

MICROARCHITECTURE CHEAT SHEET

X86 CPUs & Performance

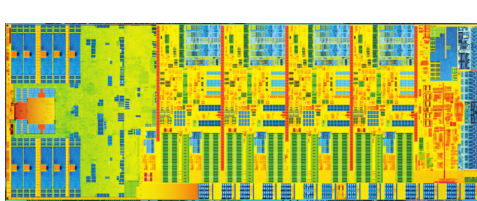


Photo of a quad-core CPU. Copyright of Intel

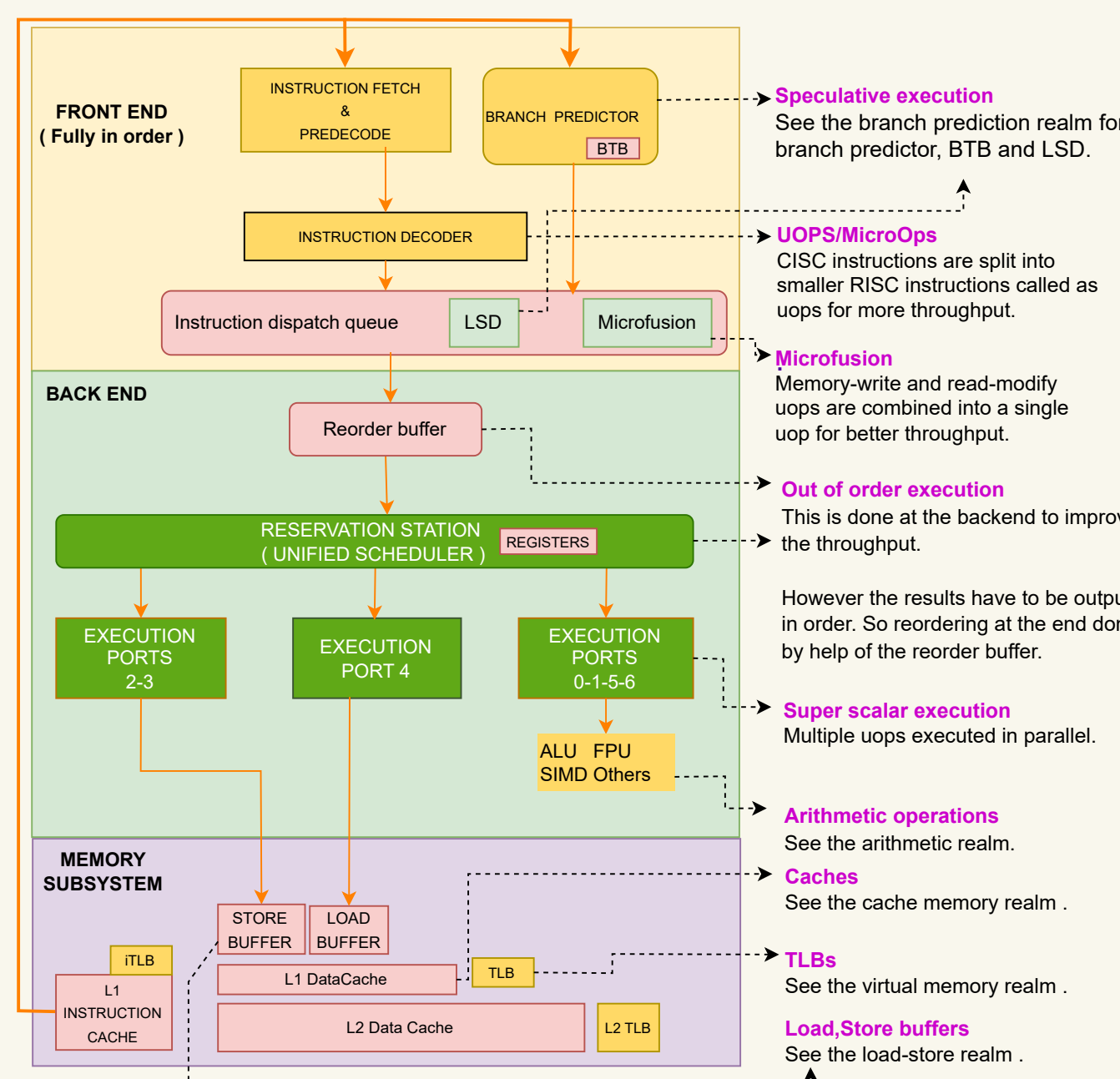
LAST UPDATE DATE : 18 JULY 2023

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

AUTHOR: AKIN OCAL akin_ocal@hotmail.com

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/en/skylake)

AMD pipelines : The main difference in AMD architectures is that there are parallel pipelines for integers and floating points : [AMD Zen2 pipeline diagram on en.wikichip.org](https://en.wikichip.org/en/amd-zen2)

- Speculative execution**
See the branch prediction realm for branch predictor, BTB and LSD.
- UOPS/MicroOps**
CISC instructions are split into smaller RISC instructions called as uops for more throughput.
- Microfusion**
Memory write and read-modify ops are combined into a single uop for better throughput.
- Out of order execution**
This is done at the backend to improve the throughput.
However the results have to be output in order. So reordering at the end done by help of the reorder buffer.
- Super scalar execution**
Multiple uops executed in parallel.
- Arithmetic operations**
See the arithmetic realm.
- Caches**
See the cache memory realm.
- L1.Bs**
See the virtual memory realm.
- Load-Store buffers**
See the load-store realm.

PIPELINE PARALLELISM & PERFORMANCE

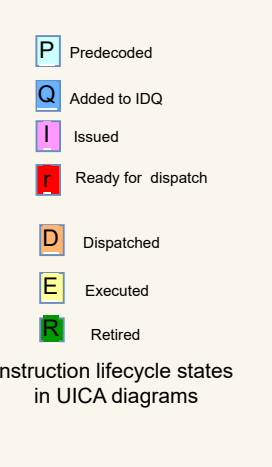
Pipeline diagrams : The diagrams below in the following topics are outputs from an online microarchitecture analysis tool [wikichip](https://en.wikichip.org/en/pipeline) and they represent parallel execution through cycles.

Rows are multiple instructions being executed at the same time.

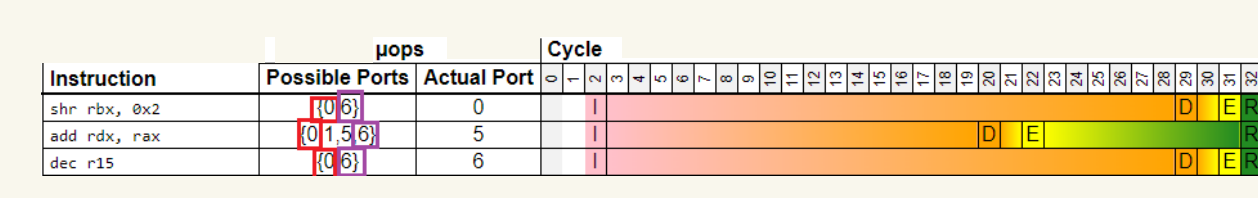
Columns display how instruction state changes through cycles.

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 fetched. Therefore a high rate of retired instructions indicates low branch prediction rate.



CONTENTION FOR EXECUTION PORTS IN THE PIPELINE

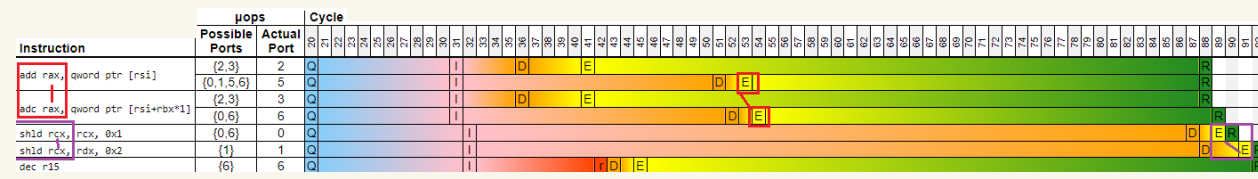


In the example above, all instructions are working on different registers, but SHR, ADD, DEC instructions are competing for ports 0 and 6. SHR and DEC are getting executed after ADD instruction.

Also notice that there is longer time between (Executed) and (Retired) states of instruction ADD as retirement has to be done in-order whereas execution is out-of-order.

Reference : [Denis Bakhtyalov's article](https://en.wikichip.org/en/denis-bakhtyalov)

INSTRUCTION STALLS DUE TO DATA DEPENDENCY



In the example above, there are 2 dependency chains, each marked with a different colour. In the first red coloured one, 2 instructions are competing for RAX register and notice that the second instruction gets executed after the first one. And the same applies to the 2nd purple pair.

Reference : [Denis Bakhtyalov's article](https://en.wikichip.org/en/denis-bakhtyalov)

RDTSSCP INSTRUCTION FOR MEASUREMENTS

TSC (Time Stamp Counter) is a special register that counts CPU cycles. RDTSSCP can be used to read the TSC value which then can be used for measurements. It can also avoid out-of-order execution effects to a degree.

It does not wait until all previous instructions have executed and all previous loads are globally visible.

(From [Intel Software Developer's Manual Volume 4.3](https://en.wikichip.org/en/intel-software-developer-manual-volume-4-3), April 2022)

Intel's [How to benchmark code execution times](https://en.wikichip.org/en/how-to-benchmark-code-execution-times) whitepaper has details of using RDTSSCP instruction.

AMD Programmers Manual Vols states : RDTSSCP forces all older instructions to retire before reading the time-stamp counter.

ESTIMATING INSTRUCTION LATENCIES

You can use Agner Fog's [Instruction Latencies](https://en.wikichip.org/en/instruction-latencies) to find out instructions' regional throughputs (clock cycle per instruction). As an example, response throughput of instruction RDTSSCP is 32 on a Skylake microarchitecture :

- > 1 cycle (84.5GHz) (highest frequency in Skylake) is 0.22 nanoseconds
- > 32*0.22=7.04 nanoseconds

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

HYPERTHREADING / SIMULTANEOUS MULTITHREADING

Hyperthreading name is used by Intel and it is called as "Simultaneous multithreading" by AMD. In both resources incases and execution units are shared.

Reference : Agner Fog's [microarchitecture book](https://en.wikichip.org/en/microarchitecture-book) has "multithreading" sections for each of Intel and AMD microarchitectures.

Regarding using it, if your app is data-intensive, halved caches won't help. Therefore it can be disabled in BIOS settings. In general, it moves the control of resources from software to hardware and that is usually not desired for performance critical applications.

DYNAMIC FREQUENCIES

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

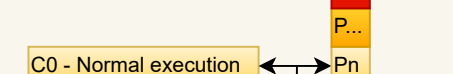
ACPI : ACPI defines multiple power states and modern CPUs implement those. P-states are for performance and C-states are for energy efficiency.

Intel has various tunability options and the most well known is TurboBoost. On AMD side there is [TurboCore](https://en.wikichip.org/en/turbo-boost).

You can use those to maximize the CPU usage.

Number of active cores & SIMD AVX2/AVX512 on Intel CPUs : Intel's power management policies are complex. See the arithmetic and the microarchitecture realms as number of active cores and some of AVX2/AVX512 extensions also may affect the frequency while in TurboBoost.

Varying max clock speeds on AMD CPUs : Some AMD CPUs' cores have slightly varying max frequencies. Reference : [AMD's GDC2 presentation page46](https://en.wikichip.org/en/amd-gdc2-presentation-page46)



It is better to switch to P-states. C-states has to be brought in to C0 state.

Number of active cores & SIMD AVX2/AVX512 on Intel CPUs

Number of active cores & SIMD AVX2/AVX512 on Intel CPUs

Number of active cores & SIMD AVX2/AVX512 on Intel CPUs

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.wikichip.org/en/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. This issue is described in [en.wikichip.org/wiki/Memory_disambiguation](https://en.wikichip.org/en/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.

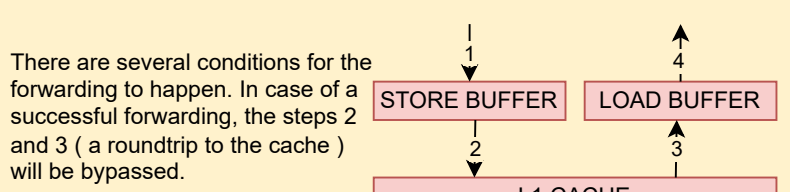
An example store and load sequence :

```
mov eax,ecx ; STORE. Write the value of ECX register to the memory address which is stored in EAX register
mov ecx,ecx ; 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](https://en.wikichip.org/en/intel-optimization-manual) 3.6.4, 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.wikichip.org/en/wiki/L1-load-Hit-Store>



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

The conditions for a successful forwarding and latency penalties in case of non-forwarding can be found in Agner Fog's [microarchitecture book](https://en.wikichip.org/en/microarchitecture-book).

What would happen without forwarding? : In the past, game consoles PlayStation3 and Xbox360 had PowerPC based processors which used in-order execution rather than out-of-order execution. Therefore developers had to separately handle LHS by using [atomic](https://en.wikichip.org/en/atomic) keyword and other methods : [Elan Rusakov's article](https://en.wikichip.org/en/Elan-Rusakov's-article)

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 Latencies](https://en.wikichip.org/en/instruction-latencies) & Skylake architecture on 64 bit registers.

Bitwise operations, integer arithmetic: 0.25 to 1 clock cycle
Floating point multiplication: about 3 clock cycles
Integer division: 24-30 clock cycles

FLOATING POINTS

X86 uses IEEE 754 standard for floating points. A 32 bit floating point consists of 3 parts in the memory layout. Below you can see all bits of 1234.5678 FP number. Used [pnlabs.org/en/pipeline](https://en.wikichip.org/en/pipeline) as visualizer :



A floating point value is calculated as : mantissa * 2^exponent

IEEE754 also defines **denormal numbers**. They are very small / near zero numbers.

As floating points are approximations, denormal numbers are needed to avoid an underflow case of: zero, but a-b=0

Without denormalize the code to the right would involve a divide-by-zero exception.

Reference : [Bruce Danielsen's article](https://en.wikichip.org/en/Bruce-Danielsen's-article)

Based on Agner Fog's [microarchitecture book](https://en.wikichip.org/en/microarchitecture-book), Intel CPUs have a penalty for denormal numbers, for ex: 129 clock cycles on Skylake. They also can be turned off on Intel CPUs.

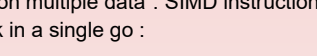
As for AMD side, the recent Zen architecture CPUs seemingly don't have the same performance degradation.

X86 EXTENSIONS

X86 extensions are specialised instructions. They have various categories from [cpuid](https://en.wikichip.org/en/cpuid) to [bugcheck](https://en.wikichip.org/en/bugcheck) network operations.

Intel [Internals Guide](https://en.wikichip.org/en/Intel-Internals-Guide) is a good page to explore those extensions.

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 go.



In the example above, an array 4 integers (11 to 44) are added to another array of integers (1 to 4). The result is also an array of same (x1 to x4). In this example, 4 add operations are executed by a single instruction.

They play key role in compilers' vectorisation optimisations : [GDC auto vectorisation](https://en.wikichip.org/en/GDC-auto-vectorisation)

Apart from arithmetic operations, they can be utilised for string operations as well : Daniel Lemire's SIMD based JSON parser : [https://github.com/daniellemire/simdjson](https://en.wikichip.org/en/https://github.com/daniellemire/simdjson)

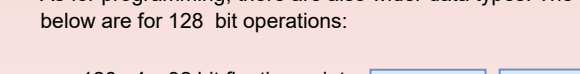
X86 EXTENSIONS : SIMD DETAILS

The most recent SIMD instruction sets for Intel CPUs are :

- AVX : Up to 256 bits
- AVX2 : Up to 256 bits
- AVX512 : Up to 512 bits

Recent AMD CPUs support AVX & AVX2. Only the latest Zen4 architecture supports AVX512.

As for programming, there are also data type details. The data type diagrams below are for 128 bit operations.



__m128i, 4 x 32 bit floating points

__m128d, 2 x 64 bit doubles

__m128t, 4 x 32 bit ints

__m128i, 4 x 32 bit long ints

Note that as SIMD instructions require more power, therefore usage of some AVX512 extensions may introduce downclocking. They should be benchmarked. For details : [Daniel Lemire's article](https://en.wikichip.org/en/Daniel-Lemire's-article)

BRANCH PREDICTION REALM

BRANCH PREDICTION BASICS

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

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

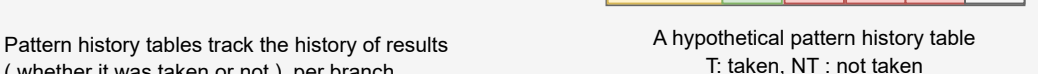
Penalty in case of misprediction : If the prediction was wrong, that prefetch will be a waste and the cost will be flushing the pipeline.

What are branch instructions? : Unconditional ones (jmp), conditional ones (eg: jne), call/ret.

How : There are auxiliary hardware buffers.

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

Pattern history tables track the history of results (whether it was taken or not) per branch.



A hypothetical pattern history table
T: taken, NT: not taken

CMOV (Conditional move) instruction also computes the conditions for some additional time. Therefore they don't introduce extra load to the branch prediction mechanism. They can be used to eliminate branches.

Reference : [Intel Optimization Manual](https://en.wikichip.org/en/Intel-Optimization-Manual) 3.4.1.1

BP METHODS : 2-LEVEL ADAPTIVE BRANCH PREDICTION

Saturating counter as a building block
Because cumulative memory requirement of multiple processes in an OS can be more than the system capacity.

Wherever a branch is taken it goes stronger.

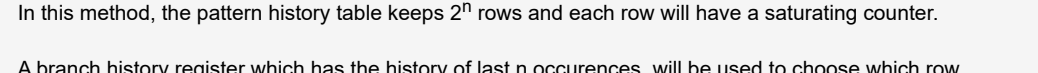
And whenever a branch is not taken it goes weaker.

2 level adaptive predictor

In this method, the pattern history table keeps 2^rows and each row will have a saturating counter.

A branch history register which has the history of last n occurrences, will be used to choose which row will be used from the pattern history table.

Reference : Agner Fog's [microarchitecture book](https://en.wikichip.org/en/microarchitecture-book) 3.1.



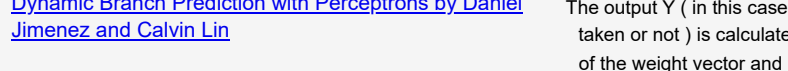
BP METHODS : AMD PERCEPTIONS

A **perception** is basically the simplest form of machine learning. It 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://en.wikichip.org/en/microarchitecture-book) 3.1.2.

For details of perception based branch prediction: [Dynamic Branch Prediction with Perceptions by Daniel Jimenez and Gavin Liu](https://en.wikichip.org/en/Dynamic-Branch-Prediction-with-Perceptions-by-Daniel-Jimenez-and-Gavin-Liu)

The output Y (in the case whether a branch taken or not) is calculated by the dot product of the weight vector and the input vector.



INTEL LSD (LOOP STREAM DETECTOR)

Intel LSD will detect a loop and stop fetching instructions to improve the frontend bandwidth. Several conditions mentioned in [Intel Optimization Manual](https://en.wikichip.org/en/Intel-Optimization-Manual) 3.4.2.1.

- Loop body size is 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 ~20 iterations.

Note that LSD is disabled on Skylake Server CPUs. Reference : [https://www.intel.com/content/www/us/en/processors/skylake-server/skylake-server-features-and-specifications.html](https://en.wikichip.org/en/https://www.intel.com/content/www/us/en/processors/skylake-server/skylake-server-features-and-specifications.html)

DISABLING SPECULATIVE EXECUTION PATCHES

You can consider disabling system patches for speculative execution related vulnerabilities such as Meltdown and Spectre for performance, if that is doable in your system. Those patches are not only microcode updates but they also need OS support.

Kernel.org documentation : [https://www.kernel.org/doc/Documentation/bugs/bugs-meltdown-spectre.txt](https://en.wikichip.org/en/https://www.kernel.org/doc/Documentation/bugs/bugs-meltdown-spectre.txt)

Red Hat Enterprise documentation : [https://access.redhat.com/articles/1311001](https://en.wikichip.org/en/https://access.redhat.com/articles/1311001)

Meltdown paper : [https://arxiv.org/pdf/1516.07551v2.pdf](https://en.wikichip.org/en/https://arxiv.org/pdf/1516.07551v2.pdf)

Spectre paper : [https://spectreattack.com/spectre.pdf](https://en.wikichip.org/en/https://spectreattack.com/spectre.pdf)

ESTIMATED LIMITS : HOW MANY IFs ARE TOO MANY ?

As for max number of entries in BTBs, there are estimations made by stress testing the BTB with sequences of branch instructions.

Intel Xeon Gold E320 -> roughly 4K
AMD EPYC 7113 -> roughly 3K

Reference : [Marek Majkowski's article on CPUloadings](https://en.wikichip.org/en/Marek-Majkowski's-article-on-CPUloadings)



CACHE MEMORY REALM

CACHE MEMORY VS SYSTEM MEMORY

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

DRAM used in system memories

SRAM used in cache memories

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

Cost : Cheaper in the price due to 6 transistors

Reference : Intel Drippers' [What every programmer should know about memory](https://en.wikichip.org/en/What-every-programmer-should-know-about-memory)

CACHE ORGANISATION

Caches are organised in multiple levels. As you go upper in that hierarchy, the capacity increases. Therefore L3 term used to indicate the last level of cache.

3 level data caches are currently the most common ones. Intel [Intel Optimization Manual](https://en.wikichip.org/en/Intel-Optimization-Manual) 3.2.3.4, there is only one kind of memory reordering that can happen. Loads can be reordered with earlier stores if they use different memory locations. That reordering will happen if they use the same addresses.

Cache line size : The unit of data transfer between the cache and the system memory. It is typically 64 bytes. And the caches are organised according to the cache line size.

There is also **instruction cache** (Cache) which stores program instructions rather than data to improve throughput of CPU kernel.

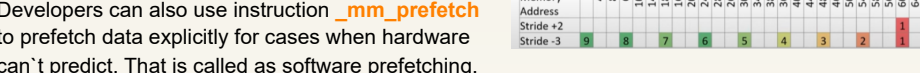
In case of a cache hit, the latency is typically single digit nanoseconds. And in case of a cache miss, we need a round trip to the system memory and total latency becomes 3 digit nanoseconds.

HARDWARE AND SOFTWARE PREFETCHING

Hardware prefetchers detect patterns like **streams** (Ex: accessing to contiguous array members) and **strides** (Ex: accessing specific members in arrays of structs) and prefetch data and instruction to cache lines automatically.

Developers can also use instruction **intel_prefetch** to predict data for cases when hardware can't prefetch data. That is called as software prefetching.

Reference for image : It is taken from [AMD's GDC2 presentation page44](https://en.wikichip.org/en/AMD's-GDC2-presentation-page44)



SYSTEM MEMORY REALM

DOR RAMs

DOR RAMs are the most common commodity hardware as system memory.

They are found in forms of DIMMs (Dual inline memory modules) / RAMsticks.

Organisation

System memory / RAM is organised as collection of ranks.

Each rank have banks which are collection of DRAM cells per bit.

DRAM refreshes

DRAM controls use capacitors which lose their charge over time. (See the cache memory realm) So RAMs have to refresh their DRAM cells periodically.

As for DORA, refreshing is rank-level which means the other banks in the same rank become inaccessible. DORs comes with [active-bank-refresh](https://en.wikichip.org/en/active-bank-refresh) feature which allows a more fine-grained bank-level refresh. Therefore it can offer a higher throughput.

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh granularity

DORs refresh granularity

DORA refresh