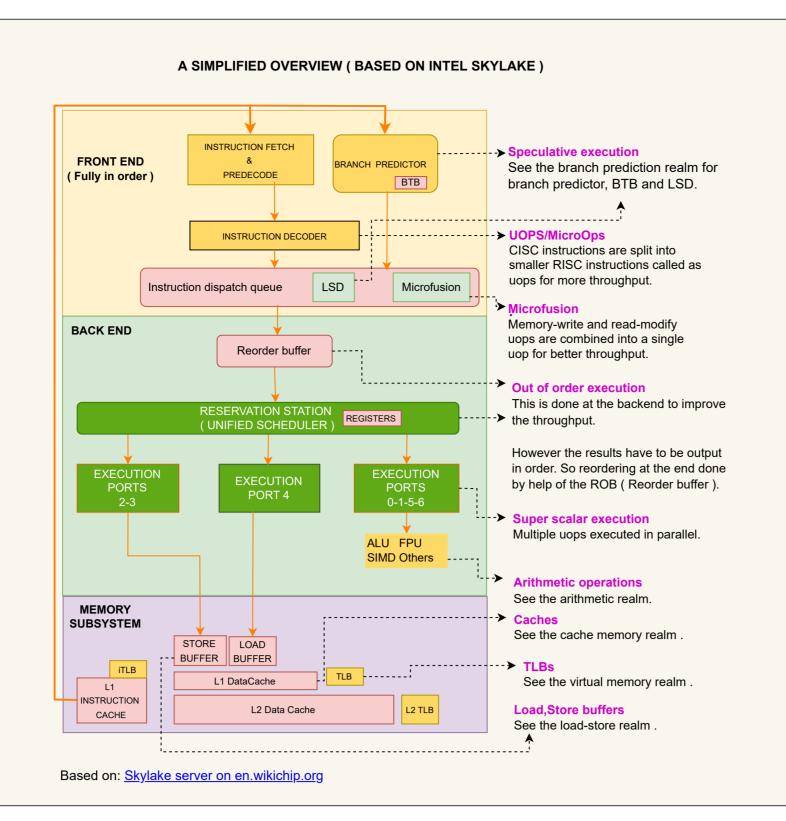
X86 CPUs & Performance

Die photo of a quadcore CPU , Copyright of Intel

LAST UPDATE DATE: 30 OCT 2022 FOR LATEST VERSION: www.github.com/akhin/microarchitecture-cheatsheet AUTHOR: AKIN OCAL akin ocal@hotmail.com

PIPELINE REALM:

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 UICA and they represent parallel execution through cycles. Q Added to IDQ I Issued Rows are multiple instructions being executed at the same time. 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 Rate of retired instructions: Apart from IPC, number of retired instructions should be checked. Retired instructions in UICA diagrams 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 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 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

<u>∞</u>. exponent

8 bits

Reference : Bruce Dawson'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 Instruction tables, RDTSCP reciprocal throughput (clock cycle per instruction) is 32 on -> 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. DYNAMIC FREQUENCIES Modern CPUs employ dynamic frequency scaling which means there is a min and max Max level ← frequency per CPU core. ACPI : ACPI defines multiple power states and C0 - Normal execution ← → Pr modern CPUs implement those. P-State's are C1 - Idle for performance and C states are for energy In order to switch o Pstates, C-state You can use Intel's <u>Turboboost</u> or AMD's nas to be brought <u>Turbocore</u> to maximise the CPU usage. to C0 level Intel CPUs number of active cores & SIMD usage: Number of active cores can introduce downclocking

regardless of SIMD usage . Therefore disabling unused cores can further improve frequency of your target

Example Xeon Gold 5120 frequency table : https://en.wikichip.org/wiki/intel/xeon_gold/5120#Frequencies

table. Regarding SIMD and downclockings: Daniel Lemire's article

Use of some SIMD instructions may introduce further downclockings. You shall refer to your CPU's frequency

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 commiting 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 syncronisation 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, 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

What would happen without forwarding?: Previous game consoles PlayStation3 and Xbox360 had PowerPC based processors which did inorder-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

no-forwarding can be found in Agner Fog's microarchitecture book.

The conditions for a successful forwarding and latency penalties in case of

ARITHMETIC REALM

Reference : Denis Bakhvalov`s article

and notice that the second instruction gets executed after the first one.

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 nteger division: 24-90 clock cycles

CACHE

MEMORY

REALM

FLOATING POINTS X86 uses IEEE 754 standard for floating points. A 32 bit floating point consists of 3 parts 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

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 undesired case of : a!=b but a-b=0 Without denormals the code to the right return 1.0f / (a - b); return 0.0f; would invoke a divide-by-zero exception.

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 x86 extensions are specialised instructions. They have various categories from <u>cryptography</u> to <u>neural network operations</u>.

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 + + + +

= = = =

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 optimisation : GCC auto vectorisation Apart from arithmetic operations, they can be utilised for string operations as well. A SIMD based JSON parser : https://github.com/simdjson/simdjson

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 int int _m128l , 2 x 64 bit long longs long long long long

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

BRANCH PREDICTION

REALM

BRANCH PREDICTION BASICS Why: CPUs proactively fetch instructions of potentially upcoming branches to utilise the pipeline as much as 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/ret How: There are auxilliary hardware buffers. 1 2 3 4 ... T T NT T Branch target buffer stores target addresses (instruction branch ... NT NT NT T pointers) of branches. AMD uses multiple level of BTBs : branch n T NT NT NT . L1 BTB, L2 BTB etc. A hypothetical pattern history table Pattern history tables track the history of results T: taken, NT : not taken (whether it was taken or not) per branch. **CONDITIONAL MOVE INSTRUCTION**

CMOV (Conditional move) instruction takes additional time in order to compute also the condition. Therefore they don't introduce extra load to branch prediction. They can be used to eliminate branches. Reference : Intel Optimisation Manual 3.4.1.1

BP METHODS: 2-LEVEL ADAPTIVE BRANCH PREDICTION Saturating counter Not taken A 2-bit saturating counter can store 4 strength states. Veaklv not taker Not taken Whenever a branch is taken it goes

Not taken

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 2ⁿ 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.

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 Agner For mentions that they are good at predicting very long branches compared to 2-level adaptive branch prediction in his microarchitecture book. The output Y (in this case whether a branch taken For details of perceptron based branch prediction:

BP METHODS: AMD PERCEPTRONS

INTEL LSD (LOOP STREAM DETECTOR)

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

Jimenez and Calvin Lin

Intel LSD will detect a loop and stop fetching instruction to improve frontend bandwidth. Several conditions mentioned in Intel Optimization Manual: \bullet Loop body size up to 60 μ ops, 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.

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 it is doable in your system. Kernel.org documentation: https://www.kernel.org/doc/html/latest/admin-guide/kernel- <u>parameters.html</u> Red Hat Enterprise documentation : https://access.redhat.com/articles/3311301 Meltdown paper: https://meltdownattack.com/meltdown.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 sequences of branch instructions:

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

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

CACHE MEMORY VS SYSTEM MEMORY System memory is made of DRAM cells. Cache memory on the other hand are made of SRAM cells which is much faster than DRAMs. On the other hand they are more expensive:

DRAM used in system memories

Access time: 50-150 nanoseonds due to capacitor Access time: Under 1 nanosecond charge/discharge times and other steps Cost: Expensive in the price due to 6 transistors

CACHE ORGANISATION

SRAM used in cache memories

Load & Store Buffers

L1 Data Cache

LLC: Last level cache

Cost: Cheaper in the price as it has less components Reference: Ulrich Drepper's What every programmer should know about memory

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. Also

upcoming AMD CPUs may come with 4 level of caches.

Lemire's article

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 e (iCache) which store program instructions rather than data to improve throughput of CPU frontend.

a round trip to the system memory and total latency becomes 3 digit nanoseconds.

HARDWARE AND SOFTWARE PREFETCHING

Intel Optimisation Manual 3.7 describes prefetching. Hardware prefetchers prefetch data and instruction to cache

lines automatically. Developers can also use instruction mm_prefetch to prefetch data explicitly. That is called as

software prefetching. However performance improvement by using software prefetcher is controversial: Daniel

In case of a cache hit, the latency is typically single digit nanoseconds. And in case of a cache miss, we need

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 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 :

OFFSET used as unique identifier used to determine used to determine the actual bytes per cacheline 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

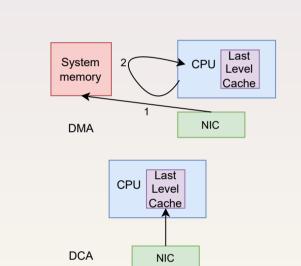
Get tag, set and offset from the address

The level of associtiavity (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 DMA mechanism doesn't require CPU involvement.

Though mechanism initiated by CPU, therefore CPU 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 REALM

or not) is calculated by dot product of the weight vector and the input vector.

VIRTUAL MEMORY ORGANISATION VIRTUAL ADDRESS SPACE SYSTEM MEMORY ADDRESS SPACE 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. - 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.

And whenever a branch is not taken it

Reference : Agner Fog`s microarchitecture book

goes weaker.

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 1 GB 1 GB supported sizes. Huge pages

TLB PRESSURE & HUGE PAGES

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 addresses. Core TLB 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 Shootdowns** See the multicore realm below. PAGE TABLE iTLB ON THE SYSTEM MEMORY Apart from data TLB, there is also iTLB for MANAGED BY OS caching addresses of instructions on both Intel and AMD architectures.

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: 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

PAGE TABLE WALKING

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.

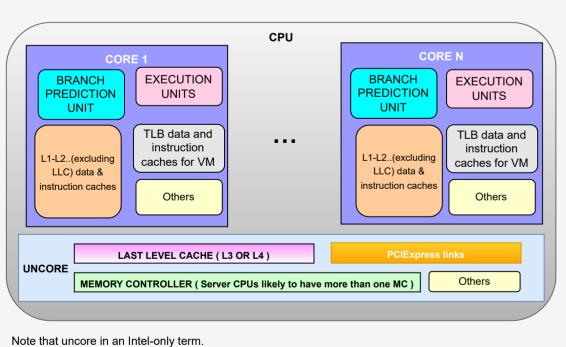
SYSTEM MEMORY REALM **DDR RAMs** DDR RAMs are the most common commodity hardware as system memory. They are found in forms of DIMMs (Dual inline memory module) / RAMSticks. A DIMM. Click for Image source **Organisation** RANK 1 BANK 1 System memory / RAM is organised as RANK ... RANK 1 collection of 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 fine-

MULTICORE REALM

Happens when the page is not on physical

memory but on the swap file which is on the

TOPOLOGICAL OVERVIEW - INTEL CPUS



Exception of E-cores: An exception to the above diagram is Intel's recent E-cores. E-cores are meant for power efficiency and paired with less resources. For ex: Alder Lake CPUs` E-cores also share L2 cache. Reference: https://www.anandtech.com/show/16959/intel-innovation-alder-lake-november-4th

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. And there is one LLC per each CCX/quad core. So practically the max number of cores competing for the LLC (without simultanenous multithreading)

An example 8 core CPU with 2 CCXs CORE 2 | CORE 5 CORE 6 CORE 4 | CORE 7 CORE 8 AMD CCX structure :There will be multiple LLCs for each group of 4 cores

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

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 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 here and in other cache lines and matches its value in system memory Invalid – cache line is unused Allowed state transitions in MESI protocol:

> N N N Y E N N N Y S N N Y Y I Y Y Y Y

> > var1 var2

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

CACHE COHERENCY: FALSE SHARING AND CACHE PING-PONGING If Core1 updates "var1" variable, that Core 1 will be touching that shared cache line. That change will then need to be

If those happen in higher rates and it Shared cache line holding var1 cache lines transerred between cores for Core1 and var2 for Core2 rapidly , that situation is called as cach ping-pong. **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

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.

INTEL'S TOPDOWN MICROARCHITECTURE ANALYSIS METHOD

IPI : Interprocessor interrupt 1. One of the cores modifies a table entry Core1 Core2 PAGE TABLE ON SYSTEM TLB TLB 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: CORE2

; x and y initially 0 ; x and y initially 0 mov [y], 1; STORE TO Y mov [x], 1; STORE TO X mov [result2], x ; LOAD FROM X mov [result1], y; LOAD FROM Y 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: 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:

> The image is taken from the online tool Compiler Explorer ATOMIC 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 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.

An atomic operation means that there will be no other operations going on during the execution

The most common type of atomic operations are RMW (read-modify-write) operations.

From point of execution, an atomic operation is indivisible and nothing can affect its execution.

TRANSACTIONAL MEMORY Transactional memory areas are programmer specified critical sections. Reads and writes in those areas are done atomically. (<u>Intel Optimization Manual</u> 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, no released AMD processors implemented it.

LIMITING CONTENTION BETWEEN CORES 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

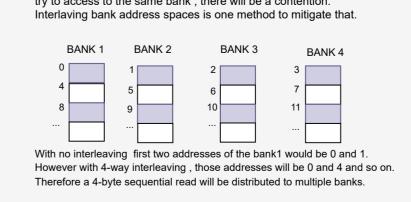
LLC cache lines shared by non LLC cache lines dedicated performance critical cores to only one core On AMD side, QOS Extensions were introduced starting from Zen2. Corresponding technologies are called as "Cache allocation enforcement" and "Code and data prioritisation": https://developer.amd.com/wp-content/resources/56375_1.00.pdf

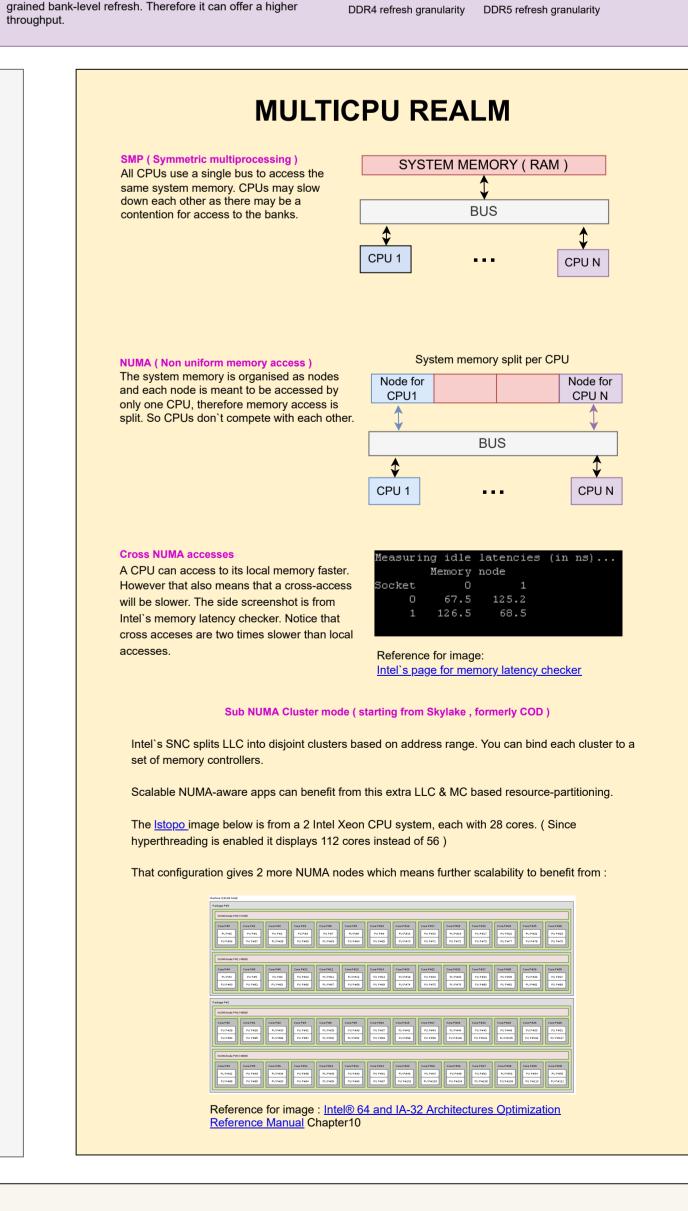
You can throttle memory bandwidth per CPU core on Intel CPUs that support MBA. Each core can be throttled with their request rate controller units. MBA: Memory bandwidth allocation, reference: Intel MBA page For AMD equivalent, QOS Extensions were introduced starting from Zen2: https://developer.amd.com/wp-content/resources/56375 1.00.pdf

MEMORY BANDWIDTH THROTTLING (SERVER CLASS CPUS)

CORE 1 CORE N PROGRAMMABL **PROGRAMMABLE** REQUEST RATE REQUEST RATE CONTROLER CONTROLER SHARED INTERCONNECT WHICH CONNECTS MULTIPLE CORES 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. Interlaving bank address spaces is one method to mitigate that. BANK 2 BANK 3

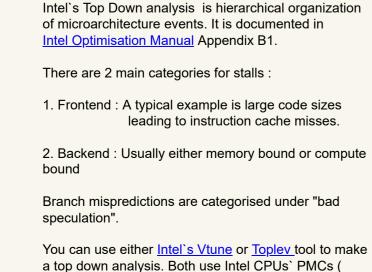




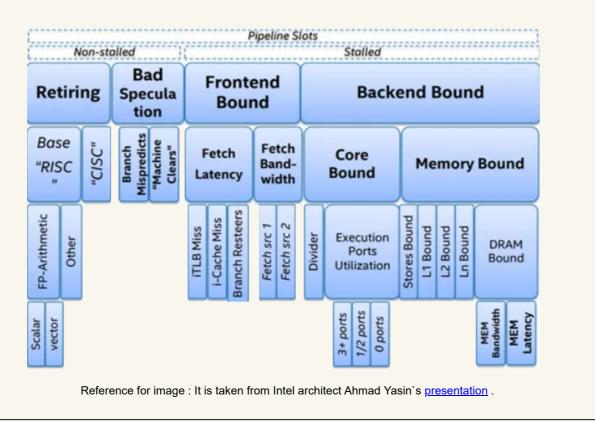
ABOUT

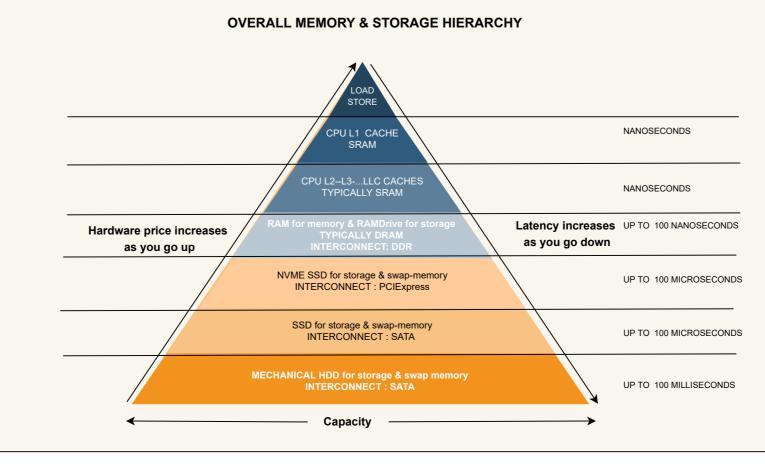
ALL

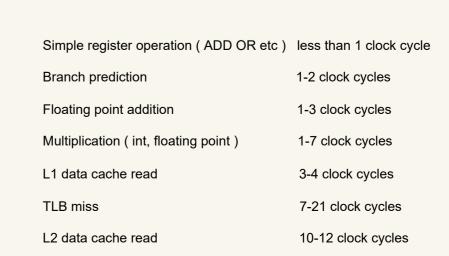
REALMS



performance monitoring counter).

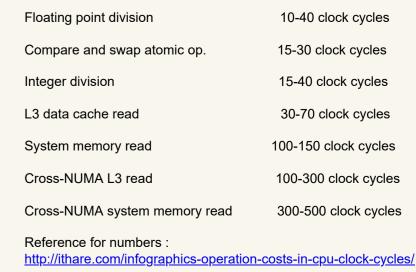






10-20 clock cycles

Branch misprediction



Note: The reference is a bit dated (2016) though it still is good to

show proportions between different event.

ESTIMATED LATENCY NUMBERS IN CLOCK CYCLES