



The Abdus Salam
International Centre
for Theoretical Physics



Overview on modern computer architectures: the software crisis

Ivan Girotto – igirotto@ictp.it

Information & Communication Technology Section (ICTS)

International Centre for Theoretical Physics (ICTP)

How fast is my CPU?!

- CPU power is measured in FLOPS
 - number of floating point operations x second
 - $\text{FLOPS} = \# \text{cores} \times \text{clock} \times \frac{\text{FLOP}}{\text{cycle}}$
- FLOP/cycle is the number of multiply-add (FMA) performed per cycle
 - architectural limit
 - depend also by the instruction set supported

The CPU Memory Hierarchy



The diagram illustrates the CPU memory hierarchy as a pyramid with three levels. The top level is a teal triangle labeled 'CPU Registers'. The middle level is a red trapezoid labeled 'CACHE'. The bottom level is a dark blue trapezoid labeled 'MAIN MEMORY'. To the right of the pyramid, there is a teal horizontal bar labeled 'COMPUTATION' at the top and a yellow horizontal bar labeled 'APPLICATION DATA' at the bottom. A large red double-headed arrow with diagonal hatching connects the 'COMPUTATION' bar to the 'APPLICATION DATA' bar, indicating bidirectional data flow.

CPU
Registers

CACHE

MAIN MEMORY

COMPUTATION

APPLICATION DATA

Cache Memory

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

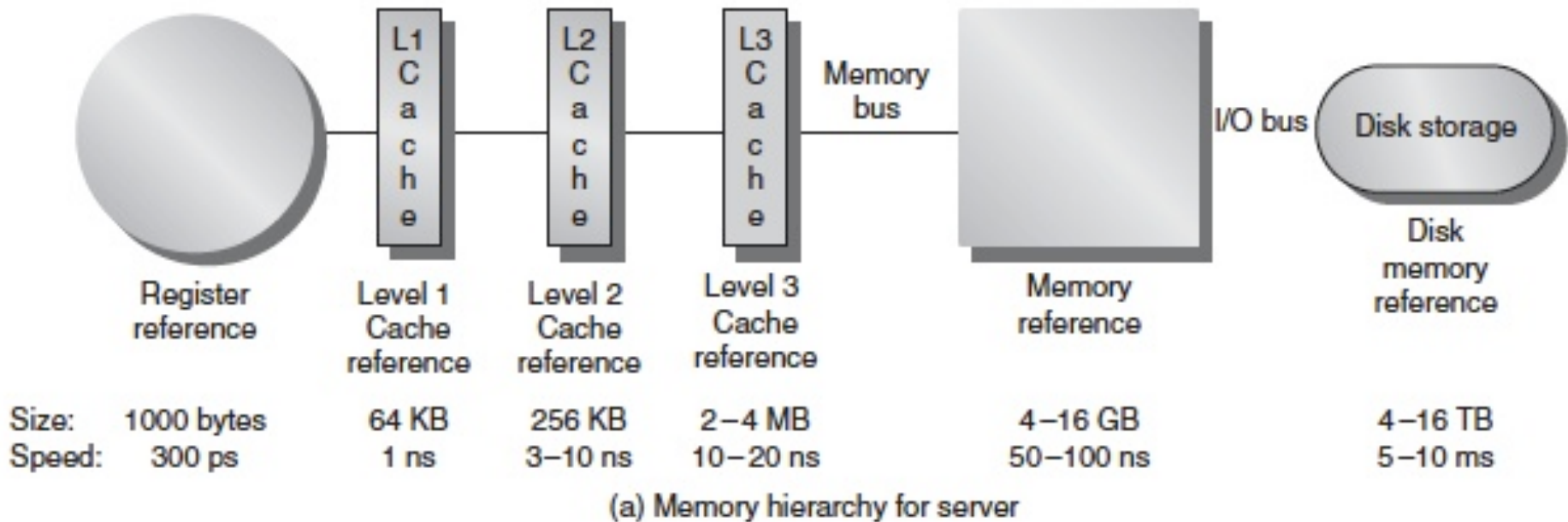
CPU
Registers

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*

The CPU Memory Hierarchy



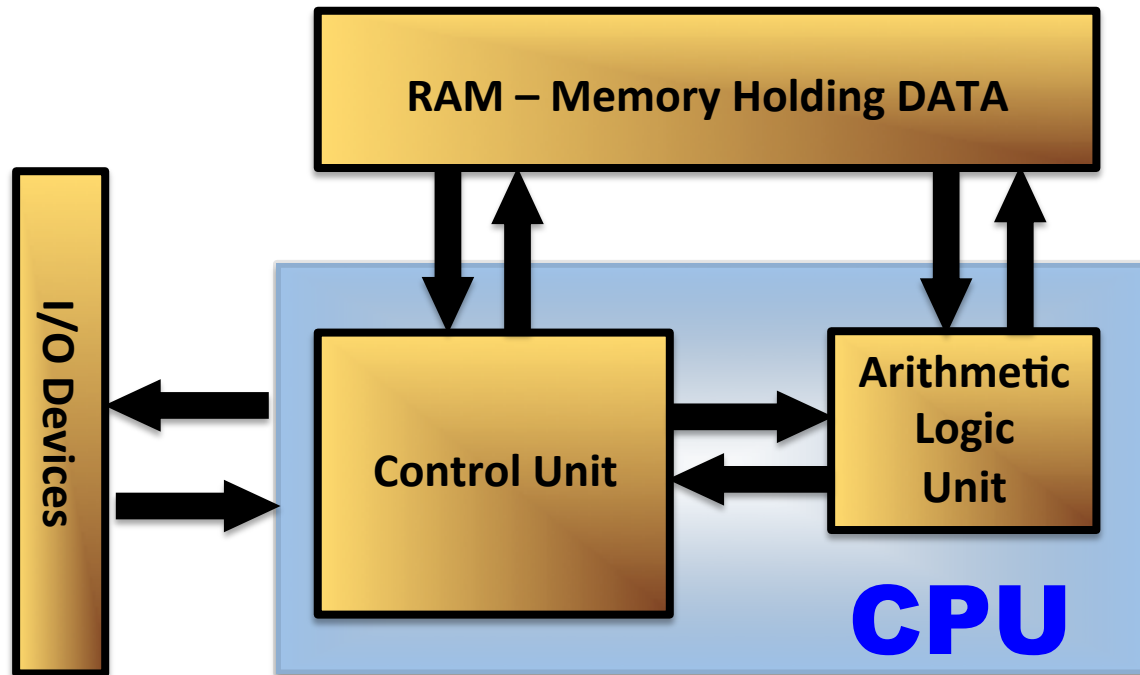
Data Memory Access

- Data ordering
- Reduce at minimum the data transfers
- Avoid complex data structure within computational intensive kernels
- Define constants and help the compiler
- Exploit the memory hierarchy



John Von Neumann

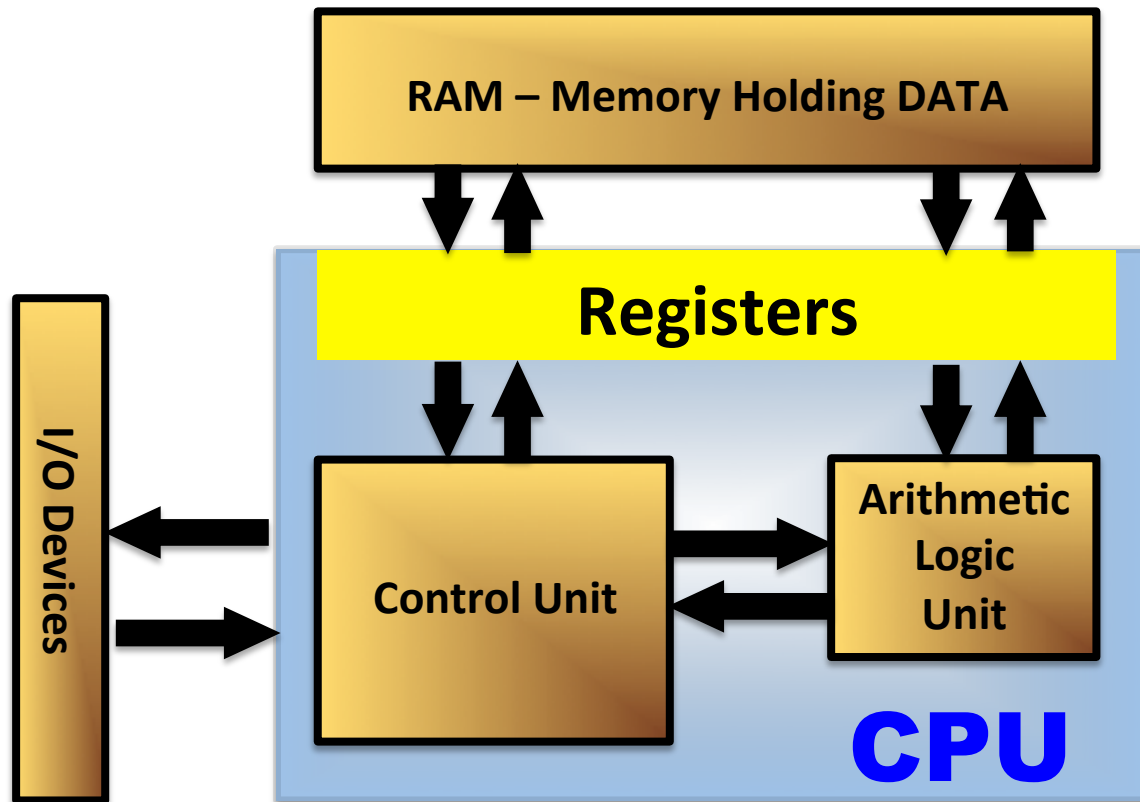
The Classical Model



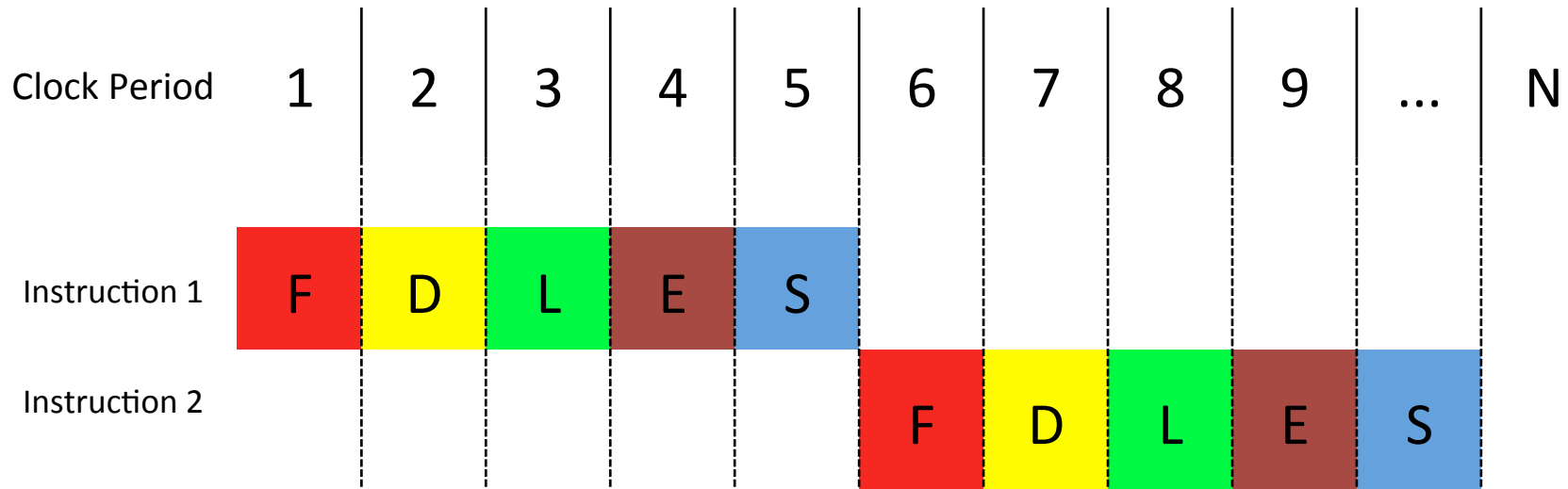


John Von Neumann

The Classical Model



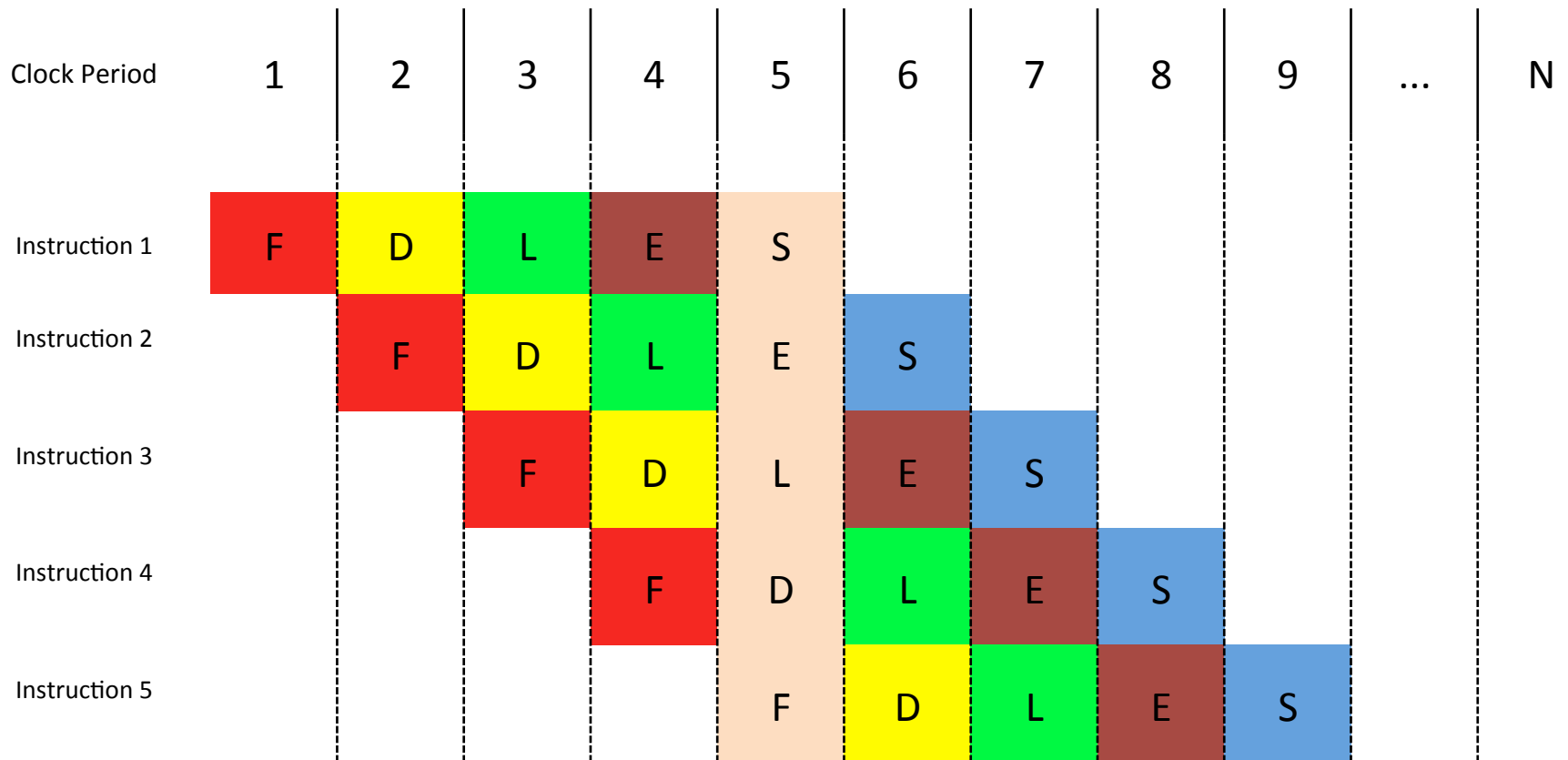
Sequential Processing



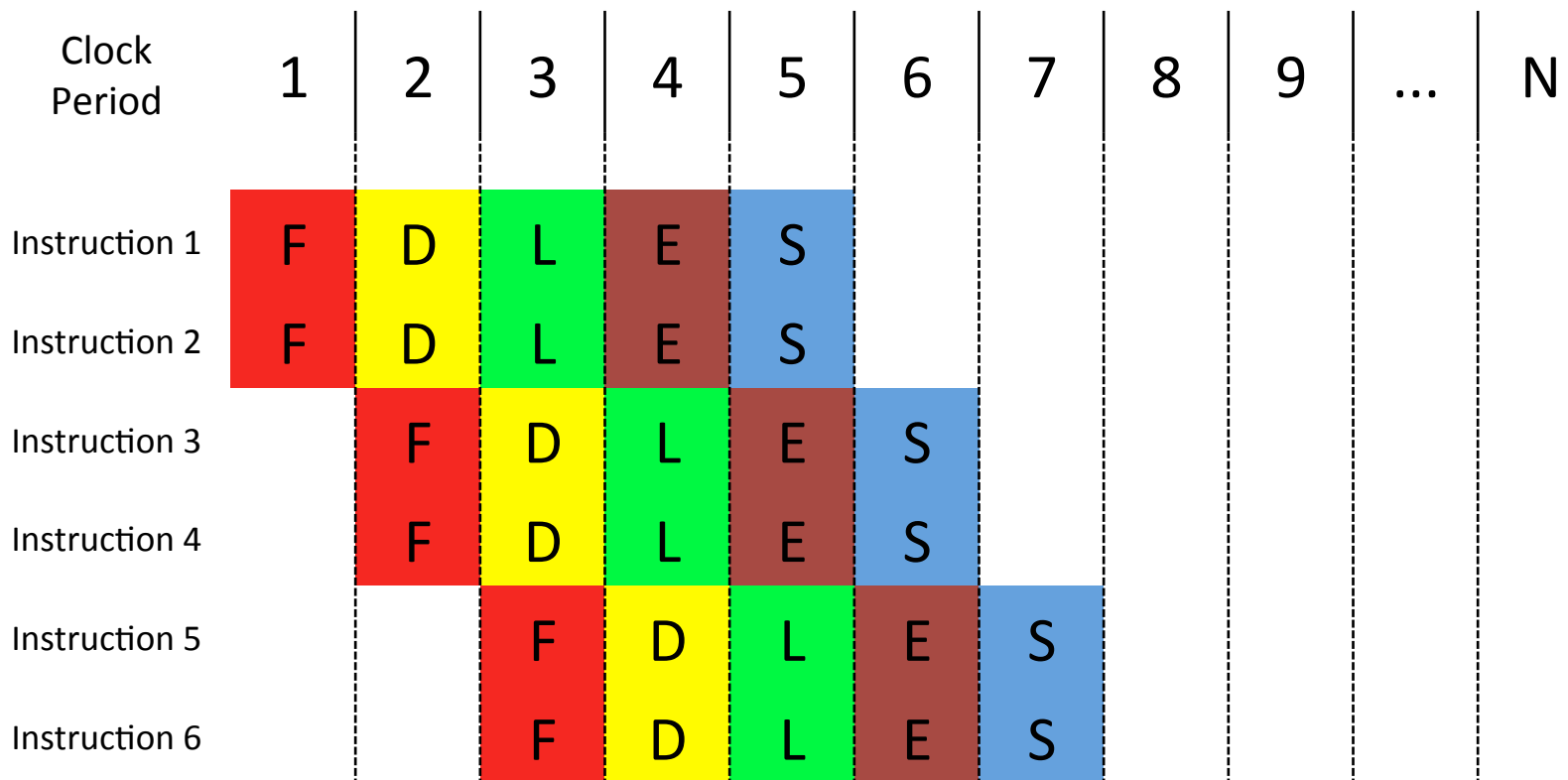
Pipelining



Pipelining

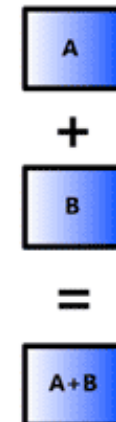


Superscalarizing



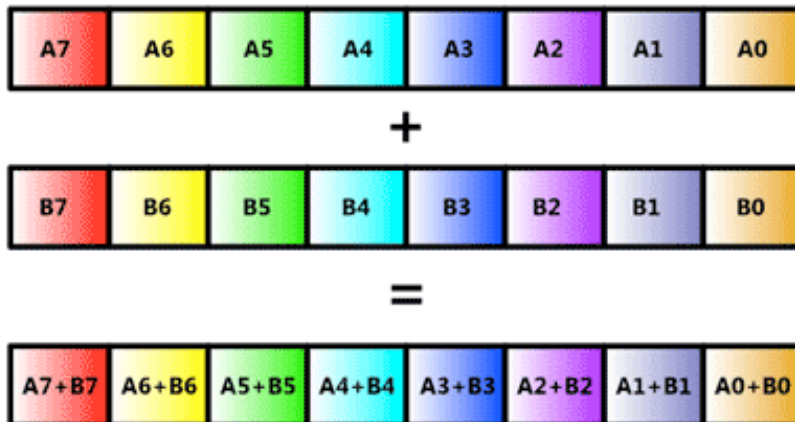
The Inside Parallelism

Scalar Mode

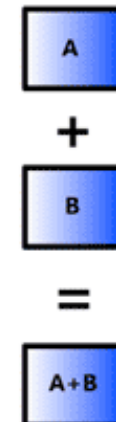


The Inside Parallelism

SIMD Mode



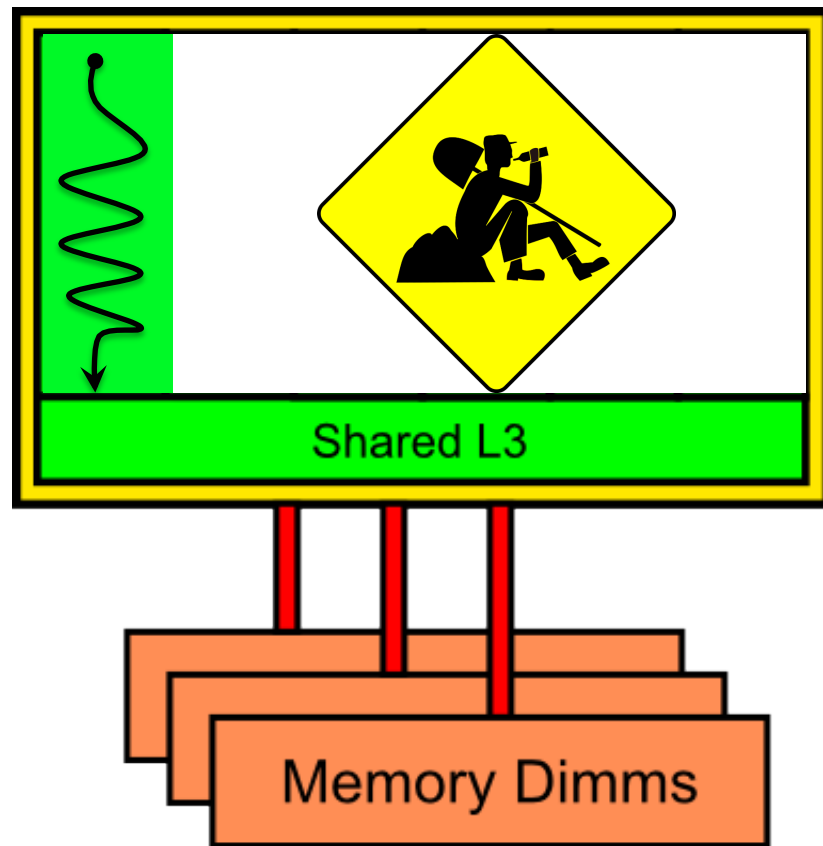
Scalar Mode



Performance Metrics

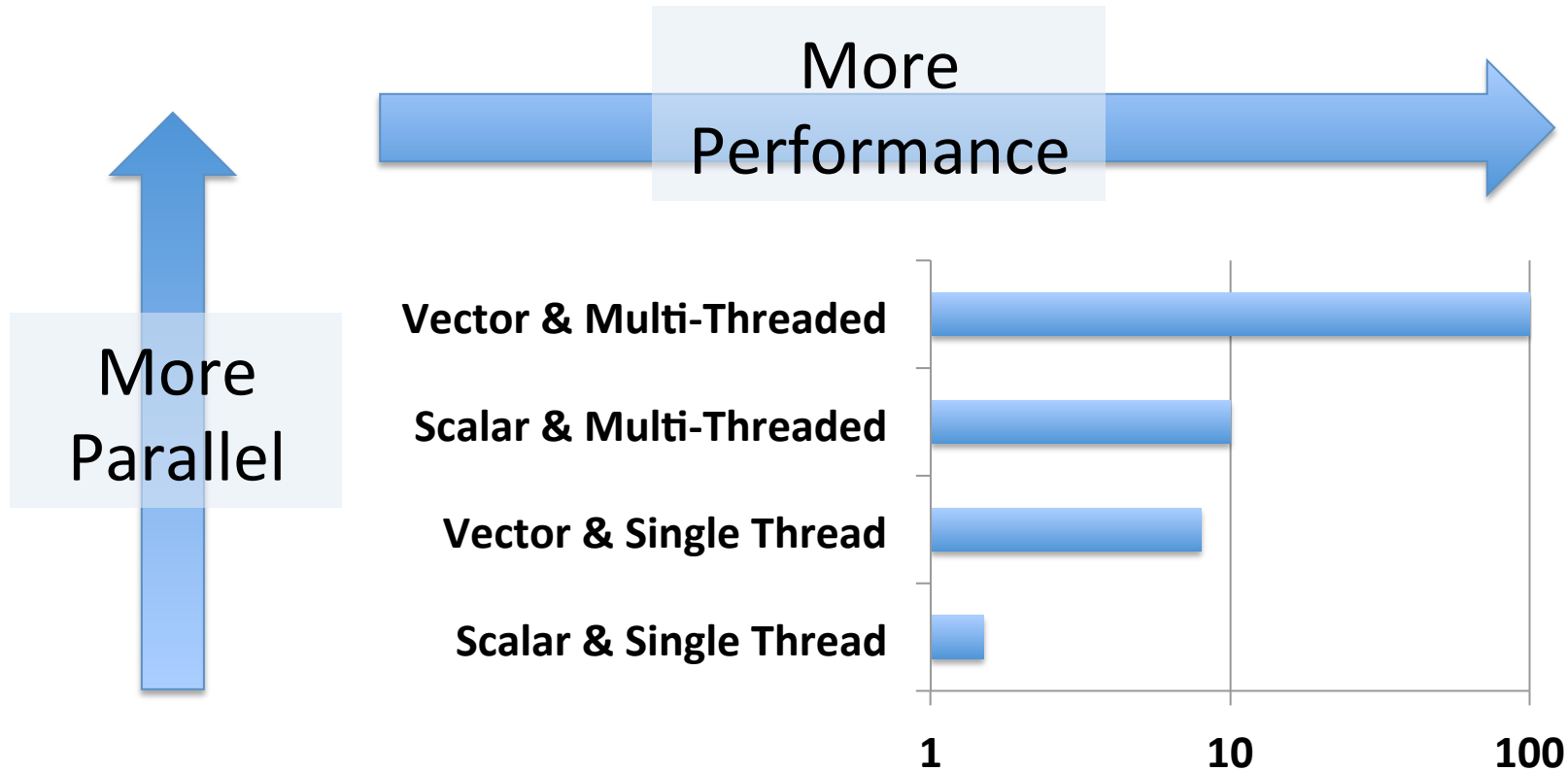
- When all CPU component work at maximum speed that is called *peak of performance*
 - Tech-spec normally describe the theoretical peak
 - Benchmarks measure the real peak
 - Applications show the real performance value
- CPU performance is measured as:
 - Floating point operations per seconds GFLOP/s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

Multi-core system Vs Serial Programming



Xeon E5650
hex-core
processors
(12GB - RAM)

Threading and Vectorization



Aid Automatic Vectorization

- Avoid Data Dependences
- Data Alignment (memalign)
- Aliasing
- Avoid Conditional Statements
- Analyze compiler report
 - `-qopt-report=5 -qopt-report-phase=vec`
- Add Compiler options:
 - `-axMIC-AVX512 (KNL)` `-march=core-avx2 (BWD)`

Conclusion

- Development of today computer architecture is based on increasing complexity: **the software crisis!!!**
- Scientific software developers for high-performance computing need to master this complexity:
 - avoid useless instructions, branches (if possible) and expensive OP
 - enhance data locality and data reuse (memory throughput)
 - aid the compiler for automatic vectorization
 - use compiler optimization (-O3)
 - make use of computer libraries for HPC