

It's not over with Meltdown

L1TF, POP SS, TLBleed and more ...

Vojtěch Pavlík

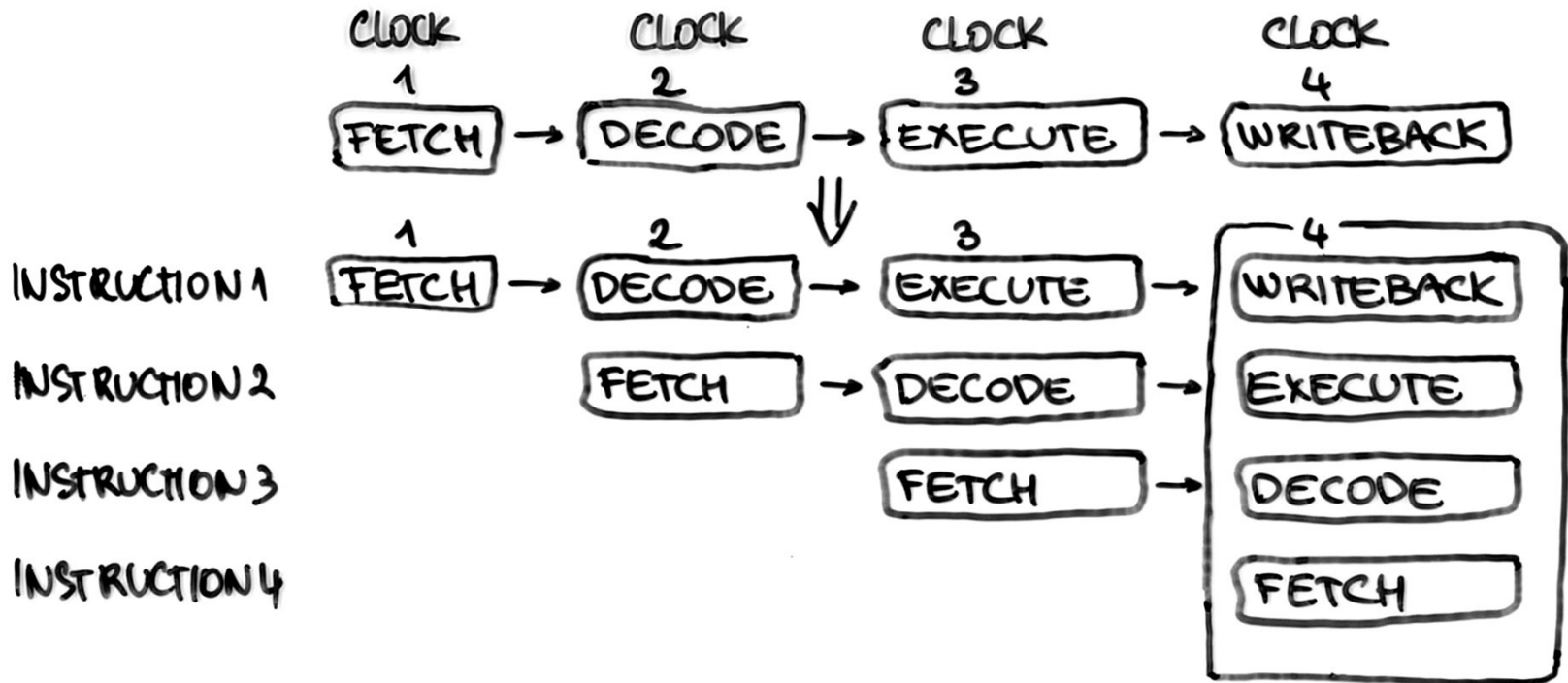
Director SUSE Labs



1967: Robert Marco Tomasulo

- Implemented a revolutionary algorithm
 - in IBM System/360 FPU
 - significantly improved execution efficiency
 - today known as Tomasulo's algorithm
- Received Eckert-Mauchly prize
 - like a Nobel prize in computer architecture
 - named after the creators of ENIAC
- But the consequences ...

Pipelining



Branch prediction

LOOP:

MUL %RAX %RBX

DEC %RAX

JNZ %LOOP

⋮

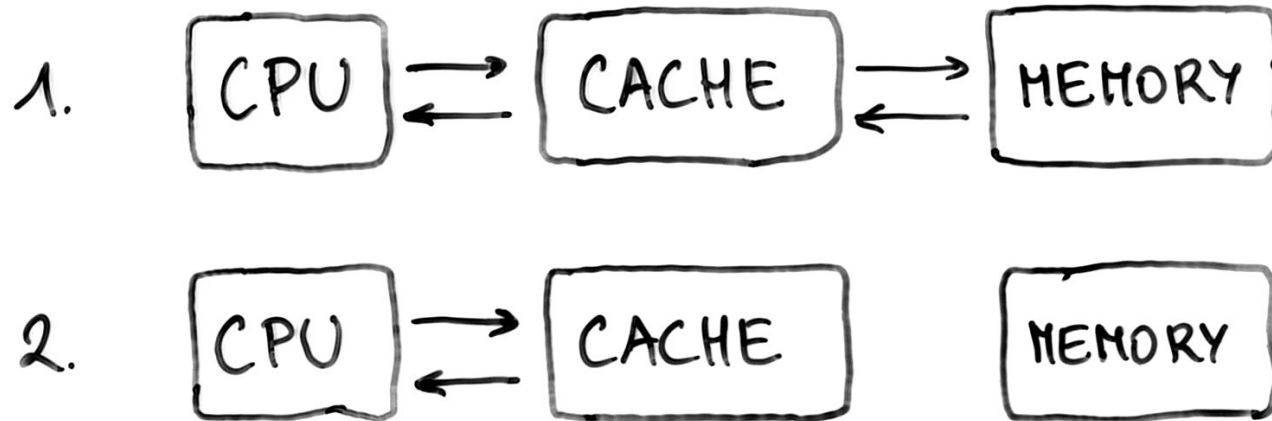
⋮



ORACLE

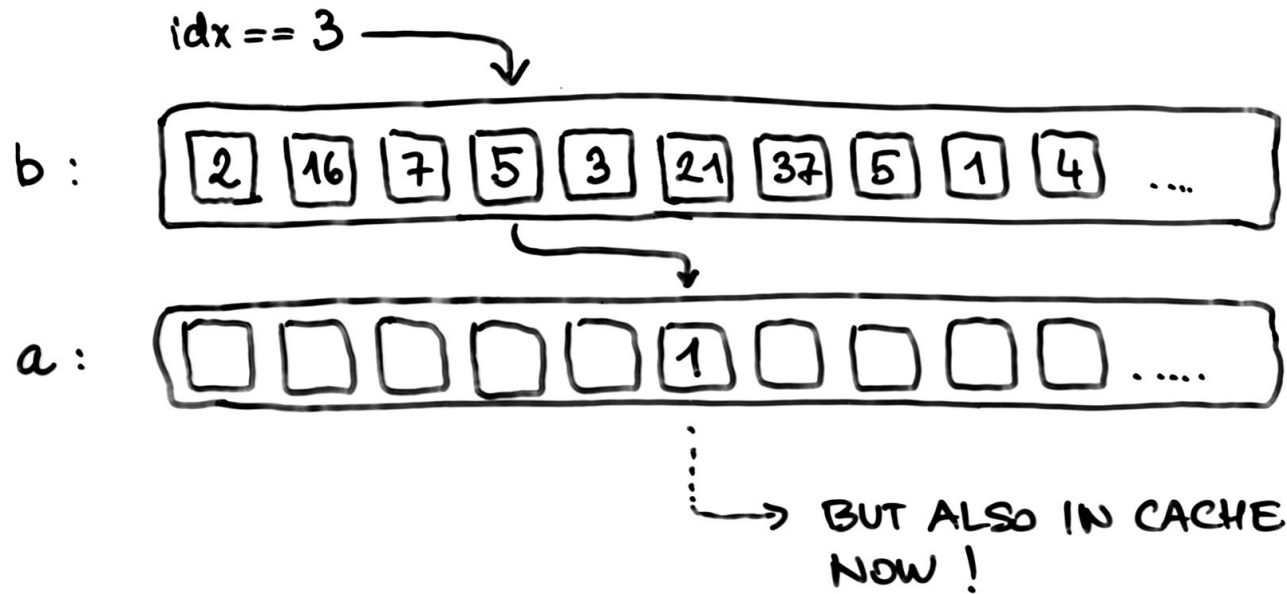
BRANCH
PREDICTOR

Cache

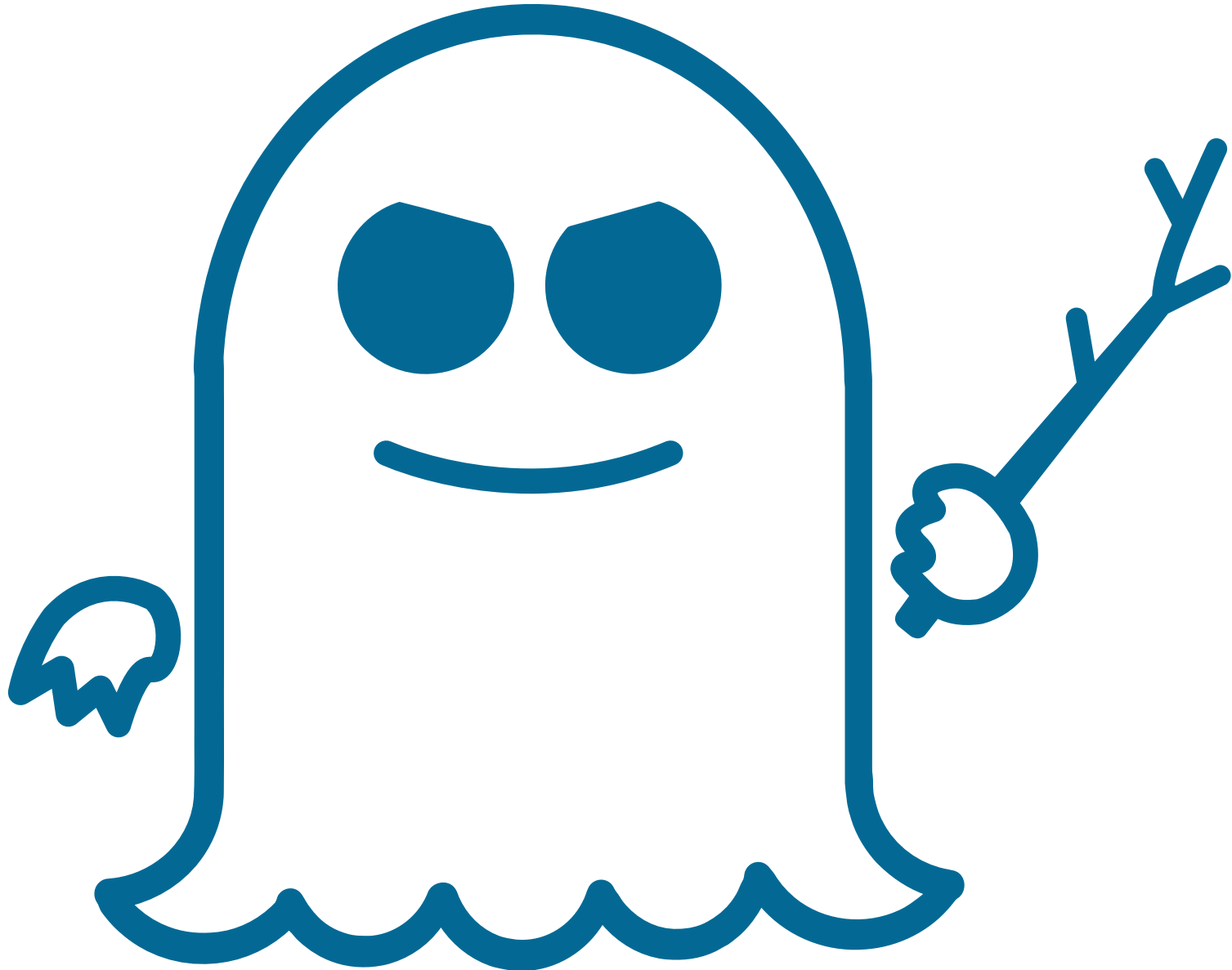


Double indirection

$a[b[idx]] = 1;$



IF WE KNOW WHICH OF A IS IN CACHE,
WE KNOW WHAT VALUE OF B IS AT IDX



Spectre v1 Attack!

```
ATTACKER
for (i=0; i < 1000; i++)
    syscall(0);

prime_cache();

syscall(1000);

for (i=0; i < CACHE_SIZE; i++)
    measure_access_time(i);
```



```
KERNEL

int syscall(int idx)
{
    if (idx < 16)
        t = a[b[idx]];

    ...
    return 0;
}
```


Arithmetics

```
int syscall (int idx)
{
    t = a[b[idx & 0xf]];
    if (idx < 16) {
        ⋮
    }
}
```

Indirect branch prediction

object → method ();

MOV %RAX OBJECT[%RSI]

JMP %RAX



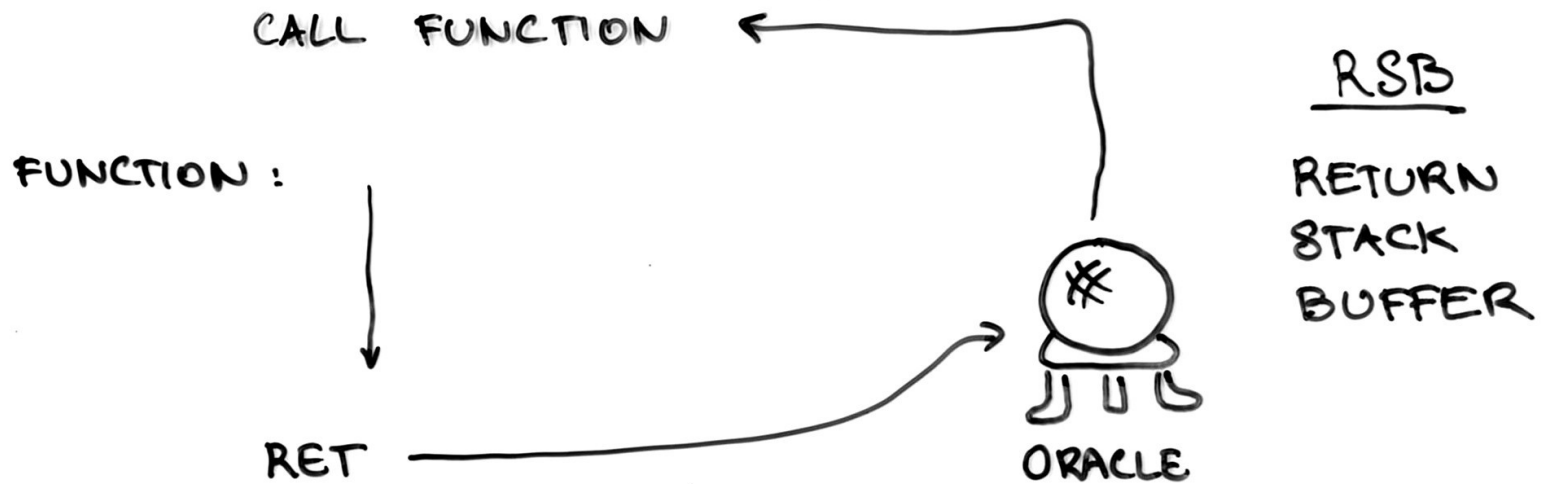
ORACLE



BTB

BRANCH
TARGET
BUFFER

Return target prediction



SpectreRSB Attack!

```
ATTACKER  
  
void stuff();  
{  
    static int cnt = 0;  
    if (cnt++ == 16)  
        goto done;  
  
    stuff();  
}  
  
stuff();  
done:
```

↓ ... kernel switches stacks ...

```
VICTIM  
  
return;  
    ↙  
    ↘  
t = a[b[idx]];
```

RSB Stuffing

- To overcome SpectreRSB the kernel has to flush the RSB contents on every context switch
- Same what the attacker did – do a sequence of calls and remove them from stack
- Costs performance
(not much of a surprise, is it?)

Lazy FPU switching

- FPU register save/restore used to be slow on 80387
- But *disabling* the FPU was fast
- So the kernel would on a context switch
 - Save the FPU state if it was used
 - Disable the FPU
- The CPU would raise an exception on FPU access
 - The kernel then restores saved FPU state
 - And remembers the FPU was used
- This avoids many costly saves and restores since most tasks don't use FPU – gaining performance!

Lazy FPU Attack!

- When speculating, the CPU doesn't check if FPU is disabled
- Combined with double-indirection, attacker can snoop FPU register contents

Eager FPU save/restore

- To mitigate, the kernel must not disable the FPU
- And save/restore FPU state on every context switch
- This costs (you knew this was coming) performance
- But fortunately there's FXSAVE and FXRSTOR
Instructions for fast FPU state saving on modern CPUs
- So it's not that bad after all

x86 Stack

8086:

SS: STACK SEGMENT



SP: STACK POINTER



STACK ADDRESS 20-BIT



POP SS Magic

POP SS

← INTERRUPT

POP SP

↓
CRASH

POP SS

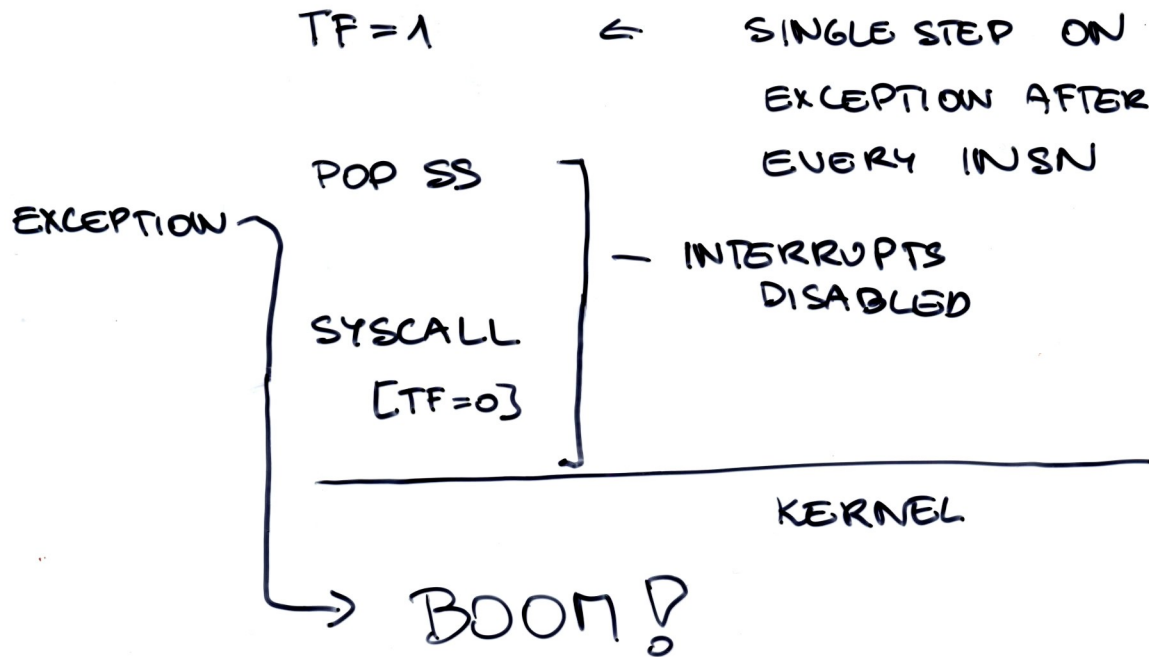
INTERRUPTS MAGICALLY

POP SP

DISABLED

POP SS Attack!

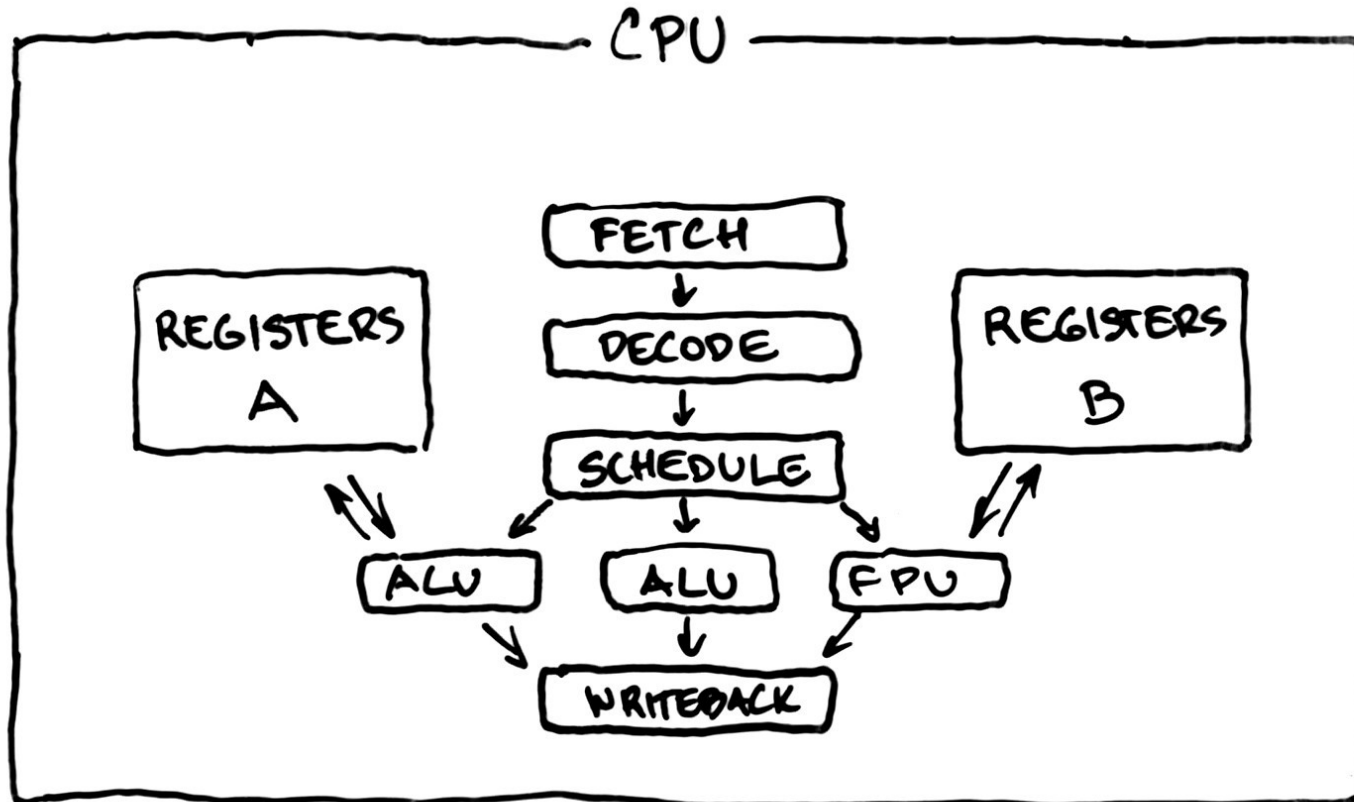
Modern x86



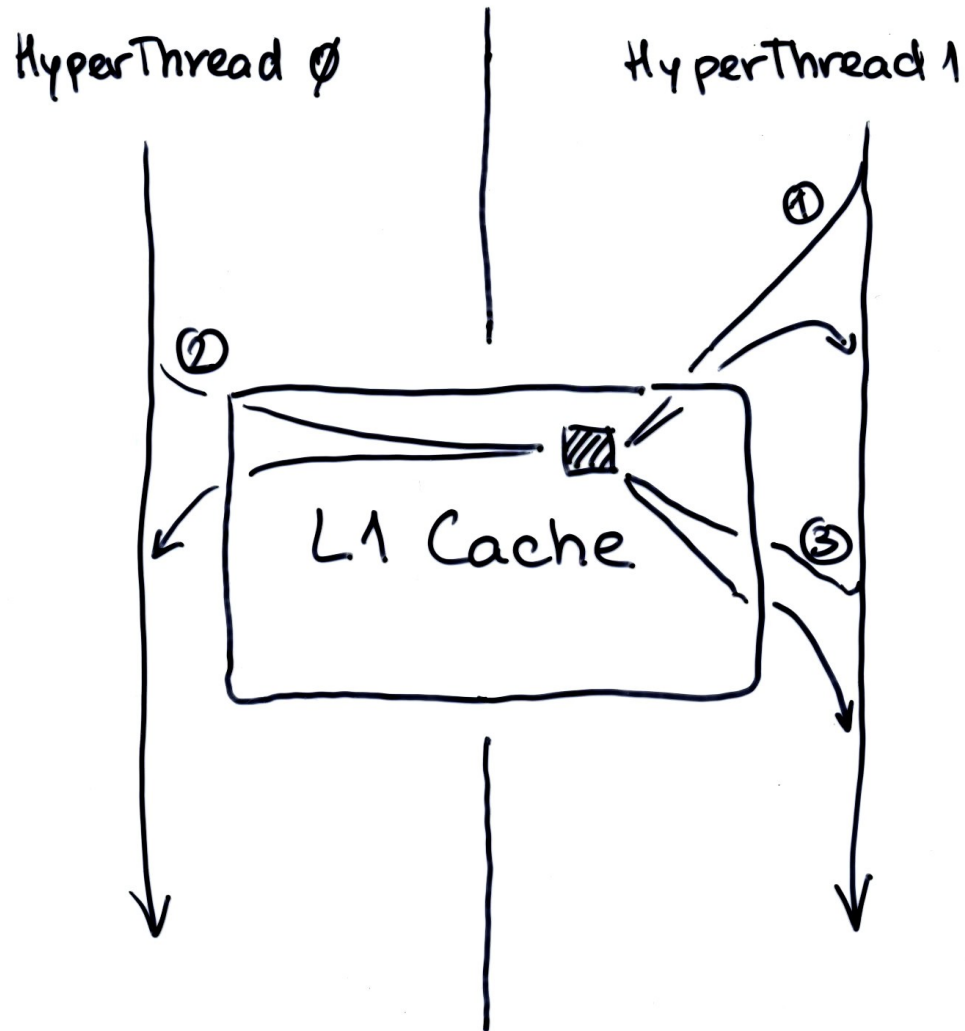
Interrupt stacks

- To overcome POP SS the kernel must use interrupt stack switching for ALL interrupts
- And when it sees an exception it needs to verify whether it may be coming from userspace
- **Costs performance**
(actually not that much, only in specific exception handlers)

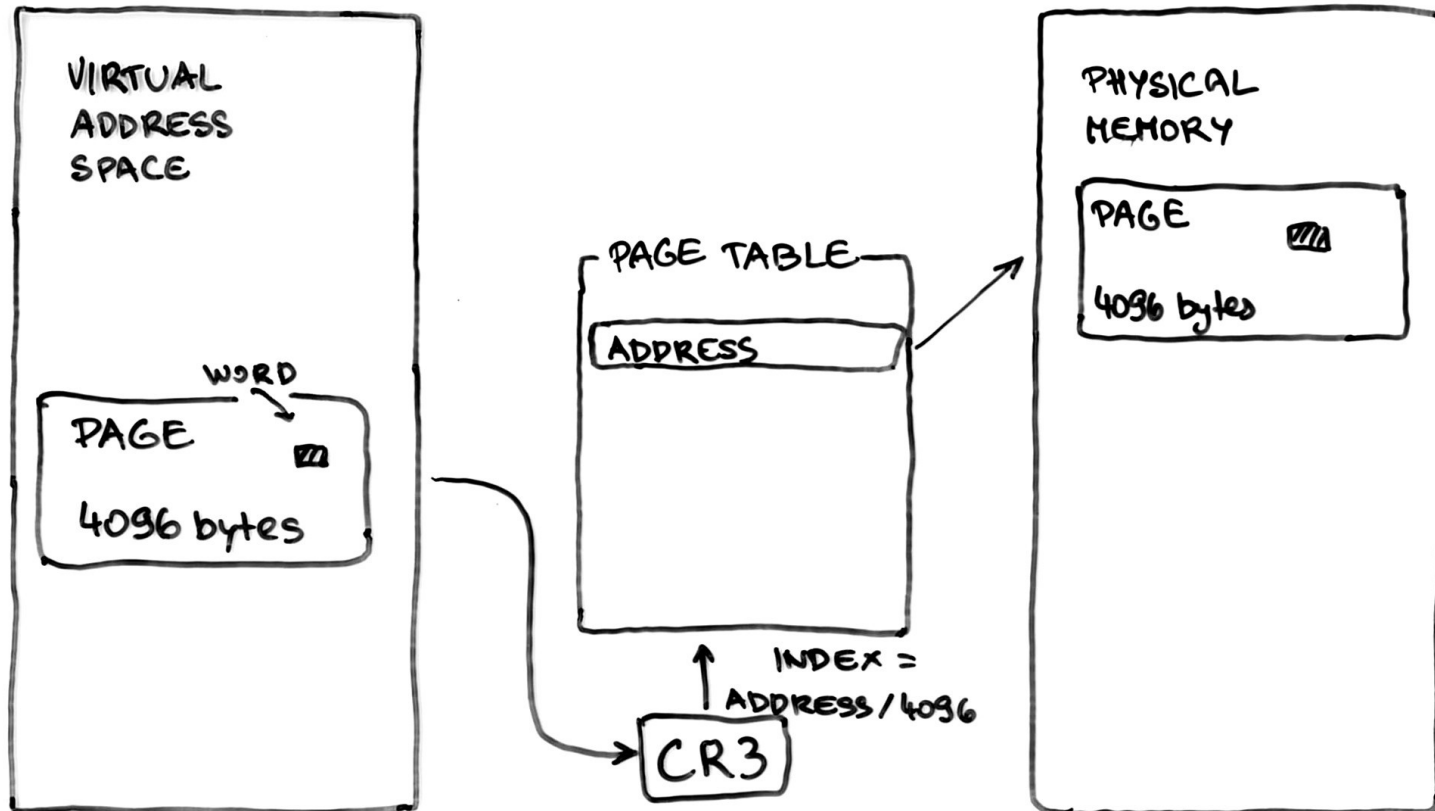
SMT / HT



Percival 2005 Attack!

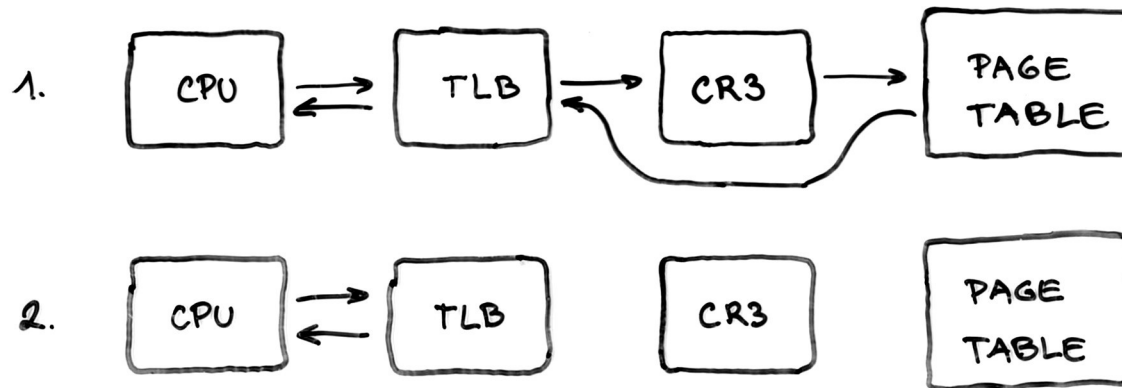


Page table and CR3



TLB

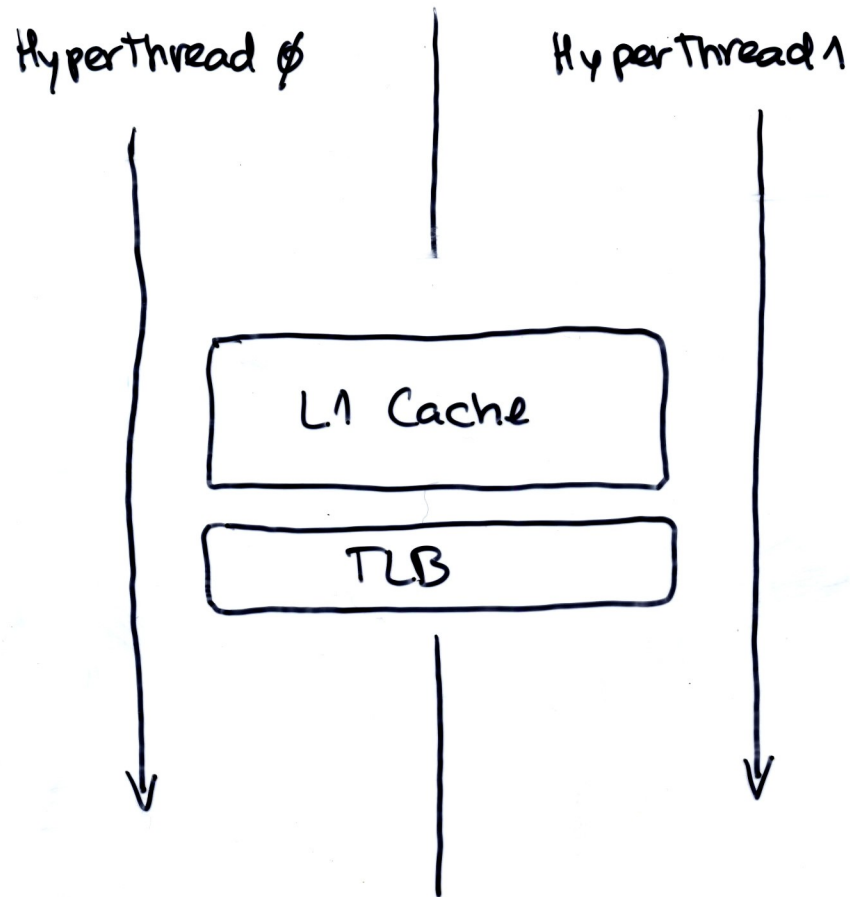
TRANSACTION LOOKASIDE BUFFER



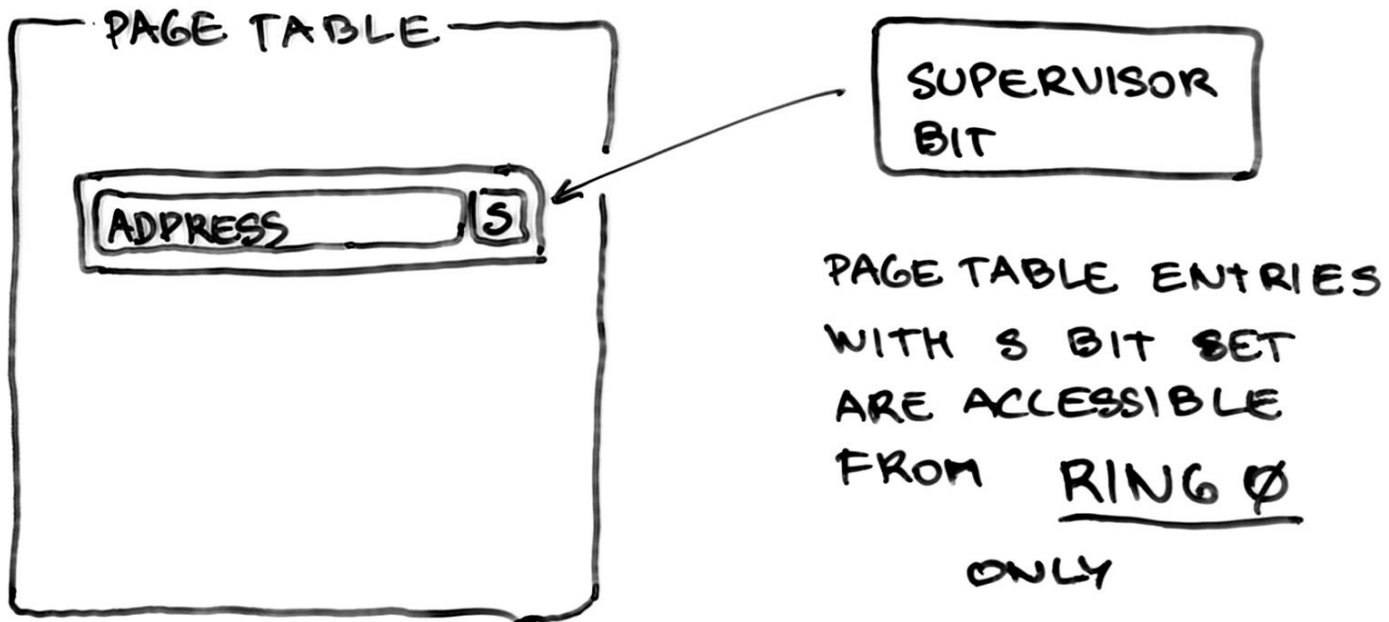
IF CR3 IS CHANGED

TLB IS EMPTIED!

TLBleed Attack!



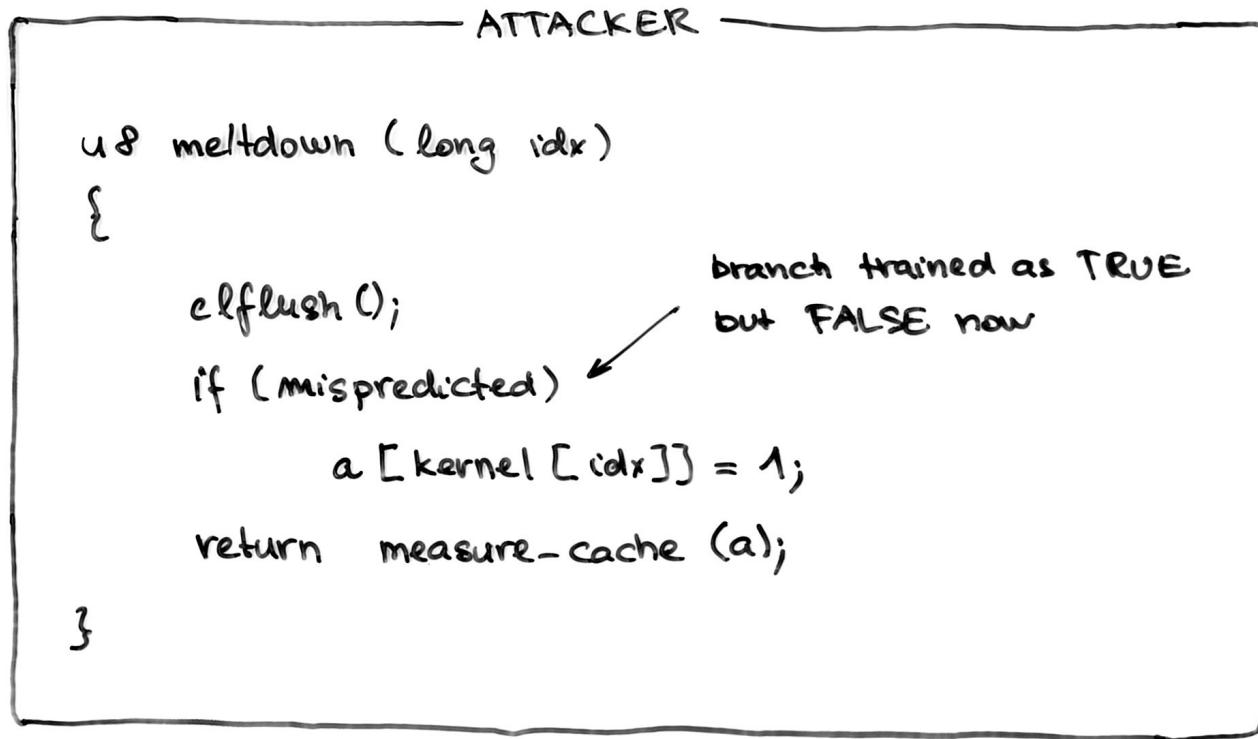
Supervisor bit





Meltdown (v3)

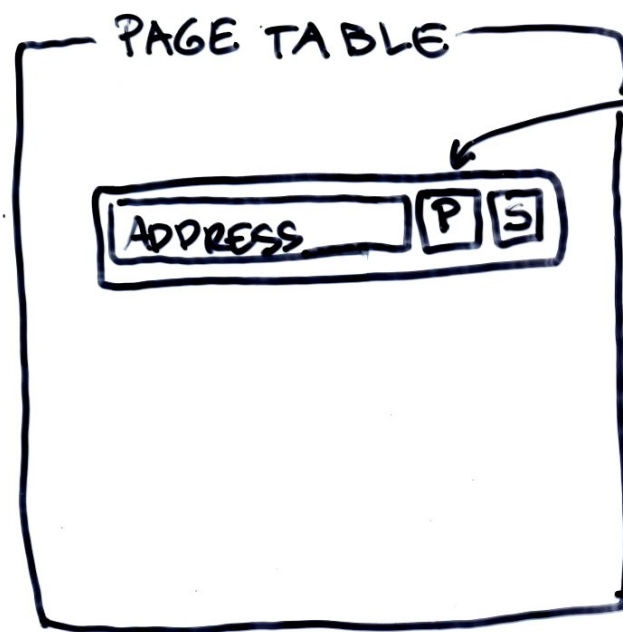
- Intel CPUs only look at the Supervisor bit once the speculation path is confirmed valid (at retirement)



KPTI

- Kernel Page Table Isolation
- Removes the "Supervisor bit" optimization
- Changes CR3 at every syscall
- Modifies stack address at every syscall
- Userspace speculation thus can't see kernel data at all and can't access them
- TLB flush at every syscall → syscall entry 4x slowdown

Present bit



PRESENT
BIT

ACCESS TO NOT PRESENT
PAGES TRIGGERS A
PAGE FAULT



L1TF Attack!

- Also known as Foreshadow
- Technique similar to Meltdown
- Can read L1D contents only
- The kernel puts stuff like swap space block numbers into non-present PTEs
- And these are small numbers
- Userspace cannot control them
- And the kernel lives on negative addresses
- So the attack is not all that useful

L1TF Mitigation

- To avoid L1TF from userspace the kernel has to sanitize all non-present PTEs to point to non-existent memory instead just low addresses
- This is done easily by putting 1's in some top PTE address bits

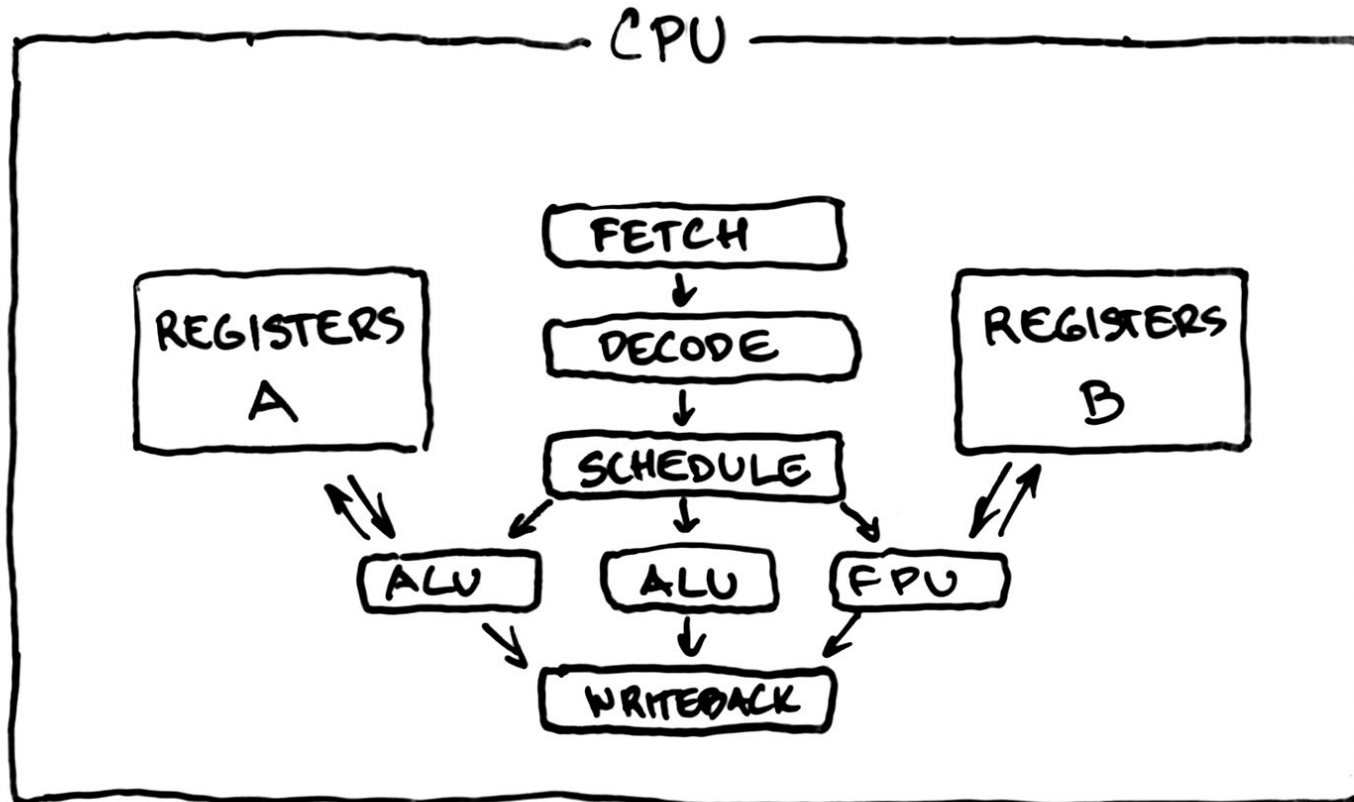
L1TF and Virtualization

- A Virtual machine is in control of page tables
- And can put ANYTHING in them
- And so it can read any address via L1D Cache
 - from other guests
 - and from the host
- And EPTs don't help (they're just ignored, too)

L1D Flushing

- To prevent L1TF from a VM, we must flush L1D before every VMENTER
- This hurts performance quite a bit

SMT / HT



VM L1TF and HyperThreading

- L1D is shared between hyperthreads
- If two VMs run on two hyperthreads, one can snoop the other
- If one hyperthread runs a VM and the other the host, the VM can snoop the host
- And there is nothing the host can do about that
→ NO MITIGATION
- The only way is to disable HT
→ 20-50% performance impact

Really nothing?

- We could disable EPT and do it in software
 - That's actually worse than disabling HT
- We could do ganged scheduling
 - That's hard and only solves the VM-VM scenario
- But it could be enough if the Host doesn't ever handle any valuable data
 - In a thin-hypervisor scenario, if no data pass through the hypervisor, the service domain (Dom0 in Xen, etc) is bound to a separate core and passes data via ring buffers to VMs, VMs never share cores, yes, then perhaps you can leave HT on.
 - Hyper-V seems to be able to do that, Xen might in the future, KVM probably never will.

