

MICROARCHITECTURE CHEAT SHEET

X86 CPUs & Performance

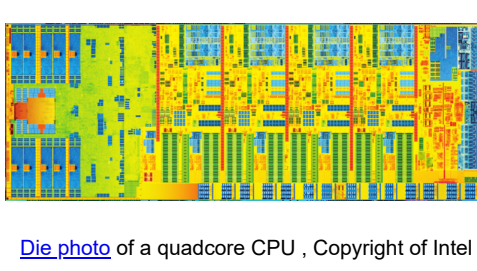


Photo of a quad-core CPU. Copyright of Intel

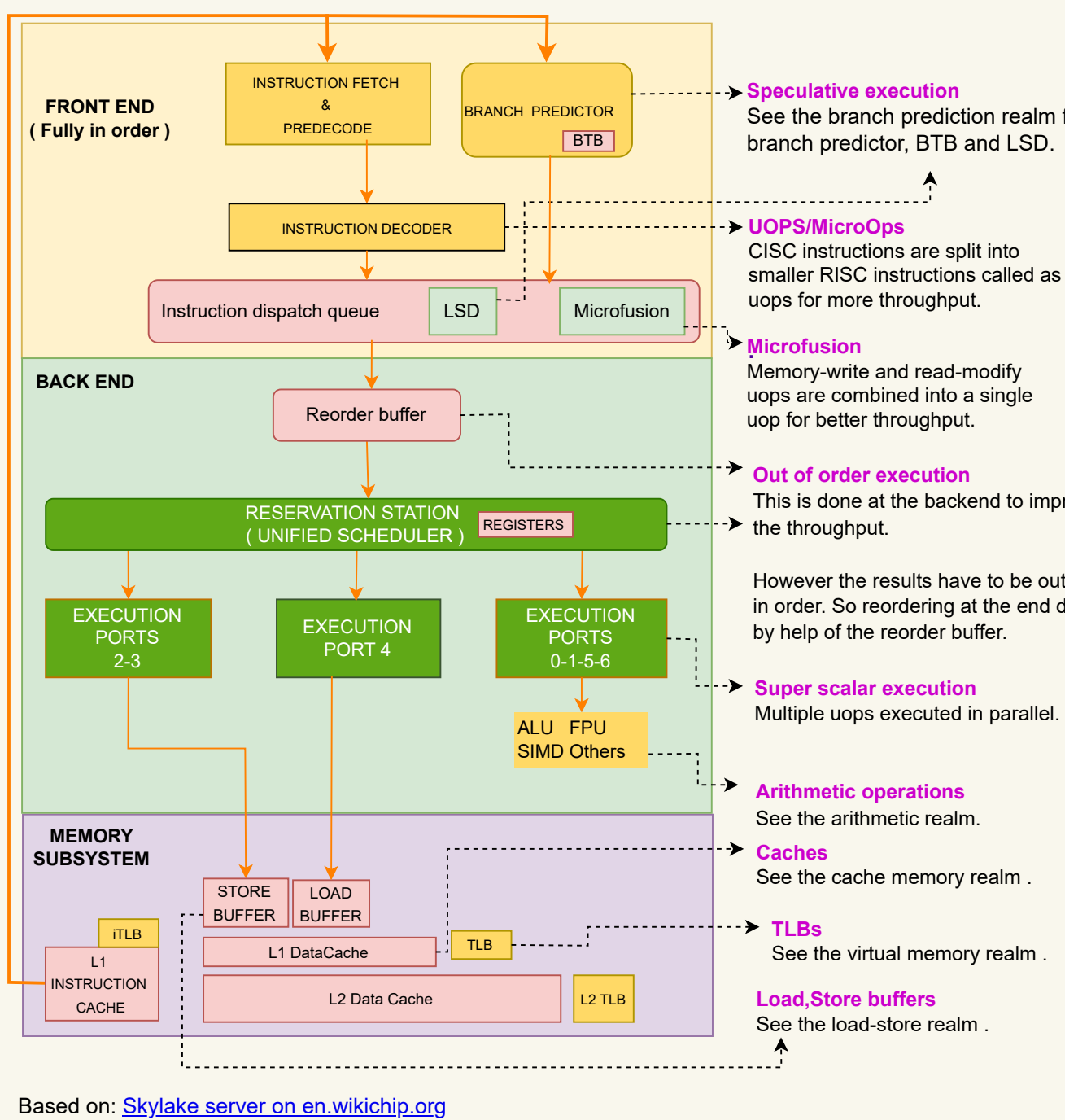
LAST UPDATE DATE : 11 NOV 2022

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

AUTHOR: AKIN OCAL skinocal@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/wiki/skylake_server_on_en.wikichip.org)

Speculative execution
See the branch predictor realm for branch predictor, BTB and LSD.

UCPIS Microfusions
CISC instructions are split into smaller RISC instructions called as UCPIs for more throughput.

Microfusion
Memory-write and read-modify UCPIs are combined into a single UCPI 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 UCPIs executed in parallel.

Arithmetic operations
See the arithmetic realm.

Caches
See the cache memory realm.

TLBs
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 [UICA](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org) and they represent parallel execution through cycles.

Rows are multiple instructions being executed at the same time.

Cumulative 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.

CONTENT FOR EXECUTION PORTS IN THE PIPELINE

Instruction	ports	Actual Port	Cycle
add rax, rax	1	1	1
add rax, rax	2	2	2
add rax, rax	3	3	3
add rax, rax	4	4	4
add rax, rax	5	5	5
add rax, rax	6	6	6
add rax, rax	7	7	7
add rax, rax	8	8	8
add rax, rax	9	9	9
add rax, rax	10	10	10
add rax, rax	11	11	11
add rax, rax	12	12	12
add rax, rax	13	13	13
add rax, rax	14	14	14
add rax, rax	15	15	15
add rax, rax	16	16	16
add rax, rax	17	17	17
add rax, rax	18	18	18
add rax, rax	19	19	19
add rax, rax	20	20	20
add rax, rax	21	21	21
add rax, rax	22	22	22
add rax, rax	23	23	23
add rax, rax	24	24	24
add rax, rax	25	25	25
add rax, rax	26	26	26
add rax, rax	27	27	27
add rax, rax	28	28	28
add rax, rax	29	29	29
add rax, rax	30	30	30
add rax, rax	31	31	31

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://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org)

INSTRUCTION STALLS DUE TO DATA DEPENDENCY

Instruction	ports	Actual Port	Cycle
add rax, rax	1	1	1
add rax, rax	2	2	2
add rax, rax	3	3	3
add rax, rax	4	4	4
add rax, rax	5	5	5
add rax, rax	6	6	6
add rax, rax	7	7	7
add rax, rax	8	8	8
add rax, rax	9	9	9
add rax, rax	10	10	10
add rax, rax	11	11	11
add rax, rax	12	12	12
add rax, rax	13	13	13
add rax, rax	14	14	14
add rax, rax	15	15	15
add rax, rax	16	16	16
add rax, rax	17	17	17
add rax, rax	18	18	18
add rax, rax	19	19	19
add rax, rax	20	20	20
add rax, rax	21	21	21
add rax, rax	22	22	22
add rax, rax	23	23	23
add rax, rax	24	24	24
add rax, rax	25	25	25
add rax, rax	26	26	26
add rax, rax	27	27	27
add rax, rax	28	28	28
add rax, rax	29	29	29
add rax, rax	30	30	30
add rax, rax	31	31	31

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://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org)

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 read out-of-order execution effects to a degree : [From Intel Software Developer's Manual Volume 4, 4.4. April 2022](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org)

How to benchmark code execution times : [whisperer](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org) has details of using RTDSCP instruction.

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

ESTIMATING INSTRUCTION LATENCIES

You can use Agner Fog's [Instruction Latency](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org) 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 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 including caches and execution units are shared Agner Fog's [microarchitecture book](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org) has the "multithreading" sections for each of Intel and AMD microarchitectures.

Regarding using it, if your app is data-intensive, "haved caches won't help. Therefore 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.

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 **variability options** and the most well known is TurboBoost. On AMD side there is **TurboCore**. You can use those to maximise the CPU usage.

Number of active cores & SIMD AVX2/SSE12 on Intel CPUs : Intel's power management policies are complex. See the arithmetic and the multicore realms as number of active cores and some of AVX2/SSE12 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.wikimedia.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. This issue is described in en.wikimedia.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 and separately handle LHS by using using keyword and other methods : Alan Rusakov's article
```

STORE-TO-LOAD FORWARDING & LHS & PERFORMANCE

Based on [Intel Optimisation Manual 3.6.4](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org), 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.

Reference : https://en.wikimedia.org/wiki/Memory_disambiguation

STORE-TO-LOAD FORWARDING

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 non-forwarding can be found in Agner Fog's [microarchitecture book](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org).

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 **using** keyword and other methods : [Alan Rusakov's article](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org)

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 Latency](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org) & 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](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org) standard for floating points. A 32 bit floating point consists of 3 parts in the memory layout. Below you can see all bits of 124 5879 FP number. [Using agner fog's ieee754-visualizer](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org) 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 a+b but a-b.

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

Based on Agner Fog's [microarchitecture book](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org), 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 [agnerfog's tutorial network operations](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org).

Intel Intrinsics 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 of integers (1 to 14) are added to another array of integers (1 to 14). The result is also an array of sums (1 to 14). In this example, 4 add operations are executed by a single instruction.

They play key role in complex vectorisation optimisations : [GCC auto vectorization](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org)

Apart from arithmetic operations, they can be utilised for string operations as well. A SIMD based JSON parser : [github.com/akhin/jsonsimd](https://en.wikichip.org/wiki/skylake_server_on_en.wikichip.org)

X86 EXTENSIONS : SIMD DETAILS

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

AVX1 : 128 bits, XMM registers
AVX2 : 256 bits, YMM registers
AVX512 : 512 bits, ZMM registers

Note about AVX512 : AMD started to implement it with Zen4.

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

__m128 : 4 x 32 bit floating points

__m128d : 2 x 64 bit doubles

__m128i : 4 x 32 bit ints

__m128i : 2 x 64 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/wiki/skylake_server_on_en.wikichip.org)

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/jmp.

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 : T T T T T
branch 2 : NT NT NT NT NT
branch 3 : T NT NT NT NT

branch 4 : T NT NT NT NT

branch 5 : T NT NT NT NT

branch 6 : T NT NT NT NT

branch 7 : T NT NT NT NT

branch 8 : T NT NT NT NT

branch 9 : T NT NT NT NT

branch 10 : T NT NT NT NT

branch 11 : T NT NT NT NT

branch 12 : T NT NT NT NT

branch 13 : T NT NT NT NT

branch 14 : T NT NT NT NT

branch 15 : T NT NT NT NT

branch 16 : T NT NT NT NT

branch 17 : T NT NT NT NT

branch 18 : T NT NT NT NT

branch 19 : T NT NT NT NT

branch 20 : T NT NT NT NT

branch 21 : T NT NT NT NT

branch 22 : T NT NT NT NT

branch 23 : T NT NT NT NT

branch 24 : T NT NT NT NT

branch 25 : T NT NT NT NT

branch 26 : T NT NT NT NT

branch 27 : T NT NT NT NT

branch 28 : T NT NT NT NT

branch 29 : T NT NT NT NT

branch 30 : T NT NT NT NT

branch 31 : T NT NT NT NT

branch 32 : T NT NT NT NT

branch 33 : T NT NT NT NT

branch 34 : T NT NT NT NT

branch 35 : T NT NT NT NT

branch 36 : T NT NT NT NT

branch 37 : T NT NT NT NT

branch 38 : T NT NT NT NT

branch 39 : T NT NT NT NT

branch 40 : T NT NT NT NT

branch 41 : T NT NT NT NT

branch 42 : T NT NT NT NT

branch 43 : T NT NT NT NT

branch 44 : T NT NT NT NT

branch 45 : T NT NT NT NT

branch 46 : T NT NT NT NT

branch 47 : T NT NT NT NT

branch 48 : T NT NT NT NT

branch 49 : T NT NT NT NT

branch 50 : T NT NT NT NT

branch 51 : T NT NT NT NT

branch 52 : T NT NT NT NT

branch 53 : T NT NT NT NT

branch 54 : T NT NT NT NT

branch 55 : T NT NT NT NT

branch 56 : T NT NT NT NT

branch 57 : T NT NT NT NT

branch 58 : T NT NT NT NT

branch 59 : T NT NT NT NT

branch 60 : T NT NT NT NT

branch 61 : T NT NT NT NT

branch 62 : T NT NT NT NT

branch 63 : T NT NT NT NT

branch 64 : T NT NT NT NT

branch 65 : T NT NT NT NT

branch 66 : T NT NT NT NT

branch 67 : T NT NT NT NT

branch 68 : T NT NT NT NT

branch 69 : T NT NT NT NT

branch 70 : T NT NT NT NT

branch 71 : T NT NT NT NT

branch 72 : T NT NT NT NT

branch 73 : T NT NT NT NT

branch 74 : T NT NT NT NT

branch 75 : T NT NT NT NT

branch 76 : T NT NT NT NT

branch 77 : T NT NT NT NT

branch 78 : T NT NT NT NT

<