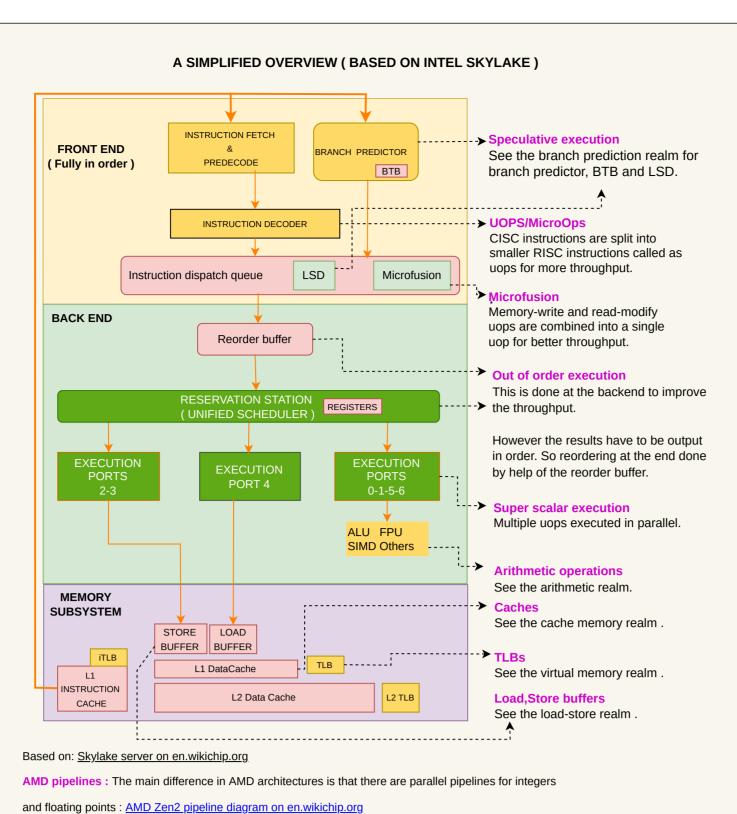
<u>Die photo</u> of a quadcore CPU , Copyright of Intel

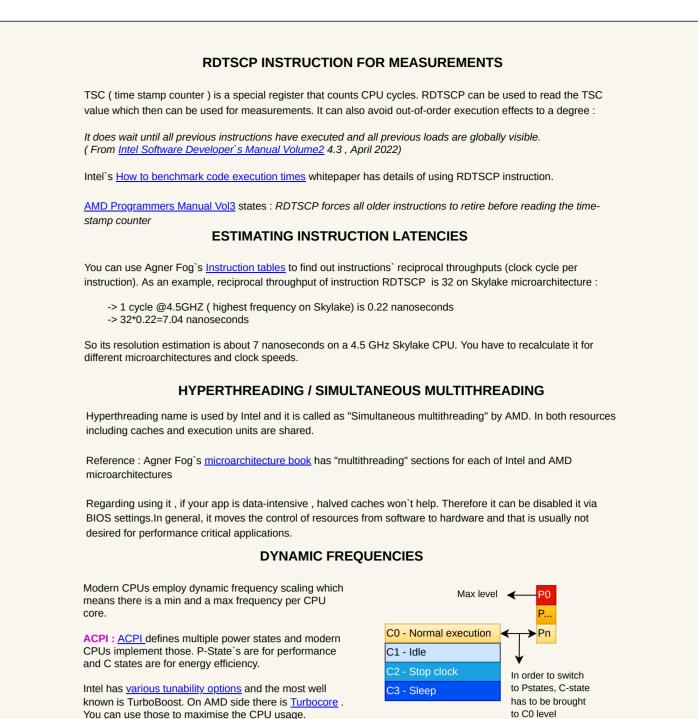
LAST UPDATE DATE: 11 JAN 2024 FOR LATEST VERSION: www.github.com/akhin/microarchitecture-cheatsheet AUTHOR: AKIN OCAL akin_ocal@hotmail.com



INSIDE INDIVIDUAL CORE



PIPELINE PARALLELISM & PERFORMANCE Pipeline diagrams: The diagrams below in the following topics are outputs from an online microarchitecture analysis P Predecoded tool <u>UICA</u> and they represent parallel execution through cycles. Q Added to IDQ Rows are multiple instructions being executed at the same time. I Issued Ready for dispatch 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 : https://perf.wiki.kernel.org/index.php/Tutorial 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 Possible Ports | Actual Port | 0 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 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. And the same applies to the 2nd purple pair. Reference: Denis Bakhvalov's article



Number of active cores & SIMD AVX2/512 on Intel CPUs: Intel's power management policies are complex.

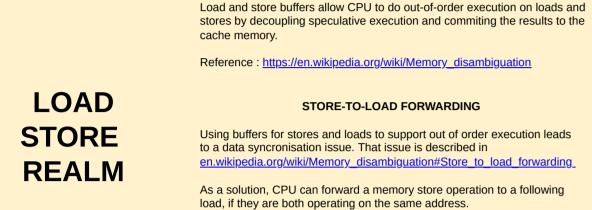
See the arithmetic and the multicore realms as number of active cores and some of AVX2/512 extensions also

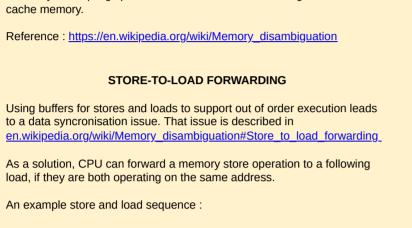
Varying max clock speeds on AMD CPUs: Some AMD CPUs` cores have slightly varying max frequencies.

__m128i , 4 x 32 bit ints

__m128l , 2 x 64 bit long longs

Therefore AMD CPUs have "preferred core" concept. Reference : <u>AMD's GDC22 presentation page19</u>





LOAD & STORE BUFFERS

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, 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) will be bypassed. L1 CACHE The conditions for a successful forwarding and latency penalties in case of no-forwarding can be found in Agner Fog`s microarchitecture book.

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 <u>restrict</u> keyword and other

methods: Elan Ruskin's article



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 <u>Instruction tables</u> & Skylake architecture on 64 bit registers Bitwise operations , integer add/sub : 0.25 to 1 clock cycle
Floating point add : 3 clock cycles

FLOATING POINTS X86 uses <u>IEEE 754</u> standard for floating points. A 32 bit floating point consists of 3 parts x86 extensions are specialised instructions. They have various categories 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 mantissa - 23 bits 8 bits

IEEE754 also defines denormal numbers. They are very small / near zero numbers. As floating points are approximations, float GetInverseOfDiff(float a, float b) denormal numbers are needed to avoid an undesired case of : a!=b but a-b=0 return 1.0f / (a - b); Without denormals the code to the right return 0.0f; would invoke a divide-by-zero exception. Reference : <u>Bruce Dawson's article</u>

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

Based on Agner Fog`s microarchitecture book, 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 The most recent SIMD instruction sets for Intel CPUs are : from <u>cryptography</u> to <u>neural network operations</u>. : Up to 256 bits 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: i1 i2 i3 i4

may affect the frequency while in Turboboost.

In the example above, an array 4 integers (i1 to i4) are added to another array of integers (j1 to j4). The result is also an array of sums (s1 to s4). In this example, 4 add operations are executed by a single instruction. They play key role in compilers' vectorisation optimisations: GCC auto vectorisation

= = = =

AVX2 : Up to 256 bits <u>AVX512</u>: Up to 512 bits Recent AMD CPUs support AVX & AVX2. Only the latest Zen4 architecture supports As for programming, there are also wider data types. The data type diagrams below are for 128 bit operations: m128 , 4 x 32 bit floating points Float Float Float __m128d , 2 x 64 bit doubles

int int

int int

long long

X86 EXTENSIONS: SIMD DETAILS

Note that as SIMD instructions require more power, therefore usage of some AVX2/512 extensions may introduce downclocking. They should be benchmarked. For details : <u>Daniel Lemire`s article</u>

OFFSET

BRANCH

REALM

Why virtual memory?

BRANCH

_1-L2..(excludin

LLC) data &

struction cacl

UNCORE

cache.

difference is CCXs.

there is one LLC per each CCX/quad core.

CORE 1

CORE 3

EXECUTION

UNITS

TLB data and

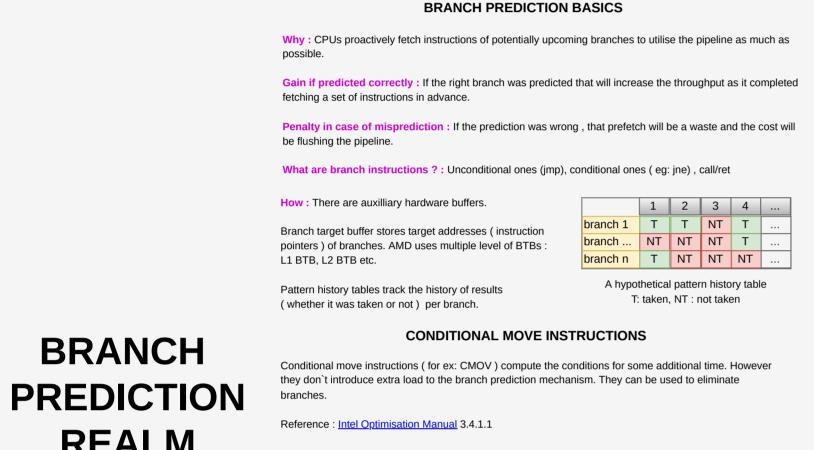
instruction

caches for VM

only term to refer to CPU functionality which are not per core.

LAST LEVEL CACHE (L3 OR L4)

Because cumulative memory requirement of



BP METHODS: 2-LEVEL ADAPTIVE BRANCH PREDICTION

Saturating counter as a building block Strongly not taken A 2-bit saturating counter can store 4 strength states. Veakly not taker Whenever a branch is taken it goes stronger. And whenever a branch is not taken it goes 2 level adaptive predictor In this method, the pattern history table keeps 2ⁿ rows and each row will have a saturating counter.

A branch history register which has the history of last n occurences, will be used to choose which row will be used from the pattern history table Reference : Agner Fog`s microarchitecture book 3.1.

BP METHODS : AMD PERCEPTRONS

A <u>perceptron</u> is basically the simplest form of machine learning. It can be considered as a linear array of 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. For details of perceptron based branch prediction: The output Y (in this case whether a branch <u>Dynamic Branch Prediction with Perceptrons by Daniel</u> taken or not) is calculated by dot product Jimenez and Calvin Lin of the weight vector and the input vector.

INTEL LSD (LOOP STREAM DETECTOR)

Intel LSD will detect a loop and stop fetching instructions to improve the frontend bandwidth. Several conditions mentioned in $\underline{\text{Intel Optimisation Manual}}\ 3.4.2.4$: • Loop body size up to 60 μops, with up to 15 taken branches, and up to 15 64-byte fetch 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://en.wikichip.org/wiki/intel/microarchitectures/skylake_(server)#Front-end DISABLING SPECULATIVE EXECUTION PATCHES

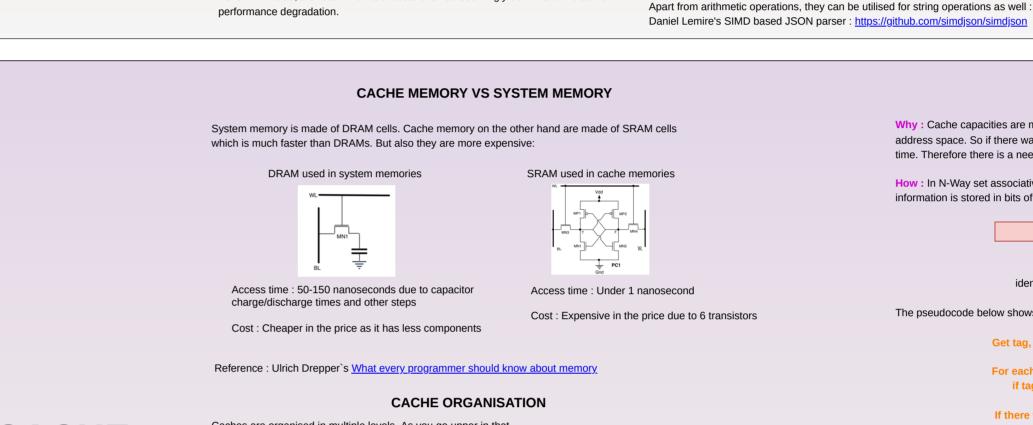
You can consider disabling system patches for speculative execution related vulnerabilities such as Meltdown and Spectre for performance, if that is doable in your system. Those patches are not only microcode updates but they also need OS support. Kernel.org documentation : <a href="https://www.kernel.org/doc/html/latest/admin-guide/kernel-or parameters.html Red Hat Enterprise documentation : https://access.redhat.com/articles/3311301

Spectre paper: https://spectreattack.com/spectre.pdf **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

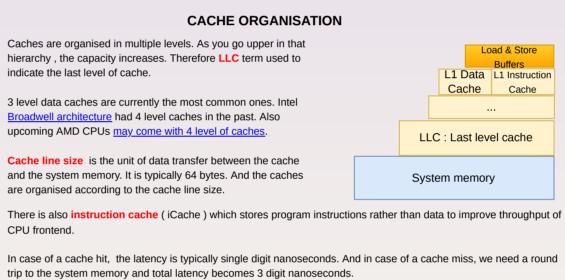
Intel Xeon Gold 6262 -> roughly 4K AMD EPYC 7713 -> roughly 3K Reference : Marek Majkovski`s article on Cloudflare blog

Meltdown paper: https://meltdownattack.com/meltdown.pdf

sequences of branch instructions :



CACHE **MEMORY** REALM



Hardware prefetchers detect patterns like streams (Ex: accessing to contiguous array members) and strides (Ex: accessing specific members in arrays of structs) and prefetch data and instruction to cache lines automatically. Developers can also use instruction _mm_prefetch to prefetch data explicitly for cases when hardware

Stride +2
Stride -3

Str can't predict. That is called as software prefetching. Reference for image: It is taken from AMD's GDC22 presentation page44

HARDWARE AND SOFTWARE PREFETCHING

N-WAY SET ASSOCIATIVITY

Why: Cache capacities are much smaller than the system memory. Moreover, software can use various regions of their 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 memory. How: In N-Way set associativity, caches are divided to groups of sets. And each set will have N cache blocks. The mapping information is stored in bits of addresses which has 3 parts:

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

SET

For each block in the current set (which we have just found out) if tag of the current block equals to tag (which we just have found out) read and return data using offset, it is a 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

Get tag, set and offset from the address

BYPASSING THE CACHE: NON-TEMPORAL STORES & WRITE-COMBINING Temporal data is data that will be accessed in a short period of time. The term non-temporal data indicates that data will not be accessed any soon. (cold data). If the amount of non-temporal data is excessive in the cache, that is called as cache

pollution. Non-temporal store instructions are introduced for this problem and they store data directly to the system memory by bypassing the cache. Write combining buffers are used with non-temporal stores. CPU will try to fill a whole cache line (typically

64byte) before committing to the system memory and only will send to the system memory when that buffer is filled. That is for reducing the load on the bus between the CPU and the system memory.

SYSTEM MEMORY REALM

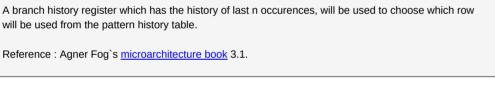
DIRECT CACHE ACCESS

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 DDIO (Direct I/O).

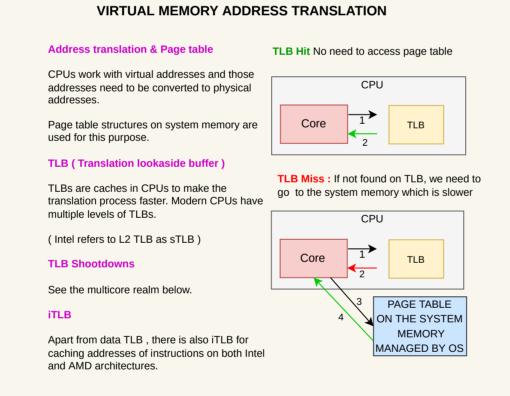
Reference : Intel documentation

Last CPU Level DCA NIC



SYSTEM MEMORY ADDRESS SPACE

VIRTUAL MEMORY REALM



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 Regular pages 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. 1 GB Note that there may be extra abstractions on top of page size For ex : Windows has 64KB page allocation granularity : Raymond Chen's article Huge pages

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. 47 39 38 30 29 21 20 4 level page table is the most common one. In the diagram above, the 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) 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.

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 DANIZ 1 System memory / RAM is organised as RANK ... BANK 1 collection of ranks. RANK N BANK 1 BANK .. Each rank have banks which are collection of BANK N — 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 throughput.

multiple processes in an OS can be more than Ok the system capacity It is basically for sharing memory resources between multiple processes. It also provides security by isolating processes. - Minimum addressable virtual space that can be requsted from OS. SWAP FILE - Typically 4KB ON DISC 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

MULTICORE REALM

TOPOLOGIES

TOPOLOGICAL OVERVIEW - INTEL CPUS

CPU

. . .

MEMORY CONTROLLER (Server CPUs likely to have more than one MC)

TOPOLOGICAL OVERVIEW - AMD CPUS

Practically the maximum number of cores competing for the LLC (without simultanenous

CPU

CORE 2 | CORE 5

CORE 4 | CORE 7

Reference: https://en.wikichip.org/wiki/amd/microarchitectures/zen#CPU Complex .28CCX.29

multithreading) is 4 in recent AMD CPUs. An example 8 core CPU with 2 CCXs:

EXECUTION

UNITS

TLB data and

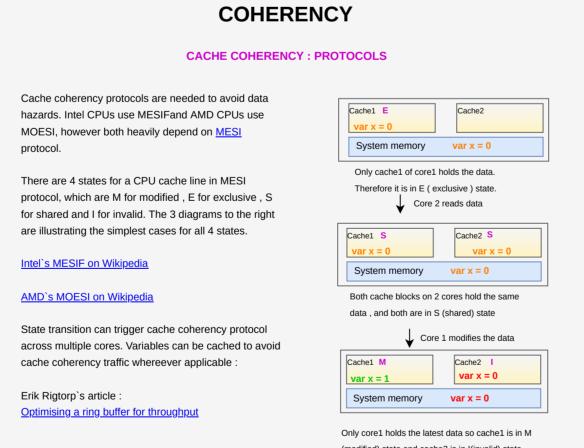
instruction

caches for VM

Others

VIRTUAL MEMORY ORGANISATION

VIRTUAL ADDRESS SPACE



(modified) state and cache2 is in I(invalid) state. CACHE COHERENCY: FALSE SHARING AND CACHE PING-PONGING In the diagram to the right, if Core1 Diagram above aims to show resource per core and shared resources. Note that uncore in an Intelchanges its var1, that change will Core 1 Core 2 need to be propagated to all other cores by the cache coherency protocol. That will lead to Hybrid topologies: An exception to the above diagram is Intel's recent E-cores. E-cores are meant invalidation of cache areas for power efficiency and paired with less resources. For ex: Alder Lake CPUs` E-cores also share L2 associated with the shared cache line across all cores, even though it var1 var2 memory Reference: https://www.anandtech.com/show/16959/intel-innovation-alder-lake-november-4th is used by only one core. Shared system memory cache line holding That situation is called false sharing. var1 for Core1 and var2 for Core2 If those happen in higher rates and if cache lines from system memory transferred between cores rapidly, Most of AMD topology is similar to Intel. However starting from Zen microarchitecture, one key that situation is called as cache ping-pong. AMD CPUs are designed as group of 4 cores which is called as CCX (Core complex) , and **VIRTUAL MEMORY PAGE TABLE COHERENCY: 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 system. IPI: Interprocessor interrupt, you can take "processor" as core in this context. 1. One of the cores modifies a table entry PAGE TABLE ON SYSTEM TLB MEMORY

MEMORY REORDERINGS & SYNCRONISATION

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 CORE1 CORE2

; x and y initially 0 ; x and y initially 0 mov [y], 1; STORE TO Y mov [x], 1; STORE TO X 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: Jeff Preshing`s article : Memory Ordering at Compile Time

INSTRUCTIONS TO AVOID REORDERINGS Reorderings can be avoided by using serialising instructions such as SFENCE, LFENCE, and MFENCE : <u>Intel Software Developer's Manual Volume3</u> 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

There is also bus locking "LOCK" prefix (<a href="Intelligent Intelligent Int which can be used as well to avoid reorderings.

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 OPERATIONS & SPLIT LOCKS If an atomic instruction is used for a memory range which is split to multiple cache lines, that will lead to

locking the whole memory bus, instead of just the cache line. Reference : <u>Detecting and handling split locks (in Linux kernel) on lwn.net</u> 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.

PAUSE INSTRUCTIONS Busy spinning applications (ex: user space spin locks) can degrade hyperthreading efficiency. There are several pause instructions (PAUSE/ TPAUSE/UMWAIT/UMMONITOR) to help that. Note that the latency for PAUSE instruction on Skylake clients is an order of magnitude slower than other architectures: Intel Optimisation Manual 2.5.4 TRANSACTIONAL MEMORY

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/ AMD equivalent is called as "Advanced Syncronisation Facility". According to Wikipedia article, there are no AMD processors using it yet.

LIMITING CONTENTION BETWEEN CORES

DISABLING UNUSED CORES TO MAXIMISE FREQUENCY (INTEL)

Therefore disabling unused cores may improve frequency for perf-critical cores,

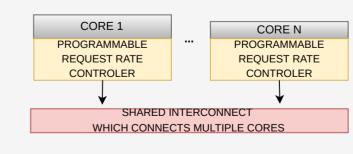
Number of active cores may introduce downclocking : Wikichip article

depending on your CPU. You shall refer to your CPU's frequency table : An example frequency table : Wikichip XeonGold5120 article ALLOCATING A PARTITION OF LLC (SERVER CLASS CPUS) 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 **CDP** (Code and data prioritisation) allows developers to allocate LLC on code basis : Intel's CDP page on supported CPUs.

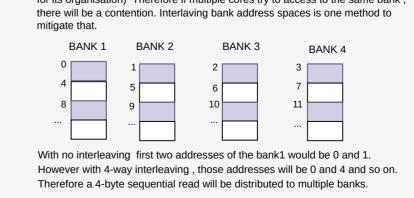
CORE 1 CORE 2 CRITICAL CORE LLC cache lines shared by non LLC cache lines dedicated performance critical cores to only one core

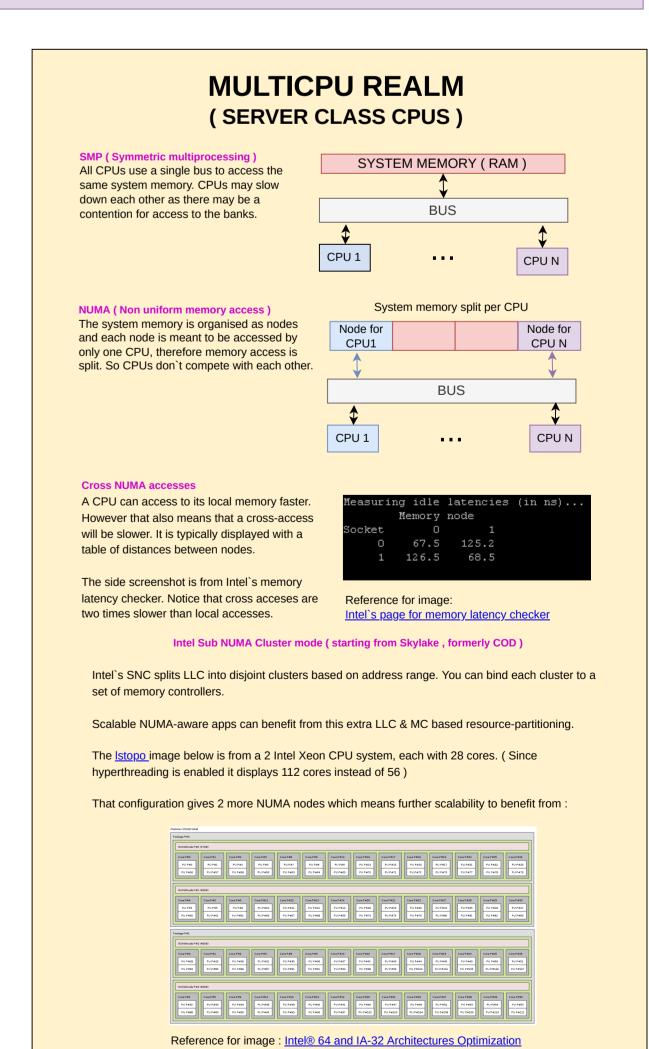
On AMD side, QOS Extensions were introduced starting from Zen2. Corresponding technologies are called as "Cache allocation enforcement" and "Code and data prioritisation": https://kib.kiev.ua/x86docs/AMD/MISC/56375 1.00 PUB.pdf MEMORY BANDWIDTH THROTTLING (SERVER CLASS CPUS)

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, QOS Extensions were introduced starting from Zen2: https://www.amd.com/system/files/TechDocs/56375 1.03 PUB.pdf



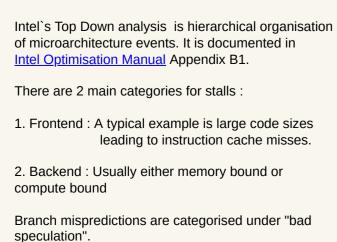
INTERLEAVING FOR REDUCING CONTENTION ON SYSTEM MEMORY Read and write requests are done at bank level. (See the system memory realm for its organisation) Therefore if multiple cores try to access to the same bank, there will be a contention. Interlaving bank address spaces is one method to





ACROSS

REALMS



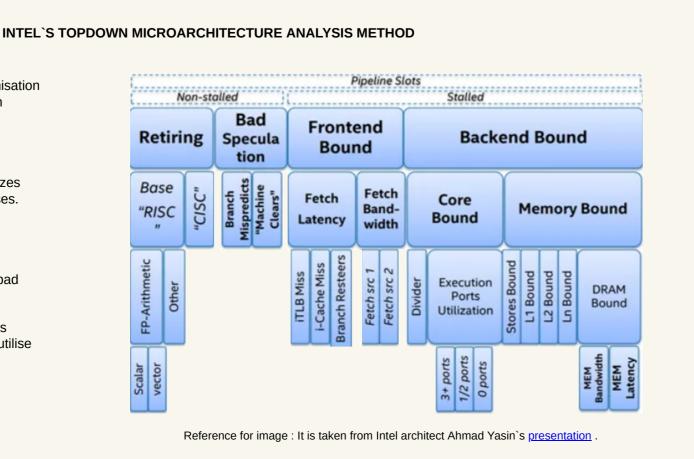
You can use either Intel's Vtune or Andi Kleen's

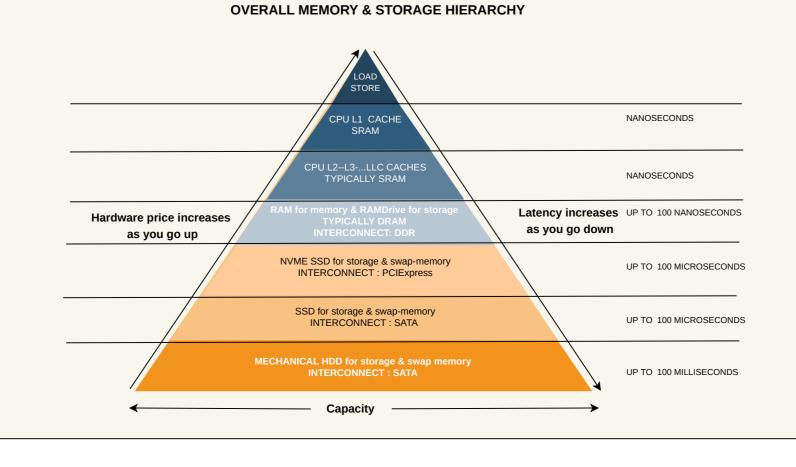
Intel CPUs` performance monitoring counters.

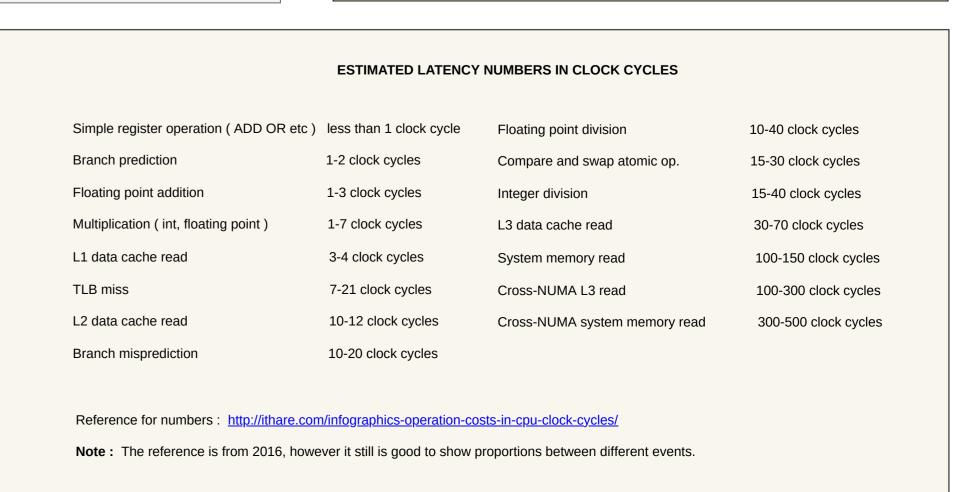
<u>Toplev</u> tool to make a top down analysis. Both utilise

CORE 6

CORE 8







Reference Manual Chapter10

AMD NUMA nodes per socket: It is the equivalent of SNC for AMD CPUs.

Reference : <u>Dell article about AMD NUMA nodes per socket feature</u>