



Wrocław
University
of Science
and Technology

Large Scale Data Processing

Lecture 1 – Basic notation, definitions

dr hab. inż. Tomasz Kajdanowicz, Piotr Bielak, Roman Bartusiak

September 27, 2020



HR EXCELLENCE IN RESEARCH



Overview

Big data – 5Vs

Types of processing

Flynn's taxonomy

Compilers



Overview

Big data – 5Vs

Types of processing

Flynn's taxonomy

Compilers



Big data – 5Vs

Value

Having access to big data is all well and good but that's only useful if we can turn it into a value.

Velocity

Speed at which data is emanating and changes are occurring between the diverse data sets

Volume

This refers to the sheer volume of data being generated every second.

5V'S OF BIG DATA

Veracity

Veracity

Data reliability and trust.
Verifying and validating the data

Variety

Variety

Can use structured as well as unstructured data.



Big data – 5Vs

- ▶ **Volume** – enormous volumes of data,
- ▶ **Velocity** – data flows in time from multiple sources and with varying speed,
- ▶ **Value** – data can be hard to obtain,
- ▶ **Veracity (wiarygodność)** – biases, noise and abnormality in data,
- ▶ **Variety** – many sources and types of data both structured and unstructured,

Sometimes this definition is extended to 7Vs:

- ▶ **Validity** – if data correct and accurate for the intended use,
- ▶ **Volatility** – how long is data valid and how long should it be stored,



Overview

Big data – 5Vs

Types of processing

Flynn's taxonomy

Compilers



Types of processing

- ▶ Sequential processing
- ▶ Distributed processing
- ▶ Parallel processing
- ▶ Concurrent processing



Sequential processing

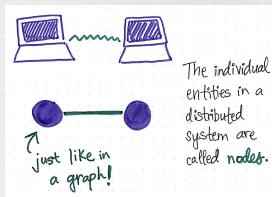
Types of processing

- ▶ processing that occurs in the order that it is received
- ▶ processor inevitably executes the same program

Distributed processing

Types of processing

- ▶ more than one computer (or processor) run an application
- ▶ memory is distributed!
- ▶ includes parallel processing in which a single computer uses more than one CPU to execute programs



- ▶ nodes run operations, that decomposes original large problem
- ▶ operations within a node are fast; communication between nodes is slow
- ▶ nodes operates on their own clocks



Parallel processing

Types of processing

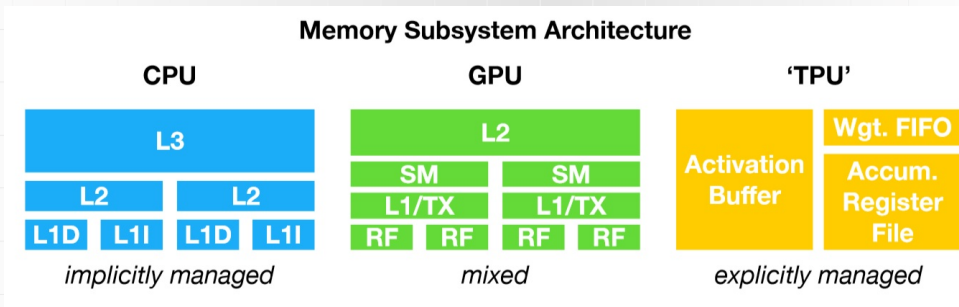
- ▶ Programs use parallel hardware to execute computation more quickly
- ▶ Possible hardware:
 - ▶ multi-core processors
 - ▶ symmetric multiprocessors
 - ▶ graphics processing unit (GPU)
 - ▶ field-programmable gate arrays (FPGAs)
 - ▶ computer clusters
- ▶ Parallel programming requires to think about:
 - ▶ How does code divide original huge problem into smaller sub-problems?
 - ▶ Which is the optimal use of parallel hardware?



CPU vs GPU vs TPU

Types of processing

► Memory Subsystem Architecture

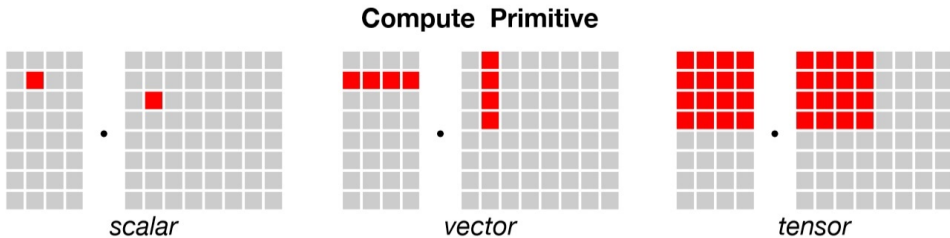




CPU vs GPU vs TPU

Types of processing

► Compute Primitive





CPU vs GPU vs TPU

Types of processing

- ▶ Dimension of data:
 - ▶ CPU: 1 X 1 data unit
 - ▶ GPU: 1 X N data unit
 - ▶ TPU: N X N data unit
- ▶ Performance
 - ▶ CPU can handle tens of operation per cycle
 - ▶ GPU can handle tens of thousands of operation per cycle
 - ▶ TPU can handle upto 128000 operations per cycle
- ▶ Purpose
 - ▶ CPU - designed to solve every computational problem in a general fashion; cache and memory optimal for any general programming problem
 - ▶ GPU - designed to accelerate the rendering of graphics
 - ▶ TPU - designed to accelerate deep learning tasks developed with TensorFlow



Concurrent processing

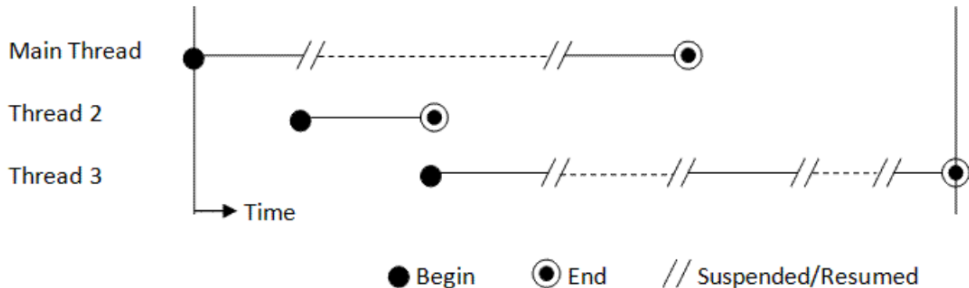
Types of processing

- ▶ concurrency is when multiple sequences of operations are run in overlapping periods of time
- ▶ task A and task B both need to happen independently of each other, and A starts running, and then B starts before A is finished
- ▶ address limits of resources
- ▶ taxonomy:
 - ▶ multitasking
 - ▶ multiprocessing
 - ▶ preemption: preemptive, cooperative



Concurrency example

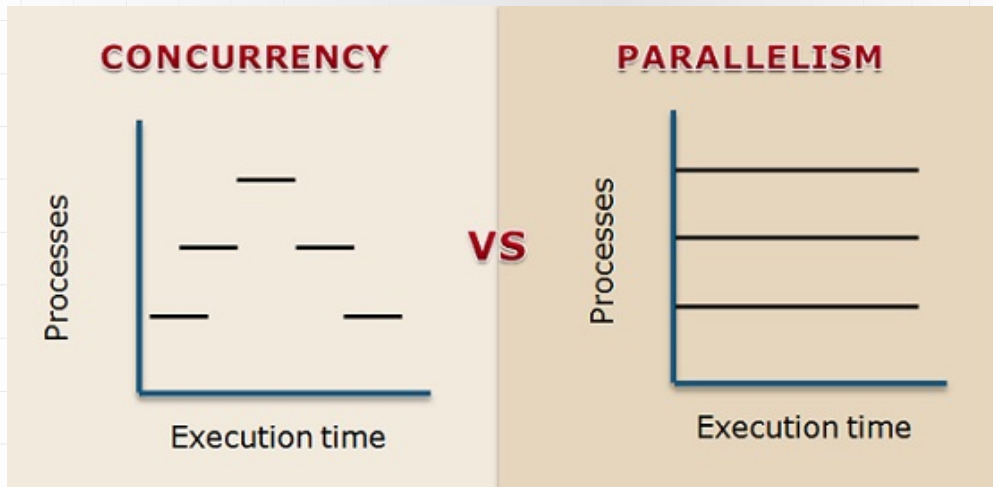
Types of processing





Concurrent vs parallel

Types of processing





Overview

Big data – 5Vs

Types of processing

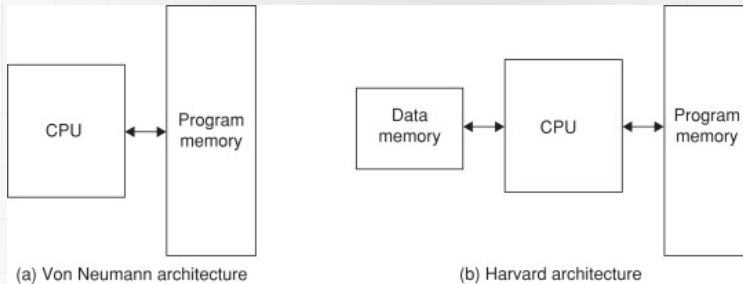
Flynn's taxonomy

Compilers



Computer architecture recap

Flynn's taxonomy



Criterion	Architecture	
	(a)	(b)
Memory/Bus	one	two
Complexity	simple	complicated
Single instruction	two clock cycles	one clock cycle
Performance	low	high (pipelining)
Cost	cheap	high

Data and Instruction streams

Flynn's taxonomy

In Flynn's taxonomy we use following criteria to define system architectures:

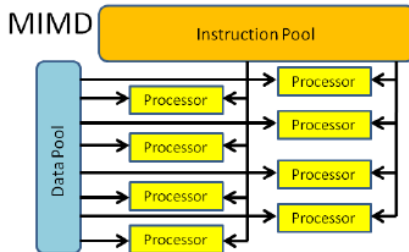
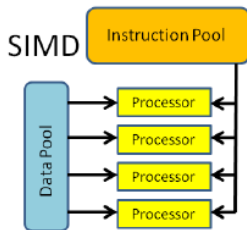
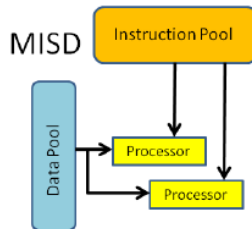
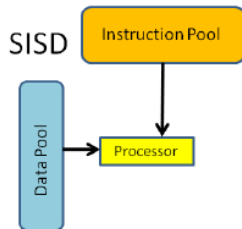
- ▶ number of **instructions** stream(s) - single or multiple,
- ▶ number of **data** stream(s) - single or multiple,

Hence we get following acronyms: **(S/M) I (S/M) D**



Architectures

Flynn's taxonomy





Examples

Flynn's taxonomy

- ▶ SISD – sequential computer; von Neumann architecture; many PCs before 2010 and mainframes
- ▶ SIMD – GPU; modern CPUs with vectorization
- ▶ MISD – systolic computer; fault-tolerant systems
- ▶ MIMD – cluster, where each processor is programmed separately; Intel Xeon Phi; multi-core superscalar processors; distributed systems



Overview

Big data – 5Vs

Types of processing

Flynn's taxonomy

Compilers

Compilation process, optimizations

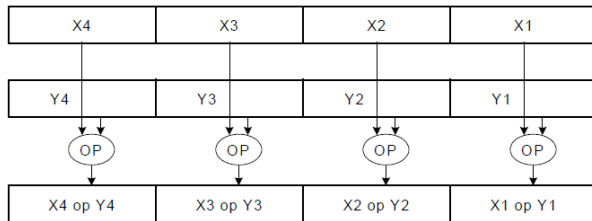
Compilers

- ▶ we won't get into the details of the compilation process,
- ▶ programming languages:
 - ▶ interpreted (e.g., Python, JavaScript),
 - ▶ compiled (e.g., C, C++, Rust),
 - ▶ mixed (e.g., Java - Bytecode+JVM, Python in some cases),
- ▶ interpreted PLs are in general slower than compiled ones (however there is JIT),
- ▶ this is caused by heavy optimizations, which are applied in the compilation process, e.g.:
 - ▶ *removal of unused code* – if the compiler detects that some variable, function etc. is declared, but is never used, then all instructions concerning that variable are removed (can be problematic in some cases like embedded systems; see: *volatile* in C/C++)
 - ▶ *unrolling loops into vector operations* – ...

Vectorization

Compilers

- ▶ 32/64-bit CPUs use general purpose registers with a capacity of 32/64 bits each,
- ▶ however there are some *special registers* with a size equal to the multiple of the architecture size (multiples of 32/64 bits),
- ▶ operations on these registers take one CPU cycle,
- ▶ hence we can speed up computations

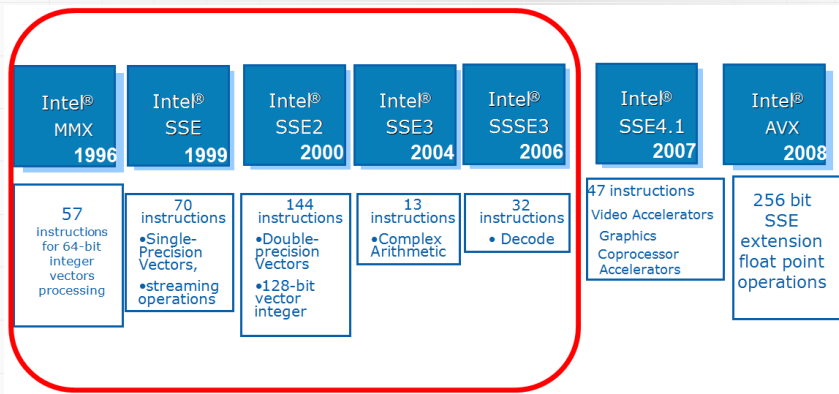


OM15148



Vector registers

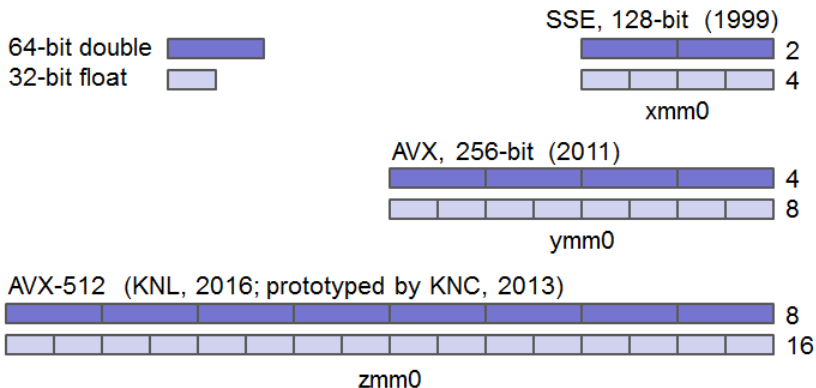
Compilers





Vector registers

Compilers





Large Scale Data Processing

Lecture 1 – Basic notation, definitions

dr hab. inż. Tomasz Kajdanowicz, Piotr Bielak, Roman Bartusiak

September 27, 2020