

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

ALBERTO DASSATTI - 2020



An instruction set, or instruction set architecture (ISA), is the part of the computer architecture related to programming, including the native data types, instructions, **registers**, addressing modes, memory architecture, interrupt and **exception handling**, and external I/O.

**Instruction set - Wikipedia, the free encyclopedia**  
[https://en.wikipedia.org/wiki/Instruction\\_set](https://en.wikipedia.org/wiki/Instruction_set)

## ISA

All information needed by the programmer for programming an architecture

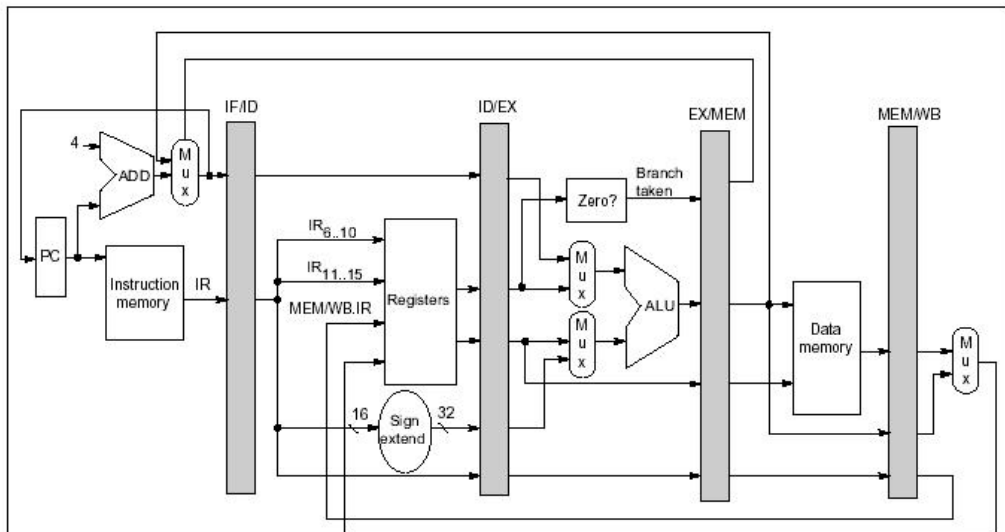
## Microarchitecture

The real hardware, CPU, memory and so on, giving (an) implementation(s) of an ISA

In computer engineering, **microarchitecture** (sometime abbreviated to *μarch* or *uarch*) is a description of the electrical circuitry of a computer, central processing unit, or digital signal processor that is sufficient for completely describing the operation of the hardware.

**Microarchitecture - Simple English Wikipedia, the free ...**  
<https://simple.wikipedia.org/wiki/Microarchitecture>

## DLX



Complex instruction set computing (**CISC** /'sɪsk/) is a processor design, where single instructions can execute several low-level operations (such as a load from memory, an arithmetic operation, and a memory store) or are capable of multi-step operations or addressing modes within single instructions.

**Complex instruction set computing - Wikipedia**

[https://en.wikipedia.org/wiki/Complex\\_instruction\\_set\\_computing](https://en.wikipedia.org/wiki/Complex_instruction_set_computing)

## CISC

PDP-11, VAX, Motorola 68k, and your desktop PCs on intel's x86 architecture.

## RISC

RISC families include DEC Alpha, AMD Am29000, ARC, ARM, Atmel AVR, Blackfin, Intel i860 and i960, MIPS, Motorola 88000, PA-RISC, Power (including PowerPC), RISC-V, etc...

Reduced instruction set computing, or **RISC** (pronounced 'risk', /rɪsk/), is a CPU design strategy based on the insight that a simplified instruction set provides higher performance when combined with a microprocessor architecture capable of executing those instructions using fewer microprocessor cycles per instruction.

**Reduced instruction set computing - Wikipedia**

[https://en.wikipedia.org/wiki/Reduced\\_instruction\\_set\\_computing](https://en.wikipedia.org/wiki/Reduced_instruction_set_computing)



# Q&A

1. Why RISC has been so succesful?

## Q&A

1. Why RISC has been so succesful?
2. What kind of ISA is X86?

## Q&A

1. Why RISC has been so succesful?
2. What kind of ISA is X86?
3. Why Intel was so succesful?

## Q&A

1. Why RISC has been so succesful?
2. What kind of ISA is X86?
3. Why Intel was so succesful?
4. How a CISC ISA is implemented today?



## SOME MORE RECENT EXAMPLES

Let's have a look at some modern microarchitectures:  
Intel [Nehalem Wikipedia](#), and [more](#) or  
AMD [Bulldozer](#) and [more](#)

# NEHALEM

1. 6 Ops: 3 mem (load, store data and store address) + 3 ALU

# NEHALEM

1. 6 Ops: 3 mem (load, store data and store address) + 3 ALU
2. Out of Order dynamically scheduled

# NEHALEM

1. 6 Ops: 3 mem (load, store data and store address) + 3 ALU
2. Out of Order dynamically scheduled
3. Speculation: branch, pre-fetch (data-flow model)

# NEHALEM

1. 6 Ops: 3 mem (load, store data and store address) + 3 ALU
2. Out of Order dynamically scheduled
3. Speculation: branch, pre-fetch (data-flow model)
4. 20 – 24 stage pipeline

# NEHALEM

1. 6 Ops: 3 mem (load, store data and store address) + 3 ALU
2. Out of Order dynamically scheduled
3. Speculation: branch, pre-fetch (data-flow model)
4. 20 – 24 stage pipeline
5.  $128\mu Ops$

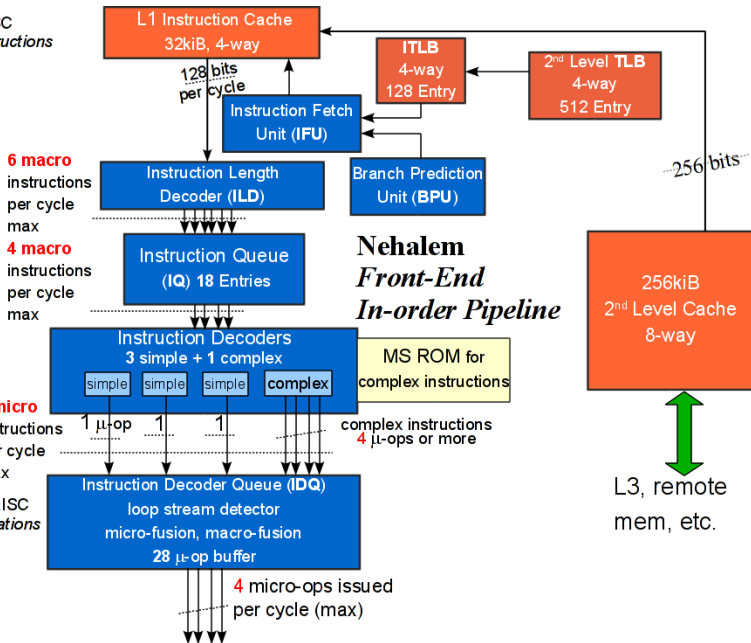
# NEHALEM

1. 6 Ops: 3 mem (load, store data and store address) + 3 ALU
2. Out of Order dynamically scheduled
3. Speculation: branch, pre-fetch (data-flow model)
4. 20 – 24 stage pipeline
5.  $128\mu Ops$
6. Very High and variable ILP

Intel64 CISC  
macro-instructions

macro-instruction decoding  
into micro-ops

Nehalem RISC  
micro-operations





Or another useful view for performance [arch](#)

## Q&A

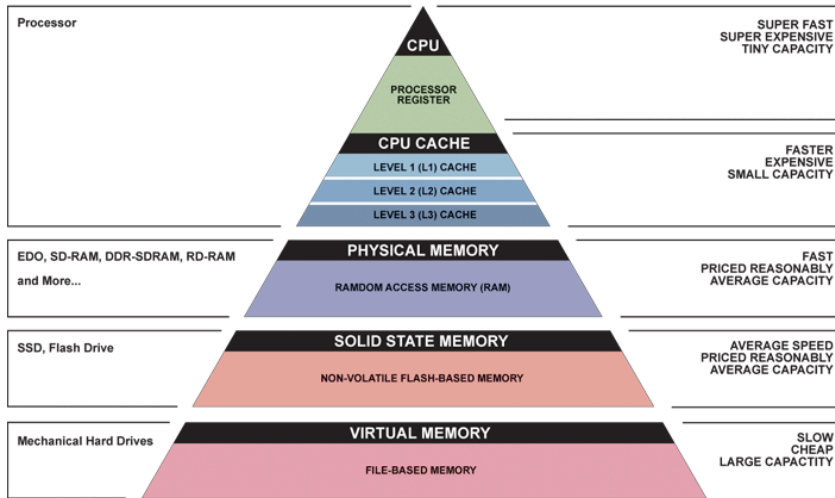
1. Why large  $\mu Ops$  buffer is important?

## Q&A

1. Why large  $\mu Ops$  buffer is important?
2. Why caches are so important in this architectures?

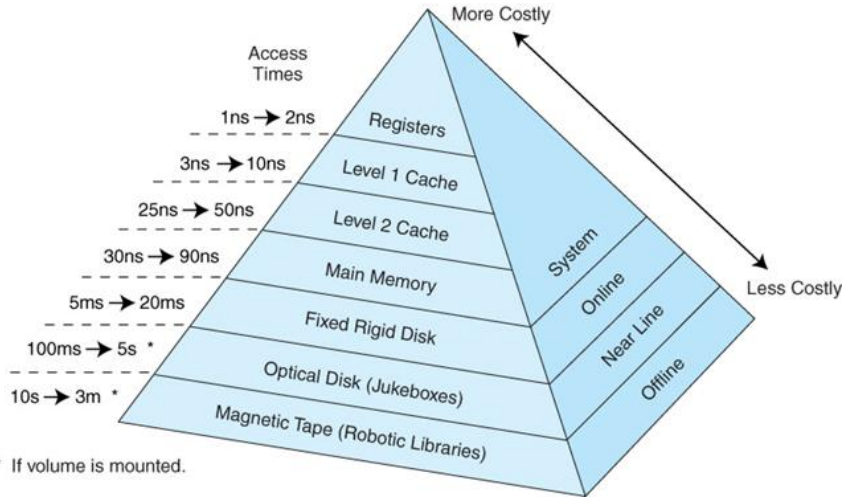
## Q&A

1. Why large  $\mu Ops$  buffer is important?
2. Why caches are so important in this architectures?
3. HiperThread technology is a compromise: when it makes sense?



▲ Simplified Computer Memory Hierarchy  
Illustration: Ryan J. Leng

source



source

# CACHE MEMORY

If you need a gentle basic introduction to cache memories see this [article](#)

## Cache Miss

- ▶ Coldness: first time the data is not there → Prefetching
- ▶ Capacity: too much data accessed → organize data and code differently
- ▶ Conflict: Multiple data mapped to the same location
- ▶ Sharing:
  - ▶ True: Thread in another processor wanted the data, it got moved to the other cache
  - ▶ False: data in the same cache line used by another core

## Latency Numbers Every Programmer Should Know

latency.txt

1	Latency Comparison Numbers					
2	-----					
3	L1 cache reference	0.5	ns			
4	Branch mispredict	5	ns			
5	L2 cache reference	7	ns			14x L1 cache
6	Mutex lock/unlock	25	ns			
7	Main memory reference	100	ns			20x L2 cache, 200x L1 cache
8	Compress 1K bytes with Zip	3,000	ns	3	us	
9	Send 1K bytes over 1 Gbps network	10,000	ns	10	us	
10	Read 4K randomly from SSD*	150,000	ns	150	us	~1GB/sec SSD
11	Read 1 MB sequentially from memory	250,000	ns	250	us	
12	Round trip within same datacenter	500,000	ns	500	us	
13	Read 1 MB sequentially from SSD*	1,000,000	ns	1,000	us	1 ms ~1GB/sec SSD, 4X memory
14	Disk seek	10,000,000	ns	10,000	us	10 ms 20x datacenter roundtrip
15	Read 1 MB sequentially from disk	20,000,000	ns	20,000	us	20 ms 80x memory, 20X SSD
16	Send packet CA->Netherlands->CA	150,000,000	ns	150,000	us	150 ms

### Notes

-----

1 ns = 10<sup>-9</sup> seconds

1 us = 10<sup>-6</sup> seconds = 1,000 ns

1 ms = 10<sup>-3</sup> seconds = 1,000 us = 1,000,000 ns

### Credit

-----

By Jeff Dean: <http://research.google.com/people/jeff/>

Originally by Peter Norvig: <http://norvig.com/21-days.html#answers>



# Q&A

1. When latency is critical?

## Q&A

1. When latency is critical?
2. When throughput is more relevant?

## Q&A

1. When latency is critical?
2. When throughput is more relevant?
3. What is the guiding principle to write performant code on these architectures?

## Q&A

1. When latency is critical?
2. When throughput is more relevant?
3. What is the guiding principle to write performant code on these architectures?
4. Where are caches introduced with respect of MMU? Why?

## EXTRA

If you are interested in the effects of branch predictors (and you should) have a look [here](#). Cache effects are very complicated to study and isolate. Read this [article](#) to have a better idea.