

Παράλληλη Επεξεργασία

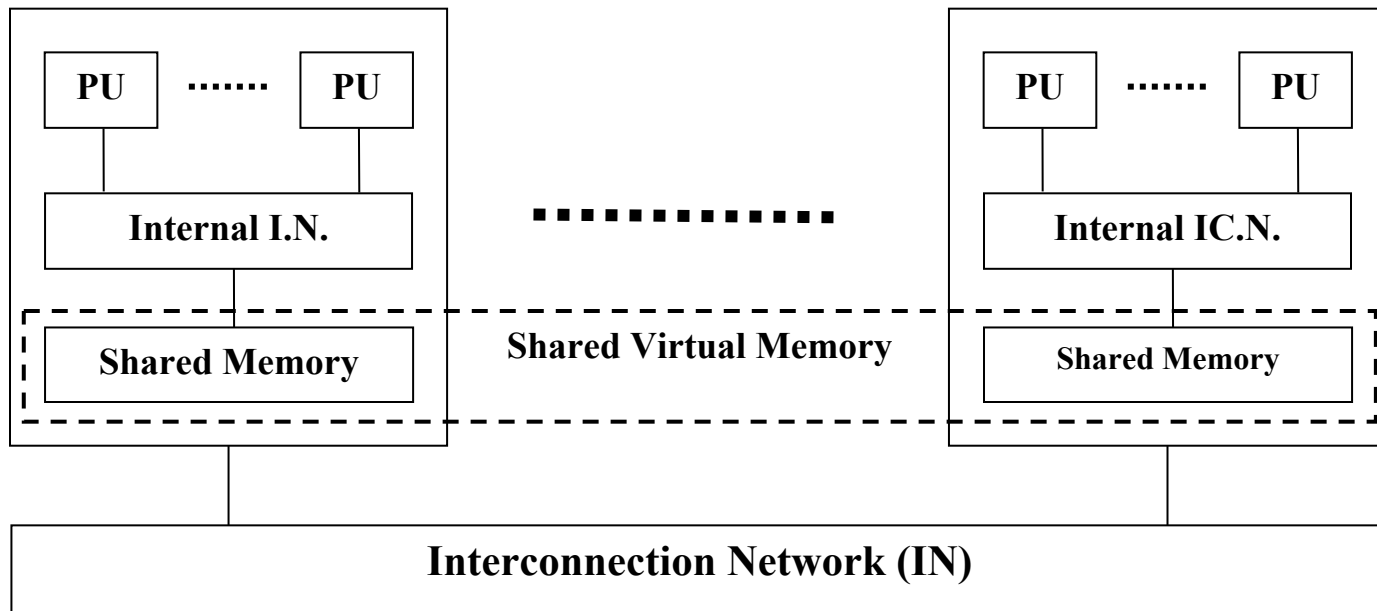
Εαρινό Εξάμηνο 2022-23
«Παράλληλες Αρχιτεκτονικές»

Παναγιώτης Χατζηδούκας, Ευστράτιος Γαλλόπουλος

Outline

- Memory organization
- Taxonomy of computer architectures
- Forms of parallelism

Clusters of Multiprocessors



- **Hybrid memory model**
 - **Distributed memory** of the nodes
 - **Shared memory** for the processors of a single node
- **Shared virtual memory:** “unification” of the memory using software techniques

Memory Organization

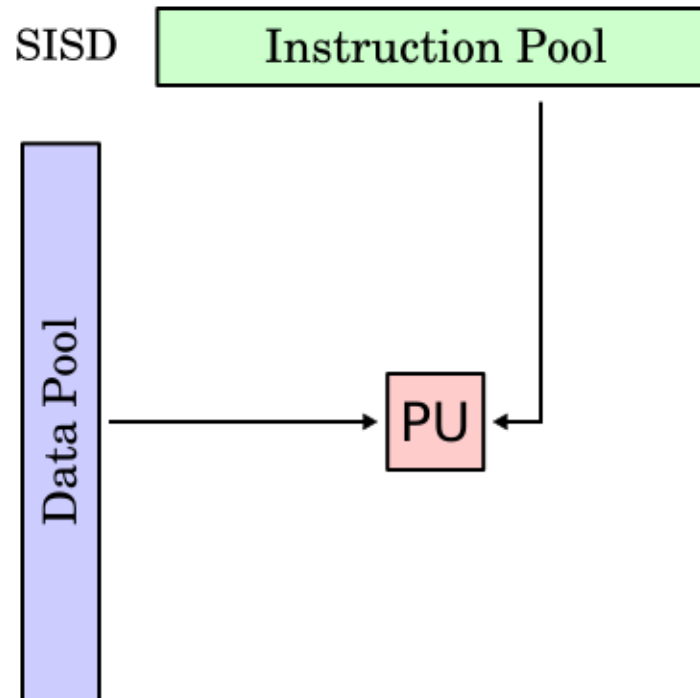
- Shared Memory (SM)
 - The processors have full access to the memory
 - Interconnection network: bus, switch
 - Memory access time can be uniform or not (UMA vs NUMA)
- Distributed Shared Memory (DSM)
 - The processors have full access to the memory
 - Each processor has faster access to some local memory module (NUMA)
- Distributed Memory (DM)
 - Private memory per node
 - Explicit communication (message passing) is needed
 - More scalable but can have high communication costs
- Shared Virtual Memory or Software DSM
 - shared memory on top of distributed memory

Computer Architectures

- Michael J. Flynn's Taxonomy (1966)
 - SISD - Single Instruction/Single Data
 - SIMD - Single Instruction/Multiple Data
 - MISD - Multiple Instruction/Single Data
 - MIMD - Multiple Instruction/Multiple Data

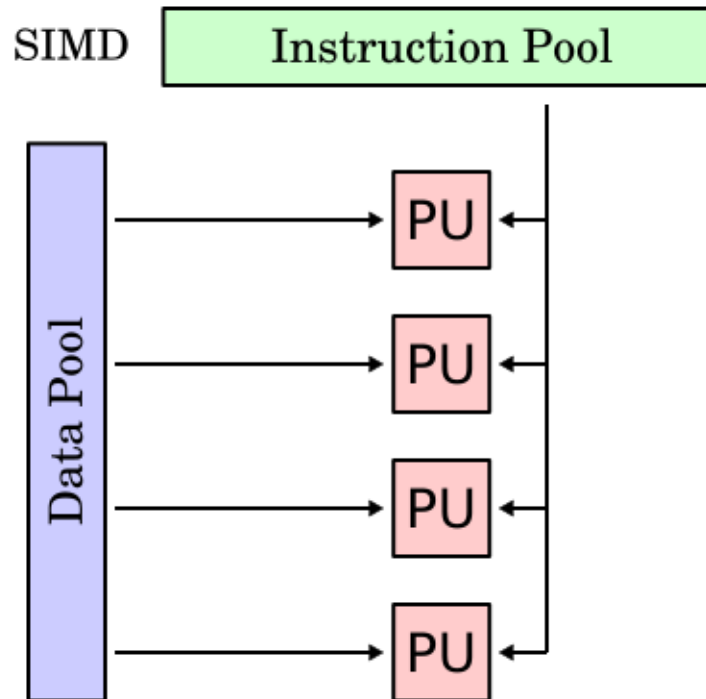
SISD

- Traditional sequential architecture
- Single processing unit
- No parallelism in the instruction and data streams



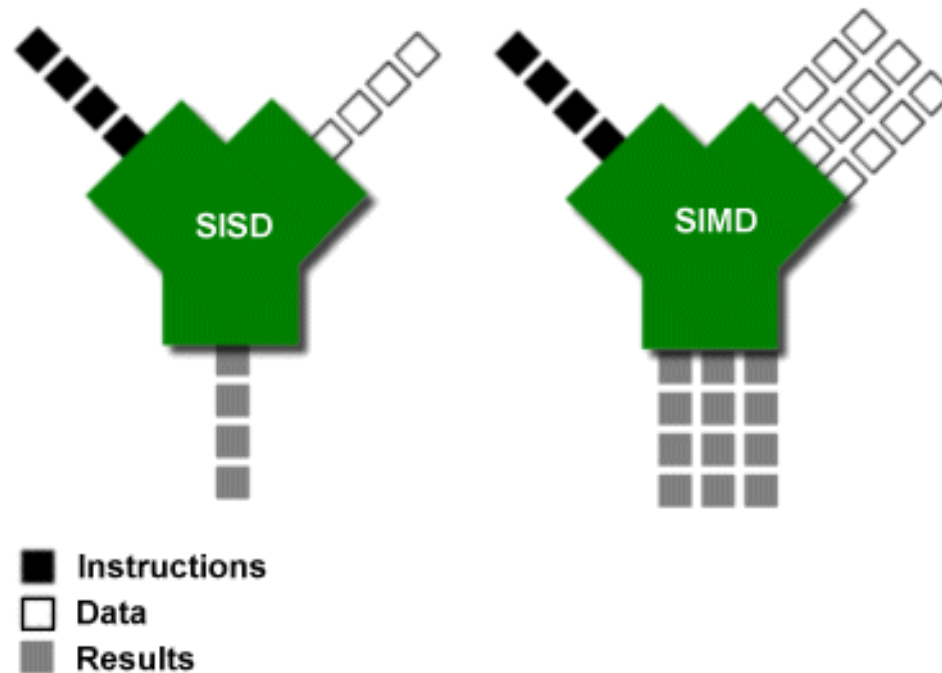
SIMD

- The processing units execute the same instruction on different data.
- Ideal for graphics and scientific computations
- Good choice for usage of multiple processing units



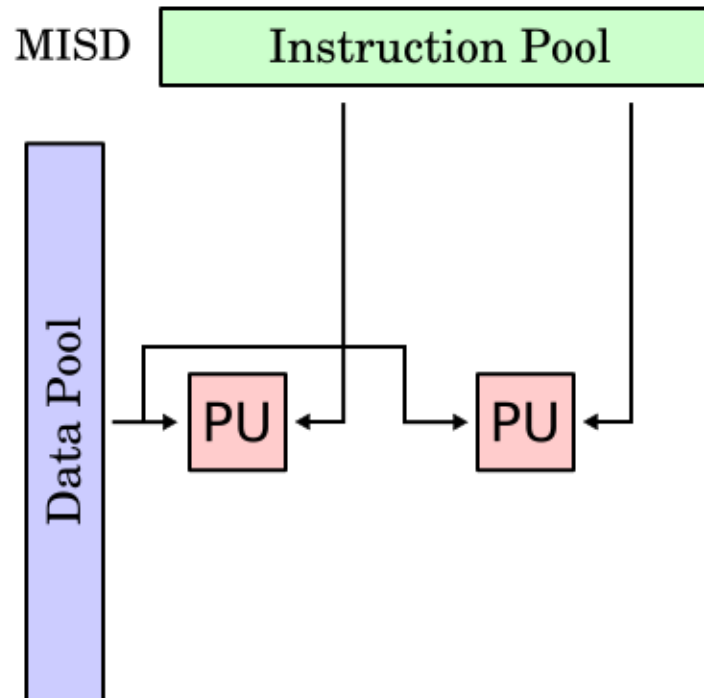
SISD-SIMD comparison

- SIMD often requires support from the compiler



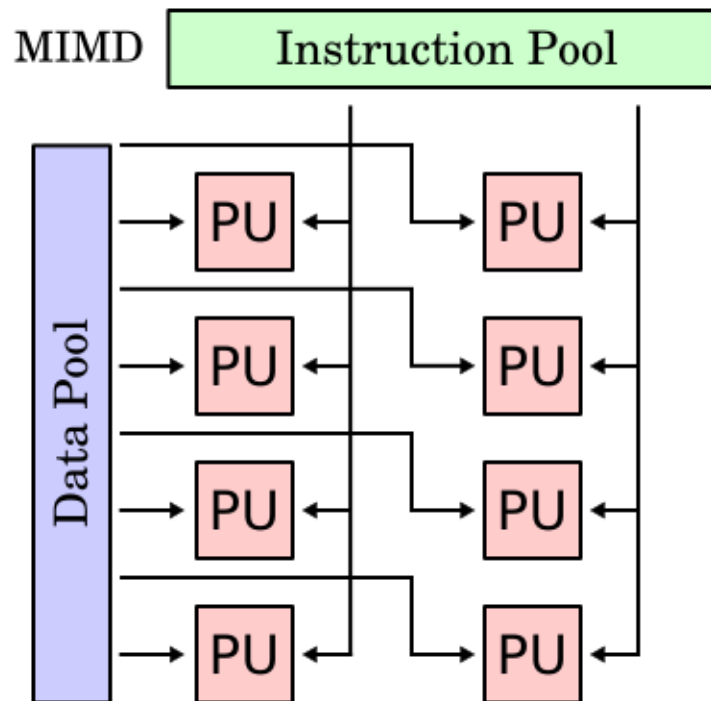
MISD

- Mostly theoretical
- Can be used for fault tolerance



MIMD

- Most multicore systems
- The processing units can also support SIMD
- Various programming models available

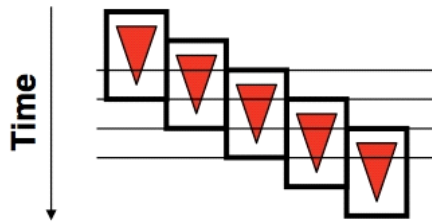
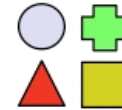


SPMD + MPMD

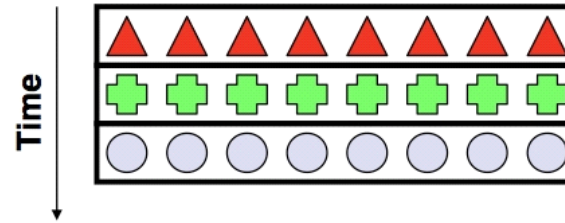
- MIMD: two main parallel execution models
- SPMD: Single program, multiple data
 - Every processing unit executes the same program
- MPMD: Multiple programs, multiple data
 - At least 2 independent programs
 - Master-worker strategy

Forms of Parallelism

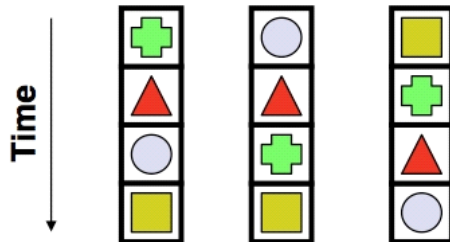
Instructions:



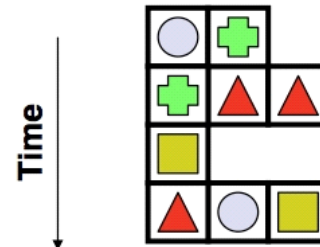
Pipelining



Data-Level Parallelism (DLP)



Thread-Level Parallelism (TLP)



Instruction-Level Parallelism (ILP)

Hardware Pipelining

- Corresponds to SISD

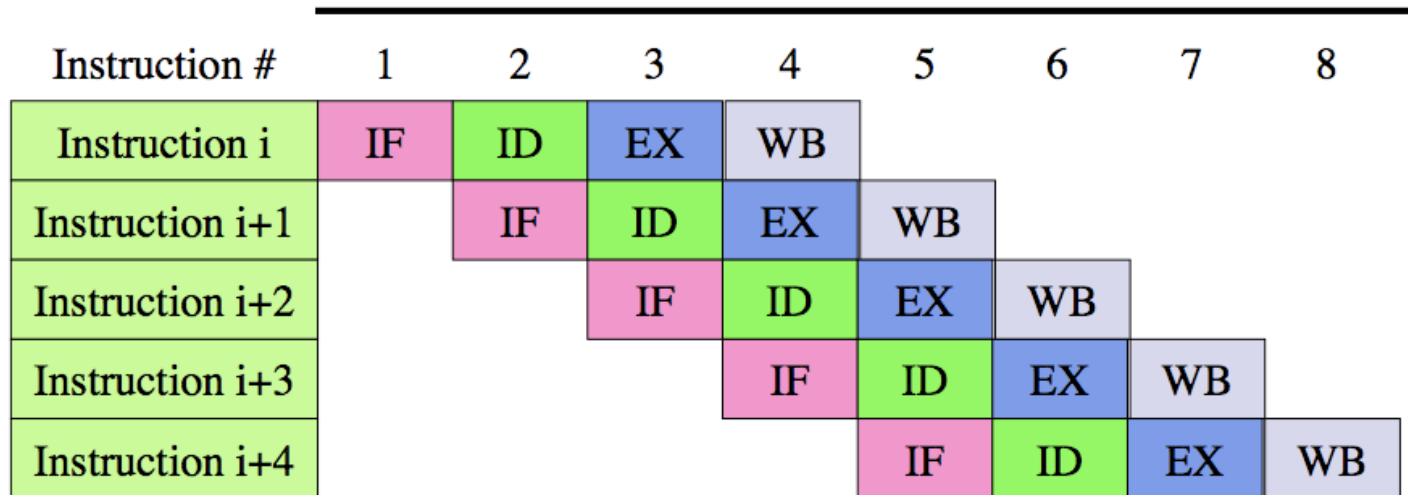
IF: Instruction fetch

ID : Instruction decode

EX : Execution

WB : Write back

Cycles



Instruction Level Parallelism (ILP)

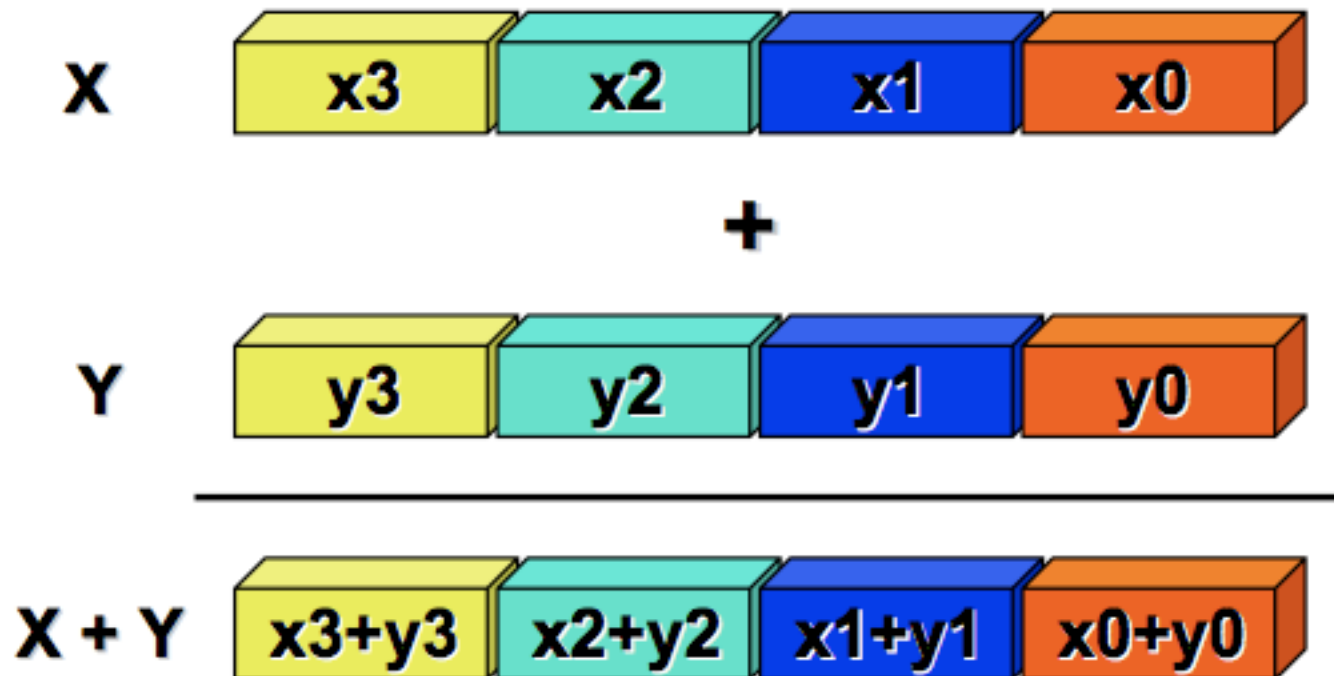
- Simultaneous execution of different instructions of a single program (see SISD)

Cycles

Instruction type	1	2	3	4	5	6	7
Integer	IF	ID	EX	WB			
Floating point	IF	ID	EX	WB			
Integer		IF	ID	EX	WB		
Floating point		IF	ID	EX	WB		
Integer			IF	ID	EX	WB	
Floating point			IF	ID	EX	WB	
Integer				IF	ID	EX	WB
Floating point				IF	ID	EX	WB

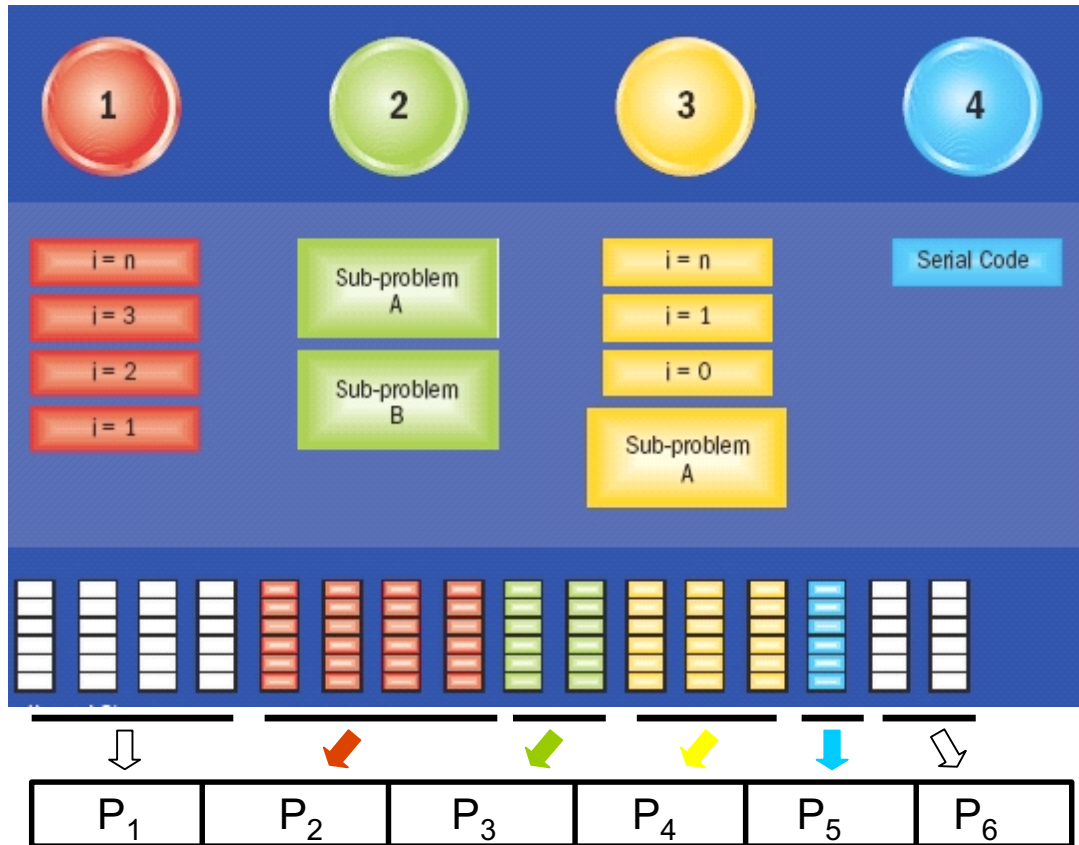
Data Level Parallelism (DLP)

- Corresponds to SIMD
- A single operation (e.g. +) produces multiple results
- X, Y, result: 1-D arrays



Thread Level Parallelism (TLP)

- Corresponds to MIMD



The program is split into threads

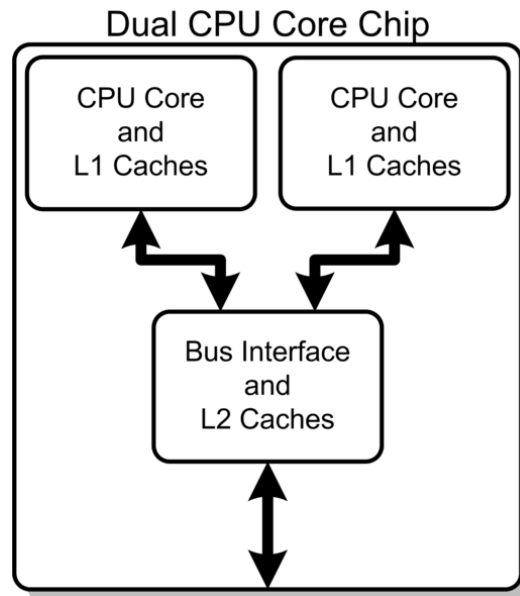
Each thread can include a set of operations

Each thread runs on a different processing unit (core)

Multicore with 6 cores.

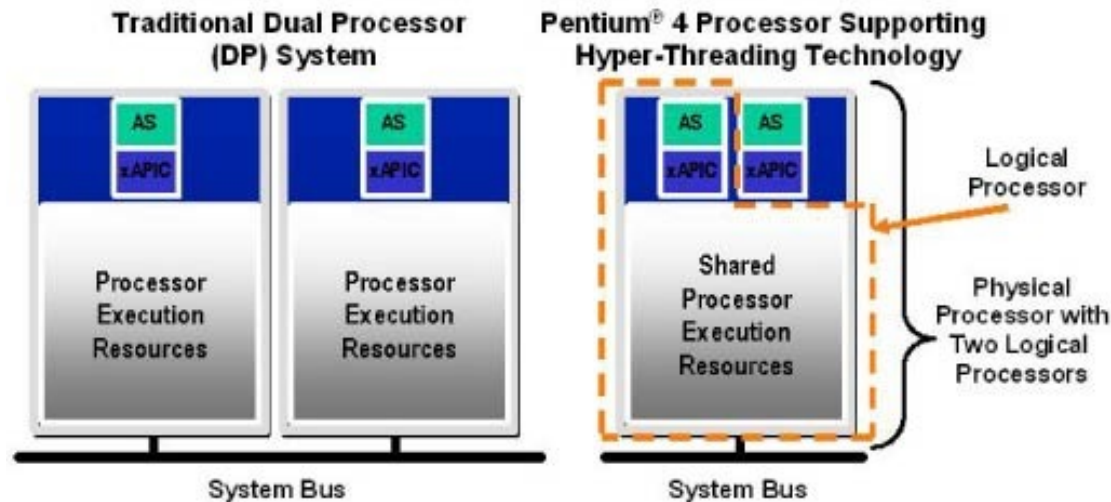
Multi-Core

- Two or more processor cores on a single chip
- Similar performance to traditional single-core multiprocessor systems



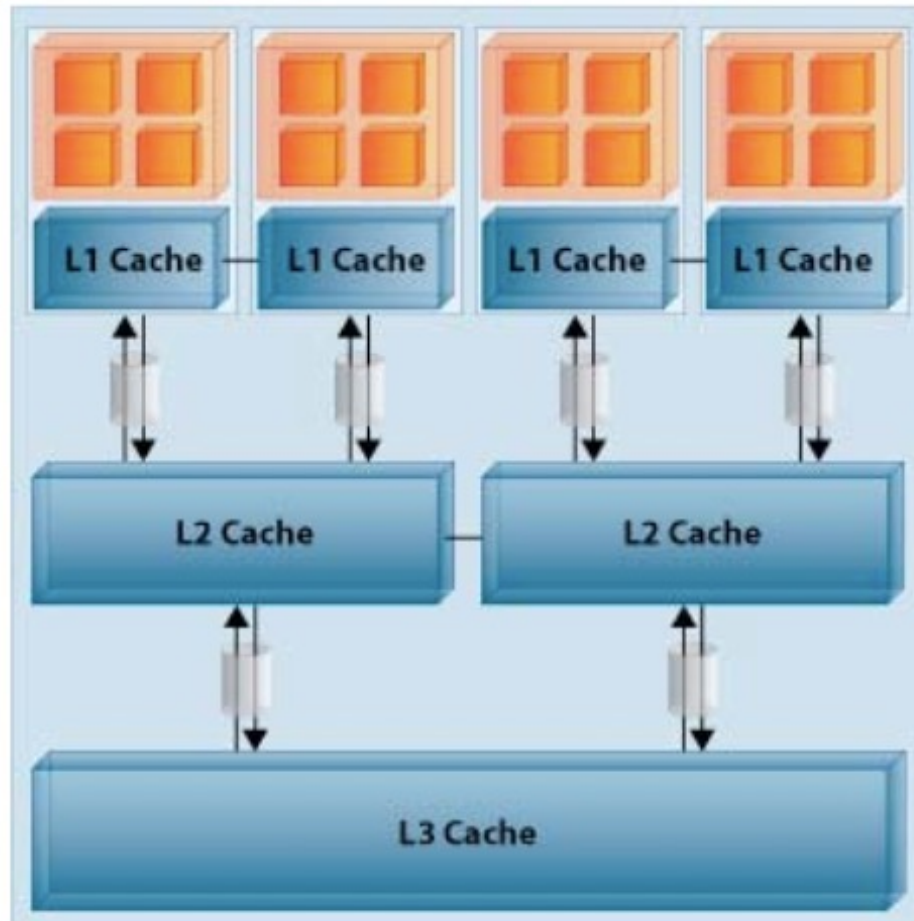
Hyper-Threading

- A physical core has two or more logical cores (hardware threads)
- Sharing of the functional execution units of the core/processor

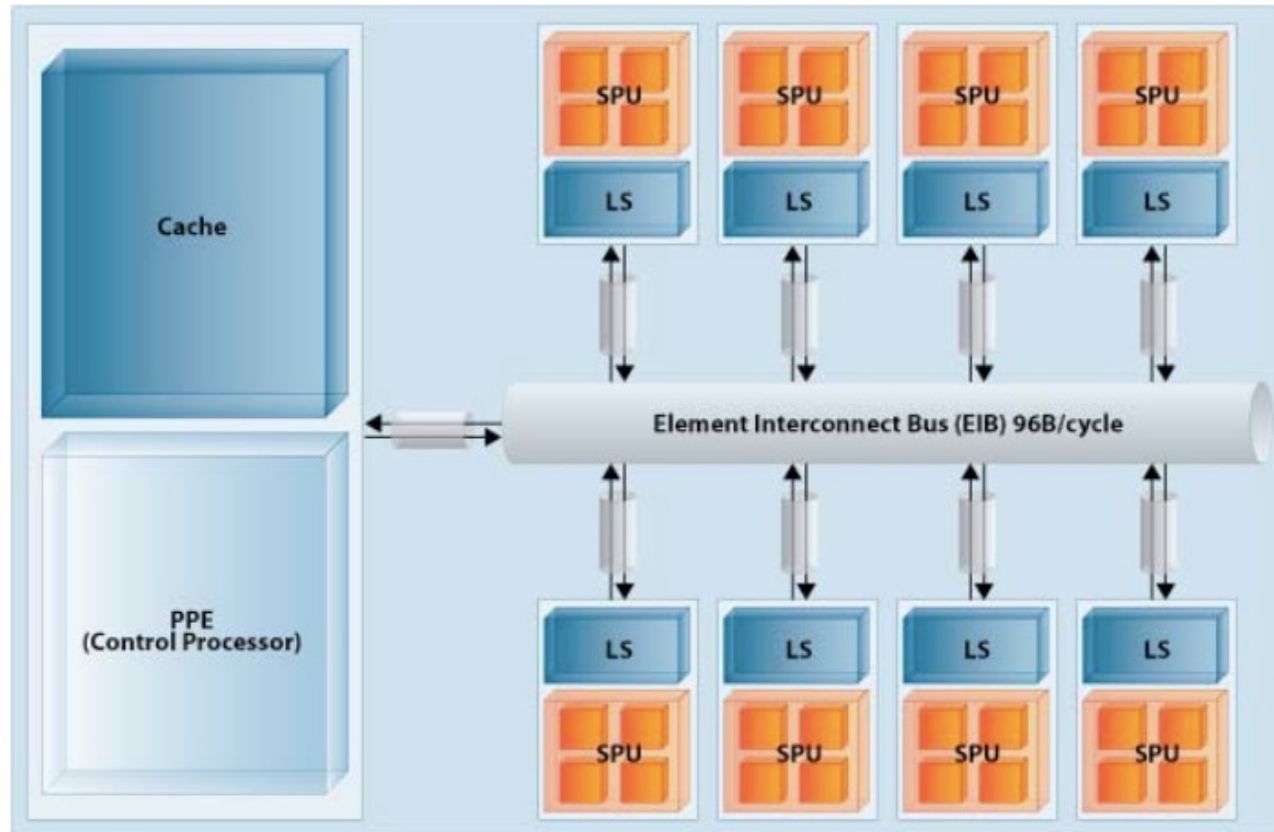


AS = Architecture State (eax, ebx, control registers, etc.)
APIC = Advanced Programmable Interrupt Controller

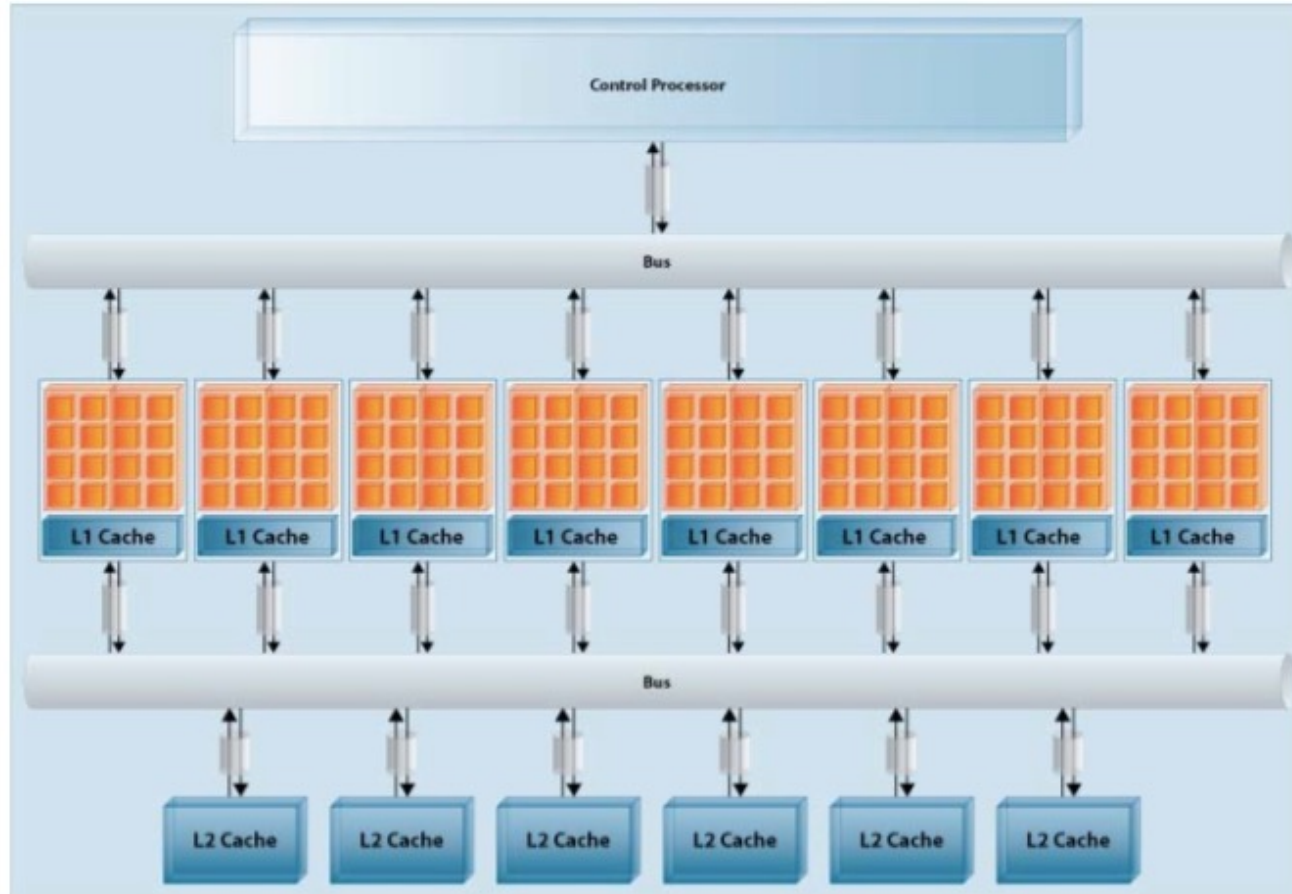
General Multicore Architecture



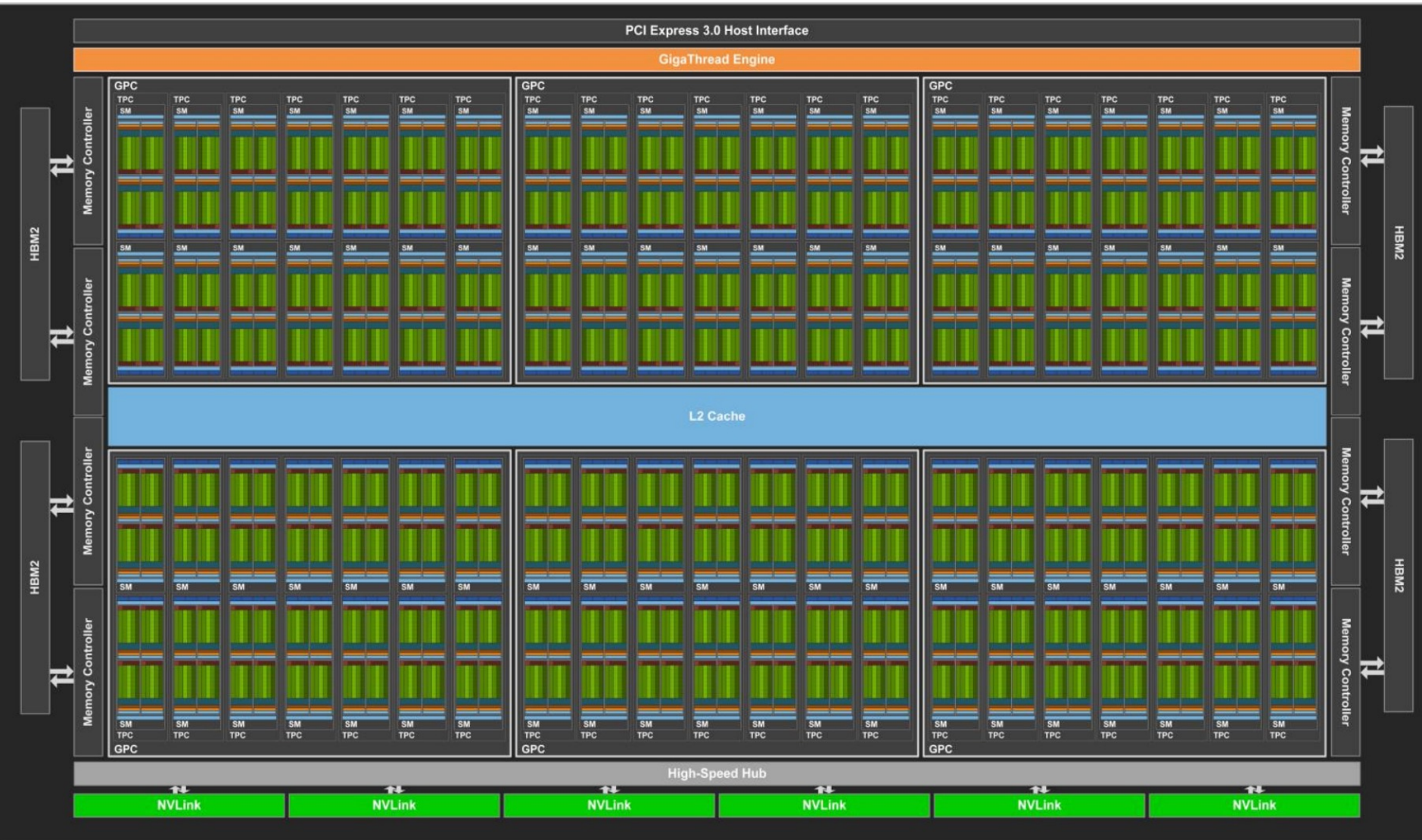
IBM Cell Broadband Engine



NVIDIA GPU G80



NVIDIA GPU V100



Τυπικές αποδόσεις (2018)

Processor	Xeon E5-2695 v4	Tesla P100
Architecture	Broadwell-EP (BDW)	Pascal
# cores	18	3584 (64 cores × 56 SMs)
Clock speed	2.1 GHz (upto 3.3 GHz)	1328 MHz (upto 1480 MHz)
Peak performance (DP)	604.8 GFLOPS	5.3 TFLOPS
Memory type & bandwidth (STREAM Triad)	DDR4 65 GB/s	HBM2 550 GB/s
Processor	Xeon Gold 6140	Xeon Phi 7150
Architecture	Skylake-SP (SKX)	Knights Landing (KNL)
# cores	18	68
Clock speed	2.3 GHz (upto 3.7 GHz)	1.4 GHz (upto 1.6 GHz)
Peak performance (DP)	1324.8 GFLOPS	3046.4 GFLOPS
Memory type (STREAM Triad) & bandwidth	DDR4 95 GB/s	MCDRAM 495 GB/s DDR4 85 GB/s