

Advanced Operating Systems (and System Security)

MS degree in Computer Engineering

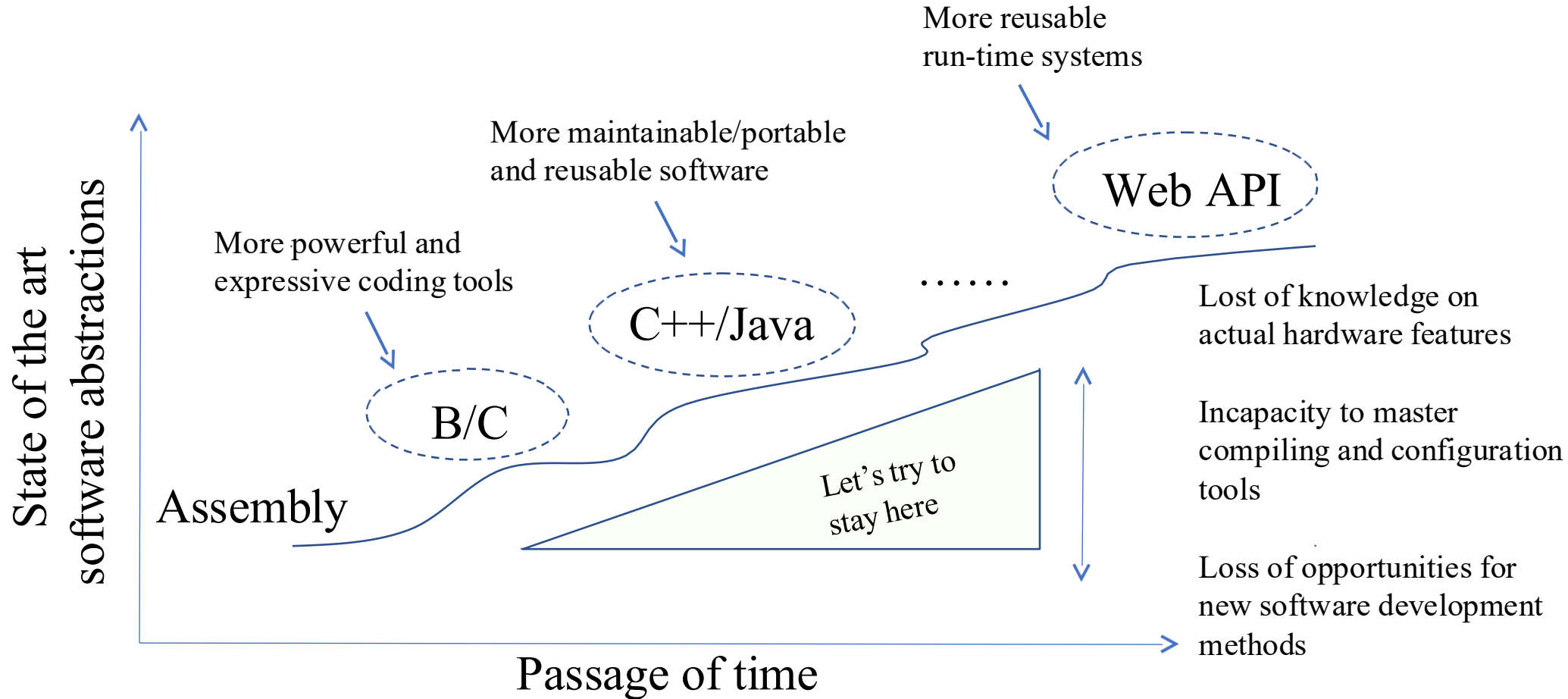
University of Rome Tor Vergata 

Lecturer: Francesco Quaglia

Hardware insights

- Pipelining and superscalar processors
- Speculative hardware
- Multi-processors and multi-cores
- Physical memory organization
- Memory coherency and consistency
- Hardware synchronization support
- Linearizability and thread coordination schemes
- Vectorization
- Hardware related security hints

The need for holistic programming



The missing piece of information

- The actual state of a program is not the “puzzle” of the states of its individual software components
- Each component sees and updates a state portion that is not trivially reflected into the view by other components
- A component even does not know whether hardware state beyond the ISA-exposed one can be affected by its execution
- a malicious component can try to catch if some footprint has been left somewhere in the system while our application components run
- As we will see these footprints represent very hard to tackle issues (e.g. for security as well as performance)

The missing piece of information

- In real systems things may occur in different ways because of
 - ✓ Compiler decisions
 - ✓ Hardware run-time decisions
 - ✓ Availability (vs absence) of hardware features
- More abstractly, there is a combination of software and hardware non-determinism
- Ideally programmers should know of all this to produce **correct, secure and efficient software**
- such a vertical vision where the programmer is aware of what really happens on the underlying hardware is mostly missing

The common fallback ... and our path

- Simply exploit what someone already did (libraries, run-time environments, algorithmic and coding approaches)
- But you should know that this still does not guarantee you are writing (or running) a program the most efficient (or even correct) way
- in the end knowing the hardware and the under-the-hood layers we really work with provides us with possibilities for better achievements in software development
- Nowadays hardware is multi-core and multi-CPU (these are hardware-threads we will look at), which is characterized by some major aspect that need to be reflected onto software programming

We could ideally study those aspects, and proper software design approaches at different levels - an OS kernel is our reference

A very trivial example - Lamport's Bakery

```
var choosing: array[1,n] of boolean;  
    number: array[1,n] of int;
```

```
repeat {  
    choosing[i] := TRUE;  
    number [i] := <max in array number[] + 1>;  
    choosing[i] := FALSE;  
    for j = 1 to n do {  
        while choosing[j] do no-op;  
        while number[j]  $\neq 0$  and (number [j],j) < (number [i],i) do no-op;  
    }  
    <critical region>;  
    number[i] := 0;  
  
}until FALSE
```

Typically no off-the-shelf machine
(unless single-core)
guarantees **globally consistent**
view of this update sequence



Entering a few details

- The machine model we have been used to think of is the von Newman's one
 - ✓ Single CPU abstraction
 - ✓ Single memory abstraction
 - ✓ Single control flow abstraction: fetch-execute-store
 - ✓ Time separated state transitions in the hardware: no more than one in-flight instruction at anytime
 - ✓ Defined memory image at the startup of any instruction
- The modern way of thinking architectures is instead not based on the flow of things as coded in a program, rather on the concept of **scheduling things** (e.g. usage of hardware components) to do something equivalent to that program flow
- Hopefully, the schedule allows **doing stuff in parallel**
- what about programs naturally made up by multiple flows? – this is exactly an OS kernel!!

Types of scheduling

- In the hardware
 - ✓ Instruction executions within a single (speculative) program flow
 - ✓ Instruction executions in parallel (speculative) program flows
 - ✓ Propagation of values within the overall memory (more generally hardware) system
- At software level
 - ✓ Definition of time frames for threads' execution on the hardware
 - ✓ Definition of time frames for activities' execution on the hardware
 - ✓ Software based synchronization supports (thread/task synchronization)

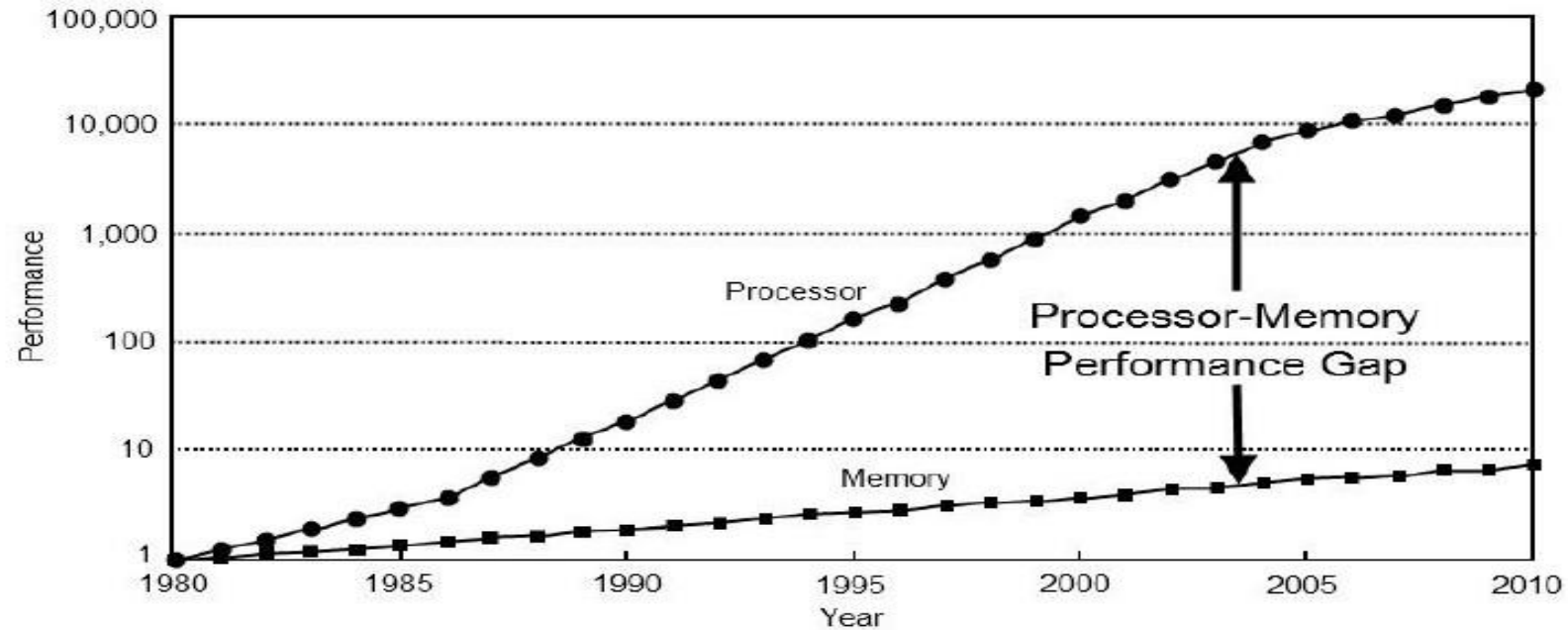
Parallelism

- Baseline hardware form – ILP (Instruction Level Parallelism):
 - ✓ The CPU is able to process 2 or more instructions of a same flow during the same cycle (even vectorized or dynamically scheduled – or both)
 - ✓ It can therefore deliver instruction commits at each individual machine-cycle, even though a single instruction can take several cycles to complete
- Software reflected form (Thread Level Parallelism):
 - ✓ A program can be thought as of the combination of multiple concurrent flows
 - ✓ Concurrency can boil down to actual wall-clock-time parallelism in multi-processing (or ILP) hardware systems

Baseline notion of computing speed

- It is typically related to the Giga Hertz (GHz) rating of a processor
- We clearly know this way of thinking is only partially correct
- There are instructions that can take long sequences of CPU-cycles just because of unpredictable factors
 - ✓ Hardware interactions
 - ✓ Asymmetries and data access patterns
- In the end we can generally think of categories of programs (or programs' blocks) that are more or less importantly affected by the clock speed
 - ✓ CPU-bound programs
 - ✓ Memory-bound programs – that is why we need to know about how to deal with memory in modern systems!!

CPU vs memory performance



CPU memory requests/sec vs DRAM maximum served requests/sec
(using 1980 as a baseline)

Overlapped processing - the pipeline

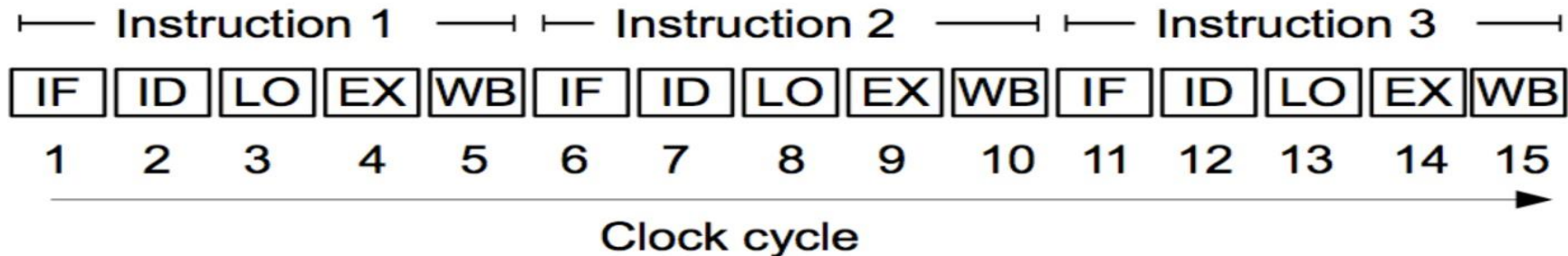
- The very baseline hardware form of overlapped processing is pipelining
- It is a **Scheduling+Parallelism** hardware-based technique
- Here we no longer have a clear temporal separation of the execution windows of different instructions (this is parallelism!!)
- What is sequenced within a program (I'm here referring to an actual executable) is not necessarily executed in that same sequence in the hardware (this is scheduling!!)
- However, causality needs to be preserved
- This is a data flow model (a source should be read based on the actual latest update along the instruction sequence)

Instruction stages

- IF – Instruction Fetch
- ID – Instruction Decode
- LO – Load Operands
- EX – Execute
- WB – Write Back

The different phases hopefully need to rely on different hardware components

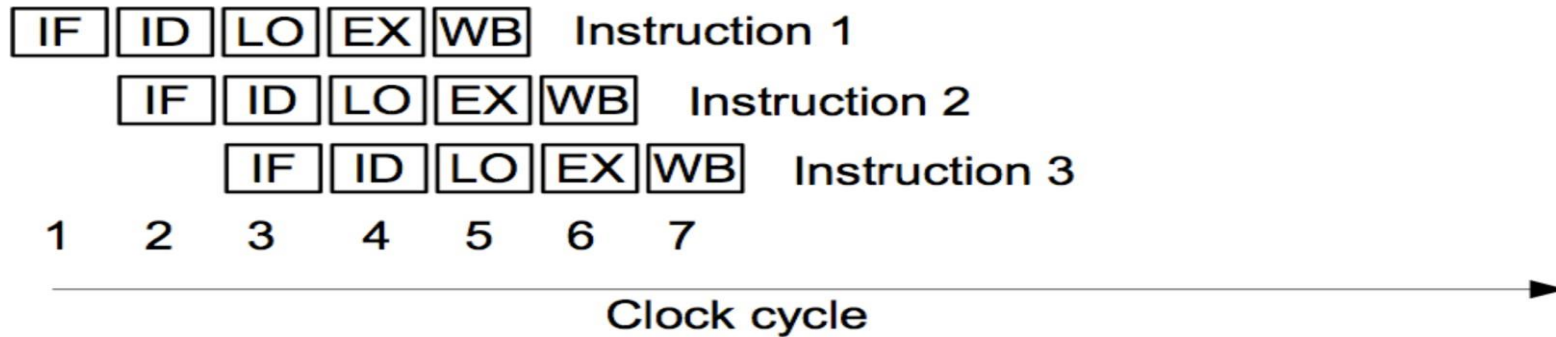
No pipeline



Overlapping stages - the pipeline

- Each instruction uses 1/5 of the resources per cycle
- We can overlap the different phases
- We can therefore get speedup in the execution of a program as compared to the non-pipeline version

Pipeline



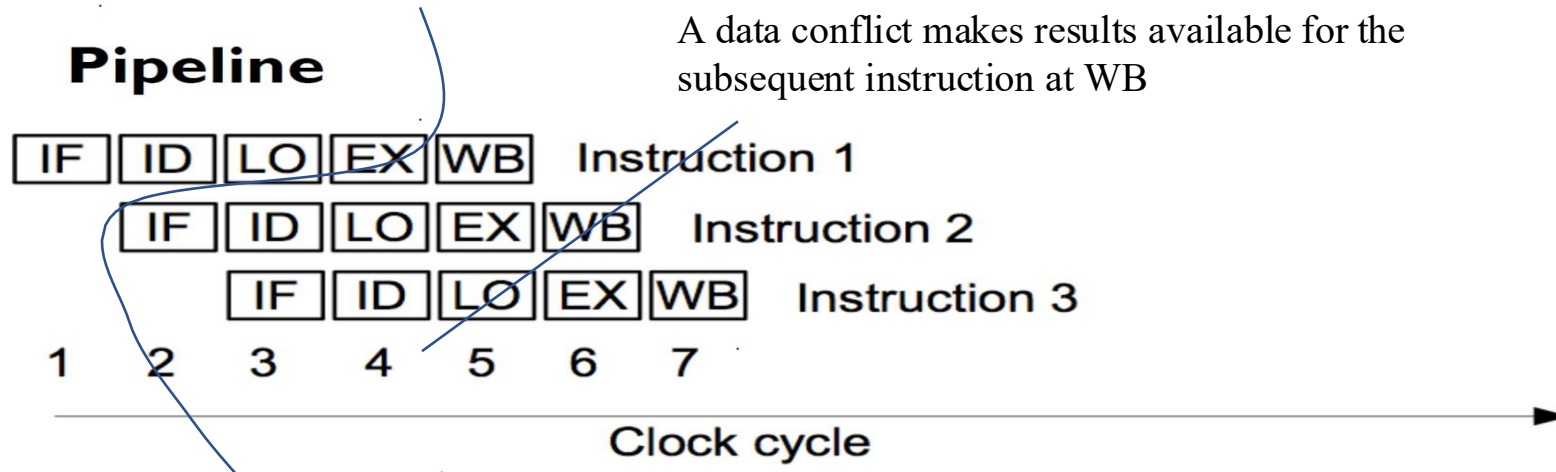
Speedup analysis

- Suppose we want to provide N outcomes (one per instruction) and we have L processing stages and clock cycle T
- With no pipelining we get $(N \times L \times T)$ delay
- With pipelining we get the order of $([N+L] \times T)$ delay
- The speedup is $(N \times L) / (N + L)$ so almost L (for large N)
- For $N = 100$ and $L = 5$ we get 4.76 speedup
- For $L = 1$ no speedup at all arises (obviously!!)
- Ideally the greater L the better the achievable performance
- But we do not live in an ideal world, in fact pipelined processors typically entail no more than the order of tens of stages (Pentium had 5 – i3/i5/i7 have 14 – ARM-11 has 8) although a few implement parts of an original instruction step

From the ideal to the real world - pipeline breaks

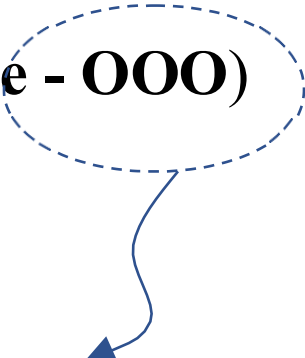
- Data dependencies
- Control dependencies

A conditional branch leads to identify the subsequent instruction at its EX stage



Handling the breaks

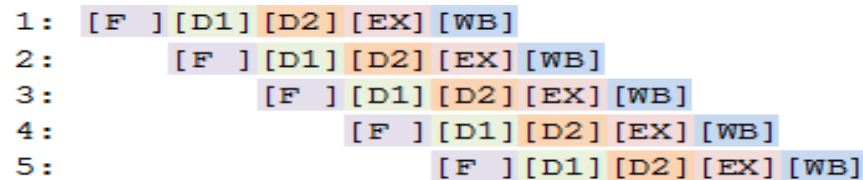
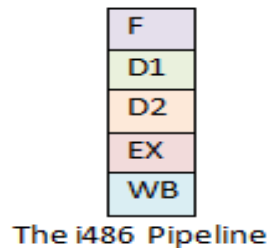
- Software stalls – compiler driven
- Software re-sequencing (or scheduling) – compiler driven
- Hardware propagation (up to what is possible)
- Hardware reschedule (**out-of-order pipeline - OOO**)
- Hardware supported hazards (for branches)



Ordering of the execution steps of the instructions is not based on how they touch ISA exposed hardware components (such as registers)

The Intel x86 pipeline

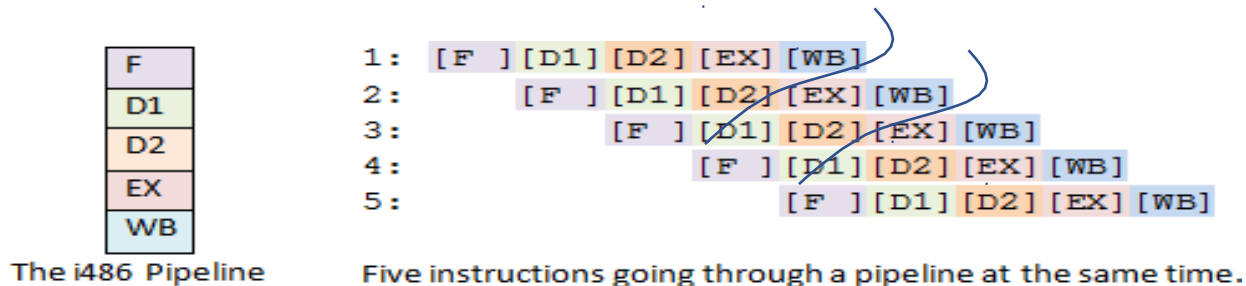
- In a broad analysis, Intel x86 processors did not change that much over time in terms of software exposed capabilities
- The 14 registers (AX, BX, .. etc) of the 8086 are still there on e.g. core-i7 processors (RAX, RBX .. etc)
- However, the 8086 was not pipelined, it processed instructions via [FETCH, DECODE, EXECUTE, RETIRE] steps in pure sequence (not in a pipeline)
- In 1999 the i486 moved to a 5-stage pipeline, with a classical organization plus 2 DECODE steps (primary and secondary – Decode/Translate)



This was for calculations like displacements in a complex addressing model

Pipelining vs software development

- Programmers cannot drive the internal behavior of a pipeline processor (that's microcode!!!!)
- However, the way software is written can hamper the actual pipeline efficiency
- An example – XOR based swap of 2 values:
 - XOR a,b – XOR b,a XOR a,b
- Each instruction has a source coinciding with a destination of the previous instruction



Some examples

- Pointer based accesses plus pointer manipulation should be carefully written
- Writing in a cycle the following two can make a non negligible difference

`a = *++p`

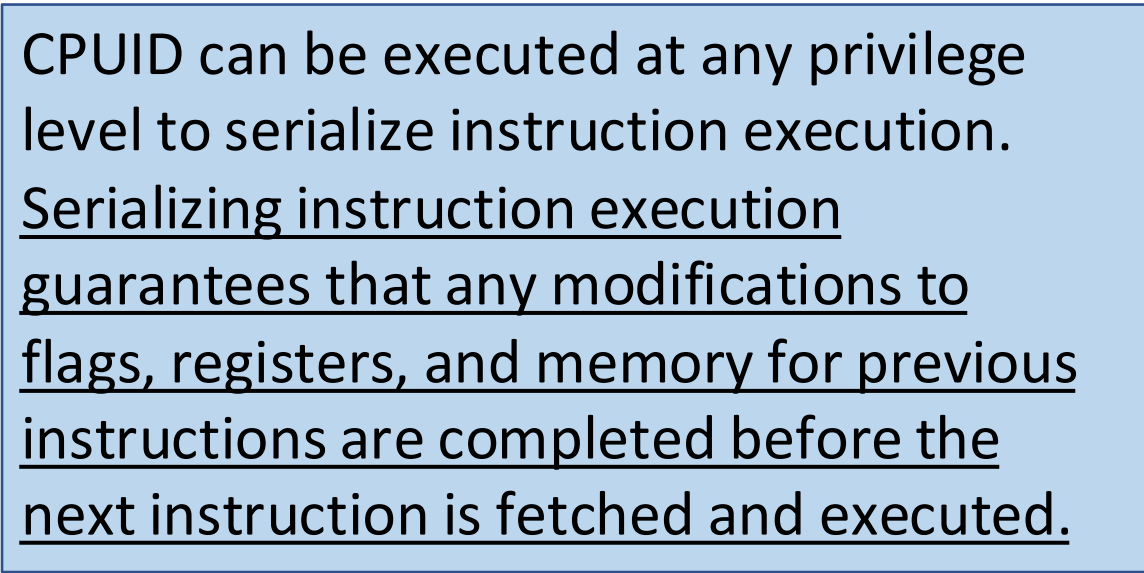
`a = *p++`

- Also, there are machine instructions which lead to flush the pipeline, because of the actual organization of the CPU circuitry
- In x86 processors, one of them is `CPUID` which gets the numerical id of the processor we are working on (and other information)
- On the other hand, using this instruction you are sure that no previous instruction in the actual program flow is still in flight along the pipeline before instructions subsequent to `CPUID` are fetched

Serializing instructions

- The mentioned CPUID instruction on x86 is also referred to as serializing instruction
- From the Intel x86 Instruction Set Reference:

Here we are referring to hardware state that correlates the execution of instructions everywhere in the hardware



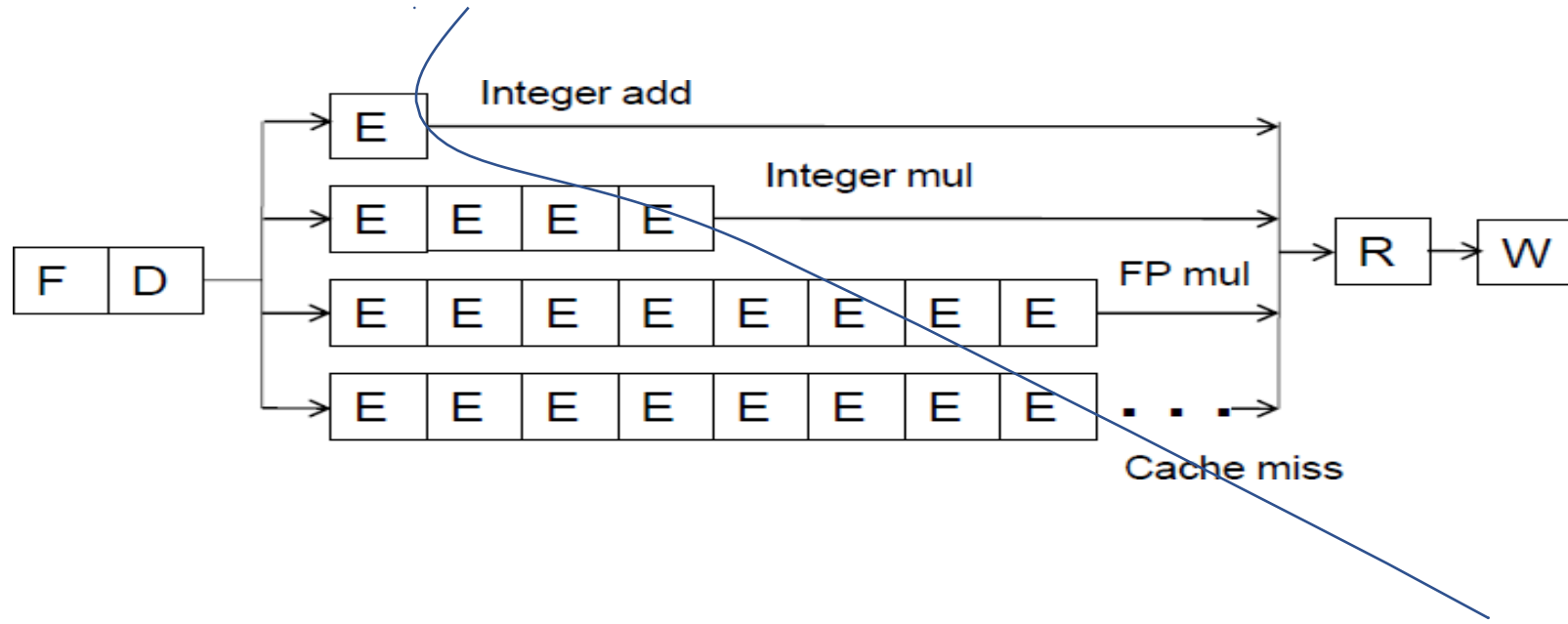
CPUID can be executed at any privilege level to serialize instruction execution. Serializing instruction execution guarantees that any modifications to flags, registers, and memory for previous instructions are completed before the next instruction is fetched and executed.

The Intel x86 superscalar pipeline

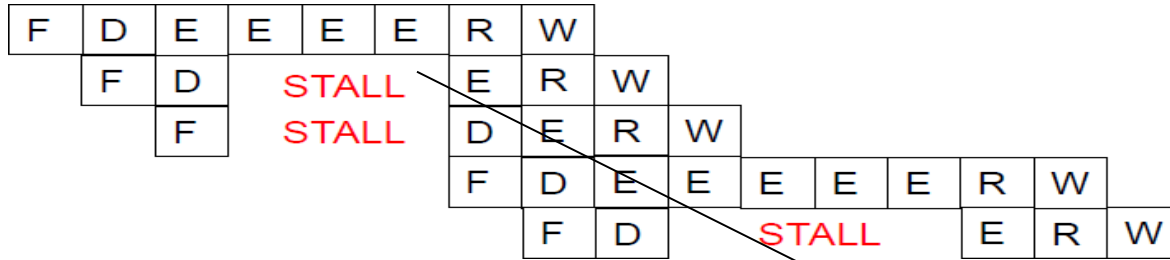
- Multiple pipelines operating simultaneously
- Intel Pentium Pro processors (1995) had 2 parallel pipelines
- EX stages could be actuated in real parallelism thanks to hardware redundancy and differentiation (multiple ALUs, differentiated int/float hardware processing support etc.)
- Given that slow instructions (requiring more processor cycles) were one major issue, this processor adopted the OOO model (originally inspired by Robert Tomasulo's Algorithm – IBM 360/91 1966)
- Baseline idea:
 - ✓ Commit (retire) instructions in program order
 - ✓ Process independent instructions (on data and resources) as soon as possible

The instruction time span problem

Delay reflected into a pipeline execution of independent instructions

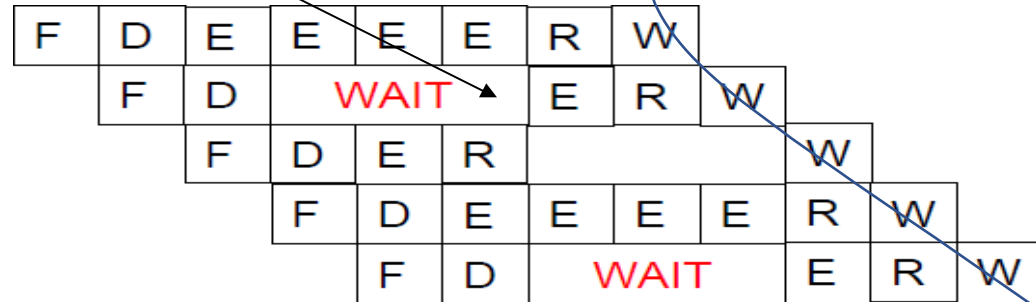


The instruction time span problem



Stall becomes
a reschedule

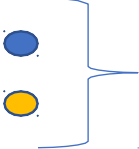
Commit order needs to be
preserved because of, e.g.
WAW (Write After Write)
conflicts



OOO pipeline - speculation

- **Emission**: the action of injecting instructions into the pipeline
- **Retire**: The action of committing instructions and making their side effects “visible” in terms of ISA exposed architectural resources
- What’s there in the middle between the two?
- An execution phase in which the different instructions can surpass each other
- Core issue (beyond data/control dependencies): **exception preserving!!!**
- OOO processors may generate **imprecise exceptions** such that the processor/architectural state may be different from the one that should be observable when executing the instructions along the original order

OOO example

 Different instructions

Program flow

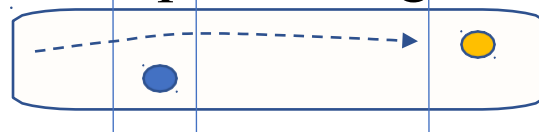


Pipeline stages



.....

Pipeline stages



.....

Program flow
speculatively subverted
No ISA exposed resource
is modified

Pipeline stages



.....

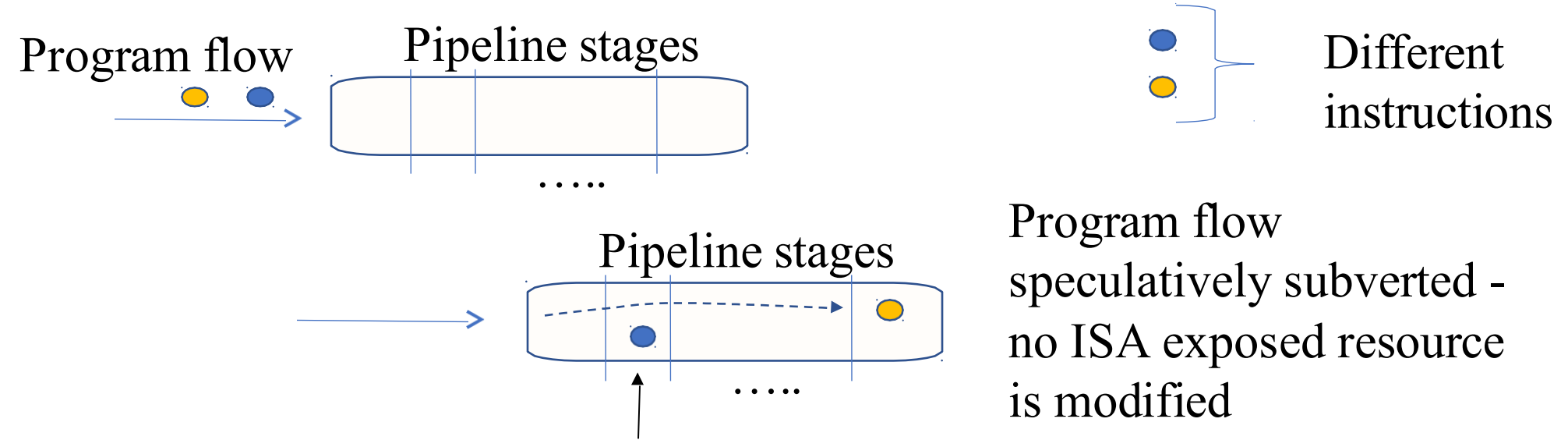
Program flow
maintained



Imprecise exceptions

- The pipeline may have already executed an instruction A that, along program flow, is located after an instruction B that causes an exception
- Instruction A may have changed the micro-architectural state, although finally not committing its actions onto ISA exposed resources (registers and memory locations updates) – the reknown Meltdown security attack exactly exploits this feature
- The pipeline may have not yet completed the execution of instructions preceding the offending one, so their ISA exposed side effects are not yet visible upon the exception
- we will be back with more details later on

A scheme



If this instruction accesses to some invalid resource (e.g. memory location, or currently un-accessible in-CPU components) that program flow is no longer valid and the other instruction cannot currently provide a valid execution, **but something in the hardware may have already happened along its processing**

Robert Tomasulo's algorithm

- Let's start from the tackled hazards – the scenario is of two instructions A and B such that $A \rightarrow B$ in program order:
- RAW (Read After Write) – B reads a datum before A writes it, which is clearly stale – this is a clear data dependency
- WAW (Write After Write) – B writes a datum before A writes the same datum – the datum exposes a stale value
- WAR (Write After Read) – B writes a datum before A reads the same datum – the read datum is not consistent with data flow (it is in the future of A 's execution)

Algorithmic ideas

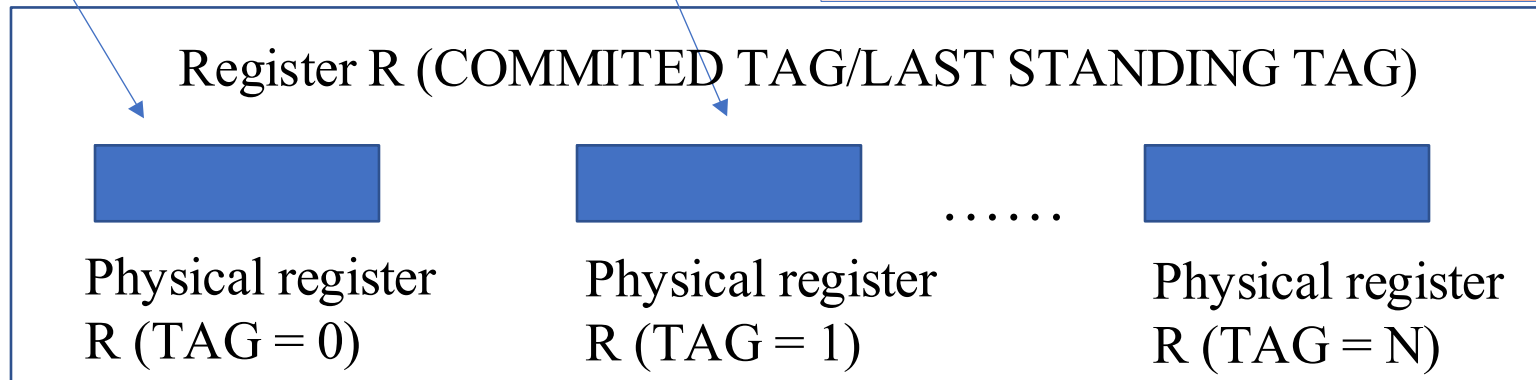
- RAW – we keep track of “when” data requested in input by instructions are ready
- Register renaming for coping with both WAR and WAW hazards
- In the renaming scheme, a source operand for an instruction can be either an actual register label, or another label (a renamed register)
- In the latter case it means that the instruction needs to read the value from the renamed register, rather than from the original register
- A renamed register materializes the concept of speculative (not yet committed) register value, made anyhow available as input to the instructions

Register renaming principles

Write instruction
generates a new
standing tag

Read instruction
gets last standing tag

Standing and commit TAGs are
re-conciliated upon instruction
retirement



Reservation stations

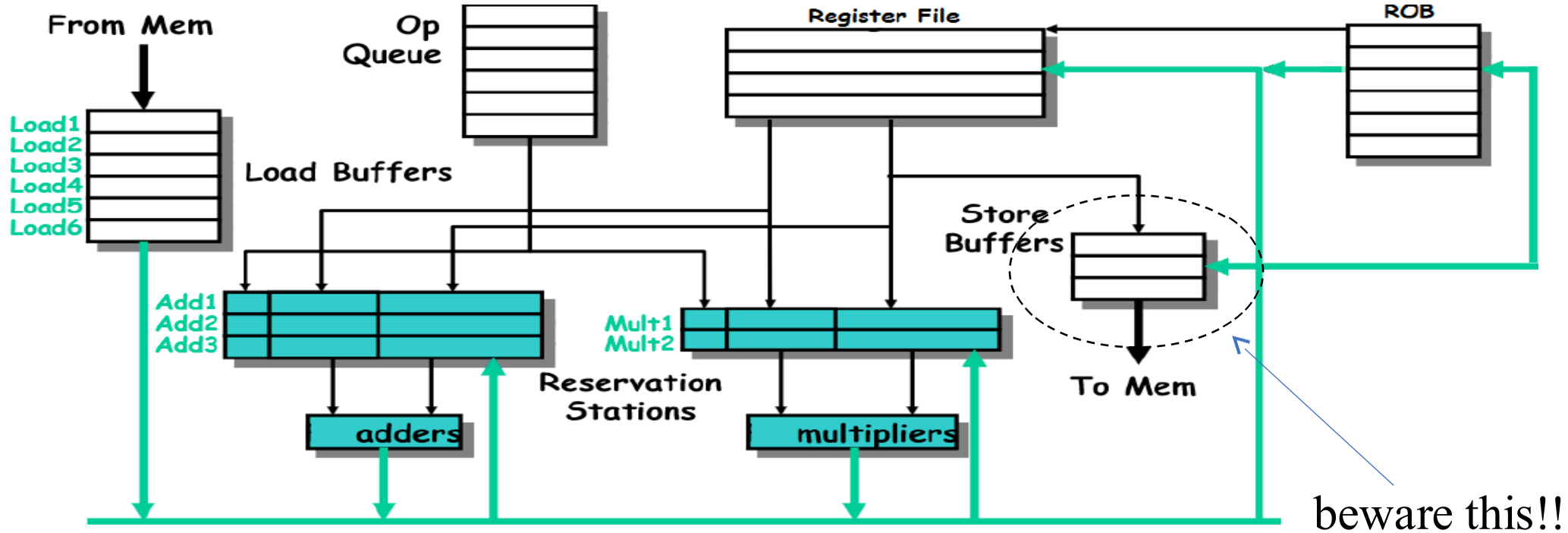
- They are buffers (typically associated with different kinds of computational resources – integer vs floating point HW operators)
- They contain:
 - OP – the operations to be executed (clearly its code)
 - Q_j, Q_k – the reservation stations that will produce the input for OP
 - Alternatively, V_j, V_k , the actual values (e.g. register values) to be used in input by OP
- By their side, registers are marked with the reservation station name Q such that it will produce the new value to be installed, if any

CDB and ROB

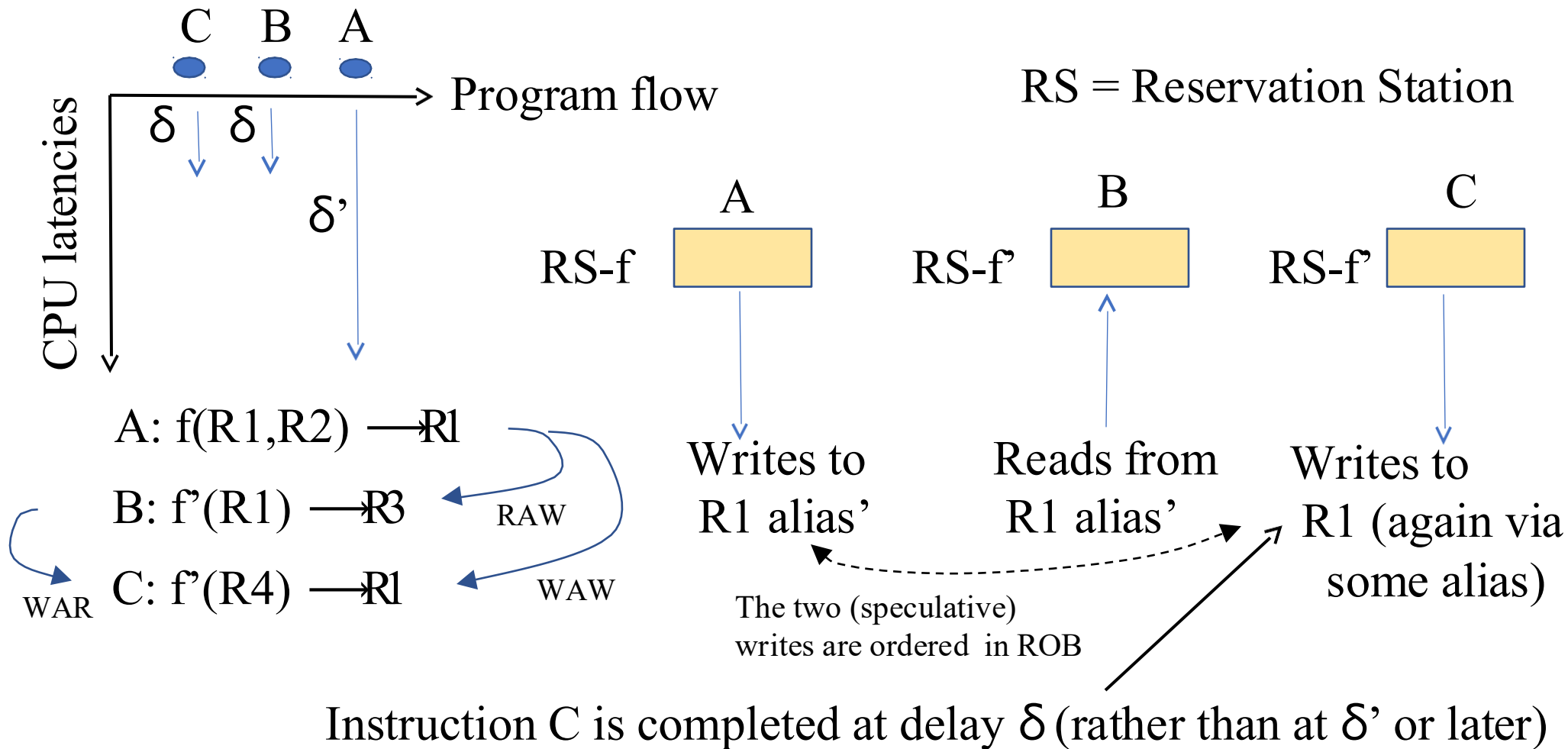
- A Common Data Bus (CDB) allows data to flow across reservation stations (so that an operation is fired when all its input data are available)
- A Reorder Buffer (ROB) keeps the metadata of the in-flight instructions
- It also acquires all the newly produced instruction values (also those transiting on CDB), and keeps them uncommitted up to the point where the instruction is retired
- This is done either directly, or via a reference to the register alias that needs to maintain the value
- ROB is also used for input to instructions that need to read from uncommitted values

An architectural scheme

ROB (Reorder Buffer) is always filled in order with program flow instructions



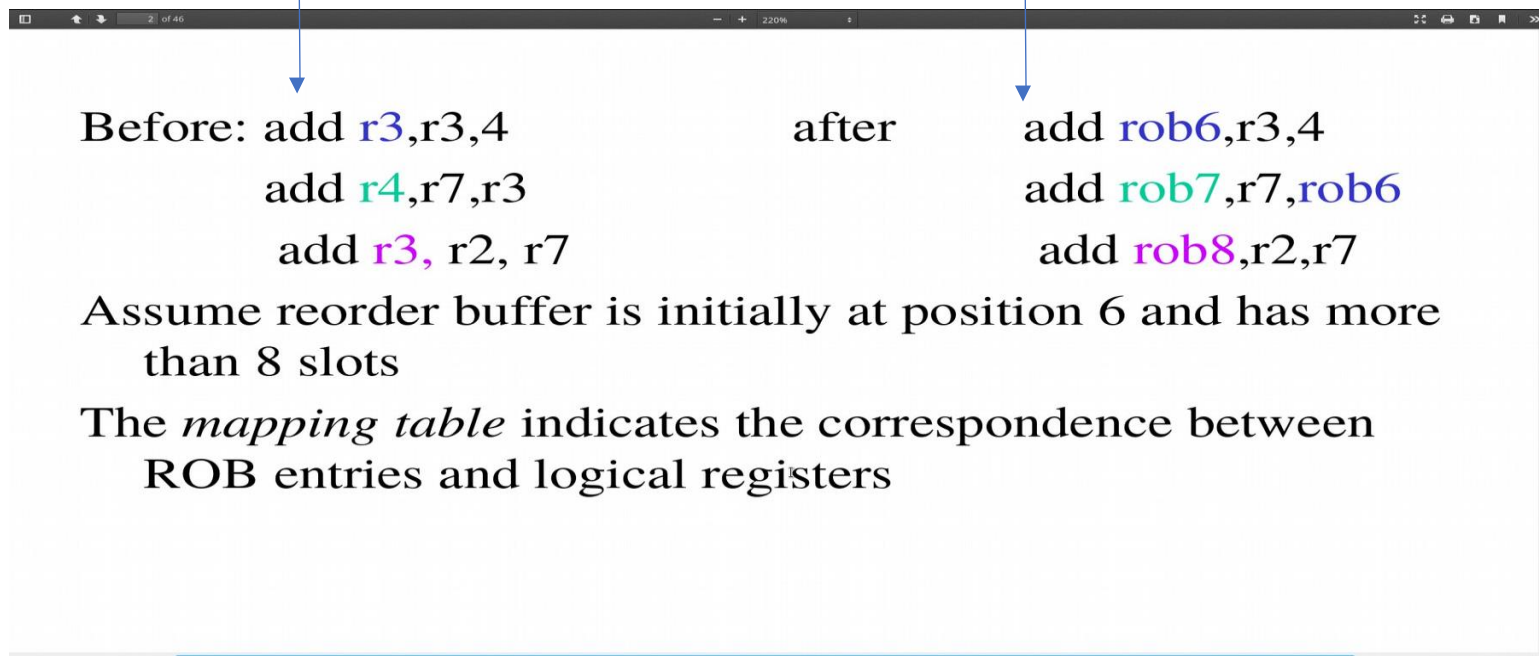
An example execution scheme



Another example with linkage to ROB

What seen at fetch time

What seen inside the OOO pipeline



Here we find data on the mapping table

`This stuff is kept by ROB

add rob6, r3, 4

add rob7, r7, rob6

add rob8, r2, r7

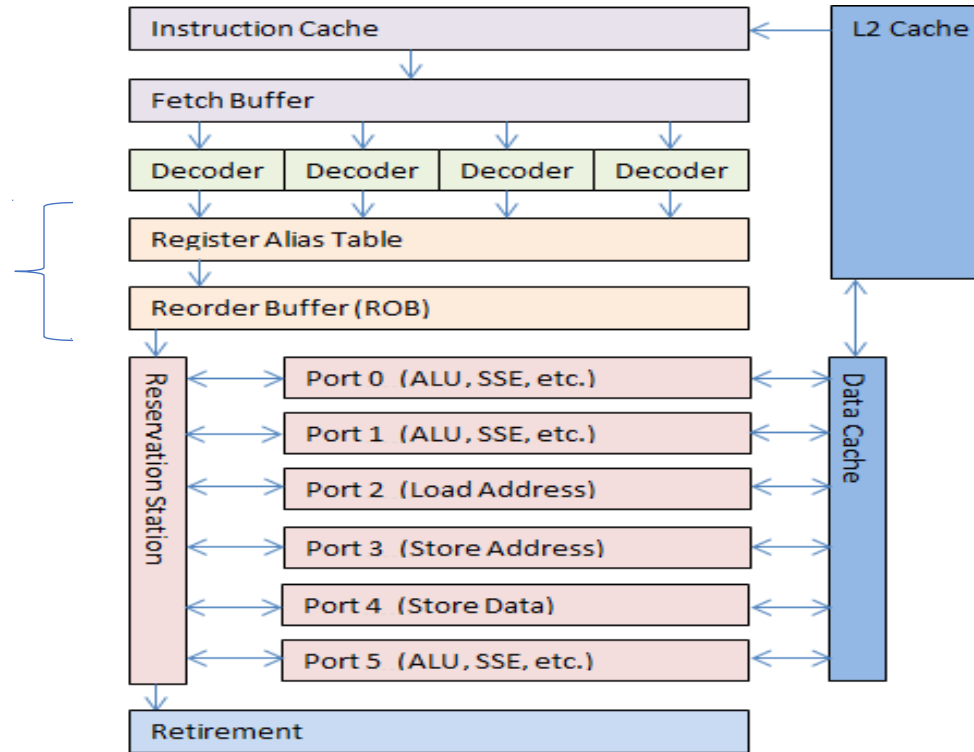
Source	Destination
-	R3 i-th renamed instance
R3 i-th renamed instance	R4 k-th renamed instance
-	R3 (i+1)-th renamed instance

Usage along history

- Tomasulo's Algorithm has been originally implemented in the IBM 360 processor (more than 40 years ago)
- Now it is commonly used (as baseline approach for handling OOO instructions execution) in almost all off-the-shelf processors
- Originally (in the IBM 360) it only managed floating point instructions
- Now we handle the whole instruction set
- Originally (in the IBM 360) it was able to handle a few instructions in the OOO window
- Now we typically deal with up to the order of 100 instructions!!

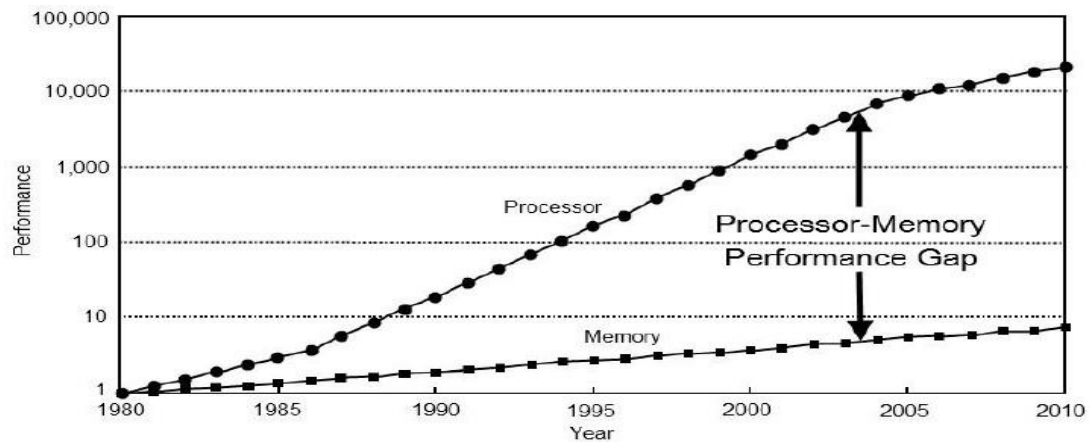
x86-64 main architectural organization

Who depends on who?

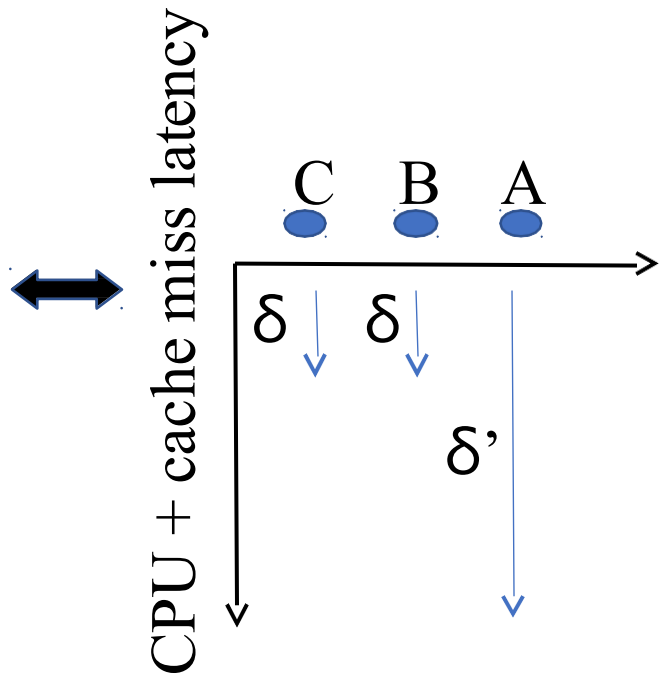


Out Of Order Core as used from 1995 to present. The color coding follows the same five stages used in previous processors. Some stages and buffers are not shown since they vary from processor to processor.

Back to the memory wall



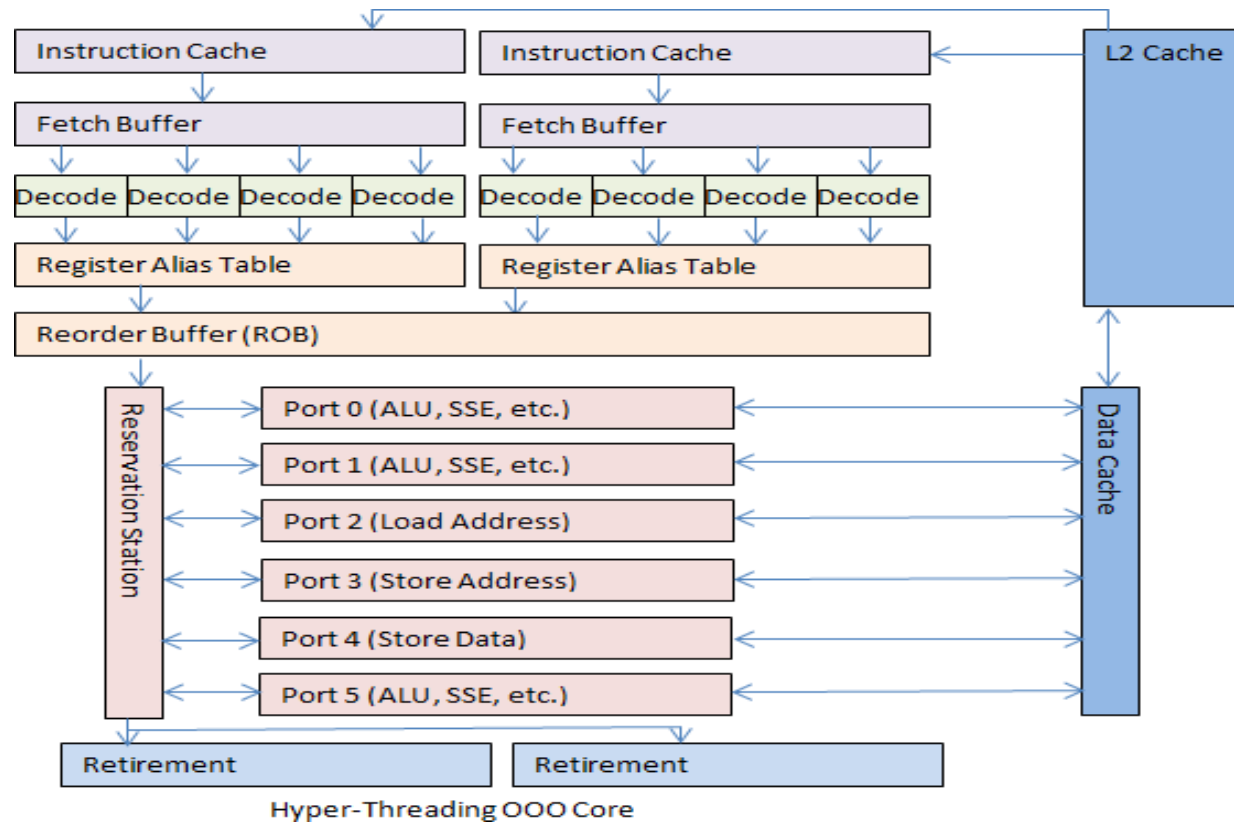
CPU memory requests/sec vs DRAM maximum served requests/sec
(using 1980 as a baseline)



Impact of OOO in x86

- OOO allowed so fast processing of instructions that room was still there on core hardware components to carry out work
- ... also because of delays within the memory hierarchy
- ... why not using the same core hardware engine for multiple program flows?
- This is called hyper-threading, and is actually exposed to the programmer at any level (user, OS etc.)
- ISA exposed registers (for programming) are replicated, as if we had 2 distinct processors
- Overall, OOO is not exposed (instructions are run as in a black box) although the way of writing software can impact the effectiveness of OOO and more generally of pipelining

Baseline architecture of OOO Hyper-threading



Nomenclature and facilities

- An hyperthread processor is typically marked with the explicit indication on how many hyperthreads it has and how many cores it has
- A classical example for a two-core and four-threads processor → **2C/4T**
- On Linux you can easily check if hyperthreading is supported on your processor by relying on the pseudo-file */proc/cpuinfo* in particular by checking the presence of the *ht* flag
- Hyperthreading can be excluded using bios
- At runtime you can check if it is active using the command “sudo dmidecode -t processor | grep Count”

Coming to interrupts

- Interrupts typically flush all the instructions in the pipeline, as soon as one is committed and the interrupt is accepted
- As an example, in a simple 4-stage pipeline IF, ID, EX, MEM residing instructions are flushed because of the acceptance of the interrupt on the WB phase of the currently finalized instruction
- This avoids the need for handling priorities across interrupts and exceptions possibly caused by instructions that we might let survive into the pipeline (no standing exception)
- Interrupts may have a significant penalty in terms of wasted work on modern OOO based pipelines
- Also, in flight instructions that are squashed may have changed the micro-architectural state on the machine

Back to exceptions - types vs pipeline stages

- Instruction Fetch, & Memory stages
 - Page fault on instruction/data fetch
 - Misaligned memory access
 - Memory-protection violation
- Instruction Decode stage
 - Undefined/illegal opcode
- Execution stage
 - Arithmetic exception
- Write-Back stage
 - *No exceptions!*

Back to exceptions - handling

- When an instruction in a pipeline gives rise to an exception, the latter is not immediately handled
- As we shall see later, such instruction in fact might even require to disappear from program flow (as an example because of miss-prediction in branches)
- It is simply marked as **offending** (with one bit traveling with the instruction across the pipeline)
- When the retire stage is reached, the exception takes place and the pipeline is flushed, resuming fetch operations from the right place in memory
- **NOTE:** micro architectural effects of in flight instructions that are later squashed (may) still stand there – see the Meltdown attack against Intel and ARM processors ...

Coming to an example

Program flow



Pipeline stages

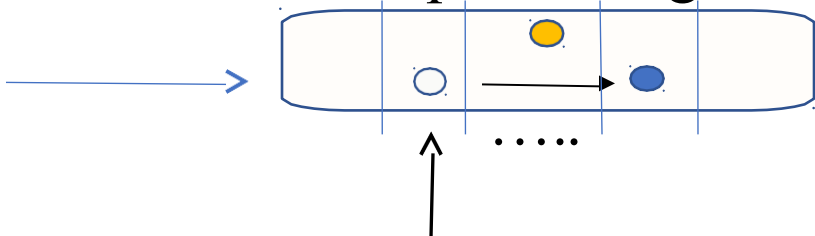


.....



Different
instructions

Pipeline stages



Program flow
speculatively subverted -
no ISA exposed resource
is modified

Offends and goes forward, and also propagates “alias” values to the other instruction, which goes forward up to the squash of the first

Meltdown primer

Flush cache

Read a kernel level byte B

Use B for displacing and reading memory

A sequence with
imprecise exception under
OOO

Offending instruction
(memory protection violation)

“Phantom” instruction with real
micro-architectural side effects

A graphical representation of what happens

Flush cache

The cache

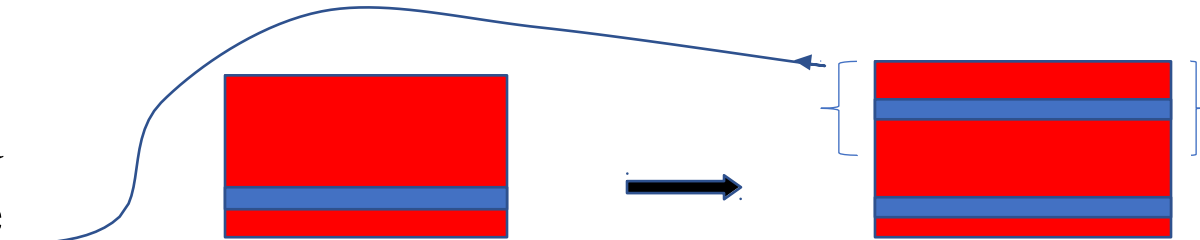


Read a **kernel level** byte B



Loading lines in cache
is not an ISA exposed
effect

Use B for displacing and
reading memory in some
known zone, say this



If we can measure the access delay for hits
and misses when reading that zone, we
would know what was the value of B

A “performance” note

- This attack is based on the processor ability to speculatively execute an access to legitime memory **after another access** to un-legitime memory
- We can expect that the faster the legitimate access can be served in the pipeline the higher the likelihood **to have it fully executed** after the un-legitimate access takes place
- hence, we can expect that the attack can more likely be carried out positively if the kernel level byte is already in cache when we try to get it

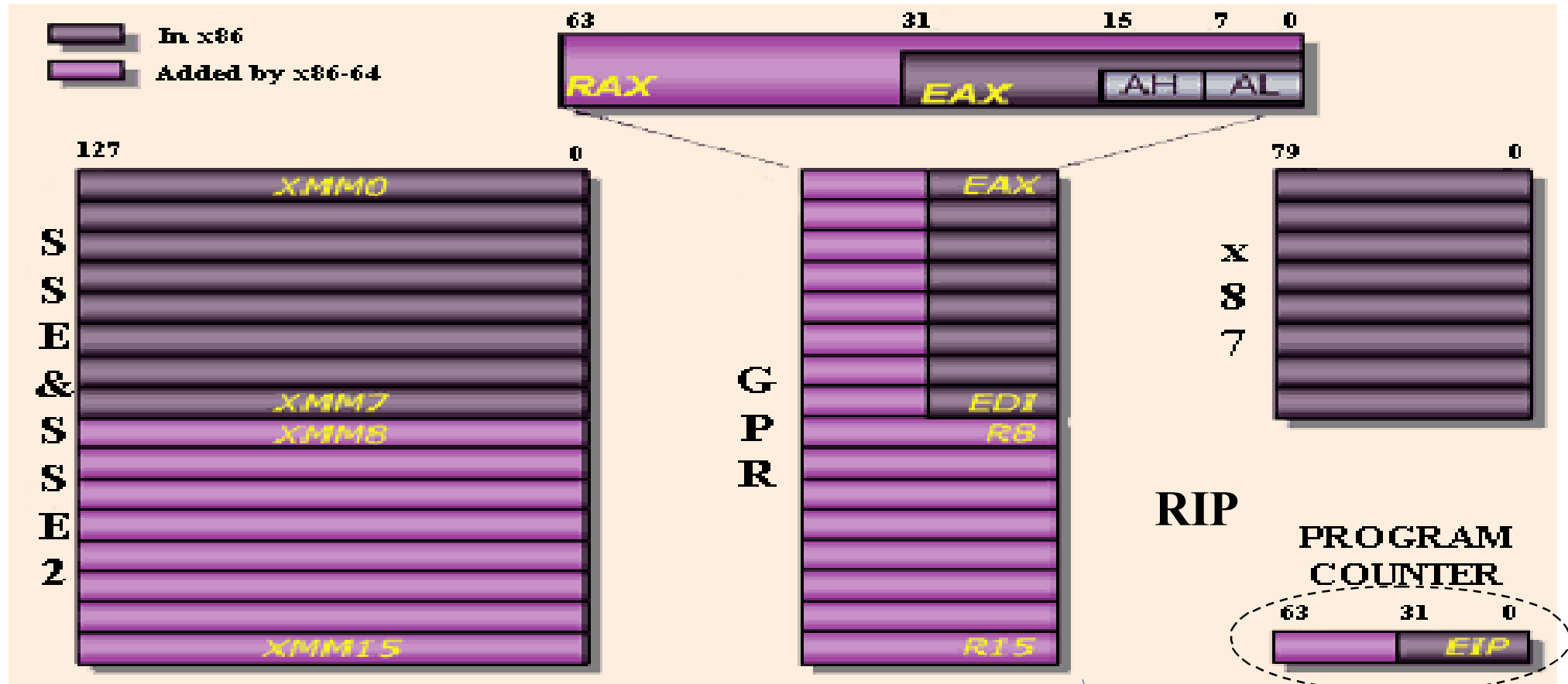
Meltdown timeline and affected processor families

- The vulnerability was independently discovered by multiple security researchers in mid-2017
- The details of Meltdown were made public on January 3, 2018, along with another major vulnerability known as Spectre
- After its disclosure, OS vendors (like Microsoft, Linux, and Apple) quickly began releasing patches to mitigate the vulnerability, though these sometimes came with performance impacts
- It affected
 - ✓ Intel: Most Intel processors from 1995–2017 are vulnerable to Meltdown
 - ✓ ARM: Some ARM Cortex-A series processors (like Cortex-A75) are vulnerable
 - ✓ AMD: Largely unaffected by Meltdown due to differences in architecture

Overall

- The cache content, namely the state of the cache can change depending on processor design and internals, not necessarily at the commitment of instructions
- Such content is in fact not directly readable in the ISA
- We can only read from logical memory, so a program flow would ideally be never affected by the fact that a datum is or is not in cache at a given point of the execution
- The only thing really affected is performance
- **But this may say something to who exactly observes performance to infer the actual state of the caching system**
- **This is a so called side-channel (or covert-channel) attack**

Additional x86 details



The instruction set - data transfer examples (AT&T syntax)

mov{b,w,l} source, dest

General move instruction

push{w,l} source

pushl %ebx # equivalent instructions
 subl \$4, %esp
 movl %ebx, (%esp)

pop{w,l} dest

popl %ebx # equivalent instructions
 movl (%esp), %ebx
 addl \$4, %esp

Variable length operands

movb \$0x4a, %al #byte
movw \$5, %ax #16-bit
movl \$7, %eax #32-bit

No operand-size specification
means (the default) 64-bit
operand on x86-64

The instruction set - linear addressing (32-bit scene)

$$\text{Offset} = \left(\begin{array}{c} \text{eax} \\ \text{ebx} \\ \text{ecx} \\ \text{edx} \\ \text{esp} \\ \text{ebp} \\ \text{esi} \\ \text{edi} \end{array} \right) + \left(\begin{array}{c} \text{eax} \\ \text{ebx} \\ \text{ecx} \\ \text{edx} \\ \text{esp} \\ \text{ebp} \\ \text{esi} \\ \text{edi} \end{array} \right) * \left(\begin{array}{c} 1 \\ 2 \\ 3 \\ 4 \end{array} \right) + \left(\begin{array}{c} \text{None} \\ 8\text{-bit} \\ 16\text{-bit} \\ 32\text{-bit} \end{array} \right)$$

Base Index scale displacement

Displacement → `movl foo, %eax`

Base → `movl (%eax), %ebx`

Base + displacement → `movl foo(%eax), %ebx`
`movl 1(%eax), %ebx`

(Index * scale) + displacement → `movl (,%eax,4), %ebx`

Base + (index * scale) + displacement → `movl foo(%ecx,%eax,4), %ebx`

The instruction set - bitwise logical instructions (base subset)

and{b,w,l} source, dest	dest = source & dest
or{b,w,l} source, dest	dest = source dest
xor{b,w,l} source, dest	dest = source ^ dest
not{b,w,l} dest	dest = 0xFF .. FF ^dest

sal{b,w,l} source, dest (arithmetic)

dest = dest << source

sar{b,w,l} source, dest (arithmetic)

dest = dest >> source

The instruction set - arithmetic (base subset)

add{b,w,l} source, dest dest = source + dest

sub{b,w,l} source, dest dest = dest - source

inc(b,w,l} dest dest = dest + 1

dec{b,w,l} dest dest = dest - 1

neg(b,w,l} dest dest = 0 - dest

cmp{b,w,l} source1, source2

source2 - source1

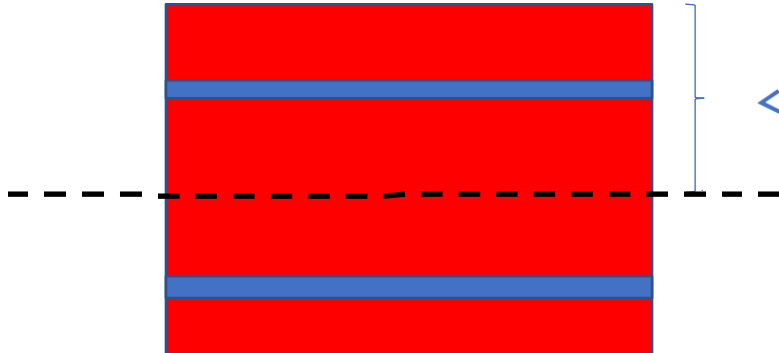
The Meltdown code example – Intel syntax (mostly reverts operand order vs AT&T)

```
1 ; rcx = kernel address
2 ; rbx = probe array
3 retry:
4 mov al, byte [rcx]
5 shl rax, 0xc
6 jz retry
7 mov rbx, qword [rbx + rax]
```

This is B

Use B as the index of a page

B becomes the displacement of a given page in an array



The target cache zone is an array of 256 pages – only the 0-th byte of the B-th page will experience a cache hit (under the assumption that concurrent actions are not using that cache zone)

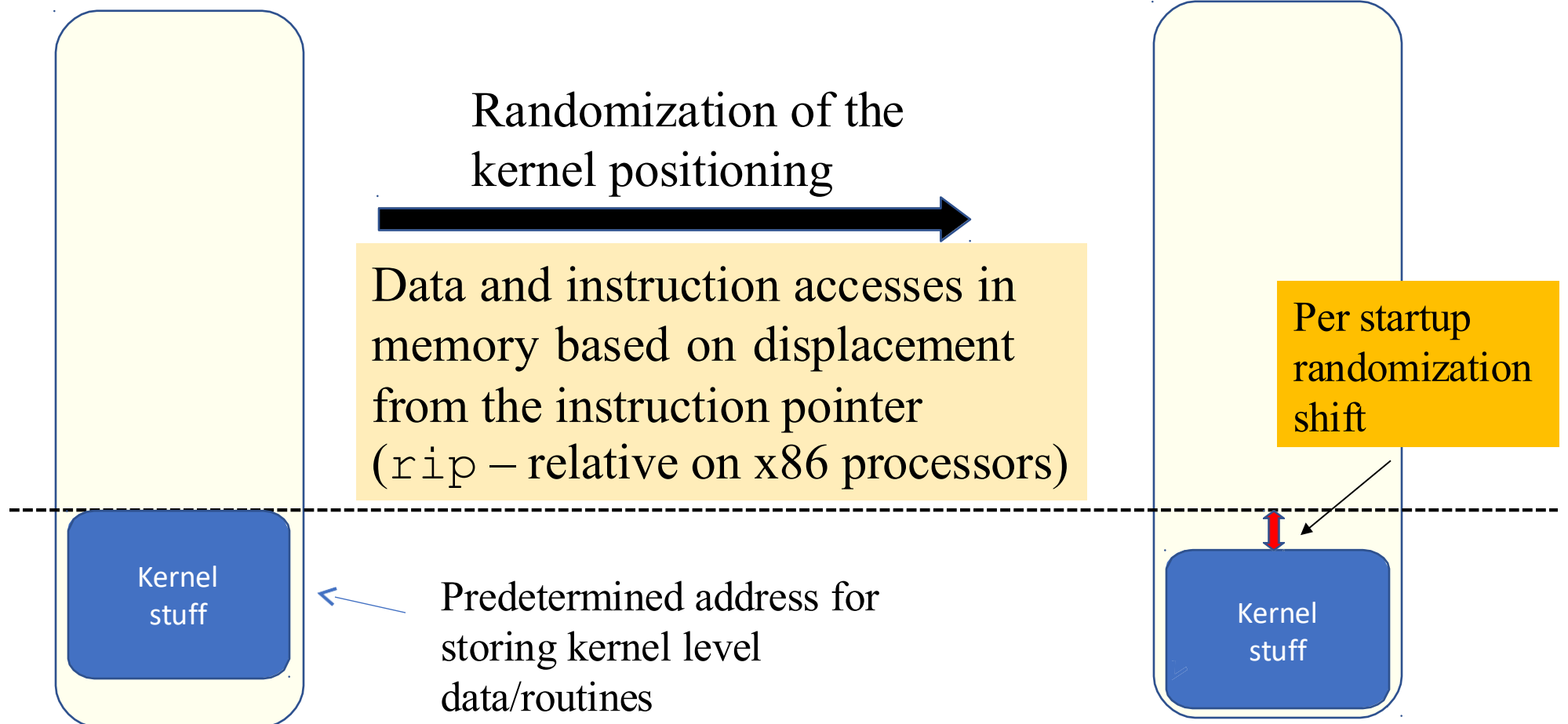
Countermeasures

- KASLR (Kernel Address Space Randomization) – limitation of being dependent on the maximum shift we apply to the logical kernel image (this varies on different Linux Kernel instances – it is typically enabled by default) - clearly this is still weak vs brute force attacks
- KAISER (Kernel Isolation in Linux) – still exposes the interrupt surface but it is highly effective
- Explicitly cache-flush at each return from kernel mode – detrimental for performance and still not fully resolving as we will discuss
- ... clearly processor patches on newer CPU releases (but we have a plethora of affected processors)

A scheme for KASLR

Classical virtual address space usage

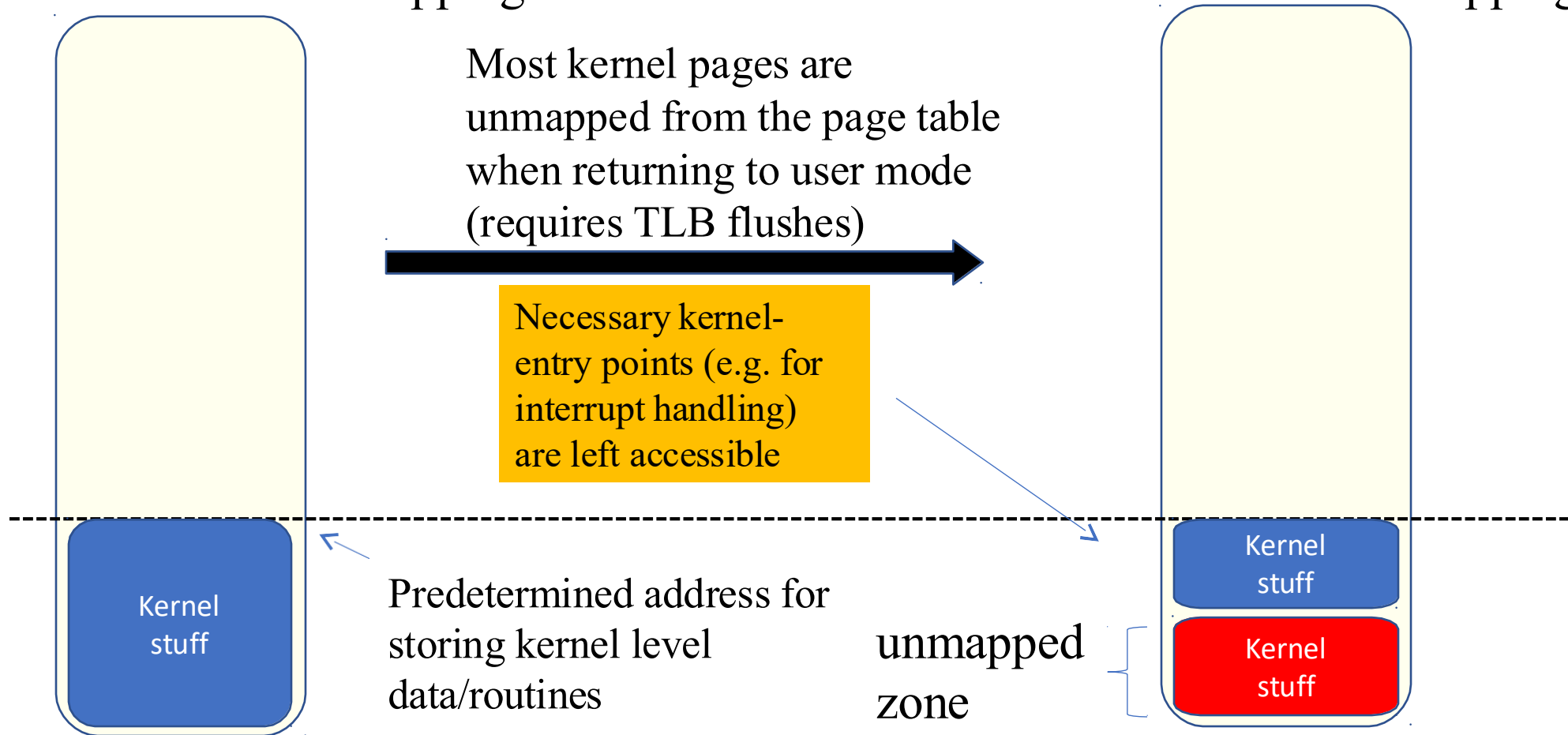
Usage with randomization



A scheme for KAISER

Classical kernel mapping

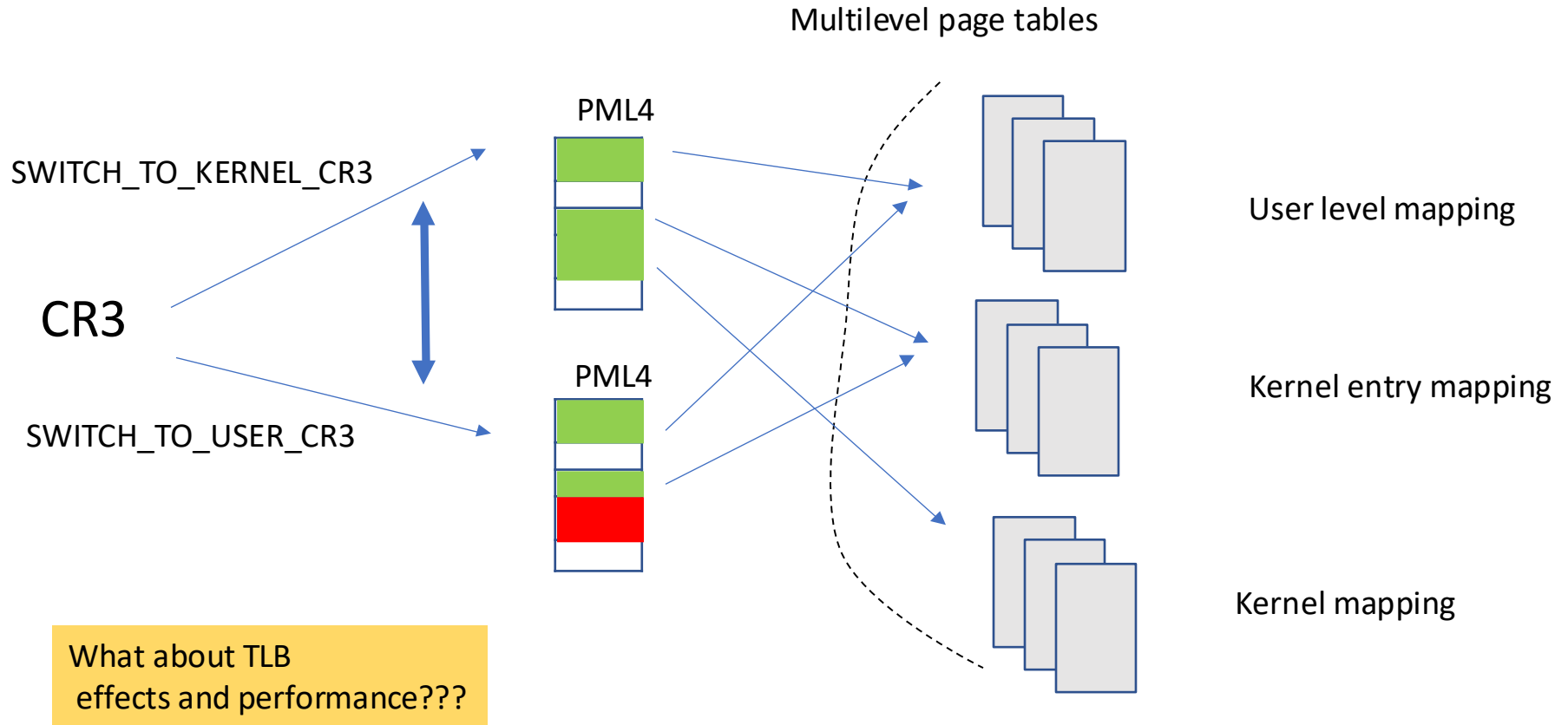
Variation of the kernel mapping



More details on the Linux world

- KAISER (Kernel Isolation in Linux) is technically denoted by the acronym PTI (Page Table Isolation)
- It is directly available in the kernel source code since kernel 4
- It clearly impacts performance, but can be disabled at kernel startup
- This can be done using the `pti=off` specification at the level of the kernel parameters (via GRUB)
- `via/sys/devices/system/cpu/vulnerabilities/meltdown` you can check if this security patch is up or not

Actual implementation in Linux (on x86)



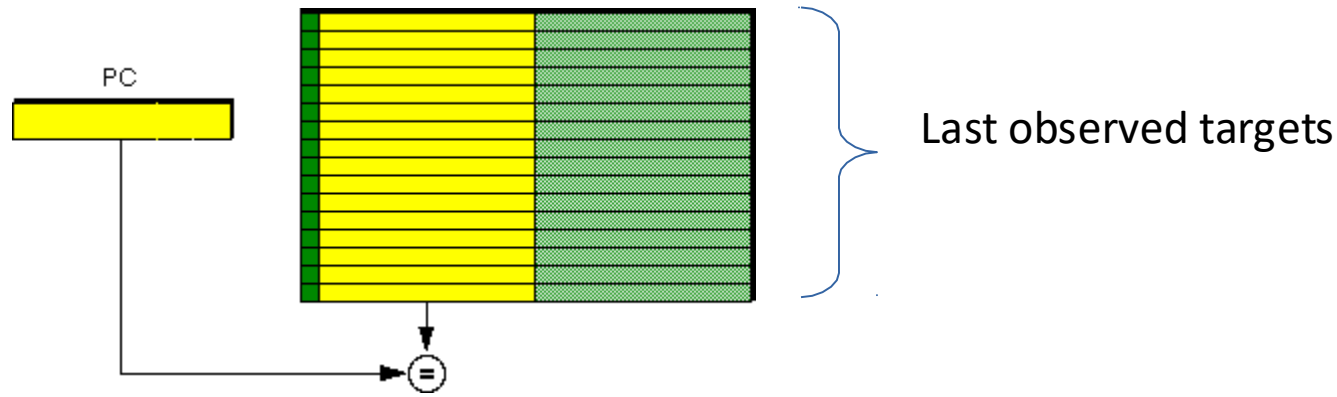
Coming to branches

Type	Direction at fetch time	Number of possible next fetch addresses?	When is next fetch address resolved?
Conditional	Unknown	2	Execution (register dependent)
Unconditional	Always taken	1	Decode (PC + offset)
Call	Always taken	1	Decode (PC + offset)
Return	Always taken	Many	Execution (register dependent)
Indirect	Always taken	Many	Execution (register dependent)

Different branch types can be handled differently

The very basic support for branches in a pipeline – Branch Target Buffer (BTB)

- It is a cache where an entry keeps the address of the branch instruction and the last target that has been observed for that branch
- Extremely useful for unconditional branches and common calls to subroutines



Coping with conditional branches

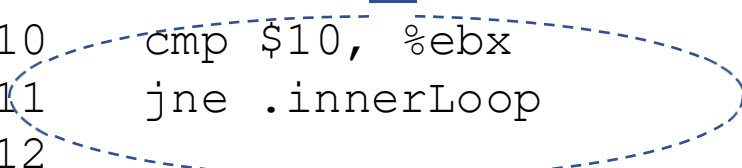
- The hardware support for improving performance under (speculative) pipelines in face of conditional branches is called *Dynamic Predictor*
- Its actual implementation consists of a Branch-Prediction Buffer (BPB) – or Branch History Table (BHT)
- The baseline implementation is based on a cache indexed by lower significant bits of branch instructions and one status bit
- The status bit tells whether the jump related to the branch instruction has been recently executed
- The (speculative) execution flow follows the direction related to the prediction by the status bit, thus following the recent behavior
- Recent past is expected to be representative of near future

Multiple-bit predictors

- One bit predictors “fail” in the scenario where the branch is often taken (or not taken) and infrequently not taken (or taken)
- In these scenarios, they leads to 2 subsequent errors in the prediction (thus 2 squashes of the pipeline)
- Is this really important? Nested loops tell yes
- The conclusion of the inner loop leads to change the prediction, which is anyhow re-changed at the next iteration of the outer loop
- Two-bit predictors require 2 subsequent prediction errors for inverting the prediction
- So each of the four states tells whether we are running with
 - ✓ YES prediction (with one or zero mistakes since the last passage on the branch)
 - ✓ NO prediction (with one or zero mistakes since the last passage on the branch)

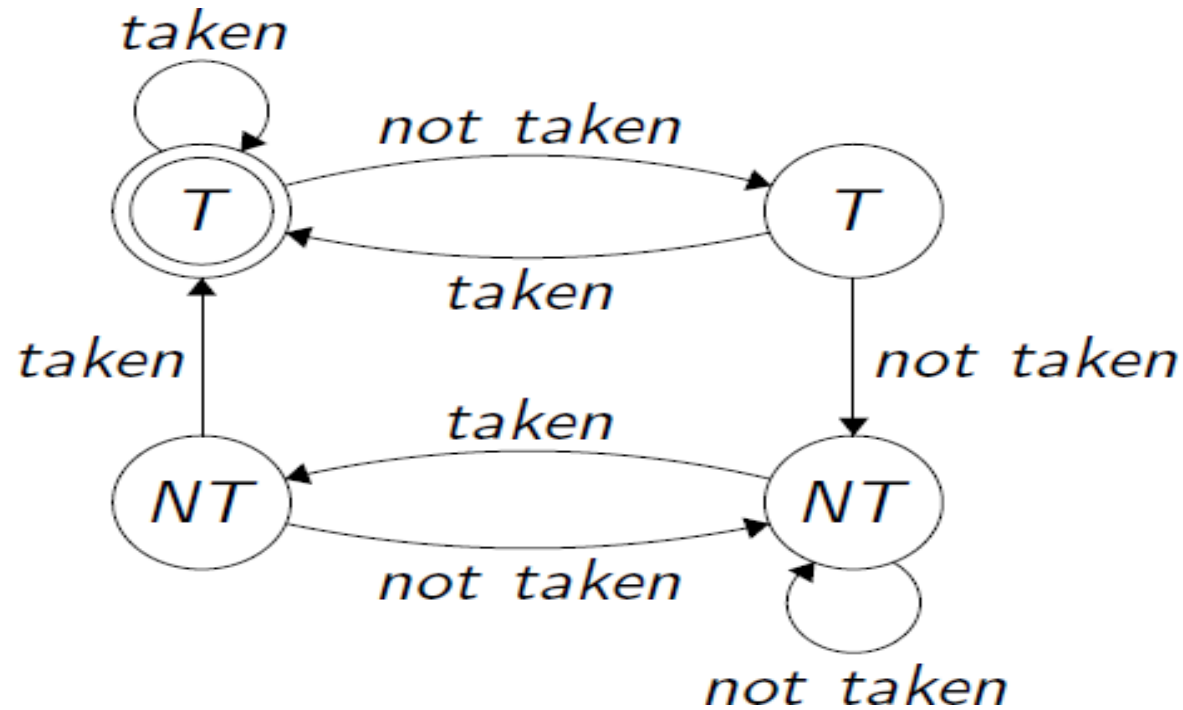
An example – AT&T syntax (we will no longer explicitly specify the actual syntax)

```
1      mov $0, %ecx
2      .outerLoop:
3      cmp $10, %ecx
4      je .done
5      mov $0, %ebx
6
7      .innerLoop:
8      ; actual code
9      inc %ebx
10     cmp $10, %ebx
11     jne .innerLoop
12
13     inc %ecx
14     jmp .outerLoop
15     .done:
```



This branch prediction is inverted at each ending inner-loop cycle

The actual two-bit predictor state machine



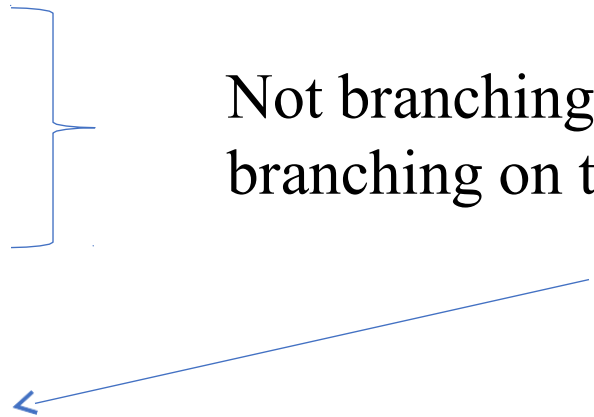
Do we need to go beyond two-bit predictors?

- Conditional branches are around 20% of the instructions in the code
- Pipelines are deeper
 - ✓ A greater misprediction penalty
- Superscalar architectures execute more instructions at once
 - ✓ The probability of finding a branch in the pipeline is higher
- The answer is clearly yes
- One more sophisticated approach offered by Pentium (and later) processors is Correlated Two-Level Prediction
- Another one offered by Alpha is Hybrid Local/Global predictor (also known as Tournament Predictor)

A motivating example

```
if (aa == VAL)
    aa = 0 ;
if (bb == VAL )
    bb = 0;
if (aa != bb) {
    //do the work
}
```

Not branching on these implies
branching on the subsequent



Idea of correlated prediction: lets' try to predict what will happen at the third branch by looking at the history of what happened in previous branches

The (m,n) two-level correlated predictor

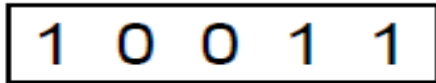
- The history of the last **m** branches is used to predict what will happen to the current branch
- The current branch is predicted with an **n**-bit predictor
- There are 2^m **n**-bit predictors
- The actual predictor for the current prediction is selected on the basis of the results of the last **m** branches, as coded into the 2^m bitmask
- A two-level correlated predictor of the form (0,2) boils down to a classical 2-bit predictor

(m,n) predictor architectural schematization

$$m = 5$$

$$n = 2$$

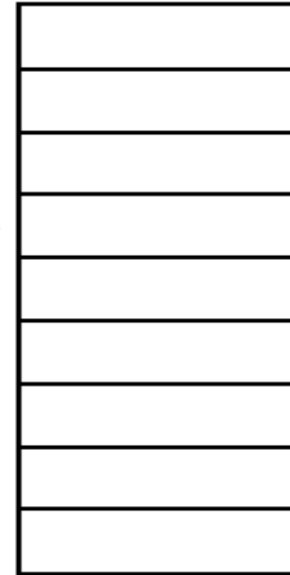
Global History Register
(shift register)



1: branch taken
0: branch not taken



Pattern History Table
(2^m entries of 2-bit counters)



This history can be speculative (hence retractable)

Tournament predictor

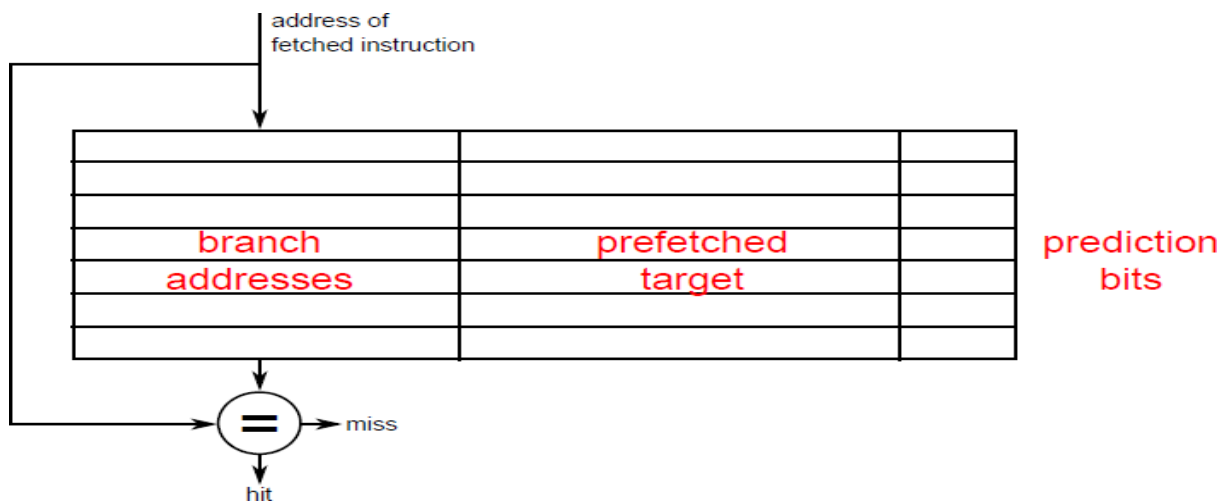
- The prediction of a branch is carried out by either using a local (per branch) predictor or a correlate (per history) predictor
- In the essence we have a combination of the two different prediction schemes
- Which of the two needs to be exploited at each individual prediction is encoded into a 4-states (2-bit based) history of success/failures
- This way, we can detect whether treating a branch as an individual in the prediction leads to improved effectiveness compared to treating it as an element in a sequence of individuals

Predicting “ret” instruction targets

- This is done via the Return Stack Buffer (RSB)
- It is simply a “stack layout cache” of addresses related to the return point of call instructions
- It is not so huge in terms of entries (typically we have the order of a few tens of entries – e.g. 32)
- Anyhow we may still experience a miss or we might (unlikely) incur a wrong prediction when return addresses into the stack area are changed by software

The very last concept on branch prediction - indirect branches

- These are branches for which the target is not known at instruction fetch time
- Essentially these are kind of families of branches (multi-target branches)
- x86 examples: `call rax` / `jmp eax`



Coming back to security

- A speculative processor can lead to micro-architectural effects by phantom instructions also in cases where the branch predictor fails, and the pipeline is eventually squashed
- If we can lead executing instructions in the phantom portion of the program flow to leave these micro-architectural effects then we can observe them via a side (covert) channel
- This is the baseline idea of Spectre attacks
- This have ben shown to kill Intel, AMD and ARM processors

Spectre primer

If (condition involving a target value X)
access array A at position
 $B[X] \ll 12$ //page size displacement

Suppose we run with
miss-prediction

The target line of A is cached
(as a side effect) and we can inspect this
via a side channel

Clearly B can be
whatever address,
so $B[X]$ is whatever
value

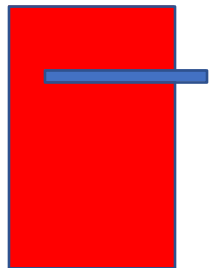
A scheme

If (condition in a target value X)
access array A at position
 $B[X] \ll 12$ //page size displacement

```
if (x < array1_size)  
    y = array2[array1[x] * 4096];
```

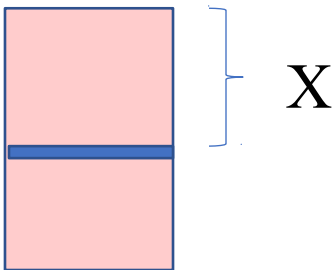
Actual code taken from
the original Spectre paper

A is cache
evicted



This brings one over 4096 bytes
(so the corresponding cache line)
into the cache and we can observe
this via a side channel

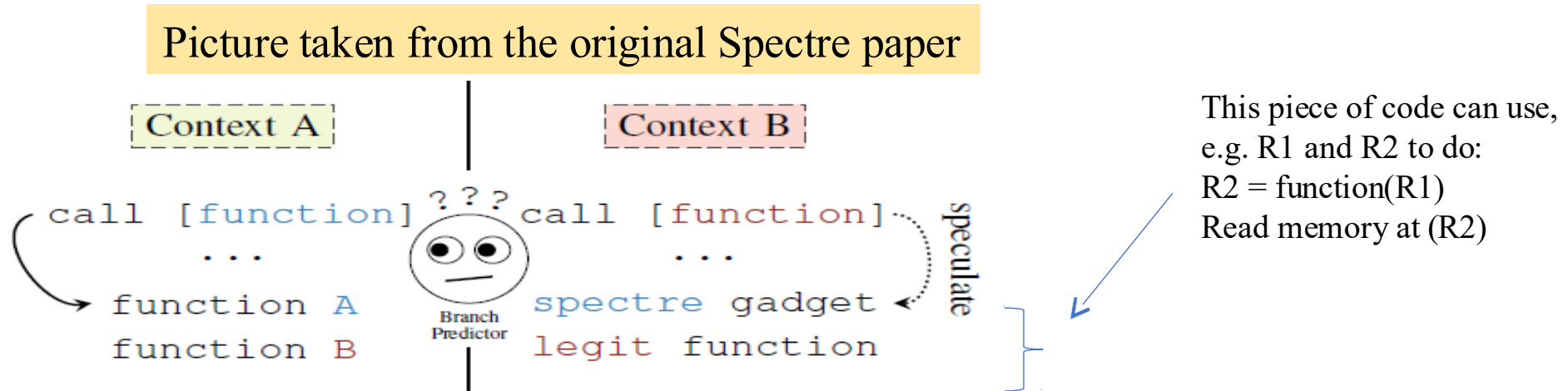
B is (partially)
a kernel zone



$B[X] \ll 12$ is used to read from A

Still Spectre - cross-context attacks

- Based on miss-predictions on indirect branches
 - ✓ train the predictor on an arbitrary address space and call-target
 - ✓ let the processor use a ‘gadget’ piece of code in, e.g. a shared library
 - ✓ somehow related to ROP (Return-Oriented –Programming), which we shall discuss



... using R1 alone in the attack

- The victim might have loaded a highly critical value into R1 (e.g. the results of a cryptographic operation)
- Then it might slide into the *call [function]* indirect branch
- The gadget might simply be a piece of code that accesses memory based on a function of R1
- **IMPORTANT NOTE:**
 - ✓ miss-training the indirect branch predictor needs to occur in the same CPU-core where the victim runs
 - ✓ while accessing the cache for latency evaluation and data leak actuation can take place on other CPU-cores, as we shall detailed see later while discussing the implementation of side/covert channels based on the caching system

Indirect branches countermeasures - Retpoline

- Retpoline (which stands for “return trampoline”) is used to perform an indirect jump using a `ret` instruction on x86 processors
- The steps in this idea are:
 - Save the target address for the jump `ADDR` onto the stack
 - Call a piece of code that removes the return PC of the call from the stack (this is the original call)
 - This piece of code then jumps to the target address via a return instruction
 - The original call has therefore no actual return hence the subsequent instructions to the call are a simple infinite loop
 - Indirect jumps branch predictor therefore will not be exploitable for tampering the control flow speculatively

Retpoline code (simplified)

```
    push target_address
1:    call retpoline_target
    //put here whatever you would like
    //typically a serializing instruction with no side effects
    jump 1b

retpoline_target:
    lea 8(%rsp), %rsp //we do not simply add 8 to RSP
                        //since FLAGS should not be modified
    ret //this will hit target_address
```

Indirect branch speculation barriers in the hardware

- Modern processors also offer a hardware based support for avoiding that the indirect branch predictor is affected by the history of previous software execution
- On x86 we have two main supports:
 - IBRS (Indirect Branch Restricted Speculation) – a MSR (Model Specific Register) can be updated in order to create an execution enclave where the history is not affected by things that are outside this enclave – useful when changing the protection mode of the processor
 - IBPB (Indirect Branch Prediction Barrier) – a MSR can be update in order to exclude indirect branch prediction affection at any point in a timeline

Checking with your Linux system

- You can type

```
grep CONFIG_RETPOLINE /boot/config-`uname -r`
```

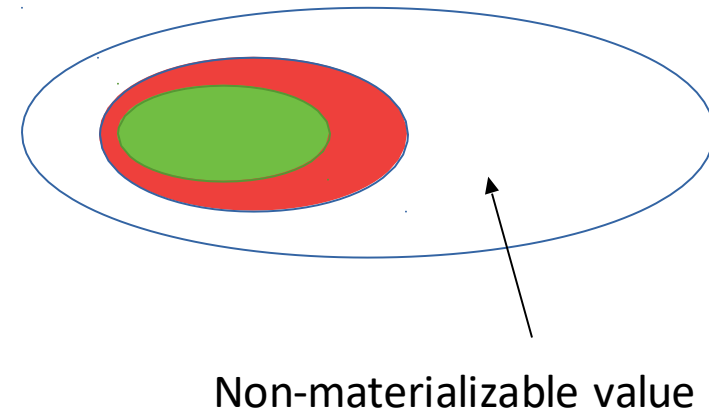
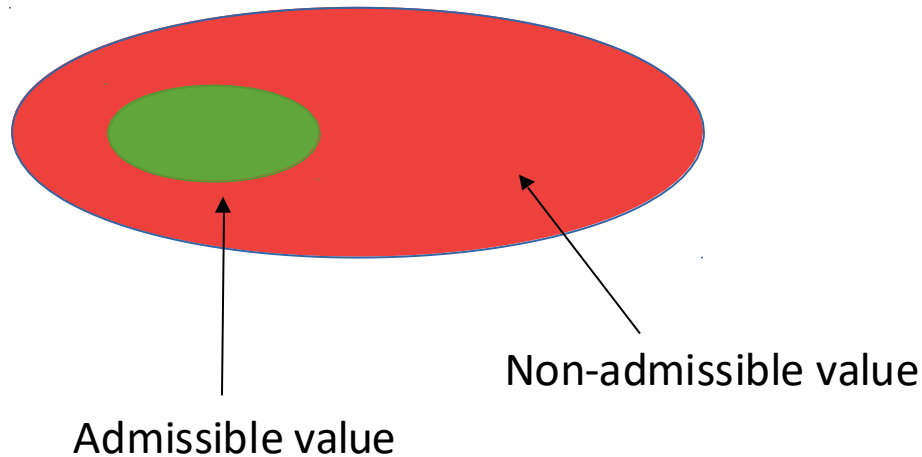
to check if your kernel has support for retpoline

- You can type

```
cat /sys/devices/system/cpu/vulnerabilities/spectre_v2
```

to determine what support your system has

Conditional branches countermeasures - sanitization



Extremely useful for table indexes

- You can type

```
cat /sys/devices/system/cpu/vulnerabilities/spectre_v1
```

to determine what support your system has

Loop unrolling

- This is a software technique that allows reducing the frequency of branches when running loops, and the relative cost of branch control instructions
- Essentially it is based on having the code-writer or the compiler to enlarge the cycle body by inserting multiple statements that would otherwise be executed in different loop iterations

```
int s=0;
```

```
for(int i=0;i<16;i++){s+=i;}
```

400545:	8b 45 fc	mov	-0x4(%rbp), %eax
400548:	01 45 f8	add	%eax, -0x8(%rbp)
40054b:	83 45 fc 01	addl	\$0x1, -0x4(%rbp)
40054f:	83 7d fc 0f	cmpl	\$0xf, -0x4(%rbp)
400553:	7e f0	jle	400545 <main+0x18>

gcc unroll directives

```
#pragma GCC push_options  
#pragma GCC optimize ("unroll-loops")
```

Region to unroll

```
#pragma GCC pop_options
```

- One may also specify the unroll factor via

```
#pragma unroll (N)
```

- In more recent *gcc* versions (e.g. 4 or later ones) it works with the `-O` directive

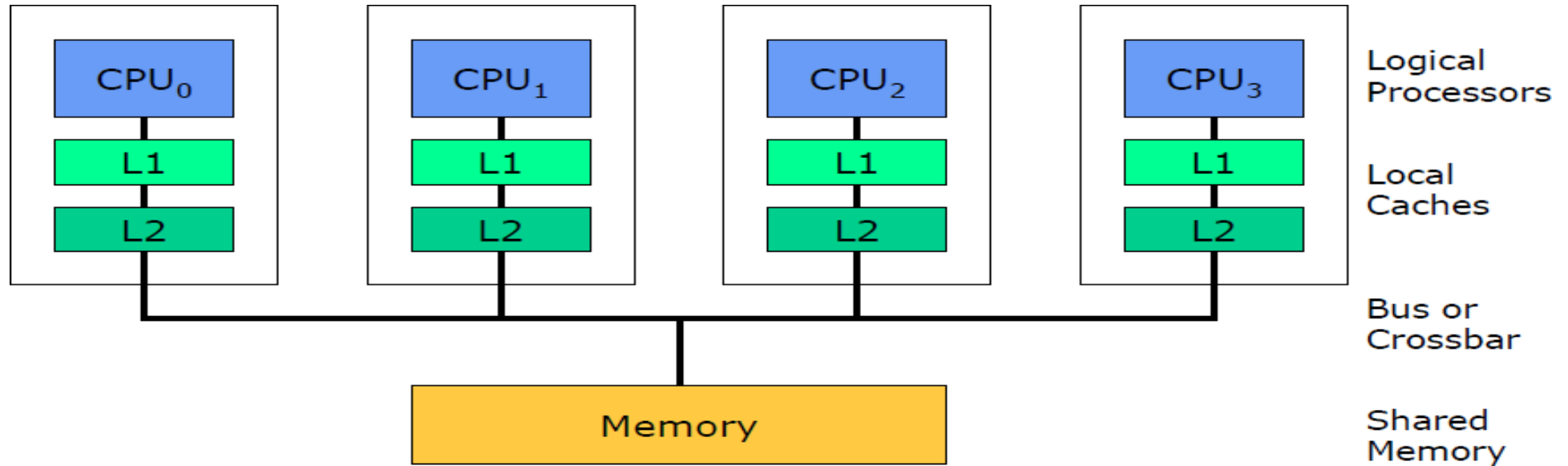
Beware unroll side effects

- It may put increased pressure on register usage leading to more frequent memory interactions
- When relying on huge unroll values code size can grow enormously, consequently locality and cache efficiency may degrade significantly
- Depending on the operations to be unrolled, it might be better to reduce the number of actual iterative steps via “vectorization”, a technique that we will look at later on

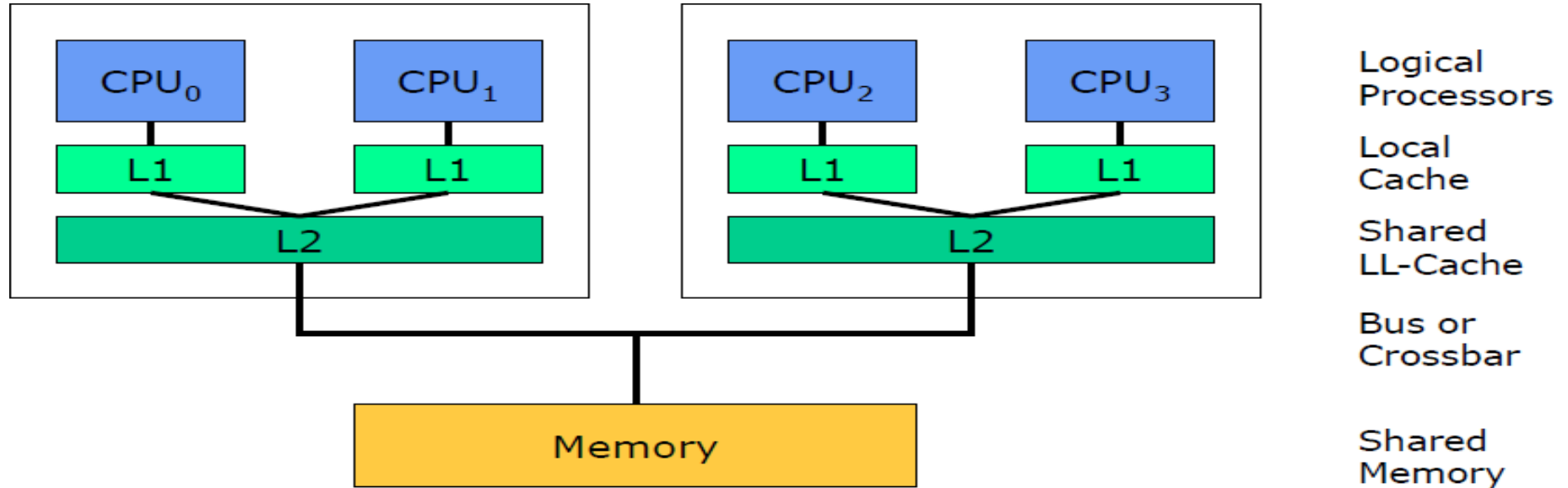
Clock frequency and power wall

- How can we make a processors run faster?
- Better exploitation of hardware components and growth of transistors' packaging – e.g. the More's law
- Increase of the clock frequency
- But nowadays we have been faced with the **power wall**, which actually prevents the building of processors with higher frequency
- In fact the power consumption grows exponentially with voltage according to the $V \times V \times F$ rule (and 130 W is considered the upper bound for dissipation)
- The way we have for continuously increasing the computing power of individual machines is to rely on parallel processing units

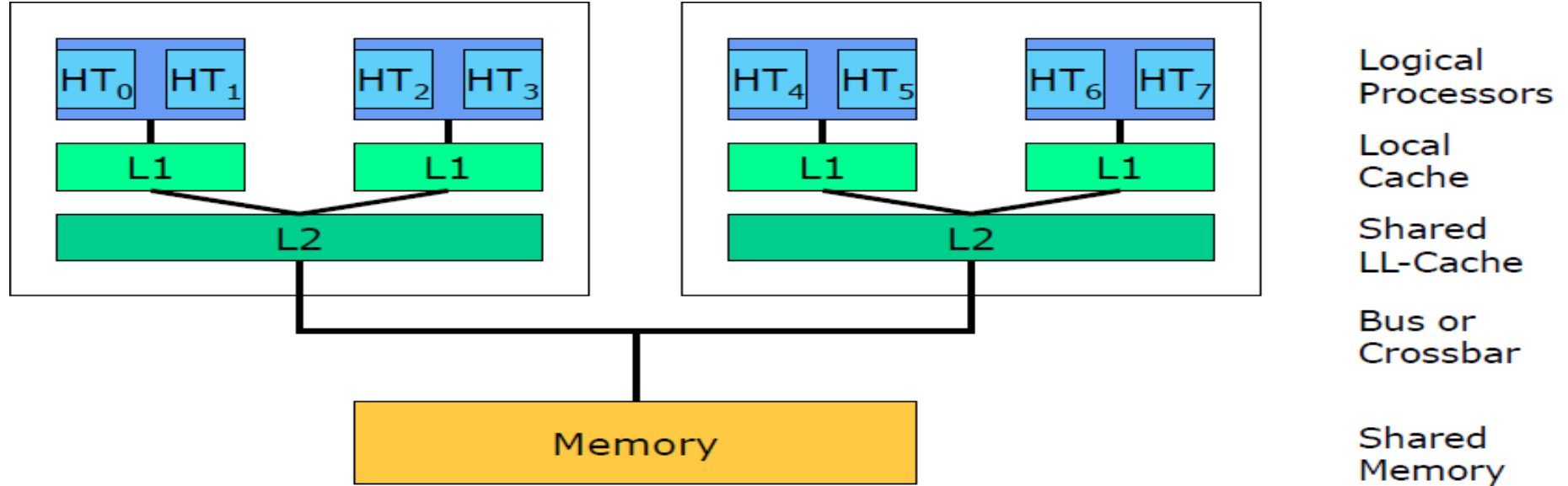
Symmetric multiprocessors



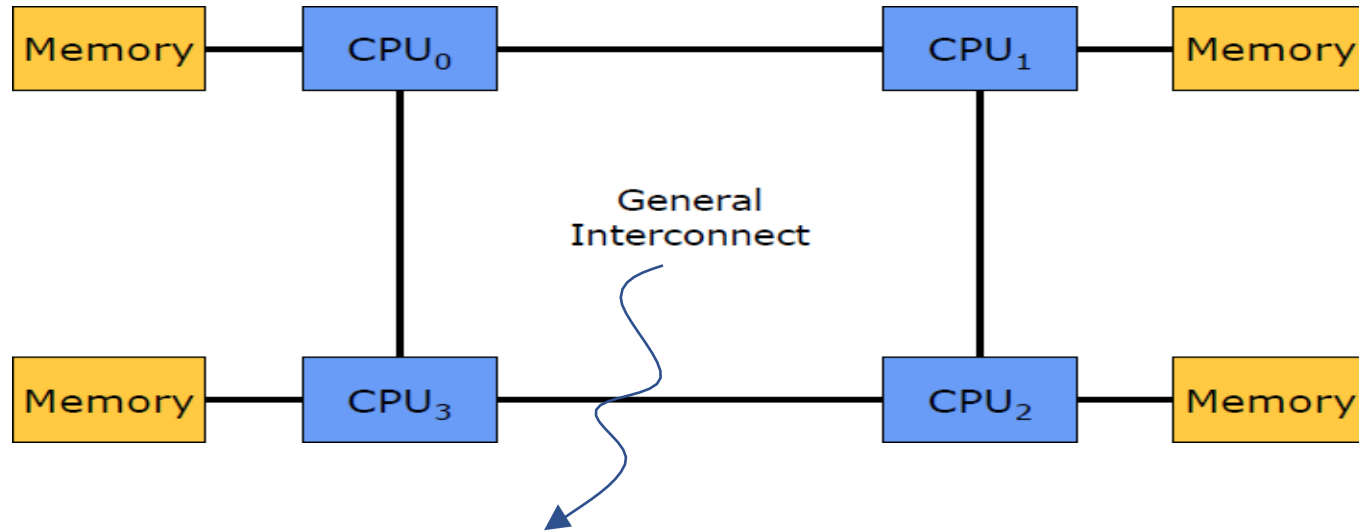
Chip Multi Processor (CMP) - Multicore



Symmetric Multi-threading (SMT) - Hyperthreading



Making memory circuitry scalable – NUMA (Non Uniform memory Access)



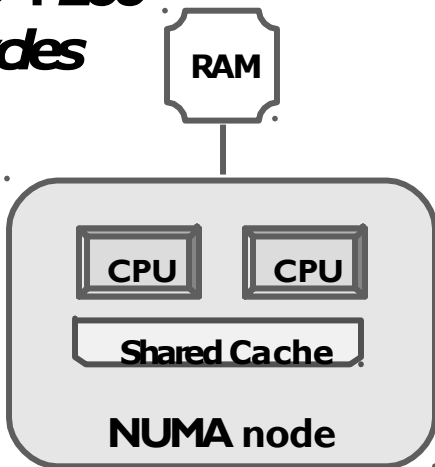
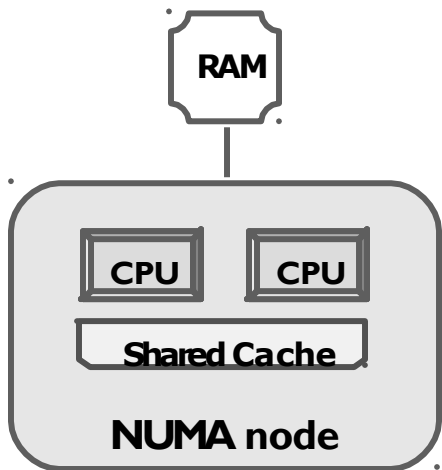
This may have different shapes
depending on chipsets

NUMA latency asymmetries

Local accesses are served by

- Inner private/shared caches
- Inner memory controllers

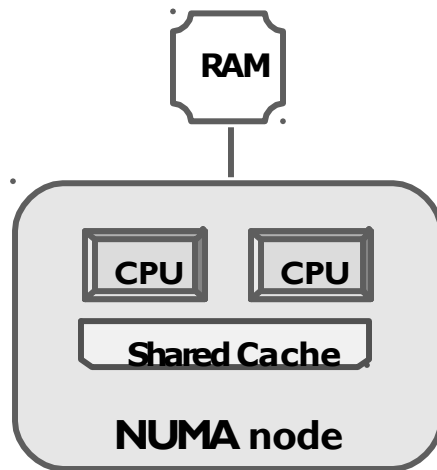
**$50 \div 200$
cycles**



Remote accesses are served by

- Outer shared caches
- Outer memory controllers

**$200 \div 300$ cycles
($1x \div 6x$)**



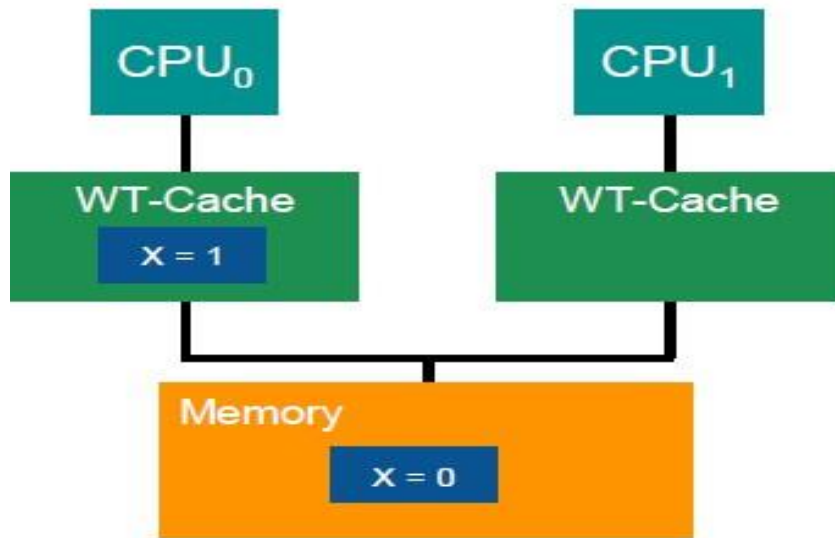
Cache coherency

- CPU-cores see memory contents through their caching hierarchy
- This is essentially a **replication system**
- The problem of defining what value (within the replication scheme) should be returned upon reading from memory is also referred to as “cache coherency”
- This is definitely different from the problem of defining when written values by programs can be actually read from memory
- The latter is in fact known to as the “memory consistency” problem, which we will discuss later on
- Overall, cache coherency is not memory consistency, but it is anyhow a big challenge to cope with, with clear implications on performance

Defining coherency

- A read from location X, previously written by a processor, returns the last written value if no other processor carried out writes on X in the meanwhile – **Causal consistency along program order**
- A read from location X by a processor, which follows a write on X by some other processor, returns the written value if the two operations are sufficiently separated along time (and no other processor writes X in the meanwhile) – **Avoidance of staleness**
- All writes on X from all processors are serialized, so that the writes are seen from all processors in a same order – **We cannot (ephemerally or permanently) invert memory updates**
- however we will come back to defining when a processor actually writes to memory!!
- Please take care that coherency deals with individual memory location operations!!!

Write through cache (and its coherency issues)

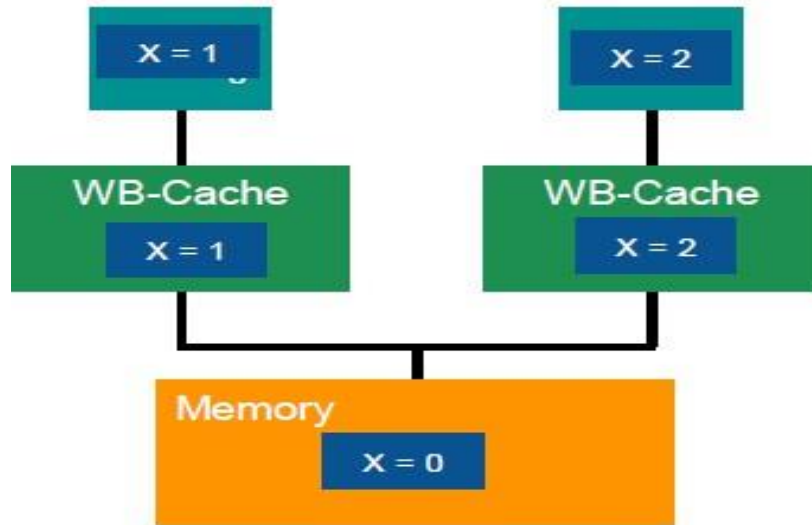


1. CPU₀ reads X from memory
 - loads X=0 into its cache
 2. CPU₁ reads X from memory
 - loads X=0 into its cache
 3. CPU₀ writes X=1
 - stores X=1 in its cache
 - stores X=1 in memory
 4. CPU₁ reads X from its cache
 - loads X=0 from its cache

Incoherent value for X on CPU₁
- CPU₁ may wait for update!

Requires write propagation!

Write back cache (and its coherency issues)



1. CPU₀ reads X from memory
 - loads X=0 into its cache
 2. CPU₁ reads X from memory
 - loads X=0 into its cache
 3. CPU₀ writes X=1
 - stores X=1 in its cache
 4. CPU₁ writes X=2
 - stores X=2 in its cache
 5. CPU₁ writes back cache line
 - stores X=2 in memory
 6. CPU₀ writes back cache line
 - stores X=1 in memory
- Later (!) store X=2 from CPU₁ lost

Requires write serialization!

Cache coherency (CC) protocols - basics

- A CC protocol is the result of choosing
 - ✓ a set of transactions supported by the distributed cache system
 - ✓ a set of states for cache blocks
 - ✓ a set of events handled by controllers
 - ✓ a set of transitions between states
- Their design is affected by several factors, such as
 - ✓ interconnection topology (e.g., single bus, hierarchical, ring-based)
 - ✓ communication primitives (i.e., unicast, multicast, broadcast)
 - ✓ memory hierarchy features (e.g., depth, inclusiveness)
 - ✓ cache policies (e.g., write-back vs write-through)
- Different CC implementations have different performance
 - ✓ Latency: time to complete a single transaction
 - ✓ Throughput: number of completed transactions per unit of time
 - ✓ Space overhead: number of bits required to maintain a block state

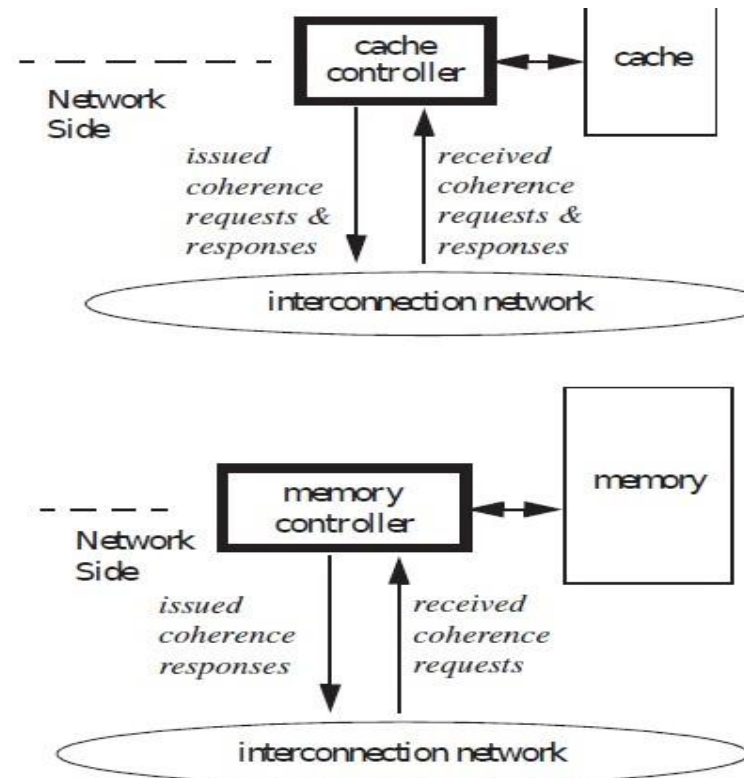
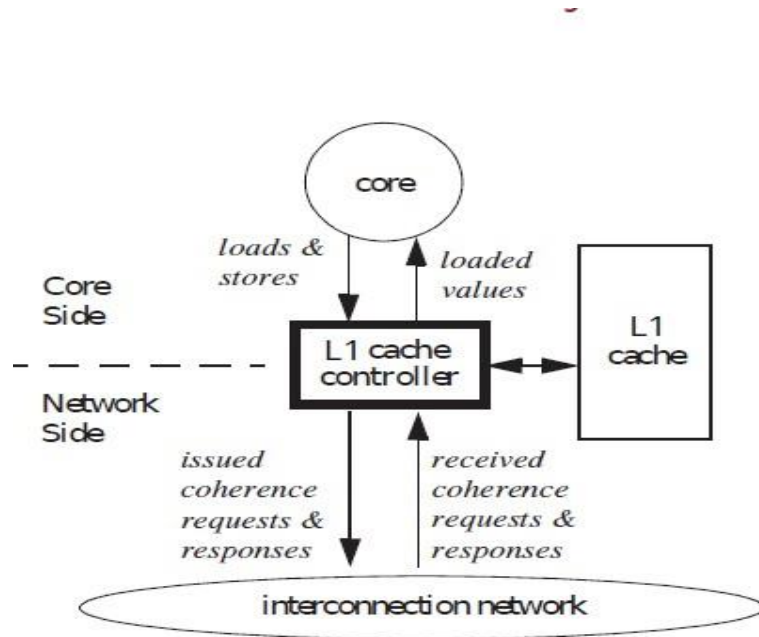
Families of CC protocols

- When to update copies in other caches?
- Invalidate protocols:
 - ✓ When a core writes to a block, all other copies are invalidated
 - ✓ Only the writer has an up-to-date copy of the block
 - ✓ Trades latency for bandwidth
- Update protocols:
 - ✓ When a core writes to a block, it updates all other copies
 - ✓ All cores have an up-to-date copy of the block
 - ✓ Trades bandwidth for latency

“Snooping cache” coherency protocols

- At the architectural level, these are based on some broadcast medium (also called network) across all cache/memory components
- Each cache/memory component is connected to the broadcast medium by relying on a controller, which snoops (observes) the in-flight data
- The broadcast medium is used to issue “transactions” on the state of cache blocks
- Agreement on state changes comes out by serializing the transactions traveling along the broadcast medium
- A state transition cannot occur unless the broadcast medium is acquired by the source controller
- State transitions are distributed (across the components), but are carried out atomically thanks to serialization over the broadcast medium

An architectural scheme



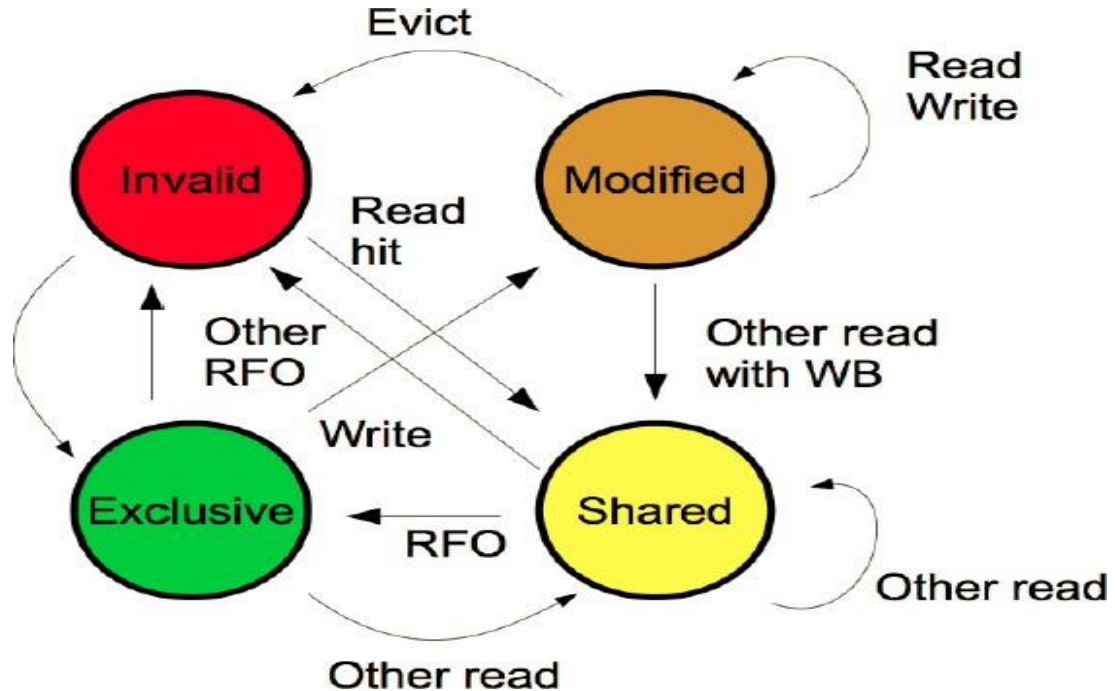
Write/read transactions with invalidation

- A write transaction invalidates all the other copies of the cache block
- Read transactions
 - ✓ Get the latest updated copy from memory in write-through caches
 - ✓ Get the latest updated copy from memory or from another caching component in write-back caches (e.g. Intel processors)
- We typically keep track of whether
 - ✓ A block is in the modified state (just written, hence invalidating all the other copies)
 - ✓ A block is in shared state (someone got the copy from the writer or from another reader)
 - ✓ A block is in the invalid state
- This is the MSI (Modified-Shared-Invalid) protocol

Reducing invalidation traffic upon writes - MESI

- Similar to MSI but includes an “exclusive” state indicating that a unique valid copy is owned, independently of whether the block has been written or not

RFO = Request
For
Ownership



MESI recap

- **Modified:** You have modified shared data
- **Exclusive:** You are the sole owner of this data and are free to modify it without a bus message
- **Shared:** You have a copy of data that another processor also has
- **Invalid:** Your copy of the data is not up to date

Transitions in response to “local” reads

State is M

- No bus transaction

State is E

- No bus transaction

State is S

- No bus transaction

State is I

- Generate bus read request (BusRd) *May force other cache operations (see later)*
- Other cache(s) signal “sharing” if they hold a copy
- If shared was signaled, go to state S
- Otherwise, go to state E

After update: return read value

Transitions in response to “local” writes

State is M

- No bus transaction

State is E

- No bus transaction
- Go to state M

State is S

- Line already local & clean
- There may be other copies
- Generate bus read request for upgrade to exclusive (BusRdX*)
- Go to state M

State is I

- Generate bus read request for exclusive ownership (BusRdX)
- Go to state M

Transitions in response to snooped bus read

State is M

- Write cache line back to main memory
- Signal “shared”
- Go to state S

State is E

- Go to state S and signal “shared”

State is S

- Signal “shared”

State is I

- Ignore

Transitions in response to snooped bus read “exclusive”

State is M

- Write cache line back to memory
- Discard line and go to I

State is E

- Discard line and go to I

State is S

- Discard line and go to I

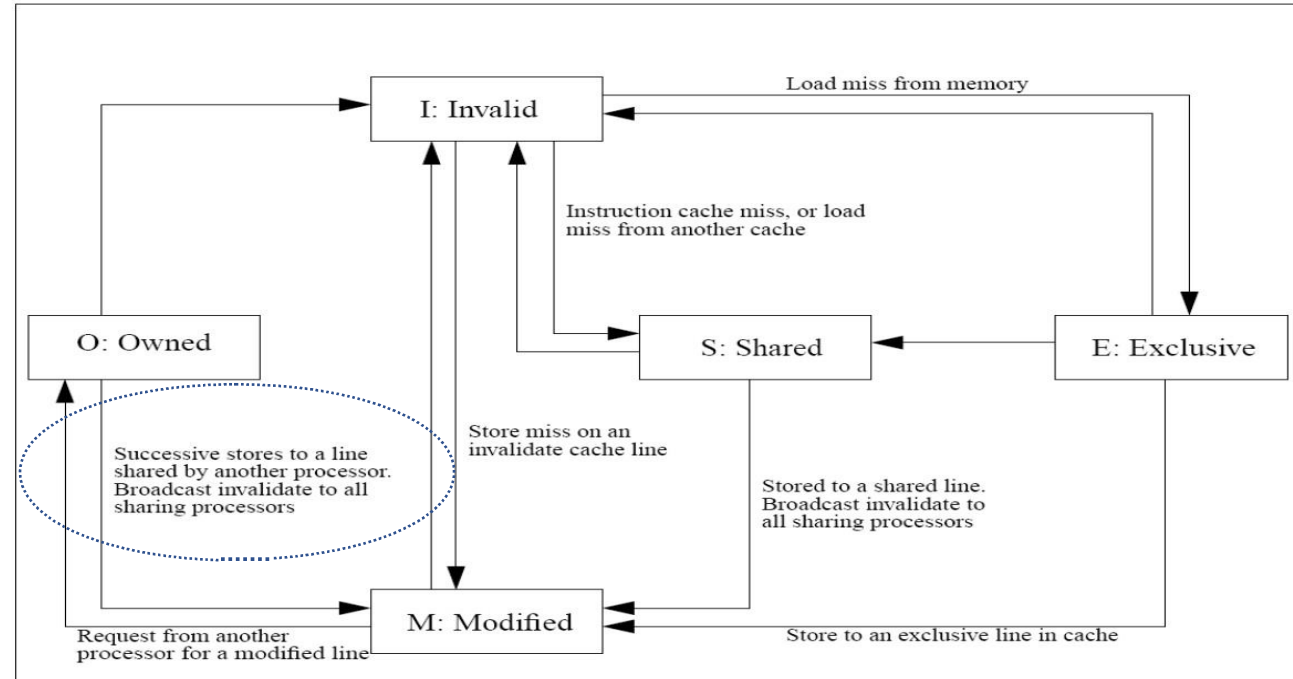
State is I

- Ignore

Further reducing bus traffic upon memory operations - MOESI

- similar to MESI but includes an “owner” state indicating that a unique owner of the master copy exists, among the updated ones that are shared

No need to pass through “exclusive” again



MOESI recap

- **Modified:** You have modified shared data
- **Owner:** Your data is shared, but you have the master copy in the cache, and can modify this data as you wish (without additional state transitions)
- **Exclusive:** You are the sole owner of this data and are free to modify it without a bus message
- **Shared:** You have a copy of data that another processor also has
- **Invalid:** Your copy of the data is not up to date

x86 Implementations

- **Intel**
 - Mostly MESI
 - Inclusive cache
 - Write back
 - L1 cache line 64 bytes
- **AMD**
 - Mostly MOESI
 - Exclusive cache (at L3)
 - Write back
 - L1 cache line 64 byte

An alternative to snooping - directory based protocols

Snooping does not scale

- Bus transactions must be *globally* visible
- Implies broadcast

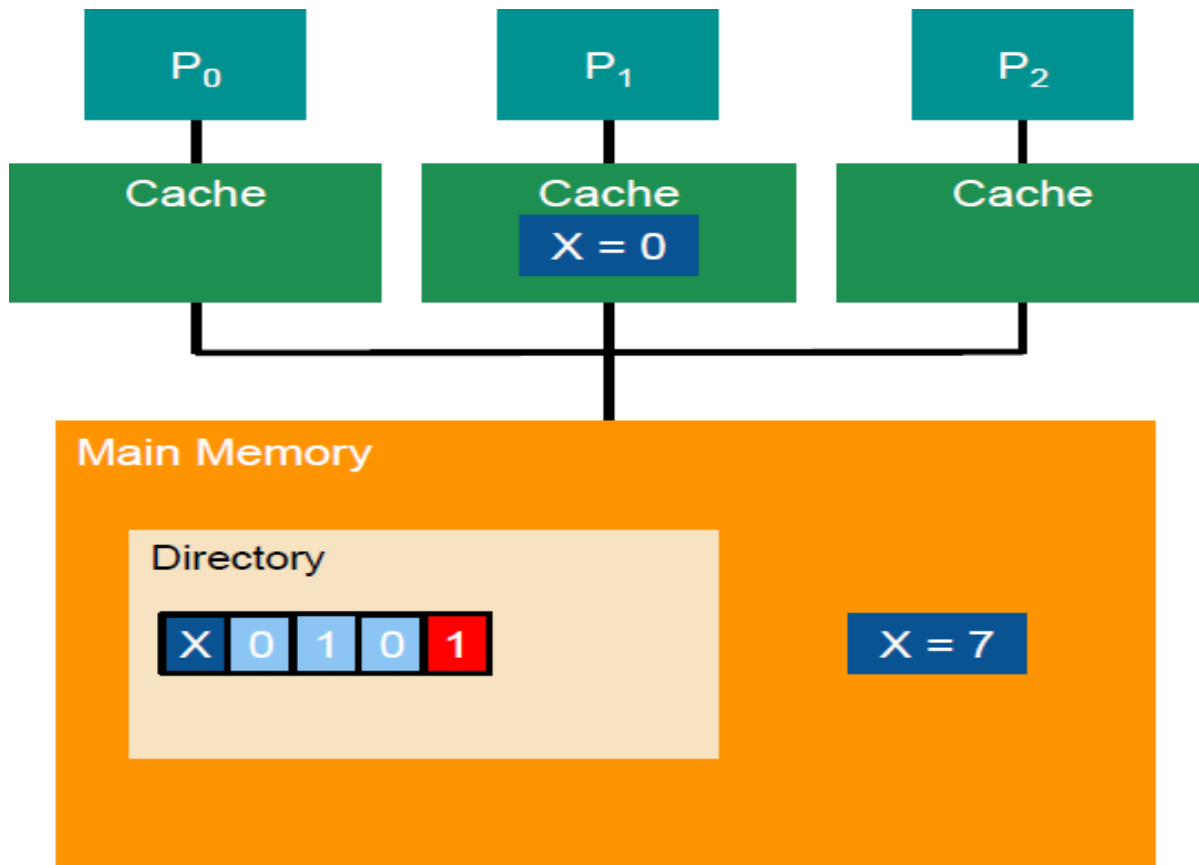
Directory-based schemes are more scalable

- Directory (entry for each CL) keeps track of all owning caches
- Point-to-point update to involved processors
 - ✓ *No broadcast*
 - ✓ *Can use specialized (high-bandwidth) network*

Basic scheme

System with N processors P_i

- For each memory block (size: cache line) maintain a directory entry
- N presence bits (light blue)
Set if block in cache of P_i
- 1 dirty bit (red)

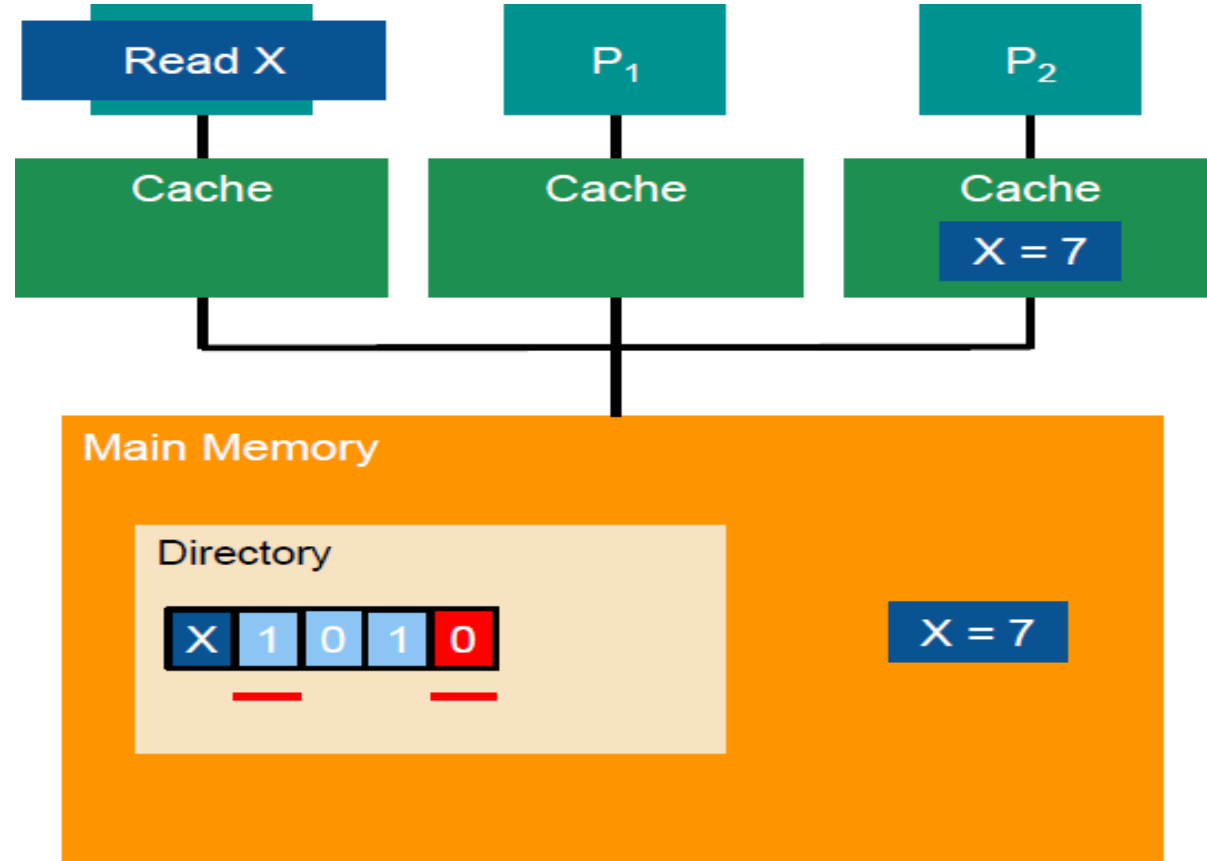


Directory based CC - read miss

P0 intends to read, misses

If dirty bit (in directory) is off

- Read from main memory
- Set presence[i]
- Supply data to reader

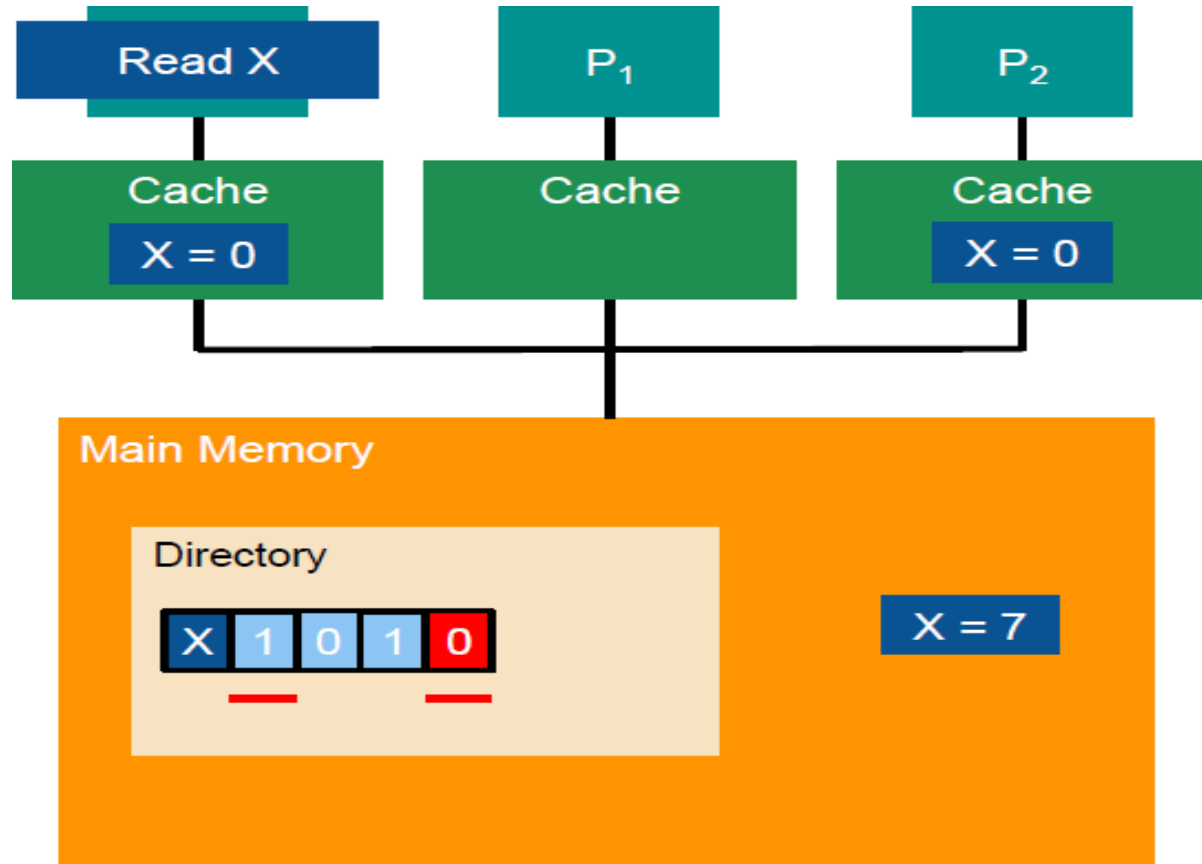


Directory based CC - read miss

P0 intends to read, misses

If dirty bit is on

- Recall cache line from P_j
(determine by presence[])
- Unset dirty bit, block shared
- Set presence[i]
- Supply data to reader and propagate to memory

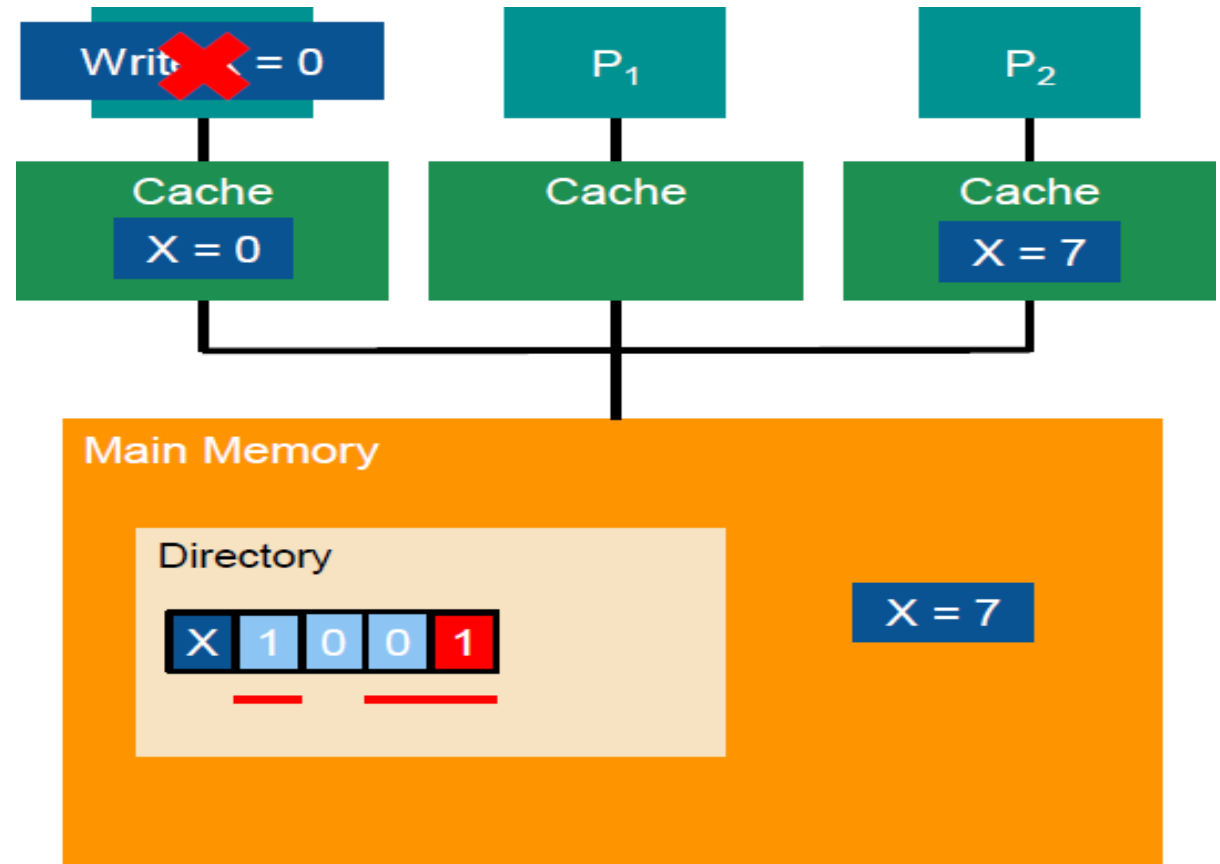


Directory based CC - write miss

P0 intends to write, misses

If dirty bit (in directory) is off

- Send invalidations to all processors P_j with presence[j] turned on
- Unset presence bit for all processors
- Set dirty bit
- Set presence[i], owner P_i

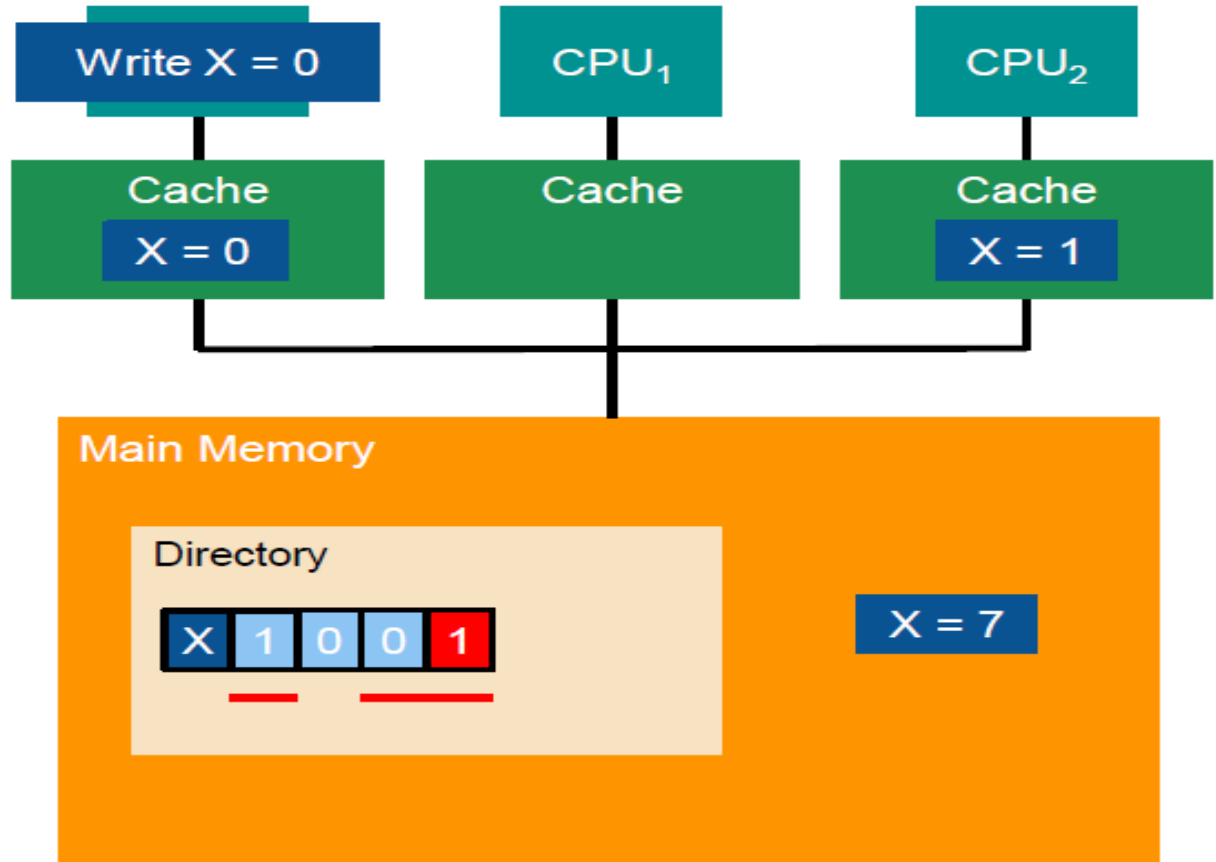


Directory based CC - write miss

P0 intends to write, misses

If dirty bit is on

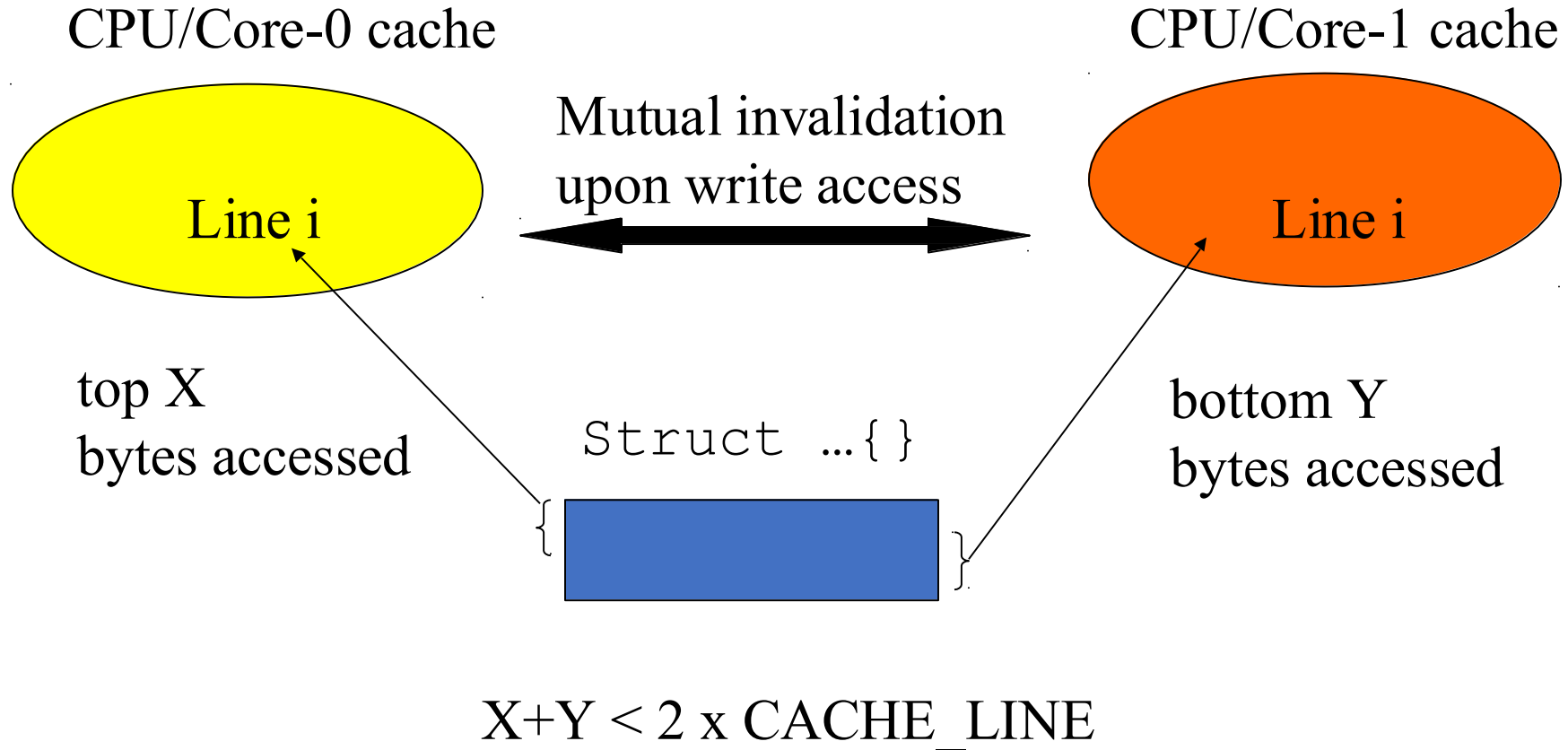
- Recall cache line from owner Pj
- Unset presence[j]
- Set presence[i], dirty bit remains set
- Acknowledge to writer



Software exposed cache performance aspects

- “Common fields” access issues
 - ✓ Most used fields inside a data structure should be placed at the head of the structure in order to maximize cache hits
 - ✓ This should happen provided that the memory allocator gives cache-line aligned addresses for dynamically allocated memory chunks
- “Loosely related fields” should be placed sufficiently distant inside the data structure so to avoid performance penalties due to false cache sharing

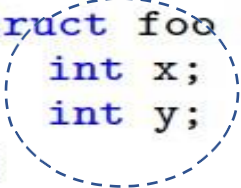
The false cache sharing problem



Example code leading to false cache sharing


Fits into a same cache line
(typically 64/256 bytes)

```
1 struct foo {  
2     int x;  
3     int y;  
4 };  
5  
6 struct foo f;
```



```
1 void inc_x(void)  
2 {  
3     int i;  
4     for(i = 0; i <  
5         1000000; i++)  
6         f.x++;  
7 }
```

```
1 int sum_y(void) {  
2     int s = 0;  
3     int i;  
4     for (i = 0; i <  
5         1000000; i++)  
6         s += f.y;  
7     return s;  
8 }
```



These reads from the cache
line find cache-invalid data, even though
the actual memory location we are reading from
does not change over time

Posix memory-aligned allocation

POSIX_MIMALIGN(3)

Linux Programmer's Manual

POSIX_MIMALIGN(3)

NAME [top](#)

`posix_memalign`, `aligned_alloc`, `memalign`, `valloc`, `pvalloc` - allocate aligned memory

SYNOPSIS [top](#)

```
#include <stdlib.h>
```

```
int posix_memalign(void **memptr, size_t alignment, size_t size);  
void *aligned_alloc(size_t alignment, size_t size);  
void *valloc(size_t size);
```

```
#include <malloc.h>
```

```
void *memalign(size_t alignment, size_t size);  
void *pvalloc(size_t size);
```

Feature Test Macro Requirements for glibc (see [feature_test_macros\(7\)](#)):

Inspecting cache line accesses

- A technique (called Flush+Reload) presented at [USENIX Security Symposium – 2013] is based on observing access latencies on shared data
- Algorithmic steps:
 - ✓ The cache content related to some shared data is flushed
 - ✓ Successively it is re-accessed in read mode
 - ✓ Depending on the timing of the latter accesses we gather whether the datum has been also accessed by some other thread
- Implementation on x86 is based on 2 building blocks:
 - ✓ A high resolution timer
 - ✓ A non-privileged cache line flush instruction
- These algorithmic steps have been finally exploited for Meltdown/Spectre attacks
- ... let's see the details

x86 high resolution timers

RDTSC

Read Time-Stamp Counter

Opcode	Mnemonic	Description
0F 31	RDTSC	Read time-stamp counter into EDX:EAX.

Description

Loads the current value of the processor's time-stamp counter into the EDX:EAX registers. The time-stamp counter is contained in a 64-bit MSR. The high-order 32 bits of the MSR are loaded into the EDX register, and the low-order 32 bits are loaded into the EAX register. The processor monotonically increments the time-stamp counter MSR every clock cycle and resets it to 0 whenever the processor is reset. See "Time Stamp Counter" in Chapter 15 of the IA-32 Intel Architecture Software Developer's Manual, Volume 3 for specific details of the time stamp counter behavior.

When in protected or virtual 8086 mode, the time stamp disable (TSD) flag in register CR4 restricts the use of the RDTSC instruction as follows. When the TSD flag is clear, the RDTSC instruction can be executed at any privilege level; when the flag is set, the instruction can only be executed at privilege level 0. (When in real-address mode, the RDTSC instruction is always enabled.) The time-stamp counter can also be read with the RDMSR instruction, when executing at privilege level 0.

The RDTSC instruction is not a serializing instruction. Thus, it does not necessarily wait until all previous instructions have been executed before reading the counter. Similarly, subsequent instructions may begin execution before the read operation is performed.

This instruction was introduced into the IA-32 Architecture in the Pentium processor.

Operation

```
if(CR4.TSD == 0 || CPL == 0 || CR0.PE == 0) EDX:EAX = TimeStampCounter;  
else Exception(GP(0)); //CR4.TSD is 1 and CPL is 1, 2, or 3 and CR0.PE is 1
```

Flags affected

None.

x86 (non privileged) cache line flush

x86 Instruction Set Reference

CLFLUSH

Flush Cache Line

Opcode	Mnemonic	Description
0F AE /7	CLFLUSH m8	Flushes cache line containing m8.

Description
<p>Invalidates the cache line that contains the linear address specified with the source operand from all levels of the processor cache hierarchy (data and instruction). The invalidation is broadcast throughout the cache coherence domain. If, at any level of the cache hierarchy, the line is inconsistent with memory (dirty) it is written to memory before invalidation. The source operand is a byte memory location.</p> <p>The availability of CLFLUSH is indicated by the presence of the CPUID feature flag CLFSH (bit 19 of the EDX register, see Section , CPUID-CPU Identification). The aligned cache line size affected is also indicated with the CPUID instruction (bits 8 through 15 of the EBX register when the initial value in the EAX register is 1).</p> <p>The memory attribute of the page containing the affected line has no effect on the behavior of this instruction. It should be noted that processors are free to speculatively fetch and cache data from system memory regions assigned a memory-type allowing for speculative reads (such as, the WB, WC, and WT memory types). PREFETCHH instructions can be used to provide the processor with hints for this speculative behavior. Because this speculative fetching can occur at any time and is not tied to instruction execution, the CLFLUSH instruction is not ordered with respect to PREFETCHH instructions or any of the speculative fetching mechanisms (that is, <u>data can be speculatively loaded into a cache line just before, during, or after the execution of a CLFLUSH instruction that references the cache line</u>).</p> <p>CLFLUSH is only ordered by the MFENCE instruction. It is not guaranteed to be ordered by any other fencing or serializing instructions or by another CLFLUSH instruction. For example, software can use an MFENCE instruction to insure that previous stores are included in the writeback.</p> <p><u>The CLFLUSH instruction can be used at all privilege levels and is subject to all permission checking and faults associated with a byte load (and in addition, a CLFLUSH instruction is allowed to flush a linear address in an execute-only segment). Like a load, the CLFLUSH instruction sets the A bit but not the D bit in the page tables.</u></p> <p>The CLFLUSH instruction was introduced with the SSE2 extensions; however, because it has its own CPUID feature flag, it can be implemented in IA-32 processors that do not include the SSE2 extensions. Also, detecting the presence of the SSE2 extensions with the CPUID instruction does not guarantee that the CLFLUSH instruction is implemented in the processor.</p>

ASM inline

- Exploited to define ASM instructions to be posted into a C function
- The programmer does not leave freedom to the compiler on that instruction sequence
- Easy way of linking ASM and C notations
- Structure of an ASM inline block of code for *gcc*

```
__asm__ [volatile][goto] (AssemblerTemplate  
    [ : OutputOperands ]  
    [ : InputOperands ]  
    [ : Clobbers ]  
    [ : GotoLabels ] );
```


Meaning of ASM inline fields

- `AssemblerTemplate` - the actual ASM code
- `volatile` – forces the compiler not to take any optimization (e.g. instruction placement effect)
- `goto` – assembly can lead to jump to any label in `GoToLabels`
- `OutputOperands` – data move post conditions
- `InputOperands` – data move preconditions
- `Clobbers` – registers involved in update by the ASM code, which require save/restore of their values (e.g. callee save registers)

C compilation directives for operands

- The = symbol means that the corresponding operand is used as an **output**
- Hence after the execution of the ASM code block, the operand value becomes the source for a given target location (e.g. for a variable)
- In case the operand needs to keep a value to be used as an **input** (hence the operand is the storage location of the value of some source location) then the = symbol does not need to be used

Main gcc supported operand specifications

r – generic register operands

m – generic memory operand (e.g. into the stack)

0-9 – reused operand index

i/I – immediate 64/32 bit operand

q - byte-addressable register (e.g. eax, ebx, ecx, edx)

A - eax or edx

a, b, c, d, S, D - eax, ebx, ecx, edx, esi, edi respectively (or al, rax etc variants depending on the size of the actual-instruction operands)

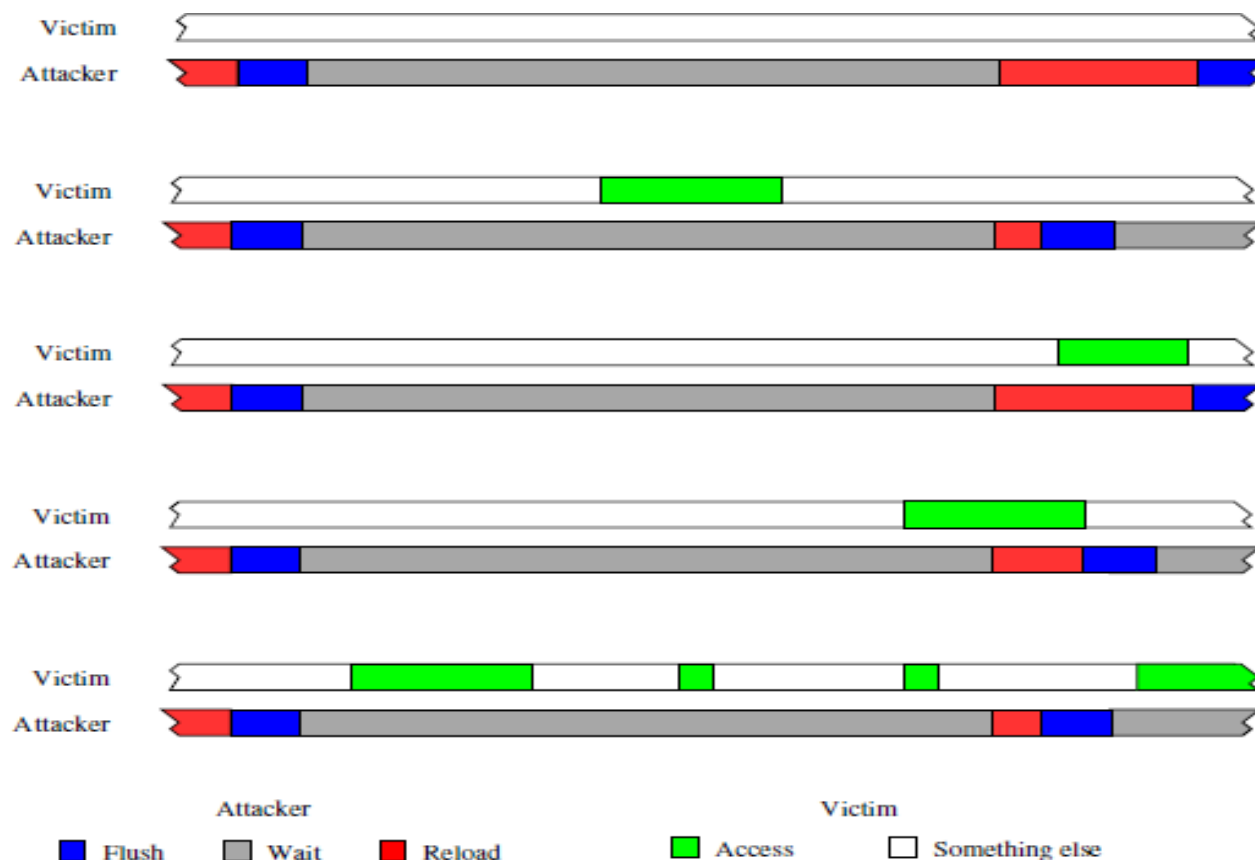
Flush+Reload - measuring cache access latency at user space

```
unsigned long probe(char* adrs){  
    volatile unsigned long cycles;  
    asm(  
        "mfence \n"  
        "lfence \n"  
        "rdtsc \n"  
        "lfence \n"  
        "movl %%eax, %%esi \n"  
        "movl (%1), %%eax \n"  
        "lfence \n"  
        "rdtsc \n"  
        "subl %%esi, %%eax \n"  
        "clflush 0(%1) \n"  
        : "=a" (cycles)  
        : "c" (adrs)  
        : "%esi" , "%edx" );  
  
    return cycles;  
}
```

A barrier on all memory accesses

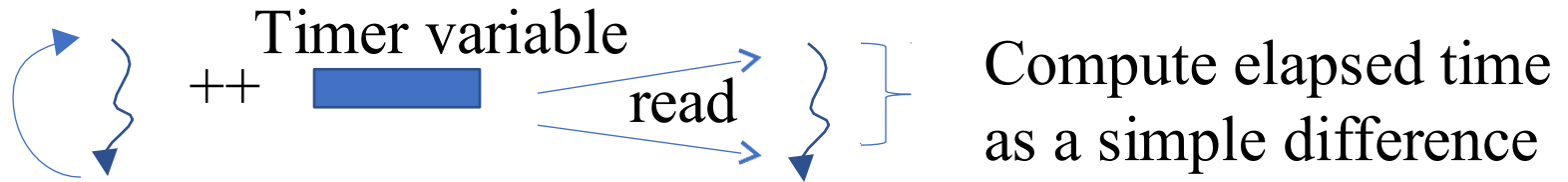
Barriers on loads

Typical Flush+Reload timelines



Can we prevent *rtdsc* to be user-space accessed?

- Yes, the processor state can be configured to run this instruction in privileged mode
- This will require passing through a system-call to get the value, which might introduce “variance” in the samples, leading to less precise measurement, a topic we will re-discuss
- In any case we can emulate a timer with a user space thread regularly incrementing a shared counter



The two threads should work on close (although different) cores

Flush+Reload vs cache inclusiveness

- An inclusive caching system (e.g. Intel) is such that a lower level caching component L_x always keeps a copy of the content cached by some upper level component L_y
- So L_y content is always included in L_x one
- A non-inclusive cache (e.g. some AMD chipsets) does not have this vincula
- So we might cache some data at $L1$ and then, e.g. after evicting from $L1$, we may load the data in, e.g. LLC
- For these caching systems, cross process Flush+Reload attacks may fail
- They can instead still be fruitfully when used, e.g., across processes running on a same CPU-core

The actual meaning of reading/writing from/to memory

- What is the memory behavior under concurrent data accesses?
 - ✓ Reading a memory location should return last value written
 - ✓ **The last value written** not clearly (or univocally) defined under concurrent access and with multi-locations as the target
- The memory consistency model
 - ✓ Defines in which order processing units perceive concurrent accesses
 - ✓ Based on ordering rules, not necessarily timing of accesses
- Memory consistency is not memory (i.e. cache) coherency!!!

Terminology for memory models

- Program Order (of a processor's operations)
 - ✓ per-processor order of memory accesses determined by program (software)
- Visibility Order (of all operations)
 - ✓ order of memory accesses observed by one or more processors (e.g. like it if they were gathered by an external observer)
 - ✓ every read from a location returns the value of the most recent write

Sequential consistency

“A multiprocessor system is sequentially consistent if the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program.” (Lamport 1979)

Program order based memory accesses cannot be subverted in the overall sequence, so they cannot be observed to occur in a different order by a “remote” observer

An example

CPU1

$[A] = 1; (a1)$

$[B] = 1; (b1)$

CPU2

$u = [B]; (a2)$

$v = [A]; (b2)$

$[A], [B] \dots$ *Memory*

$u, v \dots$ *Registers*

$a1, b1, a2, b2$

Sequentially consistent

**Visibility order does not violate
program order**

$b1, a2, b2, a1$

Not sequentially consistent

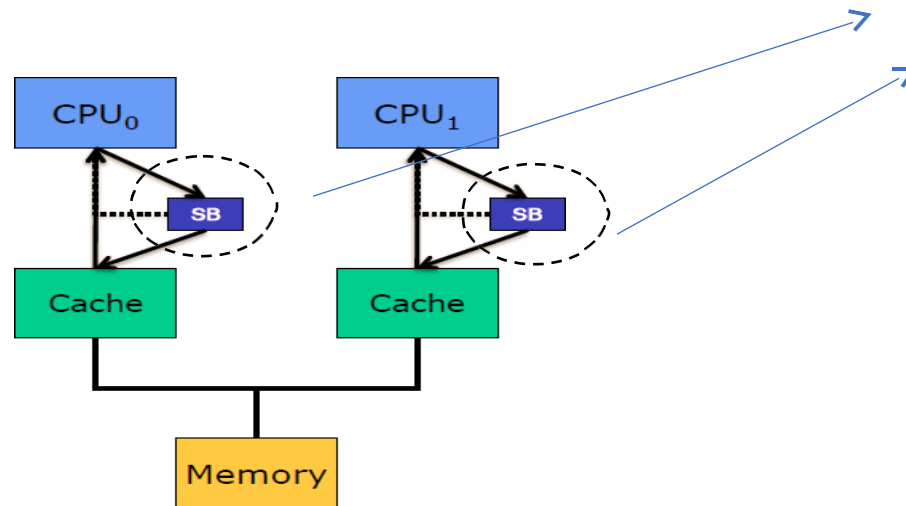
Visibility order violates program order

Total Store Order (TSO)

- Sequential consistency is “inconvenient” in terms of memory performance
- Example: cache misses need to be served “sequentially” even if they are write-operations with no currently depending instruction
- TSO is based on the idea that storing data into memory is not equivalent to writing to memory (as it occurs along program order)
- Something is positioned in the middle between a write operation (by software) and the actual memory update (in the hardware)
- A write materializes as a store when it is “more convenient” along time
- Several off-the-shelf machines rely on TSO (e.g. SPARC V8, x86)

TSO architectural concepts

- Store buffers allow writes to memory and/or caches to be saved to optimize interconnect accesses (e.g. when the interconnection medium is locked)
- CPU can continue execution before the write to cache/memory is complete (i.e. before data is stored)
- Some writes can be combined, e.g. video memory
- Store forwarding allows reads from local CPU to see pending writes in the store buffer
- Store buffer invisible to remote CPUs

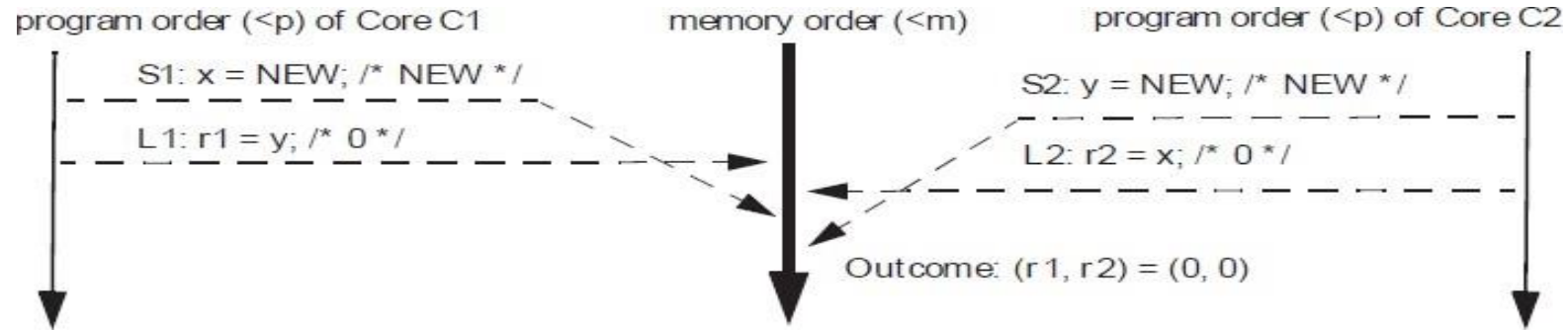


Store buffers not directly visible in the ISA

Forwarding of pending writes in the store buffer to successive read operations of the same location

Writes become visible to writing processor first

A TSO timeline



On x86 load operations may be reordered with older store operations to different locations

This breaks, e.g., Dekker's mutual exclusion algorithm

x86 memory synchronization

- x86 ISA provides means for managing synchronization (hence visibility) of memory operations
- SFENCE (Store Fence) instruction:
 - ✓ Performs a serializing operation on all store-to-memory instructions that were issued prior the SFENCE instruction. This serializing operation guarantees that every store instruction that precedes the SFENCE instruction in program order becomes globally visible before any store instruction that follows the SFENCE instruction.
- LFENCE (Load Fence) instruction:
 - ✓ Performs a serializing operation on all load-from-memory instructions that were issued prior the LFENCE instruction. Specifically, LFENCE does not execute until all prior instructions have completed locally, and no later instruction begins execution until LFENCE completes. In particular, an instruction that loads from memory and that precedes an LFENCE receives data from memory prior to completion of the LFENCE

x86 memory synchronization

- MFENCE (Memory Fence) instruction:
 - ✓ Performs a serializing operation on all load-from-memory and store-to-memory instructions that were issued prior the MFENCE instruction. This serializing operation guarantees that every load and store instruction that precedes the MFENCE instruction in program order becomes globally visible before any load or store instruction that follows the MFENCE instruction
- Fences are guaranteed to be ordered with respect to any other serializing instructions (e.g. CUID, LGDT, LIDT etc.)
- Instructions that can be prefixed by LOCK become serializing instructions
- These are ADD, ADC, AND, BTC, BTR, BTS, **CMPXCHG**, DEC, INC, NEG, NOT, OR, SBB, SUB, XOR, XAND
- CMPXCHG is used by spinlocks implementations such as

```
int pthread_mutex_lock(pthread_mutex_t *mutex);  
int pthread_mutex_trylock(pthread_mutex_t *mutex);
```

Read-Modify-Write (RMW) instructions

- More generally, CMPXCHG (historically known as Compare-and-Swap – CAS) stands in the wider class of Read-Modify-Write instructions like also Fetch-and-Add, Fetch-and-Or etc...
- These instructions perform a pattern where a value is both read and updated (if criteria are met)
- This can also be done atomically, with the guarantee of not being interfered by memory accesses in remote program flows
- In the essence, the interconnection medium (e.g. the memory bus) is locked in favor of the processing unit that is executing the Read-Modify-Write instruction
- For architecture where RMW is only allowed on cache-line aligned data, we can avoid locking the memory bus, and we can just exploit the “exclusive” state of the MESI (or derived) cache coherency protocol

gcc built-in

```
void __mm_sfence(void)
void __mm_lfence(void)
void __mm_mfence(void)
bool __sync_bool_compare_and_swap (type *ptr, type
oldval, type newval)
.....
```

- The definition given in the Intel documentation allows only for the use of the types *int*, *long*, *long long* as well as their unsigned counterparts
- gcc will allow any integral scalar or pointer type that is 1, 2, 4 or 8 bytes in length

Implementing an active-wait barrier

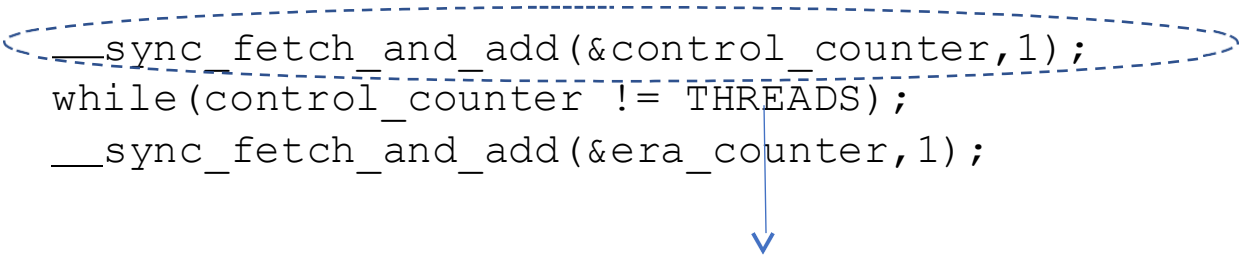
```
long control_counter = THREADS;
long era_counter = THREADS;

void barrier(void) {
    int ret;

    while(era_counter != THREADS && control_counter == THREADS);

    ret = __sync_bool_compare_and_swap(&control_counter, THREADS, 0);
    if(ret) era_counter = 0;

    __sync_fetch_and_add(&control_counter, 1);
    while(control_counter != THREADS);
    __sync_fetch_and_add(&era_counter, 1);
}
```



era_counter initial update already committed when performing this

ASM-based trylock via CMPXCHG

```
int try_lock(void * uadr) {  
    unsigned long r = 0;  
    asm volatile(  
        "xor %%rax,%%rax\n"  
        "mov $1,%%rbx\n"  
        "lock cmpxchg %%rbx, (%1) \n"  
        "sete (%0) \n"  
        : "r" (&r), "r" (uadr)  
        : "%rax", "%rbx"  
    );  
    return (r) ? 1 : 0;  
}
```

rax – eax – ax – al
are implicit registers
for *cmpxchg*

Target memory word

Set equal

If they were equal return 1

Locks vs (more) scalable coordination

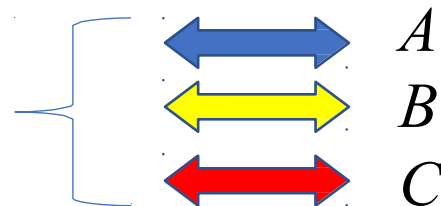
- The common way of coordinating the access to shared data is based on locks
- Up to now we understood what is the actual implementation of spin-locks
- In the end most of us never cared about hardware level memory consistency since spin-locks (and their Read-Modify-Write based implementation) **never leave pending memory updates upon exiting a lock protected critical section**
- Can we exploit memory consistency and the RMW support for achieving more scalable coordination schemes??
- The answer is yes
 - ✓ Non-blocking coordination (lock/wait-free synchronization)
 - ✓ Read Copy Update (originally born within the Linux kernel)

A recall on linearizability

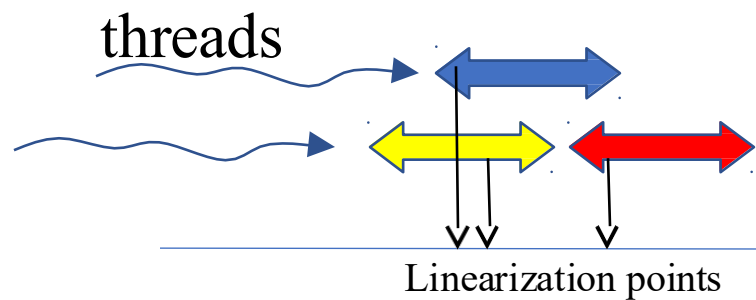
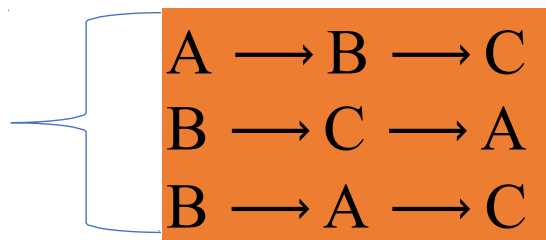
- A share data structure is “linearizable” if its operations always look to be sequentializable – we can make them equivalent to some sequential history
- This is true if
 - ✓ all its access methods/functions, although lasting a wall-clock-time period, can be seen as taking effect (materialize) at a specific point in time (they can be seen as isolated)
 - ✓ all the time-overlapping operations (if any) can be ordered based on their “selected” materialization instant
- Linearizability is a restriction of serializability since it involves operations on a single datum/object

A scheme

Operations (e.g. functions)
accessing a shared datum



Admissible
histories



$B \rightarrow C$ is a constraint

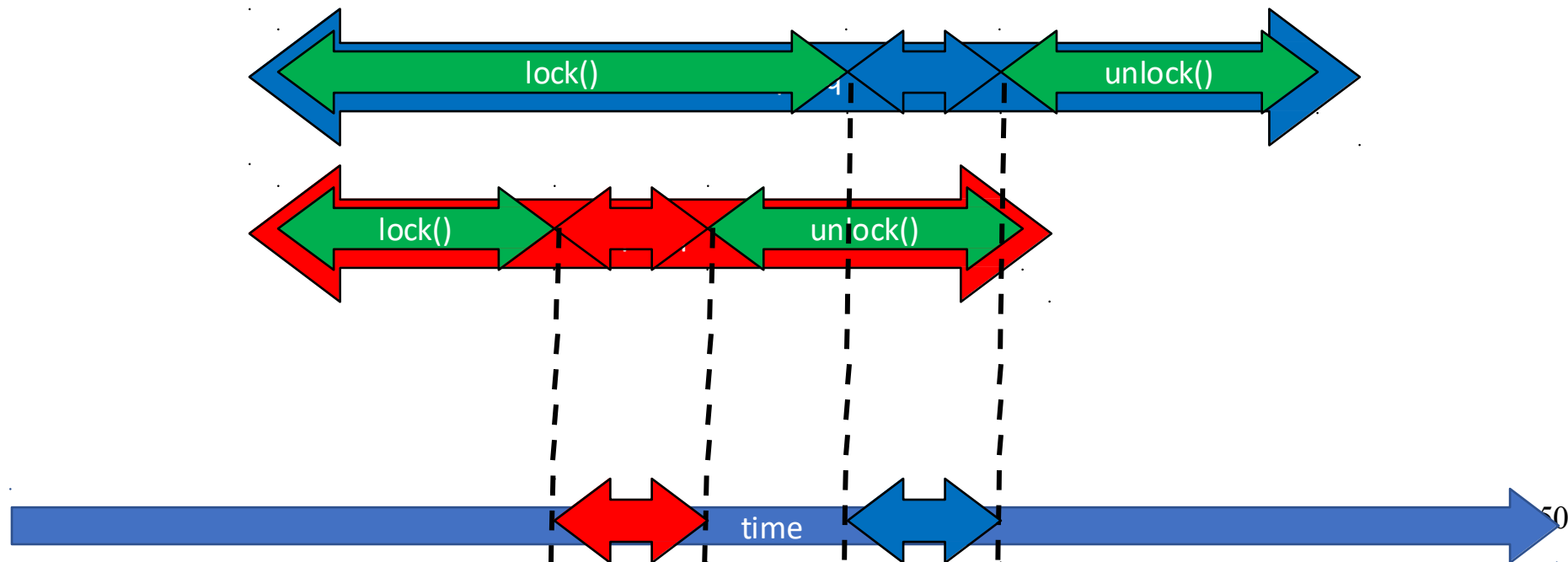
$A \rightarrow B \ \&\& \ C \rightarrow A$ would violate
linearizability

RMW vs linearizability

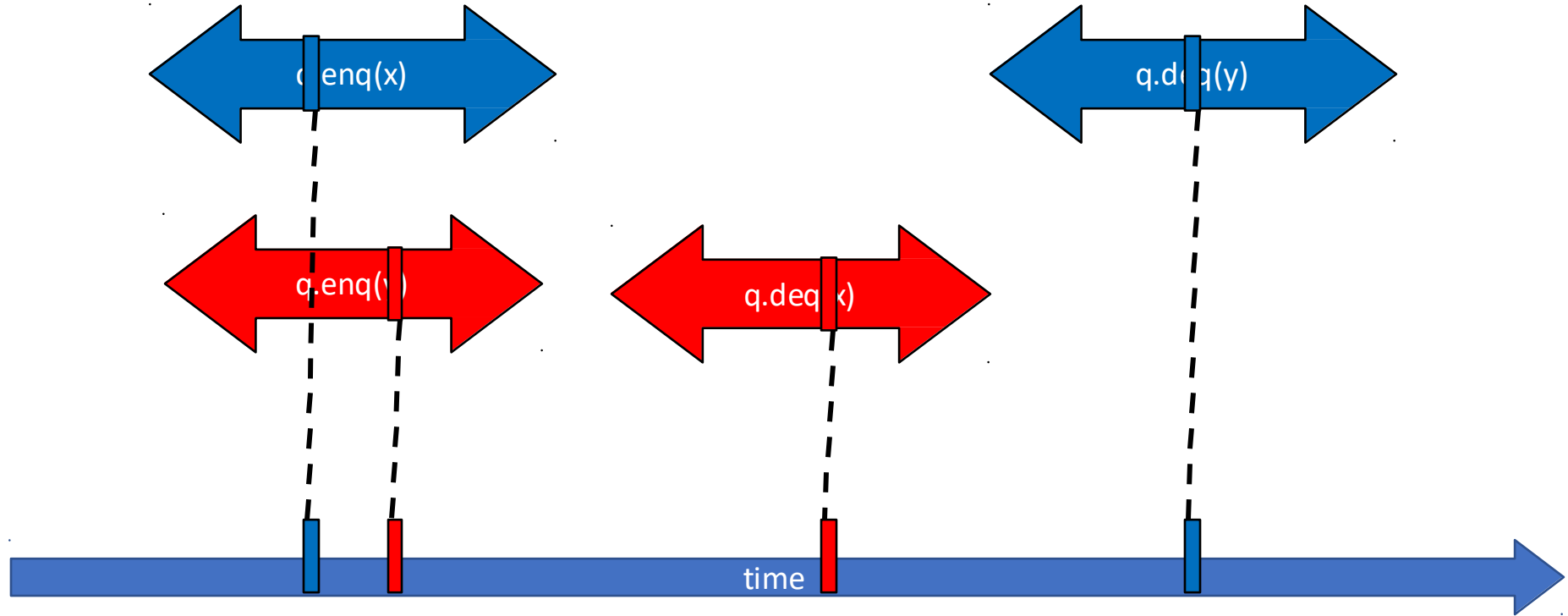
- Even though they implement non-trivial memory accesses, RMW instructions appear as atomic across the overall hardware architecture
- So they can be exploited to define linearization points of operations, thus leading to order the operations in a linearizable history
- The linearization points can be subject to differentiated execution paths (e.g. conditional branches)
- RMW instruction can fail, thus leading to drive subsequent RMW or other instructions, which can anticipate or delay the linearization point of the operations

RMW vs locks vs linearizability

- RMW-based locks can be used to create explicit wall clock time separation across operations
- We get therefore a sequential object with trivial linearization



Making RMW part of the operations



Lock-free vs wait-free synchronization

- They are both approaches belonging to **non-blocking synchronization**
- They differ in terms of the progress condition they ensure for the involved functions/methods
- **Lock-freedom**
 - ✓ Some instance of the function call successfully terminates in a finite amount of time (eventually)
 - ✓ All instances terminate successfully (or not) in a finite amount of time (eventually)
- **Wait-freedom:**
 - ✓ All instances of the function call successfully terminate in a finite amount of time (eventually)

Advantages from non-blocking synchronization

- Any thread can conclude its operations in a finite amount of time (or execution steps) independently of the other threads behavior – what if a thread crashes?!?
- This is highly attractive in modern contexts based on, e.g. CPU-stealing approaches – see Virtual Machine operations
- In classical blocking synchronization (e.g. based on spin-locks) what determines the actual number of computing steps (and time) for finalizing a given function are
 - ✓ The behavior of the lock holding thread
 - ✓ The actual sequence of lock acquisitions
- This is no longer true in non-blocking synchronization

Look-freedom aspects

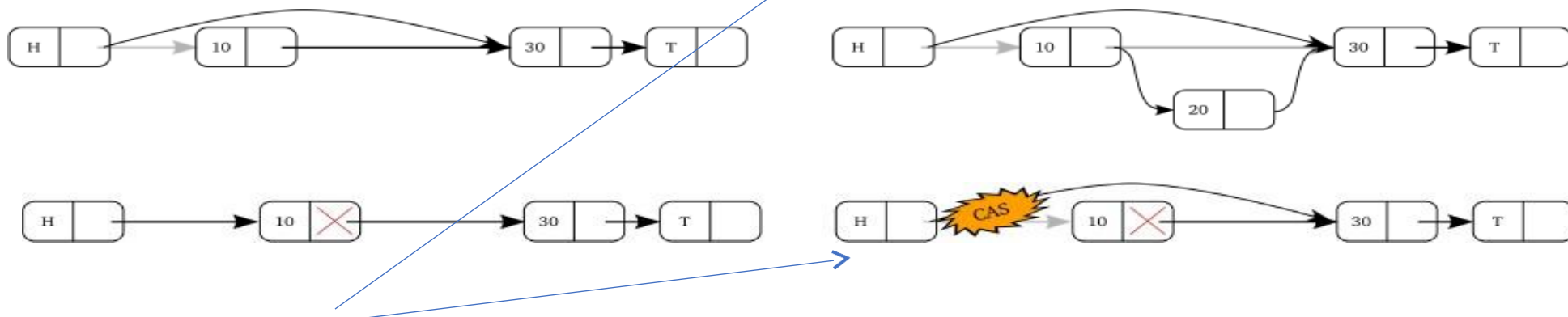
- If two ordered operations are incompatible (they conflict – possibly leading some RMW to fail), then one of them can be accepted, and the other one is refused (and maybe retried)
- Look-free algorithms are based on abort/retry
- The assumption for the effectiveness is that the aborts will occur infrequently, thus not wasting work too much
- The tradeoffs in the design are based here on devising data structures whose actual operations are somehow brought to balance operate on different parts of the data structure
- But this is not always possible!!

On the lock-free linked list example

- Insert via CAS on pointers (based on failure retries)



- Remove via CAS on node-state prior to node linkage



These CAS can fail but likely will not depending on the access pattern

On the wait-free atomic (1,N) register example

- It allows a writer to post atomically to N readers a new content
- A CAS on a pointer (with no other algorithmic step or register management logic) is not sufficient to guarantee that we can use a finite amount of memory to solve this problem
- The literature says that the lower bound on the number of buffers to use is $N+2$, and we should aim at this
 - N can be currently all read by the concurrent N readers
 - 1 can keep a new value, not yet accessed by any reader
 - 1 can be used to fill some new content

Actual register operations [ARC - TPDS journal 2018]

Algorithm 1 Register initialization.

```
1: procedure INIT(content, size)
2:   for all slot  $\in [0, N + 1]$  do
3:     register[slot].size  $\leftarrow 0$ 
4:     register[slot].r_start  $\leftarrow 0$ 
5:     register[slot].r_end  $\leftarrow 0$ 
6:   MEMCOPY(register[0].content, content, size)
7:   register[0].size  $\leftarrow$  size
8:   current  $\leftarrow N$  ▷ I1
```

Algorithm 2 The atomic register read operation.

```
1: procedure READ()
2:   index  $\leftarrow$  current  $\gg 32$  ▷ R1
3:   if last_index = index then
4:     entry  $\leftarrow$  register[last_index]
5:     return  $\langle$ entry.content, entry.size $\rangle$  ▷ R2
6:   ATOMICINC(register[last_index].r_end) ▷ R3
7:   tmp_curr  $\leftarrow$  ATOMICADDANDFETCH(current, 1) ▷ R4
8:   last_index  $\leftarrow$  tmp_curr  $\gg 32$  ▷ R5
9:   entry  $\leftarrow$  register[last_index]
10:  return  $\langle$ entry.content, entry.size $\rangle$ 
```

Unique synchronization
variable with 2 fields



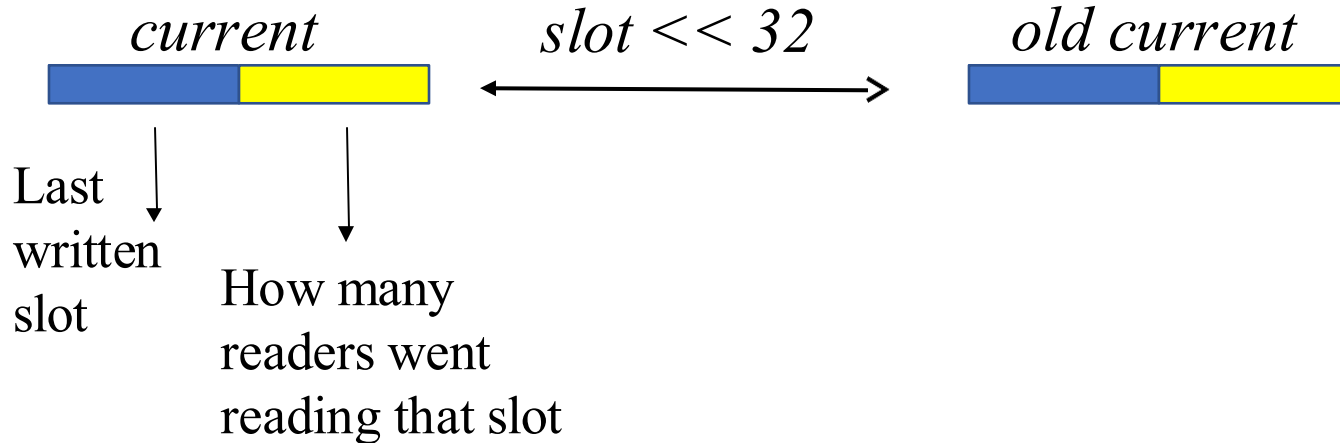
Last written
slot

How many
readers went
reading that slot

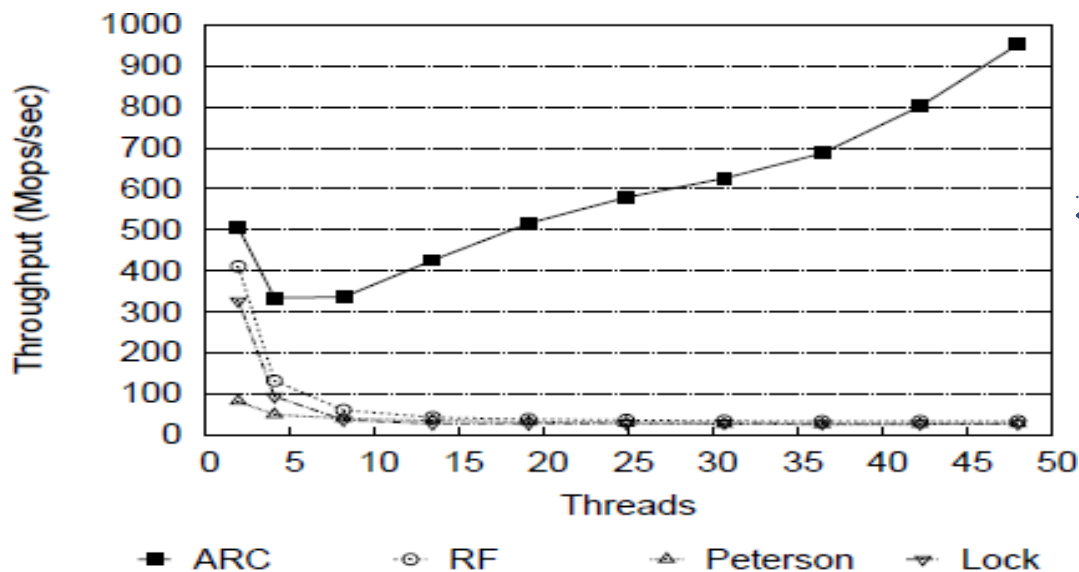
Actual register operations

Algorithm 3 The atomic register write operation.

```
1: procedure WRITE(content, size)
2:   pick slot such that  $slot \neq last\_slot \wedge register[slot].r\_start = register[slot].r\_end$ 
3:   MEMCOPY( $register[slot].content$ , content, size)
4:    $register[slot].size \leftarrow size$ 
5:    $register[slot].r\_start \leftarrow 0$ 
6:    $register[slot].r\_end \leftarrow 0$ 
7:    $old\_curr \leftarrow \text{ATOMICEXCHANGE}(current, slot \ll 32)$ 
8:    $old\_slot \leftarrow old\_curr \gg 32$ 
9:    $register[old\_slot].r\_start \leftarrow old\_curr \& (2^{32} - 1)$ 
10:   $last\_slot \leftarrow slot$ 
```



Performance



48 CPU-core
machine deploy



■ ARC

○ RF

△ Peterson

▽ Lock

(a) 128B register size

Can linearize read operations
with no RMW instruction

Always needs RMW instructions to
linearize read operations

Makes 2 data copies to
assess consistency upon reading

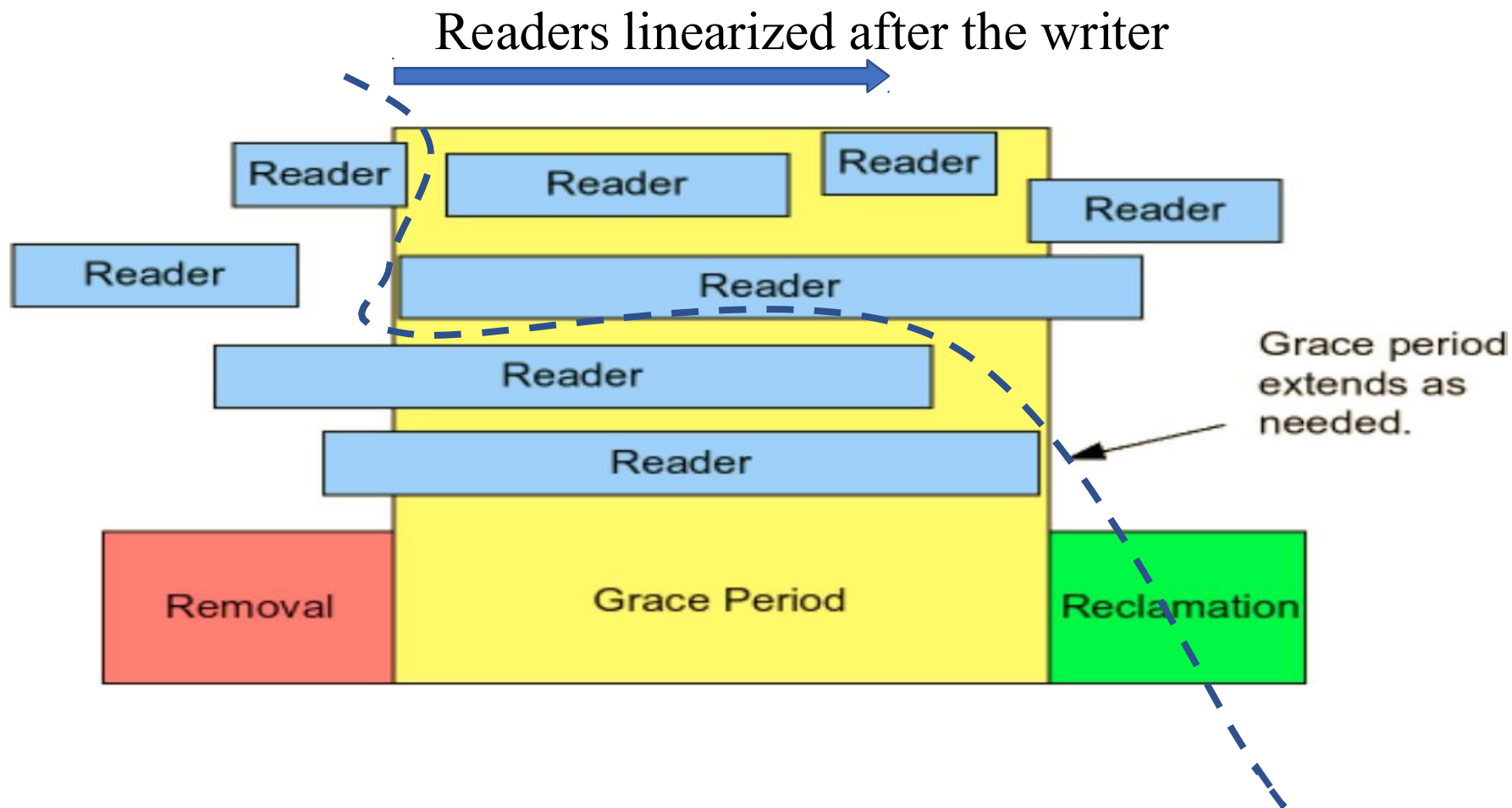
The big problem with generic data structures - buffer re-usage

- RMW based approaches allow us to understand what is the state of some (linked) data structure (still in or already out of a linkage)
- But we cannot understand if traversals on that data structure are still pending
- If we reuse the data structure (e.g. modifying its fields), we might give rise to data structure breaks
- This may even lead to security problems:
 - ✓ We traverse via a thread un-allowed pieces of information

Read Copy Update (RCU)

- Baseline idea
 - ✓ A writer at any time
 - ✓ Concurrency between readers and writers
- Actuation
 - ✓ Out-links of logically removed data structures are not destroyed prior being sure that no reader is still traversing the modified copy of the data structure
 - ✓ Buffer re-reuse (hence release) takes place at the end of a so called “**grace period**”, allowing the standing readers not linearized after the update to still proceed
- Very useful for read intensive shared data structures

General RCU timeline



RCU reads and writes

- The reader
 - ✓ Signals it is there
 - ✓ It reads
 - ✓ Then it signals it is no longer there
- The writer
 - ✓ Takes a write lock
 - ✓ Updates the data structure
 - ✓ Waits for standing readers to finish
 - ✓ **NOTE: readers operating on the modified data structure instance are don't care readers**
 - ✓ Release the buffers for re-usage

Kernel level RCU with no reader preemptability

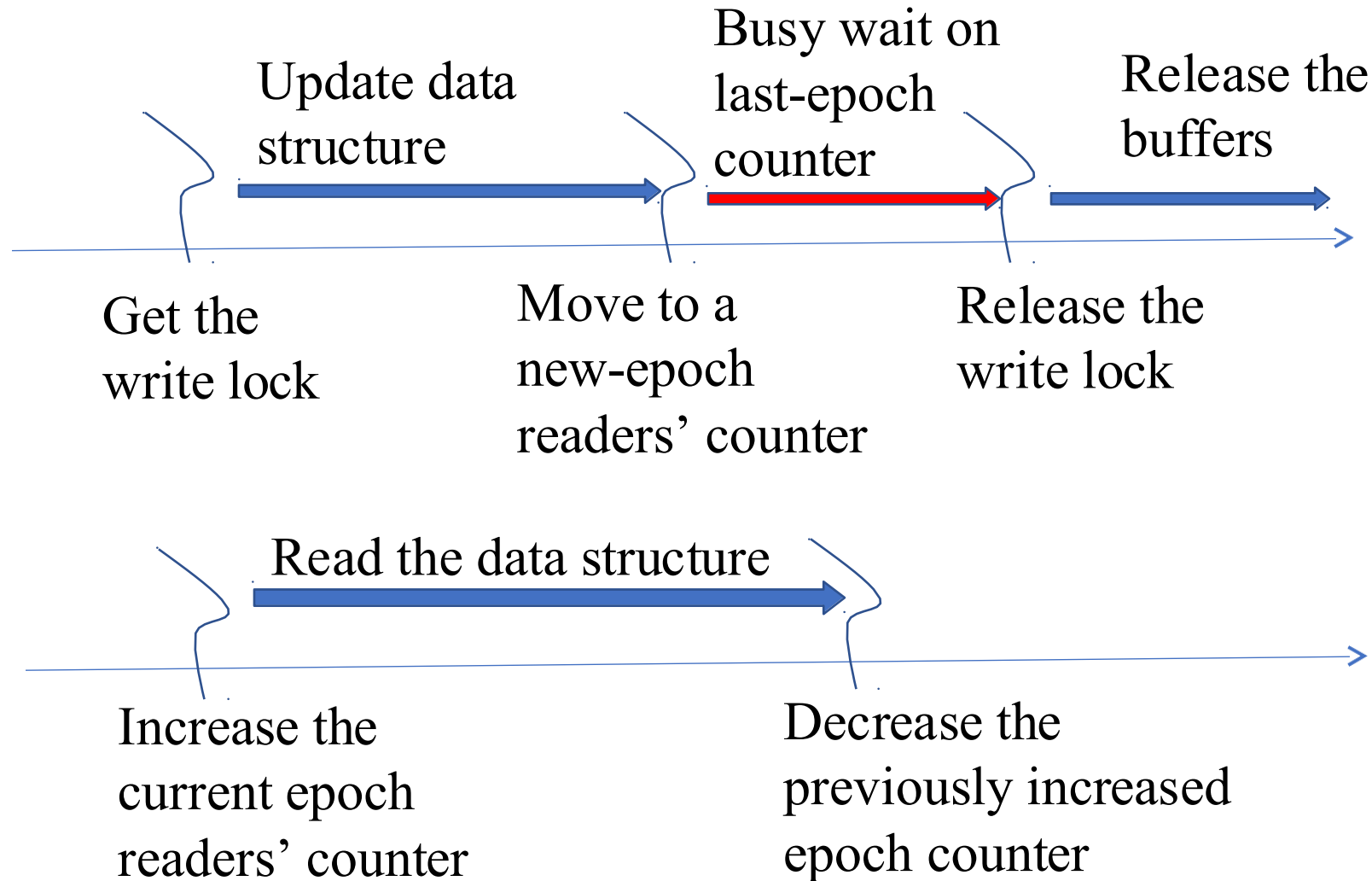
- The reader only needs to turn off preemption upon a read and re-enable it upon finishing (beware the interactions with RT scheduling) – but no call to blocking services with possible context switch is allowed
- The writer understands that no standing reader is still there thanks to its own migration to all the remote CPUs, in Linux as easily as

```
for_each_online_cpu(cpu)    run_on(cpu) ;
```
- The migrations create a context-switch leading the writer to understand that no standing reader, not linearized after the writer, is still there
- Alternatively, CPUs publish their context switch occurrence and the writer reads it

Preemptable (e.g. user level) RCU

- Discovering standing readers in the grace periods is a bit more tricky
- An atomic presence-counter indicates an ongoing read
- The writer updates the data structure and redirects readers to a new presence counter (a new epoch)
- It then waits up to the release of presence counts on the last epoch counter
- Data-structure updates and epoch move are not atomic
- However, the only risk incurred is the one of waiting for some reader that already saw the new shape of the data structure, but got registered as present in the last epoch

Preemptable CRU - example reader/writer timeline



A few kernel level API in Linux

- `rcu_read_lock()`
 - This signals that the reader has started operating (e.g. it may block preemptability, or do nothing in a kernel with `CONFIG_PREEMPT = n`)
- `rcu_read_unlock()`
 - This signals the reader has ended its activity
- `synchronize_rcu()`
 - This API implements the wait of the end of the grace period
- `call_rcu(...)`
 - This API enables the usage of a callback to be executed after the grace period has ended, but does not lock the updater thread

Additional kernel level API in Linux – sleepable RCU

```
int init_srcu_struct(struct srcu_struct *sp);
```

```
void cleanup_srcu_struct(struct srcu_struct *sp)
```

```
int srcu_read_lock(struct srcu_struct *sp);
```

```
void srcu_read_unlock(struct srcu_struct *sp, int idx);
```

```
void synchronize_srcu(struct srcu_struct *sp);
```

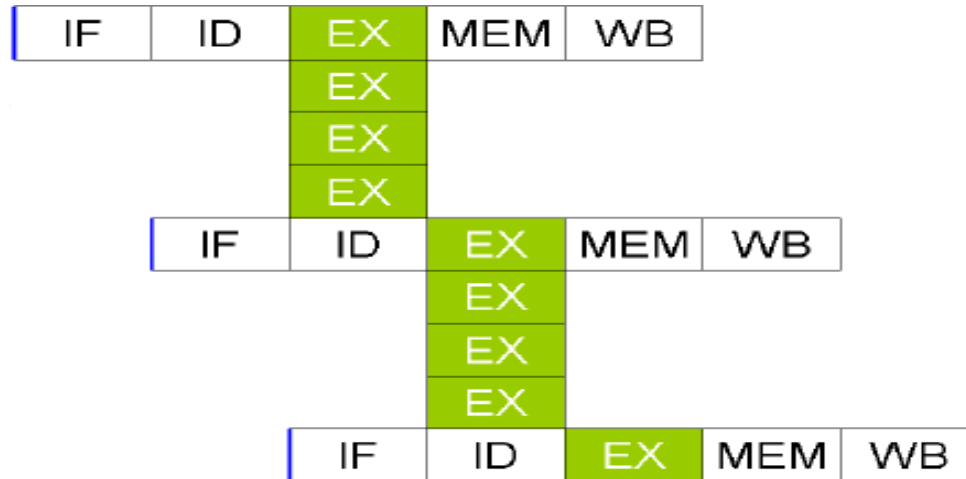
Somehow similar (in terms of functionality) to a generic preemptable implementation of RCU

Additional parallelization aspects

- This is the so called “vectorization”
- It was born way before speculative computing and multi-processor/multi-core
- Essentially it is a form of **SIMD** (Single Instruction Multiple Data) processing
- As opposed to classical **MIMD** (Multiple Instruction Multiple Data) processing of multi-processors/multi-cores
- **SIMD** is based on vectorial registers and/or vectorial computation hardware (e.g. GPUs)
- Less common is **MISD** (although somebody says that a speculative processor is MISD)
- ... **SISD** is a trivial single-core non speculative machine

The vector processor scheme

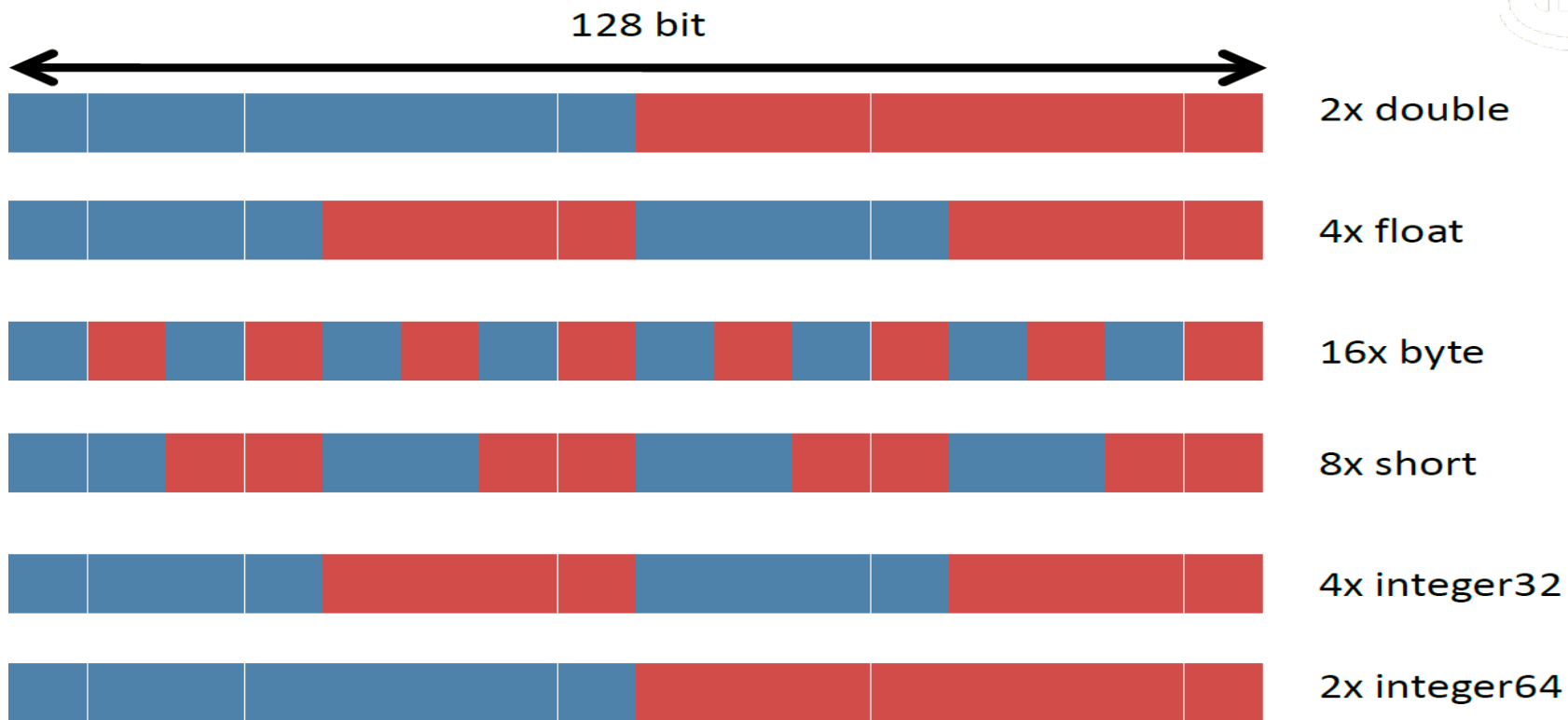
- Vector registers
- Vectorized and pipelined functional units
- Vector instructions
- Hardware data scatter/gather



x86 vectorization

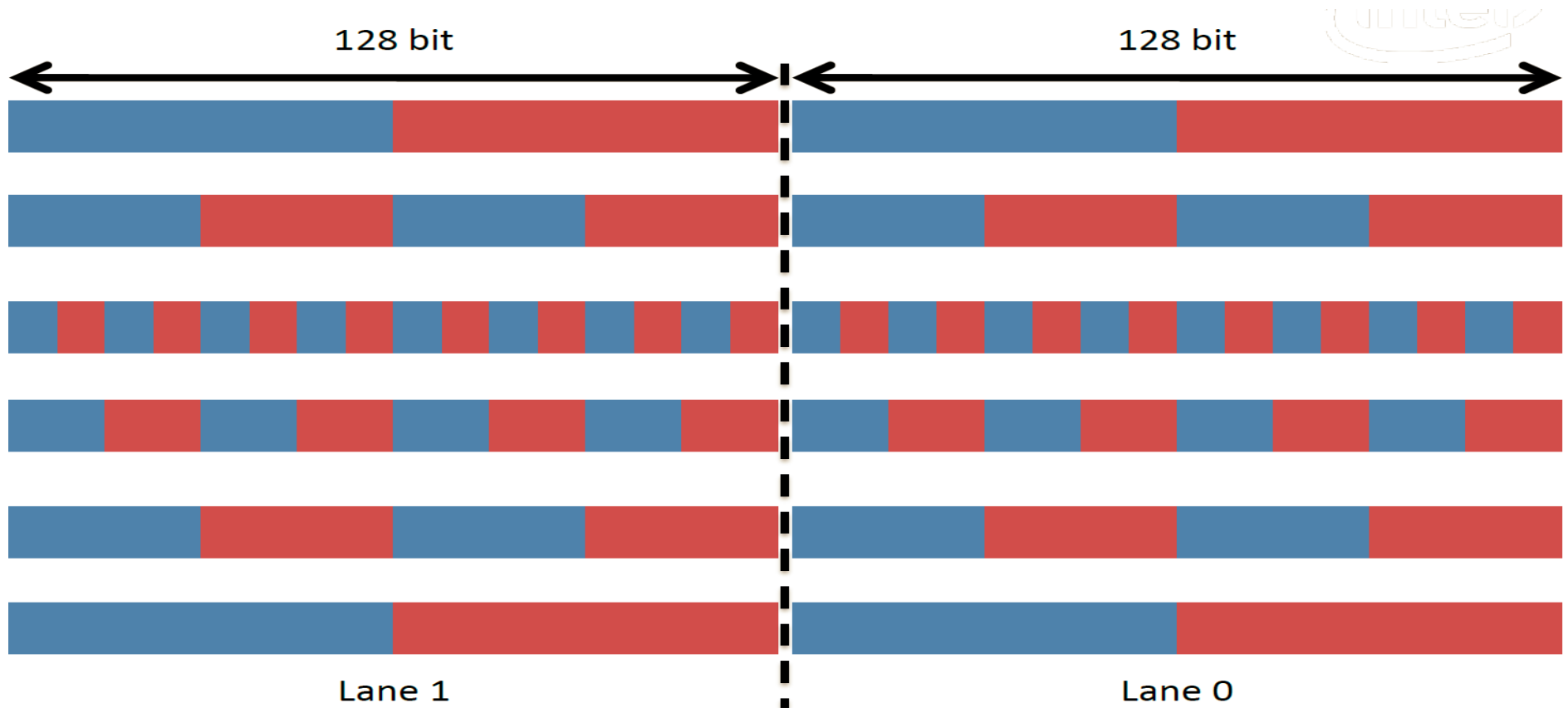
- Called SSE (Streaming SIMD Extension)
- Introduced in Pentium III processors (from 1999)
- Previously called MMX (MultiMedia eXtension or Matrix Math eXtension) on Pentium processors (based on 64-bit registers)
- In the SSE programming mode there are 8 128-bit XMM registers (16 in x86-64 SSE2) keeping vectors of
 - ✓ 2 64-bit double precision floats
 - ✓ 4 32-bit single precision floats
 - ✓ 1 128-bit integer (unsigned)
 - ✓ 2 64-bit integers (signed or unsigned)
 - ✓ 4 32-bit integers (signed or unsigned)
 - ✓ 8 16-bit integers (signed or unsigned)
 - ✓ 16 8-bit integers (signed or unsigned)

SSE data types



Sandy Bridge AVX (Advanced Vector Extensions)

- Registers are this time YMM[0-15]



Memory alignment

- Memory/XMM*/YMM* data move instructions in x86 operate with either 8/16-byte aligned memory or not
- Those with aligned memory are faster
- gcc offers the support for aligning static (arrays of) values via the `__attribute__((aligned (16)))`
- It enables compile level automatic vectorization with `-O` flags (originally `-O2`), whenever possible
- Clearly, one may also resort to dynamic memory allocation with explicit alignment
- 4KB page boundaries are intrinsically 16-bit aligned, which helps with `mmap()`
- Usage of instructions requiring alignment on non-aligned data will cause a general protection error

C intrinsics for SSE programming

- Vectorized addition - 8/16/32/64-bit integers

MMX(64-bit)	<code>d = _mm_add_pi8(a, b)</code>	<code>__m64 a, b;</code>	<code>__m64 d;</code>
SSE2(128-bit)	<code>d = _mm_add_epi8(a, b)</code>	<code>__m128i a, b;</code>	<code>__m128i d;</code>

MMX(64-bit)	<code>d = _mm_add_pi16(a, b)</code>	<code>__m64 a, b;</code>	<code>__m64 d;</code>
SSE2(128-bit)	<code>d = _mm_add_epi16(a, b)</code>	<code>__m128i a, b;</code>	<code>__m128i d;</code>

MMX(64-bit)	<code>d = _mm_add_pi32(a, b)</code>	<code>__m64 a, b;</code>	<code>__m64 d;</code>
SSE2(128-bit)	<code>d = _mm_add_epi32(a, b)</code>	<code>__m128i a, b;</code>	<code>__m128i d;</code>

MMX(64-bit)	<code>d = _mm_add_si64(a, b)</code>	<code>__m64 a, b;</code>	<code>__m64 d;</code>
SSE2(128-bit)	<code>d = _mm_add_epi64(a, b)</code>	<code>__m128i a, b;</code>	<code>__m128i d;</code>

- Vectorized addition - 32-bit floats

SSE(64-bit)	<code>d = _mm_add_ss(a, b)</code>	<code>__m128 a, b;</code>	<code>__m128 d;</code>
SSE(128-bit)	<code>d = _mm_add_ps(a, b)</code>	<code>__m128 a, b;</code>	<code>__m128 d;</code>

- Vectorized addition - 64-bit doubles

SSE2(64-bit)	<code>d = _mm_add_sd(a, b)</code>	<code>__m128d a, b;</code>	<code>__m128d d;</code>
SSE2(128-bit)	<code>d = _mm_add_pd(a, b)</code>	<code>__m128d a, b;</code>	<code>__m128d d;</code>

Additional C intrinsics

- Additional features are available for, e.g.:
 - ✓ Saturated addition
 - ✓ Subtraction
 - ✓ Saturated subtraction
 - ✓ Addition/subtraction with carry
 - ✓ Odd/even addition/subtraction
 - ✓ In vector sum reduction
- Similar functionalities are offered for the AVX case