

# The Devastation of Meltdown



J.Nider

2019-11

Background  
Virtual Memory  
OoO Execution  
Linux Memory Management  
The Exploit  
The Fix  
The Damage

# Background

- Meltdown allows an unprivileged app to read ALL memory of a victim machine
- Official name: CVE-2017-5754 “Rogue Data Cache Load” (RDCL)
- Caused by a race condition in out-of-order CPU’s
- NSA potentially knew about this since 1995

# Scope

- Affects almost all Intel processors [1990-2018]
  - IBM POWER7,8 Z
  - ARMv8 A-series
  - AMD is not vulnerable!!
- 
- All operating systems are affected (Linux, Windows, Android, etc)
  - Containers are affected (LXC, Docker, OpenVZ, etc)
  - Hardware-supported virtual machines are not (KVM, VMWare ESX, etc)

Background

**Virtual Memory**

OoO Execution

Linux Memory Management

The Exploit

The Fix

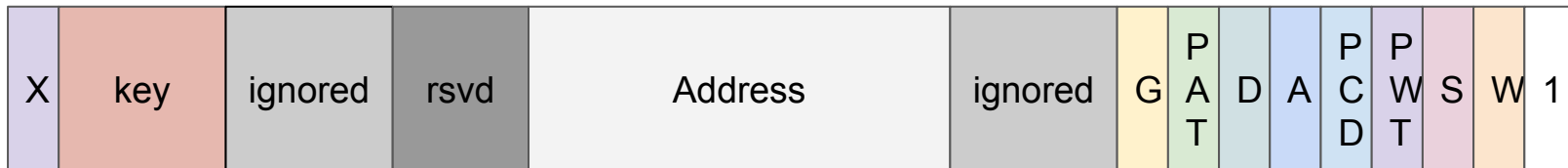
The Damage

# Virtual Memory

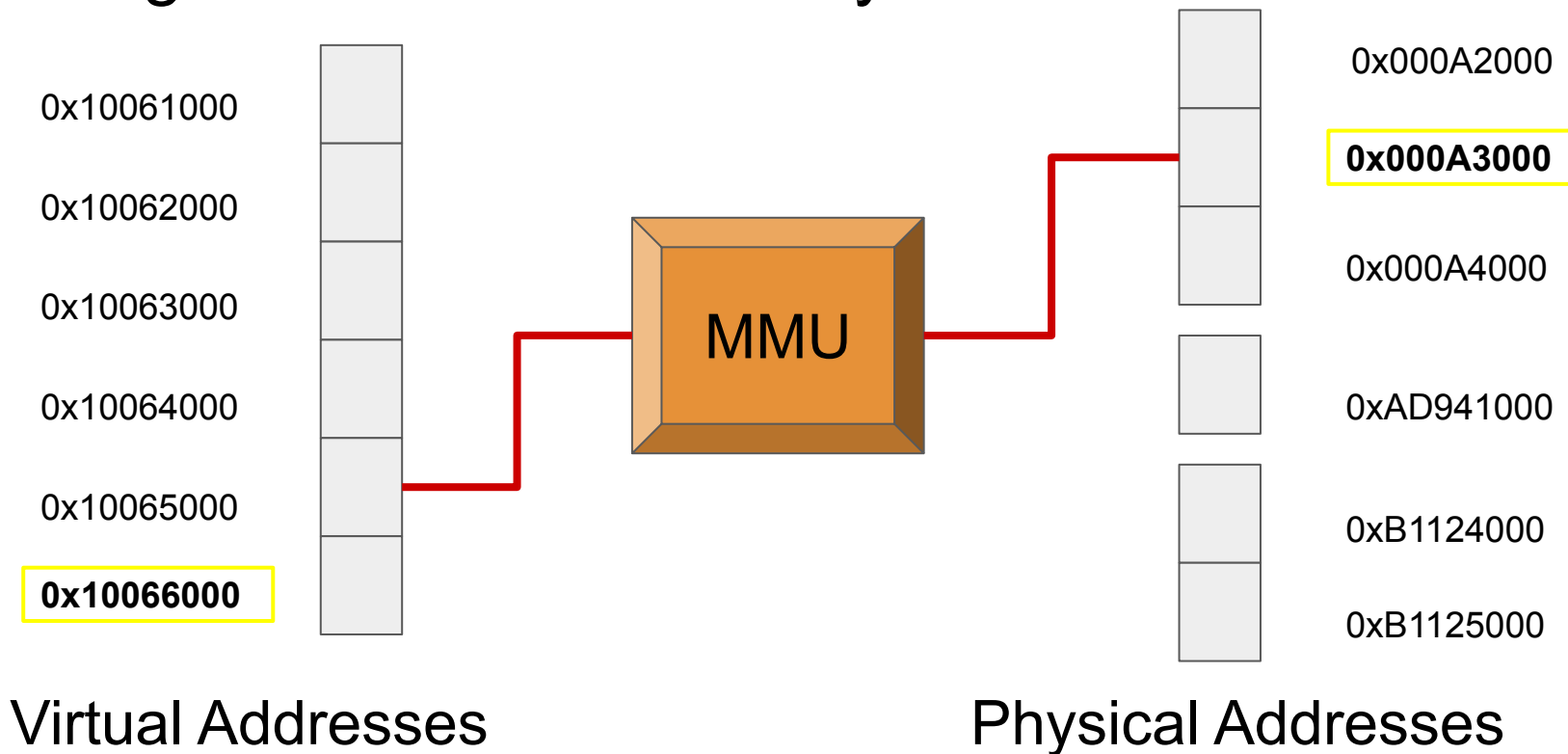
- Memory is organized into pages
  - Page sizes range from 4KB to up 1GB
- MMU maps virtual address to physical addresses
  - Usually through page tables, but there are other mechanisms (such as hashing)
- Each page has attributes
  - Describes permissions (RWX), (S) and caching (C)
- Translations are cached in TLB (translation lookaside buffer)

63

0



# Background: Virtual Memory

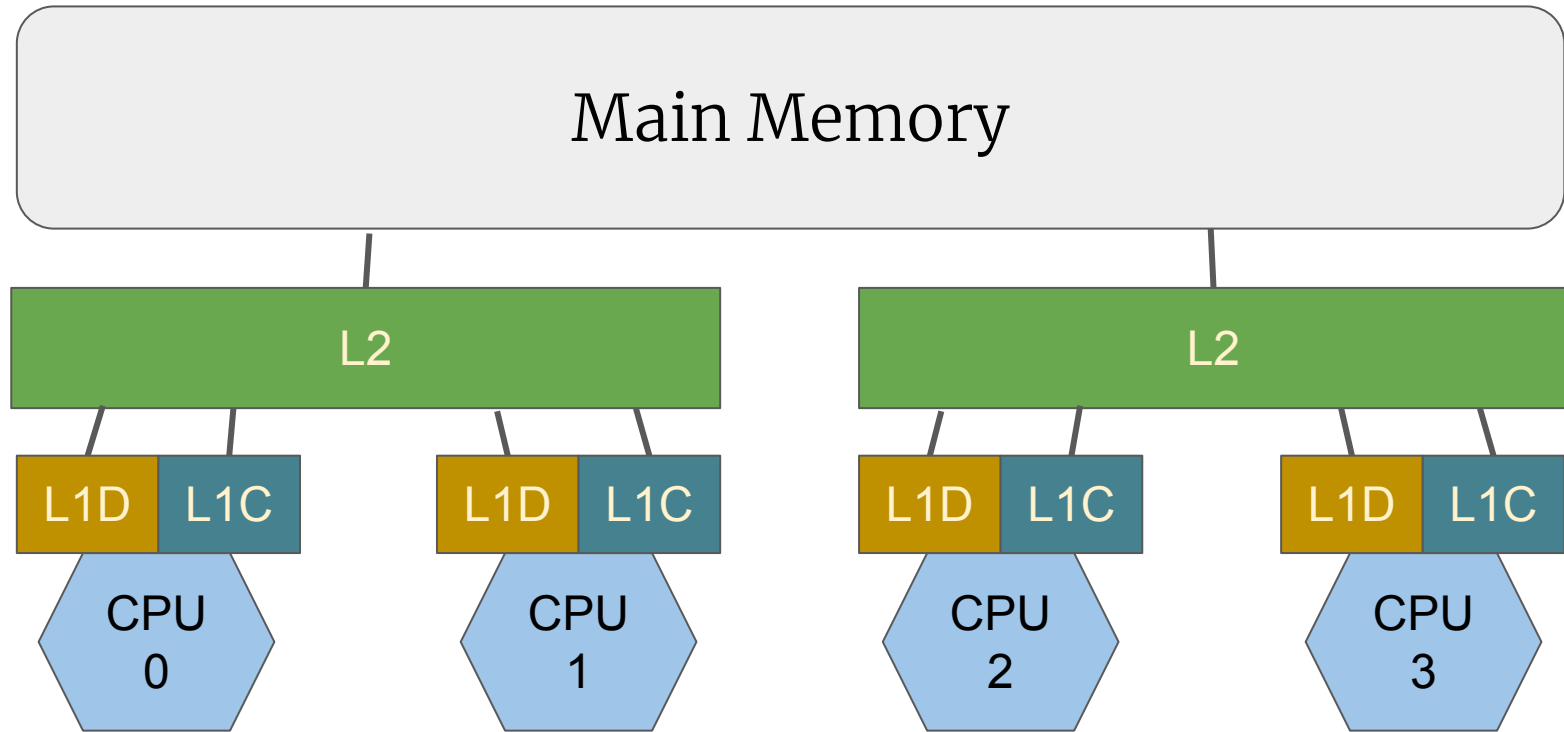


# Cache Organization

- Reading from main memory is slow!
  - In the range of 400-800ns
  - Therefore we want to avoid main memory as much as possible
- So we cache (make a copy) of any data in a smaller, faster memory
  - Much faster - in the range of 10-100ns
  - Faster memories are more expensive
- We can make a hierarchy with different attributes
  - Capacity
  - Access time
  - Mapping (Direct, Associativity)
  - Multiple ways
- **Cache is not part of the Instruction Set Architecture!**
  - **It is part of the microarchitecture**



# Cache Organization



Background

Virtual Memory

**OoO Execution**

Linux Memory Management

The Exploit

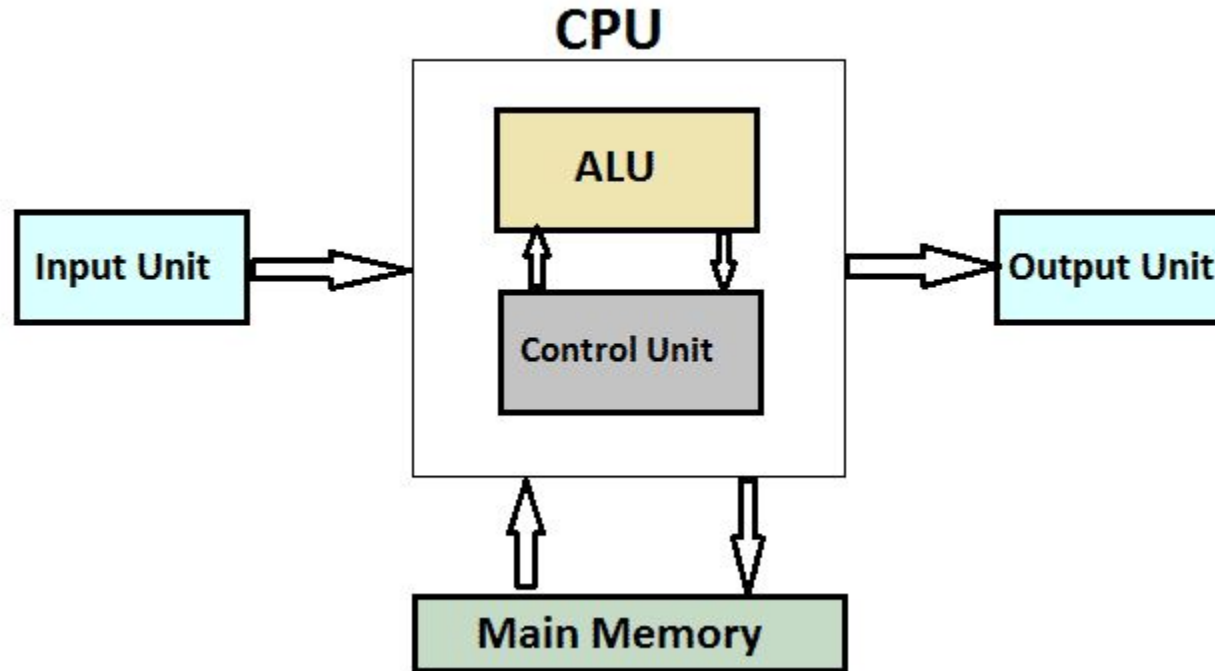
The Fix

The Damage

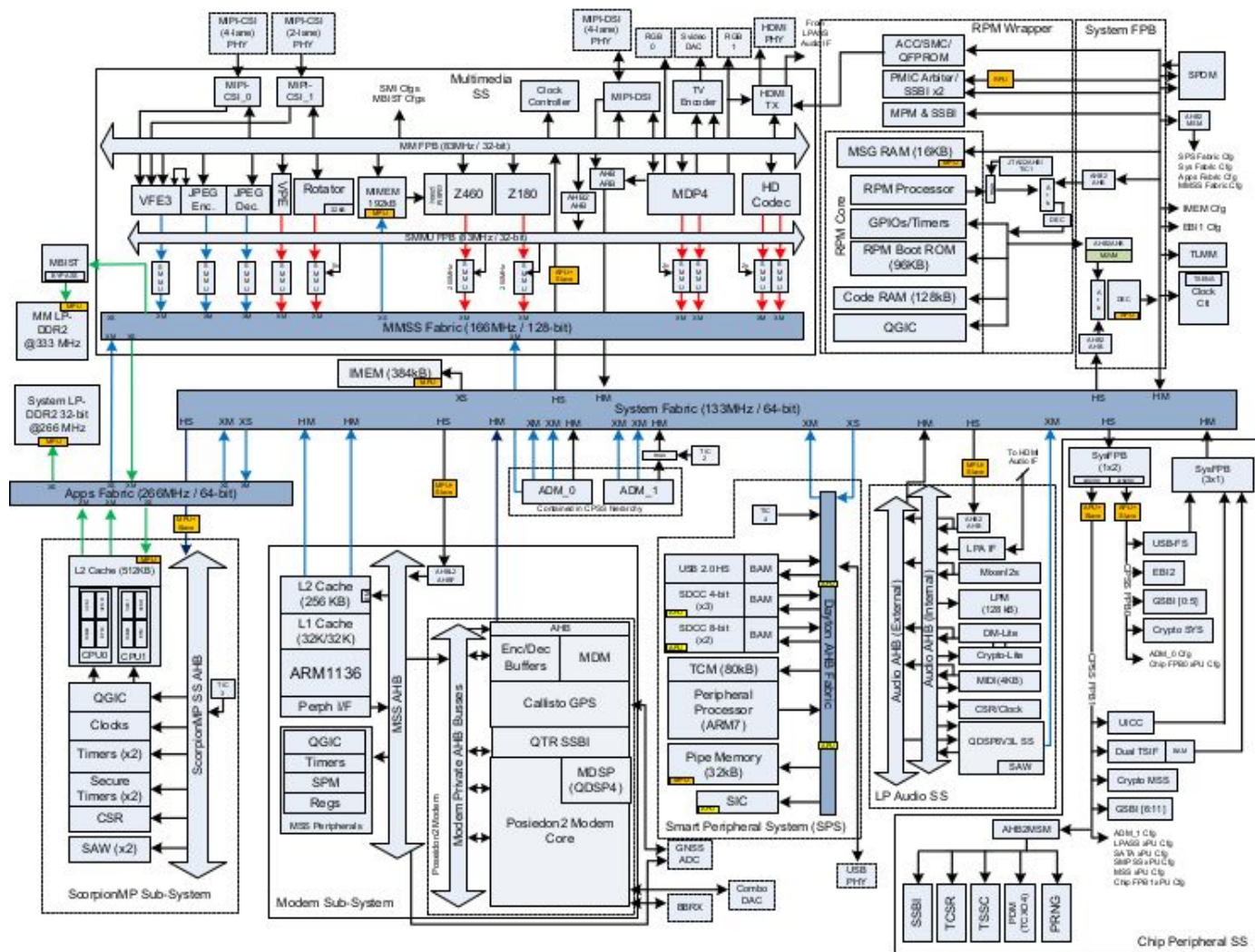
# Out-of-order Execution

- CPUs are made up of many hardware blocks
  - Integer units
  - Floating point units
  - Internal registers
  - Many more
- Not all hardware is used for each instruction
- Some instructions wait even though there are not dependencies
- We want to work as fast as possible

# CPU Architecture - Simplified

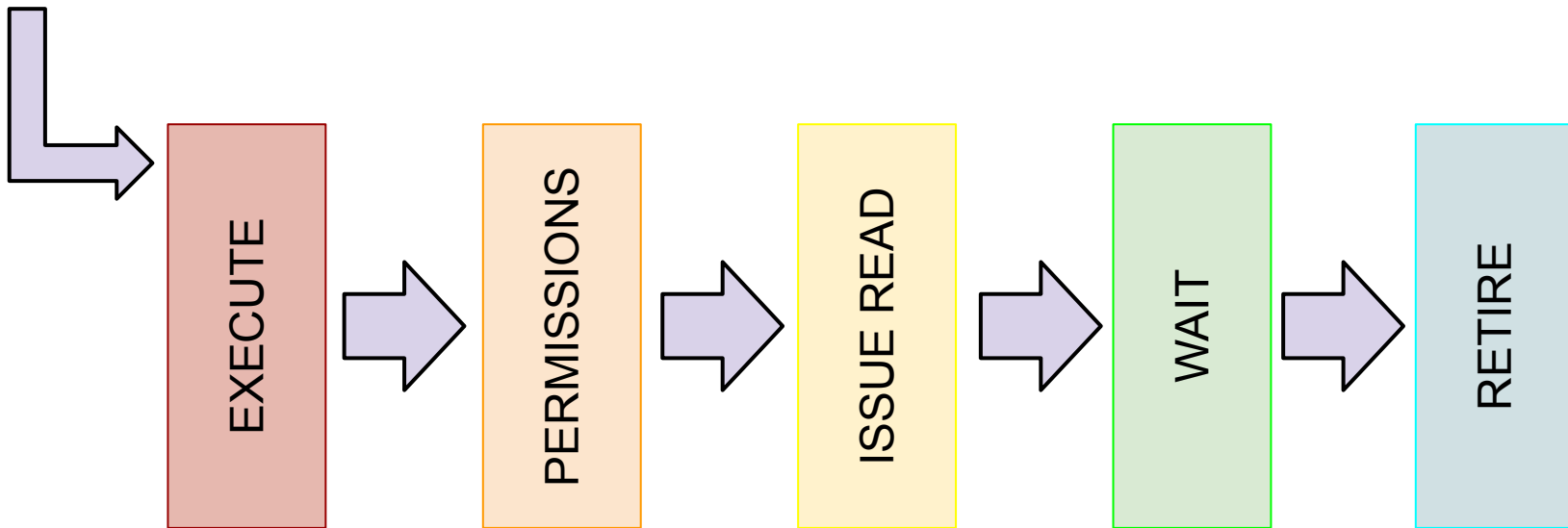


# Real Life

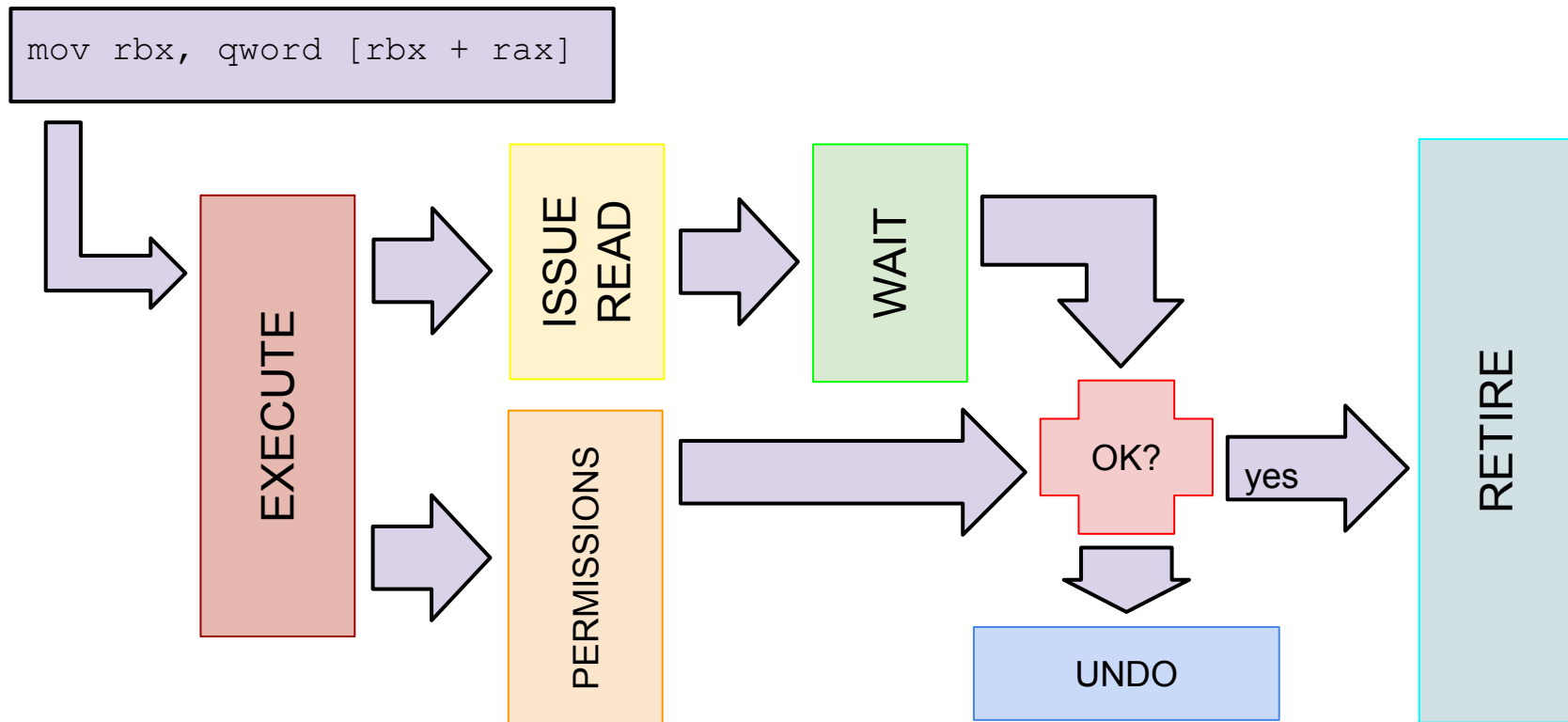


# Memory Read (sequential)

```
mov rbx, qword [rbx + rax]
```



# Memory Read (parallel)



Background

Virtual Memory

CPU Pipelining

**Linux Memory Management**

The Exploit

The Fix

The Damage

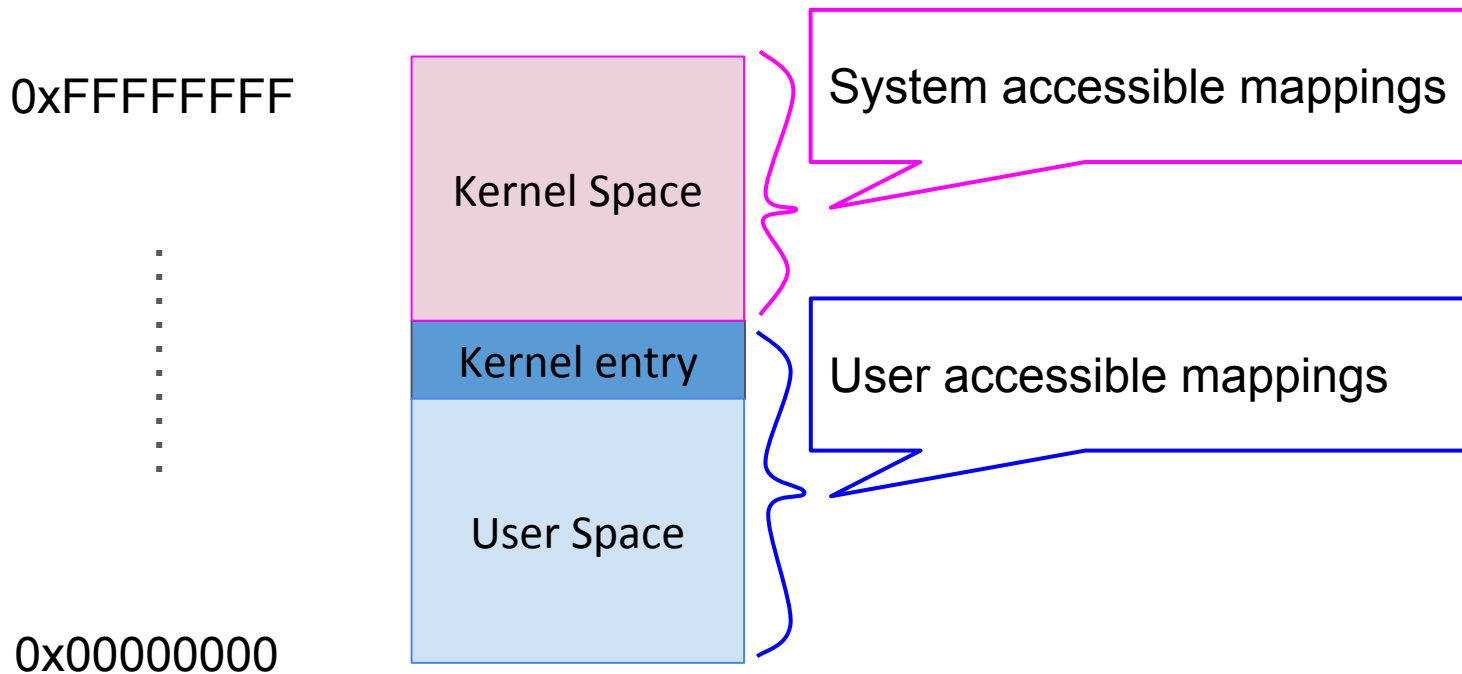


# Linux Memory Management

- The operating system uses the hardware to give us an abstraction
  - We can make assumptions without understanding the details of the hardware
- Makes sure we can get physical memory when we need it
- Use the hardware to protect our memory from other programs

Start	End	Size	Description
0000000000000000	00007fffffffffffffff	128 TB	user-space virtual memory
0000800000000000	ffff7fffffffffffffff	~16 EB	empty
ffff800000000000	ffffffffffffffffffff	128 TB	Kernel-space virtual memory

# Linux Virtual Address Space Layout (w/o KPTI)



# Linux Direct Map

- All physical memory is directly mapped in kernel virtual memory space
- Basis for `phys2virt` and `virt2phys` macros
- Used primarily for drivers and 'mm' functions
- This makes memory manipulation code small, fast and efficient
- This is also a big security risk!

Start	End	Size	Description
ffff888000000000	ffffc87ffffffffff	64 TB	Direct Map

Background  
Virtual Memory  
CPU Pipelining  
Linux Memory Management  
**The Exploit**  
The Fix  
The Damage

# The Exploit

- We want to read kernel memory - how?

## **Two conditions must hold**

1. Mapping of physical page in our virtual address space
2. Permission bit to allow unprivileged access to page

# The Exploit

- We want to read kernel memory - how?

## Two conditions must hold

1. Mapping of physical page in our virtual address space ✓
2. Permission bit to allow unprivileged access to page

# The Exploit

- We want to read kernel memory - how?

## Two conditions must hold

1. Mapping of physical page in our virtual address space ✓
2. Permission bit to allow unprivileged access to page ✗

# The Exploit in C

We could do this in C, but it is clearer in assembly language

```
unsigned long rax = 0;
char probe[4096 * 256];    the probe array
unsigned long rcx = 0xffff800000000000;  pointer to a kernel address
char val;
while (!rax) {
    rax = *(byte*)rcx;      (no permission!)
    rax <<= 12;             shift the secret value by the page size
    val = probe[rax];       secret value becomes index into probe array
}
```



# The Exploit

```
xor rax, rax
```

```
retry:
```

```
mov al, byte [rcx]
```

rcx is a pointer to a kernel address

(no permission!)

```
shl rax, 0xc
```

shift the secret value by the page size

secret value becomes index into probe array

```
jz retry
```

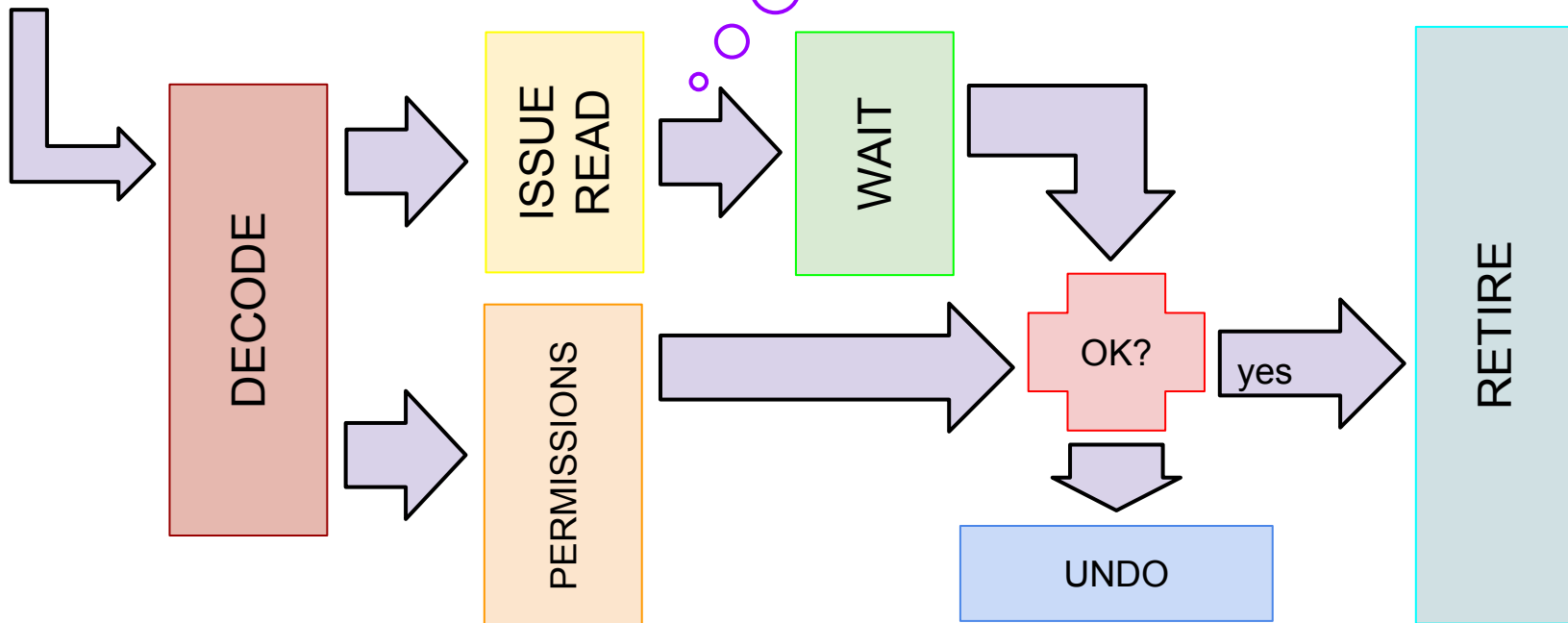
```
mov rbx, qword [rbx + rax]
```

rbx is the base address of the probe array

# The Exploit

```
mov rbx, qword [rbx + rax]
```

There is a side effect!



# The Exploit

```
xor rax, rax
```

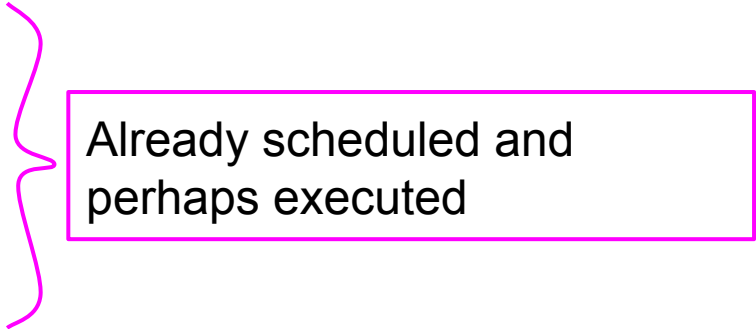
```
retry:
```

```
mov al, byte [rcx] Exception!
```

```
shl rax, 0xc
```

```
jz retry
```

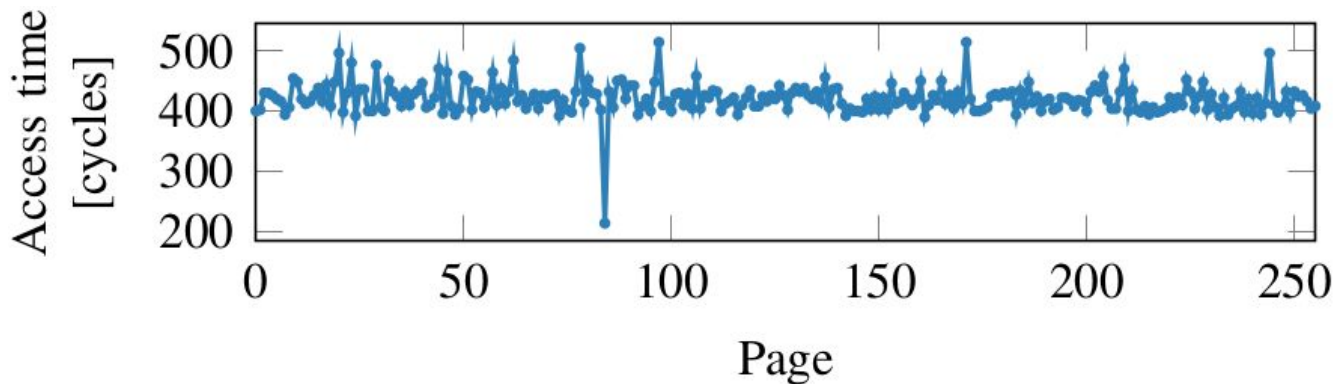
```
mov rbx, qword [rbx + rax]
```



Already scheduled and  
perhaps executed

# Flush + Reload

- Make sure the cache is empty (clflush)
- Perform attack
- Read all entries in the probe array, and measure access time
- One measurement might stand out!
- Index of cached page is the value of the secret byte



# Accessing All Memory

- Now we know how to access kernel memory!
  - Not very fast, but it works
- But how to access memory of another process?
  - Linux manages all processes (including their hierarchy) in a linked list
  - The head of this task list is stored in the `init_task` structure
- Use the direct memory map
  - Must find the page tables belonging to another process
  - Perform a page walk to find the physical page for a particular virtual address
  - Access that physical page through the direct map

# Performance

- Flush-Reload is the bottleneck of the attack
- Instead of 8 bits (=256 entries), send 1 bit (=2 entries) of information
  - Much faster
  - Less reliable (noise bias to '0')
- Can read memory at rates between 4KB/s - 500KB/s

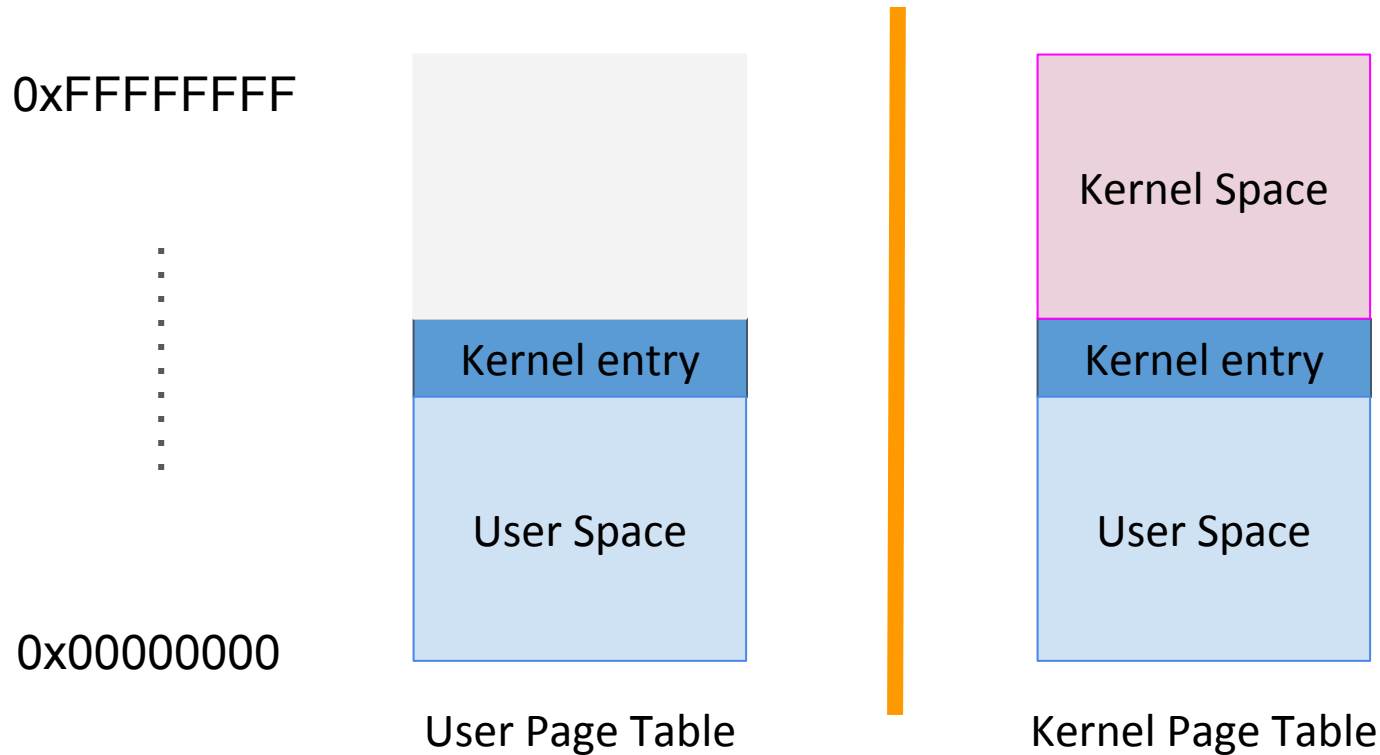
Background  
Virtual Memory  
CPU Pipelining  
Linux Memory Management  
The Exploit  
**The Fix**  
The Damage

# The Fix

- KPTI - Kernel Page Table Isolation
- Based on KAISER patches
- Removes kernel mappings from user process virtual memory
- Requires a pair of page tables for each process
  - One for user space
  - One for kernel space
- Drastically increases overhead during context switch



# Linux Virtual Address Space Layout (with KPTI)

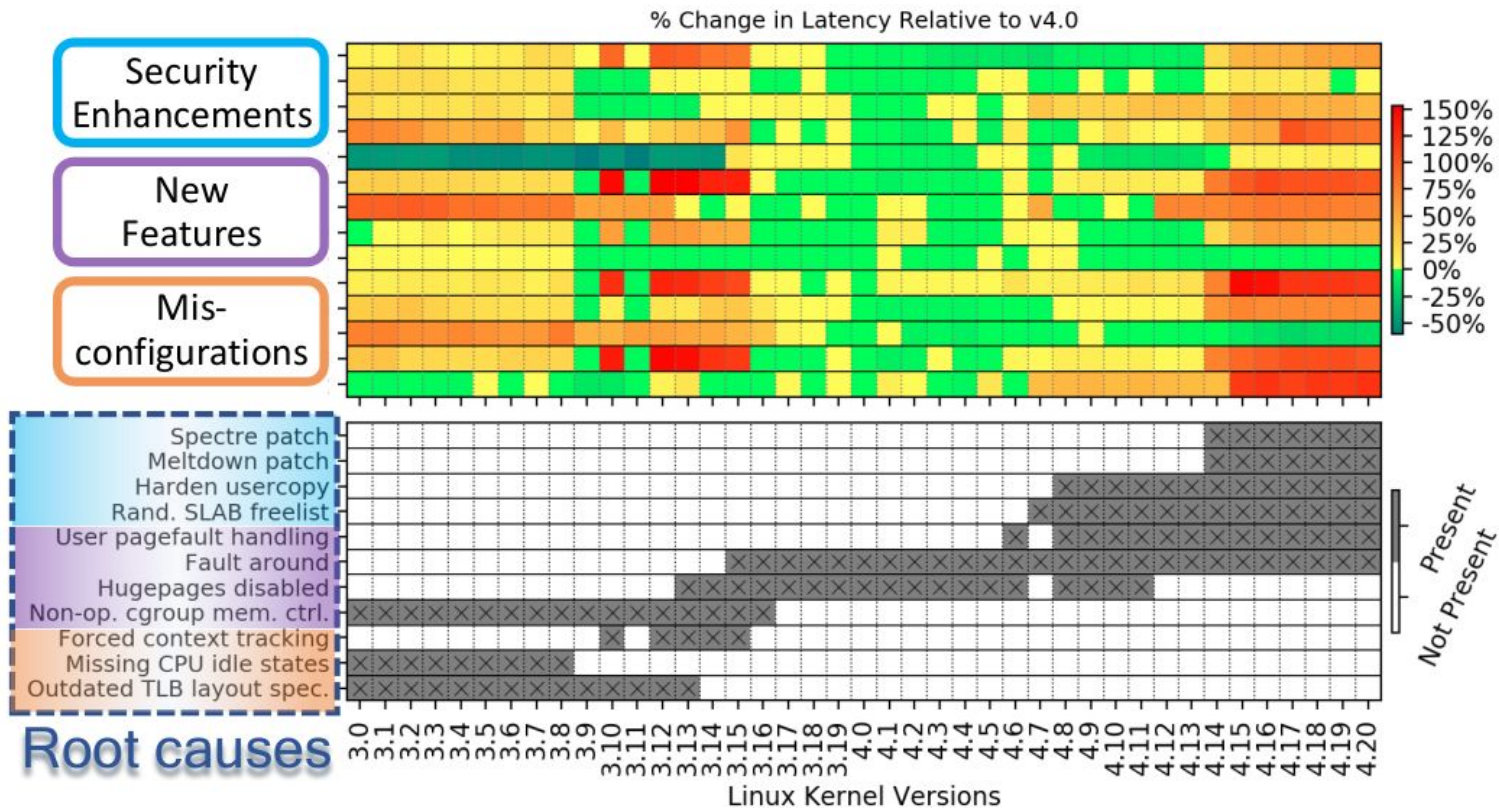


Background  
Virtual Memory  
CPU Pipelining  
Linux Memory Management  
The Exploit  
The Fix  
**The Damage**

# The Damage

- Measurements are very dependent on the number of syscalls
- The overhead was measured to be 0.28% according to KAISER's original authors
- a Linux developer measured it to be roughly 5% for most workloads and up to 30% in some cases
- for database engine PostgreSQL the impact on read-only tests on an Intel Skylake processor was 16–23% (without PCID)
- Redis slowed by 6–7%
- Linux kernel compilation slowed down by 5% on Haswell

# The Damage



# Making It Hurt Slightly Less

- PCIDs allow a logical processor to cache information for multiple linear-address spaces
  - Allows us to bypass the TLB flush on syscall entry/exit
- 
- PostgreSQL read-only tests on an Intel Skylake processor was 7–17% (or 16–23% without PCID)

# Conclusions

- Even the most commonly used, professionally made chips have bugs
- Operating systems can be used to mask these bugs
- Even so, the bugs are costly!

Meltdown

Spectre

L1TF

RIDL

Fallout

More??

# References

<https://sosp19.rcs.uwaterloo.ca/slides/ren.pdf>

<https://meltdownattack.com/meltdown.pdf>

[https://www.kernel.org/doc/Documentation/x86/x86\\_64/mm.txt](https://www.kernel.org/doc/Documentation/x86/x86_64/mm.txt)

[https://en.wikipedia.org/wiki/Kernel\\_page-table\\_isolation](https://en.wikipedia.org/wiki/Kernel_page-table_isolation)