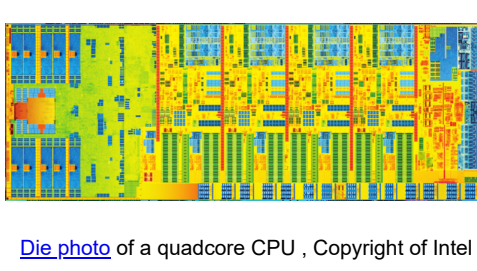


MICROARCHITECTURE CHEAT SHEET



LAST UPDATE DATE : 14 NOV 2022

FOR LATEST VERSION: www.github.com/akhin/microarchitecture-cheatsheet

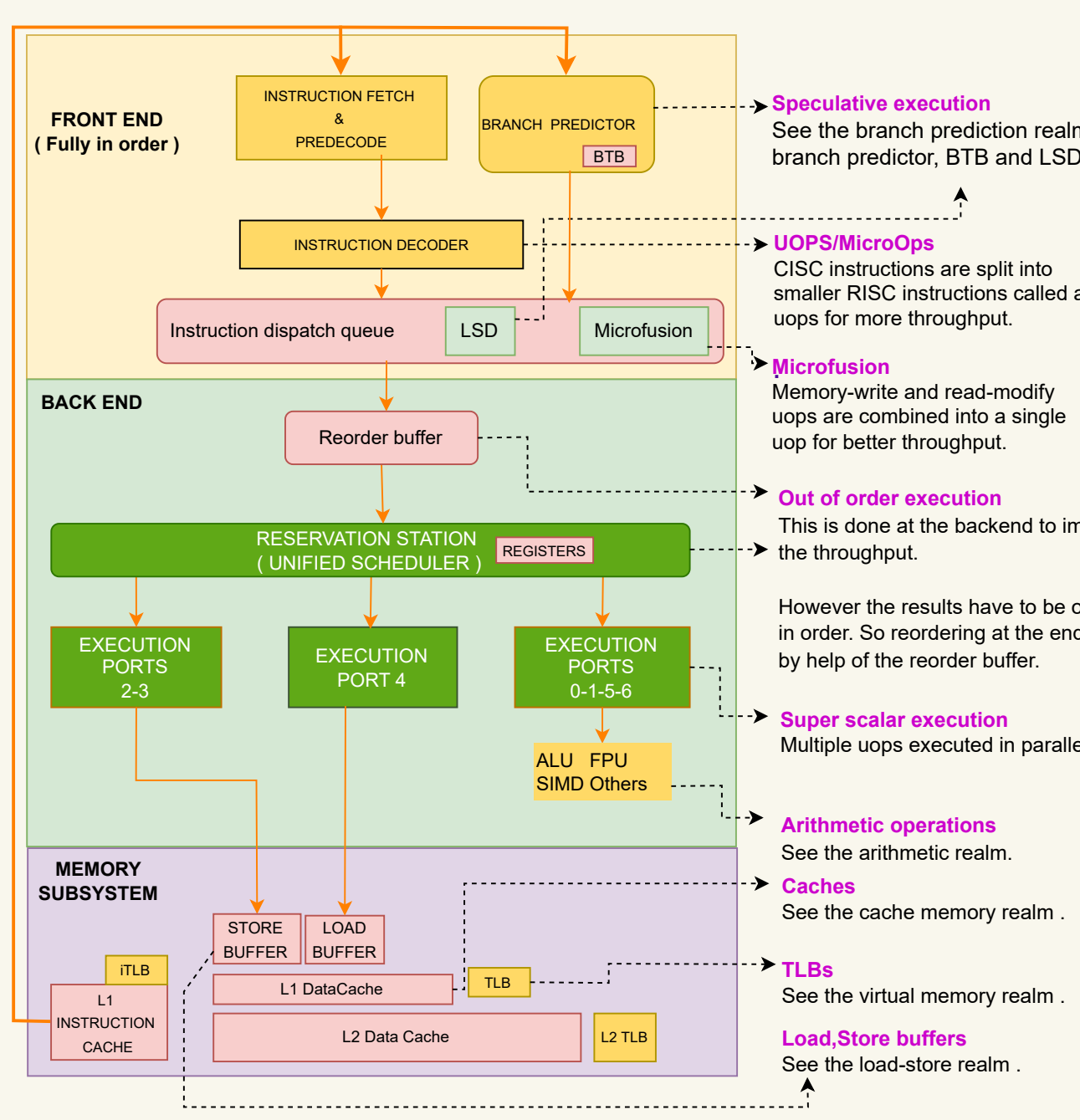
AUTHOR: AKIN OCAL skin_ocal@hotmail.com

X86 CPUs & Performance

PIPELINE REALM :

INSIDE AN INDIVIDUAL CORE

A SIMPLIFIED OVERVIEW (BASED ON INTEL SKYLAKE)



Based on: [SkyLake server on wikichip.org](https://www.wikichip.org/wiki/SkyLake_server)

AMD pipelines : The main difference in AMD architectures is that there are parallel pipelines for integers and floating points: [AMD Zen2 pipeline diagram on wikichip.org](https://www.wikichip.org/wiki/AMD_Zen2_pipeline_diagram)

PIPELINE PARALLELISM & PERFORMANCE

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

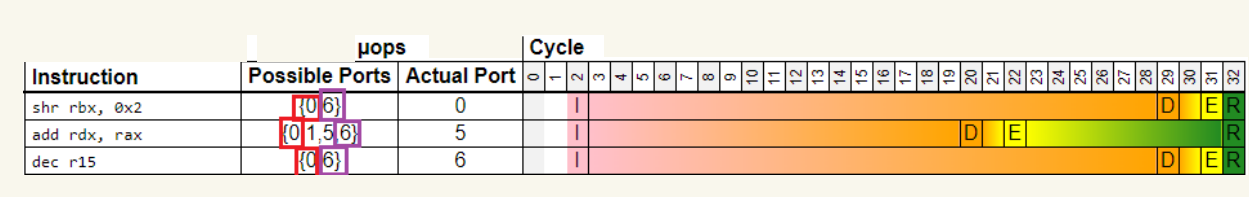
Rows are multiple instructions being executed at the same time.

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.

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, retired 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

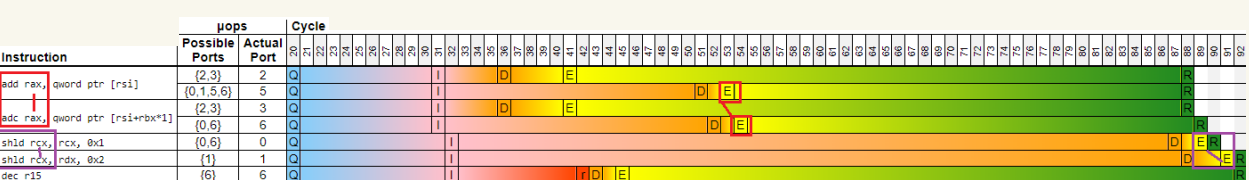


In the example above, all instructions are working on different registers, but SHR, ADD, DEC instructions are competing for ports 0 and 6. SHR and DEC are getting executed after ADD instruction.

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 Bakhtov's article](https://www.wikichip.org/wiki/PIPELINE_PARALLELISM)

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 Bakhtov's article](https://www.wikichip.org/wiki/PIPELINE_PARALLELISM)

RTDSCP INSTRUCTION FOR MEASUREMENTS

TSC (Time stamp counter) is a special register that counts CPU cycles. RTDSCP 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 previous instructions have executed and all previous loads are globally visible.

(From [Intel Software Developer's Manual Volume 4](https://www.wikichip.org/wiki/RTDSCP), 1, April 2022)

Intel's [How to benchmark code execution times](https://www.wikichip.org/wiki/RTDSCP) Wikipedia has details of using RTDSCP instruction.

AMD Programmers Manual Vol3 states : RTDSCP forces all older instructions to retire before reading the time-stamp counter

ESTIMATING INSTRUCTION LATENCIES

You can use Agner Fog's [Instruction tables](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES) to find out instructions' reciprocal throughputs (clock cycle per instruction). As an example, reciprocal throughput of instruction RTDSCP is 32 on Skylake microarchitecture :

-> 1 cycle @4.5GHz (highest frequency on 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 other microarchitectures and clock speeds.

HYPERTHREADING / SIMULTANEOUS MULTITHREADING

Hyperthreading name is used by Intel and it's called as "Simultaneous multithreading" by AMD. In both resources including caches and execution units are shared-Agner Fog's [microarchitecture book](https://www.wikichip.org/wiki/HYPERTHREADING) has "multithreading" sections for each of Intel and AMD microarchitectures.

Regarding using it, if your app is data-intensive, shared caches won't help. Therefore it can be disabled 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.

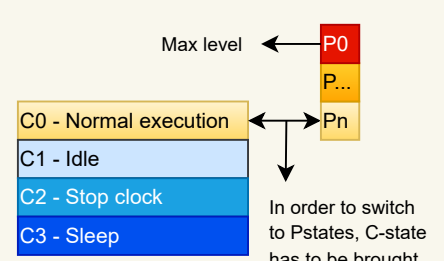
DYNAMIC FREQUENCIES

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

AGN1 : AGN1 defines multiple power states and modern CPUs implement those. P-State are for performance and C-States are for energy efficiency.

Intel has various tunability controls and the most well known is TurboBoost. On AMD side there is [TurboCore](https://www.wikichip.org/wiki/TURBOBOOST). You can use those to maximize the CPU usage.

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



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

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#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.

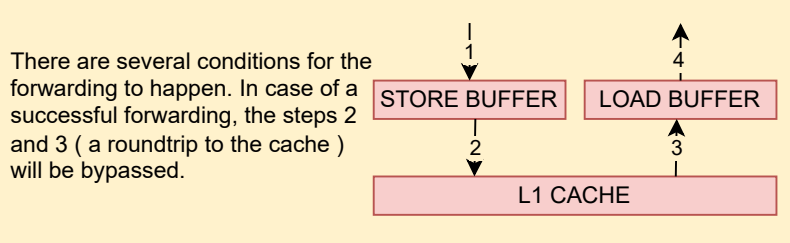
An example store and load sequence :

```
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 is stored in EAX register (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://www.wikichip.org/wiki/STORE-TO-LOAD_FORWARDING), 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/Store-to-load_forwarding



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 [PowerPC](https://www.wikichip.org/wiki/POWERPC) keyword and other methods : [Dan Ruskin's article](https://www.wikichip.org/wiki/POWERPC)

ARITHMETIC REALM

ARITHMETIC INSTRUCTION LATENCIES

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

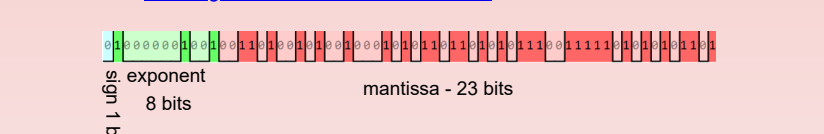
The clock cycles are based on Agner Fog's [Instruction tables](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES) & Skylake architecture on 64 bit registers.

Bitwise operations, integer arithmetic: 0.25 to 1 clock cycle
Floating point add: ~3 clock cycles
Floating point multiplication: ~about 4 clock cycles
Without denormalize the code to the right would involve a divide-by-zero exception.
Integer division: ~24-30 clock cycles

Reference : [Boris Danilov's article](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES)

FLOATING POINTS

X86 uses [IEEE 754](https://www.wikichip.org/wiki/IEEE754) standard for floating points. A 32 bit floating point consists of 3 parts in the memory layout. Below you can see all bits of 128.5879 FP number. Used [agner.fog@code42.com](https://www.wikichip.org/wiki/IEEE754) as visualizer.



A floating point's 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 undefined case of: `while (float < 0) return 1 / (float - 0);`

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

Reference : [Boris Danilov's article](https://www.wikichip.org/wiki/IEEE754)

Based on Agner Fog's [microarchitecture book](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES), 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 [cryptography](https://www.wikichip.org/wiki/X86_EXTENSIONS) to [signal network operations](https://www.wikichip.org/wiki/X86_EXTENSIONS).

[Intel Intrinsics Guide](https://www.wikichip.org/wiki/X86_EXTENSIONS) 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.

As for programming, there are also data type types. 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

__m128q, 2 x 64 bit long longs

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

They play key role in compilers' vectorisation optimisations : [GCC auto vectorization](https://www.wikichip.org/wiki/VECTORISATION)

Apart from arithmetic operations, they can be utilised for string operations as well. A SIMD based JSON parser : https://www.wikichip.org/wiki/SIMD_JSON_PARSER

X86 EXTENSIONS : SIMD DETAILS

The most recent SIMD instruction sets for Intel CPUs are :

AVX-512 : Up to 256 bits

AVX2 : Up to 256 bits

AVX128 : Up to 128 bits

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

As for programming, there are also data type types. 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

__m128q, 2 x 64 bit long longs

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

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

branch 1 2 3 4
branch 0 T NT NT NT
branch 0 T NT NT NT

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 3.4.1.1](https://www.wikichip.org/wiki/BRANCH_PREDICTION)

BP METHODS : 2-LEVEL ADAPTIVE BRANCH PREDICTION

Saturating counter as a building block
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, the pattern history table keeps 2^n ones 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 3.1](https://www.wikichip.org/wiki/BRANCH_PREDICTION).

BP METHODS : AMD PERCEPTIONS

A **perceptron** 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 3.12](https://www.wikichip.org/wiki/PERCEPTIONS).

For details of perceptron based branch prediction: [Dynamic Branch Prediction with Perceptrons by Daniel Ammerer and Calvin Lu](https://www.wikichip.org/wiki/PERCEPTIONS)

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

Intel LSD will detect a loop and stop fetching instructions to improve the frontend bandwidth. Several conditions mentioned in [Intel Optimization Manual 3.4.2.4](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES).

- Loop body size up to 60 bytes, with up to 15 taken branches, and up to 15 64-byte float 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.wikichip.org/wiki/INSTRUCTION_LATENCIES

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.

Kernel.org documentation : https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES

Red Hat Enterprise documentation : https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES

Meltdown paper : https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES

Spectre paper : https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES

ESTIMATED LIMITS : HOW MANY IFs ARE TOO MANY ?
As for number of entries in BTBs, there are estimations made by stress testing the BTB with sequences of branch instructions:

Intel Xeon Gold 6342 -> roughly 4K
AMD EPYC 7713 -> roughly 3K

Reference : [Marek Makowski's article on Cloudflare blog](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES)

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 is much faster than DRAMs. 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: Expensive in the price due to 6 transistors

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

CACHE ORGANISATION

Caches are organised in multiple levels. As you go up 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](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES) had 4 level caches in the past. Also upcoming AMD CPUs may come with 4 level of caches.

Cache line size is 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.

All the mentioned caches still have data caches. But there is also **instruction cache** (I-Cache) which store program instructions rather than data to improve throughput 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 total latency becomes 3 digit nanoseconds.

Intel Optimization Manual 3.7 describes [prefetching](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES). Hardware prefetchers prefetch data and instruction to cache lines automatically. Developers can also use instruction `prefetch` to prefetch data manually. That is called as software prefetching. However performance improvement by using software prefetcher is controversial : [Daniel Lemire's article](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES)

HARDWARE AND SOFTWARE PREFETCHING

Intel Optimization Manual 3.7 describes [prefetching](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES). Hardware prefetchers prefetch data and instruction to cache lines automatically. Developers can also use instruction `prefetch` to prefetch data manually. That is called as software prefetching. However performance improvement by using software prefetcher is controversial : [Daniel Lemire's article](https://www.wikichip.org/wiki/INSTRUCTION_LATENCIES)

The disadvantage of another level is that you can address even more space.

The disadvantage is that the time needed to walk the page table increases due to a new level of indirection.

VIRTUAL MEMORY REALM

VIRTUAL MEMORY ORGANISATION

Why virtual memory ?
Because cumulative memory requirement of multiple processes in an OS can be more than the system capacity.

It is basically for sharing memory resources between multiple processes.

It also provides security by isolating processes.

Pages :
- Minimum addressable virtual space that can be requested from OS.
- Typically 4KB

Swapping :
In case of out of memory, process memory will be evicted to the disc.

Page faults :
Happens when the page is not on physical memory but on the swap file which is on the hardware.

VIRTUAL ADDRESS SPACE

SYSTEM MEMORY ADDRESS SPACE

SWAP FILE ON DISC

VIRTUAL MEMORY ADDRESS TRANSLATION

Address translation & Page table
CPUs work with virtual addresses and those addresses need to be converted to physical addresses.

Page table structures on system memory are used for this purpose.

TLB (Translation lookaside buffer)
TLBs are caches in CPUs to make the translation process faster. Modern CPUs have multiple levels of TLBs.

(Intel refers to L2 TLB as sTLB)

TLB Shadowing
See the multicore section below.

ITLB
Apart from data TLB, there is also ITLB for caching addresses of instructions on both Intel and AMD architectures.

TLB Hit: No need to access page table

TLB Miss: If not found on TLB, we need to go to the system memory which is slower

OS support
Linux implementation refers to them as huge pages and Windows calls them as large pages.

You should check your OS and CPU in combination to find out the supported states.

PAGE TABLE ON-THE-SYSTEM MEMORY MANAGED BY OS

Each page table entry is modified by any of the cores, that particular TLB entry is invalidated in all cores via PIs. This one is not done by hardware but initiated by operating system.

PI: Interprocessor interrupt, you can take "processor" core in core context.

CACHE COHERENCY : FALSE SHARING AND CACHE PING-PONGING

In the diagram to the right, if Core1 changes its var1, that change will need to be propagated to all other cores by the cache coherency protocol. That will lead to invalidation of cache lines associated with the shared cache line across all cores, even though it is used by only one core.

That situation is called **cache ping-pong**.

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

Whenever a page table entry is modified by any of the cores, that particular TLB entry is invalidated in all cores via PIs. This one is not done by hardware but initiated by operating system.

PI: Interprocessor interrupt, you can take "processor" core in core context.

VIRTUAL MEMORY PAGE TABLE COHERENCY : TLB SHOOTDOWNS

One of the cores modifies a table entry

PAGE TABLE ON SYSTEM MEMORY

Core1 TLB Core2 TLB

var1 var2

Shared system memory cache line holding var1 for Core1 and var2 for Core2

That situation is called **cache ping-pong**.

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

Whenever a page table entry is modified by any of the cores, that particular TLB entry is invalidated in all cores via PIs. This one is not done by hardware but initiated by operating system.

PI: Interprocessor interrupt, you can take "processor" core in core context.

VIRTUAL MEMORY PAGE TABLE COHERENCY : TLB SHOOTDOWNS