# General introduction: GPUs and the realm of parallel architectures

GPU Computing Training
August 17-19th 2015

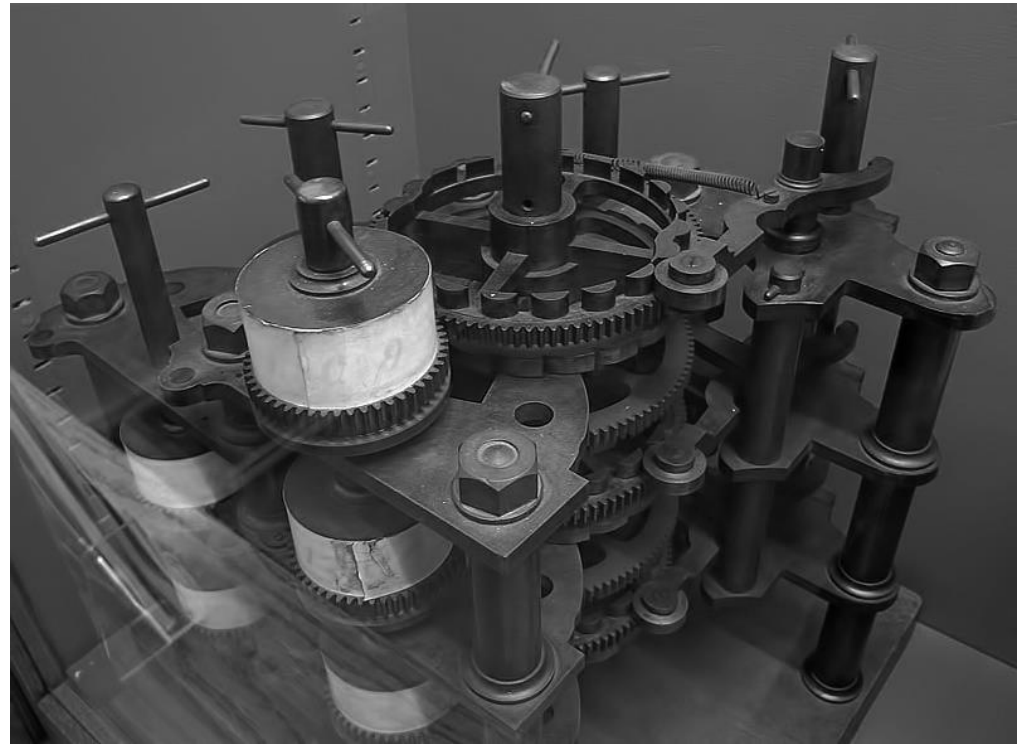# Jan Lemeire (jan.lemeire@vub.ac.be)

- Graduated as Engineer in **1994** at VUB
- Worked for **4 years** for 2 IT-consultancy companies
- **2000-2007**: PhD at the VUB while teaching as assistant
  - Subject: *probabilistic models for the performance analysis of parallel programs*
- **Since 2008**: postdoc en parttime professor at VUB, department of electronics and informatics (ETRO)
  - Teaching 'Informatics' for first-year bachelors; 'parallel systems' and 'advanced computer architecture' to masters
- **Since 2012**: also teaching for engineers industrial sciences ('industrial engineers')
- Projects, papers, phd students in *parallel processing* (performance analysis, GPU computing) & *data mining/machine learning* (probabilistic models, causality, learning algorithms)
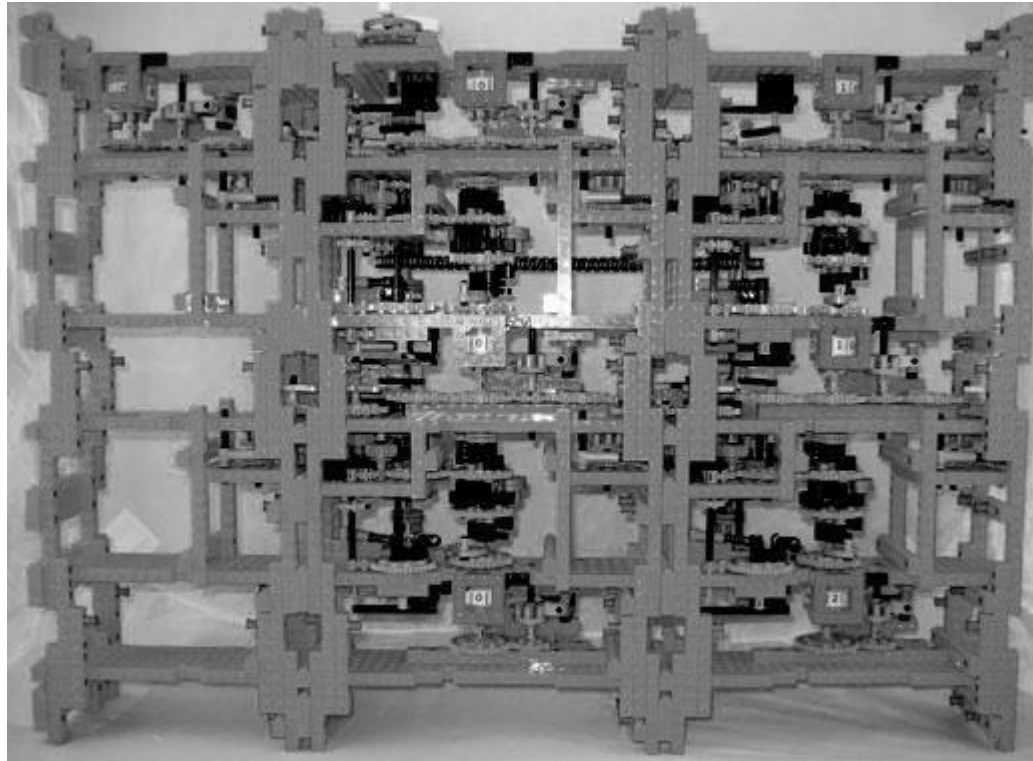
- http://parallel.vub.ac.be

# A bit of History

# The first computer, mechanical



**Charles Babbage
1791-1871**

# Babbage Difference Engine made with LEGO



- [http://acarol.woz.org/](http://acarol.woz.org/)

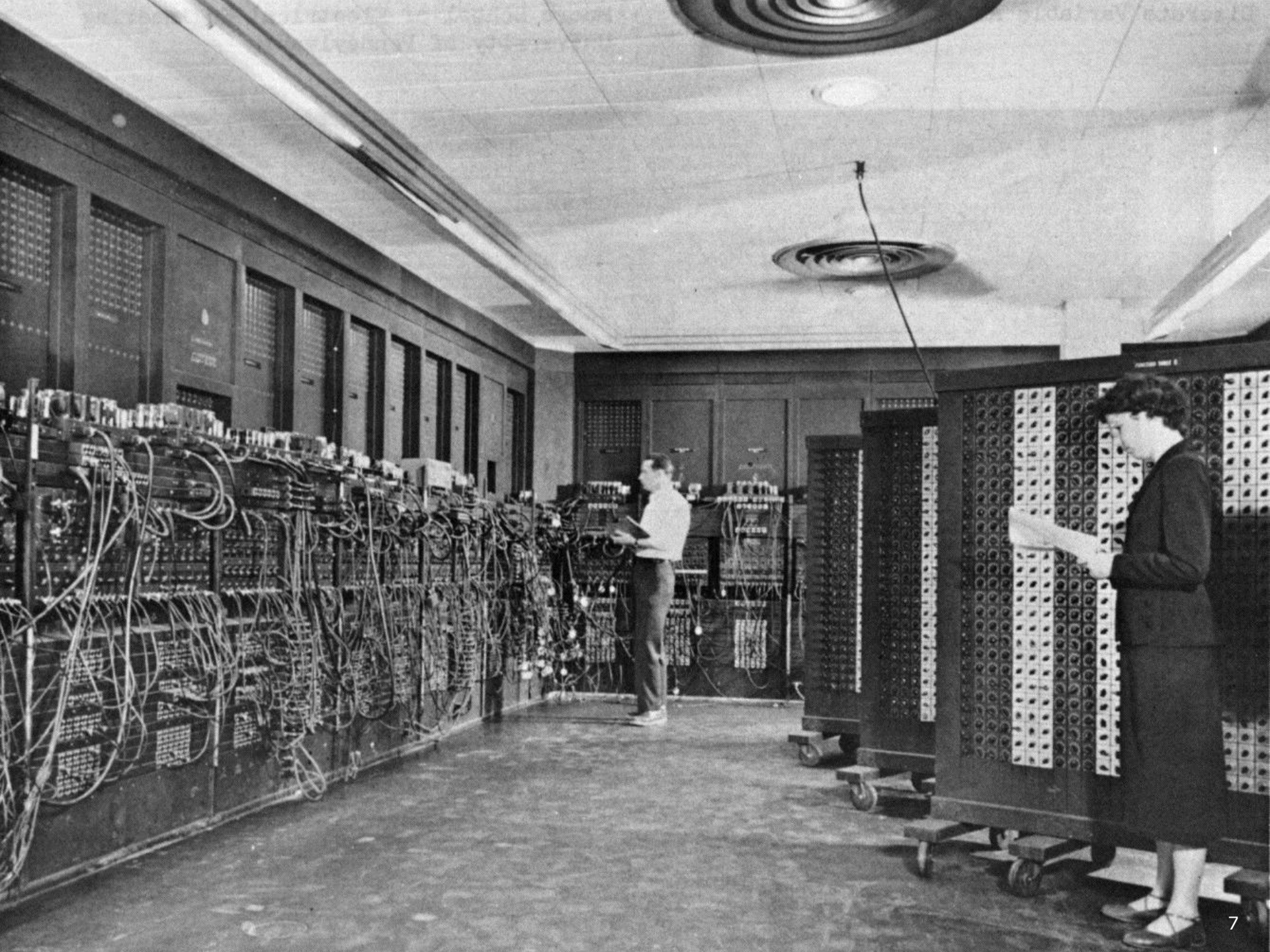This machine can evaluate polynomials of the form $Ax^2 + Bx + C$ for x=0, 1, 2, …n with 3 digit results.

# The first informatician
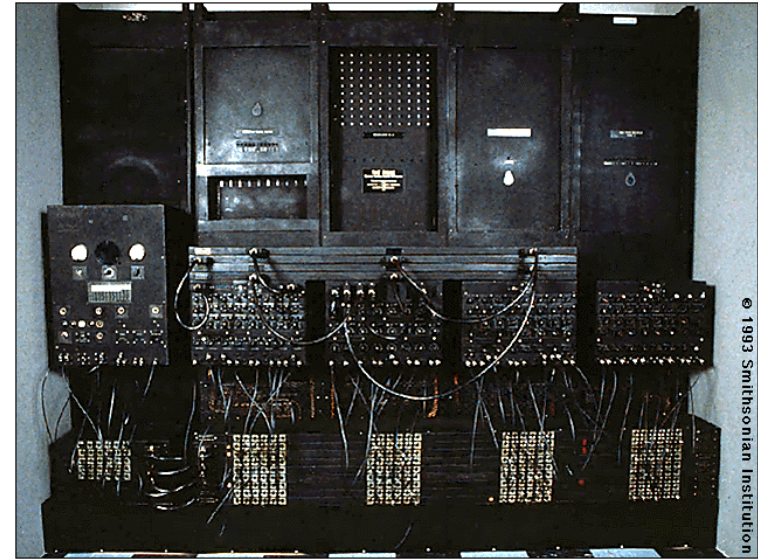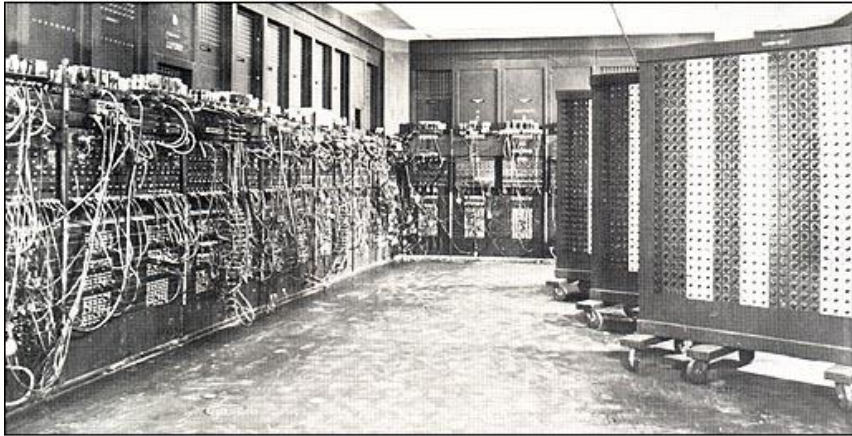
Ada Lovelace
1815 – 1852

- Describes what software is
  - She brings the insight that a computer goes beyond plain calculations
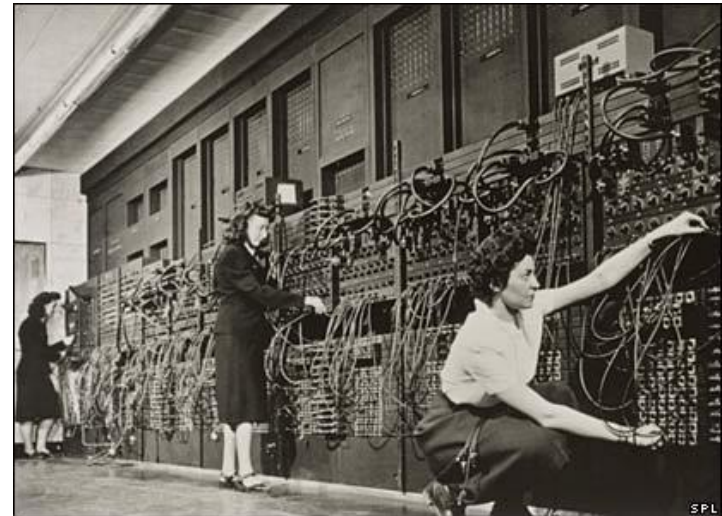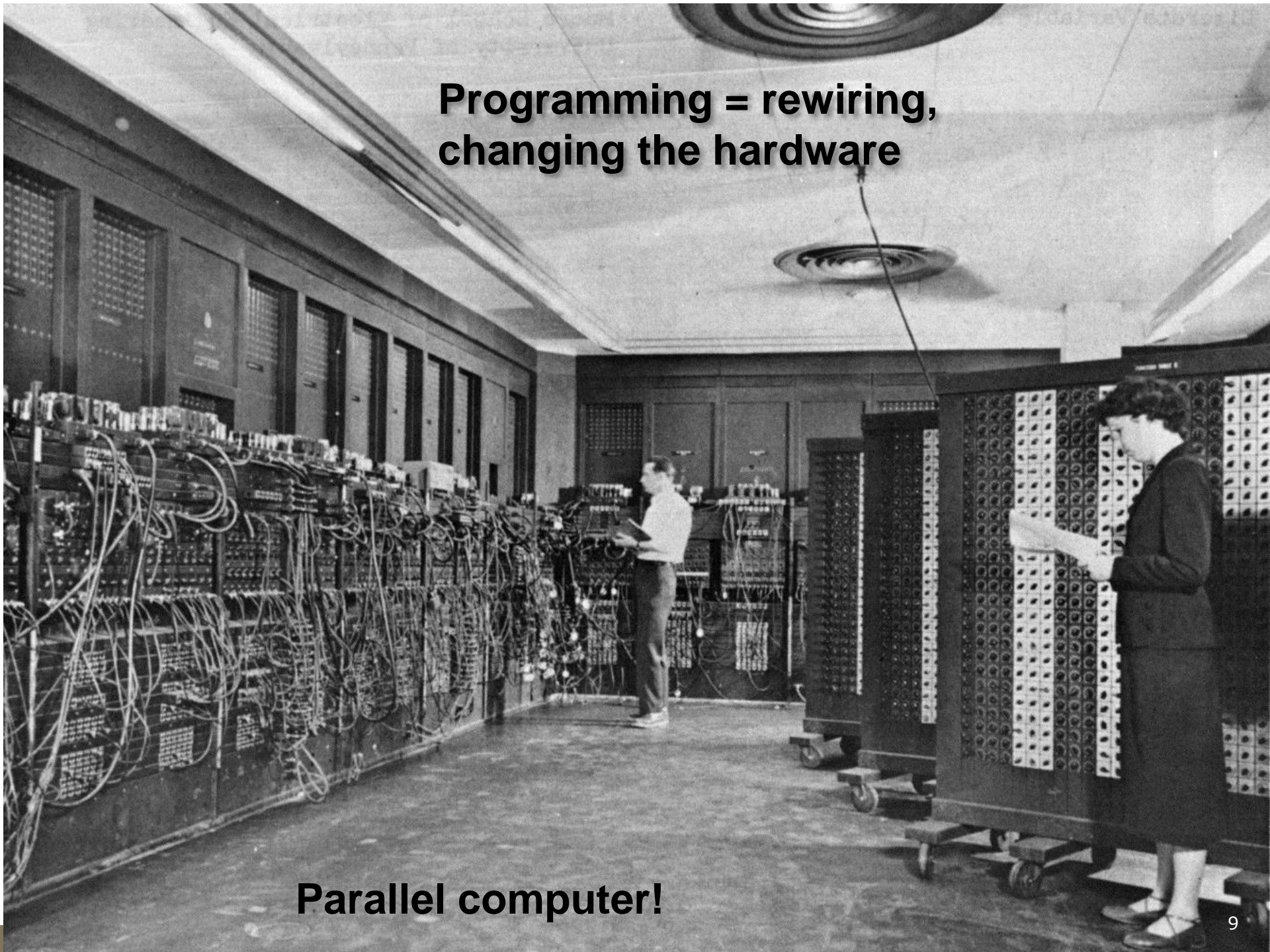- She writes the first algorithm/program

# ENIAC

First computer: WWII

John Mauchly and John Eckert, 1945

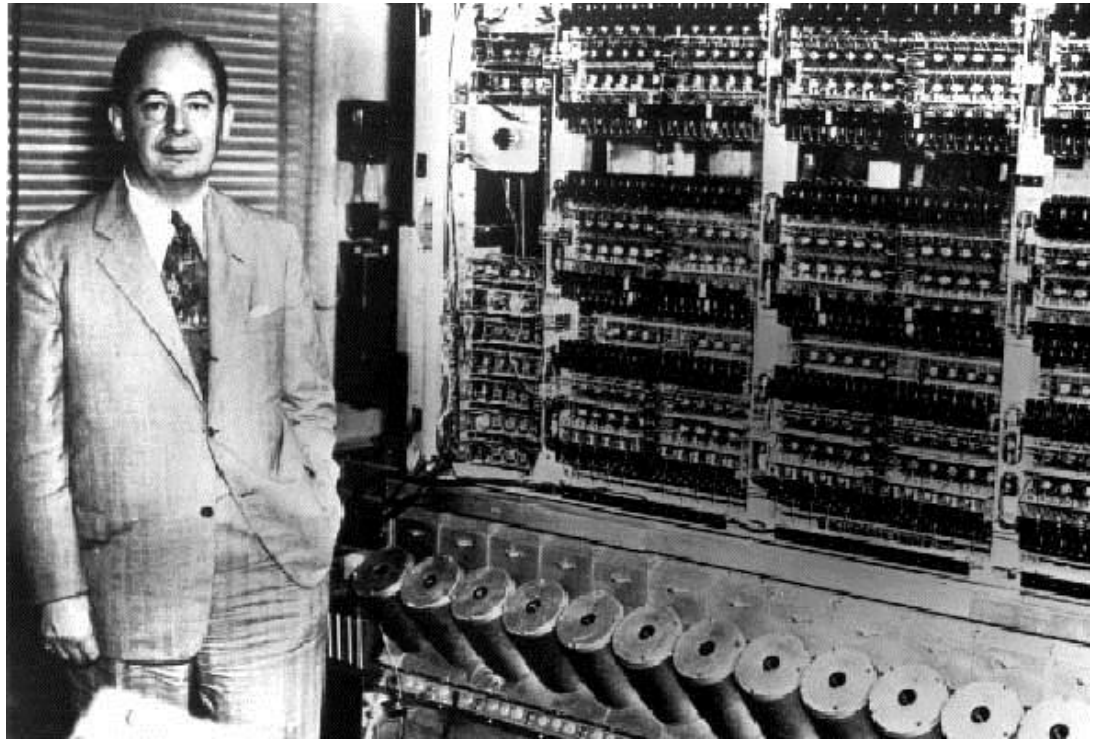**Programming = rewiring, changing the hardware**
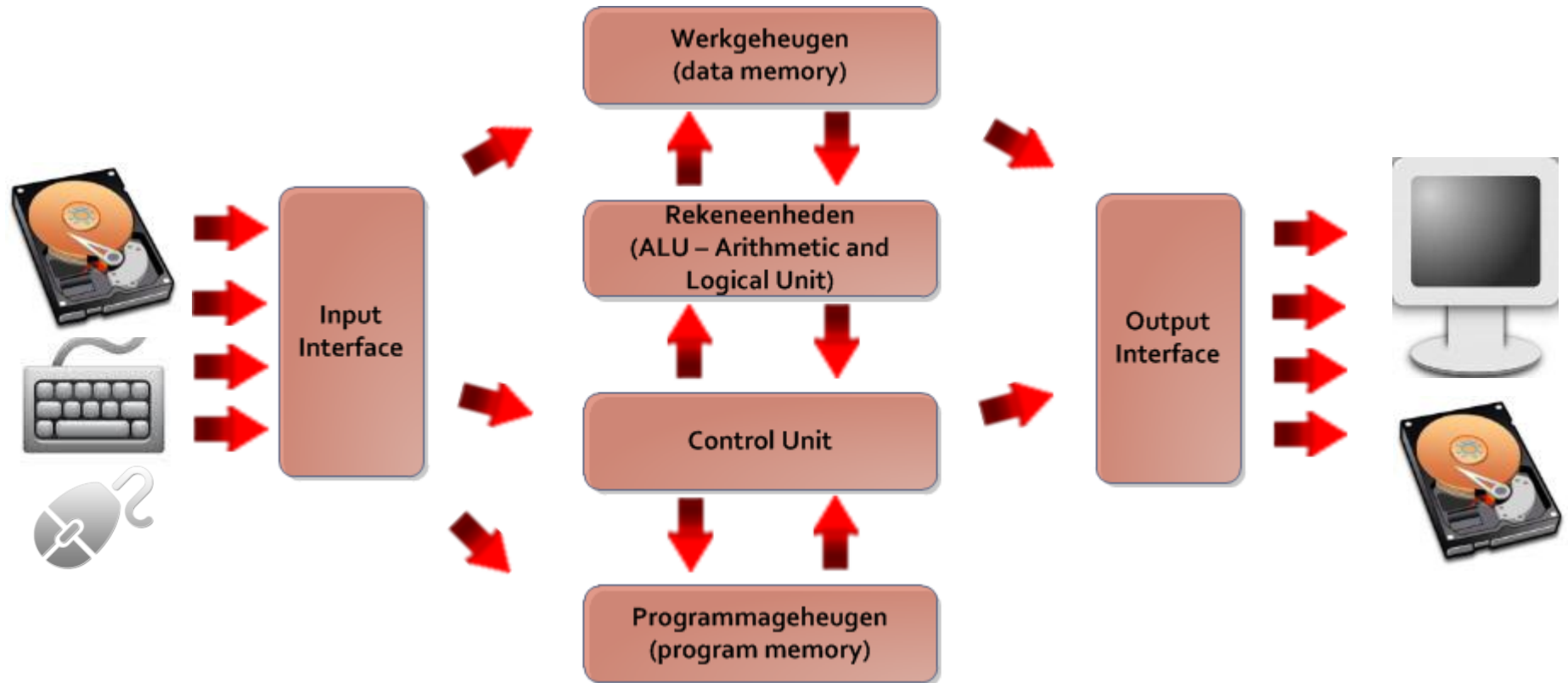
**Parallel computer!**

9

# Von Neumann rethinks the computer



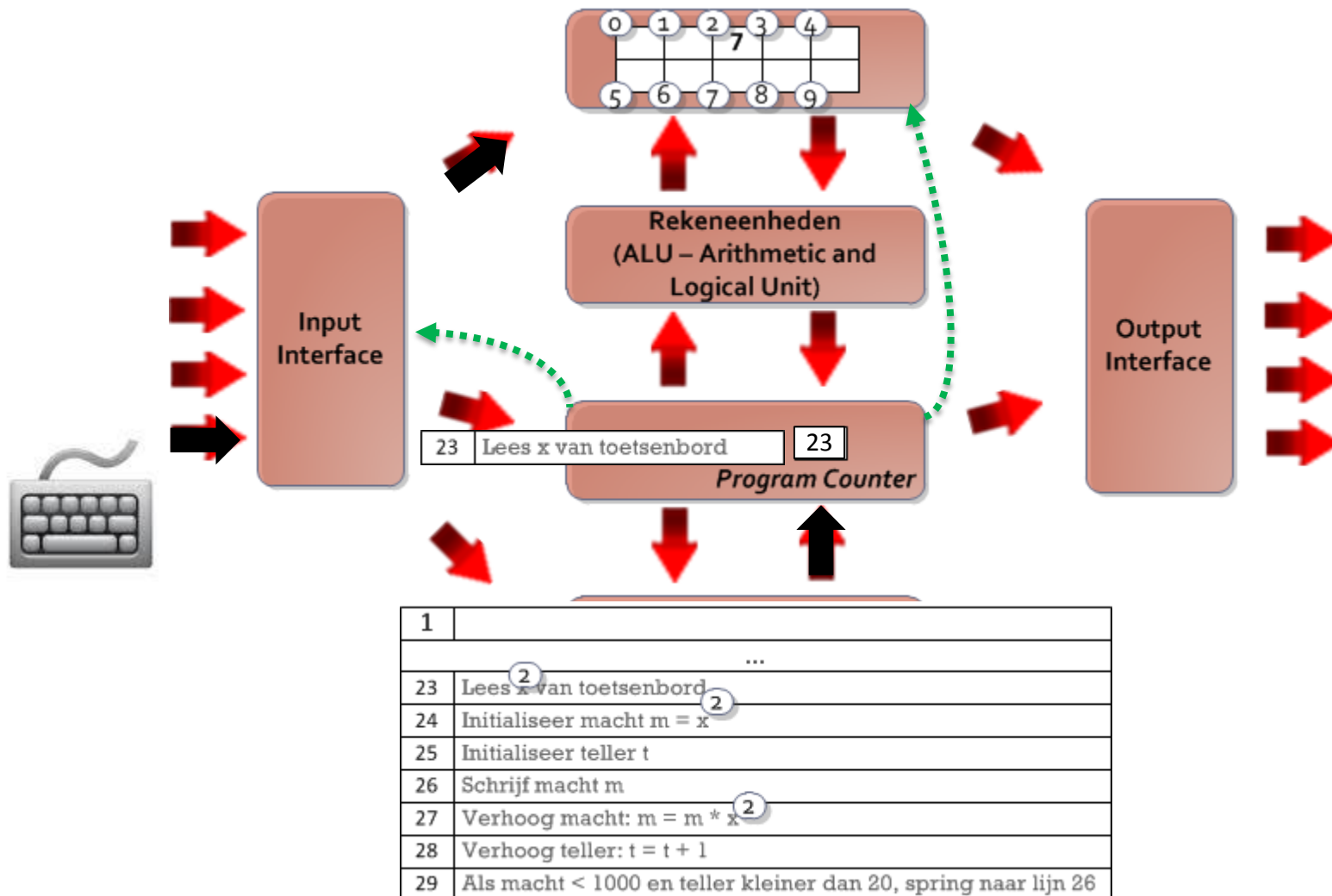**John Von Neumann**

# The Von Neumann-architecture

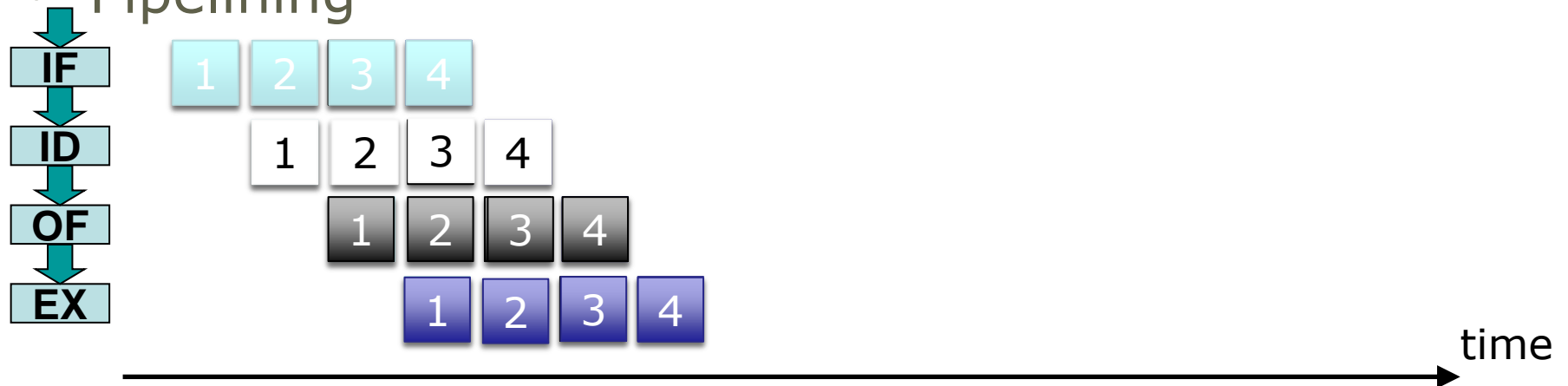# Execute program step-by-step

# For acceleration: **pipeline**

- Long operations

| 1 | 2 | 3 | 4 |

- Combination of short operations

1 1 1 1 2 2 2 2 3 3 3 3 4 4 4 4

- Pipelining

**IF** 1 2 3 4

**ID** 1 2 3 4
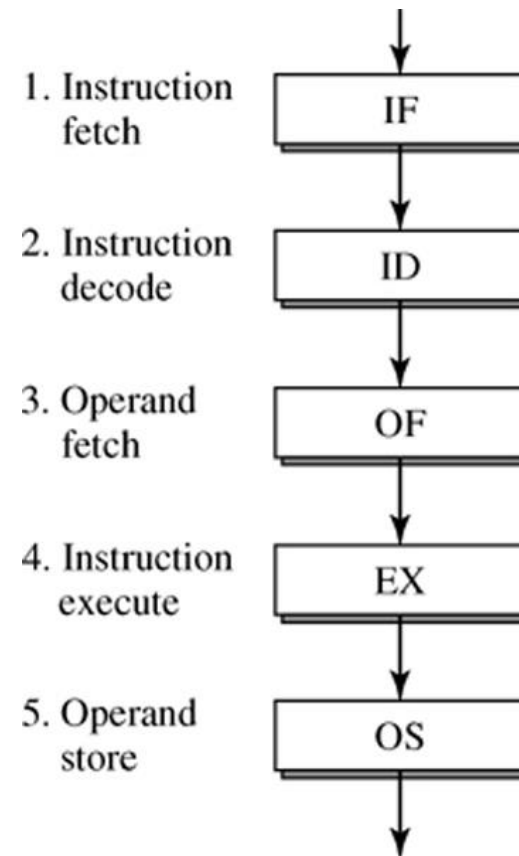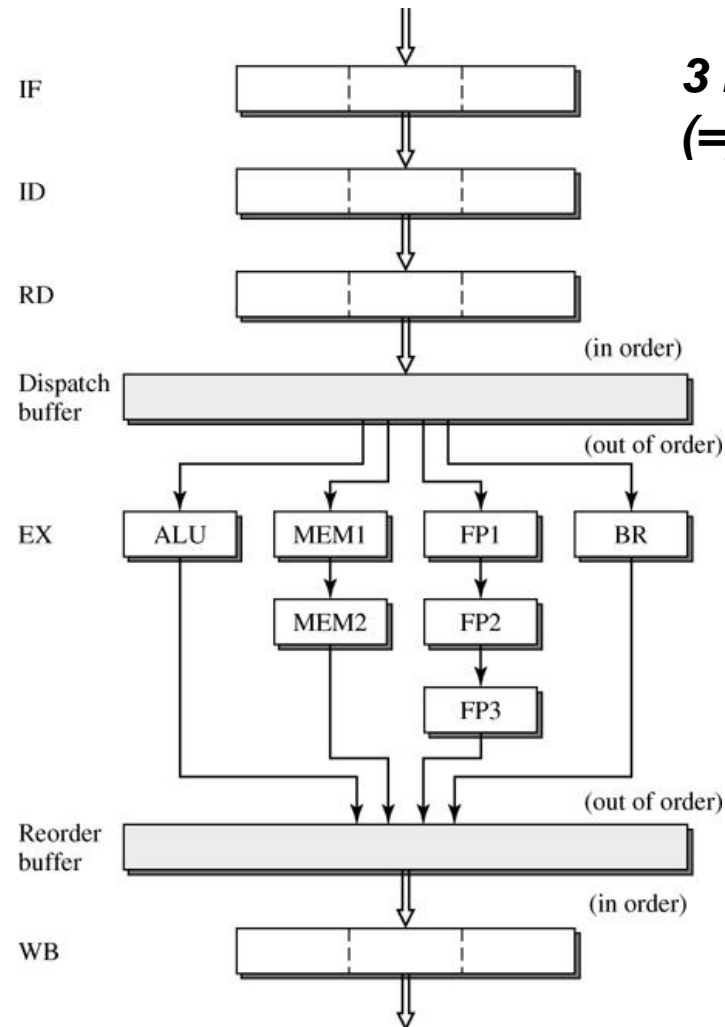
**OF** 1 2 3 4

**EX** 1 2 3 4

time

# Pipeline Design

- Typically five tasks in instruction execution
  - IF: instruction fetch
  - ID: instruction decode
  - OF: operand fetch
  - EX: instruction execution
  - OS: operand store,
    often called write-back WB



1. Instruction fetch — IF
2. Instruction decode — ID
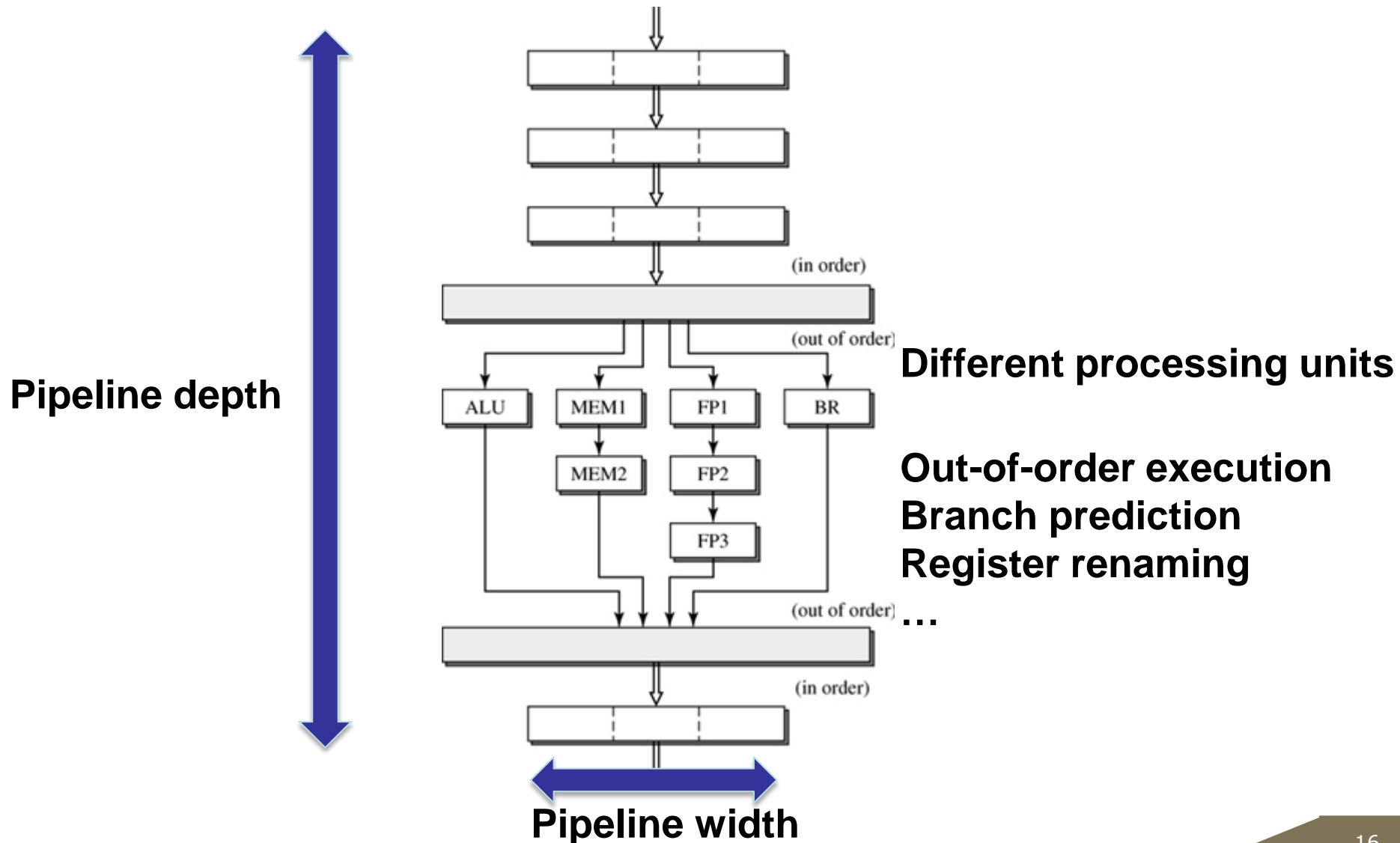3. Operand fetch — OF
4. Instruction execute — EX
5. Operand store — OS

# Superscalar out-of-order pipeline



*3 instructions simultaneously (=pipeline width)*

*Independent instructions can overtake each other*

Personal Super Computing Competence Center

# 'Sequential' processor: super-scalar out-of-order pipeline



**Pipeline depth**

(in order)

(out of order)

**Different processing units**

ALU   MEM1   FP1   BR

MEM2   FP2

FP3

**Out-of-order execution**
**Branch prediction**
**Register renaming**
**…**

(out of order)

(in order)

**Pipeline width**

# Now we are computing sequentially!

# Parallel computing?
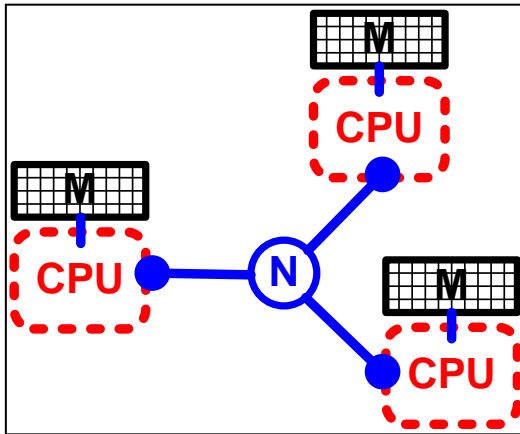
# Super computer: BlueGene/L



- IBM 2007
- 65.536 dual core nodes
  - E.g. one processor dedicated to communication, other to computation
- Each 512 MB RAM
- No 8 in Top 500 Supercomputer list (2010)
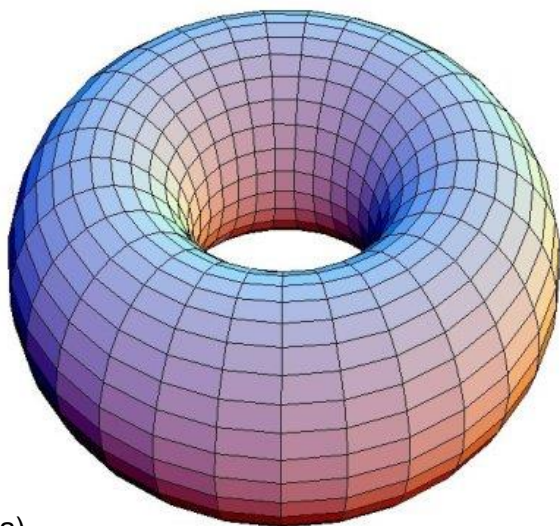  - www.top500.org

# Clusters



- Made from commodity parts
  - or blade servers
- Open-source software available

# Distributed-Memory Architectures



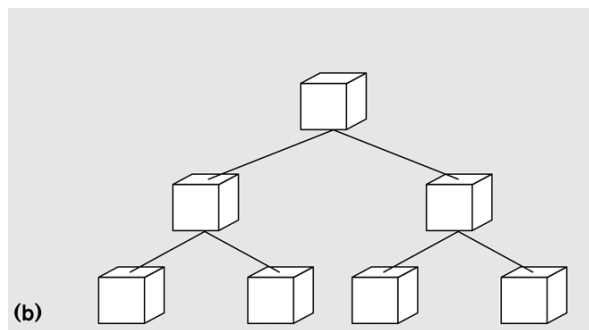- Each process got his own local memory

- Communication through messages

- Process is in control

# BlueGene/L communication networks


(a)


(b)

(a) 3D torus (64x32x32) for standard interprocessor data transfer
- Cut-through routing (see later)

(b) collective network for fast evaluation of *reductions*.
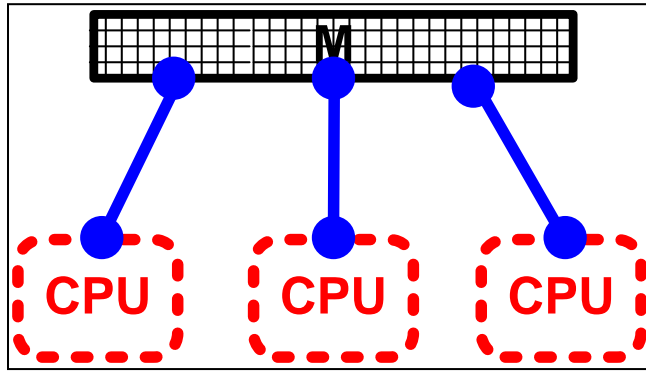
(c) Barrier network by a common wire

# Message-passing

- The ability to send and receive messages is all we need

  - void send(message, destination)

  - message receive(source)

  - boolean probe(source)
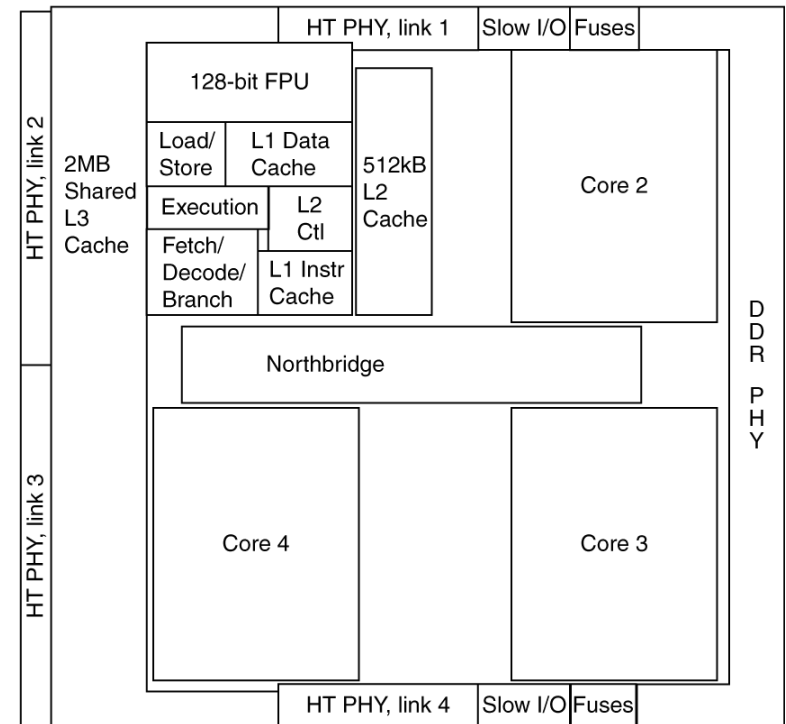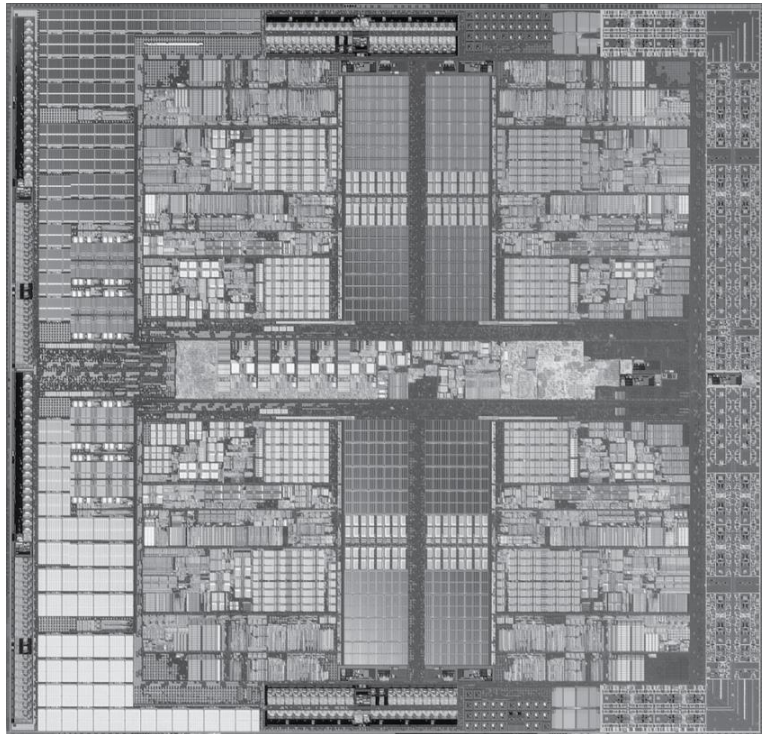
# Multicore computing

# Shared Address-space Architectures



- Example: multiprocessors

# AMD Barcelona: 4 processor cores



HT PHY, link 2

HT PHY, link 3

HT PHY, link 1 | Slow I/O | Fuses

2MB Shared L3 Cache

128-bit FPU

Load/ Store | L1 Data Cache

Execution | L2 Ctl

Fetch/ Decode/ Branch | L1 Instr Cache

512kB L2 Cache

Core 2

Northbridge

Core 4

Core 3

DDR PHY

HT PHY, link 4 | Slow I/O | Fuses
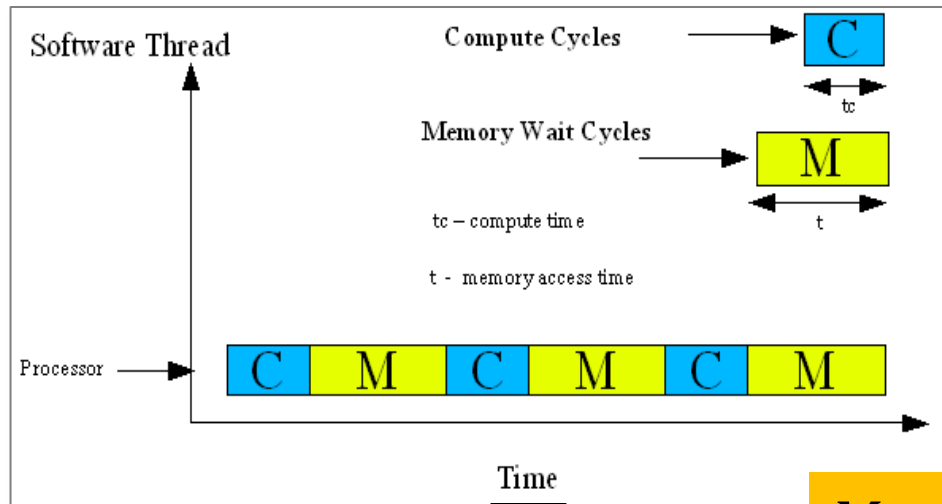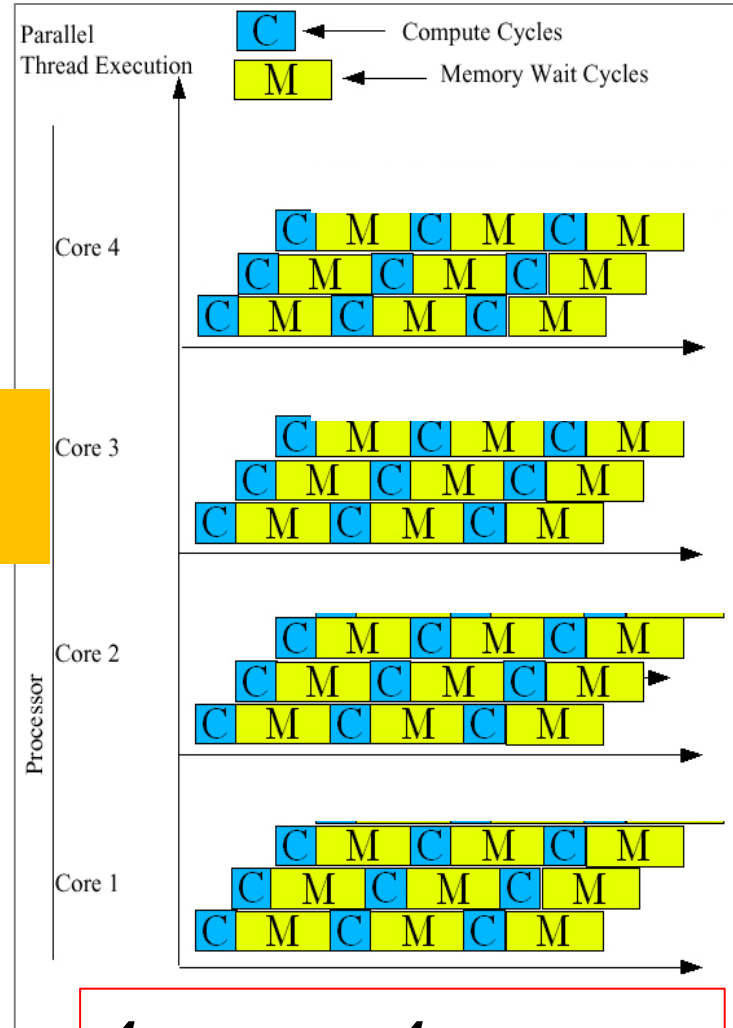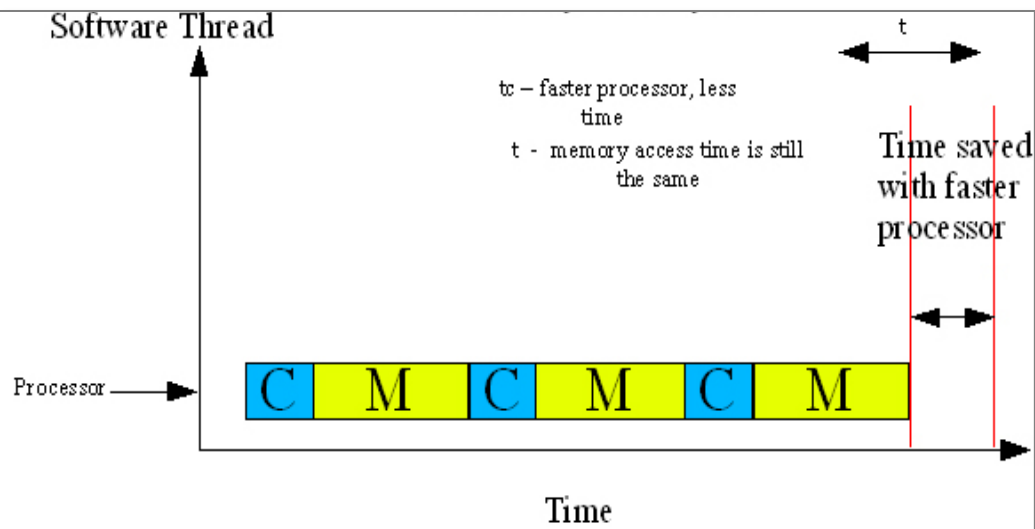
# Thread / core

- A different thread per core
- Each thread can run independently

➡ Multi-threaded programming
- Thread synchronization necessary

- Multiple threads per core also possible
  - Context switches necessary
  - Hardware threads/hyperthreading

# Latency Hiding



Faster CPU

More threads
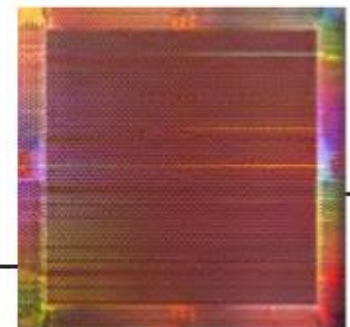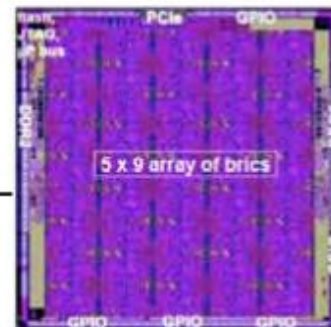
4 cores: x4
Latency hiding: x3

# Next level

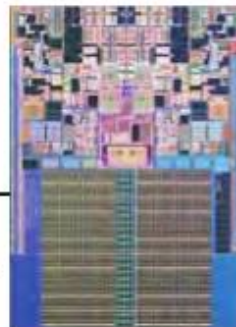# 2010

**350 Million triangles/second**
**3 Billion transistors GPU**

# 1995

**5,000 triangles/second**
**800,000 transistors GPU**

# Parallel processors



| CPUs | DSPs | Multicores | GPUs (arrays) | FPGAs |

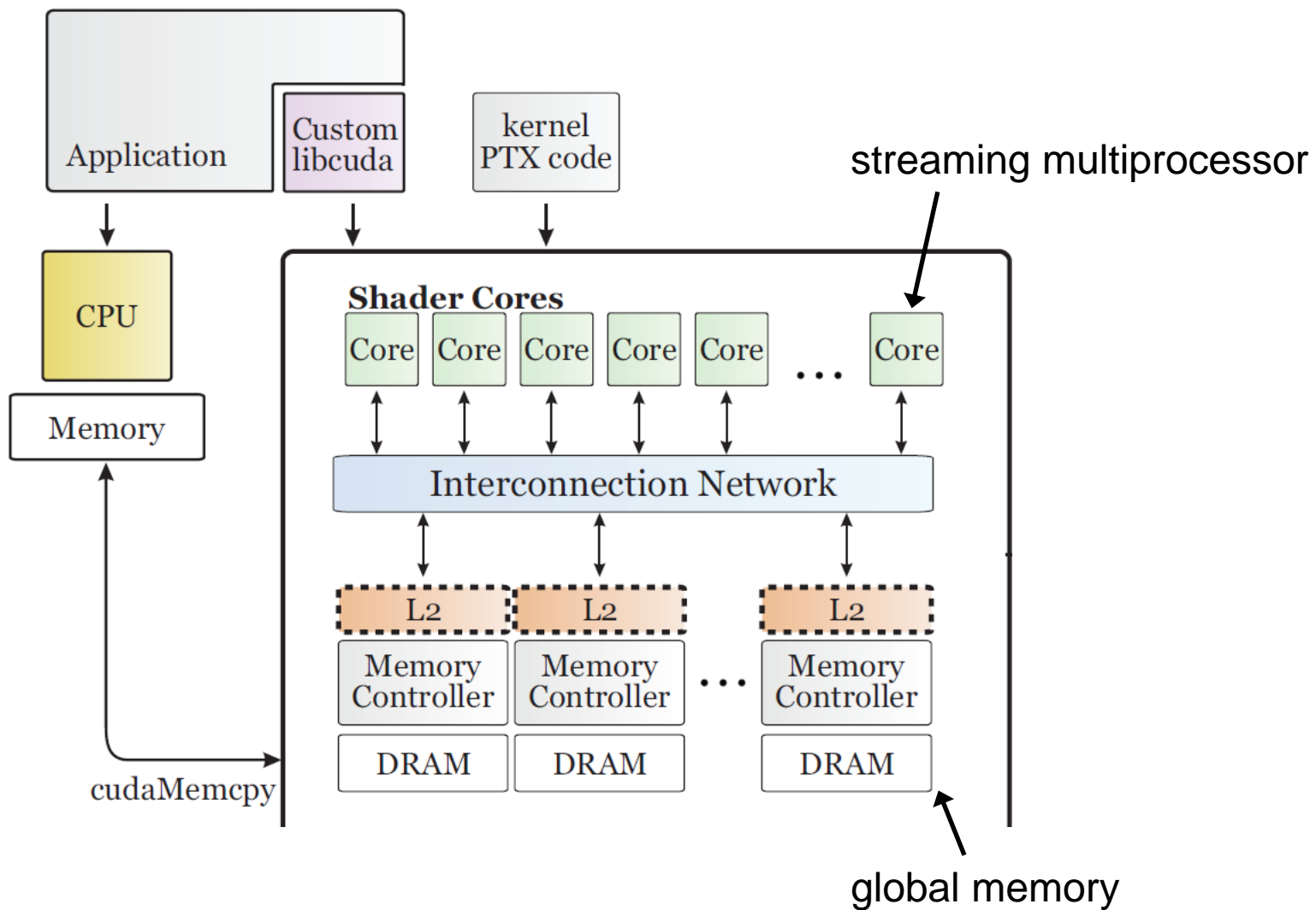**Single Cores** — **Multicores Coarse-Grained CPUs and DSPs** — **Coarse-Grained Massively Parallel Processor Arrays** — **Fine-Grained Massively Parallel Arrays**

Courtesy of ALTERA.

# GPU Architecture



Application

Custom libcuda

kernel PTX code

streaming multiprocessor

CPU

Memory

**Shader Cores**

Core Core Core Core Core ... Core

Interconnection Network

L2 L2 ... L2

Memory Controller | Memory Controller | ... | Memory Controller

DRAM DRAM DRAM

cudaMemcpy

global memory

# 1 Streaming Multiprocessor



**Shader Core**

- Thread Warp
- Thread Warp
- ...
- Thread Warp  **Scheduler**

**SIMD Pipeline**

- Fetch
- Decode
- RF ... RF RF RF

*local/global access (or L1 miss); texture or const cache miss*

- L1 tex
- L1 const
- L1 local& global
- Shared Mem.

**To interconnect**

- Thread Warp
- Thread Warp

**MSHRs**

*Data*  *All threads hit in L1?*

- Writeback

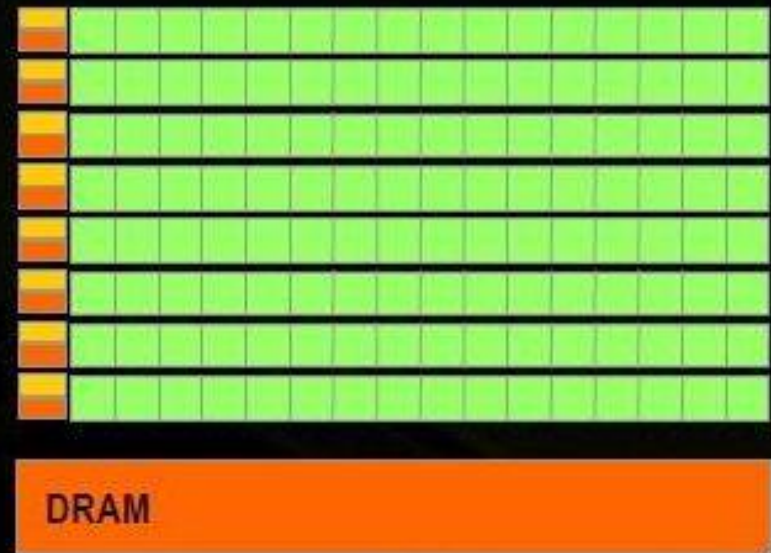The Same Instruction is executed on Multiple Data (SIMD)

width of pipeline: 8 - 32
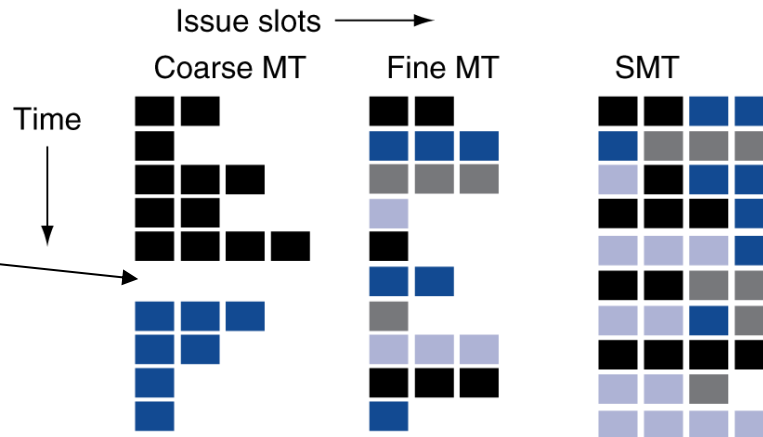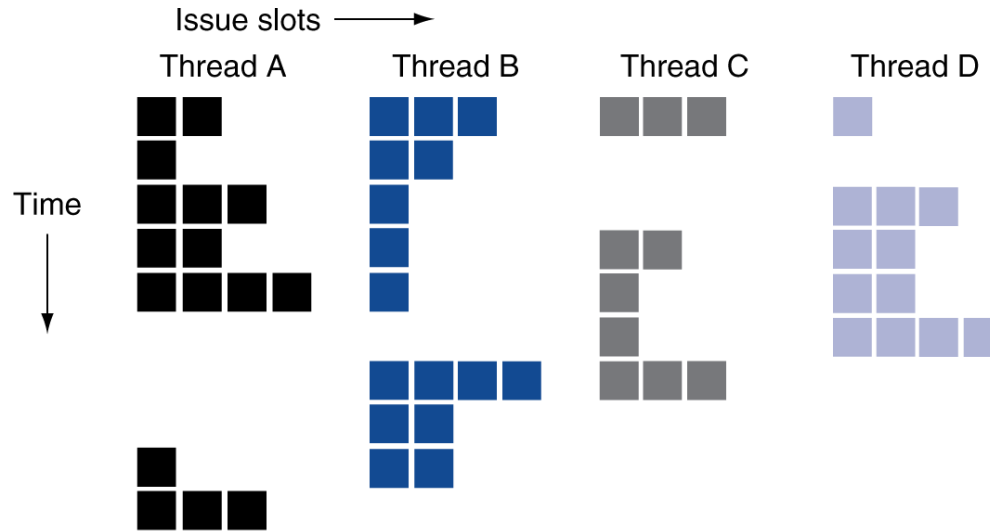
# Why are GPUs faster?



Devote transistors to… computation

# GPU processor pipeline

- ±24 stages (old), now ± 8
- in-order execution!!
- no branch prediction!!
- no forwarding!!
- no register renaming!!
- Memory system:
  - relatively small
  - Until recently no caching
  - On the other hand: much more registers (see later)
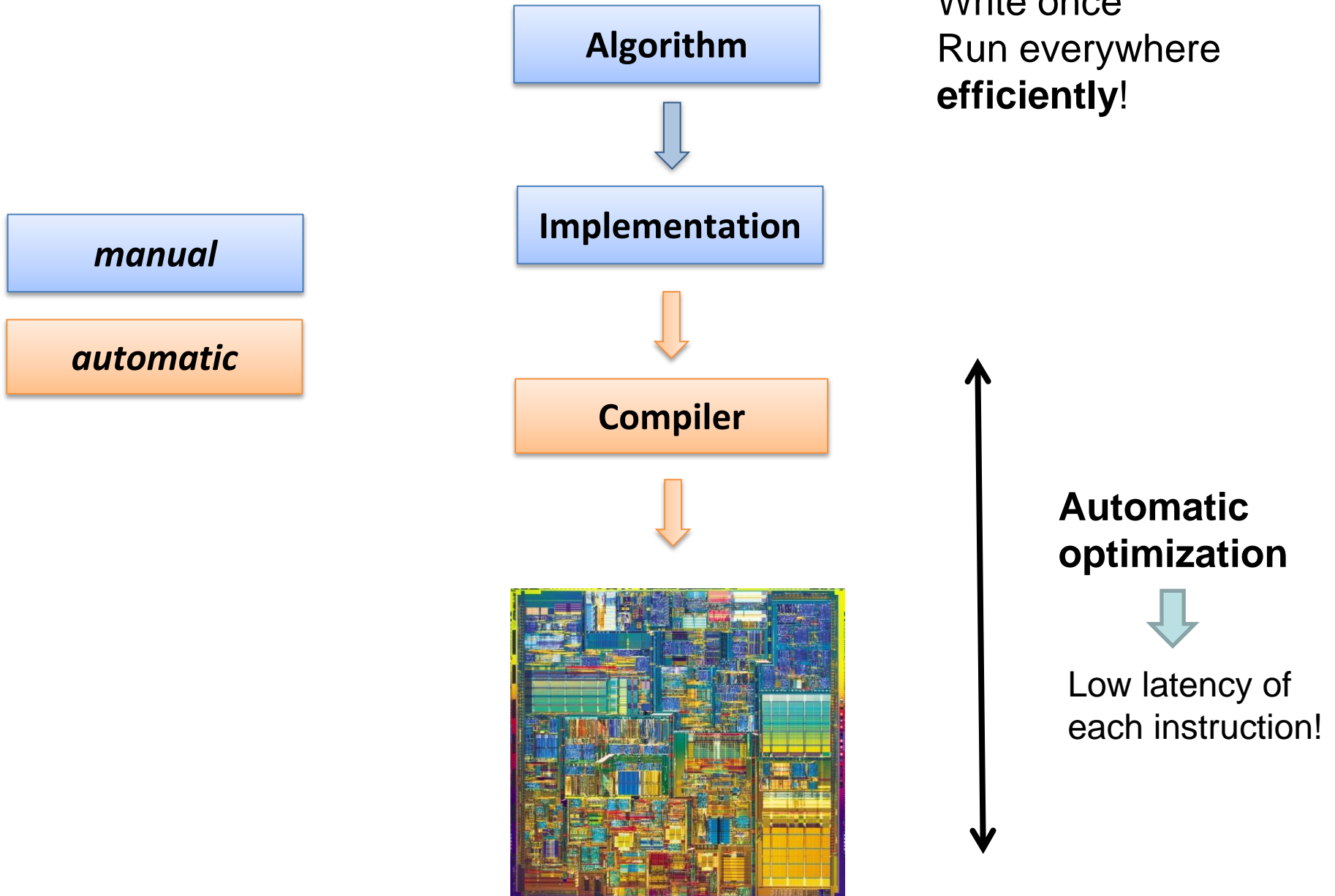
# Multi-Threading (MT) possibilities



**Context switch**

# *Processing power not for free*

# **Obstacle 1**

## Hard(er) to implement

# **Obstacle 2**

## Hard(er) to get efficiency

# CPU computing

**Algorithm**

*manual*

*automatic*

**Implementation**

**Compiler**



Write once
Run everywhere
**efficiently**!

**Automatic optimization**

Low latency of each instruction!

# Computer science is not about computers

# Programmability solutions

# Auto-parallelization

- **Key requirements**
  - A compiler must not alter the program semantics
  - If the compiler cannot determine all dependencies, it has to forego parallelization

- **Compilers sometimes need to act very conservatively**
  - Pointers make it hard for the compiler to deduce memory layout
  - Codes may produce overlapping arrays through pointer arithmetics
  - If the compiler can't tell, it does not parallelize

- **Past 30 years have shown that auto-parallelization**
  - is a tough problem in general
  - is only applicable to very regular loops
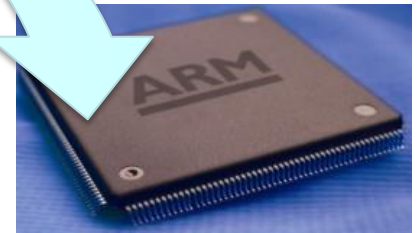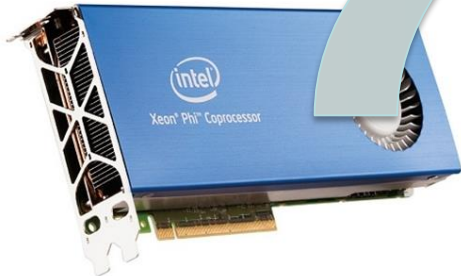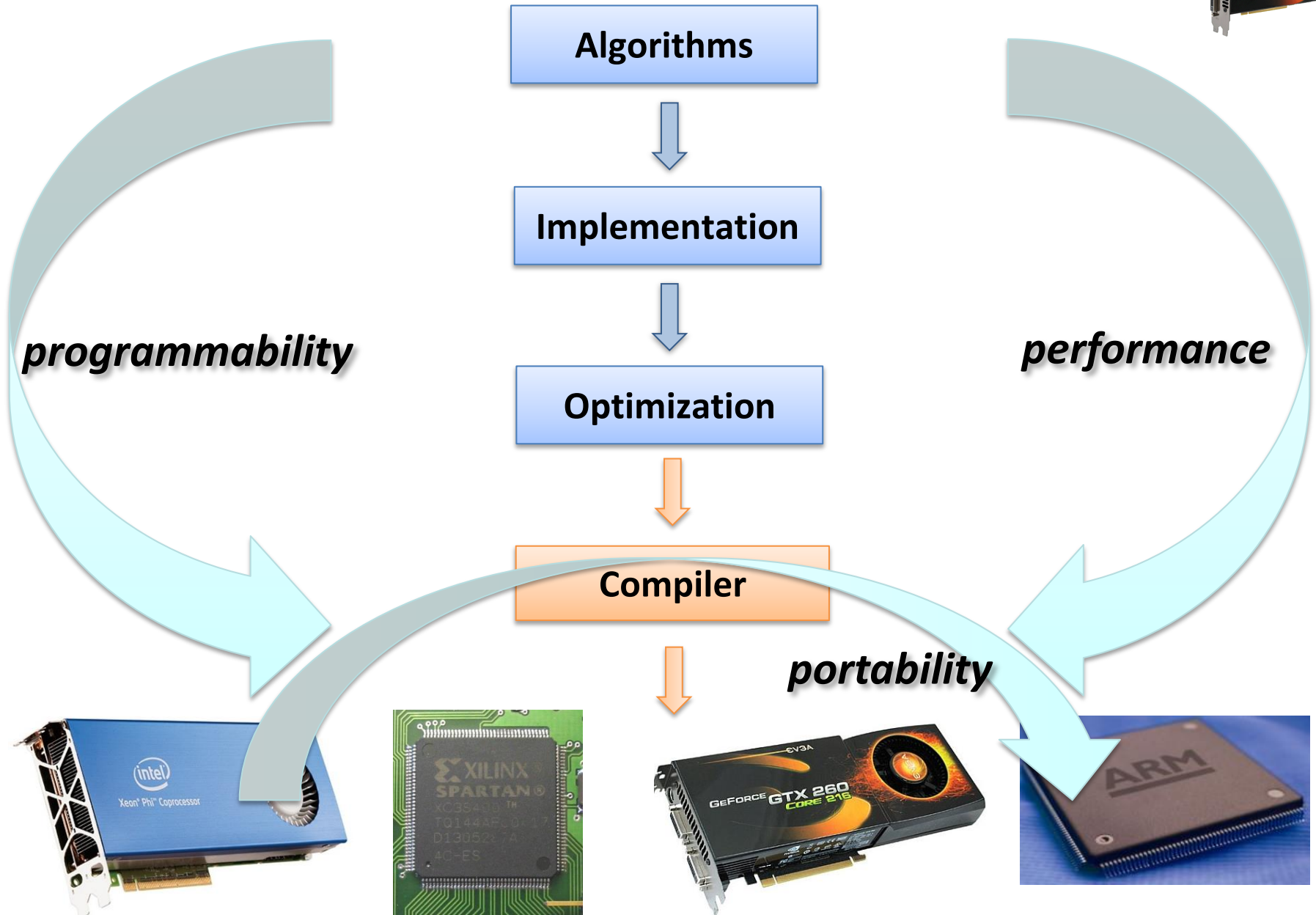  - cannot take care of manual parallelization tasks
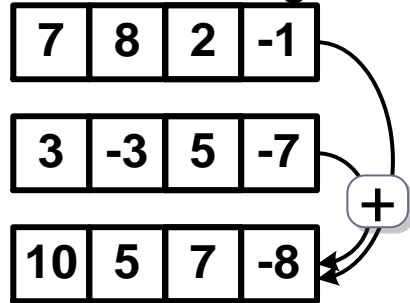
# Computer science

# is n✖t about

# computers

# Challenges of GPU computing

# Intel Xeon Phi

# Single Instruction Multiple Data (SIMD)

128-bit vector registers

| 7 | 8 | 2 | -1 |

| 3 | -3 | 5 | -7 |

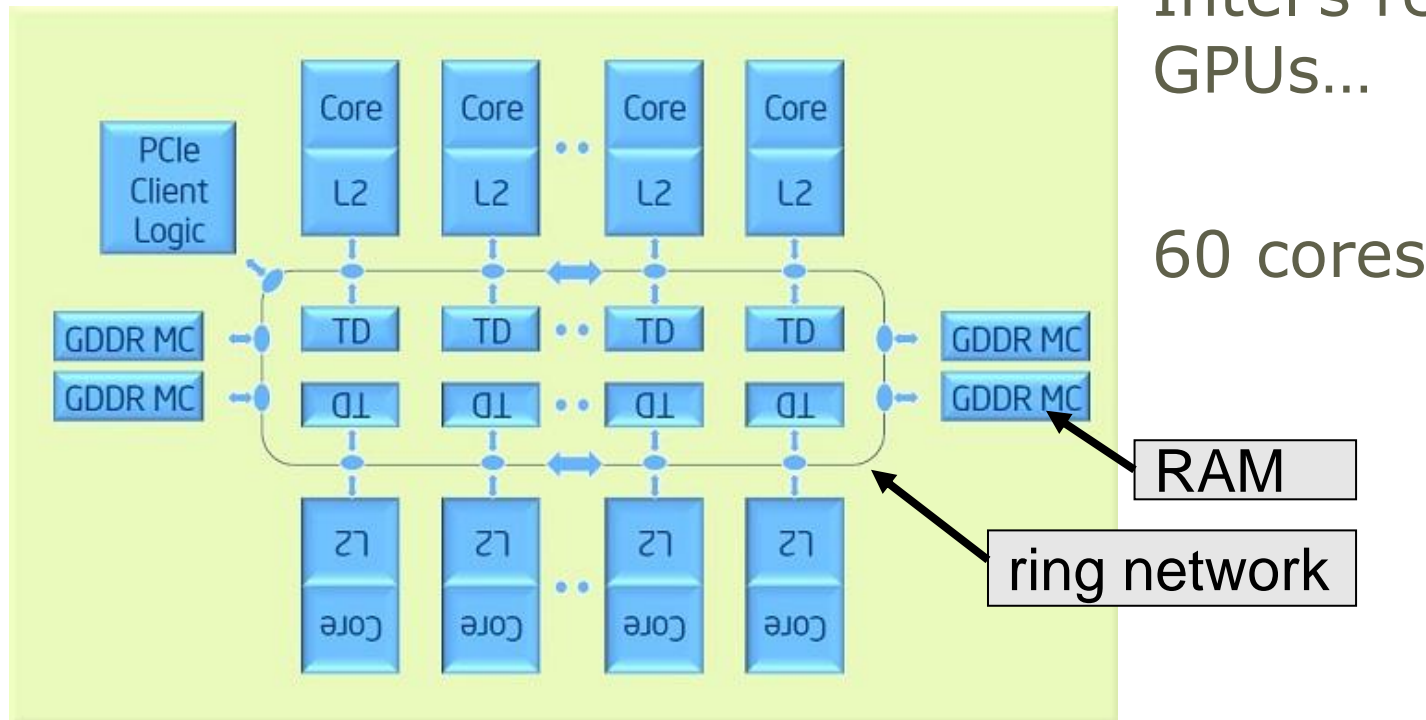**+**

| 10 | 5 | 7 | -8 |

Instructions can be performed at once on all elements of vector registers

- Operate elementwise on vectors of data
  - E.g., MMX and SSE instructions in x86
    - Multiple data elements in 128-bit wide registers
    - Data has to be moved explicitly to/from vector registers
- All processors execute the same instruction at the same time
  - Instruction has to be fetched only once

# Vector processors (SIMD)

- Highly pipelined function units
- Stream data from/to vector registers to units
  - Data collected from memory into registers
  - Results stored from registers to memory
- **Has long be viewed as the solution for high-performance computing**
  - Why always repeating the same instructions (on different data)? => just apply the instruction immediately on all data
- *However: difficult to program*
- Is SIMT (OpenCL) a better alternative??

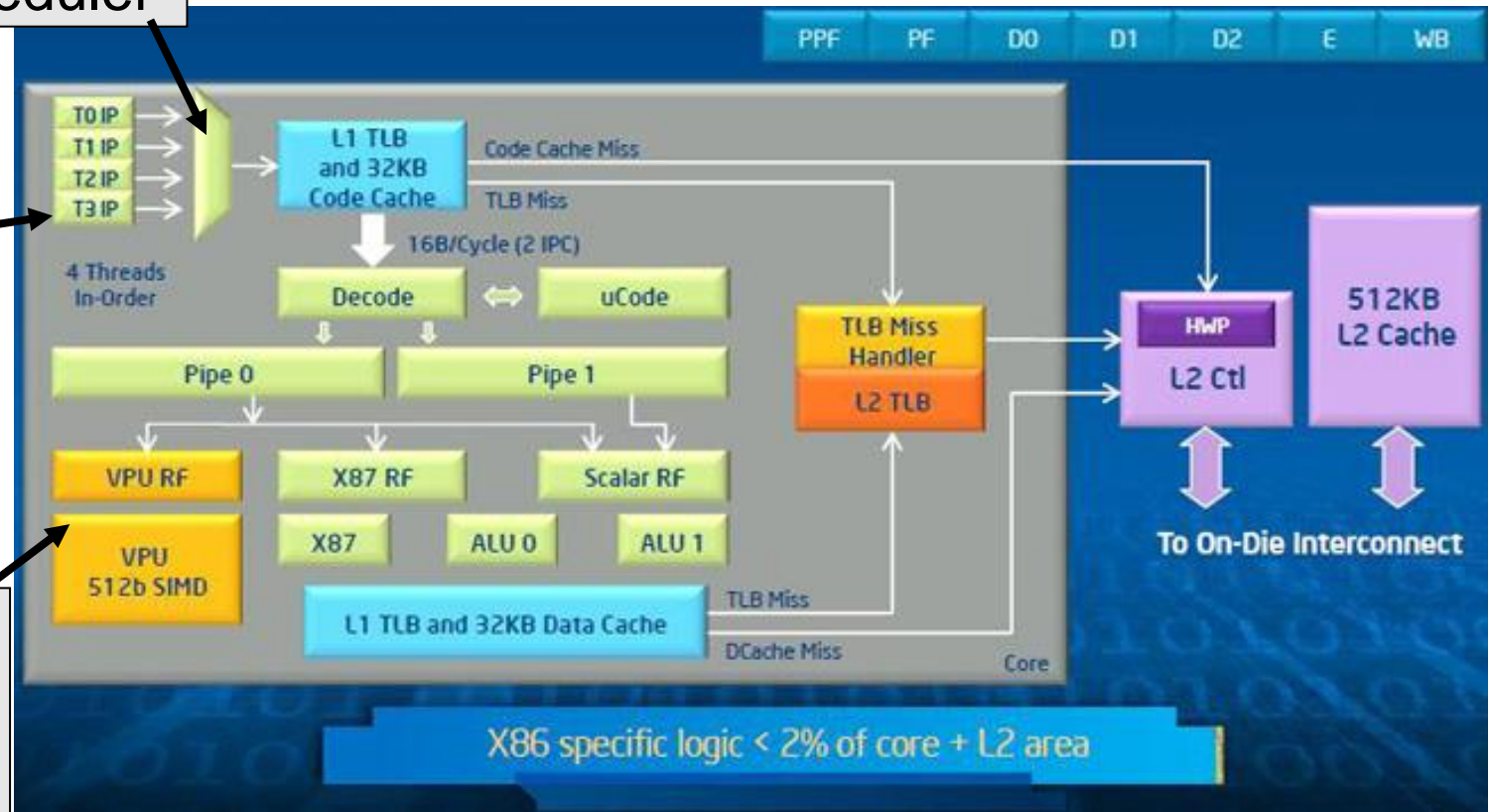# Intel's Xeon Phi coprocessor



Intel's response to GPUs…
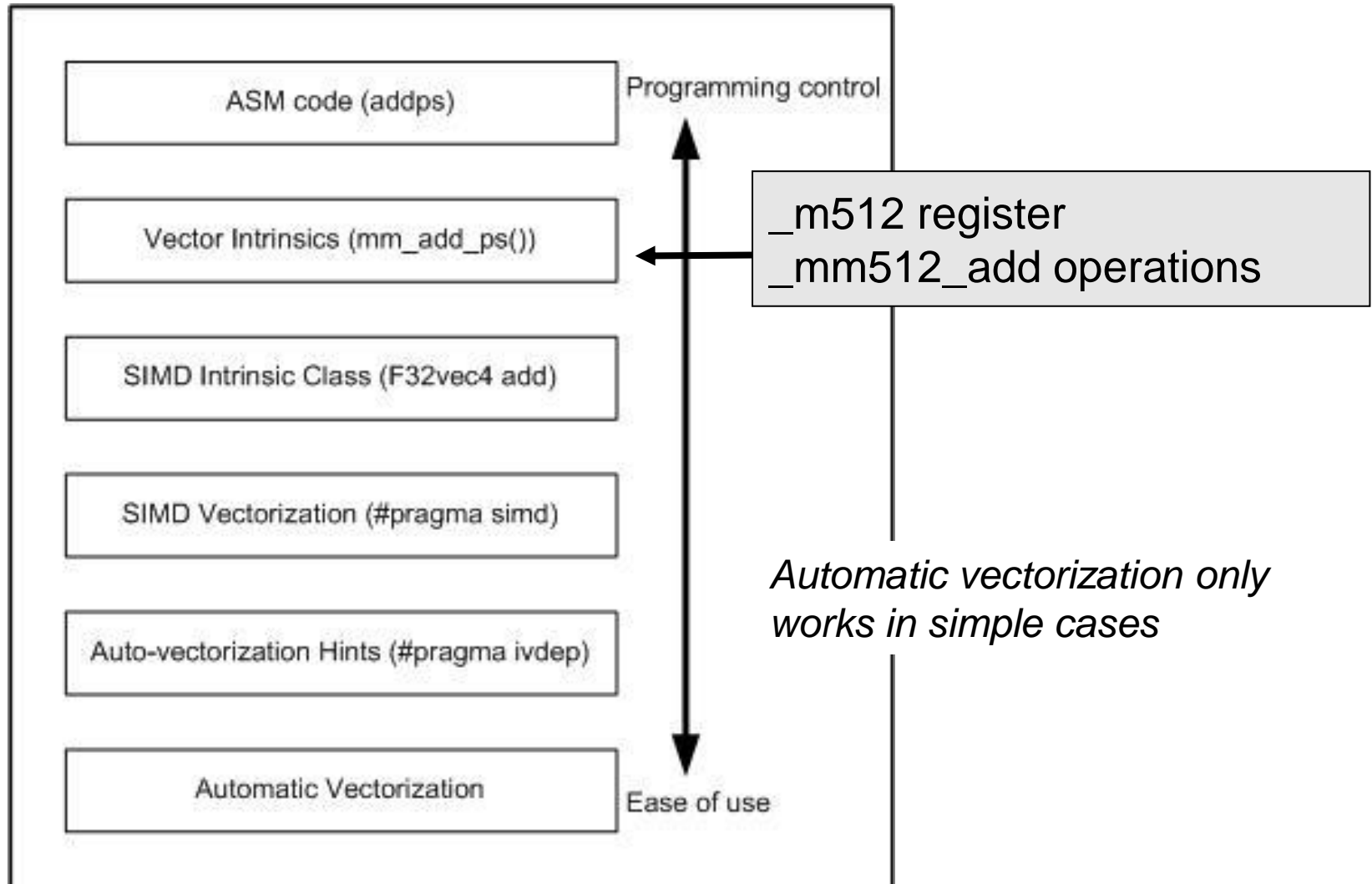
60 cores

RAM

ring network

# Intel's Xeon Phi's core

Thread scheduler

4 hardware threads

512-bit Vector unit (SIMD)

# Vectorization needed for peak performance!!



```
ASM code (addps)

Vector Intrinsics (mm_add_ps())

SIMD Intrinsic Class (F32vec4 add)

SIMD Vectorization (#pragma simd)

Auto-vectorization Hints (#pragma ivdep)

Automatic Vectorization
```

Programming control

Ease of use

_m512 register
_mm512_add operations

*Automatic vectorization only works in simple cases*

# Conclusions

**1** THEORY

**2** EXPERIMENTATION

**3** COMPUTATION

**The third pillar**
**of the scientific world**

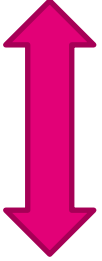# Parallel Programming Paradigms

# Hardware? Software?

- You need to have insight into the hardware!
- No universal hardware/programming model (yet)
- **Intel-approach (SIMD)**
    - Intel sticks to x86 architecture
        - That's what programmers know & they won't change
        - Vectorization necessary

- **OpenCL-approach (SIMT)**
    - Will semi-abstract model remain valid?