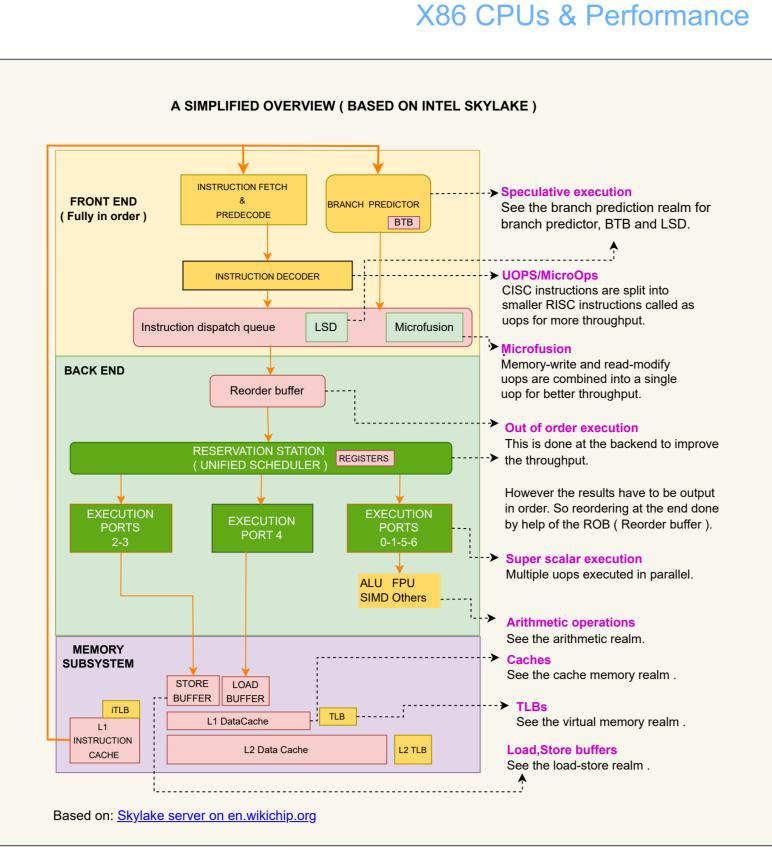
Die photo of a quadcore CPU , Copyright of Intel

LAST UPDATE DATE: 19 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 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 : 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 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

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 modern CPUs implement those. P-State's are C1 - Idle for performance and C states are for energy In order to switch to Pstates, C-state You can use Intel's <u>Turboboost</u> or AMD's

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

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 Reference: https://en.wikipedia.org/wiki/Memory_disambiguation **LOAD** STORE-TO-LOAD FORWARDING **STORE**

REALM

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

2ⁿ rows and each row has a saturating counter.

Reference : Agner Fog`s microarchitecture book

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

sequences of branch instructions:

Reference: Marek Majkovski's article on Cloudflare blog

Intel Xeon Gold 6262 -> roughly 4K

AMD EPYC 7713 -> roughly 3K

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

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

ARITHMETIC REALM

Reference: Denis Bakhvalov'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 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

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

would invoke a divide-by-zero exception. Reference : Bruce Dawson's article 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. Some typical application areas are 3D graphics and quantitative finance. As for AMD side, the recent Zen architecture CPUs seemingly don't have the same performance degradation.

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 visualizer:

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

Without denormals the code to the right

denormal numbers are needed to avoid an

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

For the list of extensions : https://en.wikipedia.org/wiki/X86#Extensions SSE (Streaming SIMD Extensions) is one of the most important ones. SIMD stands for "single instruction multiple data". SIMD instructions use wider registers to execute more work in a single go:

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

X86 EXTENSIONS

such as cryptography and neural network operations

operations are executed by a single instruction.

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

Daniel Lemire's article

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

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

has to be brought

to C0 level

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

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

OFFSET

N-WAY SET ASSOCIATIVITY

Cache capacities are much smaller than the system memory. Moreover, softwares can use various regions of their

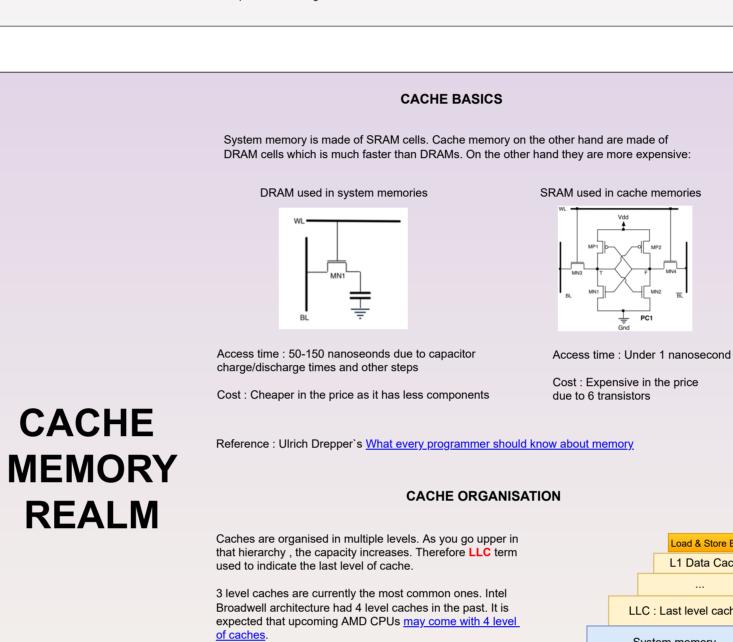
In N-Way set associativity, caches are divided to groups of sets. And each set will have N cache lines. The mapping

SET

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

BP METHODS: AMD PERCEPTRONS BRANCH PREDICTION BASICS They are used in Zen arcihtectures. Why: CPUs proactively fetch instructions of potentially upcoming branches to utilise the pipeline as much as A <u>perceptron</u> is basically the simplest form of machine learning. They can be considered as a linear array of weights. 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. Agner For mentions that they are good at predicting very long branches compared to 2-level adaptive Penalty in case of misprediction: If the prediction was wrong, that prefetch will be a waste and the cost will branch prediction in his microarchitecture book. For details of perceptron based branch prediction: The output Y (in this case whether a branch taken What are branch instructions?: Unconditional ones (jmp), conditional ones (eg: jne), call/ret or not) is calculated by dot product of the weight <u>Dynamic Branch Prediction with Perceptrons by Daniel</u> vector and the input vector. 1 2 3 4 ... How: There are auxilliary hardware buffers. INTEL LSD (LOOP STREAM DETECTOR) T T NT T Branch target buffer stores target addresses (instruction branch ... NT NT NT T Intel LSD will detect a loop and stop fetching instruction to improve frontend bandwidth. Several pointers) of branches. AMD uses multiple level of BTBs : branch n T NT NT NT conditions mentioned in Intel Optimization Manual L1 BTB, L2 BTB etc. **BRANCH** A hypothetical pattern history table • Loop body size up to 60 μops, with up to 15 taken branches, and up to 15 64-byte fetch lines. Pattern history tables track the history of results T: taken, NT: not taken (whether it was taken or not) per branch. • No mismatched stack operations (e.g., more PUSH than POP). **PREDICTION** More than ~20 iterations. **BP METHODS: 2-LEVEL ADAPTIVE BRANCH PREDICTION** Note that LSD is disabled on Skylake Server CPUs. Reference: https://en.wikichip.org/wiki/intel/microarchitectures/skylake_(server)#Front-end REALM A 2-bit saturating counter can store 4 **DISABLING SPECULATIVE EXECUTION PATCHES** strength states. You can consider disabling system patches for speculative execution related vulnerabilities such Whenever a branch is taken it goes as Meltdown and Spectre for performance, if it is doable in your system. stronger. Kernel.org documentation: https://www.kernel.org/doc/html/latest/admin-guide/kernel-And whenever a branch is not taken it Not taken <u>parameters.html</u> Meltdown paper : https://meltdownattack.com/meltdown.pdf Spectre paper : https://spectreattack.com/spectre.pdf 2 level adaptive predictor In this method, you store the history of last n occurences in a history register which is n bits. **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 Also you create a table called "pattern history table" for that branch. That pattern history table keeps



All the mentioned caches till now were data caches. But there is also in

program instructions rather than data to improve throughput of CPU frontend.

ALLOCATING A PARTITION OF LLC

CORE 2

MEMORY BANDWIDTH THROTTLING

SHARED INTERCONNECT

WHICH CONNECTS MULTIPLE CORES

BANK 3

CORE 1

CORE 1

PROGRAMMABLE

REQUEST RATE

CONTROLER

LLC cache lines shared by

non performance critical

CRITICAL

CORE

LLC cache lines

dedicated to

only one core

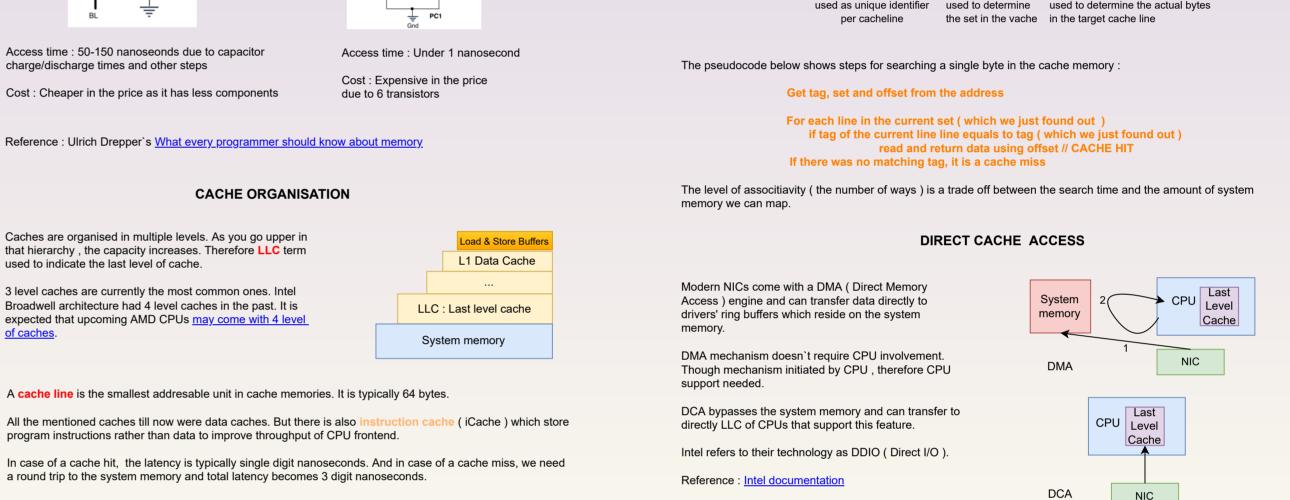
CORE N

PROGRAMMABLE

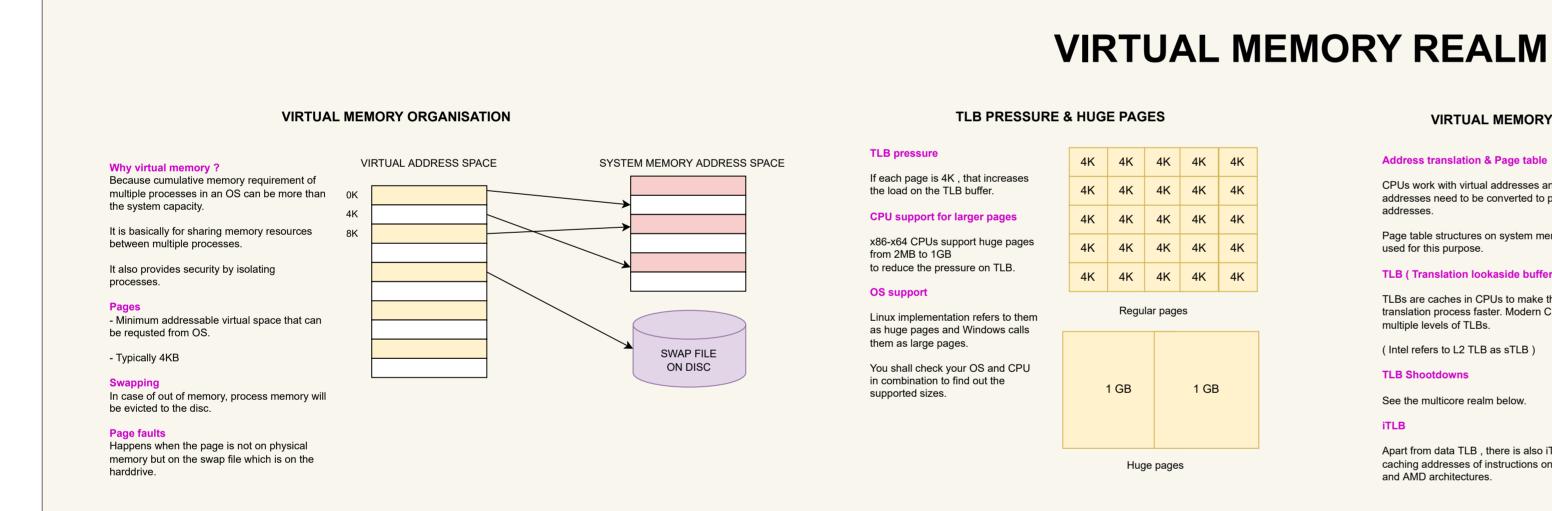
REQUEST RATE

CONTROLER

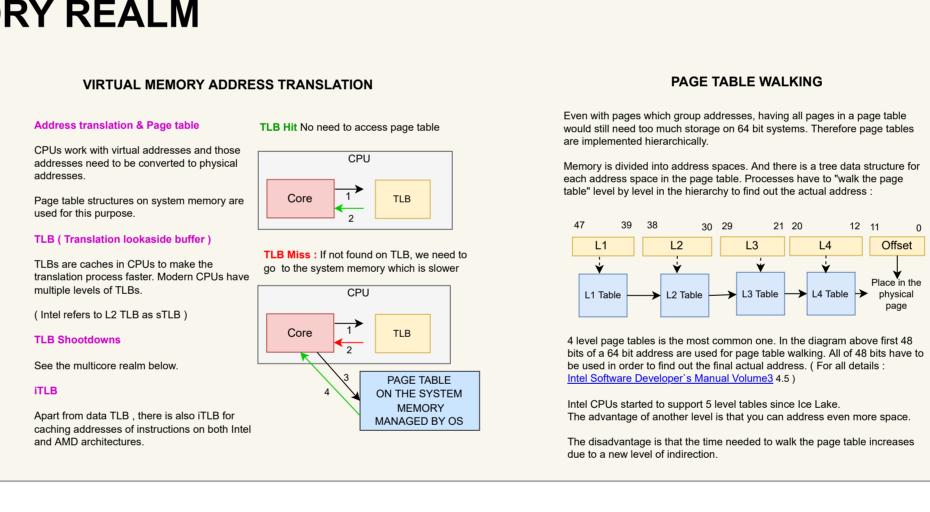
BANK 4

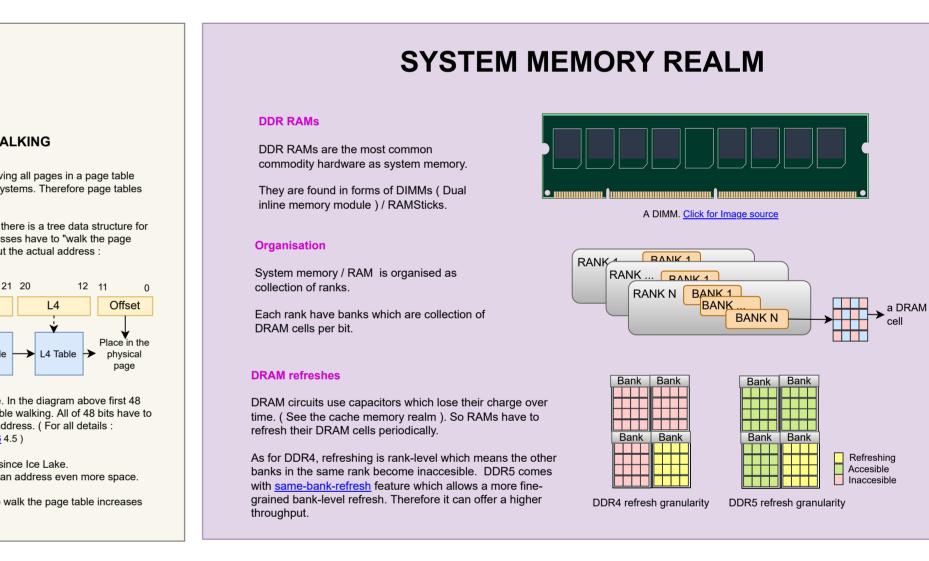


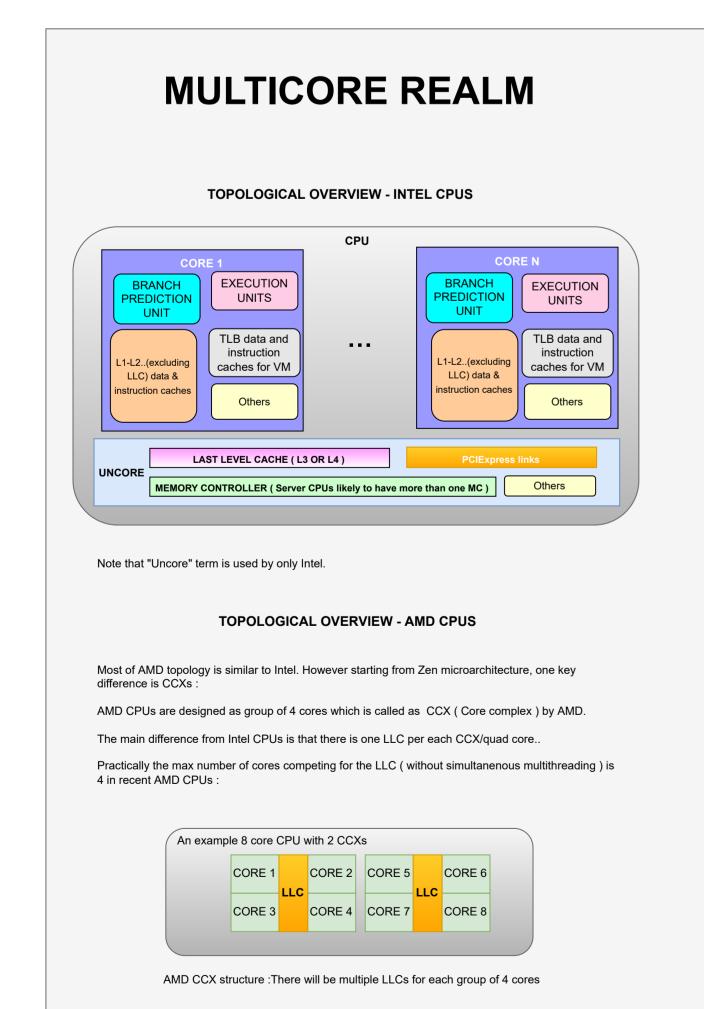
information is stored in bits of addresses. A cache address has 3 parts



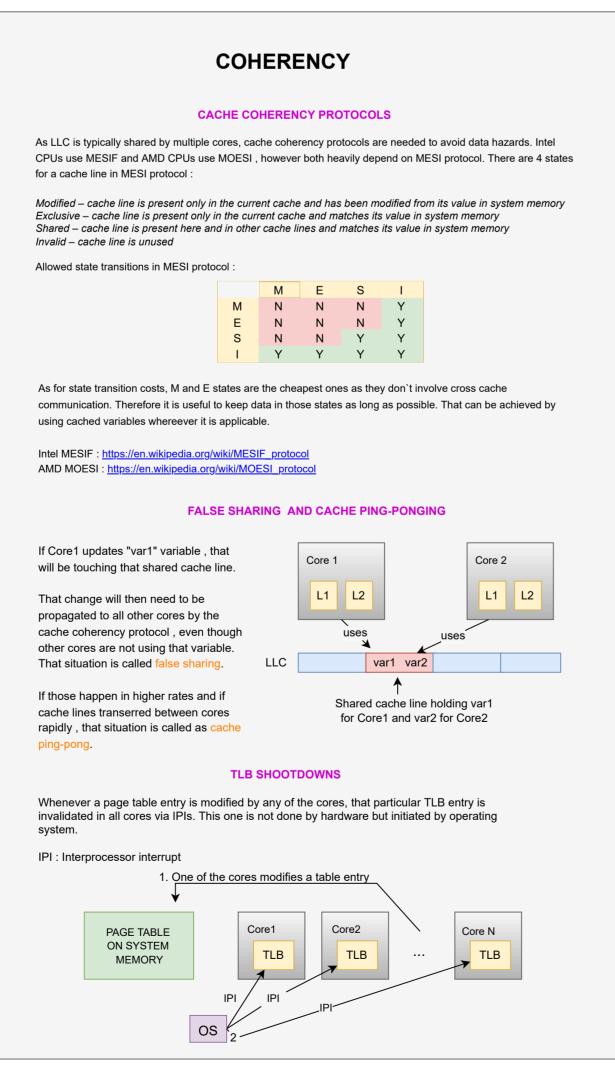
The branch history register will be used to choose which row will be used from the pattern history table.

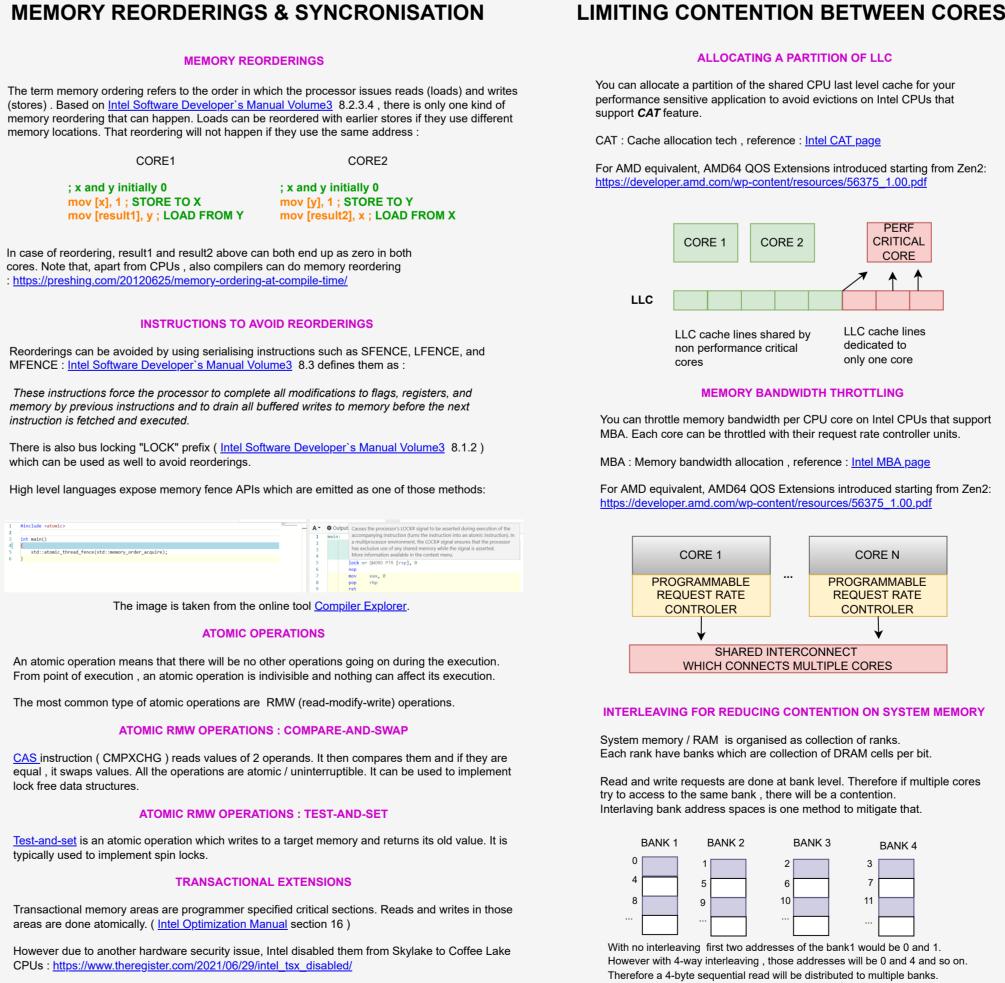


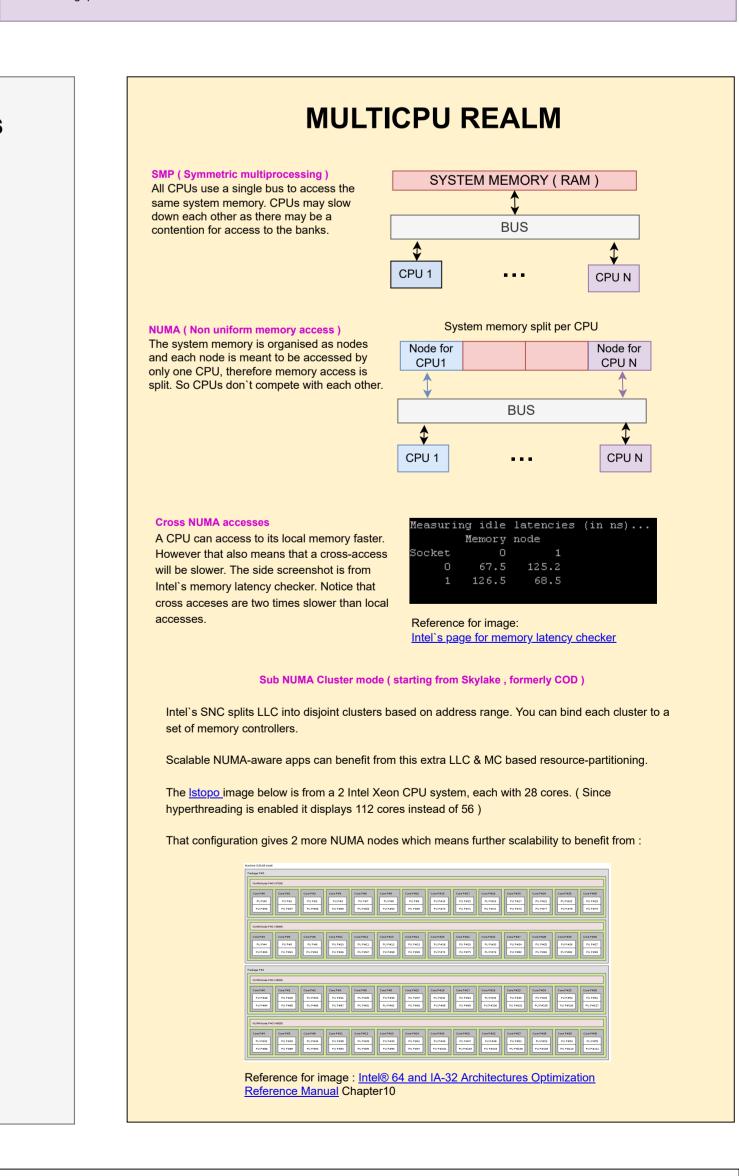




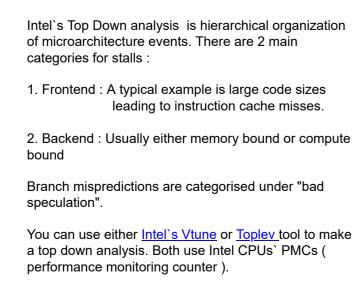
Reference: https://en.wikichip.org/wiki/amd/microarchitectures/zen#CPU_Complex_.28CCX.29

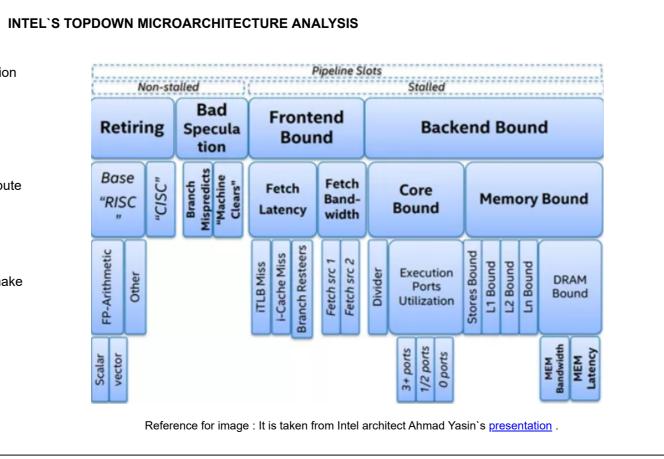


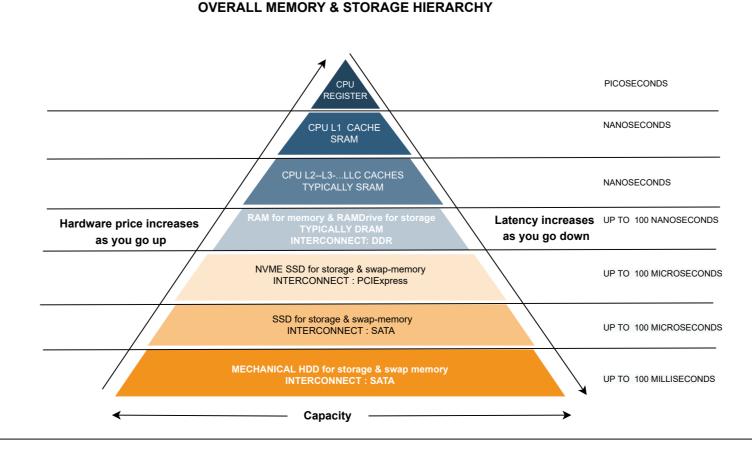




ABOUT REALMS







Floating point division 10-40 clock cycles Simple register operation (ADD OR etc) less than 1 clock cycle Compare and swap atomic op. 15-30 clock cycles Branch prediction 1-2 clock cycles Integer division 15-40 clock cycles Floating point addition 1-3 clock cycles L3 data cache read 30-70 clock cycles Multiplication (int, floating point) 1-7 clock cycles System memory read 100-150 clock cycles L1 data cache read 3-4 clock cycles Cross-NUMA L3 read 100-300 clock cycles TLB miss 7-21 clock cycles Cross-NUMA system memory read 300-500 clock cycles 10-12 clock cycles L2 data cache read Reference for numbers : 10-20 clock cycles Branch misprediction http://ithare.com/infographics-operation-costs-in-cpu-clock-cycles/

(ESTIMATED) LATENCY NUMBERS IN CLOCK CYCLES