# Introduction to HPC

Lecture 2

Jakub Gałecki

# CPU architecture – a primer

CPU = Central Processing Unit

It is the brain of the computer

But that's a bit vague…

# But that was too specific…

Let's start by placing things in context. The (modern) computer consists of:

- **CPU(s)**
- **RAM**
- **GPU(s)** (last 2 lectures)
- The hard drive (whether HDD or SSD is irrelevant to us)
- I/O devices
- …

Fundamentally, the CPU performs the following tasks:

- Fetches instructions
- Decodes instructions
- Executes instructions

How does it know where to fetch the instructions from: the program counter.

Examples of instructions:

- arithmetic operation
- read/write from/to memory
- **conditional jump**

Instructions usually operate on operands (arguments)
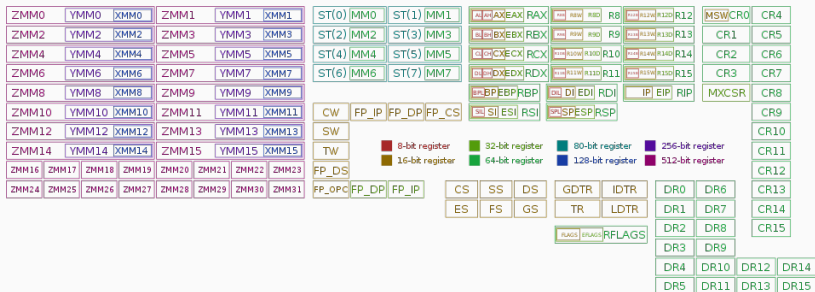
Small (e.g. 64 bit) volatile memory units

Most instructions involve data stored in registers

Registers have zero latency to access

Examples:

- General purpose
- RFLAGS
- Control
- Debug
- Vector (spoiler alert)

```
.LC0:
        .string "Hello, World!\n"
main:
        push    rbp
        mov     rbp, rsp
        mov     edi, OFFSET FLAT:.LC0
        call    puts
        mov     eax, 0
        pop     rbp
        ret
```

https://godbolt.org/z/or5Tz3cEb
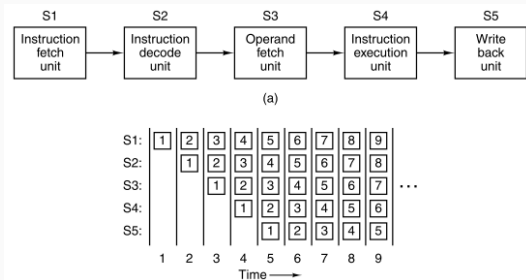
Transistor reaction speed is not instantaneous.

- Gate delay: $d$
- Desired clock rate: $f$
- Theoretical max gate chain length: $1/df$

If we want fast clock frequencies, we have a hard, physical limit on the complexity of our circuit.

The solution: pipelining

We can break the instructions down into stages and execute 1 stage per cycle. Different stages of subsequent instructions are executed concurrently!

Simplified example, 5 stage pipeline:



Problem: what happens when instruction $n + 1$ depends on the result of instruction $n$?

Pipelining introduces potential delays when subsequent operations depend on one another

- Structural hazard – resource conflict
- Data hazard – logical dependency between instructions
    - Read-after-write
    - Write-after-read
    - Write-after-write
- Control hazard – control flow depends on result of previous instruction

We should keep these in mind when programming, although the hardware and compiler do most of the heavy lifting.

What kind of hazard is this?

| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Instructions | ADD R8, R5, R5 | IF | ID | EX | MEM | WR | | | | | | | | | |
| | ADD R2, R5, R8 | | IF | Idle | | | ID | EX | MEM | WR | | | | | |
| | SUB R3, R8, R4 | | | IF | Idle | | | ID | EX | MEM | WR | | | | |
| | ADD R2, R2, R3 | | | | | | | IF | Idle | ID | EX | MEM | WR | | |

- Structural hazards: get better CPU (sorry)

- Data hazards: compiler optimization, out-of-order execution, register renaming, inline assembly if we're feeling dangerous

- Control hazards: branch prediction, write better code (stay tuned)

```asm
.LC0:
        .string "Hello, World!"
.LC1:
        .string "So many arguments :o"
main:
        sub     rsp, 8
        cmp     edi, 1
        jle     .L6
        mov     edi, OFFSET FLAT:.LC1
        call    puts
.L3:
        xor     eax, eax
        add     rsp, 8
        ret
.L6:
        mov     edi, OFFSET FLAT:.LC0
        call    puts
        jmp     .L3
```

https://godbolt.org/z/5sWsYcvTd

The problem with branching: the CPU doesn't even know which instruction to fetch until some previous instruction executes

Control hazard == "Data hazard on steroids"

The solution: take a guess and see what happens

- Correct guess: no stall, no performance penalty
- Incorrect guess: pipeline flush, undo changes – expensive

We need to try to be predictable.

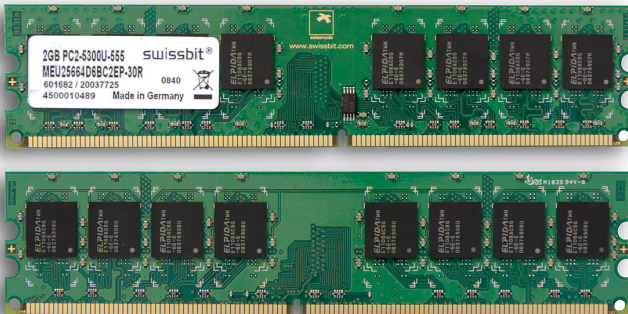Complete talk on the subject: https://youtu.be/g-WPhYREFjk

# Accessing memory

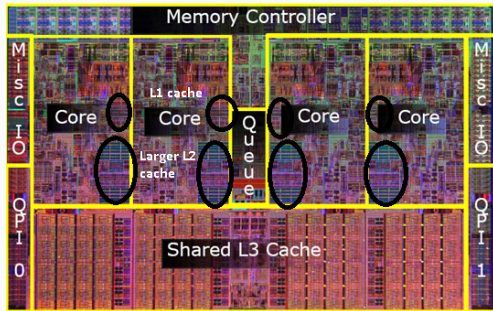DRAM == Dynamic Random Access Memory

Very large – up to hundreds of GB

Very slow to access – hundreds of cycles

Cache == fast, on-die memory

Usually several levels, nowadays: L1I, L1D, L2, L3

Smaller →faster

This is, of course, dependent on the specific hardware, but we can take a look at some reference values:

| Memory type | Size | Latency [cycles] | Bandwidth |
|---|---|---|---|
| Register | ~3KB* | 0 | - |
| L1 Cache | 32 KB | 4 | 256 GB/s |
| L2 Cache | 256 KB | 10-25 | 256 GB/s |
| L3 Cache | 8 MB | ~40 | 128 GB/s |
| Main memory | $\gg$1GB | 200+ | 17 GB/s |

The cache does not operate on individual bytes, but rather on sets of bytes, called **cachelines**.

The size of a cacheline on modern CPUs is 64B.

This has consequences:

- Aligning data to cache can increase performance
- Accessing neighboring data is faster
- Potential pitfall for concurrent programs (false sharing)

How many different slots in the cache are available for a given cacheline?

- Direct mapped – only 1 slot available
- Fully associative – any cacheline can go into any slot
- Set-associative – a cacheline can go into one of N slots

Trade-off between flexibility and power-efficiency

Modern CPUs: N-way associative, $6 \leq N \leq 16$

There is no instruction for "write N bytes from memory to LX cache"*

We have to structure our data access so that it is naturally cache-friendly

Spatial locality:

- Subsequent addresses are likely to be on the same cacheline
- The CPU can detect access patterns and prefetch our data

Temporal locality:

- Least recently used cacheline gets evicted first
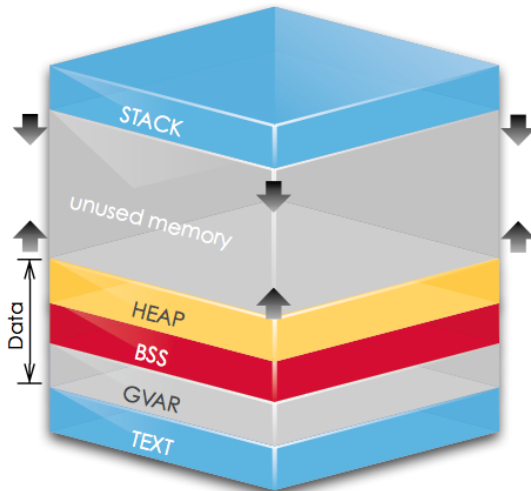- Data which was recently accessed is likely still in cache

Data is ultimately represented by electrons residing in the DRAM die – *physical* address

Our program references memory via *virtual* addresses

To de-conflict different processes, the OS *translates* virtual addresses to physical addresses

The CPU has special hardware which helps with translation

For improved efficiency, memory is divided into 4kB* *pages*

TLB = Translation Lookaside Buffer

Cache for the page translation process

TLB size: 1536 pages

TLB hit time: $\leq$1 cycle

TLB miss penalty: 10-100 cycles

*Memory thrashing* for large working sets with random memory access

Usually not an issue

We say address *i* is aligned to *a* (or has alignment *a*) iff

$$i \mod a = 0$$

where *a* must be a power of 2. For example:

- `0xa0` is aligned to 16
- `0x0777b2` is aligned to 2

CPUs are much better at accessing data which is aligned to its natural alignment, i.e., a multiple of its size.

For usual cases, this is handled by the compiler with padding: `https://godbolt.org/z/39aWbGoKW`.

We can use `alignas` or aligned allocation to override the defaults. We will soon see why this may be desired.
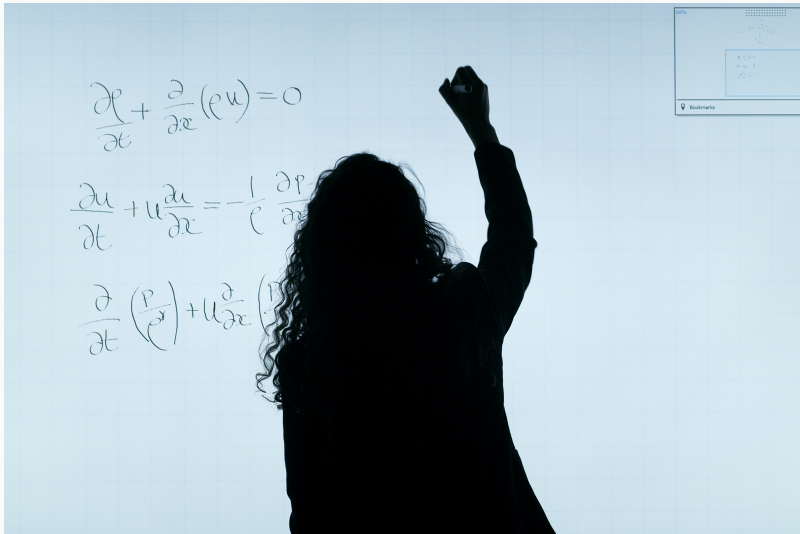
Author: Kazushige Goto (early 2000's)

Matrix-matrix multiply algorithm explicitly catering to the 3 level cache memory hierarchy

Slice & dice approach

General structure: simple, no CS PhD required

Micro-kernel: detailed knowledge of the CPU architecture is required

Fantastic explanation: `https://youtu.be/07SMaudtH6k`

- CPU architecture 101

- Assembly 101

- CPUs are pipelined

- Avoid unpredictable branches

- Cache is king

→ Want performance? Know your hardware!

→ The speed of feeding the data to the CPU is equaly as important as the speed of processing the data

→ Break down the problem, optimize the kernel