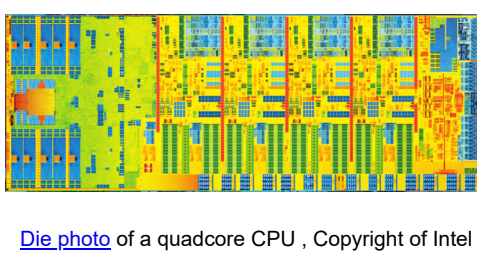


# MICROARCHITECTURE CHEAT SHEET

## X86 CPUs & Performance



Die photo of a quadcore CPU. Copyright of Intel

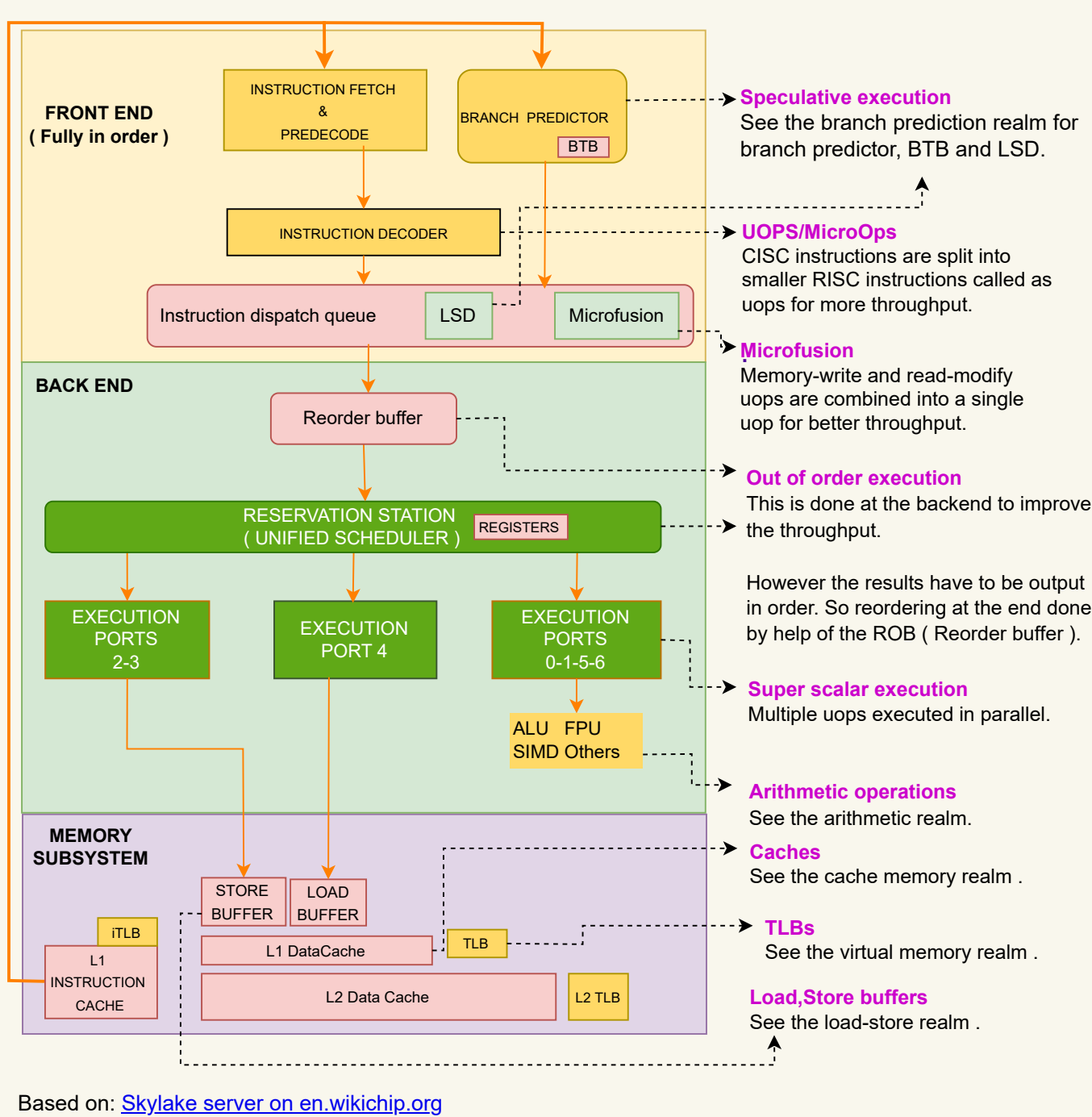
LAST UPDATE DATE : 17 OCT 2022

FOR LATEST VERSION: [www.github.com/ahin/microarchitecture-cheatsheet](https://github.com/ahin/microarchitecture-cheatsheet)

AUTHOR: AKIN OCAL [akin\\_ocal@hotmail.com](mailto:akin_ocal@hotmail.com) [https://twitter.com/akin\\_ocal\\_dev](https://twitter.com/akin_ocal_dev)

## PIPELINE REALM : INSIDE AN INDIVIDUAL CORE

### A SIMPLIFIED OVERVIEW (BASED ON INTEL SKYLAKE)



Based on: [Skylake server on en.wikichip.org](https://en.wikichip.org/wiki/skylake_server)

### PIPELINE PARALLELISM & PERFORMANCE

**Pipeline diagrams** : The diagrams below in the following topics are outputs from an online microarchitecture analysis tool [UCCA](https://en.wikichip.org/wiki/analysis), and they represent parallel execution through cycles.

Rows are multiple instructions being executed at the same time.

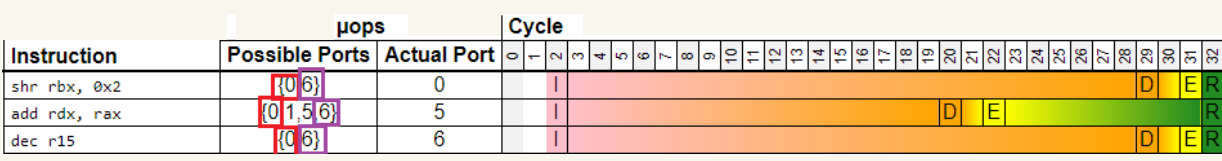
Columns display how instruction state changes through cycles.

**IPC** : As for pipeline performance, typically IPC is used. It stands for "Instructions per cycle". A higher IPC value usually means a better throughput.

You can measure IPC with perf: <https://perf.wiki.kernel.org/index.php/Tutorial>

**Rate of retired instructions** : Apart from IPC, number of retired instructions should be checked. Retired instructions are not committed/fetched as they were wrongly speculated. On the other hand, executed instructions are the ones which were finalised. Therefore a high rate of retired instructions indicates low branch prediction rate.

### CONTENTION FOR EXECUTION PORTS IN THE PIPELINE

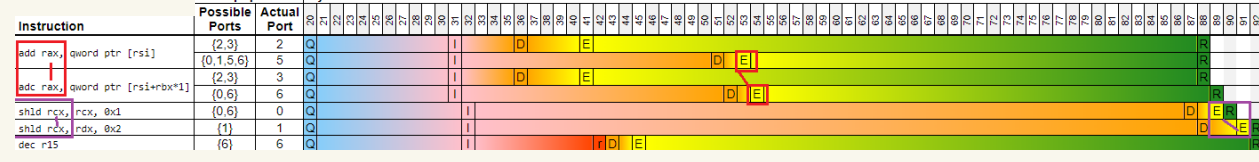


In the example above, all instructions are working on different registers, but SHR, ADD, DEC instructions are competing for ports 0 and 6. SHR and DEC are getting executed after ADD instruction.

Also notice that there is longer time between E[xecuted] and R[etired] states of instruction ADD as retirement has to be done in-order whereas execution is out-of-order.

Reference : [Denis Bashkatov's article](https://en.wikichip.org/wiki/skylake_server)

### INSTRUCTION STALLS DUE TO DATA DEPENDENCY



In the example above, there are 2 dependency chains, each marked with a different colour. In the first chain, 2 instructions are competing for RAX register and notice that the second instruction gets executed after the first one.

Reference : [Denis Bashkatov's article](https://en.wikichip.org/wiki/skylake_server)

### RDTSCTP INSTRUCTION FOR MEASUREMENTS

RDTSCTP instruction can flush the pipeline to discard the instructions prior to the measurement and read the TSC value of the CPU.

TSC : timestamp counter

You can use CPUID and RDTSCTP combination in older systems that don't support RDTSCTP.

### ESTIMATING INSTRUCTION LATENCIES

Based on Agner Fog's [Instruction tables](https://www.agner-fog.com/instruction-tables/), RDTSCTP reciprocal throughput (clock cycle per instruction) is 32 on Skylake microarchitecture.

-> 1 cycle @4.5GHz = 0.22 nanoseconds  
-> 32\*0.22=7.04 nanoseconds

So its resolution estimate is about 7 nanoseconds on a 4.5 GHz Skylake microarchitecture. You have to recalculate it for different microarchitectures and clock speeds.

### HYPERTHREADING / SIMULTANEOUS MULTITHREADING

Based on [Intel Software Developer's Manual Volume3](https://www.intel.com/content/www/us/en/developer/articles/technical/hyper-threading.html), it is implemented by 2 virtual cores that share resources including cache memory, branch prediction resources and execution ports. And AMD seems to use the resources in the same way based on Agner Fog's [microarchitecture book](https://www.agner-fog.com/instruction-tables/).

For ex if you app is data-intensive, halved caches won't help. It can be disabled it via BIOS settings.

In general, it moves the control of resources from software to hardware and that is usually not desired for performance critical applications.

Note : Its generic name is simultaneous multithreading. Hyperthreading name used by Intel only.

### DYNAMIC CLOCK SPEEDS

Modern CPUs employ dynamic frequency scaling which means there is a min and max frequency per CPU core.

Also [ACPI](https://en.wikichip.org/wiki/skylake_server) defines multiple power states and modern CPUs implement those. P-State is for performance and C states are for energy efficiency.

You can use Intel's [TurboBoost](https://www.intel.com/content/www/us/en/developer/articles/technical/hyper-threading.html) or AMD's [TurboCore](https://www.amd.com/en/tech/hyperthreading) to maintain the CPU usage.

Note that SSE usage may also introduce downclocking, therefore they should be used carefully : [Daniel Lemire's article](https://en.wikichip.org/wiki/skylake_server)

## ARITHMETIC REALM

### ARITHMETIC INSTRUCTION LATENCIES

You can see a set of arithmetic operations from fast to slow below.

The clock cycles are based on Agner Fog's [Instruction tables](#) & Skylake architecture on 64 bit registers.

Bitwise operations , integer add/sub : 0.25 to 1 clock cycle
Floating point add : 3 clock cycles
Floating point multiplication : about 5 clock cycles
Floating point division : about 14-16 clock cycles
Integer division : 24-50 clock cycles

### FLOATING POINTS

X86 uses [IEEE 754](https://en.wikichip.org/wiki/skylake_server) standard for floating points. A 32 bit floating point consists of 3 parts in the memory layout. Below you can see all bits of 32bit IEEE 754 FP number. Used [https://en.wikichip.org/wiki/skylake\\_server](https://en.wikichip.org/wiki/skylake_server) as visualizer :



A floating point's value is calculated as :  $smantissa \times 2^{exponent}$

IEEE754 also defines [denormal numbers](https://en.wikichip.org/wiki/skylake_server). They are very small / near zero numbers.

As floating points are approximations, denormal numbers are needed to avoid an undefined case of x/0, but x=0.

If (a != b) return 1; if (a < b) return 1; if (a > b) return 1;

Without denormal the code to the right would make a divide-by-zero exception. Reference : [Blue Dawson's article](https://en.wikichip.org/wiki/skylake_server)

Based on Agner Fog's [microarchitecture book](https://www.agner-fog.com/instruction-tables/), Intel CPUs have a penalty for denormal numbers, for ex 128 clock cycles on Skylake. They also can be turned off on Intel CPUs.

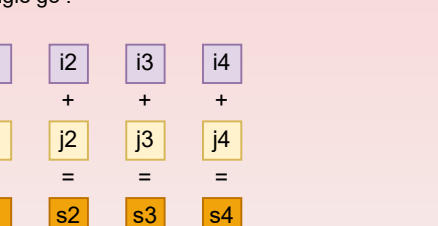
As for AMD side, the recent Zen architecture CPUs seemingly don't have the same performance degradation.

### X86 EXTENSIONS

x86 extensions are specialised instructions. They have various categories such as [conditional](https://en.wikichip.org/wiki/skylake_server) and [register-related](https://en.wikichip.org/wiki/skylake_server) operations.

For the list of extensions : [https://en.wikichip.org/wiki/skylake\\_server](https://en.wikichip.org/wiki/skylake_server)

**SSE** (Streaming SIMD Extensions) is one of the most important ones. **SIMD** stands for "single instruction multiple data". SIMD instructions use wider registers to execute more work in a single op :



In the example above, an array of 4 integers (11 to 44) are added to another array of integers (11 to 44). The result is also an array of sums (11 to 44). In this example, 4 additions are executed by a single instruction.

Some typical application areas are 3D graphics and quantitative finance.

Apart from arithmetic operations, they can be utilised for string operations as well : [https://en.wikichip.org/wiki/skylake\\_server](https://en.wikichip.org/wiki/skylake_server)

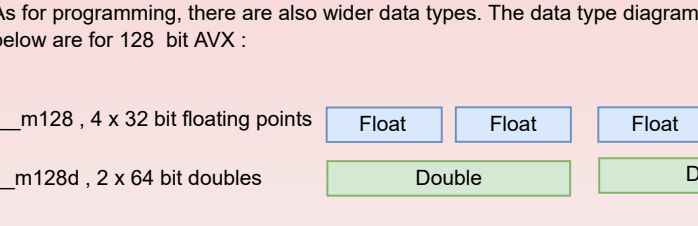
A SIMD based JSOM parser [https://en.wikichip.org/wiki/skylake\\_server](https://en.wikichip.org/wiki/skylake_server)

### X86 EXTENSIONS : SIMD DETAILS

The most recent SIMD instruction sets and their corresponding registers are :

AVX : 128 bits, XMM registers  
AVX2 : 256 bits, YMM registers  
AVX512 : 512 bits, ZMM registers

As for programming, there are also wider data types. The data type diagrams below are for 128 bit AVX :



## LOAD STORE REALM

### LOAD & STORE BUFFERS

Load and store buffers allow CPU to do out-of-order execution on loads and stores by decoupling speculative execution and committing the results to the cache memory.

Reference : [https://en.wikipedia.org/wiki/Memory\\_disambiguation](https://en.wikipedia.org/wiki/Memory_disambiguation)

### STORE-TO-LOAD FORWARDING

Using buffers for stores and loads to support out of order execution leads to a data synchronization issue. That issue is described in [en.wikipedia.org/wiki/Memory\\_disambiguation](https://en.wikipedia.org/wiki/Memory_disambiguation)Store-to-load forwarding.

As a solution, CPU can forward a memory store operation to a following load, if they are both operating on the same address.

An example store and load sequence :

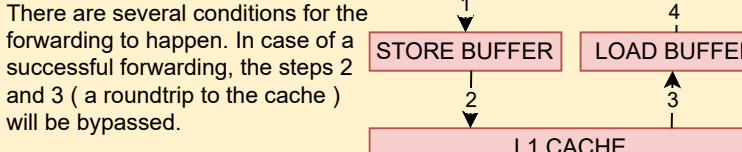
```
mov eax,ecx ; STORE : Write the value of ECX register to the memory  
; address which is stored in EAX register  
mov ecx,[eax] ; LOAD: Read the value from that memory address  
; (which was just used) and write it to ECX register
```

### STORE-TO-LOAD FORWARDING & LHS & PERFORMANCE

Based on [Intel Optimization Manual 3.6.4](https://en.wikichip.org/wiki/skylake_server), Store-to-load forwarding may improve combined latency of those 2 operations. The reason is not specified however it is potentially LHS (Load-Hit-Store) problem in which the penalty is a round trip to the cache memory.

<https://en.wikipedia.org/wiki/Load-Hit-Store>

There are several conditions for the forwarding to happen. In case of a successful forwarding, the steps 2 and 3 (a roundtrip to the cache) will be bypassed.



The conditions for a successful forwarding and latency penalties in case of no-forwarding can be found in Agner Fog's [microarchitecture book](https://www.agner-fog.com/instruction-tables/).

Previous game consoles PlayStation3 and Xbox360 had PowerPC based processors which did in-order execution rather than out-of-order execution.

Therefore developers had to separately handle LHS by using [flushing](https://en.wikichip.org/wiki/skylake_server) keyword and other methods : [Alan Rusakov's article](https://en.wikichip.org/wiki/skylake_server)

## BRANCH PREDICTION REALM

### BRANCH PREDICTION BASICS

Why : CPUs proactively fetch instructions of potentially upcoming branches to utilise the pipeline as much as possible.

Goal if predicted correctly : If the right branch was predicted that will increase the throughput as it completed fetching a set of instructions in advance.

Penalty in case of misprediction : If the prediction was wrong, that prefetch will be a waste and the cost will be flushing the pipeline.

What are branch instructions instructions ? : Unconditional ones (jmp), conditional ones (eg: jne) : called

How : There are hardware auxiliary buffers.

Branch target buffer stores target addresses (instructions) of branches. AMD uses multiple level of BTBs : L1 BTB, L2 BTB etc.

Pattern history table tracks the history of results (whether it was taken or not) per branch (jpe, jne etc.) :

BP METHODS : 2-LEVEL ADAPTIVE BRANCH PREDICTION

A 2-bit saturating counter can store 4 strength states.

Whenever a branch is taken it goes stronger.

And whenever a branch is not taken it goes weaker.

2 level adaptive predictor

In this method, you store the history of last n occurrences in a history register which is n bits.

Also you create a table called "system history table" for that branch. That pattern history table keeps 2^n ones and each row has a saturating counter.

The branch history register will be used to choose which row will be used from the pattern history table.

Reference : Agner Fog's [microarchitecture book](https://www.agner-fog.com/instruction-tables/)

### BP METHODS : AMD PERCEPTIONS

They are used in Zen architectures.

A [perception](https://en.wikichip.org/wiki/skylake_server) is basically the simplest form of machine learning. They can be considered as a linear array of weights.

Agner Fog mentions that they are good at predicting very long branches compared to 2-level adaptive branch prediction in his [microarchitecture book](https://www.agner-fog.com/instruction-tables/).

For details of perception based branch prediction : [Dynamic Branch Prediction with Perceptions by Daniel Lemire and Calvin Lin](https://en.wikichip.org/wiki/skylake_server)

Intel LBD (LOOP STREAM DETECTOR)

Intel LBD will detect a loop and stop fetching instruction to improve frontend bandwidth. Several conditions mentioned in [Intel Optimization Manual](https://en.wikichip.org/wiki/skylake_server).

- Loop body size up to 60 uops, with up to 15 taken branches, and up to 15 64-byte fetch lines.

- No CALL or RET.

- No mismatched stack operations (e.g., more PUSH than POP).

- More than 10 iterations.

Note that LBD is disabled on Skylake Server CPUs. Reference : [https://en.wikichip.org/wiki/skylake\\_server](https://en.wikichip.org/wiki/skylake_server)

Disabling speculative execution patches

You can consider disabling system patches for speculative execution related vulnerabilities such as Meltdown and Spectre for performance, if it is disabled in your system.

Kernel.org documentation : <https://www.kernel.org/doc/html/latest/admin-guide/known-vulnerabilities.html>

Meltdown paper : <https://meltdownattack.com/meltdown.pdf>

Spectre paper : <https://spectreattack.com/spectre.pdf>

ESTIMATED LIMITS : HOW MANY IFs ARE TOO MANY ?

As for max number of writes in BTBs, there are estimations made by stress testing the BTB with conditions of branch instructions :

Intel Xeon Gold 6262 -> roughly 4K  
AMD EPYC 7713 -> roughly 3K

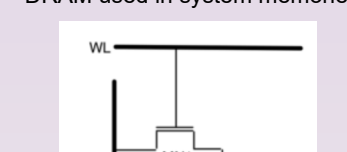
Reference : [Mark Malozki's article on Cloudflare blog](https://en.wikichip.org/wiki/skylake_server)

## CACHE MEMORY REALM

### CACHE BASICS

System memory is made of SRAM cells. Cache memory on the other hand are made of DRAM cells which is much faster than DRAMs. On the other hand they are more expensive.

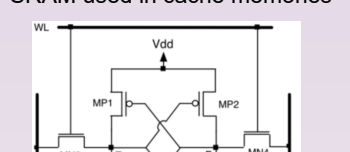
DRAM used in system memories



Access time : 50-150 nanoseconds due to capacitor charge/discharge times and other steps

Cost : Cheaper as it has less components

SRAM used in cache memories



Access time : Under 1 nanosecond

Cost : Expensive due to 6 transistors

Reference : Ulrich Drepper's [What every programmer should know about memory](https://en.wikichip.org/wiki/skylake_server)

### CACHE ORGANISATION

Caches are organised in multiple levels. As you go upper in that hierarchy, the capacity increases. Therefore LLC term used to indicate the last level of cache.

3 level caches are currently the most common ones. Intel Broadwell architecture had 4 level caches in the past. It is expected that upcoming AMD CPUs may come with 4 level of caches.

A **cache line** is the smallest addressable unit in cache memories. It is typically 64 bytes.

All the mentioned caches still now were data caches. But there is also **instruction cache** (I-Cache) which store program instructions rather than data to improve bandwidth of CPU frontend.

In case of a cache hit, the latency is typically single digit nanoseconds. And in case of a cache miss, we need a round trip to the system memory and local latency becomes 3 digit nanoseconds.

Used & Store buffers

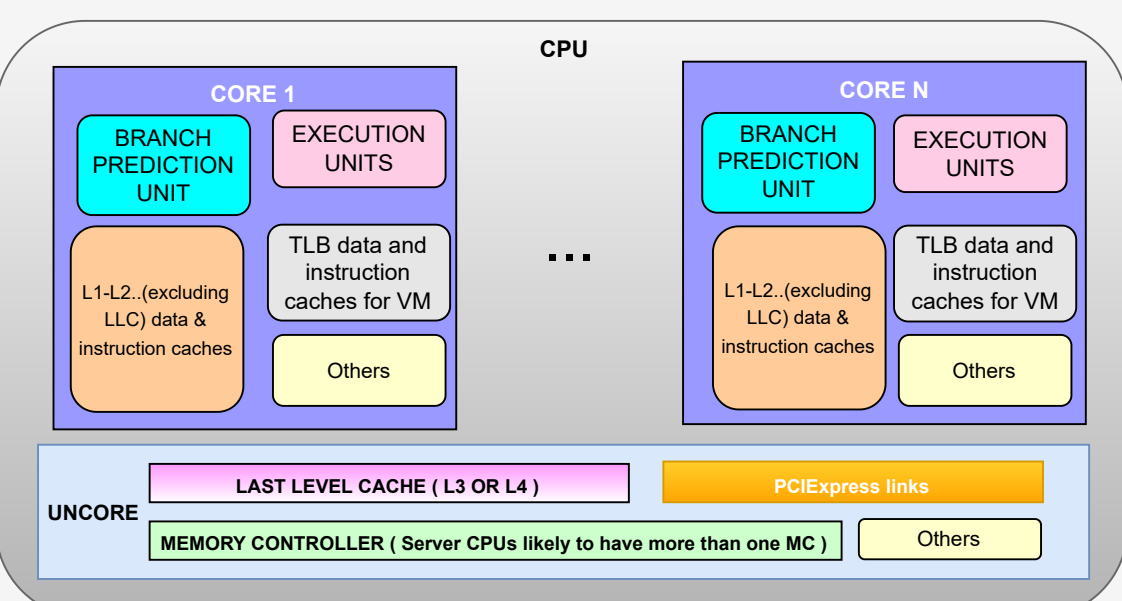
L1 Data Cache

LLC : Last level cache

System memory

## MULTICORE REALM

### TOPOLOGICAL OVERVIEW - INTEL CPUs



Note that "Uncore" term is used by only Intel.

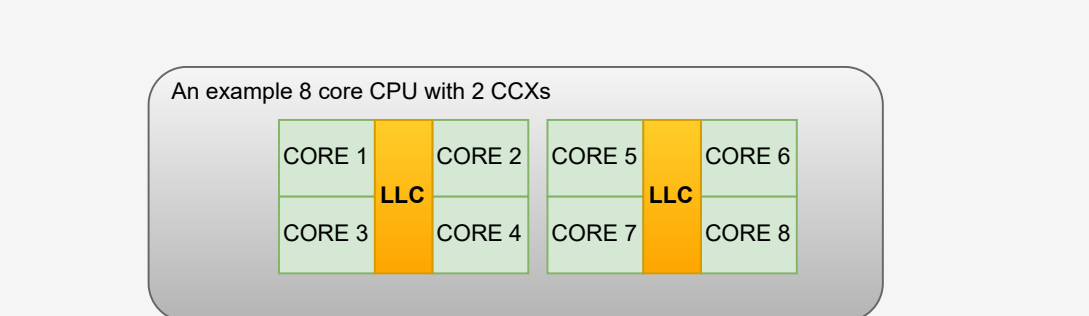
### TOPOLOGICAL OVERVIEW - AMD CPUs

Most of AMD topology is similar to Intel. However starting from Zen microarchitecture, one key difference is CDOs.

AMD CPUs are designed as group of 4 cores which is called as CCX (Core complex) by AMD.

The main difference between Intel CPUs is that there is one LLC per each CCX/quad core.

Practically the max number of cores competing for the LLC (without simultaneous multithreading) is 4 in recent AMD CPUs.



AMD CCX structure : There will be multiple LLCs for each group of 4 cores

Reference : [https://en.wikichip.org/wiki/amd/microarchitectures/zen2/CPU\\_Complex\\_28CCX\\_29](https://en.wikichip.org/wiki/amd/microarchitectures/zen2/CPU_Complex_28CCX_29)

## COHERENCY

### CACHE COHERENCY PROTOCOLS

As LLC is typically shared by multiple cores, cache coherency protocols are needed to avoid data hazards. Intel CPUs use MESI and AMD CPUs use MOESI, however both heavily depend on MESI protocol. There are 4 states for a cache line in MESI protocol :

Modified - cache line is present only in the current cache and has been modified from its value in system memory

Exclusive - cache line is present only in the current cache and matches its value in system memory

Shared - cache line is present here and the other cache line matches its value in system memory

Invalid - cache line is unused

Allowed state transitions in MESI protocol :

As for state transition costs, M and S states are the cheapest ones as they don't involve cross cache communication. Therefore it is useful to keep data in those states as long as possible. That can be achieved by using cached variables whenever it is applicable.

Initial MESI : [https://en.wikichip.org/wiki/amd/microarchitectures/zen2/CPU\\_Complex\\_28CCX\\_29](https://en.wikichip.org/wiki/amd/microarchitectures/zen2/CPU_Complex_28CCX_29)

AMD MOESI : [https://en.wikichip.org/wiki/amd/microarchitectures/zen2/CPU\\_Complex\\_28CCX\\_29](https://en.wikichip.org/wiki/amd/microarchitectures/zen2/CPU_Complex_28CCX_29)

FALSE SHARING AND CACHE PING-PONGING

If Core1 updates "var1" variable, that will be touching that shared cache line.

That change will then need to be propagated to all other cores by the cache coherency protocol, even though other cores are not using that variable.

That situation is called **cache sharing**.

If those happen in higher rates and if cache lines are transferred between cores rapidly, that situation is called as **cache ping-pong**.

IP1 : Interprocessor interrupt

1 : One of the cores modifies a table entry

PAGE TABLE ON SYSTEM MEMORY

OS

TLB SHOOTDOWNS

Whenever a page table entry is modified by any of the cores, that particular TLB entry is invalidated in all cores by sending them IP1s.

IP1 : Interprocessor interrupt

1 : One of the cores modifies a table entry

PAGE TABLE ON SYSTEM MEMORY

OS

TLB SHOOTDOWNS

Whenever a page table entry is modified by any of the cores, that particular TLB entry is invalidated in all cores by sending them IP1s.

IP1 : Interprocessor interrupt

1 : One of the cores modifies a table entry

PAGE TABLE ON SYSTEM MEMORY

OS

TLB SHOOTDOWNS

Whenever a page table entry is modified by any of the cores, that particular TLB entry is invalidated in all cores by sending them IP1s.

IP1 : Interprocessor interrupt

1 : One of the cores modifies a table entry

PAGE TABLE ON SYSTEM MEMORY

OS

TLB SHOOTDOWNS

Whenever a page table entry is modified by any of the cores, that particular TLB entry is invalidated in all cores by sending them IP1s.

IP1 : Interprocessor interrupt

1 : One of the cores modifies a table entry

PAGE TABLE ON SYSTEM MEMORY

OS

## ABOUT ALL REALMS

Intel's Top Down analysis is hierarchical organization of microarch events. There are 2 main categories for stalls :

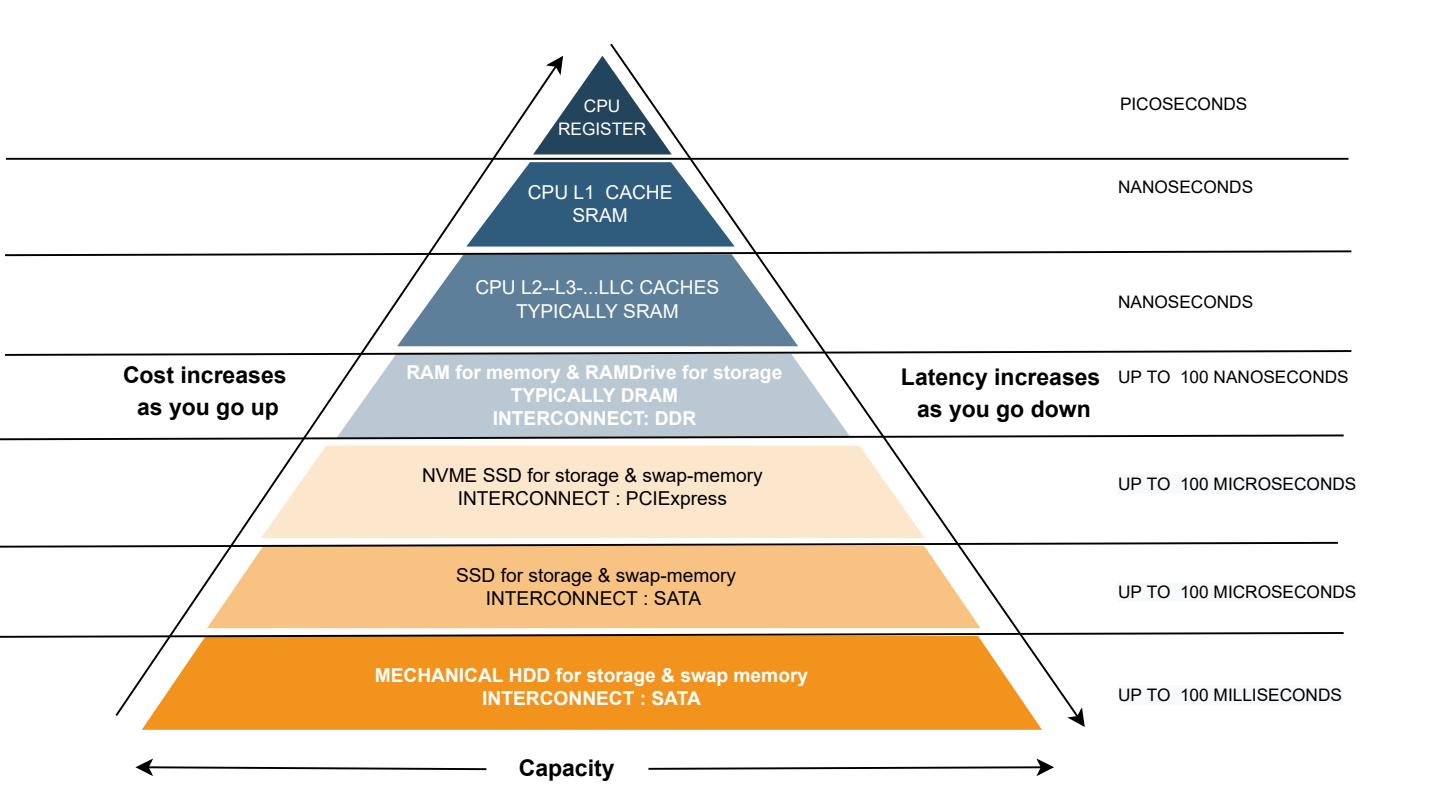
1. Frontend : A typical example is large code sizes leading to instruction cache misses.

2. Backend : Usually either memory bound or compute bound

Branch mispredictions are categorised under "bad speculation".

You can see after Intel's [Ytting](https://en.wikichip.org/wiki/skylake_server) or [Tooley](https://en.wikichip.org/wiki/skylake_server) talk to make a top down analysis. Both use Intel CPUs' PMCs (performance monitoring counter)

### OVERALL MEMORY & STORAGE HIERARCHY



### (ESTIMATED) LATENCY NUMBERS IN CLOCK CYCLES