

Systems Programming and Computer Architecture

Robin Bacher, Janis Hutz

<https://github.com/janishutz/eth-summaries>

January 29, 2026

Rust and the like have an unsafe block... C's equivalent to this is

```
1 int main( int argc, char *argv[] ) {  
2     // Unsafe code goes here  
3 }
```

i.e. ***YOU are the one that makes C code safe!***

"If you are using CMake to solve the exercises... First off, sorry that you like CMake"

- Timothy Roscoe, 2025

HS2025, ETHZ

Summary of the Lectures and Lecture Slides

Quotes

"An LLM is a lossy index over human statements"

- Professor Buhmann, Date unknown

"If you are using CMake to solve the exercises... First off, sorry that you like CMake"

"You can't have a refrigerator behave like multiple refrigerators"

"Why is C++ called C++ and not ++C? It's like you don't get any value and then it's incremented, which is true"

- Timothy Roscoe, 2025

Contents

1	Introduction	5
2	x86 Assembly	6
2.1	The syntax	7
2.1.1	Registers	7
2.1.2	Instructions	7
2.2	Data types	8
2.2.1	Arrays	8
2.2.2	Structures	8
2.2.3	Nested / Multidimensional arrays, Struct arrays	9
2.2.4	Multi-Level arrays	9
2.2.5	Unions	9
2.3	Operations	10
2.3.1	Arithmetic Operations	10
2.3.2	Condition Codes	10
2.3.3	Jumping	11
2.4	Control Flow	12
2.4.1	Conditional statements	12
2.4.2	While Loops	12
2.4.3	For Loops & Switch	13
2.5	The Stack	14
2.5.1	Calling Conventions	14
2.5.2	Examples	15
2.6	Unorthodox Control Flow	16
2.7	Coroutines	18
3	The C Programming Language	20
3.1	Basics	20
3.1.1	Control Flow	21
3.1.2	Declarations	22
3.1.3	Operators	25
3.1.4	Arrays	26
3.1.5	Strings	26
3.1.6	Integers in C	27
3.1.7	Pointers	28
3.2	The C preprocessor	30
3.3	Memory	31
3.3.1	Dynamic Memory Allocation	32
3.3.2	Garbage Collection	34
3.3.3	Common pitfalls	34
3.4	Variadic functions	35
3.5	Code Vulnerabilities	36
3.5.1	Buffer overflow	36
3.5.2	Return-oriented Programming	37
3.6	Floating Point	38
3.6.1	Fractional Binary Numbers	38
3.6.2	Floating Point Representation	38
3.6.3	Properties	39
3.6.4	Rounding	40
3.6.5	Operations	41
3.6.6	Mathematical Properties	41
3.6.7	Floating Point in C	41
4	The gcc toolchain	42
4.1	Compiler optimizations	43
4.1.1	Optimization blockers	43

4.2	Linking	44
4.2.1	Symbol Resolution	44
4.2.2	Relocation	44
4.2.3	Packaging Libraries	45
4.3	File types	46
4.3.1	Executable and Linkable Format (ELF)	46
5	Hardware	47
5.1	Code Optimizations	47
5.2	Vector Operations	48
5.3	Caches	48
5.3.1	Performance Metrics	48
5.3.2	Cache misses	49
5.3.3	Cache organization	49
5.3.4	Determining cache set from memory address	49
5.3.5	Memory reads	50
5.3.6	Memory writes	50
5.3.7	Writing fast code	50
5.3.8	Cache Addressing Schemes	50
5.4	Virtual Memory	51
5.4.1	Address Translation	51
5.4.2	Multilevel Page Tables	52
5.4.3	x86 Virtual Memory	52
5.5	Exceptions	54
5.5.1	Synchronous Exceptions	54
5.5.2	Asynchronous Exceptions	54
5.6	Multi-Core	55
5.6.1	Background	55
5.6.2	Limitations	55
5.6.3	Coherency and Consistency	56
5.6.4	MSI protocol	57
5.6.5	MSI extensions	58
5.6.6	Relaxing Sequential Consistency	59
5.6.7	Multicore synchronization	60
5.6.8	Symmetric Multiprocessing	61
5.6.9	Non-Uniform Memory Access	61
5.6.10	Performance and Optimization	62
5.7	Devices	63
5.7.1	Device Registers	63
5.7.2	Direct Memory Access	64
5.7.3	Device Drivers	64

1 Introduction

This summary tries to summarize everything that is important to know for this course. It aims to be a full replacement for the slides, but as with all my summaries, there may be missing or incorrect information in here, so use at your own risk. You have been warned!

The summary does *not* follow the order the lecture does. This is to make related information appear more closely to each other than they have in the lecture and the summary assumes you have already seen the concepts in the lectures or elsewhere (or are willing to be thrown in the deep end).

The target semester for this summary is HS2025, so there might have been changes in your year. If there are changes and you'd like to update this summary, please open a pull request in the summary's repo at

<https://github.com/janishutz/eth-summaries>

2 x86 Assembly

Definition: (*Architecture*) Also known as ISA (Instruction Set Architecture) is “The parts of a processor design that one needs to understand to write assembly code”. It includes for example the definition of instructions (and their options) and what registers are available. Notable examples are x86, RISC-V (this one is open-source!), MIPS, ARM, etc

x86-Assembly files usually use the file extension `.s` or `.asm`

Most CPU architectures (including x86, most ARM and base RISC-V implementations) are little endian (also sometimes called right-to-left or right-aligned) and thus have the LSB at the lowest address.

Definition: (*Microarchitecture*) The implementation of the ISA. It defines the actual hardware layout and how the individual instructions are actually implemented and thus also defines things such as core frequency, cache layout and more.

Thus, the ISA is more or less precisely on the boundary of the software/hardware interface.

Definition: Complex Instruction Set (CISC):

- Stack oriented instruction set: Uses it to pass arguments, save program counter and features explicit push and pop instructions for the stack.
- Arithmetic instructions can access memory
- Condition codes set side effect of arithmetic and logical instructions.
- Design Philosophy: Add new instructions for typical tasks.

Definition: Reduced Instruction Set (RISC):

- Fewer, simpler instructions, commonly with fixed-size encoding. As a result, we might need more to get a given task done. On the other hand, we can execute them with small and fast hardware
- Register-oriented instruction set with many more registers that are used for arguments, return pointers, temporaries, etc.
- Load-Store architecture, i.e. only load and store instructions can access memory
- Thus: No Condition codes

What to choose? Both have advantages that the other has as disadvantage: Compiling for CISC is usually easier and usually results in smaller code size. For RISC however, compiler optimization can give a huge performance uplift and it can run fast with even a simple chip design.

Today, the choices are made based on outside constraints usually. For desktops and servers, there is enough compute to make anything run fast. For embedded systems though, the reduced complexity of RISC makes more sense, but for how long still?

What matters most today are non-technical factors such as existence of code for one ISA or licensing costs (and of course, Geopolitics)

Example A C function that simply adds 2 arguments might be compiled (unoptimized) to this:

File: `code-examples/01_asm/01_sum.s`

```

1  sum:  # Label
2      endbr64                # Indirect Branch target (No effect on code's logic)
3      pushq %rbp             # Preserve caller stack frame
4      movq %rsp, %rbp        # Set up new stack frame
5      movl %edi, -20(%rbp)    # Arg1 register -> stack
6      movl %esi, -24(%rbp)    # Arg2 register -> stack
7      movl -20(%rbp), %edx    # stack[Arg1] -> gp register
8      movl -24(%rbp), %eax    # stack[Arg2] -> gp register
9      addl %edx, %eax         # Add Arg1 + Arg2 -> gp register
10     movl %eax, -4(%rbp)     # result -> stack
11     movl -4(%rbp), %eax     # stack[result] -> return register
12     popq %rbp              # Restore caller stack frame
13     ret                    # jump back to caller

```

2.1 The syntax

There are two common styles: AT&T syntax (common on UNIX) and Intel syntax (common on Windows). The most obvious difference between the two is that they invert the order of operands, i.e. where the AT&T syntax has the destination as the second argument, the Intel syntax puts it first.

The state that is visible to us is:

- PC (Program Counter) that contains the address of the next instruction
- Register file that contains the most used program data
- Condition codes that store status information about most recent arithmetic operation and are used for conditional branching

To view what C code looks like in assembly, we can use `gcc -O0 -S code.c`, which produces `code.s` which contains assembly code.

2.1.1 Registers

x86 assembly is a bit particular with register naming (register names all start in %). The initial 16-bit version of x86 had the following registers (sub registers are registers that can be used to access the high (h suffix) or low (l suffix) half of the register. Only registers ending in x feature these sub registers. They, as well as %si and %di are general purpose):

Name	Sub-registers	Description
%ax	%ah, %al	accumulate
%cx	%ch, %cl	counter
%dx	%dh, %dl	data
%bx	%bh, %bl	base
%si	-	Source index
%di	-	Destination index
%sp	-	Stack pointer
%bp	-	Base pointer
%ip	-	Instruction pointer
%sr	-	Status (flags)

When the architecture was extended to 32-bit, all registers previously available were retained and a 32 bit version of each was introduced with the prefix `e`. In other words, any 16 bit code would still work as previously, as e.g. the %ax register was simply now the lower 16 bits of the %eax register.

The same happened again when extending to 64-bit, only this time the `r` prefix was used. So, the register %eax was now the lower 32 bits of %rax. Additionally, the following registers are also available, with X to be substituted with 8 through 15: %rX and the lower 32 bits %rXd

2.1.2 Instructions

Instructions usually have a 3 letter mnemonic with a one letter postfix that indicates the number of bytes. The following postfixes are available: b (byte, 1 byte), w (word, 2 bytes), l (long word, 4 bytes) and q (quad, 8 bytes).

The following options can be passed for source and destination: Registers,

Immediates To use a constant value (aka Immediate) in an instruction, we prefix the number with \$ (following number is decimal). To use hex, we can use \$0x, for binary, we can use \$0b, etc.

Memory addresses To treat a register as a memory address, use parenthesis, e.g. (%rax) interprets the value of %rax as a memory address. The instruction will then read the number of bytes, as specified by the postfix of the instruction.

The full syntax for memory address modes is `D(Rb, Ri, S)`, where

- D: Displacement (constant offset), can be 0, 1, 2 or 4 bytes (not bits, if you are confused as I was)
- Rb: Base register (to which offsets, etc are added). Can be any of the 16 integer registers
- Ri: Index register: Any, except for %rsp (and %rbp is also rarely used)
- S: Scale factor (1, 2, 4 or 8, to correct offsets)

The computation that happens is the following: `Mem[Reg[Rb] + S * Reg[Ri] + D]`. Using the `lea src, dest` instruction, we can get the address computed into the dest register. Can be abused for similar arithmetic expressions.

2.2 Data types

Assembly supports the following integer types (where GAS stands for GNU Assembly). If they are signed or unsigned does not matter (as we will see in section 3.1.6), so it's up to you to interpret them as one or the other

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

These integer types are also used for pointer addresses. Assembly also supports floating point numbers. They are stored and operated on in floating point registers.

Intel	GAS	Bytes	C
single	s	4	float
double	l	8	double
extended	t	16	long double

Assembly does not support any aggregate types (such as arrays, structs, etc) natively. You can however (obviously) make your own. In the following section we will cover how C datatypes are compiled into assembly. Do note that the `sizeof` function in C returns the number of bytes.

2.2.1 Arrays

Arrays of type `T` and length `L` are allocated as a contiguous region of memory with size `L * sizeof(T)` bytes. We then also store a reference / identifier `A` to the array (i.e. similar to variable name in C), that holds the address of the first element of the array and can then be used in conjunction with “assembly pointer arithmetic”.

Array loops that are written as for-loops in code are usually transformed into do-while loops by the compiler to save one condition check in the beginning, except of course, it might be possible that the loop is never executed.

2.2.2 Structures

We again allocate a contiguous region of memory. Only now, the number of bytes required isn't as straightforward to compute anymore, but still relatively simple: We simply sum up the sizes of all members and that will be our required sizes, so for the n members x_i of struct `my_struct`, we have $\text{sizeof}(\text{my_struct}) = \sum_{i=0}^{n-1} \text{sizeof}(x_i)$.

However, the size of a struct may be different to fulfill alignment requirements set forth by the ISA or operating system. This could mean that the struct takes $n * K$ bytes, where K is the alignment of the largest element

For alignment on x86-64 we have:

- 1 byte (no restrictions)
- 2 bytes (LSB must be 0)
- 4 bytes (2 LSB must be 00)
- 8 bytes (3 LSB must be 000)
- 16 bytes (4 LSB must be 0000)

On 64 bit Linux, using `gcc`, structs are compiled according to this concept and the size is $n * k$ bytes.

Another issue is accessing members. The solution to this is however easy and efficient, as at compile time, the offsets are pre-determined and compiled into the setter and/or getter code for the struct.

2.2.3 Nested / Multidimensional arrays, Struct arrays

All of these arrays have similar underlying concepts in the way they are allocated, yet all are a bit different

Common ideas Each of the array's elements are allocated in contiguous regions of memory, with the elements also in contiguous regions of memory. (Imagine it as lining up all elements on a band, i.e. as going through the array in a nested loop and printing all the elements into a single line.) The size of the array is determined by $n * \text{sizeof}(T)$, where T is the type of the elements of the array (or outer array). This is what is different for the lot (as well as accessing elements):

Nested array T is another array. We thus have a recursive definition, where $\text{sizeof}(T)$ resolves to $n * \text{sizeof}(T1)$, etc. Accessing element i, j, k is handled as follows: $o = i * \text{sizeof}(T) + j * \text{sizeof}(T1) + k * \text{sizeof}(T2)$, with $T1$ and $T2$ the types of the nested arrays

Struct arrays T is a struct.

2.2.4 Multi-Level arrays

In comparison to multidimensional arrays, we have arrays of pointers that contain either more arrays of pointers, (normal) arrays or pointers to other data types. The size of such a Multi-Level array is determined by: $n * \text{sizeof}(\text{ptr})$, where $\text{sizeof}(\text{ptr})$ is the platform-specific size of a pointer and n is the number of elements in the array

To do an access, we need to do two (or more) memory reads, which we can again do using address computations.

The benefit of these kinds of arrays is that we can store arbitrary data types together in an array, giving us more flexibility.

2.2.5 Unions

Since unions can hold any of the elements listed (but only one at a time), we allocate based on the size of the largest element.

2.3 Operations

Assembly operations include performing arithmetic or logic functions on registers or memory data, transferring data between memory and registers and transferring control (conditional or unconditional jumps).

Note that move instructions *cannot* move data directly from memory to memory.

The following instruction formats are commonly used:

File: code-examples/01_asm/00_operations.s

```
1 main:
2     movq %rax, %rdx # rax is src, rdx is dest
3     cmp  %rax, %rdx # rax is src2, rdx is src1
4     jmp  func      # func is label
5     notq %rax      # rax is dest (and src)
6     ret  # No argument
```

2.3.1 Arithmetic Operations

Arithmetic / logic operations need a size postfix (replace X below with b (1B), w (2B), l (4B) or q (8B)).

Mnemonic	Format	Computation
addX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} + \text{Src}$
subX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} - \text{Src}$
imulX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} * \text{Src}$
salX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} \ll \text{Src}$
sarX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} \gg \text{Src}$ (arithmetic)
shrX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} \gg \text{Src}$ (logical)
xorX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} \wedge \text{Src}$
andX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} \& \text{Src}$
orX	Src, Dest	$\text{Dest} \leftarrow \text{Dest} \mid \text{Src}$
incX	Dest	$\text{Dest} \leftarrow \text{Dest} + 1$
decX	Dest	$\text{Dest} \leftarrow \text{Dest} - 1$
negX	Dest	$\text{Dest} \leftarrow -\text{Dest}$
notX	Dest	$\text{Dest} \leftarrow \sim \text{Dest}$

2.3.2 Condition Codes

Any arithmetic operation (that is truly part of the arithmetic operations group, so not including lea for example) implicitly sets the **condition codes**. The following condition codes were covered in the lecture (operation: $t = a + b$):

- CF (Carry Flag): Set if carry out from MSB (unsigned overflow)
- ZF (Zero Flag): Set if $t == 0$
- SF (Sign Flag): Set if $(a - b) < 0$ (for signed)
- OF (Overflow Flag): Set if two's complement overflow (i.e. $(a > 0 \ \&\& \ b > 0 \ \&\& \ t < 0) \mid (a < 0 \ \&\& \ b < 0 \ \&\& \ t \geq 0)$)

Explicit computation In the below explanations, we always assume $\text{src2} = b$ and $\text{src1} = a$

To explicitly compute them, we can use the `cmpX src2, src1` (i.e. for easier understanding, `cmpX b, a`) instruction (with X again any of the size postfixes), that essentially computes $(a - b)$ without setting a destination register.

When we execute that instruction, CF is set if $a < b$ (unsigned), ZF is set if $a == b$, SF is set if $a < b$ (signed) and OF is set as above, where $t = a - b$.

Another instruction that is used is `testX src2, src1` (X again a size postfix, easier: `testX b, a`), and acts like computing $a \& b$ and where ZF is set if $a \& b == 0$ and SF is set if $a \& b < 0$.

Example: Assume $x = \text{src1}$, $y = \text{src2}$ in `cmpl y, x`. Then, when you use the condition codes with `cmovle`, the instruction is executed only if $x \leq y$ (i.e. the operands of the compare instructions in the AT&T syntax are flipped)

Zeroing register We can use a move instruction, but that is less efficient than using `xorl reg, reg`, where `reg` is the 32-bit version of the reg we want to zero. This works because on 32-bit operations, the upper 32 bit of the 64 bit register will be zeroed.

Reading condition codes To read condition codes, we can use the `setC` instructions, where the `C` is to be substituted by an element of table 1

2.3.3 Jumping

To jump, use `jmp <label>` (unconditional jump) or the `jC <label>` instructions, with `C` from table 1

setC	Condition	Description
e	ZF	Equal / Zero
ne	\sim ZF	Not Equal / Not Zero
s	SF	Negative
ns	\sim SF	Nonnegative
g	\sim (SF \wedge OF) $\&$ \sim ZF	Greater (signed)
ge	\sim (SF \wedge OF)	Greater or equal (signed)
l	SF \wedge OF	Less (signed)
le	(SF \wedge OF) \mid ZF	Less or equal (signed)
a	\sim CF $\&$ \sim ZF	Above (unsigned)
b	CF	Below (unsigned)

Table 1: Condition code postfixes for jump and set instructions

Conditional Moves

Similar to `jC`, the same postfixes can be applied to `cmovC`, for example:

```
cmpl    %eax, %edx # computes (edx - eax) without overwriting edx
cmovle  %edx, %eax # only copies edx into eax if edx <= eax
```

Due to the computed condition flags, this will move `%edx` into `%eax`, only if `%edx` is less than or equal (`le`) to `%eax`.

This can be used to, for example, compile ternary expressions from `C` to Assembly. However, this requires evaluating both possible expressions, which may lead to a (significant) performance overhead.

2.4 Control Flow

Control flow structures from C like if/else or for are compiled into assembly mainly using jumps and conditional move.

By the nature of Assembly and thanks to compilers optimizing aggressively, there is no *single* definitive translation of the C control structures: The compiler may translate it very differently depending on the context of the program.

2.4.1 Conditional statements

A function using an if/else construct to choose the maximum of 2 numbers might compile like this:

File: code-examples/01_asm/02_max.s

```

1 max:
2     cmpl %esi, %edi    # Set condition flags
3     jle .IF           # Conditional jump if %edi <= %esi
4     movl %edi, %edx    # %edi -> return register
5     jmp .ELSE
6 .IF:
7     movl %esi, %edx    # %esi -> return register
8 .ELSE:
9     ret

```

A function computing the absolute difference $|x - y|$ using an if/else construct, might use a conditional move instead:

File: code-examples/01_asm/03_absdiff.s

```

1 absdiff:
2     movl %edi, %eax
3     subl %esi, %eax    # arg2 - arg1 -> eax
4     movl %esi, %edx
5     subl %edi, %edx    # arg1 - arg2 -> edx
6     cmpl %esi, %edi    # Set condition flags
7     cmovle %edx, %eax  # edx -> eax, only if eax <= edx
8     ret

```

2.4.2 While Loops

A recursive factorial function using a do while loop may be compiled like this:

File: code-examples/01_asm/04_factorial.s

```

1 factorial:
2     movl $1, %eax      # Setup
3 .AGAIN:
4     imull %edi, %eax    # %eax *= %edi
5     subl $1, %edi      # %edi--
6     cmpl $1, %edi      # 1 < %edi ?
7     jg .AGAIN          # go back, if yes
8     ret

```

The same function, using a while loop instead may lead to this:

File: code-examples/01_asm/05_factorial.s

```
1 factorial:
2     jmp .COMP          # Check condition
3 .LOOP:
4     imull %edx, %eax    # %eax *= %edx
5     decl %edx          # %edx--
6 .COMP:
7     cmpl $1, %edx      # 1 < %edx ?
8     jg .LOOP           # if yes, go to loop
9     ret
```

2.4.3 For Loops & Switch

for loops follow the same idea as while loops, albeit with a few more jumps.

switch statements are implemented differently depending on size and case values: Sparse switch statements are compiled as decision trees, whereas compact switch statements may become *jump tables*.

These jump tables are usually stored in the `.rodata` section with 8-byte alignment. The jump uses offsets:

```
jmp *.LABEL(, %rsi, 8)
```

and we jump to the effective address of `.LABEL + rsi * 8`

2.5 The Stack

The Stack is the main way to dynamically allocate memory in Assembly, i.e. for temporary values. This process is completely manual: Allocating/Deallocating memory on the stack is entirely explicit. The stack grows *downwards*: Addresses decrease as the stack grows. The Stack pointer `%rsp` points to the current top of the stack.

Using the Stack `pushX`, `popX` exist for $x = 1, 2, 4, 8$, each corresponding to a size prefix that is set with X

Stack push `pushX src`: Fetch operand at `src`, decrement `%rsp` by x , then writes the operand at address of `%rsp`

Stack pop `popX dest`: Fetch operand at address of `%rsp`, increment `%rsp` by x , then writes the operand into `dest`

So intuitively, `pushX` and `popX` do exactly what is expected.

Procedure call / return Use `call LABEL`. This pushes the return label to the stack and jumps to `LABEL`. After this instruction, we also may use the `pushX` instruction to store further registers. Just remember to `pop` in the correct order with the correct size again!

The `ret` instruction is the return instruction and it will jump back to the caller and execution will continue there.

2.5.1 Calling Conventions

Even though Assembly will never stop you from using registers in any way, **System V ABI** specifies a few conventions, standardizing especially how functions are allowed to use registers.

Caller/Callee The callee is the function that is called and the caller is the code / function that calls the function.

- `%rax` and `%eax` are caller saved (usually used as return)
- `%rdi`, `%rsi`, `%rdx`, `%rcx` are caller saved (usually used for arguments)
- `%rsp` should not be modified manually
- `%rbp` is callee saved and used as frame pointer: usually set to `%rsp` at start of procedure and used to access frame elements (should always point to the start of the frame during function)

Register Conventions

Name	Description
<code>%rax</code>	Return value, #variable args
<code>%rbx</code>	Base pointer, Callee saved
<code>%rcx</code>	Argument 4
<code>%rdx</code>	Argument 3 (and return 2)
<code>%rsi</code>	Argument 2
<code>%rdi</code>	Argument 1
<code>%rsp</code>	Stack pointer
<code>%rbp</code>	Frame pointer, Callee saved

Name	Description
<code>%r8</code>	Argument 5
<code>%r9</code>	Argument 6
<code>%r10</code>	Static chain pointer
<code>%r11</code>	Temporary
<code>%r12</code>	Callee saved
<code>%r13</code>	Callee saved
<code>%r14</code>	Callee saved
<code>%r15</code>	GOT pointer, callee saved

If we have more than 6 arguments to be passed, we can use the stack for this. If we can do all accesses to the stack relative to the stack pointer, we do not need to update `%rbp` and not even `%rbx`, or we can use it for other purposes.

Manual stack management We can also allocate the entire stack frame directly by decrementing `%rsp` to the final position and then storing data relative to it. To deallocate a stack frame, simply increment the stack pointer to the value usually stored by `%rbp`

2.5.2 Examples

The stack is commonly used for recursive functions. A recursive factorial function might compile like this:

Note how `%rbx` is saved on the stack, since `rbx` is callee-saved by convention.

`%eax` is used directly, since it is caller-saved by convention.

File: `code-examples/01_asm/06_factorial.s`

```

1 factorial:
2     pushq %rbx                # Preserve frame pointer
3     movl %edi, %ebx
4     movl $1, %eax
5     cmpl $1, %edi
6     jle .QUIT                # Base case reached: quit
7     leal -1(%rdi), %edi       # Prepare args for next function call
8     call factorial
9     imull %ebx, %eax          # Use result of function call
10 .QUIT:
11     popq %rbx                # Restore frame pointer
12     ret

```

A more complex example, passing addresses as arguments:

This function swaps 2 array elements (using a swap function) and adds the first value to an accumulator.

File: `code-examples/01_asm/07_swap_and_sum.s`

```

1 swap_and_sum:
2     movq %rbx, -16(%rsp)      # Save %rbx
3     movslq %esi, %rbx         # Save i (and extend)
4     movq %r12, -8(%rsp)       # Save %r12
5     movq %rdi, %r12           # Save a
6     leaq (%rdi,%rbx,8), %rdi  # &a[i] -> %rdi (arg 1)
7     subq $16, %rsp            # Allocate stack frame
8     leaq 8(%rdi), %rsi        # &a[i+1] -> %rsi (arg 2)
9     call swap
10    movq (%r12,%rbx,8), %rax   # a[i]
11    addq %rax, sum(%rip)       # sum += a[i]
12    movq (%rsp), %rbx          # Restore %rbx
13    movq 8(%rsp), %r12         # Restore %r12
14    addq $16, %rsp             # Deallocate stack frame
15    ret

```

2.6 Unorthodox Control Flow

In C, the `setjmp.h` header file can be included, which gives us access to `setjmp` and `longjmp`.

To use them, we first need to declare a `jmp_buf` somewhere, usually as a static variable.

The `setjmp(jmp_buf env)` function stores the current stack / environment in the `jmp_buf` and returns 0.

The `longjmp(jmp_buf env, int val)` function causes a second return, which returns `val`, to the `setjmp` invocation and jumps back to that place.

File: `code-examples/01_asm/08_unorthodox-controlflow.c`

```
1  #include <setjmp.h>
2  #include <stdio.h>
3
4  static jmp_buf buf;
5
6  void second( void ) {
7      printf( "second\n" );
8      longjmp( buf, 1 );
9  }
10
11 void first( void ) {
12     second();
13     printf( "first\n" ); // Never executed
14 }
15
16 int main() {
17     if ( !setjmp( buf ) ) // returns 0 initially
18         first();
19     else
20         printf( "main\n" ); // 1 is returned when longjmp is executed
21     return 0;
22 }
```

What the above code outputs is: `second` followed by `main`.

They are implemented in Assembly as follows. Nothing really surprising for the implementation there. The assembly code is from the Musl C library

File: code-examples/01_asm/09_setjmp-longjmp.s

```
1  # Copyright 2011-2012 Nicholas J. Kain, licensed under standard MIT license
2  setjmp:
3      mov %rbx, (%rdi)
4      mov %rbp, 8(%rdi)
5      mov %r12, 16(%rdi)
6      mov %r13, 24(%rdi)
7      mov %r14, 32(%rdi)
8      mov %r15, 40(%rdi)
9      lea 8(%rsp), %rdx
10     mov %rdx, 48(%rdi)
11     mov (%rsp), %rdx
12     mov %rdx, 56(%rdi)
13     xor %eax, %eax
14     ret
15
16 longjmp:
17     xor %eax, %eax
18     cmp $1, %esi      # CF = val ? 0 : 1
19     adc %esi, %eax    # eax = val + !val
20     mov (%rdi), %rbx  # rdi is the jmp_buf, restore regs from it
21     mov 8(%rdi), %rbp
22     mov 16(%rdi), %r12
23     mov 24(%rdi), %r13
24     mov 32(%rdi), %r14
25     mov 40(%rdi), %r15
26     mov 48(%rdi), %rsp
27     jmp *56(%rdi)    # goto saved address without altering rsp
```

2.7 Coroutines

Coroutines are functions that call each other when they are done or they need new data to work on. An example is a decompressor that calls a parser when it has finished compressing parts of the file and that parser then again calls the decompressor when it has finished parsing.

We can implement that either by rewriting the functions into a single function, which often is a bit clumsy. A way around this is to use **Continuations**, where the first function saves its state and the context is switched to the other function. That function can then load its state and continue where it left off, runs until it finishes its task, then saves its state and the context switches back to the original function.

File: code-examples/01_asm/10_coroutine.h

```
1  #include <setjmp.h>
2  typedef void( co_start_fn )( void * );
3
4  struct coroutine {
5      void *stack;           // The call stack
6      jmp_buf env;          // The saved context
7      co_start_fn *start;    // Function to call
8      void *arg;             // Argument to the function
9  };
10 struct coroutine *co_new( co_start_fn *start, void *ctxt );
11 void co_free( struct coroutine *self );
12 void *co_switchto( struct coroutine *next );
13 void co_init( void );
```

File: code-examples/01_asm/10_coroutine.c

```

1  #include "10_coroutine.h"
2  #include <setjmp.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5  #define CORO_STACK_SIZE 128
6  static struct coroutine *cur_co;
7  static struct coroutine *main_co;
8
9  void co_init() {
10     main_co = (struct coroutine *) calloc( 1, sizeof( struct coroutine ) );
11     cur_co = main_co;
12     co_switchto( main_co );
13 }
14
15 void *co_switchto( struct coroutine *next, void *arg ) {
16     if ( setjmp( cur_co->env ) == 0 ) {
17         cur_co = next;
18         cur_co->arg = arg;
19         longjmp( cur_co->env, 1 );
20     }
21     return cur_co->arg;
22 }
23
24 static void start_cl( void ) {
25     ( cur_co->start )( cur_co->arg );
26     co_switchto( main_co );
27
28     printf( "Error: returned from coroutine start closure.\n" );
29     exit( -1 );
30 }
31
32 struct coroutine *co_new( co_start_fn *start, void *ctxt ) {
33     struct coroutine *co = (struct coroutine *) calloc( 1, sizeof( struct coroutine ) );
34     co->stack = calloc( 1, CORO_STACK_SIZE + 16 );
35     co->start = start;
36     co->arg = ctxt;
37     setjmp( co->env );
38     co->env[ 0 ].__jmpbuf[ 6 ] = ( (uint64_t) ( co->stack ) + CORO_STACK_SIZE );
39     co->env[ 0 ].__jmpbuf[ 7 ] = ( (uint64_t) ( start_cl ) ); // Both lines machine specific
40 }

```

As you can see, the setjmp.h functions are the foundation of all concurrent programming.

3 The C Programming Language

I can clearly C why you'd want to use C. Already sorry in advance for all the bad C jokes that are going to be part of this section

C is a compiled, low-level programming language, lacking many features modern high-level programming languages offer, like Object Oriented programming, true Functional Programming (like Haskell implements), Garbage Collection, complex abstract datatypes and vectors, just to name a few. (It is possible to replicate these using Preprocessor macros, more on this later).

On the other hand, it offers low-level hardware access, the ability to directly integrate assembly code into the .c files, as well as bit level data manipulation and extensive memory management options, again just to name a few.

This of course leads to C performing excellently and there are many programming languages whose compiler doesn't directly produce machine code or assembly, but instead optimized C code that is then compiled into machine code using a C compiler. This has a number of benefits, most notably that C compilers can produce very efficient assembly, as lots of effort is put into the C compilers by the hardware manufacturers.

There are many great C tutorials out there, a simple one (as for many other languages too) can be found [here](#)

3.1 Basics

C uses a very similar syntax as many other programming languages, like Java, JavaScript and many more...to be precise, it is *them* that use the C syntax, not the other way around. So:

File: code-examples/00_c/00_basics/00_intro.c

```

1 // This is a line comment
2 /* this is a block comment */
3 #include "01_func.h" // Relative import
4
5 int i = 0; // This allocates an integer on the stack
6
7 int main( int argc, char *argv[] ) {
8     // This is the function body of a function (here the main function)
9     // which serves as the entrypoint to the program in C and has arguments
10    printf( "Argc: %d\n", argc );    // Number of arguments passed, always >= 1
11                                    // (first argument is the executable name)
12    for ( int i = 0; i < argc; i++ ) // For loop just like any other sane programming language
13        printf( "Arg %d: %s\n", i, argv[ i ] ); // Outputs the i-th argument from CLI
14
15    get_user_input_int( "Select a number" ); // Function calls as in any other language
16    return 0;                               // Return a POSIX exit code
17 }
```

In C we are referring to the implementation of a function as a **(function) definition** (correspondingly, *variable definition*, if the variable is initialized) and to the definition of the function signature (or variables, without initializing them) as the **(function) declaration** (or, correspondingly, *variable declaration*).

C code is usually split into the source files, ending in .c (where the local functions and variables are declared, as well as all function definitions) and the header files, ending in .h, usually sharing the filename of the source file, where the external declarations are defined. By convention, no definition of functions are in the .h files, and neither variables, but there is nothing preventing you from putting them there.

File: code-examples/00_c/00_basics/01_func.h

```

1 #include <stdio.h> // Import from system path
2                 // (like library imports in other languages)
3
4 int get_user_input_int( char prompt[] );
```

3.1.1 Control Flow

Many of the control-flow structures of C can be found in the below code snippet. A note of caution when using goto: It is almost never a good idea (can lead to unexpected behaviour, is hard to maintain, etc). Where it however is very handy is for error recovery (and cleanup functions) and early termination of multiple loops (jumping out of a loop). So, for example, if you have to run multiple functions to set something up and one of them fails, you can jump to a label and have all cleanup code execute that you have specified there. And because the labels are (as in Assembly) skipped over during execution, you can make very nice cleanup code. We can also use continue and break statements similarly to Java, they do not however accept labels. (Reminder: continue skips the loop body and goes to the next iteration)

File: code-examples/00_c/00_basics/01_func.c

```

1  #include "01_func.h"
2  #include <stdio.h>
3
4  int get_user_input_int( char prompt[] ) {
5      int input_data;
6      printf( "%s", prompt );           // Always wrap strings like this for printf
7      scanf( "%d", &input_data );       // Get user input from CLI
8      int input_data_copy = input_data; // Value copied
9
10     // If statements just like any other language
11     if ( input_data )
12         printf( "Not 0" );
13     else
14         printf( "Input is zero" );
15
16     // Switch statements just like in any other language
17     switch ( input_data ) {
18         case 5:
19             printf( "You win!" );
20             break; // Doesn't fall through
21         case 6:
22             printf( "You were close" ); // Falls through
23         default:
24             printf( "No win" ); // Case for any not covered input
25     }
26
27     while ( input_data > 1 ) {
28         input_data -= 1;
29         printf( "Hello World\n" );
30     }
31
32     // Inverted while loop (executes at least once)
33     do {
34         input_data_copy -= 1;
35         printf( "Bye World\n" );
36         if ( input_data_copy == 10 )
37             goto this_is_a_label;
38     } while ( input_data_copy > 1 );
39
40     this_is_a_label:
41     printf( "Jumped to label" );
42     return 0;
43 }
```

3.1.2 Declarations

We have already seen a few examples for how C handles declarations. In concept they are similar (and scoping works the same) to most other C-like programming languages, including Java.

File: code-examples/00_c/00_basics/02_declarations.c

```

1  int my_int;                // Allocates memory on the stack.
2                               // Variable is global (read / writable by entire program)
3  static int my_local_int;  // only available locally (in this file)
4  extern const char *var;   // Defined in some other file
5  const int MY_CONST = 10;  // constant (immutable), convention: SCREAM_CASE
6
7  enum { ONE, TWO } num;    // Enum. ONE will get value 0, TWO has value 1
8
9  enum { O = 2, T = 1 } n;  // Enum with values specified
10
11 // Structs are like classes, but contain no logic
12 struct MyStruct {
13     int el1;
14     int el2;
15 };
16
17 // Like structs, but can only hold one of the values!
18 union MyUnion {
19     int ival;
20     float fval;
21     char *sval;
22 };
23
24 int fun( int j ) {
25     static int i = 0;        // Persists across calls of fun
26     short my_var = 1;        // Block scoped (deallocated when going out of scope)
27     int my_var_dbl = (int) my_var; // Explicit casting (works between almost all types)
28     return i;
29 }
30
31 int main( int argc, char *argv[] ) {
32     if ( ( my_local_int = fun( 10 ) ) ) {
33         // Every c statement is also an expression, i.e. you can do the above!
34     }
35     struct MyStruct test;      // Allocate memory on stack for struct
36     struct MyStruct *test_p = &test; // Pointer to memory where test resides
37     struct MyStruct test2;
38     union MyUnion my_uval;    // Work exactly like structs for access
39     test.el1 = 1;              // Direct element access
40     test_p->el2 = 2;            // Via pointer
41     test2 = test;              // Copies the struct
42     return 0;
43 }

```

A peculiarity of C is that the bit-count is not defined by the language, but rather the hardware it is compiled for.

C data type	typical 32-bit	ia32	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	4	8	8
long double	8	10/12	16

Table 2: Comparison of byte-sizes for each datatype on different architectures

Type format Be however aware that this table uses the LP64 format for the x86-64 sizes and this is the format all UNIX-Systems use (i.e. Linux, BSD, Darwin (the Mac Kernel)). 64 bit Windows however uses LLP64, i.e. `int` and `long` have the same size (32) and `long long` and pointers are 64 bit.

Integers By default, integers in C are signed, to declare an unsigned integer, use `unsigned int`. Since it is hard and annoying to remember the number of bytes that are in each data type, C99 has introduced the extended integer types, which can be imported from `stdint.h` and are of form `int<bit count>_t` and `uint<bit count>_t`, where we substitute the `<bit count>` with the number of bits (have to correspond to a valid type of course).

Literaly Literaly may easily lead to obscure errors, since they are typed.

```

1 int x = 16;           // 16 is interpreted as signed by default
2 unsigned int ux1 = 16; // casts 16 (signed int) to unsigned
3 unsigned int ux2 = 16U // Declare 16 as unsigned explicitly (safer)
4 float fx1 = 16;       // casts 16 (signed int) to float
5 float fx2 = 16f;      // Declare 16 as float explicitly (safer)
6 float fx3 = 16F;
7 float fx4 = 16.0;     // casts 16.0 (double) to float
8 double dx1 = 16.0     // No casting, 16.0 is treated as double by default
9 double dx2 = 16.0L    // Explicitly declare as "Long Float" (Double)

```

Note that integer expressions (without suffix) are given the smallest fitting integer type, but always at least `int`.

Number Systems By default, decimal, hexadecimal and octal are supported.

```

1 int x = 16;
2 int xh = 0xf;           // Hexadecimal, declared by leading 0x
3 int xo = 010;          // Octal, declared by leading 0
4 int xb = 0b1000;       // Binary, only supported as of C23, so this is not recommended

```

The same suffixes from above can be used for any representation.

Booleans Another notable difference of C compared to other languages is that C doesn't natively have a boolean type, by convention a `short` is used to represent it, where any non-zero value means `true` and 0 means `false`. Since boolean types are quite handy, the `!` syntax for negation turns any non-zero value of any integer type into zero and vice-versa. C99 has added support for a `bool` type via `stdbool.h`, which however is still an integer.

Implicit casts Notably, C doesn't have a very rigid type system and lower bit-count types are implicitly cast to higher bit-count data types, i.e. if you add a `short` and an `int`, the `short` is cast to `short` (bits 16-31 are set to 0) and the two are added. Explicit casting between almost all types is also supported. Some will force a change of bit representation, but most won't (notably, when casting to and from `float`-like types, minus to `void`)

Expressions Every C statement is also an expression, see above code block for example.

Void The `void` type has *no* value and is used for untyped pointers and declaring functions with no return value

Structs Are like classes in OOP, but they contain no logic. We can assign copy a struct by assignment and they behave just like everything else in C when used as an argument for functions in that they are passed by value and not by reference. You can of course pass it also by reference (like any other data type) by setting the argument to type `struct mystruct * name` and then calling the function using `func(&test)` assuming `test` is the name of your struct

Typedef To define a custom type using `typedef <type it represents> <name of the new type>`.

You may also use `typedef` on structs using `typedef struct <struct tag> <name of the new alias>`, you can thus instead of e.g. `struct list_el my_list; write list my_list;`, if you have used `typedef struct list_el list;` before. It is even possible to do this:

```
1 typedef struct list_el {
2     unsigned long val;
3     struct list_el *next;
4 } list_el;
5
6 struct list_el my_list;
7 list_el my_other_list;
```

Namespaces C has a few different namespaces, i.e. you can have one of the same name in each namespace (i.e. you can have `struct a`, `int a`, etc). The following namespaces were covered:

- Label names (used for `goto`)
- Tags (for `struct`, `union` and `enum` names)
- Member names (one namespace for each `struct`, `union` and `enum`)
- Everything else mostly (types, variable names, etc, including `typedef`)

3.1.3 Operators

The list of operators in C is similar to the one of Java, etc. In Table 3, you can see an overview of the operators, sorted by precedence in descending order. You may notice that the `&` and `*` operators appear twice. The higher precedence occurrence is the address operator and dereference, respectively, and the lower precedence is bitwise and multiplication, respectively.

Additionally, `+` and `-` also appear twice, the higher precedence one being a unary plus or minus, respectively, which is used to denote positive or negative numbers, and it can also be used to do a sort-of assertion that we have an arithmetic type using the preprocessor macro

```
#define CHECK_ARITHMETIC(x) +(x)
```

which will cause a compiler error if `x` is for example a pointer. Of course, the lower precedence `+` and `-` is addition and subtraction, respectively.

Very low precedence belongs to boolean operators `&&` and `||`, as well as the ternary operator and assignment operators.

Operator	Associativity
<code>() [] -> .</code>	Left-to-right
<code>! ~ ++ -- + - * & (type) sizeof</code>	Right-to-left
<code>* / %</code>	Left-to-right
<code>+ -</code>	Left-to-right
<code><< >></code>	Left-to-right
<code>< <= >= ></code>	Left-to-right
<code>== !=</code>	Left-to-right
<code>& (logical and)</code>	Left-to-right
<code>^ (logical xor)</code>	Left-to-right
<code> (logical or)</code>	Left-to-right
<code>&& (boolean and)</code>	Left-to-right
<code> (boolean or)</code>	Left-to-right
<code>? : (ternary)</code>	Right-to-left
<code>= += -= *= /= %= &= ^=== <<= >>=</code>	Right-to-left
<code>,</code>	Left-to-right

Table 3: C operators ordered in descending order by precedence

In an expression like `x = (expr1, expr2)` the first expression is evaluated, its result is discarded, the second expression is executed and its result is returned. It is also possible to chain them, i.e. to use multiple without parenthesis, the last expression's value is ultimately returned. Of note is that the expressions have to be encased in parenthesis.

Associativity

- Left-to-right: $A + B + C \mapsto (A + B) + C$
- Right-to-left: $A += B += C \mapsto A += (B += C)$

As it should be, boolean and, as well as boolean or, support early termination.

The ternary operator works as in other programming languages `result = expr ? res_true : res_false;`

As previously touched on, every statement is also an expression, i.e. the following works

```
printf("%s", x = foo(y)); // prints output of foo(y) and x has that value
```

Pre-increment (`++i`, new value returned) and post-increment (`i++`, old value returned) are also supported by C.

C has an `assert` statement, but do not use it for error handling. The basic syntax is `assert(expr);`

3.1.4 Arrays

C compiler does not do any array bound checks! Thus, always check array bounds. Unlike some other programming languages, arrays are *not* dynamic length.

The below snippet includes already some pointer arithmetic tricks. The variable `data` is a pointer to the first element of the array.

File: `code-examples/00_c/00_basics/03_arrays.c`

```

1  #include <stdint.h>
2  #include <stdio.h>
3
4  int main( int argc, char *argv[] ) {
5      int data[ 10 ];           // Initialize array of 10 integers
6      data[ 5 ] = 5;           // element 5 is now 5
7      *data = 10;              // element 0 is now 10
8      printf( "%d\n", data[ 0 ] ); // print element 0 (prints 10)
9      printf( "%d\n", *data );    // equivalent as above
10     printf( "%d\n", data[ 5 ] ); // print element 5 (prints 5)
11     printf( "%d\n", *( data + 5 ) ); // equivalent as above
12     int multidim[ 5 ][ 5 ];     // 2-dimensional array
13                                 // We can iterate over it using two for-loops
14     int init_array[ 2 ][ 2 ] = {
15         {1, 2},
16         {3, 4}
17     };                          // We can initialize an array like this
18     int empty_arr[ 4 ] = {};    // Initialized to 0
19     return 0;
20 }
```

3.1.5 Strings

C doesn't have a string data type, but rather, strings are represented (when using ASCII) as char arrays, with length of the array $n + 1$ (where n is the number of characters of the string). The extra element is the termination character, called the null character, denoted `\0`. To determine the actual length of the string (as it may be padded), we can use `strlen(str, maxlen)` from `string.h`

File: `code-examples/00_c/00_basics/04_strings.c`

```

1  #include <stdio.h>
2  #include <string.h>
3
4  int main( int argc, char *argv[] ) {
5      char hello[ 6 ] = "hello";           // Using double quotes
6      char world[ 6 ] = { 'w', 'o', 'r', 'l', 'd', '\0' }; // As array
7
8      char src[ 12 ], dest[ 12 ];
9      strncpy( src, "ETHZ", 12 );           // Copy strings (extra elements will be set to \0)
10     strncpy( dest, src, 12 );             // Copy strings (last arg is first n chars to copy)
11     if ( strcmp( src, dest, 12 ) ) // Compare two strings. Returns 1 if src > dest
12         printf( "Hello World" );
13     strcat( dest, " is in ZH", 12 ); // Concatenate strings
14     return 0;
15 }
```

3.1.6 Integers in C

As a reminder, integers are encoded as follows in big endian notation, with x_i being the i -th bit and w being the number of bits used to represent the number:

- **Unsigned:** $\sum_{i=0}^{w-1} x_i \cdot 2^i$
- **Signed:** $-x_{w-1} \cdot 2^{w-1} + \sum_{i=0}^{w-1} x_i \cdot 2^i$ (two's complement notation, with x_{w-1} being the sign-bit)

The minimum number representable is 0 and -2^{w-1} , respectively, whereas the maximum number representable is $2^w - 1$ and $2^{w-1} - 1$. `limits.h` defines constants for the minimum and maximum values of different types, e.g. `ULONG_MAX` or `LONG_MAX` and `LONG_MIN`

We can use the shift operators to multiply and divide by two. Shift operations are usually *much* cheaper than multiplication and division. Left shift (`u << k` in C) always fills with zeros and throws away the extra bits on the left (equivalent to multiplication by 2^k), whereas right shift (`u >> k` in C) is implementation-defined, either arithmetic (fill with most significant bit, division by 2^k . This however rounds incorrectly, see below) or logical shift (fill with zeros, unsigned division by 2^k).

Signed division using arithmetic right shifts has the issue of incorrect rounding when number is < 0 . Instead, we represent $s/2^k = s + (2^k - 1) >> k$ for $s < 0$ and $s/2^k = s >> k$ for $s > 0$

In expressions, signed values are implicitly cast to unsigned

This can lead to all sorts of nasty exploits (e.g. provide -1 as the argument to `memcpy` and watch it burn, this was an actual exploit in FreeBSD)

Addition & Subtraction

A nice property of the two's complement notation is that addition and subtraction works exactly the same as in normal notation, due to over- and underflow. This also obviously means that it implements modular arithmetic, i.e.

$$\text{Add}_w(u, v) = u + v \bmod 2^w \quad \text{and} \quad \text{Sub}_w(u, v) = u - v \bmod 2^w$$

Multiplication & Division

Unsigned multiplication with addition forms a commutative ring. Again, it is doing modular arithmetic and

$$\text{UMult}_w(u, v) = u \cdot v \bmod 2^w$$

3.1.7 Pointers

On loading of a program, the OS creates the virtual address space for the process, inspects the executable and loads the data to the right places in the address space, before other preparations like final linking and relocation are done.

Stack-based languages (supporting recursion) allocate stack in frames that contain local variables, return information and temporary space. When a procedure is entered, a stack frame is allocated and executes any necessary setup code (like moving the stack pointer, see 2.5). When a procedure returns, the stack frame is deallocated and any necessary cleanup code is executed, before execution of the previous frame continues.

In C a pointer is a variable whose value is the memory address of another variable

Of note is that if you simply declare a pointer using type `* p;` you will get different memory addresses every time. The (Linux)-Kernel randomizes the address space to prevent some common exploits.

File: code-examples/00_c/00_basics/05_pointers.c

```

1  #include "01_func.h" // See a few pages up for declarations
2  #include <assert.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  void a_function( int ( *func )( char * ), char prompt[] ) {
7      ( *func )( prompt ); // Call function with arguments
8  }
9
10 int main( int argc, char *argv[] ) {
11     int x = 0;
12     int *p = &x;           // Get x's memory address
13     printf( "%p\n", p );   // Print the address of x
14     printf( "%d\n", *p );  // Dereference pointer (get contents of memory location)
15     *p = 10;               // Dereference assign
16     int **dbl_p = &p;      // Double pointer (pointer to pointer to value)
17     int *null_p = NULL;    // Create NULL pointer
18     *null_p = 1;           // Segmentation fault due to null pointer dereference
19
20     // pointer arithmetic
21     int arr[ 3 ] = { 2, 3, 4 };
22     char c_arr[ 3 ] = { 'A', 'B', 'C' };
23     int *arr_p = &arr[ 1 ];
24     char *c_arr_p = &c_arr[ 1 ];
25     c_arr_p += 1; // Now points to c_arr[2]
26     arr_p -= 1;   // Now points to arr[0]
27
28     char *arr_p_c = (char *) arr_p; // Cast to char pointer (points to first byte of arr[0])
29     printf( "%d", *( arr_p - 5 ) ); // No boundary checks (can access any memory)
30     assert( arr == &( arr[ 0 ] ) ); // Evaluates to true
31     int new_arr[ 3 ] = arr;          // Compile time error (cannot use other array as
        ↪ initializer)
32     int *new_arr_p = &arr[ 0 ];      // This works
33     void *t;
34
35     a_function( &get_user_input_int, c_arr );
36
37     return EXIT_SUCCESS;
38 }
```

Some pointer arithmetic has already appeared in section 3.1.4, but same kind of content with better explanation can be found here

Pointer Arithmetic Note that when doing pointer arithmetic, adding 1 will move the pointer by `sizeof(type)` bits. Pointer arithmetic with a `void` pointer is thus not allowed by standard C, as the compiler does not know the size of the data type. However, `gcc` does allow it and assumes the size of `void` is 1 byte.

You may use pointer arithmetic on whatever pointer you'd like (as long as it's not a null pointer). This means, you *can* make an array wherever in memory you'd like. The issue is just that you are likely to overwrite something, and that something might be something critical (like a stack pointer), thus you will get **undefined** behaviour! (This is by the way a common concept in C, if something isn't easy to make more flexible (example for `malloc`, if you pass a pointer to memory that is not the start of the `malloc`'d section, you get undefined behaviour), in the docs mention that one gets undefined behaviour if you do not do as it says, so...RTFM!)

As already seen in the section arrays (section 3.1.4), we can use pointer arithmetic for accessing array elements. The array name is treated as a pointer to the first element of the array, except when:

- it is operand of `sizeof` (return value is $n \cdot \text{sizeof}(\text{type})$ with n the number of elements)
- its address is taken (then `&a == a`)
- it is a string literal initializer. If we modify a pointer `char *b = "String";` to string literal in code, the "String" is stored in the code segment and if we modify the pointer, we get undefined behaviour

Fun fact : `A[i]` is always rewritten `*(A + i)` by compiler.

Function arguments Another important aspect is passing by value or by reference. You can pass every data type by reference, you can not however pass an array by value (as an array is treated as a pointer, see above).

Body-less loops

```
1 int x = 0;
2 while ( x++ < 10 ); // This is (of course) not a useful snippet, but shows the concept
```

Function pointers A function can be passed as an argument to another function using the typical address syntax with the `&` symbol is annotated as argument using type `(* name)(type arg1, ...)` and is called using `(*func)(arg1, ...)`.

3.2 The C preprocessor

To have gcc stop compilation after running through cpp, the C preprocessor, use `gcc -E <file name>`.

Imports Imports in C are handled by the preprocessor, that for each `#include <file1.h>`, the preprocessor simply copies the contents of the file recursively into one file. Note that this can easily lead to errors caused by multiple definitions. Using the cpp directive `#ifndef` can be used to avoid this.

Depending on if we use `#include <file1.h>` or `#include "file1.h"` the preprocessor will search for the file either in the system headers (on Linux typically in `/usr/include/`) or in the project directory. Be wary of including files twice, as the preprocessor will recursively include all files (i.e. it will include files from the files we included) It is however possible to prevent that from becoming an issue using this (called an include guard):

File: `code-examples/00_c/01_preprocessor/00_include-guards.h`

```

1  #ifndef MYLIB_H
2      #define MYLIB_H
3
4      // all code goes now here
5  #endif

```

The name you pick for the preprocessor macro used in the check, by convention usually has the form `NAME_H`. It should not start in an underscore (or multiple) and it should not contain more than one consecutive underscore.

The C preprocessor gives us what are called preprocessor macros, which have the format `#define NAME SUBSTITUTION`.

File: `code-examples/00_c/01_preprocessor/01_macros.c`

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #define FOO      BAZ
4  #define BAR( x ) ( x + 3 )
5  #define SKIP_SPACES( p )          \
6      do {                          \
7          while ( p > 0 ) { p--; }    \
8      } while ( 0 )
9  #define COMMAND( c ) { #c, c##_command } // Produces { "<val(c)>", "<val(c)>_command" }
10
11 #ifdef FOO // If macro is defined, ifndef for if not defined
12     #define COURSE "SPCA"
13 #else
14     #define COURSE "Systems Programming and Computer Architecture"
15 #endif
16
17 #if 1
18     #define OUT_HELLO // if statement
19 #endif
20
21 int main( int argc, char *argv[] ) {
22     int i = 10;
23     SKIP_SPACES( i );
24
25     printf( "%s", COURSE );
26
27     return EXIT_SUCCESS;
28 }

```

To avoid issues with semicolons at the end of preprocessor macros that wrap statements that cannot end in semicolons,

we can use a concept called semicolon swallowing. For that, we wrap the statements in a `do ... while(0)` loop, which is removed by the compiler on compile, also taking with it the semicolon.

There are also a number of predefined macros:

- `__FILE__`: Filename of processed file
- `__LINE__`: Line number of this usage of macro
- `__DATE__`: Date of processing
- `__TIME__`: Time of processing
- `__STDC__`: Set if ANSI Standard C compiler is used
- `__STDC_VERSION__`: The version of Standard C being compiled
- ... many more

In headers, we typically use `#ifndef __FILENAME_H_` followed by a `#define __FILENAME_H_` or the like to check if the header was already included before

3.3 Memory

In comparison to most other languages, C does not feature automatic memory management, but instead gives us full, manual control over memory. This of course has both advantages and disadvantages.

File: `code-examples/00_c/02_memory/00_memory.c`

```

1  #include <stdlib.h>
2
3  int main( int argc, char *argv[] ) {
4      long *arr = (long *) malloc( 10 * sizeof( long ) ); // Allocate on heap
5      if ( arr == NULL )                                // Check if successful
6          return EXIT_FAILURE;
7      arr[ 0 ] = 5;
8
9      long *arr2;
10     if ( ( arr2 = (long *) calloc( 10, sizeof( long ) ) ) == NULL )
11         return EXIT_FAILURE; // Same as above, but fewer lines and memory zeroed
12
13     // Reallocate memory (to change size). Always use new pointer and do check!
14     if ( ( arr2 = (long *) realloc( arr2, 15 * sizeof( long ) ) ) == NULL )
15         return EXIT_FAILURE;
16
17     free( arr );    // Deallocate the memory
18     arr = NULL;    // Best practice: NULL pointer
19     free( arr2 );  // *Can* omit NULLing pointer because end
20
21     long arr3[5];  // Allocate on Stack
22                   // Deallocated automatically when returning
23
24     return EXIT_SUCCESS;
25 }
```

Notably, the argument `size_t sz` for `malloc`, `calloc` and `realloc` is an unsigned integer of some size and differs depending on hardware and software platforms.

`malloc` keeps track of which blocks are allocated. If you give `free` a pointer that isn't the start of the memory region previously `malloc`'d, you get undefined behaviour.

Memory corruption There are many ways to corrupt memory in C. The below code shows off a few of them:

File: code-examples/00_c/02_memory/01_mem-corruption.c

```

1  #include <stdlib.h>
2
3  int main( int argc, char **argv ) {
4      int a[ 2 ];
5      int *b = (int *) malloc( 2 * sizeof( int ) ), *c;
6      a[ 2 ] = 5;           // assign past the end of an array
7      b[ 0 ] += 2;          // assume malloc zeroes out memory
8      c = b + 3;           // mess up your pointer arithmetic
9      free( &(amp; a[ 0 ] ) ); // pass pointer to free() that wasn't malloc'ed
10     free( b );
11     free( b ); // double-free the same block
12     b[ 0 ] = 5; // use a free()'d pointer
13     // any many more!
14     return 0;
15 }
```

Memory leaks If we allocate memory, but never free it, we use more and more memory (old memory is inaccessible)

Dynamic data structures We build it using structs that have a pointer to another struct inside them. We have to allocate memory for each element and then add the pointer to another struct. For a generic dynamic data structure, make the element a void pointer. This in general is the concept used for functions operating on any data type.

3.3.1 Dynamic Memory Allocation

Memory allocated with malloc is typically 8- or 16-byte aligned.

Explicit vs. Implicit In explicit memory management, the application does both the allocation *and* deallocation memory, whereas in implicit memory management, the application allocates the memory, but usually a *Garbage Collector* (GC) frees it.

For some languages, like Rust, one would assume that it does implicit allocation, but Rust is a language using explicit management, it's just that the *compiler* and not the programmer decides when to allocate and when to deallocate.

Assumption in this course: Memory is **word** addressed (= 8 Bytes on 64-bit platform).

Goals The allocation should have the highest possible throughput and at the same time the best (i.e. lowest) possible memory utilization. This however is usually conflicting, so we have to balance the two.

Definition: **Aggregate payload** P_k : All malloc'd stuff minus all free'd stuff

Definition: **Current heap size** H_k : Monotonically non-decreasing. Grows when sbrk system call is issued.

Definition: **Peak memory utilization** $U_k = (\max_{i < k} P_i) / H_k$

A bit problem for the free function is to know how much memory to free without knowing the size of the to be freed block. This is just one of many other implementation issues:

- (1) How much memory to free? → Headers
- (2) How do we keep track of the free blocks? I.e. where and how large are they? → Free lists
- (3) What do we do with the extra space of a block when allocating a smaller block? → Coalescing
- (4) How do we pick a block? → Placement policies
- (5) How do we reinsert a freed block into the heap? → When to coalesce

This all leads to an issue known as **fragmentation**

Definition: **Internal Fragmentation**: If for a given block the payload (i.e. the requested size) is smaller than the block size. This depends on the pattern of previous requests and is thus easy to measure

Definition: **External Fragmentation**: There is enough aggregate heap memory, but there isn't a single large enough free block available This depends on the pattern of future requests and is thus hard to measure.

Header : Stores size of block and is usually placed in the word that precedes the allocated block (standard method).

Free lists

M1 **Implicit list** using length: Links all blocks and uses a low-order bit to indicate free / allocated, as for aligned blocks, a / some low-order bit(s) are always 0).

M2 **Explicit list** among free blocks using a pointer in the first (and possibly second) word of the block

M3 **Segregated free list**: different free lists for different size classes

M4 **Blocks sorted by size**: Using a balanced tree with pointers within each free block and the length is the key

Definition: Coalescing Connecting two (or more) (free) blocks to form a larger (free) block.

We can do this efficiently in one direction with just a header, however in both directions requires what are referred to as boundary tags. They are simply headers on both sides of the block and this allows us to traverse backwards.

We can do coalescing in constant time, by looking at the previous block's footer and next block's header to check if they are free or not.

- If the previous block is free, we can coalesce it by updating its header to include the length of the to be freed block and the middle two words. Same update has to happen to the to be freed block's footer.
- If the next block is free, update its footer to the size of the to be freed block plus the free block's size plus the two words in the middle. Do the same update to the to be freed block's header.

If both blocks are free, then of course we can do this step in one go for both.

Using the headers or boundary tags is just one option to do it and it can be optimized.

Implicit Free List If we use the size that is stored in the header, we know where the next block is going to be already. To know if a block is allocated, we use a low-order bit. This is possible, since if the blocks are aligned, then some of the low-order bits are always 0.

Allocating Blocks: We traverse the list, and choose depending on the placement policy: *first fit*, *next fit* or *best fit* policies. When a block is picked, we might want to split it by adding a header to the remaining part.

Freeing Blocks: We can simply clear the allocated flag, which might lead to fragmentation. Thus: Coalesce the freed blocks.

This is the simplest approach, the price we pay is that allocating is linear w.r.t *all blocks* managed by the allocator.

Explicit Free List To improve on this, here we maintain pointers to the next and preferably previous free block(s). This means that all free blocks form a (doubly) linked list, and we don't traverse already allocated blocks. Since the free blocks aren't used by the program, the list pointers can be stored inside the free blocks.

Allocating Blocks: We remove the element from the list by updating the pointers and we can again use the *first fit*, *next fit* or *best fit* policies.

Freeing Blocks: Requires updating the list pointers, which depends on the list ordering. If boundary tags are used for coalescing, we also have to clear the allocated bit in the footer and header.

The list ordering may be LIFO (last-in-first-out) or address-ordered. LIFO (with *first fit*) allows constant time freeing, while address-ordered provides better memory utilization (linear time freeing).

To prevent fragmentation, we want to use coalescing again, which can be implemented using boundary tags (again leads to linear time coalescing, if the list is FIFO). This time, the pointers need to be updated too. How this works (and how fast) depends on the list ordering.

Segregated Free List An issue with the previous approaches is finding a *large enough* block, which is why finding a fit remains in linear time. Here, we keep separate lists for different *size classes* to make this faster. For each *size class* a separate free list is maintained. For example, the classes might be:

$$\{1\}, \{2\}, \{3, 4\}, \{5 - 8\}, \dots, \{1025 - 2048\}, \{2049 - 4096\}, \{4097 - \infty\} \quad (\text{Partitions by Powers of 2})$$

Allocating Blocks: We first check the list for the requested size class, and if no fitting block is found, move up to the next list. If we have reached the last list and there is no suitable block, request new memory using `sbrk()`.

Freeing Blocks: Requires adding the newly freed block into the respective free list, potentially after coalescing.

This leads to an increased throughput and better memory utilization as compared to the previous two. Additionally, this structure opens up many optimization options.

The way splitting, coalescing, etc. is implemented varies a lot by implementation.

For example, we might choose to only coalesce if the length of a certain list has reached a certain threshold.

Similarly, if a list of some size class is empty, we might either request new memory and create a new list, or split a block from a higher size class.

Some options used are: *Simple Segregated Storage*, *Segregated Fits*, *Buddy System*

Fun fact : the GNU malloc package uses a Segregated Free List, with *Segregated fits*.

Placement policies

- **First fit** Search the list from the beginning, pick first free block that fits. This will usually cause “splinters” at the beginning of the list and can take linear time in the total number of blocks (allocated and free)
- **Next fit** Like first fit, but start at point the previous search finished at. This should be faster, however leads to worse fragmentation.
- **Best fit** Searches the list and chooses the *best* free block that fits and has fewest bytes left over. Leads to lower fragmentation, but is slower than first fit.

3.3.2 Garbage Collection

The memory manager must somehow be able to tell what memory can be freed. In general, we cannot know if memory is going to be used or not, except if there exists no pointer to it anymore. Garbage collectors use graphs to track pointer availability. In other words, a block is reachable if there exists a path from a root node to it.

An easy GC algorithm is called **Mark and Sweep**. It has an extra bit in the header called the *mark bit* and can be built on top of malloc/free. The concept is to use malloc until we “run out of space” and to then run these steps:

- **Mark**: Starts at each root node and sets a mark bit on each reachable block.
- **Sweep**: Scan all blocks and free all blocks that are unmarked.

3.3.3 Common pitfalls

- Dereferencing bad pointers (e.g. passing an `int` to a function expecting a pointer)
- Reading uninitialized memory (memory allocated with `malloc` should be considered garbage)
- Overwriting memory (if you mess up pointer arithmetic or don't do boundary checks)
- Referencing nonexistent variables (variables go out of scope on function returns, except `static`)
- Freeing blocks multiple times (can corrupt the heap)
- Referencing freed blocks (always `NULL` pointers after using `free()`)
- Failing to free blocks (memory leak incoming and make sure to free the ENTIRE data structure!)

Some of these bugs (especially bad references) can usually be found using a debugger.

Substitute `malloc` with a `malloc` that has extra checking code (like UToronto CSRI `malloc` to detect memory leaks)

Another option is using `valgrind` (a memory debugger). Or, simply don't bother with C and use Rust.

3.4 Variadic functions

Variadic functions take a variable number of arguments and use the `...` syntax in C. A notable example of such a function is `printf`

File: `code-examples/00_c/03_others/00_variadic.c`

```
1  #include <stdarg.h> // Variadic function utilities
2  #include <stdio.h>
3
4  void print_int( unsigned int num_ints, char *msg, ... ) {
5      va_list ap; // keeps track of current argument, similar (in concept) to iterator
6      va_start( ap, msg ); // Initialize the iterator based on last fixed arg
7      for ( int i = 0; i < num_ints; i++ )
8          printf( "int %d = %d\n", i, va_arg( ap, int ) ); // Returns next arg cast to int
9      va_end( ap ); // Free up iterator
10 }
```

3.5 Code Vulnerabilities

A brief interjection on some code vulnerabilities.

System-level protections

- **Compiler-inserted checks** on functions
- **Randomized stack offsets**: Allocate *random* amount on stack before running the program
- **Nonexecutable segments**: Memory needs a special *execute* permission

3.5.1 Buffer overflow

Buffer overflows are a method for code injection on vulnerable code with specific buffer size-checking deficiencies. There are 2 ways to do this:

1. Change a function call or return address
2. Push malicious assembly onto the stack

For example, consider this code:

File: code-examples/00_c/05_vulnerabilities/01_buffer_overflow_echo.c

```

1  #include <stdio.h>
2
3  void echo() {
4      char buf[4];    // Limited size
5      gets(buf);      // Assumes size matches, does not check!
6      puts(buf);
7  }
8
9  int main()
10 {
11     printf("Type a string:"); // No size check enforced!
12     echo();
13     return 0;
14 }
```

This is a problem, since echo may be compiled to something similar to this:

File: 02_buffer_overflow_echo_asm.s

```

1  echo:
2      subq $24, %rsp    # Allocate stack space for buf
3      movq %rsp, %rdi
4      call gets
5      movq %rsp, %rdi
6      call puts
7      addq $24, %rsp
8      ret
```

Since buf is on the stack, and there is no size-enforcement when writing to buf, malicious input can write *before* %rsp, since the Stack grows downwards. This means stack memory that the program is intending to use again can be modified.

However, inserting executable assembly like this usually does not work, since the stack may not be executable due to missing system permission.

The vulnerability above could be fixed by using fgets(buf, 4, stdin) instead, which checks the size.

Heap overflow On the heap, buffer overflows work differently, as the heap contains no return addresses. However, the heap stores function pointers, which can be modified. Further, sophisticated attacks can use buffer overflow to potentially modify pointers in dynamically allocated memory.

3.5.2 Return-oriented Programming

Return-oriented Programming is a more sophisticated exploit, which does not rely on injecting any new code.

The key idea is: Overwrite return addresses and jump to *specific* machine instruction sequences *already present* in process memory.

3.6 Floating Point

Floating point numbers are a representation of real numbers.

Though there are many ways to accomplish this, *IEEE Standard 754* is used practically everywhere, also in x86. This standard is a little more complicated than fractional binary numbers, but has a few numeric advantages, especially for representing very large (very small) numbers.

float.exposed is an excellent website to understand floating point by example.

3.6.1 Fractional Binary Numbers

We can represent any real number (with a finite decimal representation) as:

$$d = \sum_{i=-n}^m 10^i \cdot d_i \quad \underbrace{d_m d_{m-1} \cdots d_1 d_0 . d_{-1} d_{-2} \cdots d_{-(n-1)} d_{-n}}_{d_i \text{ is the } i\text{-th digit of } d \text{ (neg. indices indicate decimals)}}$$

We can use the same idea for Base 2 as well:

$$b = \sum_{i=-n}^m 2^i \cdot b_i \quad b_m b_{m-1} \cdots b_1 b_0 . b_{-1} b_{-2} \cdots b_{-(n-1)} b_{-n}$$

To get an intuition for this representation, looking at some examples is helpful:

A few observations:

1. Shifting the dot right: Division by 2
2. Shifting the dot left: Multiply by 2
3. Numbers of the form $0.111\dots$ are just below 1.0
4. Some numbers representable in finite Base 10 are infinite in Base 2, e.g. $\frac{1}{5} = 0.20_{10}$

Binary	Fraction	Decimal
0.0	$\frac{0}{2}$	0.0
0.01	$\frac{1}{4}$	0.25
0.010	$\frac{2}{8}$	0.25
0.0011	$\frac{3}{16}$	0.1875
0.00110	$\frac{6}{32}$	0.1875
0.001101	$\frac{13}{64}$	0.203125
0.0011010	$\frac{26}{128}$	0.203125
0.00110101	$\frac{51}{256}$	0.19921875

A major issue with this representation is that very large (respectively very small) numbers require a large representation. E.g. $a_{10} = 5 \cdot 2^{100}$ has the representation $a_2 = 101 \underbrace{0000000000000000\dots}_{100 \text{ Zeros}}$. Floating Point is designed to address this.

3.6.2 Floating Point Representation

Floating point numbers instead use the representation:

$$a = \underbrace{(-1)^s}_{\text{Sign}} \cdot \underbrace{M}_{\text{Mantissa}} \cdot \underbrace{2^E}_{\text{Exponent}}$$

Single precision and Double precision floating point numbers store the 3 parameters in separate bit fields s, e, m :

Single Precision:

31: Sign	30 – 23: Exponent	22 – 0: Mantissa
----------	-------------------	------------------

Bias: 127, Exponent range: $[-126, 127]$

Double Precision:

63: Sign	62 – 52: Exponent	51 – 0: Mantissa
----------	-------------------	------------------

Bias: 1023, Exponent range: $[-1022, 1023]$

Most of the extra precision in 64b floating point numbers is associated to the mantissa. Note how double precision is necessary to represent all 32b signed Integers, and not all 64b signed Integers can be represented in either format.

The way these bitfields are interpreted *differs* based on the exponent field e :

1. **Normalized Values:** Exponent bit field e is neither all 1s nor all 0s.
In this case, E is read in *biased* form: $E = e - b$. The bias is $b = 2^{k-1} - 1$, where k is the number of bits reserved for e . This produces the exponent ranges $E \in [-(b-1), b]$.
The mantissa field m is interpreted as $M = 0.m_{n-1} \dots m_1 m_0 + 1$, where n is the amount of bits reserved for m .
2. **Denormalized Values:** Exponent bit field e is all 0s.
In this case, E is read in *biased* form $E = 1 - b$. (Instead of $E = e - b$)
The mantissa field m is interpreted as $M = 0.m_{n-1} \dots m_1 m_0$ (without adding 1)
3. **Special Values:** Exponent bit field e is all 1s.
 $m = 0$ represents infinity, which is signed using s .
 $m \neq 0$ is NaN, regardless of what is in m or s .

Why is the Bias chosen this way? It allows smooth transitions between normalized and denormalized values.

3.6.3 Properties

The advantage of having denormalized values is that 0 can be represented as the bit-field with all 0s. Further, this enforces equidistant points for values close to 0, whereas normalized values increase in distance as they move further from 0.

Example 8b Floating Point table to visualize the different cases.

8b precision Floating Point: $\underbrace{0}_s \underbrace{0000}_e \underbrace{000}_m$

Case	s	e	m	E	Value
Denormalized	0	0000	000	-6	0
	0	0000	001	-6	$\frac{1}{8} \cdot \frac{1}{64} = \frac{1}{512}$
	0	0000	010	-6	$\frac{2}{8} \cdot \frac{1}{64} = \frac{2}{512}$
			\vdots		\vdots
	0	0000	110	-6	$\frac{6}{8} \cdot \frac{1}{64} = \frac{6}{512}$
	0	0000	111	-6	$\frac{7}{8} \cdot \frac{1}{64} = \frac{7}{512}$
Normalized	0	0001	000	-6	$\frac{8}{8} \cdot \frac{1}{64} = \frac{8}{512}$
	0	0001	001	-6	$\frac{9}{8} \cdot \frac{1}{64} = \frac{9}{512}$
			\vdots		\vdots
	0	0110	110	-1	$\frac{14}{8} \cdot \frac{1}{2} = \frac{14}{16}$
	0	0110	111	-1	$\frac{15}{8} \cdot \frac{1}{2} = \frac{15}{16}$
	0	0111	000	0	$\frac{8}{8} \cdot 1 = 1$
	0	0111	001	0	$\frac{9}{8} \cdot 1 = \frac{9}{8}$
	0	0111	010	0	$\frac{10}{8} \cdot 1 = \frac{10}{8}$
			\vdots		\vdots
	0	1110	110	7	$\frac{14}{8} \cdot 128 = 224$
Special	0	1110	111	7	$\frac{15}{8} \cdot 128 = 240$
	0	1111	000	n/a	∞

3.6.4 Rounding

The basic idea of Floating Point operations is:

1. Compute exact result
2. Round, so it fits the desired precision

IEEE Standard 754 specifies 4 rounding modes: *Towards Zero*, *Round Down*, *Round Up*, *Nearest Even*.

The default used is *Nearest Even*¹, which rounds up/down depending on which number is closer, like regular rounding, but picks the nearest even number if it's exactly in the middle.

Rounding can be defined using 3 different bits from the *exact* number: G, R, S

$$a = 1.B_1B_2 \dots B_{n-2}B_{n-1} \underbrace{G}_{\text{Guard}} \underbrace{R}_{\text{Round}} \underbrace{X_1X_2 \dots X_{k-1}X_k}_{\text{Sticky}}$$

where n is the number of bits in the mantissa of the format (e.g. 3 as in the above example of an 8bit floating point number). Do keep in mind that when computing them in C in the standard way, they are inverted, i.e. we need a `!= 0` at the end to correct.

1. **Guard Bit** G is the least significant bit of the (rounded) result (i.e. it is B_n)
2. **Round Bit** R is the 1st bit cut off after rounding
3. **Sticky Bit** S is the logical OR of all remaining cut off bits X_i . The mask can be computed using $(1 \ll n) - 1$

Based on these bits the rounding can be decided (we increment the rounded part if the expression evaluates to true):

Round up: $R \wedge S$

Round to even: $G \wedge R \wedge \neg S$

It is notable that for round to even, the special condition only applies if the sticky bit is not set. If it is set, the round up condition is to be used. An easy way to implement the condition is as follows

```
(sticky && round) || (!sticky && round && guard)
```

This will be ever so slightly more efficient than a different order, as the computation will be stopped shorter if a condition is not fulfilled

Example Rounding 8b precise results to 8b precision floating point (4b mantissa):

Value	Fraction	GRS	Incr?	Rounded
128	1.000 0000	000	N	1.000
13	1.101 0000	100	N	1.101
17	1.000 1000	010	N	1.000
19	1.001 1000	110	Y	1.010
138	1.000 1010	011	Y	1.001
63	1.111 1100	111	Y	10.000

Post-Normalization: Rounding may cause overflow. In this case: Shift right once and increment exponent.

¹Changing the rounding mode is usually hard to do without using Assembly.

3.6.5 Operations

Multiplication is straightforward, all 3 parameters can be operated on separately:

$$(-1)^{s_1} M_1 \cdot 2^{E_1} \cdot (-1)^{s_2} M_2 \cdot 2^{E_2} = (-1)^{s_1 \oplus s_2} (M_1 \cdot M_2) 2^{E_1 + E_2}$$

Post-Normalization:

1. If $M \geq 2$, shift M right and increment E
2. If E out of range, overflow (set to ∞)
3. Round M to fit desired precision.

Addition is more complicated: (Assumption: $E_1 \geq E_2$)

$$(-1)^{s_1} M_1 \cdot 2^{E_1} + (-1)^{s_2} M_2 \cdot 2^{E_2} = (-1)^{s'} M' \cdot 2^{E_1}$$

s', M' are the result of a signed align & add.

This means $(-1)^{s_1} M_1$ is shifted left by $E_1 - E_2$, and then $(-1)^{s_2} M_2$ is added.

Post-Normalization:

1. if $M \geq 2$, shift M right, increment E
2. if $M \leq 1$, shift M left k , decrement E by k
3. Overflow E if out of range (set to ∞)
4. Round M to desired precision

3.6.6 Mathematical Properties

Floating point is *almost* an Abelian Group.

- **Closed** under Addition, Multiplication (But may generate NaN, $\pm\infty$)
- **Commutative**
- **Not Associative** (Overflow & Rounding)
- **0, 1 are Identity**
- **Additive Inverse** (Except $\pm\infty$ and NaN)
- **Monotonicity** (Except for $\pm\infty$ and NaN)
- **Not Distributive** (Overflow & Rounding)

3.6.7 Floating Point in C

C99 guarantees float and double, and long double is usually interpreted as quadruple precision.

Casting/Conversion between Integer types and float, double *changes* the bit representation in most cases (e.g. 0 stays the same)

4 The gcc toolchain

The GNU Compiler Collection, or short (which is also its executable name) `gcc`, is a C compiler toolchain that is commonly used to compile C code on UNIX platforms.

It includes all the necessary tools to compile a C program:

- A preprocessor, called `cpp`
- The actual C compiler, called `cc` (though on the slides it states it is called `cc1`, but at least in the Arch Linux package, `cc1` is not a thing)
- An x86 assembler, called `as`
- A static linker, called `ld`

However, the individual parts are usually not called individually, but using the toolchain command `gcc` (and that is usually again abstracted away using CMake or Make, which is in turn commonly called via a build system like Meson to automatically build packages for distribution).

`gcc` has (as of GCC 15.2.1 20260103 on Arch Linux) about 1000 CLI arguments that can be passed. Below is a list of the most important flags that can be passed, as discussed in the lectures:

Flag	Description
<code>-E</code>	Stop after the preprocessor (output is a <code>.i</code> file)
<code>-S</code>	Stop after the compiler (output is assembly in <code>.s</code> file)
<code>-c</code>	Stop after the assembler (output is <code>.o</code> file)
<code>-o</code>	Specify the executable name
<code>-DNDEBUG</code>	Removes all assert statements
<code>-OX</code>	Optimization level where X can be one of 0, 1, 2, 3
<code>-g</code>	Compile with debugging information
<code>-Wall</code>	Enable common warnings
<code>-Wextra</code>	Enable more warnings
<code>-Werror</code>	Makes all warnings errors
<code>-march=XXX</code>	Optimize for the architecture (can be e.g. <code>native</code> , <code>x86-64-v4</code> , ...)
<code>-fno-tree-vectorize</code>	Do not vectorize code (<code>-O3</code> commonly enables vectorization)

Table 4: Command line flags for GCC

4.1 Compiler optimizations

While the compiler can do quite a bit to speed up code, it can't rework the core logic, as it has to guarantee that the executable does do what was specified in the code.

So, it is really important to not only consider asymptotic runtime (as $100n$ and $5n$ are both $\mathcal{O}(n)$, but obviously the latter is 20 times faster). We thus need to optimize the algorithms, data representations, loops, etc and for that, we need to properly understand how programs are compiled, executed and how the hardware works.

When using `gcc`, it is usually a good idea to compile a final build with the `-O2` or `-O3` flags.

The `-march` flag was already mentioned in table 4 and can be used if you want to go above and beyond, as it will optimize for the specific hardware. The values that can be passed to `-march` are listed [here](#) and even include a specific CPU microarchitecture. For example, to compile for Intel Alderlake (12000 series), you can specify `-march=alderlake`

To understand what you need to optimize, you need to understand what the compiler is good at:

- Register allocation
- Scheduling (i.e. code selection and ordering)
- Dead code elimination
- Eliminating minor (!) inefficiencies

and what it is not good at:

- Improving Asymptotic efficiency (compiler can't turn BubbleSort into e.g. QuickSort)
- Improving the constant factor (if your implementation is slow, it likely won't magically become faster, though some bad practices can be eliminated)
- Overcoming other optimization blockers such as memory aliasing and procedure side-effects

Code motion is a compiler technique, where it moves certain computations out of loops that always produce the same result. However: Always remember that the compiler will be *conservative*, i.e. it will always err on the side of caution.

Strength reduction is a compiler technique, where e.g. sequences of products are turned into cheaper additions in each iteration. An example is that if you have an operation such as $n * i$ in a loop, the compiler might replace that with a variable `ni` that is incremented by `n` in each iteration. Similarly, it might replace $16 * x$, or even worse still, $x / 16$ with `x << 4` or `x >> 4`, respectively

Common sub-expressions can be extracted into pre-computations and then only use cheaper operations on the individual steps. A good example is if you are using similar multiplications that then only require one addition or subtraction to get to a result close by.

4.1.1 Optimization blockers

A sure-fire way to make your code slow is by using a large number of procedure calls. They are among the slowest operations in C. And, the compiler cannot safely extract the function in a for loop like this:

```

1  int i;
2  for (i = 0; i < strlen(s); i++) {
3      if (s[i] >= 'A' && s[i] <= 'Z') {
4          s[i] -= ('A' - 'a');
5      }
6  }
```

The compiler can't safely remove `strlen(s)` from the loop, as it may have side-effects, i.e. may modify other program content other than simply returning a value. Thus, only ever call functions in the loop condition when you need the side-effects and otherwise, pre-compute it and simply use a variable to check against. You can declare a function *side-effect free* using `__attribute__((pure))` or `__attribute__((const))` (this is more strict, as the function is also not allowed to read global memory) in the function declaration. The compiler may then extract `strlen(s)` from the loop.

Another common blocker is memory aliasing. This happens when two pointers point to the same address and of course, since we can do pointer arithmetic, it is very easy to do that in C. The easiest way to prevent this from happening is to use local variables where possible, such that they do not need to be passed in using a pointer.

Normally the compiler assumes there can be another pointer that accesses the memory pointed to by this pointer. If you use the `restrict` keyword on the variable (i.e. in a function declaration, we have `void test(double restrict *a)`), the compiler will assume that for the lifetime of this pointer, there are no other pointers that will be used to access the memory to which it points.

Another technique to improve throughput for something like matrix multiplications is to do it in blocks due to the way caching works. Since the compiler doesn't *understand* your code, it can't do this for you (as it assumes associativity of the operation)

4.2 Linking

Linking is the final step in the compilation pipeline: separately compiled object files are combined into an executable.

The advantages of using Linkers are clear:

1. **Separate Compilation:** Changing one source file requires only recompiling that file.
2. **Space Optimization:** Executable code only contains functions (e.g. from libraries) that are actually used.

4.2.1 Symbol Resolution

The first step during Linking is Symbol Resolution.

In the context of Linking, all variables and functions are considered *Symbols*. Compilers store all symbol definitions in a *Symbol Table*. The linker associates symbol references with *exactly one* definition.

Definition: Symbol types

- **Global Symbols** can be referenced by other modules (e.g. `non-static` in C)
- **External Symbols** are referenced globals defined elsewhere
- **Local Symbols** are defined and referenced exclusively in one module (e.g. `static` in C)

Note: Local linker symbols and local program variables are *not* the same.

Definition: Symbol strength

Duplicate symbols either lead to linking errors (`-fno-common`, the default) or compile (`-fcommon`)

- **Strong Symbols** are procedure names and initialized globals
- **Weak Symbols** are uninitialized globals (on `-fcommon`)

in C, function symbols can explicitly be declared weak using:

```
#pragma weak func
__attribute__((weak)) __ func()
```

Duplicate Handling

The linker uses these definitions to handle duplicates:

1. Given multiple strong symbols are illegal
2. Given a strong symbol and multiple weak symbols, pick the strong symbol
3. Given multiple weak symbols, choose an *arbitrary* one

4.2.2 Relocation

The second step during Linking is Relocation.

Code and data sections of separate sources are combined, and symbols are relocated from relative locations (in `.o` files) to absolute locations (in executable files)

Command line order matters for this, since the Linker will scan `.o` and `.a` files in the order they are appear in the CLI arguments. In general, libraries should therefore be linked *last*.

4.2.3 Packaging Libraries

Using just the Linker, there are only 2 inconvenient ways to package libraries:

1. All functions into 1 file \mapsto linking unnecessarily big objects.
2. One function per file \mapsto Requires linking a lot of files, annoying for programmer.

Static Libraries solve this: The linker looks for functions inside the static library, and only links matching archive *members* into the executable. However, these come with issues too:

1. Duplication in stored executables (e.g. libc.a functions)
2. Duplication in running executables
3. Any fix in a library requires importing applications to explicitly relink

Shared Libraries solve this: These are linked at load-time or during run-time. Another advantage is that *multiple* processes can use the same shared library simultaneously. This is how, for example, libc is packaged.

During runtime, shared libraries can be loaded using dlopen:

File: code-examples/00_c/04_toolchain/01_dynamic_linking.c

```

1  #include <stdio.h>
2  #include <dlfcn.h> // contains addvec
3
4  int x[2] = {1, 2}; int y[2] = {3, 4}; int z[2];
5
6  int main(int argc, char *argv[])
7  {
8      void *handle;
9      void (*addvec)(int *, int *, int *, int); // Declaration
10     char *error;
11
12     handle = dlopen("./libvector.so", RTLD_LAZY); // Load .so, makes addvec usable
13     if (!handle) {
14         fprintf(stderr, "%s\n", dlerror());
15         exit(1);
16     }
17
18     addvec = dlsym(handle, "addvec"); // get a pointer to advec
19     if ((error = dlerror()) != NULL) {
20         fprintf(stderr, "%s\n", error);
21         exit(1);
22     }
23
24     addvec(x, y, z, 2); // Now callable like any other function
25     printf("z = [%d %d]\n", z[0], z[1]);
26
27     if (dlclose(handle) < 0) { // Unload shared library
28         fprintf(stderr, "%s\n", dlerror());
29         exit(1);
30     }
31     return 0;
32 }
```

4.3 File types

The most common file types used during compilation:

- **Source Code File** (.c) Uncompiled source code in C.
- **Relocatable Object File** (.o) Code & Data in a format ready for linking.
- **Shared Object File** (.so) Special object file, can be loaded & linked dynamically: at load or run time.
- **Executable File** Code & Data in a format that can be directly copied into memory & run.
- **Archive files** (.a) concatenate related .o files into one .a, with an index.

Alternate names On Windows, .dll files are used instead of .so and are called *Dynamic Link Libraries*.

4.3.1 Executable and Linkable Format (ELF)

The standard unified format for all object files (.exe, .o, .so) in use since UNIX.

Section	Content
ELF header	contains basic information: Word size, byte ordering, file type, machine type
Segment header table	page size, virtual address memory segments, segment sizes
.text	actual code
.rodata	(Read-only data) for example, jump tables
.data	initialized global variables
.bss	uninitialized global variables
.symtab	symbol table, procedure and static variable names, section names & locations.
.rel.text	relocation info for .text, e.g. addresses of instructions that need modifying
.rel.data	relocation info for .data, e.g. addresses of pointers that need modifying
.debug	info for symbolic debugging (gcc -g)

5 Hardware

Hardware was only touched on briefly in this course, and knowledge of a previous Computer Architecture course is assumed. This section is similarly incomplete, only focusing on the specific topics covered in the lecture.

5.1 Code Optimizations

By understanding the underlying hardware, simple code optimizations can be made to improve performance significantly.

To analyze program performance, a solid benchmark and performance metrics need to be chosen.

Common metrics are:

1. **Latency** i.e. how long does a request take?
2. **Throughput** i.e. how many requests (per second) can be processed?

For pipelined processors, throughput is more tricky to specify. Usually, the steady-state throughput is considered.

Definition: **Cycles per Element (CPE)** is used as a performance metric for operations on vectors/lists. The execution time t can then be formulated as: $t = \text{CPE} \cdot n + \text{Overhead}$. (n is the vector/list size)

Definition: **Cycles per Instruction (CPI)** is the inverse of *Instructions per Cycle*

CPI can be further divided into $\text{CPI} = \text{CPI}_{\text{Base}} + \text{CPI}_{\text{Stall}}$ assuming a pipelined processor, which may stall for many reasons: Data hazards, control hazards, memory latency ...

Definition: **Clock Cycle Time (CCT)** is the inverse of *Clock Frequency*

Using these we can define Program Execution Time: $t = n \cdot \text{CPI} \cdot \text{CCT}$ (n is the instruction count)

Definition: **Instruction Level Parallelism (ILP)** is the natural parallelism occurring through independent instructions.

Superscalar Processors capitalize on ILP by executing multiple instructions per cycle. The instructions are (usually) scheduled dynamically. This is particularly great since it requires no effort from the programmer.

Different instructions generally may take vastly different amounts of time, though this can be reduced via pipelining and multiple execution units.

Example This table is specific to an Intel Haswell CPU.

Instruction	Latency	Cycles/issue
Load/store	4	1
Int Add	1	1
Int Multiply	3	1
Int/Long Divide	3–30	3–30
FP Multiply	5	1
FP Add	3	1
FP Divide	3–15	3–15

Performance of these operations is **Latency bound** if they are sequential, **Throughput bound** if they can run parallel.

Example **Loop unrolling** can increase performance for latency bound operations. (Optimization flags *should* do this)

Unrolling can be done to an arbitrary degree, but there are diminishing returns at some point. An ideal unrolling factor must be found experimentally.

Example **Reassociation** can also improve performance for parallelizable operations, if it breaks some sequential dependency.

Example **Separate Accumulators** can improve performance for latency bound operations too, as separate load/store units may be used.

5.2 Vector Operations

Extreme performance gains beyond the results of the previous section can be gained using hardware vector registers on supported CPUs.

Example In Intel AVX2, 256b vector registers like %ymm0, %ymm1 can be used to perform component-wise single/double precision FP operations.

```
vaddsd %ymm0, %ymm1, %ymm1    # Comp.-wise 32b FP add
vaddsd %ymm0, %ymm1, %ymm1    # Comp.-wise 64b FP add
```

5.3 Caches

Processors generally improve quicker than memory speed does. Therefore, optimizing memory is necessary, and caches are one such method.

Example: The importance of caches can quickly be seen when looking at the memory hierarchy:

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	4 kB page	Main memory (RAM)	100	Hardware + OS
Buffer cache	4 kB 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

5.3.1 Performance Metrics

- **Miss Rate** is the fraction of memory references that are not found in the cache. Defined as $\frac{\text{misses}}{\text{accesses}} = 1 - \text{hit rate}$ and is typically 3-10% in L1 caches and less than 1% in L2 caches.
- **Hit Time** is the time required to deliver a line in the cache to the processor and typically is 1-2 cycles for L1 cache and 5-20 cycles for L2 caches.
- **Miss penalty** is the additional time required when a cache miss occurs. Typically around 50-200 cycles

Judging by these numbers, it makes a huge difference if we hit the cache and the speed difference can easily exceed a factor of 100x. Of note is as well that a 99% hit rate is twice as good as a 97% hit rate with a miss penalty of 100 cycles and a cache hit time of 1 cycle:

- **97% hits:** 1 cycle + 0.03 · 100 cycles = 4 cycles
- **99% hits:** 1 cycle + 0.01 · 100 cycles = 2 cycles

Thus, we always use *miss rate* instead of hit rate.

For a multi-level cache, we start with the last level cache and compute its miss penalty and combine that with the next higher level and so on (example with 2 level cache)

$$\begin{aligned}\text{MissPenaltyL2} &= \text{DRAMaccessTime} + \frac{\text{BlockSize}}{\text{Bandwidth}} \\ \text{MissPenaltyL1} &= \text{HitTimeL2} + \text{MissRateL2} \cdot \text{MissPenaltyL2} \\ \text{AverageMemoryAccessTime} &= \text{HitTimeL1} + \text{MissRateL1} \cdot \text{MissPenaltyL1}\end{aligned}$$

5.3.2 Cache misses

- **Compulsory / Cold miss** Occurs on the first access of a block (there can't be any data there yet)
- **Conflict miss** The cache may be large enough, but multiple lines may map to the current block, e.g. referencing blocks 0, 8, 0, 8, ... would miss every time if they are both mapped to the same cache line. This is the typical behaviour for most caches.
- **Capacity miss** The number of active cache blocks is larger than the cache
- **Coherency miss** See in section 5.6. They happen if cache lines have to be invalidated due in multiprocessor scenarios to preserve sequential consistency, etc

5.3.3 Cache organization

Structure Caches can be defined using S, E, B such that $S \cdot E \cdot B = \text{Cache Size}$.

$S = 2^s$ is the set count, $E = 2^e$ is the lines per set, and $B = 2^b$ is the number of bytes per cache block / cache line².

Each block has the following structure where v is the valid bit:

Cache line:

v	tag	data ($B = 2^b$ bytes per block)
-----	-----	-----------------------------------

Address The cache address can be separated into fields which dictate the cache location:

Address:

tag	set index	block offset
-----	-----------	--------------

Since we have $S = 2^s$ sets and $B = 2^b$ bytes per block, we need s bits for the set index, b bits for block offset. The remaining part (tag) is stored with the cache block and needs to match for a cache hit.

Do note that the cache address is the same as the physical memory address, we just refer to it as cache address when talking about the interpretation the cache uses for it. See section 5.3.4 for more information

Higher cache associativity helps reduce the number of conflict misses (often) by increasing the number of lines available for each block. It however, at the same cache size, reduces the number of available cache sets, which in turn may increase the number of conflict misses.

Definition: Direct-mapped i.e. $E = 1$ (1 cache line per set only).

Definition: N -way Set-Associative i.e. $E = N$ (N cache lines per set, in this course we primarily covered $N = 2$).

5.3.4 Determining cache set from memory address

The cache set a memory address is mapped to does not depend (directly) on the associativity of the cache³ To be able to compute the cache block that a memory request maps to, we need to know the number of sets and the line size.

We can then compute the set it maps to using $x \% S$, where $x = \text{addr} \gg b$ is the block number in the memory.

Thus, in the cache address, the tag corresponds to $y = x / S$, so it indicates which actual memory location is stored in the cache (as multiple different memory block numbers will map to the same cache set).

²The terms *cache block* and *cache line* will be used interchangeably and sometimes shortened to just *block* or *line*

³The only dependence it has on the associativity is that with the same cache size, the number of sets is reduced if the associativity is increased

5.3.5 Memory reads

In Direct-Mapped caches

1. We find the cache set and check if the valid bit is set. If not, skip to step 3, false
2. Check if the tag is equal to the requested tag.
3. If true, return the block at the correct offset, else evict the line and fetch the correct line (will be B bytes) and return the block at the correct offset.

In 2-way Set-Associative caches

1. Find the corresponding cache set and compare the tag to both blocks.
2. If one matches, check its valid bit. If none match, go to step 5
3. If valid, return the block at the correct offset.
4. If invalid, evict the line, fetch the correct one. Return to step 3
5. If no match, evict one of the two (choose using a replacement policy like LRU (Least Recently Used) or randomly if the other block is not invalid), fetch requested line, go to step 3.

5.3.6 Memory writes

Memory writes are just as slow, if not sometimes slower than memory reads. Thus, they also have to be cached to improve performance. Again, there are a few options to handle write caching and we will cover the two most prevalent ones:

- **Write-through** Here, the data is immediately written to main memory. The obvious benefit is that the data is always up-to-date in the memory as well, but we do not gain any speed from doing that, it is thus very slow.
- **Write-back** We defer the write until the line is replaced (or sometimes other conditions or an explicit cache flush). The obvious benefit is the increased speed, as we can benefit from the low cache access latency. We do however need a *dirty bit* to indicate that the cache line differs from main memory. This introduces additional complexity, especially in multi-core situations, which is why in that case, often a write-through mode is enabled for the variables that need atomicity.

Another question that arises is what to do on a *write-miss*?

- **Write-allocate** The data is loaded into the cache and is beneficial if more writes to the location follow suite. It is however harder to implement and it may evict an existing value from the cache. This is commonly seen with write-back caches.
- **No-write-allocate** This writes to the memory immediately, is easier to implement, but again slower, especially if the value is later re-read. This is commonly seen with write-through caches.

5.3.7 Writing fast code

Improving the speed of the code can be done by optimizing the code's locality (both temporal and spatial locality).

Computing a matrix product or iterating over a list can be much faster depending on the way you access the array and how it is stored.

A common tactic to improve the throughput for these kinds of operations is to make sure that the order of operations is correct (i.e. for a row-major matrix to iterate over the elements of a row in the inner loop), or to do block multiplication where you multiply small blocks at once.

5.3.8 Cache Addressing Schemes

The cache can see either the virtual or physical address, and the tag and index do *not* need to both use the physical/virtual address. If this seems confusing, first read and understand section 5.4

Indexing	Tagging	Code
Virtually Indexed	Virtually Tagged	VV
Virtually Indexed	Physically Tagged	VP
Physically Indexed	Virtually Tagged	PV
Physically Indexed	Physically Tagged	PP

5.4 Virtual Memory

Conceptually, Assembly operations treat memory as a very large contiguous array of memory: Each byte has an individual address.

```
movl    (%rcx), %eax    # Refers to a Virtual Address
```

In truth of course, this is an abstraction for the memory hierarchy. Actual allocation is done by the compiler & OS.

The main advantages are:

- Efficient use of (limited) RAM: Keep only active areas of virtual address space in memory
- Simplifies memory management for programmers
- Isolates address spaces: Processes can't interfere with other processes

The reason virtual memory is even feasible is that most programs have great locality. The performance is good as long as the total virtual memory that is actively being used does not exceed the available physical memory. If that happens, we speak of **Thrashing**, where the performance degrades significantly due to the large number of swaps occurring.

Another benefit of virtual memory is that we can use the automated virtual to physical mapping to simplify memory allocation and management, since we don't have to manually seek out free physical pages anymore. Additionally, that serves as protection, as the OS can choose to allow certain processes to share memory, whilst it can disallow others to access the data by updating the page tables correctly. For that, the page table entries (PTEs) are extended to also include some permissions.

As touched on already, this also allows for memory sharing, which can become useful for dynamic linking. These are just some of the benefits of virtual memory

5.4.1 Address Translation

Address translation happens in a dedicated hardware component: The Memory Management Unit (MMU). Thus, the CPU can use the virtual memory addresses and does not have to worry about translating to the physical addresses.

Virtual and Physical Addresses share the same structure, but the VPN (Virtual Page Number) is usually far longer than the PPN (Physical Page Number), since the virtual space is far bigger. The offsets match however.

Virtual:

V. Page Number	V. Page Offset
----------------	----------------

Physical:

P. Page Number	P. Page Offset
----------------	----------------

The Page Table (PT) (Located at a special Page Table Base Register (PTBR)) contains the mapping $VPN \mapsto PPN$. Page Table Entries (PTE) are cached in the L1 cache like any other memory word.

The Translation Lookaside Buffer (TLB) is a small hardware cache inside the MMU that is used to accelerate the PT lookups, which is typically faster than an L1 cache hit. The PT is usually stored in memory⁴ and contains for each address the corresponding physical address, as well as a valid bit, which indicates if the page is in memory.

If a page is not in memory a page fault is triggered, which transfers control to the OS, which then loads the page into memory, updates the page table and returns control back to the process. This and the inverse (i.e. unloading pages from memory onto the disk) is often referred to as *swapping*.

Due to the slowness of disks, page sizes are usually fairly large (4-8KB) typically, in some cases up to 4MB. The replacement policy algorithms are highly sophisticated and too complicated to be implemented in hardware and are thus usually handled by the operating system.

Address Translation with page hit The CPU requests a virtual memory address from the MMU. It fetches the PTE from memory and sends the physical address to the memory system, which sends the data to the CPU.

Address Translation with page fault During the check of the valid bit of the page table, the MMU find that it is not set. It thus triggers a page fault exception, which is then handled and a victim page (if necessary) is then picked and evicted (and if the dirty flag is set, it is paged out to disk). The handler then loads the new page into memory and updates the PT and the original instruction is then restarted on the CPU and the address translation will then succeed.

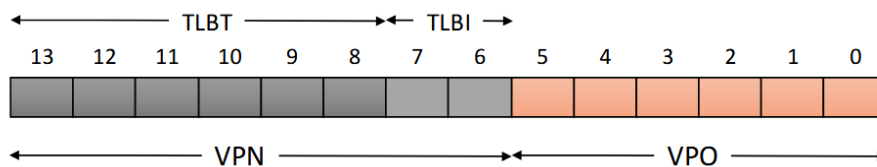
As already touched on, the TLB can be used to speed up translation. **Address Translation with TLB Hit** When checking the TLB for the entry, it find the entry and we save one memory access. With the PPN retrieved, the memory system sends the data to the CPU.

⁴In practice, most address translations actually hit the TLB.

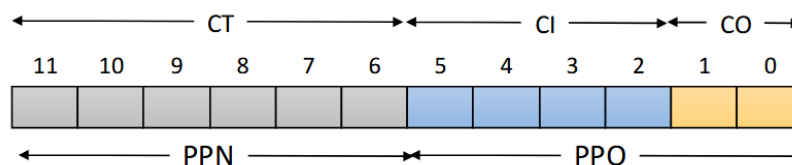
Address Translation with TLB Miss This works similar to the case when there is no TLB, as the TLB returns a miss signal for the request. Only that the PTE that is returned is inserted into the TLB via replacement policy (if applicable). The data is then fetched from the physical address and sent to the CPU.

Example We consider $N = 14$ bit virtual addresses and $M = 12$ bit physical addresses. The offset takes 6 bits.⁵

If we assume a TLB with 16 entries, and 4 way associativity, the VPN translates like this: (where TLBT = Tag and TLBI = Set)



Similarly, if we assume a direct-mapped 16 line cache with 4 byte blocks: (where CT = Tag, CI = Set and CO = Offset)



5.4.2 Multilevel Page Tables

Motivation For a 48-bit Virtual Address Space with 4KB ($= 2^{12}$ bytes) page size, the number of bits required is $2^{48} / 2^{12} \cdot 2^3 = 2^{39}$ bytes (that is 512 GB). The 2^3 bytes is the size of the page table entry (8 bytes).

Multi-Level page tables add further steps to this process: Instead of a PT we have a Page Directory Table (PDE) which contains the addresses of separate Page Tables. The top of the VPN is used to index into each of these, which technically allows any depth of page tables.

5.4.3 x86 Virtual Memory

In x86-64 Virtual Addresses are 48 bits long, yielding an address space of 256TB.

Physical Addresses are 52 bits, with 40 bit PPNs, yielding a page size of 4KB (we thus have 64 bit PTEs).

On the slides, they are again using (as far as we can tell) a Skylake CPU (Core 6000 series, could also be Kaby Lake, Core 7000 series).

On that architecture, the TLB contained the 40 bit PPN, a 32 bit TLB Tag, as well as

- a valid bit (V)
- a global bit (G, copied from PDE / PTE and prevents eviction)
- a supervisor-only bit (S, i.e. only accessible to OS, copied from PDE / PTE)
- a writable bit (W, page is writable, copied from PDE / PTE)
- a dirty bit (D, PTE has been marked dirty (i.e. modified vs memory))

There are a number of flags set and there are also a significant number of bits available for systems programmers to use on x86. Since they are highly unlikely to be exam-relevant, we will only point out that there are a lot of them (including setting supervisor mode, read/write mode, dirty, etc). To view them all, find them in the lecture slides of lecture 20, pages 87 through 90.

For many years, cache sizes have been stagnant, this was due to the limited number of bits that could be used to efficiently determine if a line is available in the cache. Today, there are techniques to overcome that limitation and we have seen fairly substantial increases in cache sizes since (primarily from Team Red, starting with the AMD Ryzen 7 5800X3D).

This was caused by the fact that only 6 bits from the PPO were used to determine the set and not more to improve performance, as the cache indexing could occur during address translation.

⁵The images in this example are from the SPCA lecture notes for FS25.

Addressing Schemes Revisited Returning to the Addressing Schemes from section 5.3.8, it becomes evident that this is the key to solving the issue just touched on. If we virtually tag and virtually index the address, we have the issue that there may exist multiple multiple PA for each VA (i.e. it is context dependent). To circumvent that issue, an ASID (Address Space Identifier) is added to the tag.

The Virtually Indexed, physically tagged scheme is what we have just seen and is commonly used for L1 caches.

The Physically Tagged, Physically Indexed is the solution to the cache size restriction. It however suffers from slower access times, as the address translation has to complete before the cache line can be identified.

Write buffers It is also common to have write buffers (which act like FIFO queue). It enables slower cache operations to complete that are typically associated with writing.

Large pages Simply a lower number of bits remain in the PPN and some bits are ignored (we increased the “offset-portion” to 21 bits = 2MB). We can increase that to a page size of 1GB, if we increase the “offset-portion” further, to 30 bits to be precise.

5.5 Exceptions

Control flow is mainly dictated by program state, and manipulated using jumps, branches, calls and returns. To react to changes in system state, exceptional control flow is used instead.

Low level mechanisms

- Hardware Exceptions
- Exceptions via combination of Hardware and OS software

High level mechanisms

- Process context switch
- Signals
- Nonlocal jumps
- Language-level exceptions (e.g. Java)

Generally, on an exception, control is transferred to a handler specific to the type of exception, which investigates the situation and returns control upon success. Mostly, this is handled via a *Exception Table* which is allocated on boot. On exception, this table is indexed depending on the type of exception to locate the corresponding handler. This causes a switch to Kernel Mode.

Definition: **Exception:** A control transfer to the OS in response to an event

- **Synchronous:** result of executing some instruction
- **Asynchronous:** result of an event 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

5.5.1 Synchronous Exceptions

Definition: **Trap** is an intentional exception that transfers control back to the next instruction

For example, opening a file in C executes a trap via a system call.

Definition: **Fault** is an unintentional, possibly recoverable exception. Either re-executes faulty instruction or aborts

For example, page faults, protection faults, floating point exceptions

Definition: **Abort** is unintentional and unrecoverable. Always aborts the program.

For example, a machine error.

5.5.2 Asynchronous Exceptions

Asynchronous Exceptions are indicated by setting the processor's (physical) interrupt pin.

For example,

- **Interrupts** are actions like network data arrival or hitting a key on the keyboard
- **Hard Reset Interrupts** are executed by hitting the system reset button
- **Soft Reset Interrupts** are caused by, for example, hitting CTRL+ALT+DEL (on Windows)

5.6 Multi-Core

5.6.1 Background

In the early days of computer hardware it was fairly easy to get higher performance due to the rapid advances in transistor technology. However, today, what is known as Moore's Law (i.e. that the transistor count of integrated circuits doubles every two years). However, due to power constraints and the slowing down of advances in transistor technology, the transistor count growth has slowed down quite a bit since the beginning of the century and is predicted to further stagnate as time goes on.

This leads to various issues, among others, the performance of CPUs isn't going up as quickly anymore as it used to. Additionally, due to power constraints, building faster and faster single-core CPUs is not possible and the advances in that field have slowed to a crawl.

To mitigate and offset these issues, manufacturers started to add multiple cores to parallelize operations. This however brings a whole host of new issues with it, for example, how do you make sure that no data races occur, how do you schedule, etc? These questions have mostly been answered in the course Parallel Programming, so we will not cover that here.

The only reason transistor count is still growing at a seemingly constant rate today is that manufacturers manage to cram more and more cores into a CPU. But even that has slowed down in recent years.

While in 2019 a highest core count AMD EPYC CPU (i.e. the EPYC 7742 from the ROME family) had 64 Zen 2 cores, in 2025 the highest core count EPYC CPU (i.e. the EPYC 9965 from the EPYC Turin Dense Family) had 192 Zen 5c cores, where the highest full core CPU was the EPYC 9755 (from the EPYC Turin family), which had 128 Zen 5 cores.

The way they manage this while not hitting the power wall is by making the CPUs physically larger. While a consumer Ryzen 9 9950X3D (the fastest consumer CPU at the time of writing) easily fits into the palm of even a small hand, an EPYC Turin CPU is so large that it covers most of even a big hand.

5.6.2 Limitations

The Power Wall More and more transistors need more and more power, thus leading to power delivery and dissipation becoming an issue. To compute the power dissipation, use the formula $P_{diss} = P_{dyn} + P_{leak} + P_{short}$, where $P_{dyn} = CV^2f$ (with C the capacitance, V the supply voltage and f the processor frequency) is the dynamic power, P_{leak} the leakage power (see DDCA) and P_{short} the short circuit power while switching.

At some point the chip becomes almost impossible to cool. A great example of a CPU series that suffers from this is the Intel Rocket Lake CPUs. The Intel Core i9-14900K is notoriously hot-running, using almost 300 watts for a very small chip and thus runs very hot.

Thus, to further increase performance, chip designers are trying to make the hardware more efficient, which allows them to further boost performance with extra power headroom.

The Memory Wall Between 1985 and 2005, CPU performance has increased on average by 55% a year, whereas memory throughput has only increased by roughly 10% a year. Thus, performance has more and more become limited by memory performance rather than pure CPU performance and to this day is the largest overhead in most applications.

The ILP Wall While it is possible to improve single core performance using instruction-level parallelism, this has been thoroughly exhausted and is not a feasible way to significantly improve CPU performance.

Around 2003, all of these walls were hit simultaneously, as they hit a power wall and thus could not clock the processors any higher, the memory access times were the limiting factors and ILP was almost completely exhausted, as not enough parallel instructions existed in code.

Current trends are a reduction in clock frequency in favour of more parallelism in the hardware, e.g. by providing more cores, or better caching, branch prediction, etc.

5.6.3 Coherency and Consistency

Definition: Coherency The values in cache all match each other and the processors all see a coherent view of the memory

Definition: Consistency The order in which changes are seen by different processors is consistent

Most modern system's CPU cores are caches coherent, i.e. it behaves as if all cores access a single memory array. This leads to one big advantage: It is easy to program, however is hard to implement in hardware and memory is also slower as a result.

Memory consistency on the other hand is not standardized across manufacturers. We are asking questions like what happens if several processors read and write data, which value is read by each one of them? That question is not easy to answer and there is also more than one "correct" answer. The key though is to *have* an answer.

Definition: Program order is the order in which a program on a processor *appears* to issue reads and writes. This only refers to local reads and writes and even on a uniprocessor (in other words, a single core processor), it does not correspond to the order in which the CPU issues them

Definition: Visibility order is the order in which all reads and writes are seen by one or more processors. This refers to all operations on the machine and might not be the same for all processors. Each processor then reads the value written by the last write in visibility order.

Definition: Sequential consistency Operations from a processor appears to all others in **program order** and every processor's visibility order is the same *interleaving* of all the program orders. For that to work, every processor has to issue memory operations in program order, the RAM then has to totally order all operations and the operations have to be globally atomic.

You can imagine that to work as though the each processor issues a memory operation and the memory picking a (random) processor to process the request from, which it completes fully and then chooses another processor.

Advantages include being easily understandable for programmers, it being easier to write correct code for and it makes code more easily automatically analyzable. On the other hand, it is hard to make it fast, as we cannot reorder reads or writes (which can often speed up processing) and we cannot combine writes to the same cache line.

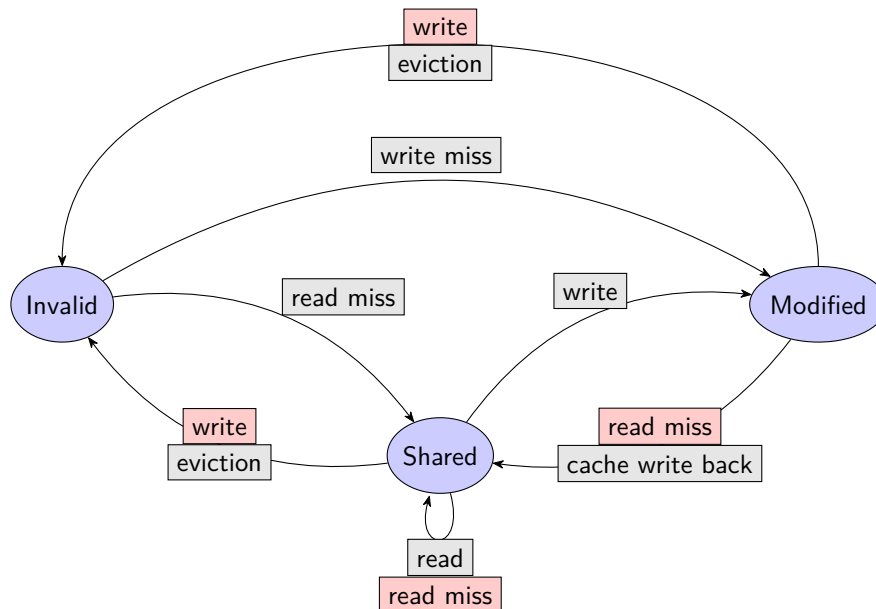
In Hardware it can be challenging to maintain proper sequential consistency, as multiple caches could hold invalid lines. This can however be worked around in hardware.

Snoopy caches are caches that "snoop" on reads and writes from other processors and thus, if a line that is valid in the local cache is written by another processor, the local line is invalidated. A write-through cache makes life a bit easier, but it can also work with a write-back cache, if cache lines can be marked as dirty (i.e. modified). It also requires a cache coherency protocol. A simple example is the MSI protocol, where a line can have three states (modified, shared, invalid). It basically forms a finite state machine that looks a bit like this:

5.6.4 MSI protocol

In MSI, a cache line may be in one of 3 states:

1. **Modified:** This is the only valid copy in any cache, newer than RAM. (Dirty)
2. **Shared:** This copy is coherent with RAM, other caches may have it too. (Clean)
3. **Invalid:** Block is unused or contains invalid data.



In these protocols communication happens through messages passed via the bus:

Code	Meaning
PrRd	Processor Read
PrWr	Processor Write
BudRd	Bus Read
BusRdX	Bus Read (Exclusive)

There are 3 main types of transitions happening:

1. **Local Read Miss:** Request read from State I, broadcasts BusRd. If any cache has this line set to M, this line is flushed to RAM. Transitions to S.
2. **Local Write Miss:** Request write from State I, broadcasts BusRdX. Other caches having this line in S or M must set to I (invalidate). Transitions to M
3. **Writing to Shared:** Request write from State S, broadcasts BusRdX. Other caches having this line in S must set to I. Transitions to M.

As nice as MSI is, as basically everything that is simple, it comes with issues, primarily here that it introduces unnecessary broadcasts. For instance, imagine a single-threaded program: Each time this program operates on a variable, broadcasts are sent over the bus even though no other core has this variable. MESI fixes this issue by introducing the Exclusive state.

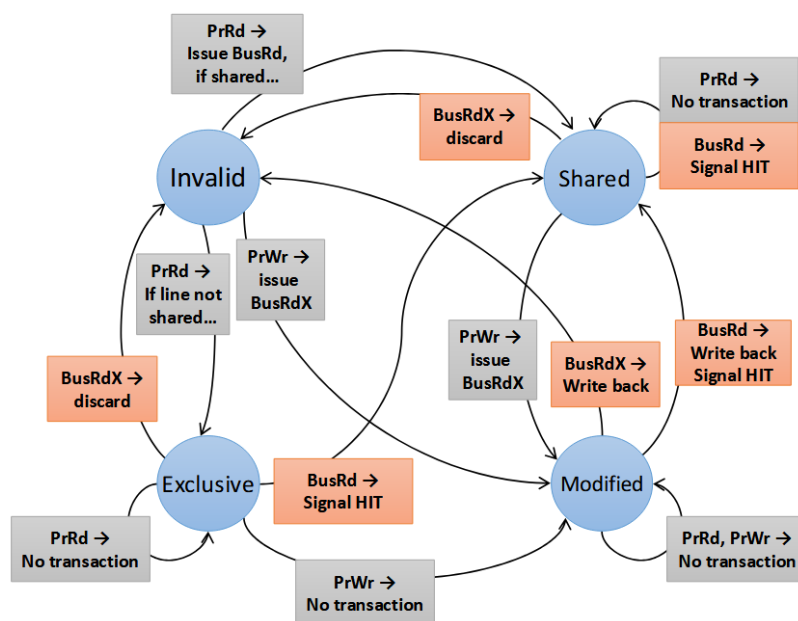
5.6.5 MSI extensions

MESI is an extension to the MSI protocol in which the processor gets to know that it is the only reader of a block. It has four states:

- **Modified:** This is the only copy, but it's modified.
- **Exclusive:** This is the only copy and it is not modified
- **Shared:** This might be one of several copies, all clean
- **Invalid**

When accessing cache, it signals a remote processor that it has hit the local cache. The cache can then load a block in either *shared* or *exclusive* states depending on whether or not the block is a HIT in the remote processor cache. The advantage is that writes to blocks in Exclusive state need *no* broadcasts.

This finite state machine is much more complex:⁶



Here, gray boxes are processor-initiated while orange boxes represent snoops.

MOESI AMD then added an owner state, in which the line can be modified, but there exist dirty copies in other caches. It has the benefit of being more quickly readable, by using the owner's cache. This of course is only beneficial if the latency to the remote cache is lower than to main memory, which in case of AMD CPUs, it is starting with the Zen 3 architecture and thus the Vermeer series of desktop CPUs (i.e. Ryzen 5000 series), as they are the first AMD CPUs with a unified L3 cache for each CCX, compared to unified cache for only four cores.

MESIF Intel added a forward state, in which cache requests are forwarded to the most recent cache line. Again, we only benefit from this if cache latency is lower than main memory latency and thus, Alder Lake (12000 series) and later benefit from the more than previous generations. (technically, also Intel 4004 in 1971 until the first Pentium, but they are hardly relevant today)

⁶This state transition diagram is from the SPCA lecture notes for HS25.

5.6.6 Relaxing Sequential Consistency

As we have outlined, sequential consistency may not be desirable when trying to build a high-performance system. We thus may want to relax sequential consistency. A primary reason to this is out-of-order execution giving a massive speed boost, as we do not have to wait for slow memory accesses to finish and can already compute what we have ready. Luckily, there are plenty of ways to work around this. For example, we can make later writes bypass earlier writes, later reads bypass earlier writes, break write atomicity (i.e. the order is not fixed anymore) or we cannot make any ordering guarantees at all.

x86 introduces specific instructions for synchronizing, for example `lfence` (load fence), `sfence` (store fence), `mfence` (memory fence) and others. It is typical for an x86-64 processor to implement relaxed sequential consistency and this is sometimes referred to as Total Store Ordering (TSO).

As a general rule, the weaker the consistency model is, the quicker it runs and the cheaper it is in hardware. However too weak a consistency model and some algorithms will simply stop working correctly.

Definition: **Barriers / Fences** are synonyms and are used to stop either the compiler or the CPU from reordering instructions or statements for order-critical operations.

Compiler barriers If we use gcc, we can use the following compiler intrinsic to stop it from reordering visible loads and stores:

```
__asm__ __volatile__ ("" ::: "memory");
```

The above intrinsic will also apply a memory barrier, which in terms of assembly code will use the `mfence` instruction. This instruction stops the CPU reordering past it, i.e. any instruction before the fence cannot happen after one behind the fence. However, any instructions past the fence are fair game and can be reordered (i.e. two instructions behind the fence can be reordered).

If we only need this for stores or loads, we can use `lfence` or `sfence`, respectively.

5.6.7 Multicore synchronization

There are two main ways to synchronize, which are:

1. **Atomic operations** such as TAS, CAS, etc. It does still have ordering constraints specified in the memory model
2. **Interprocessor interrupts** (IPIs) This invokes the interrupt handler on remote CPU, but is VERY slow (500+ cycles) and thus often avoided, except in the OS

TAS (Test-and-Set) We can only set to the memory location using TAS if said location is 0. It can thus be used for a mutex, with a simple spinlock, which is simple to implement and often the fastest if the lock isn't held for long. Since we most commonly do not read a value of 0 in the lock memory location, we can use a TATAS (Test And Test-and-Set) lock to reduce the performance overhead.

File: code-examples/03_hw/01_tas.c

```

1 void acquire( int *lock ) {
2     while ( TAS( lock ) == 1 );
3 }
4
5 void acquire_tatas( int *lock ) {
6     do {
7         while ( *lock == 1 );
8     } while ( TAS( lock ) == 1 );
9 }
10
11 void release( int *lock ) {
12     *lock = 0;
13 }

```

A word of caution: Do not use TAS to check if a value has changed outside a lock. It will most likely not work in C and almost certainly not in Java or any higher level languages. In systems code however, it works well, as TAS is a hardware instruction. The processor won't reorder instructions past a TAS and the compiler also should not. They are told with memory fences and the use of the volatile keyword.

However, be aware that the code depends on the processor architecture and its memory consistency model.

CAS (Compare-and-Swap) first loads the current value, then checks if it is the same as the "old" value and only if it is sets it to the "new" value. The typical function signature of a CAS function is CAS(location, old, new) and it typically returns the old value.

In comparison to TAS, CAS is commonly used for lock-free programming, whereas TAS is almost an integral part to programming with locks. As already covered in parallel programming, we can use CAS to read a value, do a computation with it, then commit it back, only if the value was unchanged, or else restart our computation.

Another option is to use a ref counter, which by default is set to 1 on an "untouched" data structure. The global pointer points to the current memory location and should only be modified using CAS.

When a reader reads the data, it is incremented (ref1 = 2). Then, if a writer comes along, it reads and increments the ref counter (ref1 = 3), then copies the data and decrements the ref counter of the original data (ref1 = 2, ref2 = 1 by convention). It then uses CAS to swap the global pointer to the new data and decrements the ref count (ref1 = 1). When the reader now finishes, it decrements the ref counter of the original data again (ref1 = 0) and the data is deleted.

ABA problem occurs if one process (or processor) does not see an intermediate change B because a second change wrote back the same value, A , as it was before the intermittent change. This problem is caused by the fact that CAS doesn't tell you if the variable had been overridden or not, but only if the value is different.

An easy way to solve this is to simply use more bits and make e.g. the upper few bits an always incrementing counter.

DCAS (Double CAS) is an alternative option for CAS, where you compare two memory locations and you only update if both match. However, this proved to not useful over the normal CAS and was subsequently dropped by most ISAs, as everything you could achieve with DCAS could also be achieved really fast with just CAS.

In x86 There are many options, the easiest one being to prefix an instruction with `lock`, which locks the data buses. If you only want to use TAS, you can use the `xchg` instruction. The `lock xadd` instruction executes an atomic fetch and add. The `CMPXCHG`, `CMPXCHG8B` and `CMPXCHG16B` implement 32 bit, 64 bit and 128 bit CAS, respectively, where the 128 bit version is exclusive to x86-64

5.6.8 Symmetric Multiprocessing

SMP allows multiple cores to access the same memory. It does still allow each core to have a separate cache, which is a de-facto requirement for it to work, as otherwise the memory becomes an even more serious bottleneck.

However, even with all the cache optimizations, the memory is still the bottleneck in SMP. The MOESI protocol can alleviate some of the slowness by enabling reads to be serviced by other caches, but that can again slow that cache down. Additionally, memory accesses can stall the current processor, as well as other processors while accessing data.

So, to reduce idle times, we would like to issue instructions to the Functional Units (FUs), but Instruction Level Parallelism (ILP) is limited due to data dependencies.

This is where SMT (Simultaneous Multithreading) comes into play. As of 2025, only AMD's consumer hardware features SMT, Intel has abandoned it with Arrow Lake (Core Ultra series). On most platforms, especially Intel platforms up to there, SMT was often referred to as Hyperthreading. It uses the fact that there are likely operations that can be served by cache in other threads and we can keep the FUs busy.

There are three main ways of achieving SMT:

- **Thread IDs** There are multiple independent instruction streams with each having their own thread IDs for later reassignment.
- **Fine-grained MT** On every instruction dispatch, a thread to dispatch from is picked
- **Coarse-grained MT** When a memory stall is encountered, switch to another thread

Of note is that CPU cores with multithreading appear to the OS as multiple cores, and in the common case of what AMD and Intel (up to Arrow Lake) do (did), they would appear as two cores, which is also where the "Thread" count comes from on CPU spec sheets.

As nice and good of an idea as this sounds, it isn't necessarily faster, as threads might be competing for cache, which is why Intel has abandoned SMT altogether and they are seeing quite competitive performance. But then again, it is very cheap to implement in terms of extra transistors used, but the performance gain is in the 10-20% range typically, if even.

Finally, the performance gain strongly depends on the workload. Scientific computing rarely benefits from it, as the compute is often the limiting factor, however, web servers for example strongly benefit from it, as they are commonly severely memory constrained.

5.6.9 Non-Uniform Memory Access

In a typical early SMP architecture, more cores were provided, but not necessarily more cache, or even completely shared cache between the cores.

This is where the NUMA concept comes in, restricting each core, or group of cores to a subset of the memory. This is specifically advantageous in large data centers, where there might be hundreds of CPUs with Terabytes of memory.

Initially, NUMA was constricted to data center uses, but with the Zen microarchitecture, AMD has brought the concept to the consumer market, where larger and more performant CPUs with larger core counts are split into multiple CCXs (Core Complexes), each with cache per core and cache per CCX. Then, there is the Infinity Fabric, which is a CCX interconnect, allowing the cores from the other CCX(s) to access the data in the cache of the current CCX.

However, that is not a full NUMA implementation, as they still share one memory bus. In a multi CPU deployment (i.e. with multiple CPU sockets containing a CPU), they often have their own memory and memory controller, as well as an interconnect, allowing them to communicate with the other sockets to access their data in the memory. In such a deployment, a CPU socket with its own memory is referred to as a NUMA Node.

Let's look at an example with two AMD EPYC 7742 CPUs, each having 64 cores. Each of these CPUs features 128 PCIe Gen 4 lanes, 64 of which could be used for the Infinity Fabric CPU to CPU interconnect (which uses the PCIe Interface). It is comparatively slow, maxing out approximately at a theoretical 128 GB/s in one direction. Compare that to the roughly 205 GB/s bandwidth of DDR4 3200 MT/s

Cache coherence We can no longer snoop on the bus, as it is no bus anymore. A solution is to emulate a bus, which enables something akin to snooping, but without the shared bus. Each node sends a message to all other nodes and waits for a reply from all the other nodes before proceeding. This is the way AMD's Infinity Fabric works (or used to work)

Another option to circumvent the issue is to use a Cache Directory, where we store the data, the node ID of which one is the one it originated from, plus one bit per node indicating if the line is present in that node. It is primarily efficient if the lines are not widely shared or if there are lots of NUMA nodes.

5.6.10 Performance and Optimization

Depending on the interconnect layout, accessing data from certain nodes is *much* slower than from others. This especially happens if there is no direct path to the other node from the current node and the access needs to pass through other nodes.

It is also of utmost importance to keep threads on the same CPU to improve locality, which improves cache coverage and try to pack data into a single cache line. Otherwise, the cache "ping-pongs".

MCS locks There are many ways to improve performance on multiprocessor systems. The MCS lock is one of the best locking system for multiprocessors. It tries to solve the issue that if a cache line contains a lock, it is continuously invalidated and that dominates the intercom traffic. The solution to the problem is that a processor enqueues itself on a list of waiting processors and then spins on its *own entry* in the list.

File: code-examples/03_hw/02_mcs.c

```

1  #include <stdlib.h>
2
3  struct qnode {
4      struct qnode *next;
5      int locked;
6  };
7  typedef struct qnode *lock_t;
8
9  void acquire( lock_t *lock, struct qnode *local ) {
10     local->next = NULL;
11     struct qnode *prev = XCHG( lock, local );
12     if ( prev ) { // queue was non-empty
13         local->locked = 1;
14         prev->next = local;
15         while ( local->locked )
16             ; // spin
17     }
18 }
19
20 void release( lock_t *lock, struct qnode *local ) {
21     if ( local->next == NULL ) {
22         if ( CAS( lock, local, NULL ) )
23             return;
24         while ( local->next == NULL )
25             ; // spin
26     }
27     local->next->locked = 0;
28 }

```

5.7 Devices

From a programmer's perspective a Device can be seen as:

- Hardware assessible via software
- Hardware occupying some bus location
- Hardware mapping to some set of registers
- Source of interrupts
- Source of direct memory transfers

5.7.1 Device Registers

Sometimes devices expose register: The CPU can load from these to obtain e.g. status info or input data. The CPU can store to device registers to e.g. set device state or write output.

Device registers can be addressed in 2 different ways:

1. **Memory Mapped:** Registers *appear* as memory locations, access via `movX`
2. **IO Instructions:** Special ISA instructions to work with devices

It's important to note: despite *appearing* as memory, device registers behave differently: State may change without CPU manipulation and writes may trigger actions. The specific way this interaction works is device-specific.

Example A very simple device driver in C may look like this:

File: `code-examples/03_hw/00_driver.c`

```

1  #define UART_BASE 0x3f8                // Registers specified by
2  #define UART_THR (UART_BASE + 0)      // the device
3  #define UART_RBR (UART_BASE + 0)
4  #define UART_LSR (UART_BASE + 5)
5
6  void serial_putc(char c) {
7      while( (inb(UART_LSR) & 0x20) == 0); // Wait until FIFO can hold more chars
8      outb(UART_THR, c);                  // Write character to FIFO
9  }
10
11 char serial_getc() {
12     while( (inb(UART_LSR) & 0x01) == 0); // Wait until there is a char to read
13     return inb(UART_RBR);               // Read from the receive FIFO
14 }
```

Of course, a proper driver would also include error handling, initialization and wouldn't spin to wait. To avoid waiting, interrupts are usually used.

Caches For Device Registers, the cache must be bypassed. Memory mapped IO causes a lot of issues for caching:

1. Reads can't be cached (Value may change independent of CPU)
2. Write-back doesn't work exactly (Device controls when the write happens)
3. Reads and writes can't be combined into 1 cache-line

5.7.2 Direct Memory Access

Direct Memory Access (DMA) requires a dedicated DMA controller, which is generally built-in nowadays. DMA allows bypassing the CPU entirely: Data is transferred directly between IO device and memory. This is especially useful for large transfers, but not small transfers due to the induced overhead.

The key advantage is that data transfer and processing are decoupled.

The CPU never needs to deal with copying between device and memory, and the CPU cache is never polluted.

The key disadvantage is that Memory is inconsistent with the CPU cache. This is addressed in various ways:

1. CPU may mark DMA buffers as *non-cacheable*
2. Cache can *snoop* DMA bus transactions (Doesn't scale well, only for small systems)
3. OS can explicitly flush/invalidate cache regions. (Usually done by a device driver)

Another issue is that DMA addresses are *Physical*. The OS (via device drivers) must *manually* translate these to virtual addresses. Some systems also contain a dedicated component called IOMMU that deals with this.

5.7.3 Device Drivers

Device drivers are programs used by the OS to communicate with devices. In a nutshell, the driver is the only program that *directly* interacts with the device, any other program talks to the device *through* the driver (which ideally abstracts away a lot of the process).

Intuitively, both the driver and device can be thought of as state machines, which affect each other.

Definition: Descriptor Ring is a type of buffer commonly used to interact with devices. The datastructure is a looped queue (ring): Device reads from the head, OS writes at the tail. The space beyond is then "owned" by the device/OS.

This can either be implemented as contiguous memory or using pointers (which is mainly what is done in practice, for flexibility). Overruns (Device has no buffers for received packets) and Underruns (CPU has read all received packets) are usually handled sensibly: i.e. the CPU waits for an interrupt or the device will simply wait.

Parallel Programming: These are producer/consumer queues! But these use messages instead of mutexes and monitors.