# Modern Computer Architectures

**Ivan Girotto** – **igirotto@ictp.it**

Information & Communication Technology Section (ICTS)

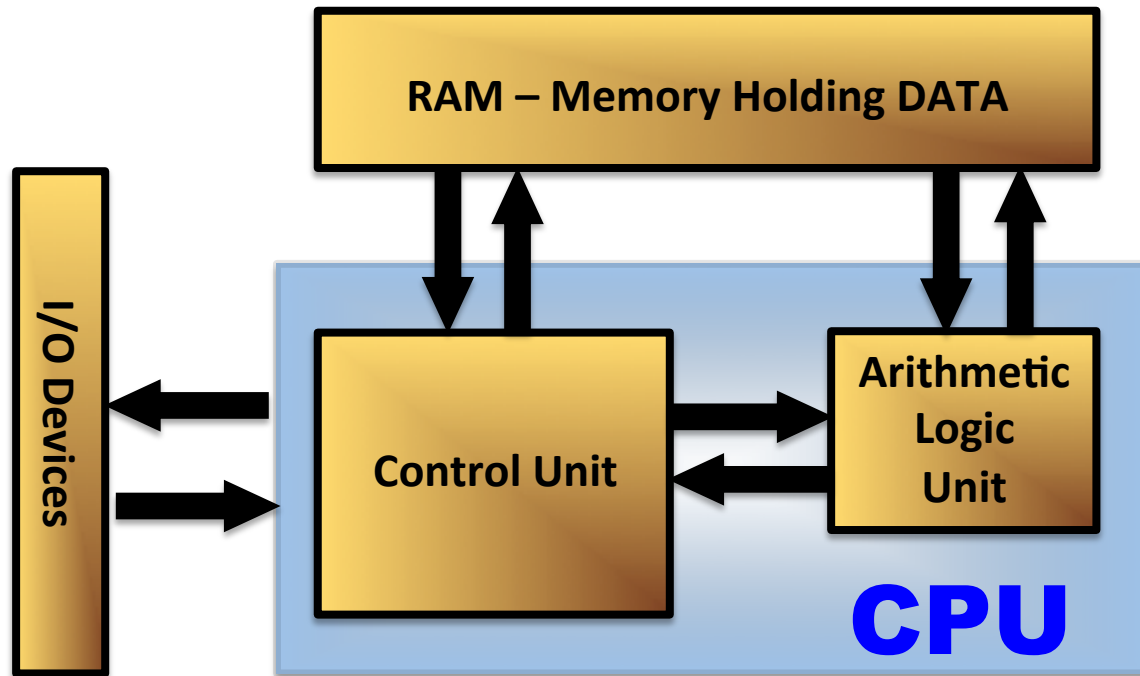International Centre for Theoretical Physics (ICTP)

# Performance Metrics

- When all CPU component work at maximum speed that is called *peak of performance*
  - Tech-spec normally describe the theoretical peak
  - Benchmarks measure the real peak
  - Applications show the real performance value
- CPU performance is measured as:
  - Floating point operations per seconds FLOP/s
- The real performance is in many cases mostly related to the memory bandwidth (Bytes/s) and the exploitation of the parallelism within the CPU

# The Classical Model
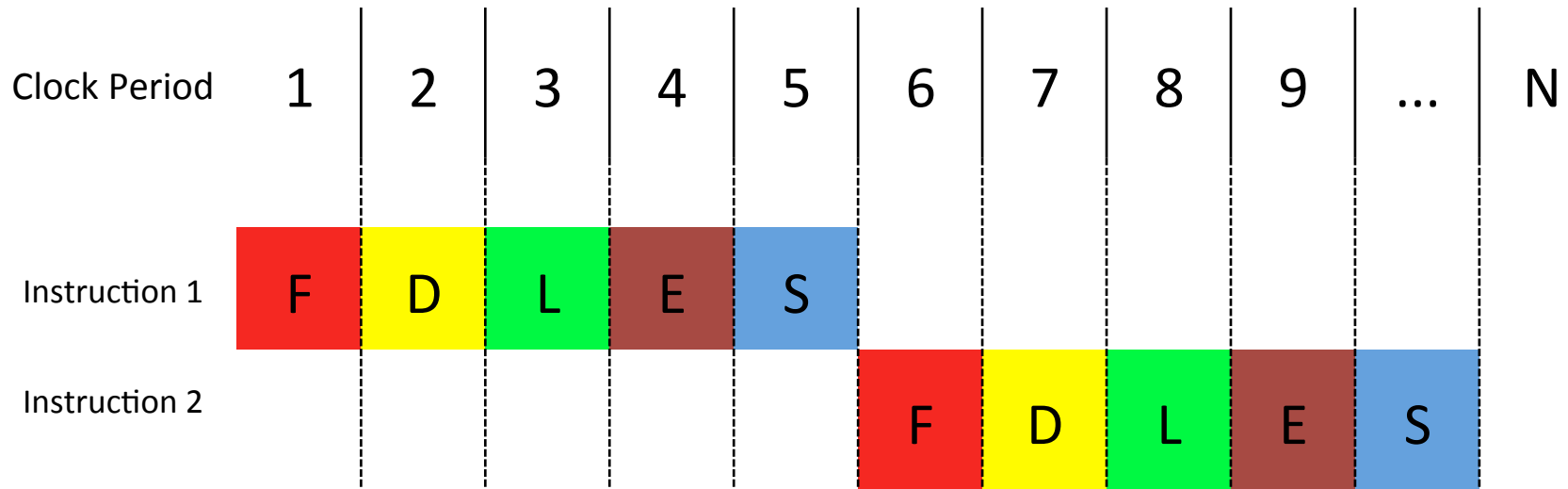
**John Von Neumann**

RAM – Memory Holding DATA

I/O Devices

Control Unit

Arithmetic Logic Unit

**CPU**

# The Instruction Processing Cycle

- Fetch: read the next instruction from memory

    - 001000 00001 00010 000000100001000

- Decode: operands and operation are decoded

    - add, $r1, $r2, 10

- Load: retrieve the data from memory to registers

- Execute: execute the instruction
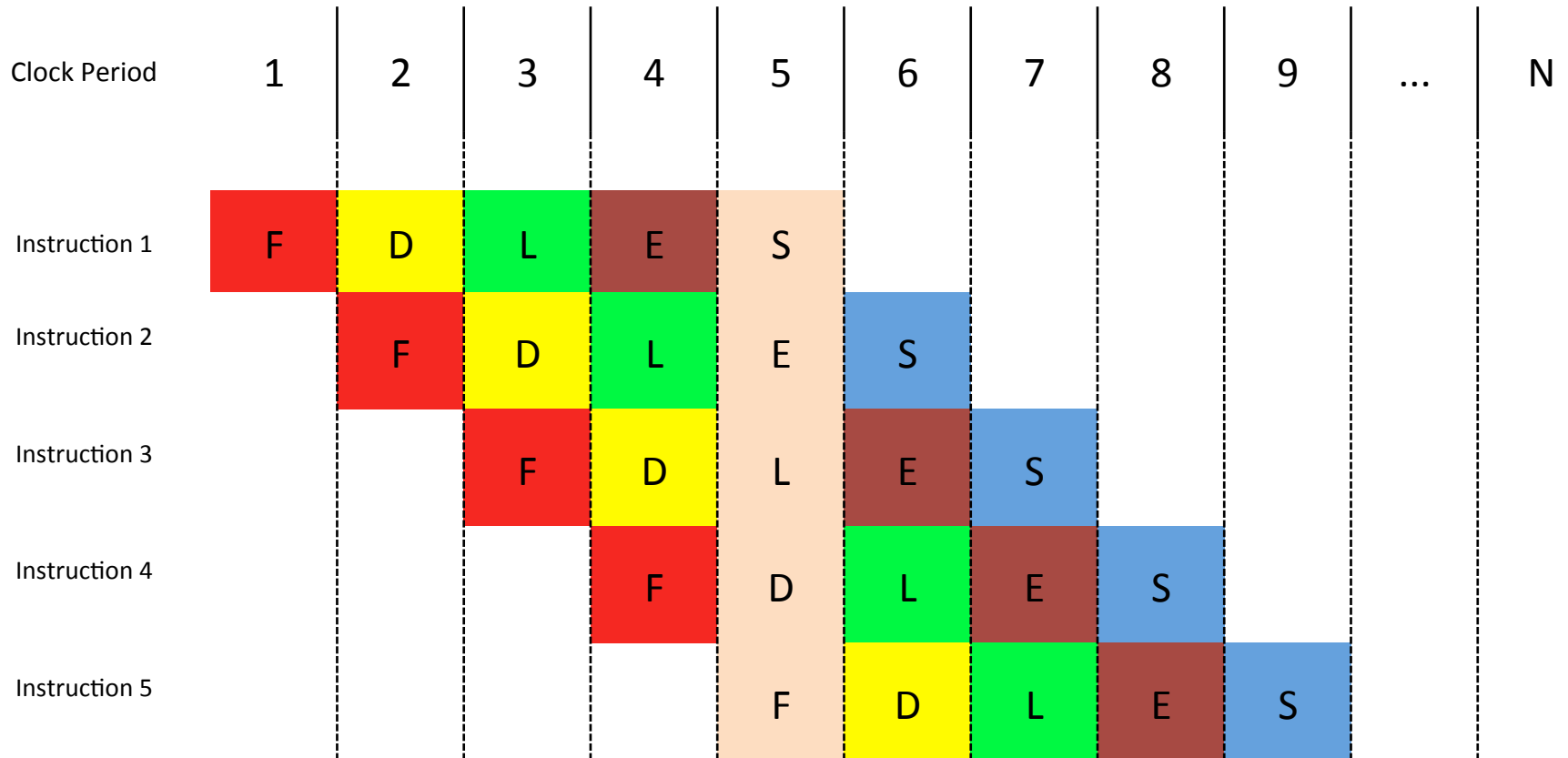
    - $r1 = 4500 + 10

- Store: store the results
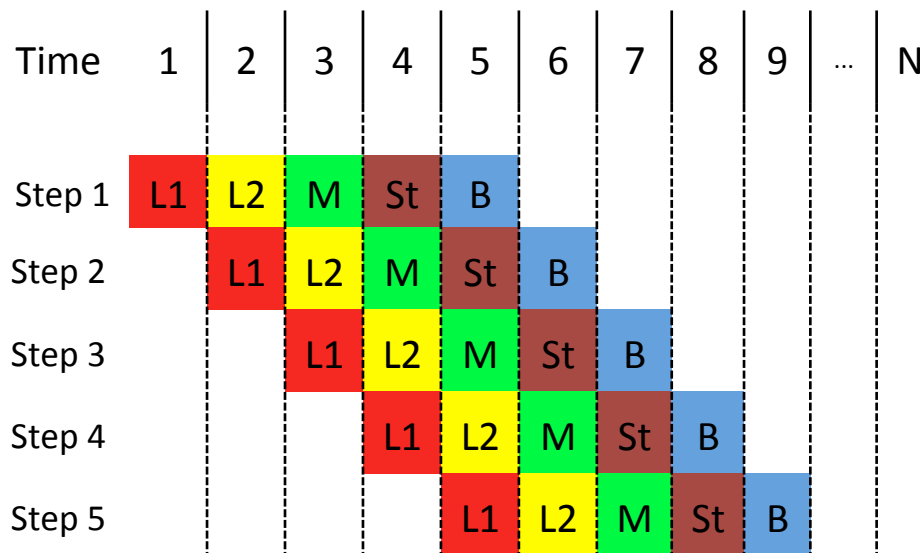
# Sequential Processing

# Pipelining

| Clock Period | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... | N |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Instruction 1 | F | D | L | E | S | | | | | | |
| Instruction 2 | | F | D | L | E | S | | | | | |
| Instruction 3 | | | F | D | L | E | S | | | | |
| Instruction 4 | | | | F | D | L | E | S | | | |
| Instruction 5 | | | | | F | D | L | E | S | | |

Ivan Girotto
igirotto@ictp.it

# Pipelining

| Clock Period | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... | N |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Instruction 1 | F | D | L | E | S | | | | | | |
| Instruction 2 | | F | D | L | E | S | | | | | |
| Instruction 3 | | | F | D | L | E | S | | | | |
| Instruction 4 | | | | F | D | L | E | S | | | |
| Instruction 5 | | | | | F | D | L | E | S | | |

# Superscalaring

| Clock Period | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... | N |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Instruction 1 | F | D | L | E | S | | | | | | |
| Instruction 2 | F | D | L | E | S | | | | | | |
| Instruction 3 | | F | D | L | E | S | | | | | |
| Instruction 4 | | F | D | L | E | S | | | | | |
| Instruction 5 | | | F | D | L | E | S | | | | |
| Instruction 6 | | | F | D | L | E | S | | | | |

# Loops and Pipeline

```
for( i = 0; i < N; i += 1 )
{
    A[i] = s * A[i]
}
```

| Time | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | … | N |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Step 1 | L1 | L2 | M | St | B | | | | | | |
| Step 2 | | L1 | L2 | M | St | B | | | | | |
| Step 3 | | | L1 | L2 | M | St | B | | | | |
| Step 4 | | | | L1 | L2 | M | St | B | | | |
| Step 5 | | | | | L1 | L2 | M | St | B | | |

```
Loop:   load r1, A(i)
        load r2, s
        mult r3, r2, r1
        store A(i), r3
        branch => loop
```

Ivan Girotto
igirotto@ictp.it

# The CPU Memory Hierarchy



**CPU Registers**

**CACHE**

**MAIN MEMORY**

**COMPUTATION**

**APPLICATION DATA**

# Cache Memory

- Expensive (SRAM) high-speed memory

- Relatively low-capacity in regards to RAM

- Cache Memory are for Instructions (i.e., L1I) and for Data (i.e., L1D)

- Modern CPU are designed with several levels of cache memories

# Cache Memory

```
Loop: load r1, A(i)
      load r2, s
      mult r3, r2, r1
      store A(i), r2
      branch => loop
```
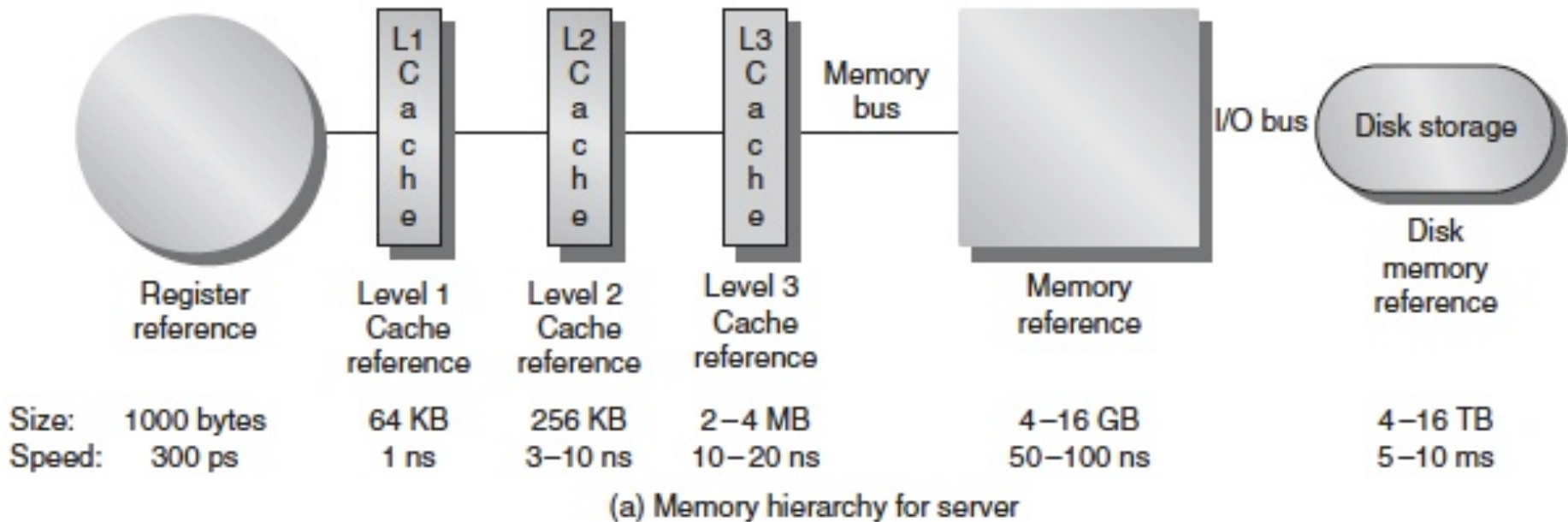
**CPU Registers**

**CACHE**

**MAIN MEMORY**

- Designed for temporal/spatial locality

- Data is transferred to cache in blocks of fixed size, called *cache lines*.

- Operation of LOAD/STORE can lead at two different scenario:
  - *cache hit*
  - *cache miss*

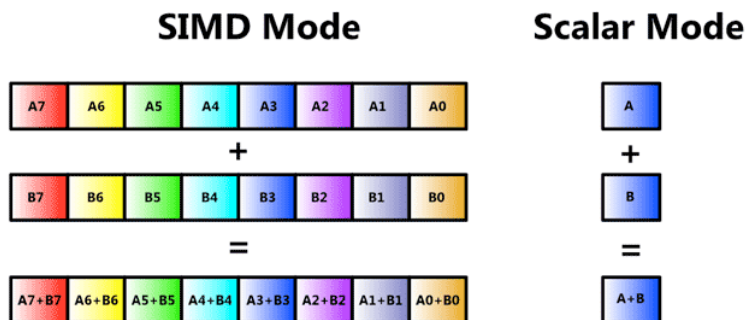# Caches

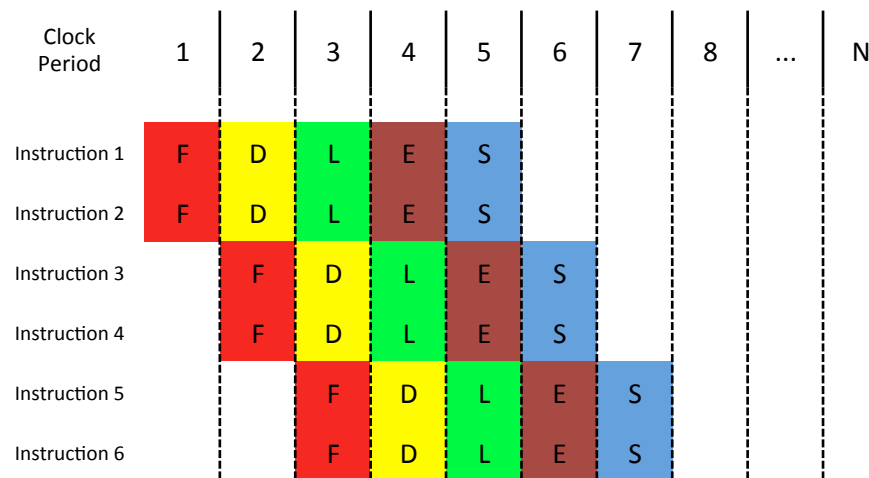- Fast memory to exploit spatial and temporal locality!

# The CPU Memory Hierarchy



(a) Memory hierarchy for server

# HPC Trend and Moore's Law

# To the Extreme - Parallel Inside



**SIMD Mode**     **Scalar Mode**

Vector Units for processing multiple data in //



| Clock Period | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ... | N |
|---|---|---|---|---|---|---|---|---|---|---|
| Instruction 1 | F | D | L | E | S | | | | | |
| Instruction 2 | F | D | L | E | S | | | | | |
| Instruction 3 | | F | D | L | E | S | | | | |
| Instruction 4 | | F | D | L | E | S | | | | |
| Instruction 5 | | | F | D | L | E | S | | | |
| Instruction 6 | | | F | D | L | E | S | | | |

Pipelined/Superscalar design: multiple functional units operate concurrently

# Few basic rules for optimized codes

- Do less work!!
  - Elimination of common sub-expressions
- Avoid expensive operations
  - Reduce your math to cheap operations
  - Avoid branches
- Think as a the compiler works
  - Enhance the compiler
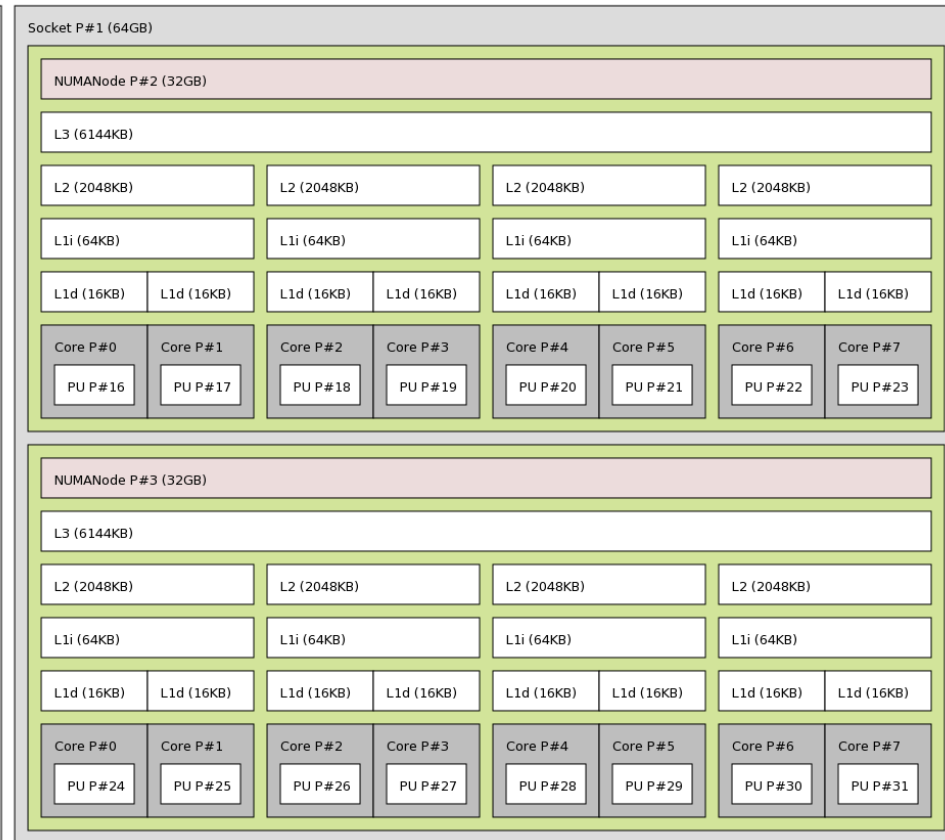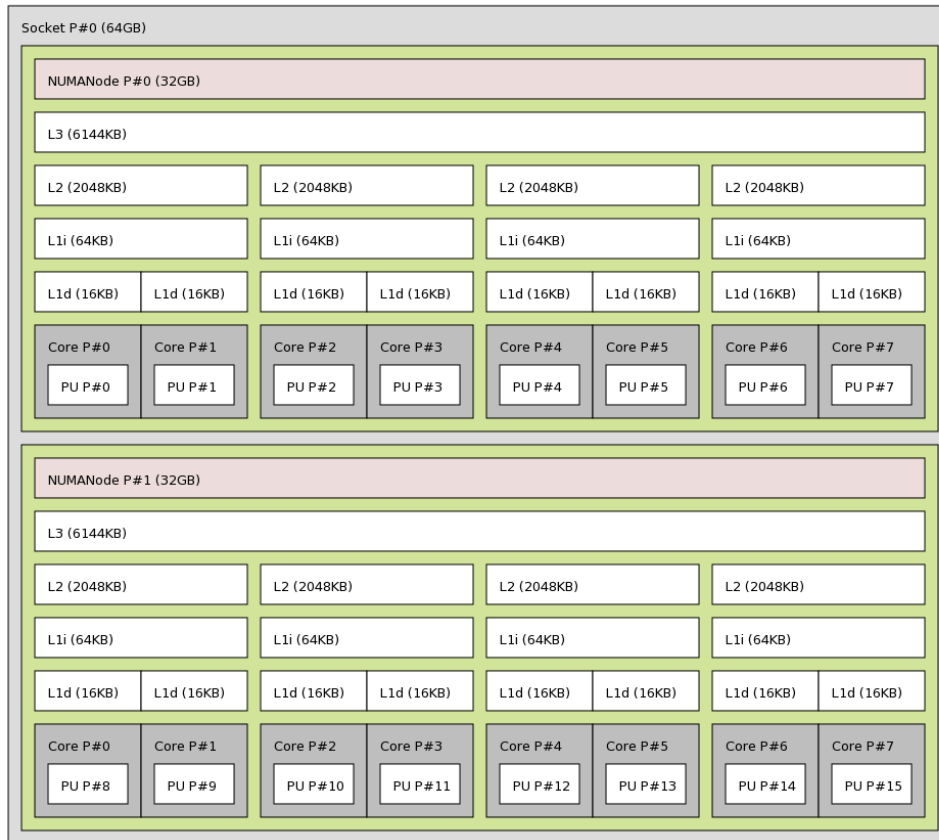
# Symmetric Multiprocessors (SMP)
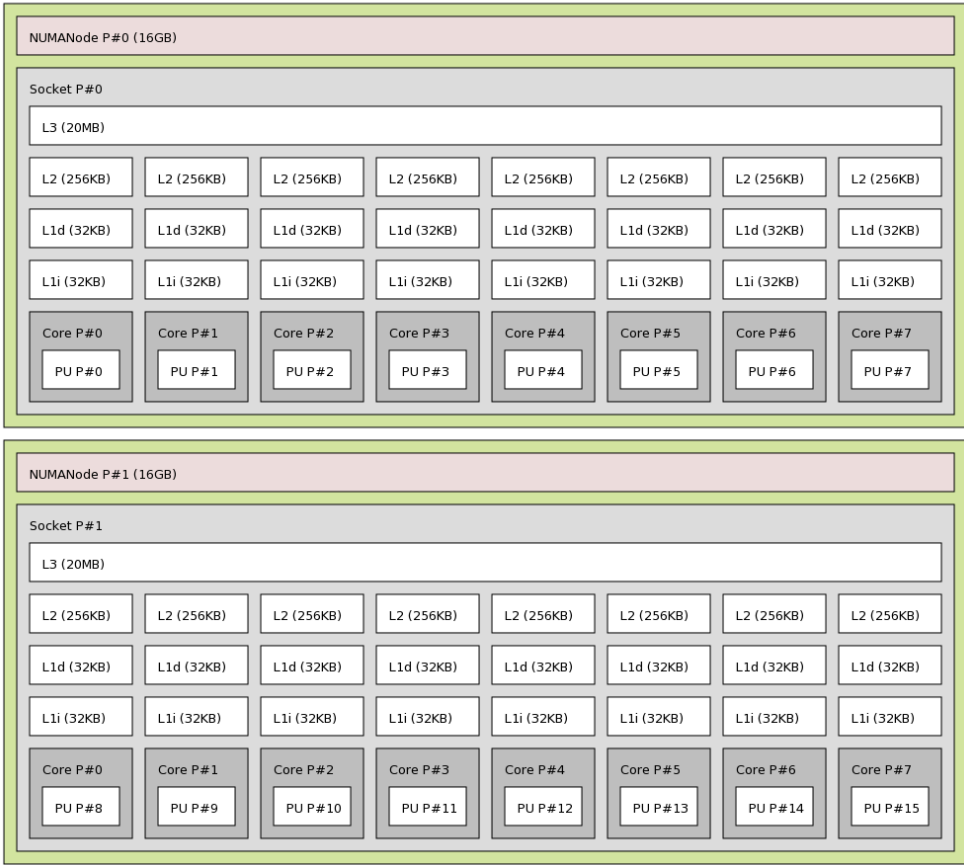
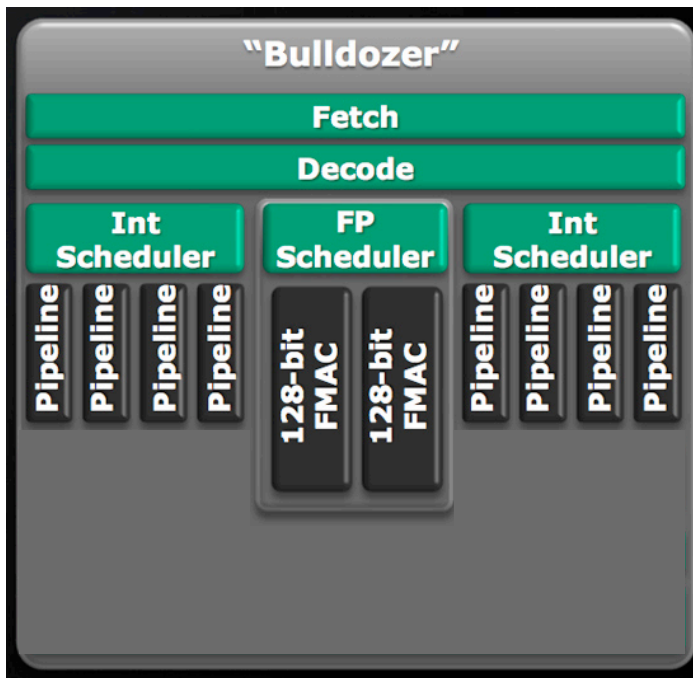# Modern NUMA Multicores

# The AMD Opteron 6380 Abu Dhabi 2.5GHz
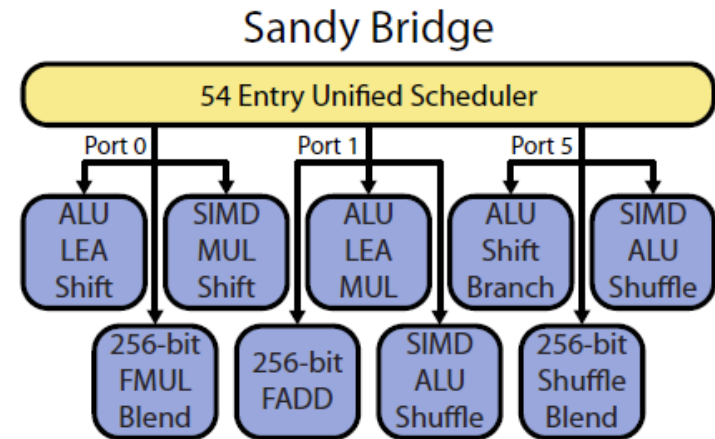
# The Intel Xeon E5-2665 Sandy Bridge-EP 2.4GHz

# State of the art

- AMD



- Intel

# Threading and Vectorization