1)*k]*b[(k+1)*n+(j+64)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+64)]
c[i*(n+1)*n+(j+65)] = a[i*(n+1)*k]*b[(k+1)*n+(j+65)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+65)]
c[i*(n+1)*n+(j+66)] = a[i*(n+1)*k]*b[(k+1)*n+(j+66)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+66)]
c[i*(n+1)*n+(j+67)] = a[i*(n+1)*k]*b[(k+1)*n+(j+67)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+67)]
c[i*(n+1)*n+(j+68)] = a[i*(n+1)*k]*b[(k+1)*n+(j+68)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+68)]
c[i*(n+1)*n+(j+69)] = a[i*(n+1)*k]*b[(k+1)*n+(j+69)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+69)]
c[i*(n+1)*n+(j+70)] = a[i*(n+1)*k]*b[(k+1)*n+(j+70)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+70)]
c[i*(n+1)*n+(j+71)] = a[i*(n+1)*k]*b[(k+1)*n+(j+71)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+71)]
c[i*(n+1)*n+(j+72)] = a[i*(n+1)*k]*b[(k+1)*n+(j+72)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+72)]
c[i*(n+1)*n+(j+73)] = a[i*(n+1)*k]*b[(k+1)*n+(j+73)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+73)]
c[i*(n+1)*n+(j+74)] = a[i*(n+1)*k]*b[(k+1)*n+(j+74)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+74)]
c[i*(n+1)*n+(j+75)] = a[i*(n+1)*k]*b[(k+1)*n+(j+75)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+75)]
c[i*(n+1)*n+(j+76)] = a[i*(n+1)*k]*b[(k+1)*n+(j+76)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+76)]
c[i*(n+1)*n+(j+77)] = a[i*(n+1)*k]*b[(k+1)*n+(j+77)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+77)]
c[i*(n+1)*n+(j+78)] = a[i*(n+1)*k]*b[(k+1)*n+(j+78)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+78)]
c[i*(n+1)*n+(j+79)] = a[i*(n+1)*k]*b[(k+1)*n+(j+79)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+79)]
c[i*(n+1)*n+(j+80)] = a[i*(n+1)*k]*b[(k+1)*n+(j+80)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+80)]
c[i*(n+1)*n+(j+81)] = a[i*(n+1)*k]*b[(k+1)*n+(j+81)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+81)]
c[i*(n+1)*n+(j+82)] = a[i*(n+1)*k]*b[(k+1)*n+(j+82)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+82)]
c[i*(n+1)*n+(j+83)] = a[i*(n+1)*k]*b[(k+1)*n+(j+83)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+83)]
c[i*(n+1)*n+(j+84)] = a[i*(n+1)*k]*b[(k+1)*n+(j+84)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+84)]
c[i*(n+1)*n+(j+85)] = a[i*(n+1)*k]*b[(k+1)*n+(j+85)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+85)]
c[i*(n+1)*n+(j+86)] = a[i*(n+1)*k]*b[(k+1)*n+(j+86)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+86)]
c[i*(n+1)*n+(j+87)] = a[i*(n+1)*k]*b[(k+1)*n+(j+87)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+87)]
c[i*(n+1)*n+(j+88)] = a[i*(n+1)*k]*b[(k+1)*n+(j+88)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+88)]
c[i*(n+1)*n+(j+89)] = a[i*(n+1)*k]*b[(k+1)*n+(j+89)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+89)]
c[i*(n+1)*n+(j+90)] = a[i*(n+1)*k]*b[(k+1)*n+(j+90)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+90)]
c[i*(n+1)*n+(j+91)] = a[i*(n+1)*k]*b[(k+1)*n+(j+91)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+91)]
c[i*(n+1)*n+(j+92)] = a[i*(n+1)*k]*b[(k+1)*n+(j+92)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+92)]
c[i*(n+1)*n+(j+93)] = a[i*(n+1)*k]*b[(k+1)*n+(j+93)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+93)]
c[i*(n+1)*n+(j+94)] = a[i*(n+1)*k]*b[(k+1)*n+(j+94)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+94)]
c[i*(n+1)*n+(j+95)] = a[i*(n+1)*k]*b[(k+1)*n+(j+95)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+95)]
c[i*(n+1)*n+(j+96)] = a[i*(n+1)*k]*b[(k+1)*n+(j+96)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+96)]
c[i*(n+1)*n+(j+97)] = a[i*(n+1)*k]*b[(k+1)*n+(j+97)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+97)]
c[i*(n+1)*n+(j+98)] = a[i*(n+1)*k]*b[(k+1)*n+(j+98)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+98)]
c[i*(n+1)*n+(j+99)] = a[i*(n+1)*k]*b[(k+1)*n+(j+99)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+99)]
c[i*(n+1)*n+(j+100)] = a[i*(n+1)*k]*b[(k+1)*n+(j+100)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+100)]
c[i*(n+1)*n+(j+101)] = a[i*(n+1)*k]*b[(k+1)*n+(j+101)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+101)]
c[i*(n+1)*n+(j+102)] = a[i*(n+1)*k]*b[(k+1)*n+(j+102)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+102)]
c[i*(n+1)*n+(j+103)] = a[i*(n+1)*k]*b[(k+1)*n+(j+103)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+103)]
c[i*(n+1)*n+(j+104)] = a[i*(n+1)*k]*b[(k+1)*n+(j+104)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+104)]
c[i*(n+1)*n+(j+105)] = a[i*(n+1)*k]*b[(k+1)*n+(j+105)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+105)]
c[i*(n+1)*n+(j+106)] = a[i*(n+1)*k]*b[(k+1)*n+(j+106)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+106)]
c[i*(n+1)*n+(j+107)] = a[i*(n+1)*k]*b[(k+1)*n+(j+107)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+107)]
c[i*(n+1)*n+(j+108)] = a[i*(n+1)*k]*b[(k+1)*n+(j+108)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+108)]
c[i*(n+1)*n+(j+109)] = a[i*(n+1)*k]*b[(k+1)*n+(j+109)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+109)]
c[i*(n+1)*n+(j+110)] = a[i*(n+1)*k]*b[(k+1)*n+(j+110)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+110)]
c[i*(n+1)*n+(j+111)] = a[i*(n+1)*k]*b[(k+1)*n+(j+111)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+111)]
c[i*(n+1)*n+(j+112)] = a[i*(n+1)*k]*b[(k+1)*n+(j+112)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+112)]
c[i*(n+1)*n+(j+113)] = a[i*(n+1)*k]*b[(k+1)*n+(j+113)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+113)]
c[i*(n+1)*n+(j+114)] = a[i*(n+1)*k]*b[(k+1)*n+(j+114)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+114)]
c[i*(n+1)*n+(j+115)] = a[i*(n+1)*k]*b[(k+1)*n+(j+115)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+115)]
c[i*(n+1)*n+(j+116)] = a[i*(n+1)*k]*b[(k+1)*n+(j+116)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+116)]
c[i*(n+1)*n+(j+117)] = a[i*(n+1)*k]*b[(k+1)*n+(j+117)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+117)]
c[i*(n+1)*n+(j+118)] = a[i*(n+1)*k]*b[(k+1)*n+(j+118)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+118)]
c[i*(n+1)*n+(j+119)] = a[i*(n+1)*k]*b[(k+1)*n+(j+119)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+119)]
c[i*(n+1)*n+(j+120)] = a[i*(n+1)*k]*b[(k+1)*n+(j+120)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+120)]
c[i*(n+1)*n+(j+121)] = a[i*(n+1)*k]*b[(k+1)*n+(j+121)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+121)]
c[i*(n+1)*n+(j+122)] = a[i*(n+1)*k]*b[(k+1)*n+(j+122)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+122)]
c[i*(n+1)*n+(j+123)] = a[i*(n+1)*k]*b[(k+1)*n+(j+123)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+123)]
c[i*(n+1)*n+(j+124)] = a[i*(n+1)*k]*b[(k+1)*n+(j+124)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+124)]
c[i*(n+1)*n+(j+125)] = a[i*(n+1)*k]*b[(k+1)*n+(j+125)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+125)]
c[i*(n+1)*n+(j+126)] = a[i*(n+1)*k]*b[(k+1)*n+(j+126)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+126)]
c[i*(n+1)*n+(j+127)] = a[i*(n+1)*k]*b[(k+1)*n+(j+127)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+127)]
c[i*(n+1)*n+(j+128)] = a[i*(n+1)*k]*b[(k+1)*n+(j+128)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+128)]
c[i*(n+1)*n+(j+129)] = a[i*(n+1)*k]*b[(k+1)*n+(j+129)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+129)]
c[i*(n+1)*n+(j+130)] = a[i*(n+1)*k]*b[(k+1)*n+(j+130)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+130)]
c[i*(n+1)*n+(j+131)] = a[i*(n+1)*k]*b[(k+1)*n+(j+131)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+131)]
c[i*(n+1)*n+(j+132)] = a[i*(n+1)*k]*b[(k+1)*n+(j+132)] + a[i*(n+1)*k+1]*b[(k+1)*n+(j+132)]
c[i*(n+1)*n+(j+1
```

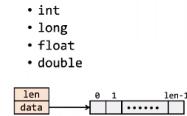
## Chapter 15

# Architecture and Optimization

In order to measure performance we use benchmarks. The benchmark should be as general as possible testing different datatypes e.g.:

### Benchmark example: data type for vectors

- Data types: different declarations for `data_t`:



```
/* data structure for vectors */
struct vec {
    size_t len;
    data_t *data;
};

/* retrieves vector_element
   parameter v: val */
int get_vec_element
    (struct vec* v, size_t idx, data_t *val)
{
    if (idx > v->len)
        return 0;
    *val = v->data[idx];
    return 1;
}
```

### Benchmark computation

- Data types (`data_t`):
- Operations: definitions of OP and IDENT
- Operations: definitions of OP and IDENT

```
void combine1(struct vec *v, data_t *dest)
{
    long int i;
    *dest = IDENT;
    for (i = 0; i < vec_length(v); i++) {
        data_t val;
        get_vec_element(v, i, &val);
        *dest = *dest OP val;
    }
}
```

Compute sum or product of vector elements

A way to express performance of a program that operates on vectors or lists is Cycles per Element (CPE). The **Execution time** is given by:

$$CPE \cdot n + \text{overhead}$$

We can increase performance by using the optimization flags, or other basic code optimizations:

### Benchmark Performance

```
void combine1(struct vec *v, data_t *dest)
{
    long int i;
    *dest = IDENT;
    for (i = 0; i < vec_length(v); i++) {
        data_t val;
        get_vec_element(v, i, &val);
        *dest = *dest OP val;
    }
}
```

Compute sum or product of vector elements

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Operation	22.68	20.02	19.98	20.18
Combine1 unoptimized	10.12	10.12	10.17	11.14

### Basic Optimizations

- Move `vec_length` out of loop
- Avoid bounds check on each cycle, instead access array directly (no procedure call in loop)
- Accumulate in temporary

```
void combined(struct vec *v, data_t *dest)
{
    long i;
    long length = vec_length(v);
    data_t *d = get_vec_start(v);
    data_t t = IDENT;
    for (i = 0; i < length; i++)
        t = t OP d[i];
    *dest = t;
}
```

### Effect of Basic Optimizations

- Eliminates sources of overhead in loop

```
void combine4(struct vec *v, data_t *dest)
{
    long i;
    long length = vec_length(v);
    data_t *d = get_vec_start(v);
    data_t t = IDENT;
    for (i = 0; i < length; i++)
        t = t OP d[i];
    *dest = t;
}
```

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Operation	22.68	20.02	19.98	20.18
Combine1 -O1	10.12	10.12	10.17	11.14

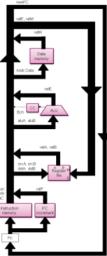
But can we do better?

## 15.1 modern processor design

### 15.1.1 Sequential processor stages

Sequential processor stages

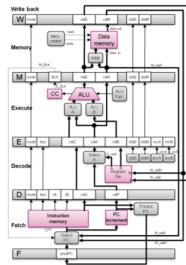
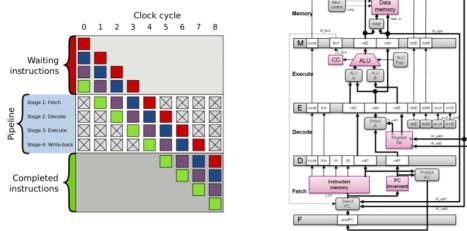
- Fetch
  - Read instruction from instruction memory
- Decode
  - Read program registers
- Execute
  - Compute value or address
- Memory
  - Read or write data
- Write Back
  - Write program registers
- PC
  - Update program counter



In practice this is too slow. Not all instructions need each step. And there will be steps which are idle. The solution is using pipelined hardware.

### 15.1.2 Pipelined hardware

Pipelined hardware



The following limitations for pipelining exist:

- **Data Hazards:** Hazards caused by data dependencies i.e we need a result from another calculation. These can be solved by implementing data forwarding or stalling the pipeline
- **Control Hazards:** Hazards caused by branches i.e we dont know for sure wether or not to take a branch. This can be solved by rollback or stalling the fetch stage until the branch is known.

### 15.1.3 Performance

Performance



$$\bullet \text{ Program Execution Time} = IC \cdot CPI \cdot CCT$$

- IC = instruction count
- CPI = cycles per instruction ( $= 1/IPC$ )
- CCT = clock cycle time ( $= 1/Frequency$ )

#### Limits to increasing pipeline depth:

- Delay of pipeline registers
- Inequalities in work per stage
  - Cannot break up work into stages at arbitrary points
- Clock skew
  - Clocks to different registers may not be perfectly aligned

$$\bullet \text{ Program Execution Time} = IC \cdot CPI \cdot CCT$$

- IC = instruction count
- CPI = cycles per instruction ( $= 1/IPC$ )
- CCT = clock cycle time ( $= 1/Frequency$ )

$$\bullet \text{ Cycles per instruction: } CPI = CPI_{base} + CPI_{stalls}$$

- Stalls due to data hazards
  - Read-after-write, Write-after-read, Write-after-write dependencies
- Stalls due to control hazards
  - Resolving jumps and branches
- Stalls due to memory latency
  - Large memories are slow

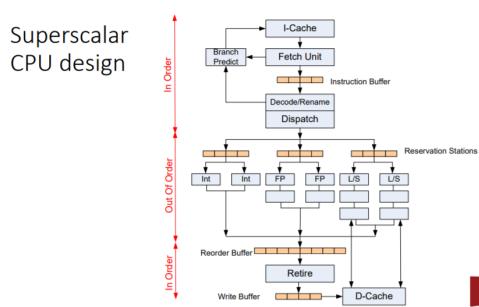
*How to improve CPI<sub>base</sub>?*

*Data forwarding, branch prediction, caching*

We can optimize IC by optimizing our code. We can decrease CCT, by pipelining. When there are more stages, the signal doesn't need to traverse as far and hence we can increase the clock frequency. What we can do to improve the CPI is widen the Processor i.e use a superscalar processor

#### 15.1.4 Superscalar processor

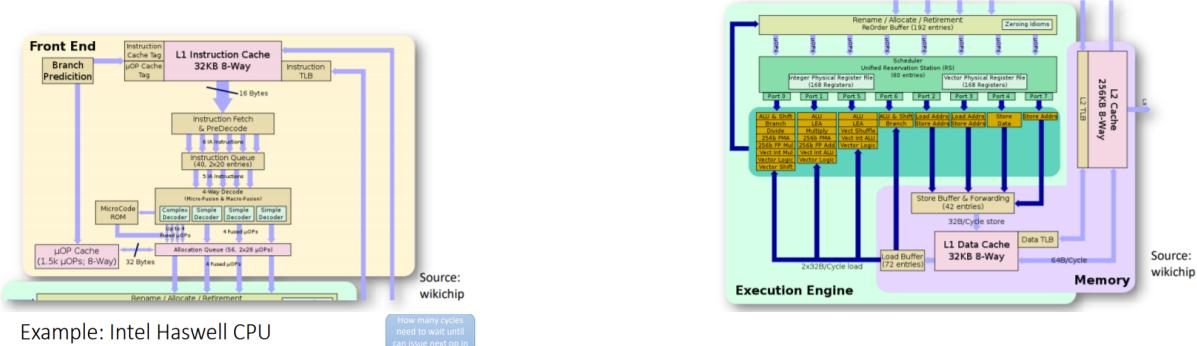
- **Definition:** A superscalar processor can issue and execute *multiple instructions in one cycle*. The instructions are retrieved from a sequential instruction stream and are usually scheduled dynamically.
  - Benefit: without programming effort, superscalar processor can take advantage of the *instruction level parallelism* that most programs have
  - Most CPUs since about 1998 are superscalar.
  - Intel: since Pentium Pro



We need to fetch instructions in the order they are being executed in. But when we are actually executing them then we can execute any instruction which has no dependencies. We do need to preserve the semantics of the program hence we need to commit the instructions in the same order as they have been fetched, this is done with the reorder buffer.

## 15.2 Superscalar processor performance

### 15.2.1 Intel Haswell CPU schematics



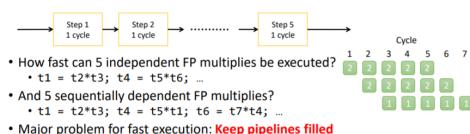
- 8 total functional units
  - Multiple instructions can execute in parallel
    - 2 load, with address computation
    - 1 store, with address computation
    - 4 integer
    - 2 FP multiply
    - 1 FP add
    - 1 FP divide
  - Some instructions take > 1 cycle, but can be pipelined

Instruction	Latency	Cycles/issue
Load/store	4	1
Int Add	1	1
Int Multiply	3	1
Int/Long Divide	3-30	3-30
Sngl/Dbl FP Multiply	5	1
Sngl/Dbl FP Add	3	1

Cycles/issue means how many cycles we need to wait before we can issue a new instruction on that EU.

### 15.2.2 Latency vs Throughput

- Last slide:  
  FP Multiply:      *latency*      *cycles/issue*



If we have independent multiplications we can finish the 5 multiplications in 5 cycles. If there are sequential dependencies for each of the multiplications then we can only do one multiply every 5 cycles and the whole calculation will then take 25 cycles.

### 15.2.3 Data Hazards

- Data dependencies for instruction  $j$  following instruction  $i$ 
    - **Read after Write (RAW) (true dependence)**
      - Instruction  $j$  tries to read before instruction  $i$  tries to write it
    - **Write after Write (WAW) (output dependence)**
      - Instruction  $j$  tries to write an operand before  $i$  writes its value
    - **Write after Read (WAR) (anti dependence)**
      - Instruction  $j$  tries to write a destination before it is read by  $i$
    - No such thing as a Read after Read (RAR) hazard since there is never a problem reading twice

- No such thing as a Read after Read (RAR) hazard since there is never a problem reading twice

Can avoid hazard with  
register renaming

RAW is a true dependancy because we really need to read the value that instruction i writes. WAW, WAR can be avoided because j is only writing a value and does not need to wait for anything. Register renaming allows this by using extra registers to store the write value of j and hence not overwrite the value of i.

### 15.2.4 Register Renaming

- Use more registers than in the ISA
  - Map *architectural* registers to a larger pool of *physical* registers
  - Give each new value produced its own physical register
  - Avoid WAW and WAR hazards
- Example: Before & after renaming
 

$R1 = R2 + R3$ $R4 = R1 + R5$ $R1 = R6 + R7$ $R6 = R1 + R3$	$R1 = R2 + R3$ $R4 = R1 + R5$ $R33 = R6 + R7$ $R34 = R33 + R3$
----------------------------------------------------------------------	-------------------------------------------------------------------------

(RAW hazard)

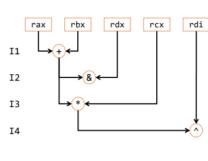
The ISA defines how many registers are available the "architectural registers". The processor itself can have a larger pool of physical registers. We dont add more registers to the ISA because we want to keep backwards compatibility. There are sometimes add ons to the ISA but in general we dont want to change the contract between software and hardware.

### 15.2.5 Instruction Execution

Traditional view of instruction execution

- Imperative view
  - Registers are fixed storage locations
    - Individual instructions read & write them
  - Instructions must be executed in specified sequence to guarantee proper program behavior

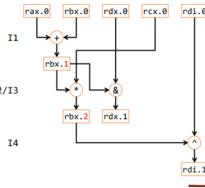
```
addq %rax, %rbx # I1
andq %rbx, %rdx # I2
mulq %rcx, %rbx # I3
xorq %rbx, %rdi # I4
```



Dataflow view of instruction execution

- Functional view
  - View each write as creating new instance of value
  - Operations can be performed as soon as operands available
  - No need to execute in original sequence

```
addq %rax, %rbx # I1
andq %rbx, %rdx # I2
mulq %rcx, %rbx # I3
xorq %rbx, %rdi # I4
```



### 15.2.6 Meaning for Program performance

What is the performance bound?

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Operation				
Combine4	1.27	3.01	3.01	5.01
Unroll 2x1	1.01	3.01	3.01	5.01
Latency bound	1.00	3.00	3.00	5.00
Throughput bound	0.50	1.00	1.00	0.50

Instructions that can execute in parallel:

- 2 load, with address computation
- 1 store, with address computation
- 4 add
- 1 integer multiply
- 2 FP multiply
- 1 FP add
- 1 FP divide

But how can we make our program throughput-bound?  
What can we do about the sequential dependency?



We get the latency bound by looking at the latency entries of the respective operation from the table on the right. In the original Combine4 program the add was still slower, this was because we could exploit it some more using loopunrolling. The value of the throughput bound is dependant not only on the number of functional units but also on the number of loads that can be done. We see that the throughput bound offers us better performance and hence we will want our program to be throughput bound and not sequential, this can be done with reassociation

### 15.3 Reassociation

Loop Unrolling with Reassociation (2x1a)

- Can this change the result of the computation?
- Yes, for FP. Why?

Compare to before:  
 $x = (x \text{ OP } d[i]) \text{ OP } d[i+1];$

```
void unrollzaa_combine(struct vec *v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length;
    data_t *d = get_vec_start(v);
    data_t x = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x = x OP (d[i] OP d[i+1]);
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x = x OP d[i];
    }
    *dest = x;
}
```

- Nearly 2x speedup for Int \*, FP +, FP \*
- Reason: Breaks sequential dependency

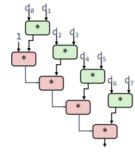
$x = x \text{ OP } (d[i] \text{ OP } d[i+1]);$

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Operation				
Combine4	1.27	3.01	3.01	5.01
Unroll 2x1	1.01	3.01	3.01	5.01
Unroll 2x1a	1.01	1.51	1.51	2.51
Latency Bound	1.00	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50

The reason that the change of brackets increase performance is because we remove the dependancy off the previous result, hence we allow for more efficient scheduling. For floating point values this would not be valid because FP's are not associative. If you however know that the the values wont be compromised then you can do it. The resulting computation structure is:

- What changed:
  - Ops in the next iteration can be started early (no dependency)
- Overall Performance
  - N elements, D cycles latency/op
  - $(N/2+1)D$  cycles:
 
$$CPE = D/2$$

$$x = x \text{ OP } (d[i] \text{ OP } d[i+1]);$$



### 15.3.1 Separate Accumulators

- Different form of reassociation

```
void unroll2a_combine(struct vec *v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x0 = IDENT;
    data_t x1 = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x0 = x0 OP d[i];
        x1 = x1 OP d[i+1];
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x0 = x0 OP d[i];
    }
    *dest = x0 OP x1;
}
```

- Int + makes use of two load units
- 2x speedup (over unroll2) for Int \*, FP +, FP \*

$$x0 = x0 \text{ OP } d[i];$$

$$x1 = x1 \text{ OP } d[i+1];$$

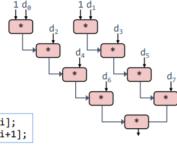
Method	Integer		Double FP	
	Add	Mult	Add	Mult
Combine4	1.27	3.01	3.01	5.01
Unroll 2x1	1.01	3.01	3.01	5.01
Unroll 2x1a	1.01	1.51	1.51	2.51
Unroll 2x2	0.81	1.51	1.51	2.51
Latency Bound	1.00	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50

#### Separate Accumulators

- What changed:
  - Two independent "streams" of operations
- Overall Performance
  - N elements, D cycles latency/op
  - Should be  $(N/2+1)D$  cycles:
 
$$CPE = D/2$$

$$x0 = x0 \text{ OP } d[i];$$

$$x1 = x1 \text{ OP } d[i+1];$$
  - CPE matches prediction!



This increases the performance for the add because we are overlapping operations. In the previous case we were doing two loads and then the respective additions, whereas now we can do the loads and additions simultaneously. It does not affect the other arithmetic operations because they have a higher latency.

## 15.4 Combining multiple accumulators and unrolling

- Idea
  - Can unroll to any degree L
  - Can accumulate K results in parallel
  - L must be multiple of K
- Limitations
  - Diminishing returns
    - Cannot go beyond throughput limitations of execution units
  - Large overhead for short lengths
    - Finish off iterations sequentially

#### Unrolling & Accumulating: Double \*

Accumulators	FP *	Unrolling Factor L									
		K	1	2	3	4	6	8	10	12	
1	5.01	5.01	5.01	5.01	5.01	5.01	5.01	5.01	5.01	5.01	
2		2.51									
3			1.67								
4				1.25							
6					0.84						
8						0.63					
10							0.51				
12								0.52			

#### Unrolling & Accumulating: Int +

Accumulators	FP *	Unrolling Factor L											
		K	1	2	3	4	6	8	10	12			
1	1.27	1.01	1.01	1.01	1.01	1.01	1.01	1.01	1.01	1.01			
2		0.81											
3			0.74										
4				0.69									
6					0.56								
8						0.54							
10							0.54						
12								0.56					

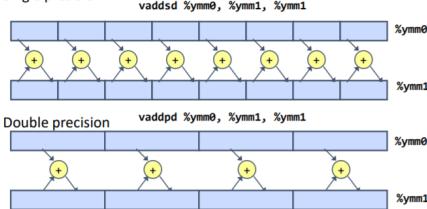
- Limited only by throughput of functional units
- Up to 42X improvement over original, unoptimized code

The compiler can do loop unrolling, but will not know which is best, hence we must make it specific in the code.

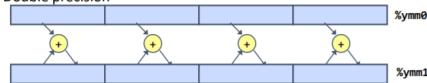
### 15.4.1 AVX2 SIMD operations

AVX2 SIMD operations (256-bit vectors)

- Single precision



- Double precision



- Make use of AVX Instructions (256 bits wide):

Method	Integer		Double FP	
	Add	Mult	Add	Mult
Operation				
Scalar Best	0.54	1.01	1.01	0.52
Vector Best	0.06	0.24	0.25	0.16
Latency Bound	0.50	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50
Vector Throughput Bound	0.06	0.12	0.25	0.12

# Chapter 16

## Caches

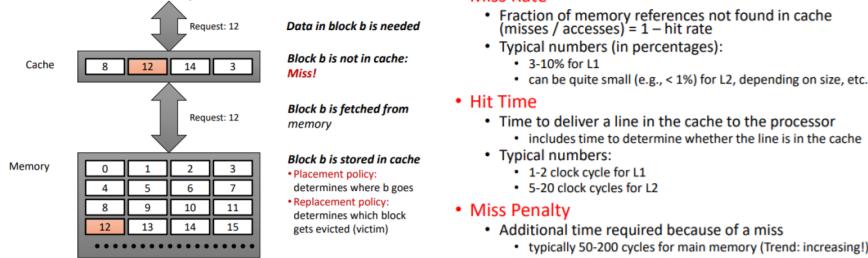
### 16.1 Intro

#### 16.1.1 Processor-memory bottleneck

The CPU performance continuously improves, but the bandwidth to access memory is evolving much slower. E.g. Haswell core can process 512 Bytes/cycle but the bandwidth is 10 Bytes/cycle, hence accessing main memory results in stalls. The solution to this problem is to use a hierarchy of caches.

#### 16.1.2 General Cache concepts

##### General cache concepts



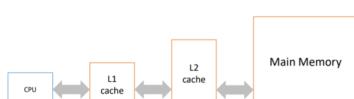
Because the huge difference between hit and miss times each percent in improved hit rate leads to enormous performance gains.

- 99% hit rate is twice as good as 97%!
  - Consider:
    - cache hit time of 1 cycle
    - miss penalty of 100 cycles
  - Average access time:
    - 97% hits: 1 cycle + 0.03 \* 100 cycles = 4 cycles
    - 99% hits: 1 cycle + 0.01 \* 100 cycles = 2 cycles
- This is why we use **miss rate** instead of hit rate

Another metric is "MPKI" which is misses per kilo instruction. MPKI is how many misses per thousand instructions.

#### 16.1.3 Two-level cache performance

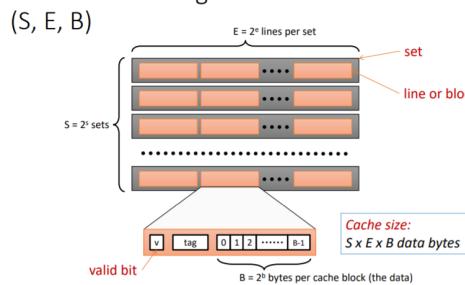
- **Average Memory Access Time** = HitTimeL1 + MissRateL1 \* MissPenaltyL1
- **MissPenaltyL1** = HitTimeL2 + MissRateL2 \* MissPenaltyL2
- **MissPenaltyL2** = DRAMaccessTime + (BlockSize/Bandwidth)



### 16.1.4 Types of cache misses

- **Cold (compulsory) miss**
  - Occurs on first access to a block
- **Conflict miss**
  - Most caches limit placement to small subset of available slots
    - e.g., block  $i$  must be placed in slot  $(i \bmod 4)$
  - Cache may be large enough, but multiple lines map to same slot
    - e.g., referencing blocks 0, 8, 0, 8, ... would miss every time
- **Capacity miss**
  - Set of active cache blocks (working set) larger than cache
- **Coherency miss**
  - Multiprocessor systems: see later in the course

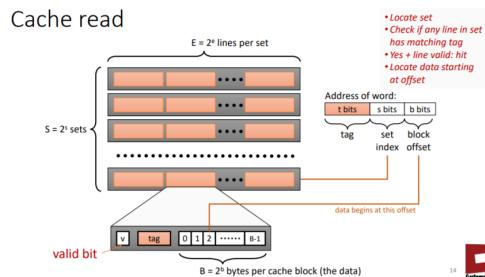
## 16.2 Cache organization



We need  $S$  bits to determine which set of the cache we want to access. Each set can have multiple entries (also called "ways"). Each block will have:

- **Valid bit:** Indicates whether the contents are a valid cache entry
- **Tag:** Used to determine a hit

## 16.3 Cache reads



the block offset indicates which byte of the block we want to access  
the set index indicates which set we want to access

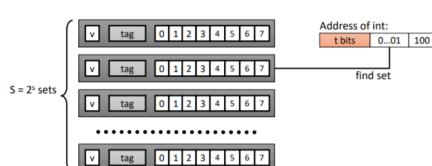
Since many locations could have the same set and block access we must compare the tags to ensure we are accessing the correct element.

### 16.3.1 Direct mapped cache

We only have one way and block per set

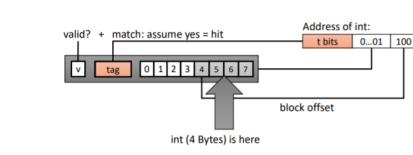
Direct mapped cache ( $E = 1$ )

Direct mapped: One line per set  
This example: cache block size 8 bytes



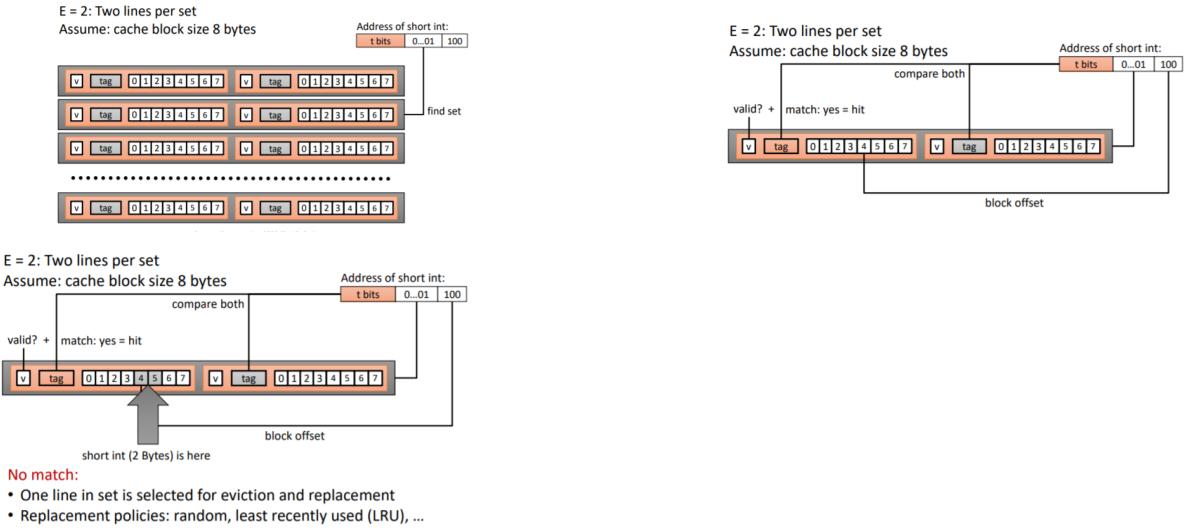
Direct mapped cache ( $E = 1$ )

Direct mapped: One line per set  
This example : cache block size 8 bytes



### 16.3.2 2-way set-associative cache

#### 2-way set-associative cache

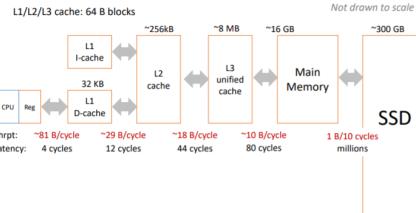


### 16.4 The memory hierarchy

#### Examples of caching in the hierarchy

#### Intel® Core™ i7-6700K (“Skylake”, 2015-)

Cache type	What is cached?	Where is it cached?	Latency (cycles)	Managed by
Registers	4/8-byte words	CPU core	0	Compiler
TLB	Address translations	On-chip TLB	0	Hardware
L1 cache	64-byte blocks	On-chip L1	1	Hardware
L2 cache	64-byte blocks	On-chip L2	10	Hardware
Virtual memory	4kB page	Main memory (RAM)	100	Hardware + OS
Buffer cache	4kB sectors	Main memory	100	OS
Network buffer cache	Parts of files	Local disk, SSD	1,000,000	SMB/NFS client
Browser cache	Web pages	Local disk	10,000,000	Web browser
Web cache	Web pages	Remote server disks	1,000,000,000	Web proxy server



Types of L1 caches:

- **I-cache:** Instruction Cache. There will be locality in the instructions because when we don't branch we use instructions that are right after each other.
- **D-cache:** Data Cache. There typically be locality in the data that is accessed (e.g. arrays)

The Disk is persistent storage, meaning if power is lost data is preserved. The Caches are all volatile storage, hence losing power means losing data. SSD are more efficient than Disks.

### 16.5 Cache writes

#### 16.5.1 Write-hit

Means we have this block in the cache and we are trying to write. There are two options:

- **Write-through**
  - Write immediately to memory
  - Memory is always consistent with the cache copy
  - Slow: what if the same value (or line!) is written several times
- **Write-back**
  - Defer write to memory until replacement of line
  - Need a **dirty** bit
    - ⇒ indicates line is different from memory
  - Higher performance (but more complex)

### 16.5.2 Write-miss

- Write-allocate (load into cache, update line in cache)
  - Good if more writes to the location follow
  - More complex to implement
  - May evict an existing value
  - Common with write-back caches
- No-write-allocate (writes immediately to memory)
  - Simpler to implement
  - Slower code (bad if value subsequently re-read)
  - Seen with write-through caches

### 16.5.3 Other hardware cache features

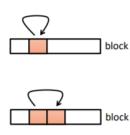
- Unified
  - Serves both instruction and data fetches
- Private
  - Only one core uses this cache
- Shared
  - Multiple cores share the cache
- Inclusive
  - Anything in this cache is **also** in every lower-level cache
- Exclusive
  - Anything in this cache is **not** in any lower-level cache

## 16.6 Cache optimizations

### 16.6.1 Locality

Why caches work: *locality*

- Programs tend to use data and instructions with addresses near or equal to those they have used recently



- Temporal locality:
  - Recently referenced items are likely to be referenced again in the near future

- Data:

- Temporal: sum referenced in each iteration
- Spatial: array  $a[i]$  accessed in stride-1 pattern

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

- Instructions:

- Temporal: cycle through loop repeatedly
- Spatial: reference instructions in sequence

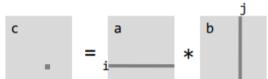
- Being able to assess the locality of code is a crucial skill for a programmer

### 16.6.2 Locality with matrix multiplication

Example: matrix multiplication

```
c = (double *) callloc(sizeof(double), n*n);

/* Multiply n x n matrices a and b */
void mmm(double *a, double *b, double *c, int n)
{
    int i, j, k;
    for (i = 0; i < n; i++)
        for (j = 0; j < n; j++)
            for (k = 0; k < n; k++)
                c[i*n+j] += a[i*n+k]*b[k*n + j];
}
```



#### Cache miss analysis

- Assume:

- Matrix elements are doubles
- Cache block = 8 doubles
- Cache size C << n (much smaller than n)



- First iteration:

- $n/8 + n = 9n/8$  misses

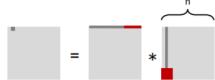
*Afterwards in cache:  
(schematic)*



#### Cache miss analysis

- Assume:

- Matrix elements are doubles
- Cache block = 8 doubles
- Cache size C << n (much smaller than n)



- Second iteration:

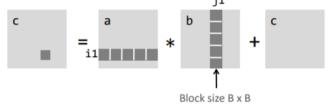
- Again:
  - $n/8 + n = 9n/8$  misses

- Total misses:

- $9n/8 * n^2 = (9/8) * n^3$

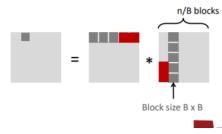
## 16.7 Blocking

## Blocked matrix multiplication



## Cache miss analysis

- Assume:
    - Cache block = 8 doubles
    - Cache size C << n (much smaller than n)
    - Three blocks  $\equiv$  fit into cache:  $3B^2 < C$
    - B is a multiple of cache block
  - Second (block) iteration:
    - Same as first iteration
    - $2n/B * B^2/B = nB/4$
  - Total misses:
    - $nB/4 * (n/B) = n^3/(4B)$



## Summary

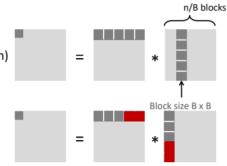
- No blocking:  $(9/8) * n^3$
  - Blocking:  $1/(4B) * n^3$
  - Suggests:  $\max(B)$  s.t.  $3B^2 < C$   
(can possibly be relaxed a bit, but there is a limit for)

- Reason for dramatic difference:
    - Matrix multiplication has inherent temporal locality:
      - Input data:  $3n^2$ , computation  $2n^3$
      - Every array elements used  $O(n)$  times!
    - But program has to be written properly
  - “Systems code” is usually much more subtle...

## Cache miss analysis

- Assume:

- Cache block = 8 doubles
  - Cache size  $C \ll n$  (much smaller than  $n$ )
  - Three blocks ■ fit into cache:  $3B^2 < C$
  - **B is a multiple of cache block**



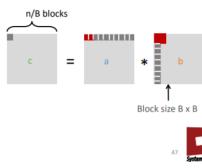
## Cache miss analysis

- Assume:

- Cache block = 8 doubles
  - Cache size  $C \ll n$  (much smaller than  $n$ )
  - Three blocks ■ fit into cache:  $3B^2 < C$

**$B \leq \text{cache block} \dots$**

Total misses:



# Chapter 17

## Exceptions

### 17.1 Intro

#### 17.1.1 Control Flow

- Processors do only one thing:
    - From startup to shutdown, a CPU simply reads and executes (interprets) a sequence of instructions, one at a time
    - This sequence is the CPU's **control flow** (or flow of control)
- Time ↓
- Physical control flow*
- ```
<startup>
inst1
inst2
inst3
...
instn
<shutdown>
```

#### 17.1.2 Altering the control flow

- Up to now: mostly two mechanisms for changing control flow:
  - Jumps and branches
  - Call and returnBoth react to changes in **program state**
- Insufficient for a useful system:  
Difficult to react to changes in **system state**
  - data arrives from a disk or a network adapter
  - instruction divides by zero
  - user hits Ctrl-C at the keyboard
  - System timer expires

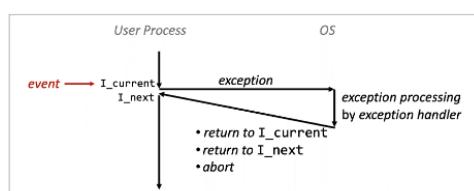
These alternative methods to alter control flow are known as Exceptional control flow

#### 17.1.3 Exceptional control flow

- Exists at all levels of a computer system
- Low level mechanisms
  - Hardware exceptions**
    - change in control flow in response to a system event (i.e., change in system state)
  - Combination of hardware and OS software
- Higher level mechanisms
  - Process context switch
  - Signals
  - Nonlocal jumps: `setjmp()`/`longjmp()`
  - Implemented by either:
    - OS software (context switch and signals)
    - C language runtime library (nonlocal jumps)
  - Language-level exceptions (Java, etc.)

#### 17.1.4 Exceptions

An exception is a transfer of control to the OS in response to some event (i.e. a change in processor state)



Examples: div by 0, machine check, page fault, I/O request completes, Ctrl-C

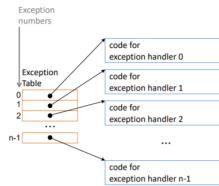
- A **synchronous** exception occurs as a result of executing an instruction.
- An **asynchronous** exception occurs as a result of events that are external to the processor.

Type of exception	Cause	Async/Sync
Interrupt	Signal from I/O device	Async
Trap	Intentional exception	Sync
Fault	Potentially recoverable error	Sync
Abort	Nonrecoverable error	Sync

## 17.2 Exception vectors and kernel mode

### 17.2.1 Exception vectors

- At boot time, OS allocates and initializes the **exception table**
  - The table's base address is stored in the **Exception Table Base Register**
- Each type of event has a unique exception number  $k$
- $k = \text{index into exception table}$  (a.k.a. interrupt vector)
- Handler  $k$  is called each time exception  $k$  occurs



Each exception table entry points to the beginning address of the code which must be executed for this particular exception.

### 17.2.2 x86 exception vectors

0	Divide error
1	Debug exception
2	Non-maskable interrupt (NMI)
3	Breakpoint
4	Overflow
5	Bounds check
6	Invalid opcode
7	Coprocessor not available
8	Double fault
9	Coprocessor segment overrun (386 or earlier)
A	Invalid task state segment
B	Segment not present
C	Stack fault
D	General protection fault
E	Page Fault
F	Reserved
10	Math fault
11	Alignment check
12	Machine check
13	SIMD floating point exception
14-1F	Reserved to Intel
20-FF	Available for external interrupts

### 17.2.3 Kernel

A Kernel is a part of the operating system. It is a part which runs in privileged mode i.e you have access to resources to registers and services which you don't have in user mode. Only time Kernel mode is used if a system call is done (e.g. when an exception is called there is a context switch into kernel mode).

## Kernel mode

- Exceptions cause a switch to **kernel mode**
  - Also: supervisor mode, privileged mode, ring 0, etc.
- Things look **very different**:
  - Access to system state (virtual memory, etc.)
  - Some exceptions are disabled
  - Some new instructions and registers
  - MMU behavior may change
- Details vary between processors
  - Even in the same architecture

## The Kernel

- Most operating systems have a **kernel**
  - = the part of the OS that runs in kernel mode
- Think of the kernel as:
  - A set of trap handling functions (always)
  - A set of threads in a special address space (sometimes)
  - Code to create the illusion of user-space processes (always)
- Many ways to structure the kernel and OS
  - Microkernels, monolithic kernels, multikernels, etc.
  - See the OS course next year...

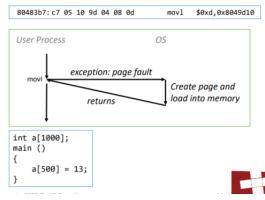
## 17.3 Synchronous Exceptions

- Caused by events that occur as a result of executing an instruction:
  - **Traps**
    - Intentional
    - Examples: **system calls**, breakpoint traps, special instructions
    - Returns control to "next" instruction
  - **Faults**
    - Unintentional but possibly recoverable
    - Examples: page faults (recoverable), protection faults (unrecoverable), floating point exceptions
    - Either re-executes faulting ("current") instruction or aborts
  - **Aborts**
    - Unintentional and unrecoverable
    - Examples: parity error, machine check
    - Aborts current program

### 17.3.1 Fault Examples

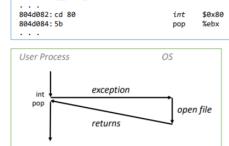
Fault example: page fault

- User writes to memory location
- That portion (page) of user's memory is currently on disk
- Page handler must load page into physical memory
- Returns to faulting instruction
- Successful on second try



Trap example: opening a file

- User calls: `open(filename, options)`
- Function open executes system call (e.g. `int` instruction)
- OS must find or create file, get it ready for reading or writing
- Returns integer file descriptor



A file descriptor is an integer which uniquely identifies a file

## 17.4 Asynchronous Exceptions (Interrupts)

- Caused by events **external** to the processor
  - Indicated by setting the processor's interrupt pin
  - Handler returns to "next" instruction
- Examples:
  - **I/O interrupts**
    - hitting Ctrl-C at the keyboard
    - arrival of a packet from a network
    - arrival of data from a disk
  - Hard reset interrupt
    - hitting the reset button
  - Soft reset interrupt
    - hitting Ctrl-Alt-Delete on a PC
- CPU **Interrupt-request line** triggered by I/O device
  - Might be edge- or level-triggered
- **Interrupt handler** receives interrupts
- **Maskable** to ignore or delay some interrupts
- Interrupt vector to dispatch interrupt to correct handler
  - Based on priority
  - Some **nonmaskable**
- Interrupt mechanism also used for exceptions

### 17.4.1 x86 interrupts

There are two interrupt pins:

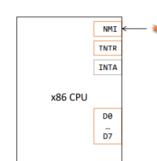
- INTR: interrupt request
- NMI: non-maskable interrupt

#### NMI

- NMI asserted ⇒
  - CPU completes current instruction
  - Issues Exception #2
  - **Always**.
- Cannot be disabled by the processor.
- Reserved for
  - Major hardware faults, e.g. memory parity error
  - "Watchdog" timer



- What caused the NMI?
- No obvious way to tell: the OS (INT 0x2 handler) must poll potential sources
- Typical machine: o(100s) devices
  - ⇒ slow reaction, inefficient
- Might get a second NMI while first is still being handled!



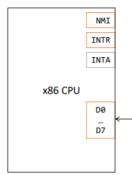
A "watchdog" timer is a timer which must be periodically reset, if this timer reaches zero, the processor enters a "panic state" (blue screen) and stops working. Exceptions are indexed from 0. Since there is only one pin there is no obvious

way to tell which device caused the interrupt, and hence it must check each device for the source of the interrupt (called polling). NMI interrupts are not so common, IRQ interrupts are more common

## IRQ

Interrupt requests are sent into the INTR pin. The CPU then acknowledges that the interrupt was received on the INTA pin (hence we decide which device we deal with) when the device confirms it sends data to the data bits. The Data bits allow the device to send us its interrupt vector, allowing us to except different kinds of interrupts.

- Interrupt request ⇒
  - Can be disabled by IE status flag (CLI/STI instructions to clear/set interrupt flag)
  - If enabled, complete current instruction, then:
- CPU **acknowledges** using INTA pin
- Interrupt **vector** is then supplied to the CPU via the data bus
- CPU issues exception # from the vector

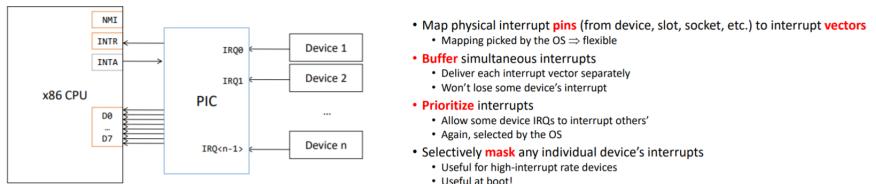


## 17.5 Interrupt controllers

### 17.5.1 Programmable Interrupt Controllers (PIC)

Instead of having Devices which directly connect to the INTR pin like we have for the NMI pin, we have the devices connected to the PIC. This allows us to manage which devices asked for the interrupts and then transmit the proper data to the CPU

#### Programmable Interrupt Controllers



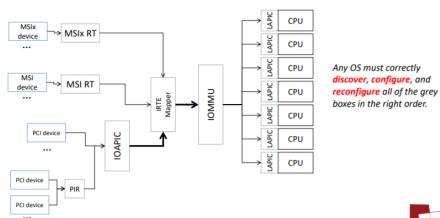
In the case of multiple interrupts we can use a buffer to store them.

### 17.5.2 Modern PICs

#### Modern PICs

- Every processor in a modern PC has a “Local APIC”
  - “Local Advanced Programmable Interrupt Controller”
  - PIC → APIC → LAPIC → xAPIC → ...
- Additional capabilities:
  - Inter-processor interrupts
  - Programmable timer
  - Sophisticated interrupt scheduling
  - Etc.
- Integrated onto same die as the core itself
- Also: IOxAPIC / MSI-X
  - Off-chip, attached to device I/O interconnect

#### Modern PCs (simplified!)

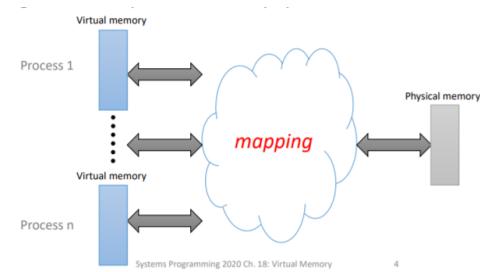


# Chapter 18

## Virtual Memory

### 18.1 Address Translation:

Each Process gets its own private memory space



#### 18.1.1 Address Spaces:

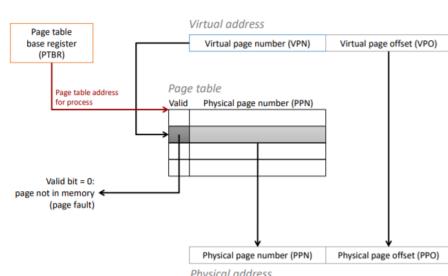
- **Linear address space:** Ordered set of contiguous non-negative integer addresses:  
 $\{0, 1, 2, 3 \dots\}$
- **Virtual address space:** Set of  $N = 2^n$  virtual addresses  
 $\{0, 1, 2, 3, \dots, N-1\}$
- **Physical address space:** Set of  $M = 2^m$  physical addresses  
 $\{0, 1, 2, 3, \dots, M-1\}$
- Clean distinction between data (bytes) and their attributes (addresses)
- Each object can now have multiple addresses
- Every byte in main memory: one physical address, one (or more) virtual addresses

One physical address can map to multiple addresses in Virtual memory.  
Virtual Memory requires the OS and hardware:

- **MMU (Memory Managing Unit):** This hardware is dedicated to translating virtual to physical addresses using a table in hardware which is managed and populated by hardware. Before memory can be accessed the virtual address needs to be translated to physical address. The memory hierarchy also applies to this.

#### 18.1.2 Address Translation with page table

**Page:** A page is a block of memory. This is the granularity of memory (corresponds to a cache block).



The page table is stored in a register in the CPU. The page table can be indexed with the virtual the virtual address. The page table tells of if a particular page is valid. The Virtual page offset specifies which byte in a page we want. The

Virtual page number needs to be translated into a physical page number. The Virtual page offset remains the same for the physical address.

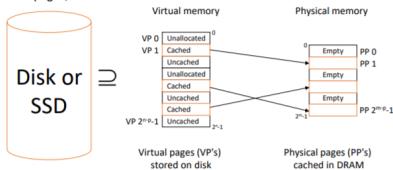
## 18.2 Uses of virtual memory

### 18.2.1 Benefits

- **Efficient use** of limited main memory (RAM)
  - Use RAM as a cache for the parts of a virtual address space
    - some non-cached parts stored on disk
    - some (unallocated) non-cached parts stored nowhere
  - Keep only active areas of virtual address space in memory
    - transfer data back and forth as needed
- **Simplifies** memory management for programmers
  - Each process gets the same full, private linear address space
- **Isolates** address spaces
  - One process can't interfere with another's memory
    - because they operate in different address spaces
  - User process cannot access privileged information
    - different sections of address spaces have different permissions

### 18.2.2 VM as a tool for caching

- Virtual memory: array of  $N = 2^n$  contiguous bytes
  - think of the array (allocated part) as being stored on disk
- Physical main memory (DRAM) = cache for allocated virtual memory
- Blocks are called pages; size =  $2^p$



Pages can be in 3 different states:

- **Unallocated**: A page which hasn't yet been used by the program. These are not present on disc.
- **Cached**: A page which has already been allocated. A Cached page resides in main memory
- **Uncached**: A page which has already been allocated. But does not reside in main memory. This page can be found on disc.

When trying to access an Uncached page, we will get a page fault. And the page fault handler will bring in that page from memory and retry the instruction, which will then be mapped to memory.

### 18.2.3 DRAM cache organization

**SRAM**: Memory type used for L1,L2,L3 caches.

- DRAM cache organization driven by the enormous miss penalty
  - DRAM is about **10x** slower than SRAM
  - Disk is about **10,000x** slower than DRAM
    - For first byte, faster for next byte
- **Consequences**
  - Large page (block) size: typically 4-8 KB, sometimes 4 MB
  - Fully associative
    - Any VP can be placed in any PP
    - Requires a "large" mapping function – different from CPU caches
  - Highly sophisticated, expensive replacement algorithms
    - Too complicated and open-ended to be implemented in hardware
  - Write-back rather than write-through
- Virtual memory works because of locality
  - At any point in time, programs tend to access a set of active virtual pages called the **working set**
    - Programs with better temporal locality will have smaller working sets
  - If (**working set size < main memory size**)
    - Good performance for one process after compulsory misses
  - If (**sum(working set sizes) > main memory size**)
    - **Thrashing**: Performance meltdown where pages are swapped (copied) in and out continuously

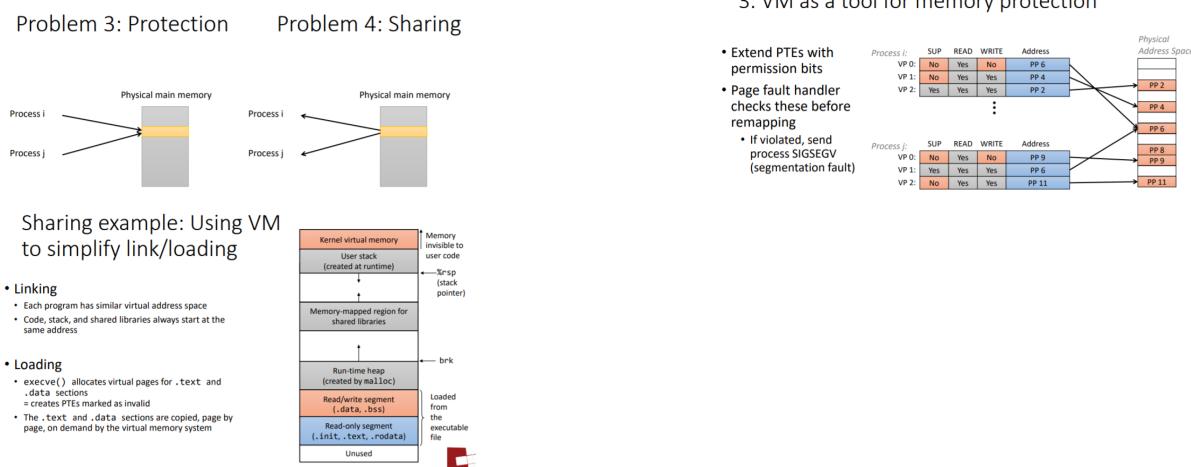
**Paging/Swapping** Bringing in pages from Disk. The reason that the page sizes are so large, is that accessing disc is so slow. So we try to compensate by bringing in more data when accessing it.

### 18.2.4 Memory Management

The idea of mapping multiple processes to physical memory.

- Key idea: each process has its own virtual address space
    - Views memory as a simple linear array
    - Mapping function scatters addresses through physical memory
    - Well-chosen mappings simplify memory allocation and management
- 
- Memory allocation
    - Each virtual page can be mapped to any physical page
    - A virtual page can be stored in different physical pages at different times
  - Sharing code and data among processes
    - Map virtual pages to the same physical page (here: PP 6)

## 18.2.5 Protection and Sharing

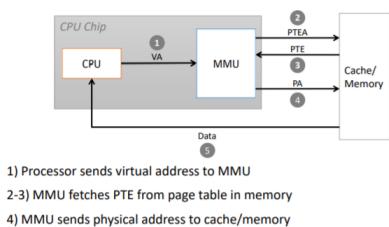


Supervisor mode permission are pages which can only be accessed by the kernel.

## 18.3 The address translation process

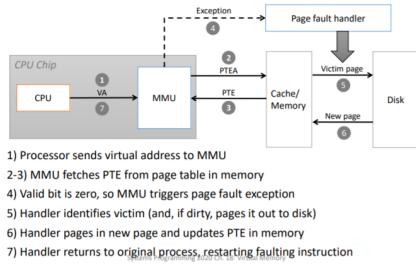
### 18.3.1 Page Hit:

Address translation: page hit



### 18.3.2 Page Fault:

Address translation: page fault



The Page fault handler is implemented in software. The cache can be indexed by virtual or physical addresses. Typically the cache is indexed by physical addresses, it is a bit slower because of translation, but is simpler to implement.

## 18.4 Translation lookaside buffer

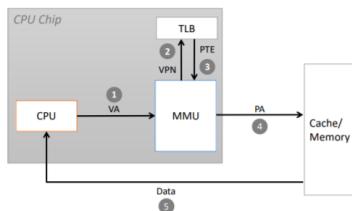
Caching for the page table translations.

Speeding up translation  
with a TLB

- Page table entries (PTEs) are cached in L1 like any other memory word
  - PTEs may be evicted by other data references
  - PTE hit in L1 still requires an extra 1 (or 2)-cycle delay
- Solution: *Translation Lookaside Buffer* (TLB)
  - Small hardware cache in MMU
  - Maps virtual page numbers to physical page numbers
  - Contains complete page table entries for small number of pages

### 18.4.1 TLB hit

TLB hit

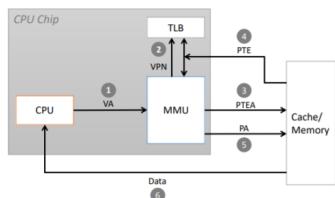


A TLB hit eliminates a memory access

When considering multiple processes we can either have identification bits for the processes. An alternative is to flush the TLB when context switching.

### 18.4.2 TLB Miss

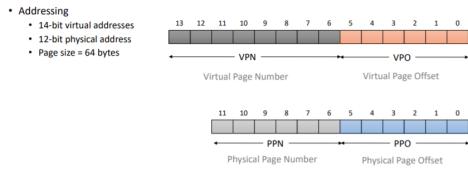
TLB miss



A TLB miss incurs an additional memory access (the PTE)  
 Fortunately, TLB misses are rare (we hope)

## 18.5 Simple Memory system Example

A simple memory system



### 18.5.1 Page Table

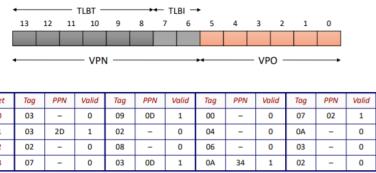
Simple memory system: page table

- Only show first 16 entries (out of 256)

VPN	PPN	Valid
00	28	1
01	-	0
02	33	1
03	02	1
04	-	0
05	16	1
06	-	0
07	-	0
08	13	1
09	17	1
0A	09	1
0B	-	0
0C	-	0
0D	2D	1
0E	11	1
0F	0D	1

### 18.5.2 TLB

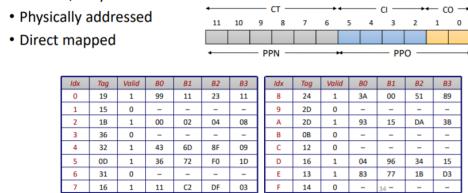
- 16 entries
- 4-way associative



- TLBT (TLB-tag):** Used to check whether the page we are actually accessing is stored in the TLB
- TLBI (TLB- index):** Used to index into the right set of the TLB

### 18.5.3 Cache

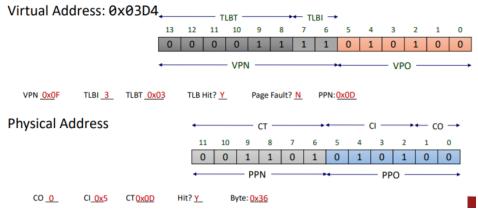
- 16 lines, 4-byte block size



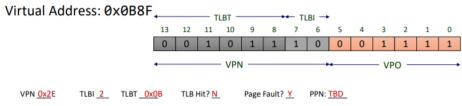
- PPN:** Physical Page Number
- PPO:** Physical Page Offset
- CO:** Cache Offset: Number of bits is determined by the cache block size
- CI:** Cache Index: Tells us which set to go to.
- CT:** Cache Tag: Tells us if we have a hit.

### 18.5.4 Address translation examples

#### Address translation example #1



#### Address translation example #2

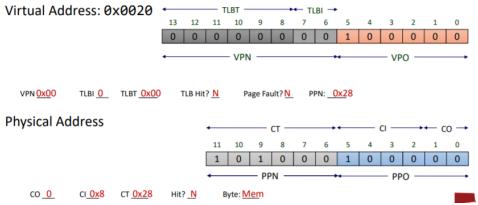


For the last example we don't have a hit in the TLB. So we check if the page table entry 00 which has PPN 28. Hence we don't have a page fault. A page fault is not the same as a miss. A fault refers to a miss in DRAM for the page itself.

## 18.6 Multi-level page tables

### 18.6.1 Terminology

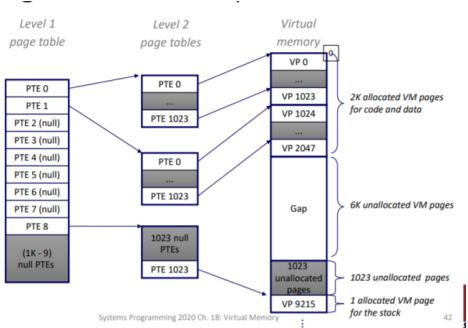
#### Address translation example #3



### 18.6.2 Linear Page Table size

- Given:
  - 4KB ( $2^{12}$ ) page size
  - 48-bit address space
  - 8-byte PTE
- How big a page table do you need?
- Answer: 512 GB!
  - $2^{48} / 2^{12} * 2^3 = 2^{39}$  bytes

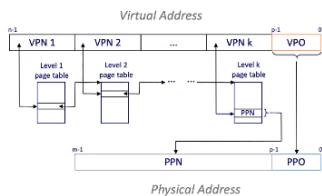
### 18.6.3 2-level page table hierarchy



We use the first bits of the VM address to index the Page directory table. This points to a different Level of smaller page tables. The next bits index the page table that was pointed to by the page directory table which in turn points to an

Address in VM. This reduces page Table size because a large range of addresses are represented by much fewer numbers and we only access allocate page table space on demand. The user tries to access a location in memory. If that hasn't been accessed before then the Page Directory Table will have a nuller pointer and hence a 2nd level page table is created. This idea can be extended to a k-level page table: We use the virtual page number to index into the multi level page tables. We take the virtual page number which is the top bits. The bottom bits are the virtual page offset. We divide the top bits by k (i.e the number of levels we have). The top 1/k bits will be used to index the first page table. The address we get from the first table will be the base of the second level table. And the next bits in the VPN will serve as a index from the second level page table. This continues until we reach the final level table which will give the PPN

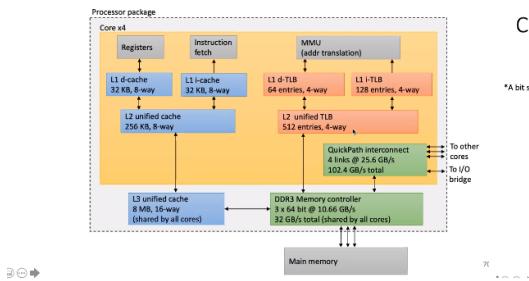
### Translating with a k-level page table



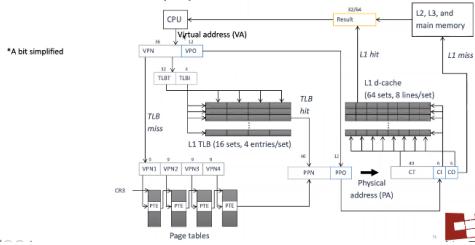
A miss on a multi level page table has a longer miss penalty called a "page lock".

## 18.7 Case study of the core i7 processor

### A typical core i7 processor



Core i7 memory system\*



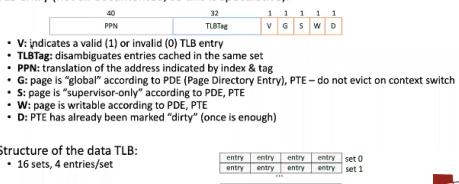
The size of the physical address space is proportional to the amount of RAM that is expected to be accessible on the machine. If one of the levels is a null pointer, then we create this level and continue this for the whole hierarchy for this address range we are trying to access.

### x86-64 paging

- 48-bit virtual address
  - 256 terabytes (TB)
  - Not yet ready for full 64 bits
    - Nobody can buy that much DRAM yet
    - Mapping tables would be huge
    - Multi-level array map may not be the right data structure
- 52-bit physical address = 40 bits for PPN
  - Requires 64-bit table entries
- Keep older x86 4KB page size (including for page tables)
  - $(4096 \text{ bytes per PT}) / (8 \text{ bytes per PTE}) = 512 \text{ entries per page}$

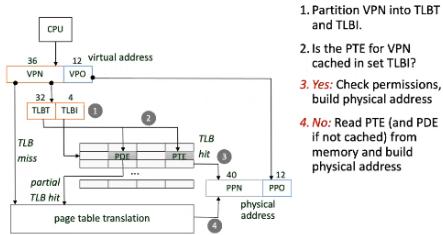
### 18.7.1 Core i7 TLB

- TLB entry (not all documented, so this is speculative):

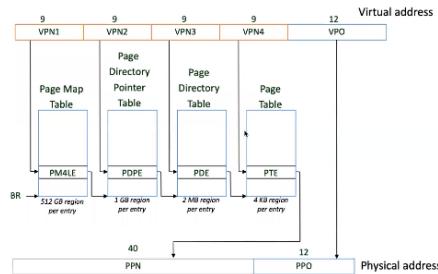


the dirty bit indicates that the value has been changed and has a different value than what's on disc. Hence when evicting this entry we also update the disc. The Page Directory refers to the higher levels of page tables. The Page Directory entry is the highest level in the page table.

## 18.7.2 Translating with the i7 TLB



## 18.7.3 x86-64 paging



Page Directory Pointer Entry (level 3)

Page-Map Level 4											
43-42	41	40	39	38	37	36	35	34	33	32	31
X	X	Available	Page-Directory Base Address (This is an architectural limit. A given implementation may support fewer bits.)								
			Page-Directory/Pointer Base Address (This is an architectural limit. A given implementation may support fewer bits.)								
			12	11	10	9	8	7	6	5	4
			A	L	M	B	C	D	E	F	G
			P	O	N	H	I	J	K	L	M
			W	R	S	T	U	V	W	X	Y
			Z	U	V	W	X	Y	Z	U	V
			P	R	S	T	U	V	W	X	Y

Page-Map Level 4

- **Page-Directory-Pointer Base Address:** 40 most significant bits of physical page table address (forces page tables to be 4KB aligned)
- **Avail:** These bits available for system programmers
- **A:** accessed (set by MMU on reads and writes, cleared by software)
- **PCD:** cache disabled (1) or enabled (0)
- **PWT:** write-through or write-back cache policy for this page table
- **U/S:** user or supervisor mode access
- **R/W:** read-only or read-write access
- **P:** page table is present in memory (1) or not (0)

Page Directory Entry (level 2)

Page-Map Level 4											
43-42	41	40	39	38	37	36	35	34	33	32	31
X	X	Available	Page-Table Base Address (This is an architectural limit. A given implementation may support fewer bits.)								
			Page-Table/Pointer Base Address (This is an architectural limit. A given implementation may support fewer bits.)								
			12	11	10	9	8	7	6	5	4
			A	L	M	B	C	D	E	F	G
			P	O	N	H	I	J	K	L	M
			W	R	S	T	U	V	W	X	Y
			Z	U	V	W	X	Y	Z	U	V
			P	R	S	T	U	V	W	X	Y

- **Page-Directory Base Address:** 40 most significant bits of physical page table address (forces page tables to be 4KB aligned)
- **Avail:** These bits available for system programmers
- **A:** accessed (set by MMU on reads and writes, cleared by software)
- **PCD:** cache disabled (1) or enabled (0)
- **PWT:** write-through or write-back cache policy for this page table
- **U/S:** user or supervisor mode access
- **R/W:** read-only or read-write access
- **P:** page table is present in memory (1) or not (0)

Page Table Entry (level 1)

Page Table Entry (level 1)											
43-42	41	40	39	38	37	36	35	34	33	32	31
X	X	Available	Physical-Page Base Address (This is an architectural limit. A given implementation may support fewer bits.)								
			Physical-Page/Pointer Base Address (This is an architectural limit. A given implementation may support fewer bits.)								
			12	11	10	9	8	7	6	5	4
			A	L	M	B	C	D	E	F	G
			P	O	N	H	I	J	K	L	M
			W	R	S	T	U	V	W	X	Y
			Z	U	V	W	X	Y	Z	U	V
			P	R	S	T	U	V	W	X	Y

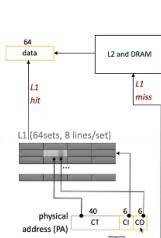
- **Page Table Base Address:** 40 most significant bits of physical page table address (forces pages to be 4 KB aligned)
- **Avail:** These bits available for system programmers
- **A:** accessed (set by MMU on reads and writes, cleared by software)
- **PCD:** cache disabled (1) or enabled (0)
- **PWT:** write-through or write-back cache policy for this page
- **U/S:** user or supervisor mode access
- **R/W:** read-only or read-write access
- **P:** page is present in physical memory (1) or not (0)

## 18.8 Core i7 cache

### 18.8.1 cache access

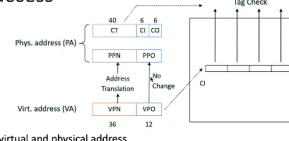
L1 cache access

- Partition physical address: CO, CI, and CT
- Use CT to determine if line containing word at address PA is cached in set CI
- **No:** check L2 (then L3...)
- **Yes:** extract word at byte offset CO and return to processor



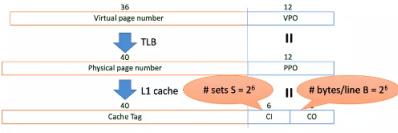
Speeding up L1 access

- Bits that determine CI **identical** in virtual and physical address
- Can index into cache **while** address translation taking place
- Generally we hit in TLB, so PPN bits (CT bits) available next
- "Virtually indexed, physically tagged"
- Cache carefully sized to make this possible



## 18.8.2 Cache size

Why the cache doesn't get bigger



- $\log_2(\text{cache size}) = \text{bits(CI)} + \text{bits(CO)} + \log_2(\text{associativity})$
- For Core i7:  $6 + 6 + 3 = 15 \Rightarrow 32\text{kB}$
- For performance:  $\text{bits(CI)} + \text{bits(CO)} \leq \text{bits(VPO)}$

As seen in the formula we can increase the cache size by increasing the associativity. This has the tradeoff that we must do more tag comparisons, and hence more hardware is needed to keep this in parallel.

## 18.9 Caches revisited

Until now we have assumed the cache only sees a virtual or physical address. But indexing and tagging can use different addresses:

- Virtually-indexed, virtually tagged (VV)
- Virtually-indexed, physically tagged (VP)
- Physically-indexed, virtually tagged (PV)
- Physically-indexed, physically tagged (PP)

Virtually-indexed means you index the cache based on bits available from the virtual address i.e before translation. Virtually tagged means that the tag you use in the cache to check if this is really a hit comes from either the virtual address in a virtually tagged cache or from the physical address in a physically tagged cache after translation.

PV is complex and rare and hence won't be discussed in further detail.

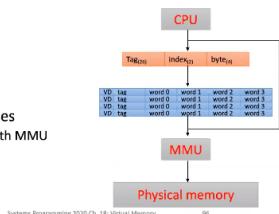
PP is the most straightforward implementation. Its disadvantage though is that we must wait for the translation to happen before we can access the tag.

### 18.9.1 Virtually-indexed, virtually tagged

#### Virtually-indexed cache issues

##### Virtually indexed, virtually tagged

- Also called:
  - Virtually addressed
- Also (misleadingly!)
  - Virtual cache
  - Virtual address cache
- Only uses virtual addresses
  - Operates concurrently with MMU



- Homonyms:** same names for different data
- VA used for indexing is context dependent
  - Same VA refers to different PAs
  - Tag does not uniquely identify data
  - OS must prevent this!
- Homonym prevention:**
  - Flush cache on context switch
  - Force non-overlapping address-space layout
  - Tag VA with address-space ID (ASID)

With a virtually indexed cache, we can access the cache at the same time as the MMU does the translation.

#### Summary: VV caches

##### Summary: VV caches with keys

- Add **address space identifier** (ASID) part of tag

On access compare with CPU's ASID register

Removes homonyms, creates synonyms (different VAs map to same PA)

Potentially better context switching performance

ASID recycling still requires a cache flush

Doesn't solve synonym problem

(but that's less serious)

Doesn't solve write-back problem

ASIDs much more common in TLBs!

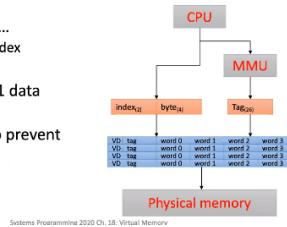
- Fastest:** don't rely on TLB for retrieving data
  - Still need TLB lookup for protection
  - Or other mechanism to provide protection
- Suffer from synonyms and homonyms**
  - Requires flush on context switch
    - Makes context switches expensive
    - May even be required on kernel  $\Rightarrow$  user switch
    - ... or guarantee of no synonyms or homonyms
- Require TLB lookup for writeback!**
- Used for i-cache on many architectures
  - Alpha, Pentium 4, etc.
- Rare as d-cache today...
  - Used on i860, ARM7/ARM9/StrongARM/Xscale...

Synonyms are a problem because we might have multiple occurrences of the same block in cache. This is a problem when writing because this would cause inconsistencies because the same block would be written and not written.

## 18.9.2 Virtually-indexed, physically tagged

Virtually indexed, physically tagged

- What we've just seen...
  - Virtual addr  $\Rightarrow$  line index
  - Physical addr  $\Rightarrow$  tag
- Commonly used for L1 data caches
- Requires alignment to prevent homonyms



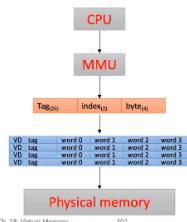
Systems Programming 2020 Ch. 18: Virtual Memory

The CPU can start indexing the cache, while the MMU does the translation, but we must wait for the MMU to generate the tag in order to start comparisons in the cache.

## 18.9.3 Physically-indexed, physically tagged

Physically indexed, physically tagged

- Only uses physical addresses
- Translation must complete before cache access can start
- Typically used off-core
- Note: page offset invariant under virtual address translation
  - Index bits  $\subseteq$  offset
  - Cache accessed without translation
  - Can be used on-core



Systems Programming 2020 Ch. 18: Virtual Memory

Slow but easy to use as there are no homonyms or synonyms. Frequently used as L2 or L3 cache.

## 18.9.4 Write Buffers

Write buffers

- Store operations take long to complete
  - E.g. if cache line must be read or allocated
- Can avoid stalling CPU by buffering writes
- **Write buffer** is FIFO queue of incomplete stores
  - also called *store buffer* or *write-behind buffer*
- Can also read intermediate values out of buffer
  - To service load of a value that is still in the write buffer
  - Avoids unnecessary stalls of load operations
- Implies that memory contents are temporarily stale
  - On a multiprocessor, CPUs see different order of writes
  - “weak store order”, to be revisited in SMP context!

