

Node-Level Performance Engineering

Jan Eitzinger

Two-day Tutorial

GWDG, Göttingen

November 19-20, 2019

<http://tiny.cc/NLPE-GWDG>

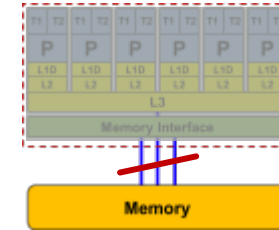
Day 1		
10:00	Welcome – Intro	JE
10:15	Computer architecture for software developers	JE
11:30	Performance engineering fundamentals	JE
12:00	Lunch	
13:00	Tools: Topology and affinity, frequency	JE
13:30	Exercise 1: Microarchitecture exploration	JE
14:30	Optimal use of parallel resources: SIMD	JE
16:00	Roofline Model: Basic	JE
17:00	End of day 1	
Day 2		
9:00	Tools: Performance counters	JE
9:45	Optimal use of parallel resources: ccNUMA	JE
10:45	Exercise 2: Dense Matrix Vector Multiplication	JE
12:00	Lunch	
13:00	Performance Engineering: Basic skills	JE
14:00	Roofline case studies: Jacobi smoother + SpMVM	JE
16:00	Exercise 3: MiniMD Analysis	JE
17:00	End of day 2	

- What does “clock frequency” mean in computers?

The “heartbeat” of the CPU. A clock cycle is the smallest unit of time on a CPU chip.
Typically $< 1\text{ns}$ $\rightarrow f \gtrsim 1\text{GHz}$

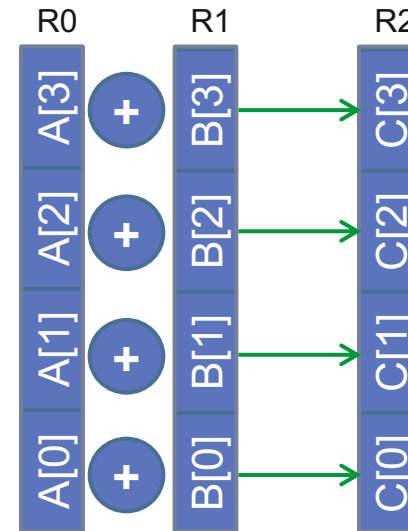
- What is “memory bandwidth”?

Rate of data transfer between main memory (RAM) and CPU chip. Typical $b_S \approx 10 \dots 100\text{GB/s}$

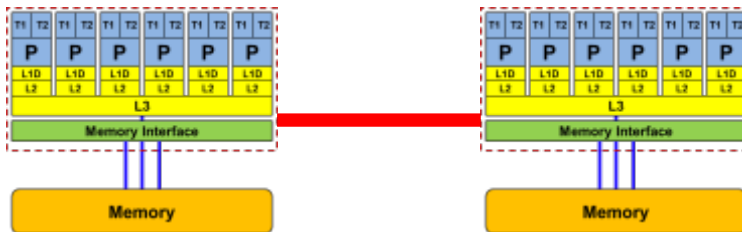


- What is SIMD vectorization?

Single Instruction **M**ultiple **D**ata.
Data-parallel load/store and execution units.



- What is ccNUMA?



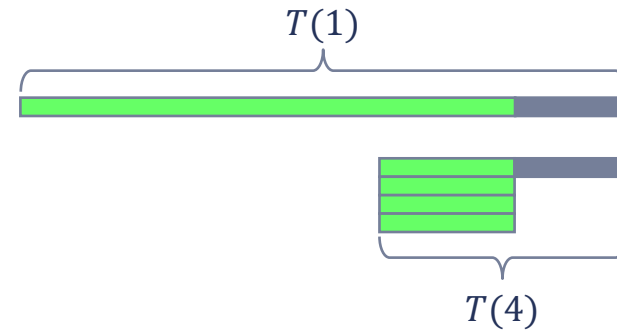
- What is a register?

A storage unit in the CPU core that can take one single value (a few values in case of SIMD). Operands for computations reside in registers.

rax	ymm0
rbx	ymm1
rcx	ymm2
rdx	ymm3
rsi	ymm4

- What is Amdahl's Law?

$$S_p = \frac{T(1)}{T(N)} = \frac{1}{s + \frac{1-s}{N}}$$



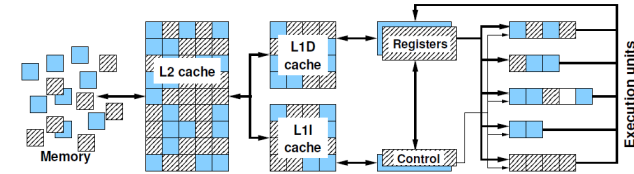
- What is a pipelined functional unit?

An instruction execution unit on the core that executes a certain task in several simple sub-steps. The stages of the pipeline can act in parallel on several instructions at once.



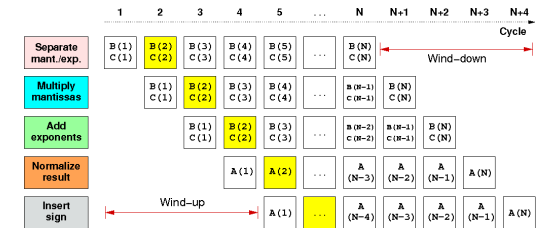
■ What is SMT?

Simultaneous Multi-Threading, a.k.a. hyper-threading. A CPU core can execute multiple threads concurrently. Such threads share all execution resources except the registers.



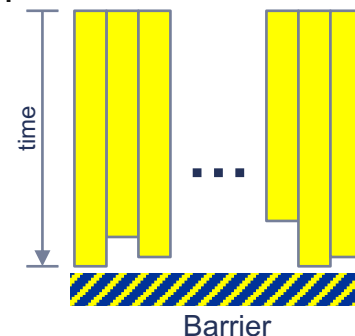
■ How long does a single double-precision floating-point multiplication take?

3-5 clock cycles, depending on the CPU



■ How long does an OpenMP barrier take?

Depending on the OpenMP runtime, the number of cores or threads, and where they are running, a barrier can take tens of thousands of cycles (typically a few 1000s on a full socket).



- 1 cycle = smallest unit of time on a CPU (“heartbeat”)
 - Clock speed of typical CPU: **2.4 Gcy/s (or GHz)**
- Basic unit of work: Floating-point operation (Flop)
 - Typical peak performance of 20-core CPU: **$P_{\text{peak}} = 768 \text{ Gflop/s}$**
 - How many Flops per cycle per core is that? $\frac{768 \cdot 10^9 \frac{\text{Flops}}{\text{s}}}{20 \text{ cores} \cdot 2.4 \cdot 10^9 \frac{\text{cy}}{\text{s}}} = 16 \frac{\text{Flops}}{\text{cy} \cdot \text{core}}$
 - Typical **duration** of a double precision **multiply: 5 cycles**
 - › How much time is that? $\frac{5 \text{ cy}}{2.4 \cdot 10^9 \frac{\text{cy}}{\text{s}}} = 2.08 \cdot 10^{-9} \text{ s} = 2.08 \text{ ns}$
- Basic unit of traffic: **Byte**
- Unit of bandwidth: **Bytes/s**
 - Typical memory bandwidth: **$62 \text{ Gbytes/s} = 6.2 \cdot 10^{10} \text{ Bytes/s}$**
 - How many bytes per cycle is that? $\frac{62 \cdot 10^9 \frac{\text{Bytes}}{\text{s}}}{2.4 \cdot 10^9 \frac{\text{cy}}{\text{s}}} = 25.8 \frac{\text{Bytes}}{\text{cy}}$