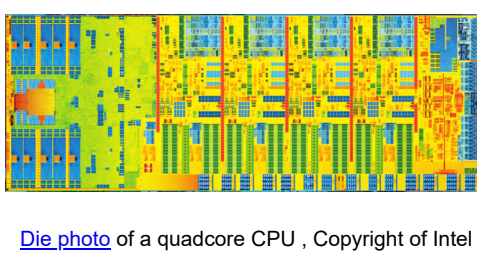


MICROARCHITECTURE CHEAT SHEET

X86 CPUs & Performance



This photo of a quad-core CPU. Copyright of Intel

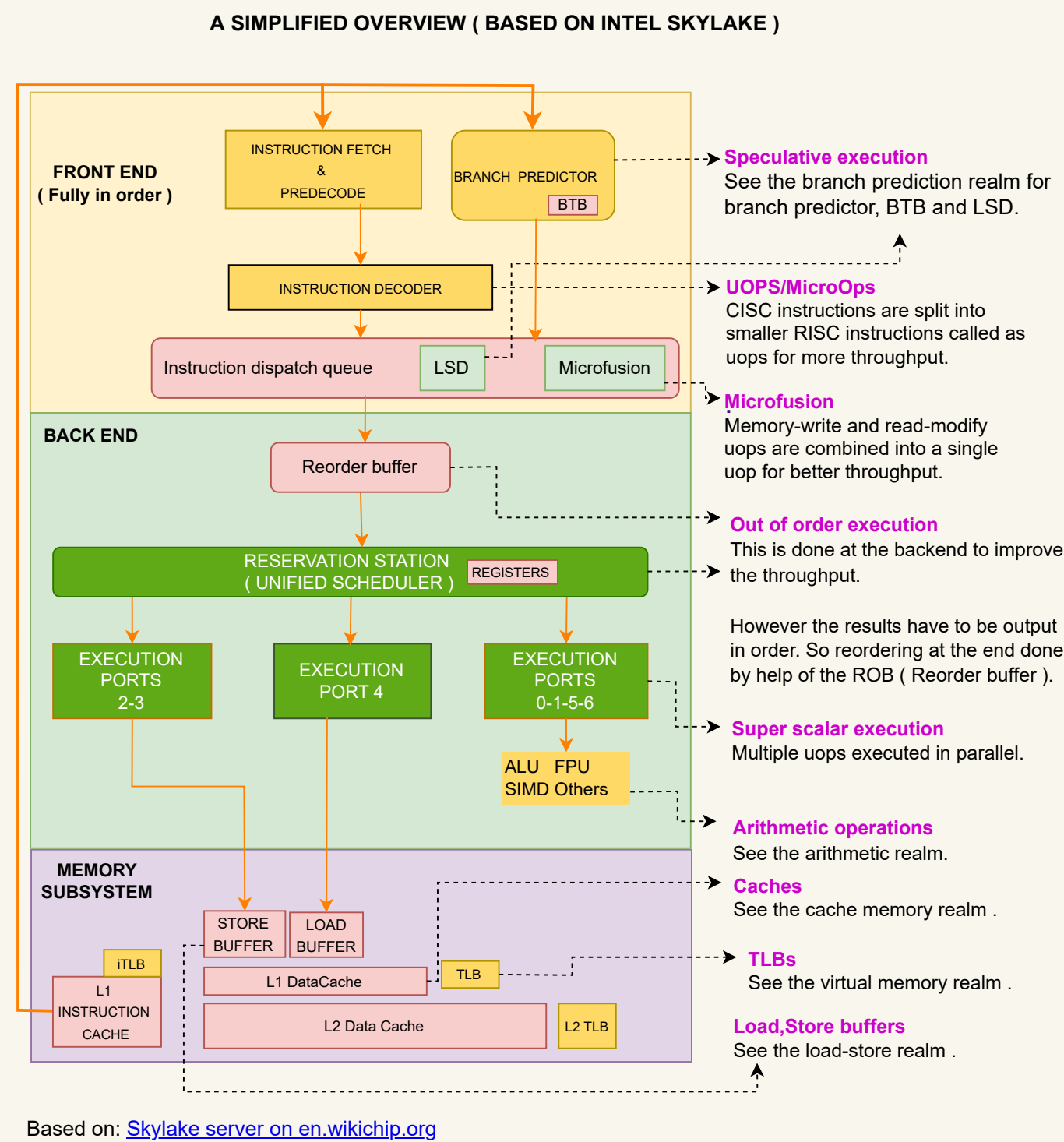
LAST UPDATE DATE : 18 OCT 2022

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

AUTHOR: AKIN OCAL akin_ocal@hotmail.com



PIPELINE REALM : INSIDE AN INDIVIDUAL CORE



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 [UCA](https://en.wikichip.org/wiki/analysis), and they represent parallel execution through cycles.

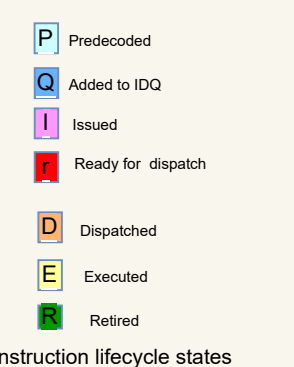
Rows are multiple instructions being executed at the same time.

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/flushed as they were wrongly speculated. On the other hand, executed instructions are the ones which were flushed. Therefore a high rate of retired instructions indicates low branch prediction rate.



Instruction lifecycle states in UCA diagrams

CONTENTION FOR EXECUTION PORTS IN THE PIPELINE

Instruction	Possible Ports	Actual Port	Cycle
mov rax, rax	1, 2, 3, 4	1	1
add rax, rax	1, 2, 3, 4	2	2
add rax, rax	1, 2, 3, 4	3	3
add rax, rax	1, 2, 3, 4	4	4

In the example above, all instructions are working on different registers, but SHR, ADD, DEC instructions are competing for ports 3 and 4. 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 Bashvalov'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 red coloured one, 2 instructions are competing for RAX register and notice that the second instruction gets executed after the first one.

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

RDTSOP INSTRUCTION FOR MEASUREMENTS

RDTSOP 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 RDTSOP combination in older systems that don't support RDTSOP.

ESTIMATING INSTRUCTION LATENCIES

Based on Agner Fog's [Instruction tables](https://www.agner-fog.com/instruction-tables), RDTSOP 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/tools/oneapi/whitepapers/2019-05-01-hyper-threading-technology-whitepaper.pdf), 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://www.intel.com/content/www/us/en/developer/tools/oneapi/whitepapers/2019-05-01-hyper-threading-technology-whitepaper.pdf) 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/tools/oneapi/whitepapers/2019-05-01-hyper-threading-technology-whitepaper.pdf) or AMD's [TurboCore](https://www.amd.com/en/techblog/2019-05-01-hyper-threading-technology-whitepaper.pdf) 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)

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 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.intel.com/content/www/us/en/developer/tools/oneapi/whitepapers/2019-05-01-hyper-threading-technology-whitepaper.pdf), 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.

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)

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).

ARITHMETIC REALM

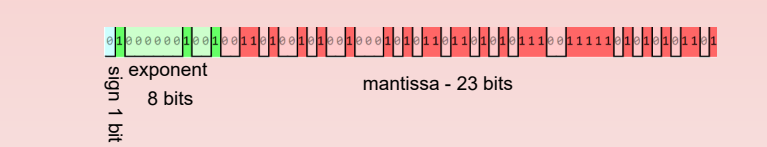
ARITHMETIC INSTRUCTION LATENCIES

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

Bitwise operations : Integer addition : 0.25 to 1 clock cycle
Floating point add : 3 clock cycles
Floating point multiplication : about 3 clock cycles
Floating point division : about 14 to 20 clock cycles
Integer division : 24-60 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 as visualizer :



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

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 $x \neq 0$, but $x < 0$.

Without denormal the code to the right would make a divide-by-zero exception. Reference : [Blaze 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, it takes 128 clock cycles on Skylake. They also can be turned off on Intel CPUs.

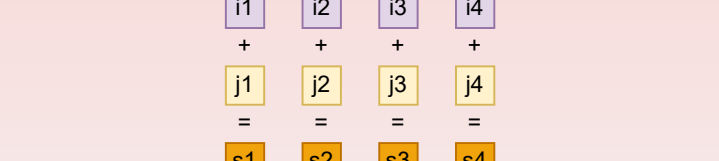
As for AMD side, the recent Zen architecture CPUs seemingly don't have the same performance degradation.

X86 EXTENSIONS

X86 extensions are specialised instructions. They have various categories such as [controlflow](https://en.wikichip.org/wiki/skylake_server) and [register-related operations](https://en.wikichip.org/wiki/skylake_server).

For the list of extensions : 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 go :



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 additions are executed by a single instruction.

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

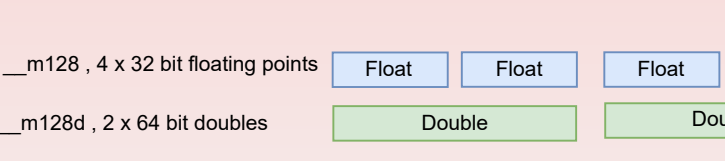
Apart from arithmetic operations, they can be utilized for string operations as well. A SIMD based JSOX parser <https://github.com/ahvin/jssox>

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 :



Note that SSE instructions require more power, therefore their usage may also introduce downclocking. They should be benchmarked : [Daniel Lemire'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) : callnt

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

BP METHODS : 2-LEVEL ADAPTIVE BRANCH PREDICTION

Saturating counter
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 rows 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** 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://en.wikichip.org/wiki/skylake_server).

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 LSD 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 jumps, with up to 15 taken branches, and up to 15 64-byte fetch lines.
- No CALL or RET.
- No non-mathematical stack operations (e.g., more PUSH than POP).
- More than 10 iterations.

Note that LSD is disabled on Skylake Server CPUs. Reference : https://en.wikichip.org/wiki/skylake_server

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

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

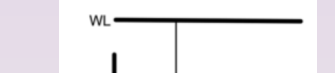
Reference : [Mark Malozki's article on Cloudflare blog](https://en.wikichip.org/wiki/skylake_server)

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 in the price as it has less components

SRAM used in cache memories



Access time : Under 1 nanosecond

Cost : Expensive in the price due to 6 transistors

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 have data caches. But there is also **instruction cache** (I-Cache) which store program instructions rather than data to improve throughput of CPU frontend.

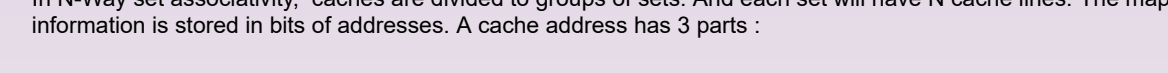
In case of a cache hit, the latency is typically single digit nanoseconds. And in case of a cache miss, we need a round trip to the system memory and local latency becomes 3 digit nanoseconds.



N-WAY SET ASSOCIATIVITY

Cache capacities are much smaller than the system memory. Moreover, softwares can use various regions of their address space. So if there were one to one mapping of a fully sequential memory that would lead to cache misses most of the time. Therefore there is a need for efficient mapping between the cache memory and the system memory.

In N-Way set associativity, caches are divided to groups of sets. And each set will have N cache lines. The mapping information is stored in bits of addresses. A cache address has 3 parts :



used as unique identifier per cache line
used to determine the set in the cache
used to determine the actual bytes in the target cache line

The pseudocode below shows steps for searching a single byte in the cache memory :

Get tag, set and offset from the address

For each line in the current set (which we just found out)
if tag of the current line equals to tag (which we just found out)
read and return data using offset of CACHE HIT
if there was no matching tag, it is a cache miss

The level of associativity (the number of ways) is a trade off between the search time and the amount of system memory we can map.

DIRECT CACHE ACCESS

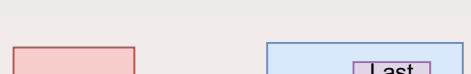
Modern NICs come with a DMA (Direct Memory Access) engine and can transfer data directly to drivers' ring buffers which reside on the system memory.

DMA mechanism doesn't require CPU involvement. Though mechanism initiated by CPU, therefore CPU support needed.

DCA bypasses the system memory and can transfer to directly LLC of CPUs that support this feature.

Intel refers to their technology as DIO (Direct IO).

Reference : [Intel documentation](https://en.wikichip.org/wiki/skylake_server)



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 removed from OS.
- Typically 4KB

Swapping
Happens when the page is not on physical memory but on the swap file which is on the harddrive.

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

TLB PRESSURE & HUGE PAGES

TLB pressure
If each page is 4K, that increases the load on the TLB buffer.

CPU support for larger pages
x86-64 CPUs support huge pages from 2MB to 1GB to reduce the pressure on TLB.

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

You shall check your OS and CPU in combination to find out the supported sizes.

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

VIRTUAL MEMORY ADDRESS TRANSLATION

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

Page table structure 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 Shootdowns
See the multicore realm below.

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

PAGE TABLE WALKING

Even with pages which group addresses, having all pages in a page table would still need too much storage on 64 bit systems. Therefore page tables are implemented hierarchically.

Memory is divided into address spaces. And there is a tree data structure for each address space in the page table. Processes have to walk the page table level by level in the hierarchy to find out the actual address.

4 level page tables is the most common one. In the diagram above first 48 bits of a 64 bit address are used for page table walking. All of 48 bits have to be used in order to find out the final actual address.

Intel CPUs is started to support 5 level tables since Ice Lake.

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

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

