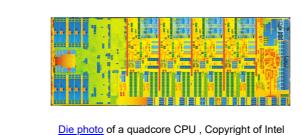
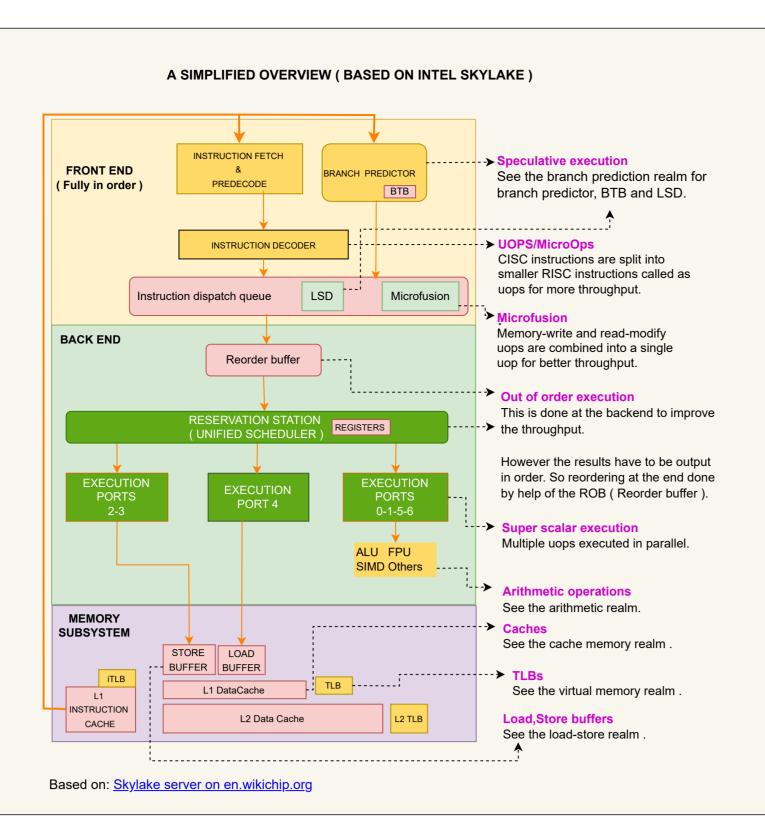
X86 CPUs & Performance



LAST UPDATE DATE: 20 OCT 2022 FOR LATEST VERSION: www.github.com/akhin/microarchitecture-cheatsheet AUTHOR: AKIN OCAL akin ocal@hotmail.com



INSIDE **INDIVIDUAL** CORE



PIPELINE PARALLELISM & PERFORMANCE Pipeline diagrams: The diagrams below in the following topics are outputs from an online microarchitecture analysis tool <u>UICA</u> 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 cyle". A higher IPC value usually means a better throughput You can measure IPC with perf : <a href="https://perf.wiki.kernel.org/index.php/Tutorial">https://perf.wiki.kernel.org/index.php/Tutorial</a> Instruction lifecycle states in UICA diagrams Rate of retired instructions: Apart from IPC, number of retired instructions should be checked. Retired instructions are not committed/finalised as they were wrongly speculated. On the other hand executed instructions are the ones which were finalised. Therefore a high rate of retired instructions indicates low branch prediction rate. CONTENTION FOR EXECUTION PORTS IN THE PIPELINE after ADD instruction.

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

Also notice that there is longer time between E(executed) and R(retired) states of instruction ADD as retirement has to be done in-order whereas execution is out-of-order. Reference : Denis Bakhvalov's article 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 Bakhvalov's article

RDTSCP INSTRUCTION FOR MEASUREMENTS RDTSCP 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 RDTSC combination in older systems that don't support RDTSCP. **ESTIMATING INSTRUCTION LATENCIES** Based on Agner Fog`s <u>Instruction tables</u>, RDTSCP reciprocal throughput (clock cycle per instruction) is 32 on Skylake microarchitecture: -> 1 cycle @4.5GHZ is 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, 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. For ex if your 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 only Intel. DYNAMIC CLOCK SPEEDS Modern CPUs employ dynamic frequency scaling which means there is a min and max Max level ◀ frequency per CPU core. Also ACPI defines multiple power states and C0 - Normal execution ← → Pn

C1 - Idle

Note that SSE usage may also introduce downclocking, therefore they should be used carefully :

modern CPUs implement those. P-State's are

for performance and C states are for energy

You can use Intel's <u>Turboboost</u> or AMD's

<u>Turbocore</u> to maximise the CPU usage.

Daniel Lemire's article

**LOAD STORE** 

**REALM** 

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 Reference: https://en.wikipedia.org/wiki/Memory\_disambiguation

**LOAD & STORE BUFFERS** 

STORE-TO-LOAD FORWARDING Using buffers for stores and loads to support out of order execution leads to a data syncronisation issue. That issue is described in en.wikipedia.org/wiki/Memory disambiguation#Store to load forwarding

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

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 :

STORE-TO-LOAD FORWARDING & LHS & PERFORMANCE Based on 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.wikipedia.org/wiki/Load-Hit-Store There are several conditions for the forwarding to happen. In case of a STORE BUFFER LOAD BUFFER successful forwarding, the steps 2 and 3 ( a roundtrip to the cache )

The conditions for a successful forwarding and latency penalties in case of no-forwarding can be found in Agner Fog's microarchitecture book.

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 restrict keyword and other methods : Elan Ruskin's article

will be bypassed.



ARITHMETIC INSTRUCTION LATENCIES You can see a set of arithmetic opertions 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 add/sub: 0.25 to 1 clock cycle

Floating point add: 3 clock cycles

Based on Agner Fog's microarchitecture book, Intel CPUs have a penalty for denormal As for AMD side, the recent Zen architecture CPUs seemingly don't have the same performance degradation.

numbers, for ex: 129 clock cycles on Skylake. They also can be turned off on Intel CPUs. Some typical application areas are 3D graphics and quantitative finance.

**FLOATING POINTS** 

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

mantissa - 23 bits

float GetInverseOfDiff(float a, float b)

return 1.0f / (a - b);

in the memory layout. Below you can see all bits of 1234.5678 FP

number. Used <u>bartaz.github.io/ieee754-visualization</u> as visualiser:

A floating point's value is calculated as: ±mantissa × 2 exponent

As floating points are approximations,

undesired case of : a!=b but a-b=0

Reference : Bruce Dawson's article

charge/discharge times and other steps

Cost: Cheaper in the price as it has less components

Without denormals the code to the right would invoke a divide-by-zero exception.

denormal numbers are needed to avoid an

integers (j1 to j4). The result is also an array of sums ( s1 to s4). In this example, 4 add

Cost: Expensive in the price

due to 6 transistors

X86 uses IEEE 754 standard for floating points. A 32 bit floating point consists of 3 parts x86 extensions are specialised instructions. They have various categories

operations are executed by a single instruction. Apart from arithmetic operations, they can be utilised for string operations as well. A SIMD based JSON parser : <a href="https://github.com/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/simdjson/si

X86 EXTENSIONS

SSE (Streaming SIMD Extensions) is one of the most important ones. SIMD stands

In the example above, an array 4 integers (i1 to i4) are added to another array of

for "single instruction multiple data". SIMD instructions use wider registers to

from cryptography to neural network operations

execute more work in a single go:

Intel Intrinsics Guide is a good page to explore those extensions.

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

m128 , 4 x 32 bit floating points Float Float Float Float Double \_\_m128d , 2 x 64 bit doubles \_\_m128i , 4 x 32 bit ints \_\_m128l , 2 x 64 bit long longs long long

In order to switch to Pstates, C-state

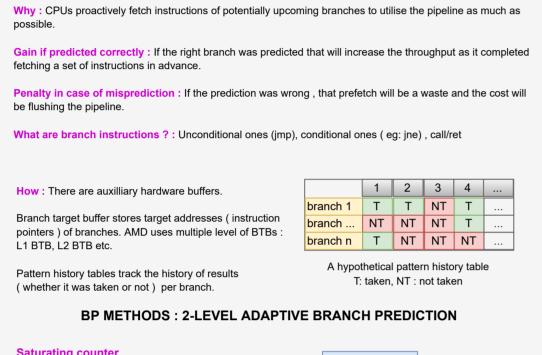
has to be brought

to C0 level

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

OFFSET

# **BRANCH PREDICTION** REALM



**BRANCH PREDICTION BASICS** 

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

2 level adaptive predictor In this method, you store the history of last n occurences in a history register which is n bits. Also you create a table called "pattern history table" for that branch. That pattern history table keeps 2<sup>n</sup> 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

#### **BP METHODS: AMD PERCEPTRONS** They are used in Zen arcihtectures. A <u>perceptron</u> is basically the simplest form of machine learning. They can be considered as a linear array of weights. Agner For mentions that they are good at predicting very long branches compared to 2-level adaptive branch prediction in his microarchitecture book. For details of perceptron based branch prediction: The output Y (in this case whether a branch taken or not ) is calculated by dot product of the weight

L1 CACHE

## <u>Dynamic Branch Prediction with Perceptrons by Daniel</u>

Note that LSD is disabled on Skylake Server CPUs. Reference:

https://en.wikichip.org/wiki/intel/microarchitectures/skylake\_(server)#Front-end

Intel LSD will detect a loop and stop fetching instruction to improve frontend bandwidth. Several conditions mentioned in Intel Optimization Manual • Loop body size up to 60 μops, with up to 15 taken branches, and up to 15 64-byte fetch lines. • No mismatched stack operations (e.g., more PUSH than POP). More than ~20 iterations.

INTEL LSD ( LOOP STREAM DETECTOR )

vector and the input vector.

**DISABLING SPECULATIVE EXECUTION PATCHES** You can consider disabling system patches for speculative execution related vulnerabilities such as Meltdown and Spectre for performance, if it is doable in your system. Kernel.org documentation: https://www.kernel.org/doc/html/latest/admin-guide/kernel-

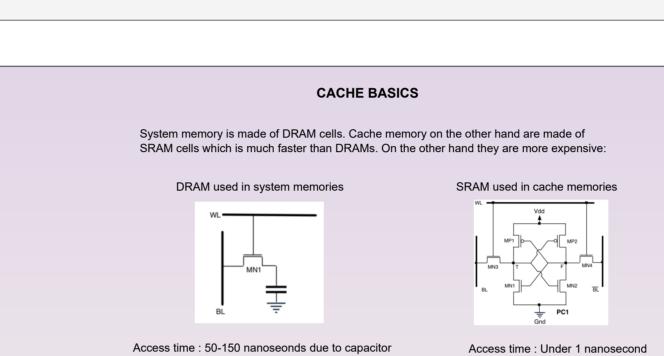
**ESTIMATED LIMITS: HOW MANY IFS ARE TOO MANY?** 

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

<u>parameters.html</u> Meltdown paper : <a href="https://meltdownattack.com/meltdown.pdf">https://meltdownattack.com/meltdown.pdf</a> Spectre paper : <a href="https://spectreattack.com/spectre.pdf">https://spectreattack.com/spectre.pdf</a>

Intel Xeon Gold 6262 -> roughly 4K AMD EPYC 7713 -> roughly 3K Reference: Marek Majkovski's article on Cloudflare blog

sequences of branch instructions:



Reference: Ulrich Drepper's What every programmer should know about memory

## **CACHE MEMORY REALM**

**CACHE ORGANISATION** Load & Store Buffers Caches are organised in multiple levels. As you go upper in that hierarchy , the capacity increases. Therefore **LLC** L1 Data Cache term used to indicate the last level of cache. 3 level caches are currently the most common ones. Intel LLC : Last level cache Broadwell architecture had 4 level caches in the past. Also upcoming AMD CPUs may come with 4 level of caches. System memory

A cache line is the smallest addresable unit in cache memories. It is typically 64 bytes. All the mentioned caches till now were data caches. But there is also in he (iCache) 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.

## N-WAY SET ASSOCIATIVITY Cache capacities are much smaller than the system memory. Moreover, softwares can use various regions of their

Get tag, set and offset from the address

address space. So if there was 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 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

SET

used as unique identifier used to determine used to determine the actual bytes the set in the cache in the target cache line The pseudocode below shows steps for searching a single byte in the cache memory:

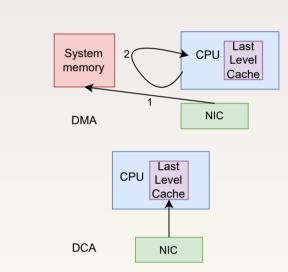
For each line in the current set ( which we just found out ) if tag of the current line line equals to tag ( which we just found out ) read and return data using offset // CACHE HIT If there was no matching tag, it is a cache miss The level of associtiavity (the number of ways) is a trade off between the search time and the amount of system

**DIRECT CACHE ACCESS** 

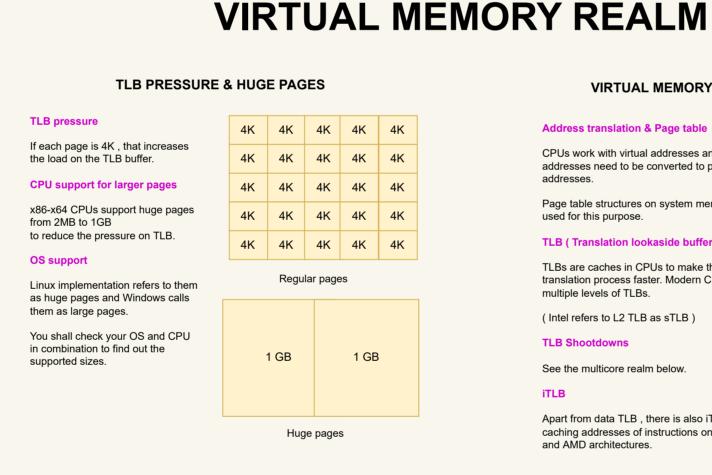
Modern NICs come with a DMA ( Direct Memory Access ) engine and can transfer data directly to drivers' ring buffers which reside on the system DMA mechanism doesn't require CPU involvement. Though mechanism initiated by CPU, therefore CPU support needed.

memory we can map.

DCA bypasses the system memory and can transfer to directly LLC of CPUs that support this feature. Intel refers to their technology as DDIO ( Direct I/O ). Reference : Intel documentation



#### VIRTUAL MEMORY ORGANISATION SYSTEM MEMORY ADDRESS SPACE VIRTUAL ADDRESS SPACE Why virtual memory? Because cumulative memory requirement of multiple processes in an OS can be more than 0K the system capacity. It is basically for sharing memory resources between multiple processes. It also provides security by isolating - Minimum addressable virtual space that can be regusted from OS. SWAP FILE - Typically 4KB ON DISC In case of out of memory, process memory will be evicted to the disc. Happens when the page is not on physical



# **VIRTUAL MEMORY ADDRESS TRANSLATION**

Address translation & Page table TLB Hit No need to access page table CPUs work with virtual addresses and those addresses need to be converted to physical Core Page table structures on system memory are used for this purpose. TLB (Translation lookaside buffer) TLB Miss: If not found on TLB, we need to TLBs are caches in CPUs to make the go to the system memory which is slower translation process faster. Modern CPUs have multiple levels of TLBs. (Intel refers to L2 TLB as sTLB) Core TLB **TLB Shootdowns** See the multicore realm below. PAGE TABLE ON THE SYSTEM

MEMORY

MANAGED BY OS

#### 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: 47 39 38 30 29 21 20 12 11 L1 Table L2 Table L3 Table L4 Table physical 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. ( For all details : Intel Software Developer's Manual Volume3 4.5)

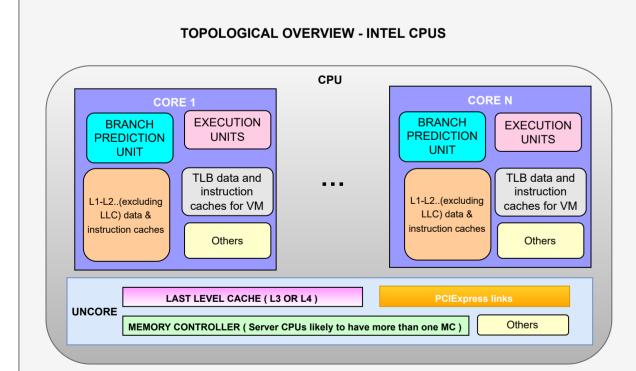
PAGE TABLE WALKING

Intel CPUs 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 table increases due to a new level of indirection.

#### SYSTEM MEMORY REALM DDR RAMs are the most common commodity hardware as system memory. They are found in forms of DIMMs ( Dual inline memory module ) / RAMSticks. A DIMM. Click for Image source RANK 1 RANK 1 System memory / RAM is organised as RANK ... PANK 1 collection of ranks. RANK N BANK 1 BANK N a DRAM BANK .. Each rank have banks which are collection of DRAM cells per bit. **DRAM** refreshes DRAM circuits use capacitors which lose their charge over time. (See the cache memory realm). So RAMs have to refresh their DRAM cells periodically. As for DDR4, refreshing is rank-level which means the other banks in the same rank become inaccesible. DDR5 comes with same-bank-refresh feature which allows a more finegrained bank-level refresh. Therefore it can offer a higher DDR4 refresh granularity DDR5 refresh granularity

## **MULTICORE REALM**

memory but on the swap file which is on the



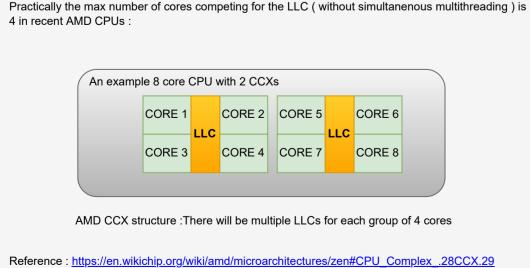
TOPOLOGICAL OVERVIEW - AMD CPUS Most of AMD topology is similar to Intel. However starting from Zen microarchitecture, one key AMD CPUs are designed as group of 4 cores which is called as CCX (Core complex) by AMD.

The main difference from Intel CPUs is that there is one LLC per each CCX/quad core.

Note that "Uncore" term is used by only Intel

**ABOUT** 

**REALMS** 



## COHERENCY

**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 AMD CPUs use MOESI , however both heavily depend on MESI protocol. There are 4 states 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 here and in other cache lines and matches its value in system memory Invalid – cache line is unused

for Core1 and var2 for Core2

Allowed state transitions in MESI protocol: N N Y E N N N Y S N N 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 whereever it is applicable. Intel MESIF: https://en.wikipedia.org/wiki/MESIF\_protocol AMD MOESI: https://en.wikipedia.org/wiki/MOESI\_protocol

Core 1 will be touching that shared cache line That change will then need to be propagated to all other cores by the cache coherency protocol, even though other cores are not using that variable. That situation is called false sharing. var1 var2 If those happen in higher rates and if Shared cache line holding var1

**FALSE SHARING AND CACHE PING-PONGING** 

ping-pong. **TLB SHOOTDOWNS** Whenever a page table entry is modified by any of the cores, that particular TLB entry is invalidated in all cores via IPIs. This one is not done by hardware but initiated by operating IPI : Interprocessor interrupt

If Core1 updates "var1" variable, that

cache lines transerred between cores

rapidly, that situation is called as cache

1. One of the cores modifies a table entry PAGE TABLE ON SYSTEM MEMORY

## MEMORY REORDERINGS & SYNCRONISATION

Apart from data TLB, there is also iTLB for

and AMD architectures

caching addresses of instructions on both Intel

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 Volume3 8.2.3.4 , 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 not happen if they use the same address CORE2

; x and y initially 0 ; x and y initially 0 mov [x], 1; STORE TO X mov [y], 1; STORE TO Y mov [result1], y ; LOAD FROM Y mov [result2], x ; LOAD FROM X In case of reordering, result1 and result2 above can both end up as zero in both cores. Note that, apart from CPUs, also compilers can do memory reordering : https://preshing.com/20120625/memory-ordering-at-compile-time/

INSTRUCTIONS TO AVOID REORDERINGS

Reorderings can be avoided by using serialising instructions such as SFENCE, LFENCE, and MFENCE: Intel Software Developer's Manual Volume3 8.3 defines them as: These instructions force the processor to complete all modifications to flags, registers, and memory by previous instructions and to drain all buffered writes to memory before the next instruction is fetched and executed. There is also bus locking "LOCK" prefix ( Intel Software Developer's Manual Volume3 8.1.2 ) which can be used as well to avoid reorderings. High level languages expose memory fence APIs which are emitted as one of those methods:

int main()
{
 std::atomic\_thread\_fence(std::memory\_order\_acquire The image is taken from the online tool Compiler Explorer. **ATOMIC OPERATIONS** An atomic operation means that there will be no other operations going on during the execution.

From point of execution, an atomic operation is indivisible and nothing can affect its execution. The most common type of atomic operations are RMW (read-modify-write) operations. ATOMIC RMW OPERATIONS: COMPARE-AND-SWAP CAS instruction ( CMPXCHG ) reads values of 2 operands. It then compares them and if they are equal, it swaps values. All the operations are atomic / uninterruptible. It can be used to implement lock free data structures.

**ATOMIC RMW OPERATIONS: TEST-AND-SET** 

Test-and-set is an atomic operation which writes to a target memory and returns its old value. It is typically used to implement spin locks. TRANSACTIONAL EXTENSIONS Transactional memory areas are programmer specified critical sections. Reads and writes in those areas are done atomically. ( Intel Optimization Manual section 16 )

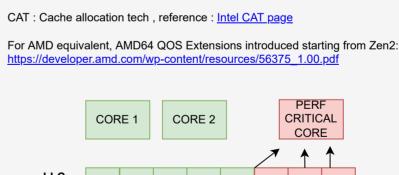
However due to another hardware security issue, Intel disabled them from Skylake to Coffee Lake CPUs: https://www.theregister.com/2021/06/29/intel\_tsx\_disabled/

### LIMITING CONTENTION BETWEEN CORES **ALLOCATING A PARTITION OF LLC**

You can allocate a partition of the shared CPU last level cache for your

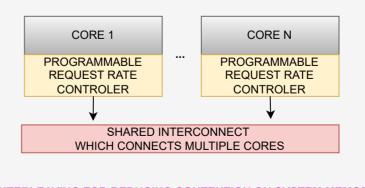
performance sensitive application to avoid evictions on Intel CPUs that

support CAT feature.

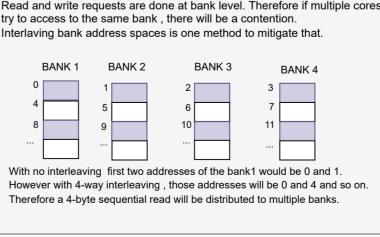


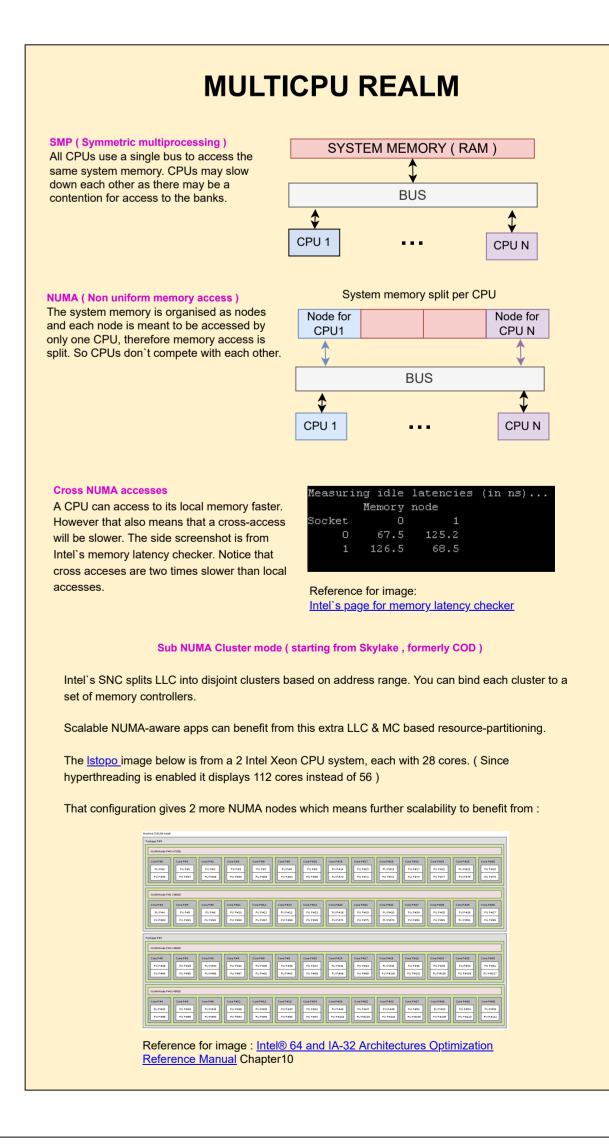
LLC cache lines LLC cache lines shared by dedicated to non performance critical only one core **MEMORY BANDWIDTH THROTTLING** You can throttle memory bandwidth per CPU core on Intel CPUs that support

MBA. Each core can be throttled with their request rate controller units. MBA: Memory bandwidth allocation, reference: Intel MBA page For AMD equivalent, AMD64 QOS Extensions introduced starting from Zen2: https://developer.amd.com/wp-content/resources/56375 1.00.pdf



INTERLEAVING FOR REDUCING CONTENTION ON SYSTEM MEMORY System memory / RAM is organised as collection of ranks. Each rank have banks which are collection of DRAM cells per bit. Read and write requests are done at bank level. Therefore if multiple cores try to access to the same bank, there will be a contention.

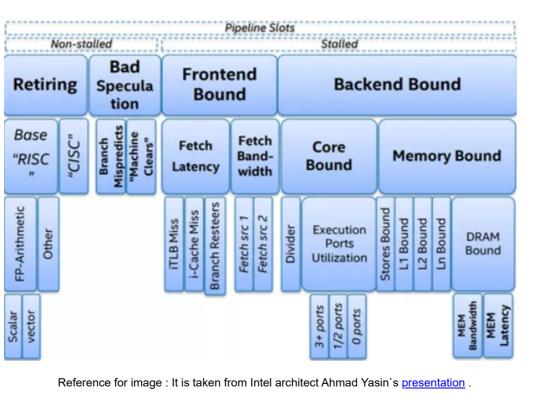


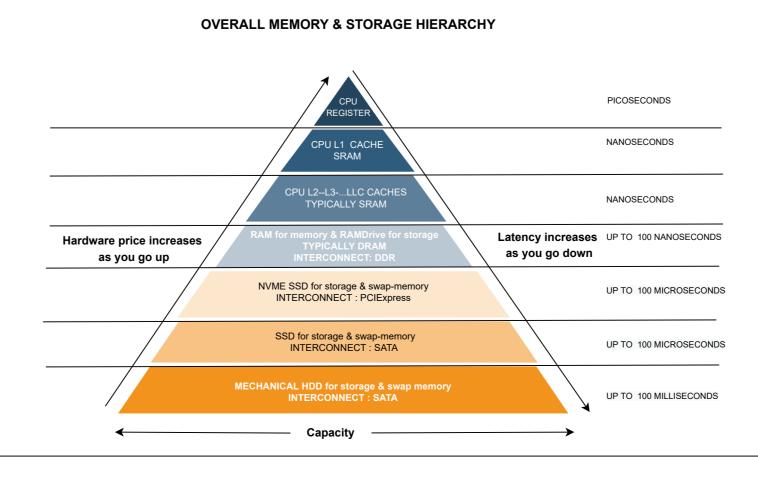


## INTEL'S TOPDOWN MICROARCHITECTURE ANALYSIS

Intel's Top Down analysis is hierarchical organization of microarchitecture events. There are 2 main categories for stalls: 1. Frontend : A typical example is large code sizes leading to instruction cache misses. 2. Backend: Usually either memory bound or compute bound

Branch mispredictions are categorised under "bad speculation". You can use either Intel's Vtune or Toplev tool to make a top down analysis. Both use Intel CPUs` PMCs ( performance monitoring counter).





## (ESTIMATED) LATENCY NUMBERS IN CLOCK CYCLES

Simple register operation ( ADD OR etc ) less than 1 clock cycle Branch prediction 1-2 clock cycles Floating point addition 1-3 clock cycles Multiplication (int, floating point) 1-7 clock cycles L1 data cache read 3-4 clock cycles TLB miss 7-21 clock cycles 10-12 clock cycles L2 data cache read 10-20 clock cycles Branch misprediction

Floating point division 10-40 clock cycles Compare and swap atomic op. 15-30 clock cycles Integer division 15-40 clock cycles L3 data cache read 30-70 clock cycles System memory read 100-150 clock cycles Cross-NUMA L3 read 100-300 clock cycles Cross-NUMA system memory read 300-500 clock cycles Reference for numbers : http://ithare.com/infographics-operation-costs-in-cpu-clock-cycles/