

LAST UPDATE DATE : 18 OCT 2022

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

AUTHOR: AKIN OCAL      [akin\\_ocal@hotmail.com](mailto:akin_ocal@hotmail.com)      

## DYNAMIC CLOCK SPEEDS

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

Also **ACPI** defines multiple power states and Modern CPUs implement those. P-State is a state for performance and C states are for energy efficiency.

You can use Intel's [TurboBoost](#) or AMD's [TurboCore](#) to maximise the CPU usage.

Note that SSE usage may also introduce downclocking, therefore they should be used carefully : [Daniel Lemire's article](#)

## STORE-TO-LOAD FORWARDING & L1S & PERFORMANCE

Based on [CISC Organization Number 3.6.4](#), store-to-load forwarding may improve combined latency of these two operations. The reason is not as self-evident however, as it potentially **L1S** (Load-Hit-Store) problem in which the penalty is a round trip to the cache memory :

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

```

graph TD
    P[Processor] -- STORE --> L1[L1 Cache]
    L1 -- HIT --> P
    P -- LOAD --> L1
    L1 -- MISS --> P
    P -- LOAD --> L2[L2 Cache]
    L2 -- HIT --> P
    
```

There are several conditions for the forwarding to happen. In case of a **STORE** followed by a **LOAD** (or **STORE** and **3**) a round trip to the cache memory will be bypassed.

The conditions for a successful forwarding and latency penalties in case of no-forwarding can be found in Agner Fog's [microarchitectures](#)

Previous game consoles PlayStation3 & Xbox360 had PowerPC based processors which did in-order-execution rather than out-of-order execution. This means that instructions are executed in the order they are using [restrict](#) keyword and other things. [Evan Raskin'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 tables](#) & Skylake architecture on 64 bit registers.

Bitwise operations	Integer divide: 0.25 to 1 clock cycle
<b>Floating point add: 0.5 clock cycles</b>	
<b>Floating point multiply: 0.5 to 1 clock cycles</b>	
<b>Floating point divide: about 14-16 clock cycles</b>	
<b>Integer divide: 24-90 clock cycles</b>	

## FLOATING POINT

X86 uses [IEEE 754](#) standard for floating points. A 32 bit floating point consists of 3 parts in the binary layout. Below you can see bits of 12345678 FP number. Used [agner fog's integer724-visualization](#) as visualizer:

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

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

As floating points are approximations, denormal numbers are needed to avoid undesired cases of  $a \cdot b = 0$  but  $a \neq 0$  and  $b \neq 0$ . Without denormals the code to the right would make a divide-by-zero exception. Reference: [Bruce Dawson's article](#).

```
float DeliversZero(Denorm float a, float b)
{
    if (a != 0)
    {
        Without denormals the code to the right
        would make a divide-by-zero exception.
        return 0.0f;
    }
}
```

Based on Agner Fog's [microarchitectures](#), Intel CPUs have a penalty for denormal numbers, as e.g. 128 clock cycles on Skylake. They 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 specialized instructions. They have various categories such as [cryptographic](#) and [neural network operations](#).

For the list of extensions: [https://en.wikipedia.org/wiki/X86\\_extensions](https://en.wikipedia.org/wiki/X86_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 wider data types. The data type diagrams below are for 128 bit AIX:

**X86 EXTENSIONS : SIMD DETAILS**

The most recent [SIMD instructions](#) and their corresponding registers are:

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

## BRANCH PREDICTION BASICS

CPUs proactively fetch instructions from potentially upcoming branches to utilise the pipeline as much as possible.

**High if predicted correctly** : if the right branch was predicted that will increase the throughput as it completes fetching of 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 having the pipeline.

**What are the branch instructions?** : Unconditional ones (jnp), conditional ones (eg: jge), called **branch instructions**.

How : There are auxiliary hardware buffers.

Branch target addresses stored in auxiliary ( instruction pointers of branch instructions). ARM uses multiple level of BTBs : L1 BTB, L2 BTB etc.

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

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

	1	2	3	4	...
branch 1	T	T	NT	T	...
branch 2	NT	NT	NT	T	...
branch n	T	NT	NT	NT	...

## BP METHODS : 2-LEVEL ADAPTIVE BRANCH PREDICTION

**Saturating counter**

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

Whenever a branch is taken it goes stronger.

And whenever a branch is not taken it goes weaker.

↑ Not taken

↑ Not taken

↑ Not taken

Strongly not taken

Weakly not taken

Weakly taken

Strongly taken

↓ Taken

↓ Taken

↓ Taken

**2 level adaptive predictor**

In this method, you store the history of not taken occurrences in a history register which is 1 bit.

Also you create a table called "pattern history table" for that branch. That pattern history table keeps 2<sup>n</sup> ones and each one has a saturating counter.

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

Reference : Armer Exp's [microprocessors book](#)

### DRAM used in system memories

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

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](#)

## 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.

Load & Store Buffers

L1 Data Cache

...

LLC : Last level cache

System memory

3 level caches are currently the most common ones. Intel Broadwell architecture has 4 level caches in **part**. It is expected that upcoming AMD CPUs **will come with 4 levels of caches**.

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

All the mentioned caches till now were data caches. But there is also **instruction cache** (Cache) which stores program instructions rather than data to improve throughput of CPU further.

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

**VIRTUAL MEMORY ORGANISATION**

The diagram shows the flow of memory access. On the left, a vertical stack of 16 yellow rectangles represents the **VIRTUAL ADDRESS SPACE**, with labels 0K, 4K, 8K, and 12K on the left. Arrows point from the 0K, 4K, 8K, and 12K addresses to the first four rows of a 4x4 grid of pink rectangles representing the **SYSTEM MEMORY ADDRESS SPACE**. An arrow from the 12K address points to a purple oval labeled **SWAP FILE ON DISC**.

**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-x64 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.

The diagram shows a 4x4 grid of 16 pink rectangles, each labeled **4K**. Below this grid is a section titled **Regular pages** showing two large yellow rectangles, each labeled **1 GB**. Below these are two more yellow rectangles, each labeled **Huge pages**.

## CACHE COHERENCY PROTOCOLS

As LLC is typically shared by multiple cores, cache coherency protocols are needed to avoid data hazards. Intel CPUs use MESIF and ARM 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 shared in other cache lines and matches its value in system memory.  
**Invalid** – cache line is not valid.

Allowed state transitions in MESI protocol:

	M	E	S	I
M				
E	N	N	N	Y
S	N	N	Y	Y
I	Y	Y	Y	Y

As for state transition costs, M and E 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 wherever it is applicable.

Intel MESIF: [https://en.wikipedia.org/wiki/MESI\\_protocol](https://en.wikipedia.org/wiki/MESI_protocol)  
 ARM MOESI: [https://en.wikipedia.org/wiki/MOESI\\_protocol](https://en.wikipedia.org/wiki/MOESI_protocol)

## MEMORY REORDERINGS

The term memory ordering refers to the order in which the processor issues reads (loads) and writes stores. Based on [Intel Software Developer's Manual Volume 3](#), 8.2.3.4, there is only one kind of memory ordering that can happen. Loads can be reordered with earlier stores if they use different memory locations. That reordering will not happen if they use the same address as follows:

CORE1		CORE2
$x$ and $y$ initially 0 $\text{mov } [x], 1; \text{ STORE to } X$ $\text{mov } [\text{result}], y; \text{ LOAD from } Y$		$x$ and $y$ initially 0 $\text{mov } [y], 1; \text{ STORE to } Y$ $\text{mov } [\text{result2}], x; \text{ LOAD from } X$

In this case of reordering, result1 and result2 above can both end up as zero in both processors. Note that, apart from CPUs, also compilers can do memory reordering. <https://wreshing.com/2010/02/25/memory-ordering-at-compile-time/>

## ALLOCATING A PARTITION OF LLC

You can allocate a partition of the shared CPU last level cache for your performance sensitive application to avoid evictions on Intel CPUs that support **CAT** feature.

CAT : Cache allocation tech. , reference : [Intel CAT page](#)

For AMD equivalent, AMD64 QOS Extensions introduced starting from [https://developer.amd.com/wp-content/resources/56375\\_1.00.pdf](#)

The diagram illustrates the allocation of a partition of the Last Level Cache (LLC) for performance-sensitive applications. It shows two cores, CORE 1 and CORE 2, sharing a common LLC. A portion of the LLC is designated as a 'PERF CRITICAL CORE' partition, which is dedicated to only one core (CORE 2).

LLC

LLC cache lines shared by non performance critical cores

PERF CRITICAL CORE

LLC cache lines dedicated to only one core

**SMP (Symmetrical multiprocessing)**  
All CPUs use a single bus to access the same system memory. CPUs may slow down each other as there may be a contention for access to the banks.

**NUMA (Non uniform memory access)**  
The system memory is organized as nodes and each node in memory is to be accessed by only one CPU, therefore memory access is split. So CPUs don't compete with each other.

**Cross NUMA access**  
A CPU can access to its local memory faster. However that also means that a cross-access will be slower. The side screenshot is from Intel's memory locality checker. Notice that cross accesses are two times slower than local accesses.

	Accessing node	Latencies (in ns)...
Socket 0	0	1
	1	2
Socket 1	0	67.5
	1	135.2
Node 1	1	126.5
	0	60.5

Reference for image:  
[Intel's paper for memory locality checker](#)

[illegible]

The diagram illustrates the hierarchy of computer storage technologies, showing how hardware cost increases as you go up the pyramid and latency increases as you go down. The technologies are categorized into five main groups, each with a corresponding performance metric.

Storage Technology	Performance Metric
CPU REGISTER	PICOSECONDS
CPU L1 CACHE SRAM	NANOSECONDS
CPU L2-L3, L1C CACHED, TYPICALLY DRAM	NANOSECONDS
RAM for memory & DRAM only for storage, TYPICALLY DRAM	UP TO 100 NANOSECONDS
INTERCONNECT: SDR	UP TO 100 MICROSECONDS
NVM SSD for storage & swap-memory INTERCONNECT: PCIe/express	UP TO 100 MICROSECONDS
SSD for storage & swap-memory INTERCONNECT: SATA	UP TO 100 MICROSECONDS
MECHANICAL HDD for storage & swap memory INTERCONNECT: SATA	UP TO 100 MILLISECONDS

Simple register operation (ADD OR etc.)	less than 1 clock cycle	Floating point addition	10-40 clock cycles
Branch prediction	1-2 clock cycles	Compare and exchange	15-30 clock cycles
Floating point addition	1-3 clock cycles	Integer division	15-40 clock cycles
Multiplication (int. floating point)	1-7 clock cycles	L3 data cache read	70-100 clock cycles
L1 data cache read	3-4 clock cycles	System memory read	100-150 clock cycles
TLB miss	7-21 clock cycles	Cross-NUMA L3 read	100-300 clock cycles
L2 data cache read	10-12 clock cycles	Cross-NUMA system memory read	300-600 clock cycles
Branch misprediction	10-20 clock cycles	Reference for numbers:	
		<a href="http://libreoffice.com/contributors-operation-clocks-in-cpu-clocks-cycles/">http://libreoffice.com/contributors-operation-clocks-in-cpu-clocks-cycles/</a>	