

# HIGH PERFORMANCE Co<sup>mpu</sup>ting<sub>ding</sub> (HPC)

ALBERTO DASSATTI - 2023



## Performance

an hard definition... for us a **measurable** quantity of interest that we can exploit for comparisons and that has some connection with a figure of merit of a system under test.

# MANY PERFORMANCES

- ▶ Development Time
- ▶ Resources requested
- ▶ Resources utilization
- ▶ Market Shares
- ▶ Maintainability
- ▶ ...

end very generally applicable

- ▶ Cost

# HPC INTERESTING PERFORMANCES

All these performances are Computer Architecture metrics used to compare different machines, here we will use them in a different way: given one machine, we will compare different codes solving the same problem and we will derive from these numbers knowledge and optimization opportunities.

- ▶ Instruction Level Parallelism
- ▶ Cache behaviour
- ▶ Latency
- ▶ Throughput
- ▶ Power

Warning: they are not independent! Trade-offs will become your daily job.

## Instruction Level Parallelism

Theory: is a measure of how many of the operations in a computer program can be performed simultaneously

Practice: is a measure of how many of the operations in a computer program are performed simultaneously

2.066862	task-clock (msec)	#	0.746 CPUs utilized
7	context-switches	#	0.003 M/sec
0	cpu-migrations	#	0.000 K/sec
105	page-faults	#	0.051 M/sec
1,825,561	cycles	#	0.883 GHz
1,167,832	stalled-cycles-frontend	#	63.97% frontend cycles idle
0	stalled-cycles-backend	#	0.00% backend cycles idle
1,424,315	instructions	#	0.78 insns per cycle
		#	0.82 stalled cycles per insn
294,760	branches	#	142.612 M/sec
13,046	branch-misses	#	4.43% of all branches

## Instruction Level Parallelism

Theory: is a measure of how many of the operations in a computer program can be performed simultaneously

Practice: is a measure of how many of the operations in a computer program are performed simultaneously

2.066862	task-clock (msec)	#	0.746 CPUs utilized
7	context-switches	#	0.003 M/sec
0	cpu-migrations	#	0.000 K/sec
105	page-faults	#	0.051 M/sec
1,825,561	cycles	#	0.883 GHz
1,167,832	stalled-cycles-frontend	#	63.97% frontend cycles idle
0	stalled-cycles-backend	#	0.00% backend cycles idle
1,424,315	instructions	#	0.78 insns per cycle
		#	0.82 stalled cycles per insn
294,760	branches	#	142.612 M/sec
13,046	branch-misses	#	4.43% of all branches

## Cache

The [hit/miss](#) ratio is very interesting.

We still have no information about what kind of miss we are facing.

28,973	L1-dcache-load-misses	
9,457	L1-dcache-store-misses	
6,106	L1-dcache-prefetch-misses	
31,855	L1-icache-load-misses	
1,387,200	instructions	# 0.81 insns per cycle
27,840	cache-references	
7,740	cache-misses	# 27.802 % of all cache refs
1,720,732	cycles	

## Cache

The [hit/miss](#) ratio is very interesting.

We still have no information about what kind of miss we are facing.

28,973	L1-dcache-load-misses	
9,457	L1-dcache-store-misses	
6,106	L1-dcache-prefetch-misses	
31,855	L1-icache-load-misses	
1,387,200	instructions	# 0.81 insns per cycle
27,840	cache-references	
7,740	cache-misses	# 27.802 % of all cache refs
1,720,732	cycles	



## Latency

Time needed to complete a task.

How to reduce it:

- ▶ change architecture
- ▶ change algorithm/implementation

## Throughput

Number of task done in parallel.

How to increase it:

- ▶ parallel execution (Pipeline, ILP, Multi-core, etc)
- ▶ it has specific obstacles we will address later in this class

# POWER

Reducing power consumption is the real motivation of HPC. Why?

# POWER

Reducing power consumption is the real motivation of HPC. Why?

$$P = \frac{1}{2}CV^2f + P_{leak}$$

# CPU C-States

	C0 HFM	C0 LFM	C1/C2	C4	C6
Core voltage					
Core clock			OFF	OFF	OFF
PLL				OFF	OFF
L1 caches			 flushed	 flushed	 off
L2 caches				 Partial flush	 off
Wakeup time	active	active			
Power					

## WHERE IS THE LIMIT

### memory bound

The task execution is limited by the amount or access time to the memory

## WHERE IS THE LIMIT

### memory bound

The task execution is limited by the amount or access time to the memory

### CPU bound

The task execution is limited by the number of instructions the system performs per second

# WHERE IS THE LIMIT

## memory bound

The task execution is limited by the amount or access time to the memory

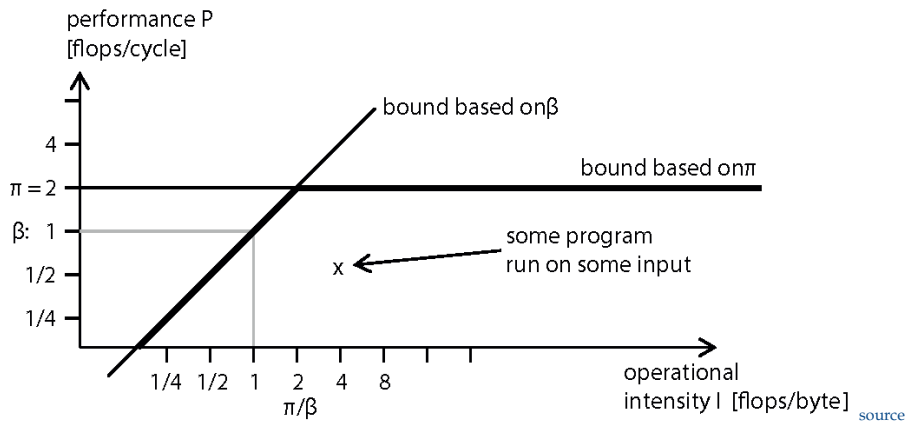
## CPU bound

The task execution is limited by the number of instructions the system performs per second

## I/O Bound

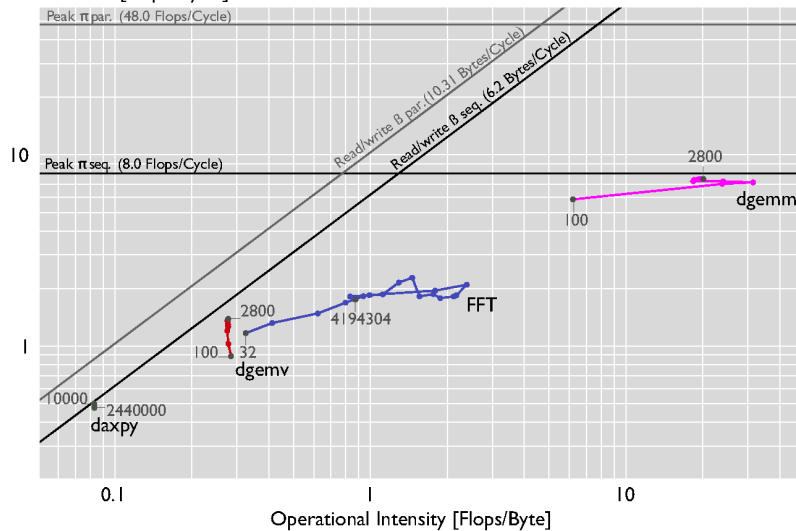
The task execution is limited by the amount or access time to external devices

# ROOFLINE MODEL





# Performance [Flops/Cycle]



source

# QUESTIONS

