



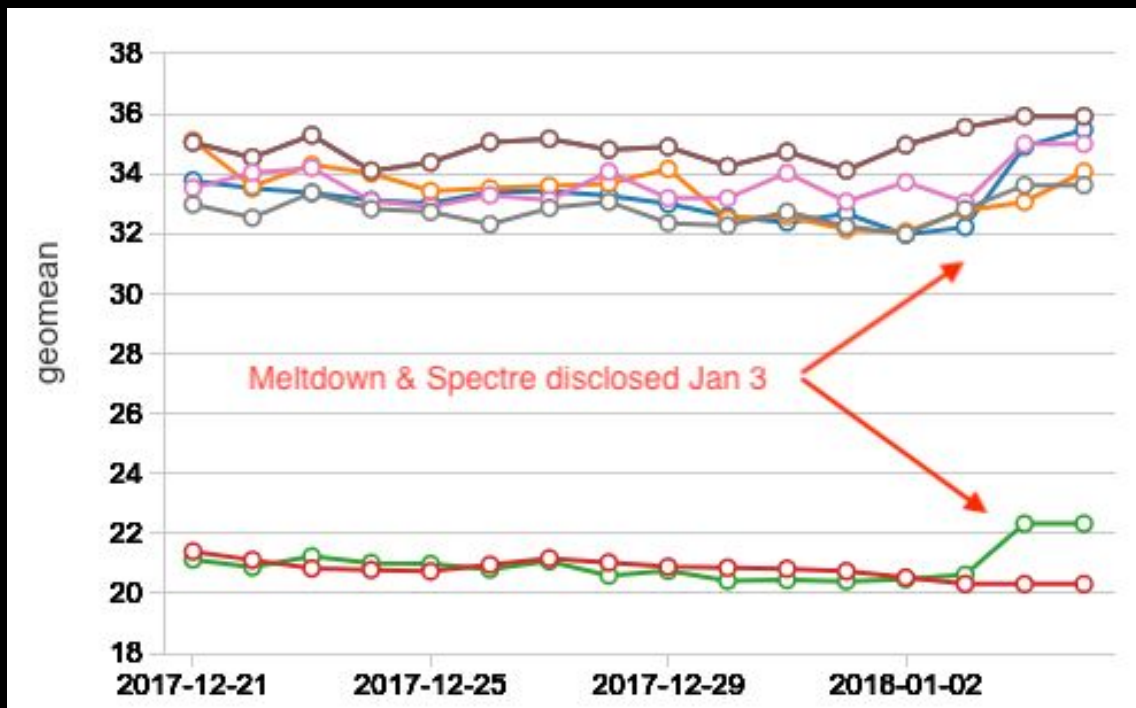
# Meltdown, Spectre and Apache Spark™ Performance

Chris Stevens

June 5, 2018



# Databricks Performance on AWS



3-5% performance degradation

# Overview

Goal: Understand the 3-5% degradation

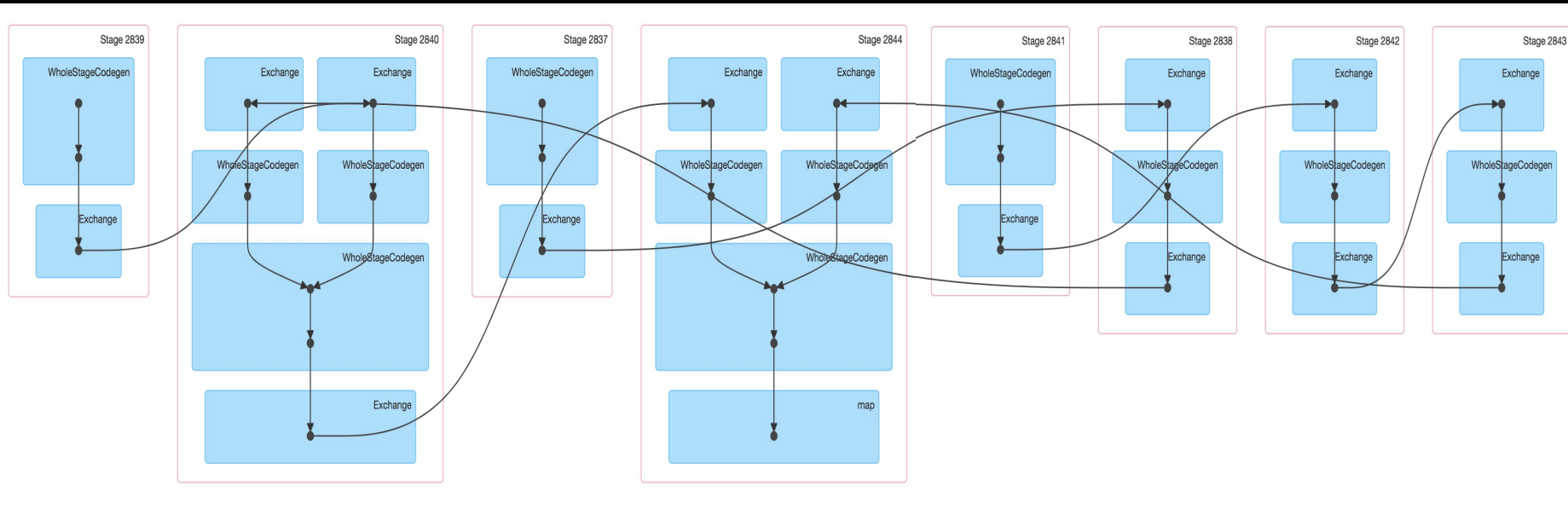
Steps:

- TPC-DS Benchmarks
- System Analysis
- Breakdown the exploits and patches
  - Meltdown
  - Spectre V1 - Bounds Check Bypass
  - Specter V2 - Branch Target Injection

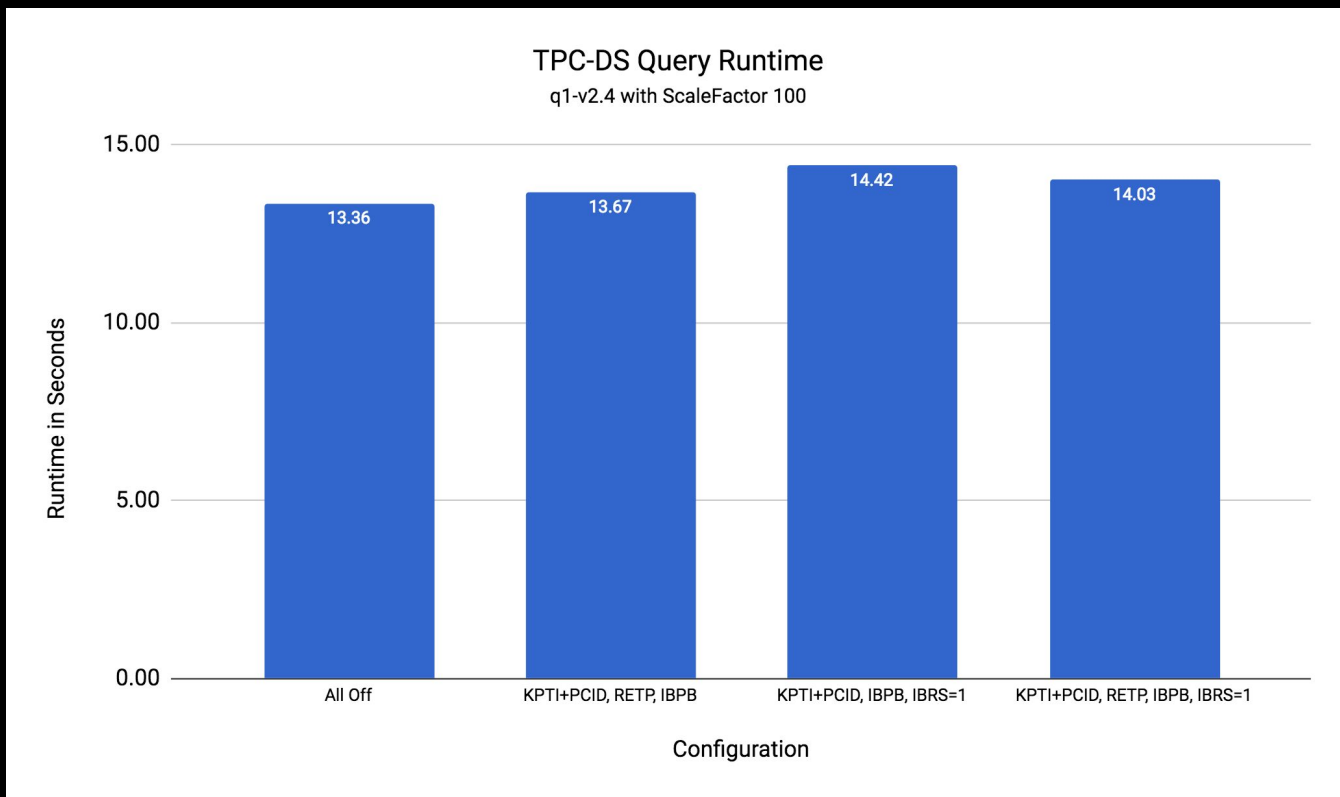
# TPC-DS: q1-v2.4

```
WITH customer_total_return AS
  (SELECT sr_customer_sk AS ctr_customer_sk,
          sr_store_sk AS ctr_store_sk,
          sum(sr_return_amt) AS ctr_total_return
   FROM store_returns, date_dim
   WHERE sr_returned_date_sk = d_date_sk
   AND d_year = 2000
   GROUP BY sr_customer_sk, sr_store_sk)
SELECT c_customer_id
  FROM customer_total_return ctr1, store, customer
 WHERE ctr1.ctr_total_return >
       (SELECT avg(ctr_total_return)*1.2
        FROM customer_total_return ctr2
        WHERE ctr1.ctr_store_sk = ctr2.ctr_store_sk)
 AND s_store_sk = ctr1.ctr_store_sk
 AND s_state = 'TN'
 AND ctr1.ctr_customer_sk = c_customer_sk
 ORDER BY c_customer_id LIMIT 100
```

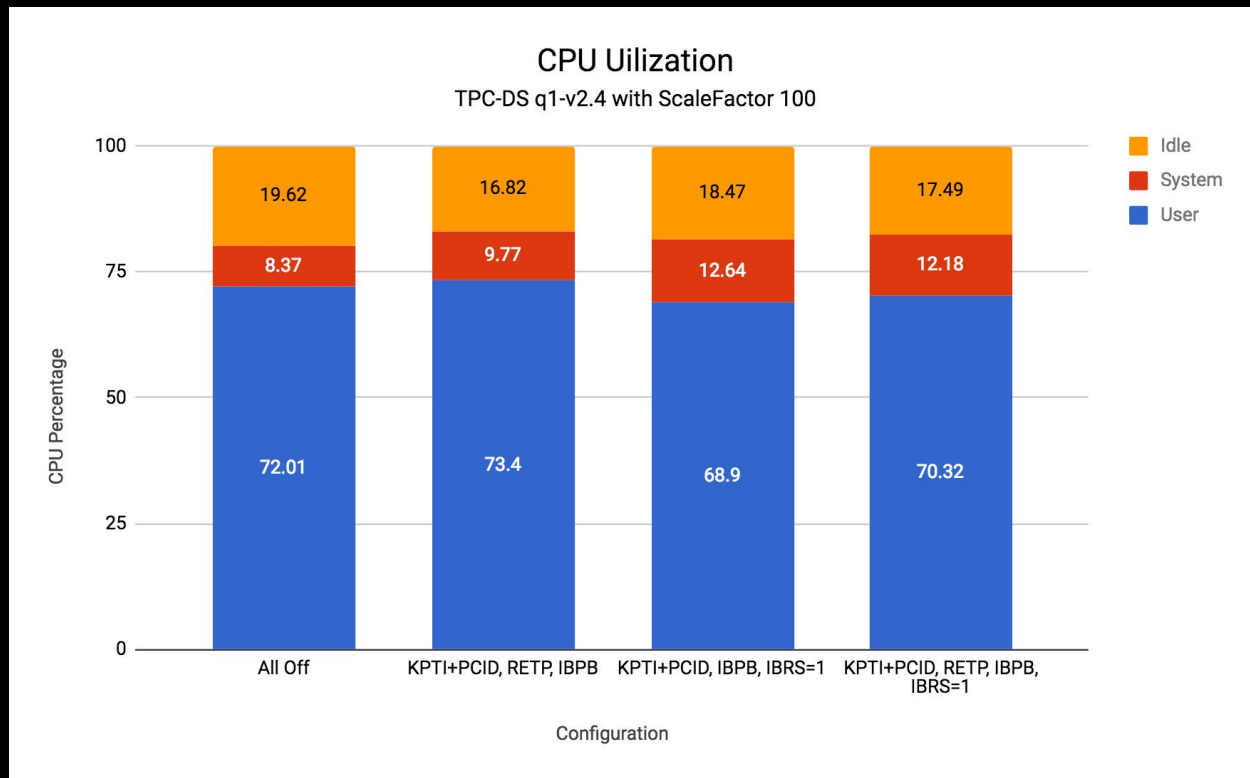
# TPC-DS: q1-v2.4



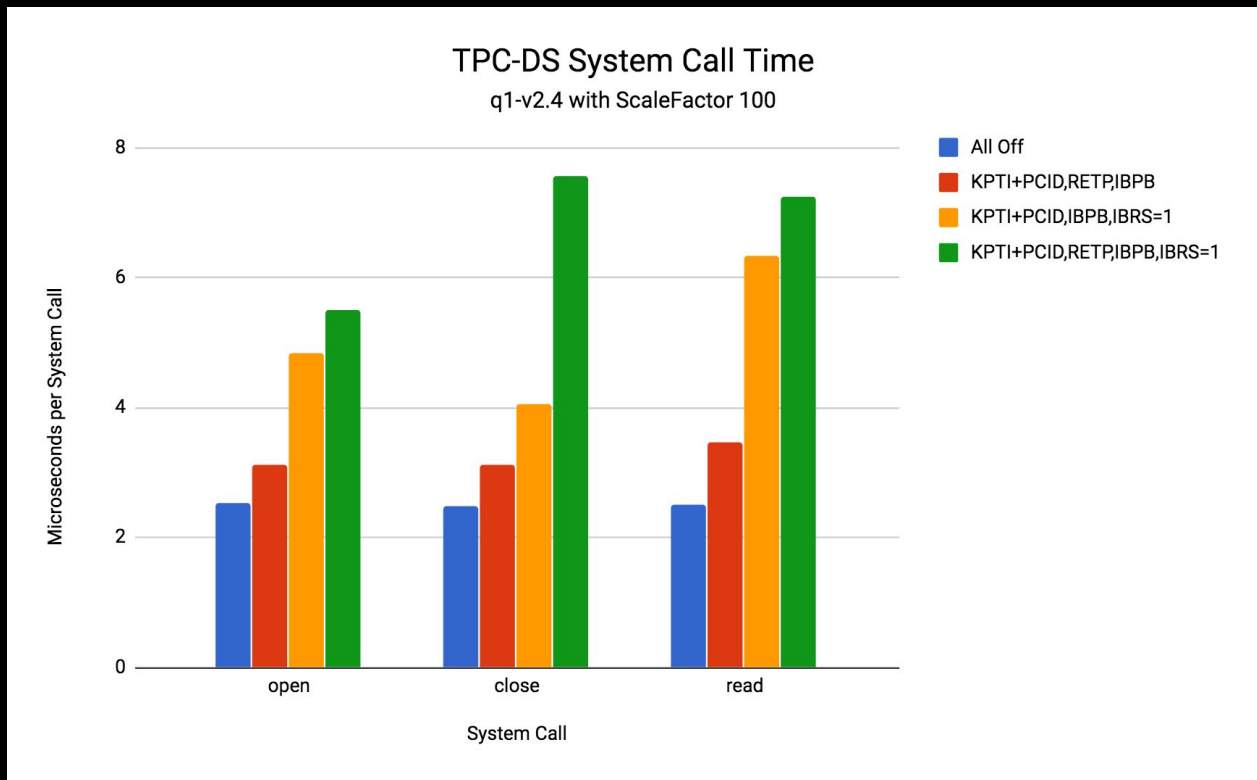
# TPC-DS: q1-v2.4



# q1-v2.4 CPU Utilization



# TPC-DS: q1-v2.4







# Exploits Background

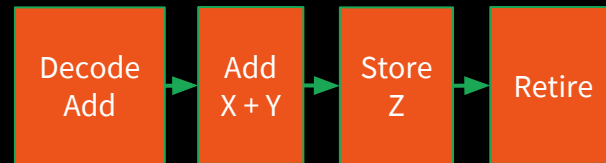
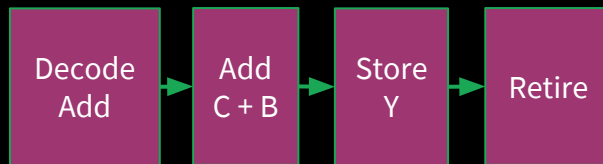
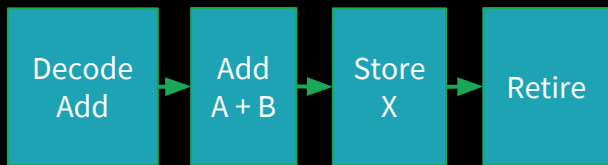
Out-of-Order Execution + Side Channel Attacks

# In-order Execution

$x = a + b$

$y = c + d$

$z = x + y$



Time

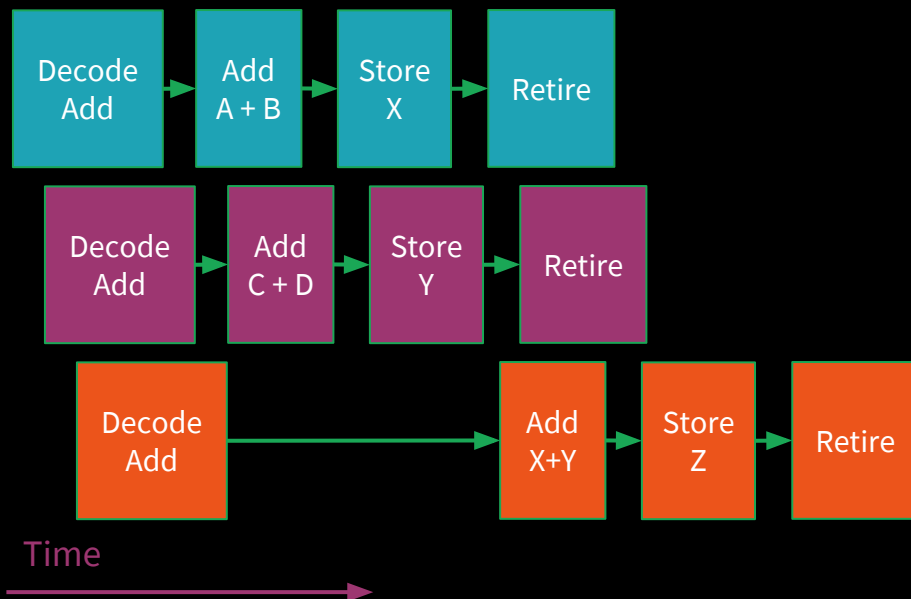


# Out-of-order Execution

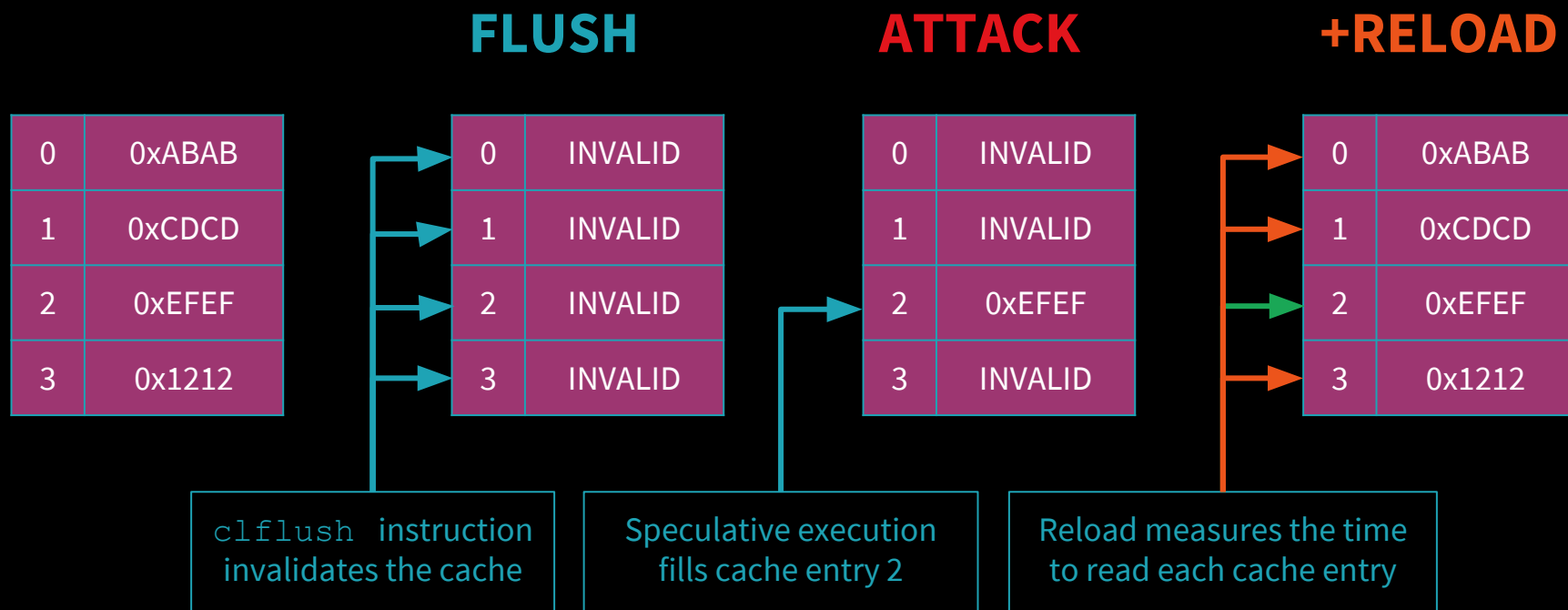
$x = a + b$

$y = c + d$

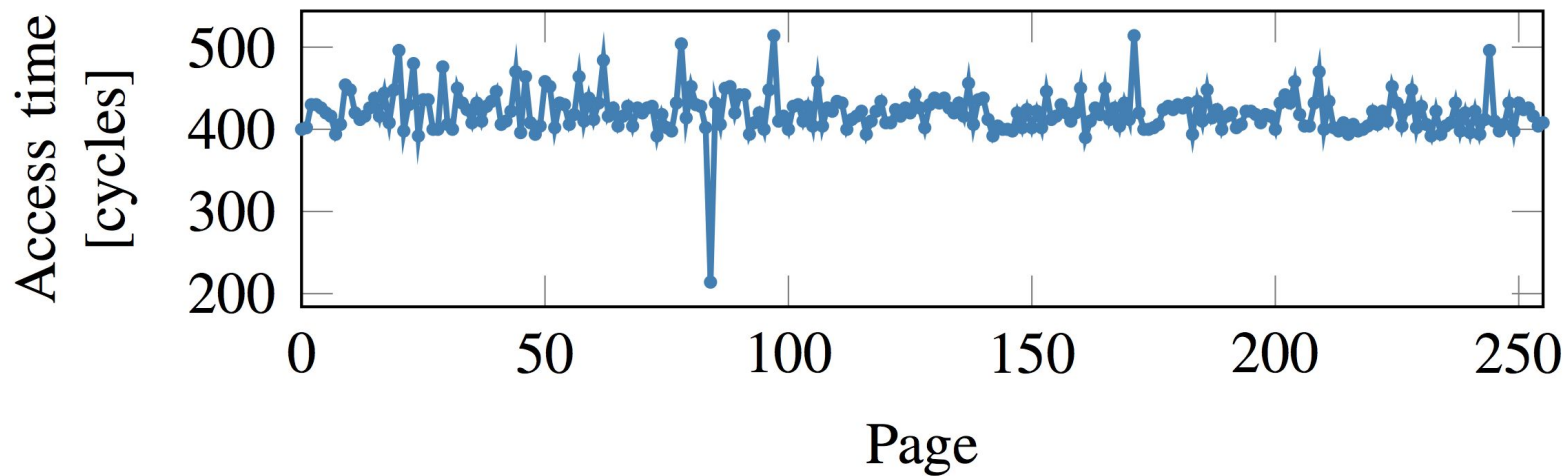
$z = x + y$



# Side-Channel Attacks



# Side-Channel Attacks



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



# Meltdown

# Meltdown

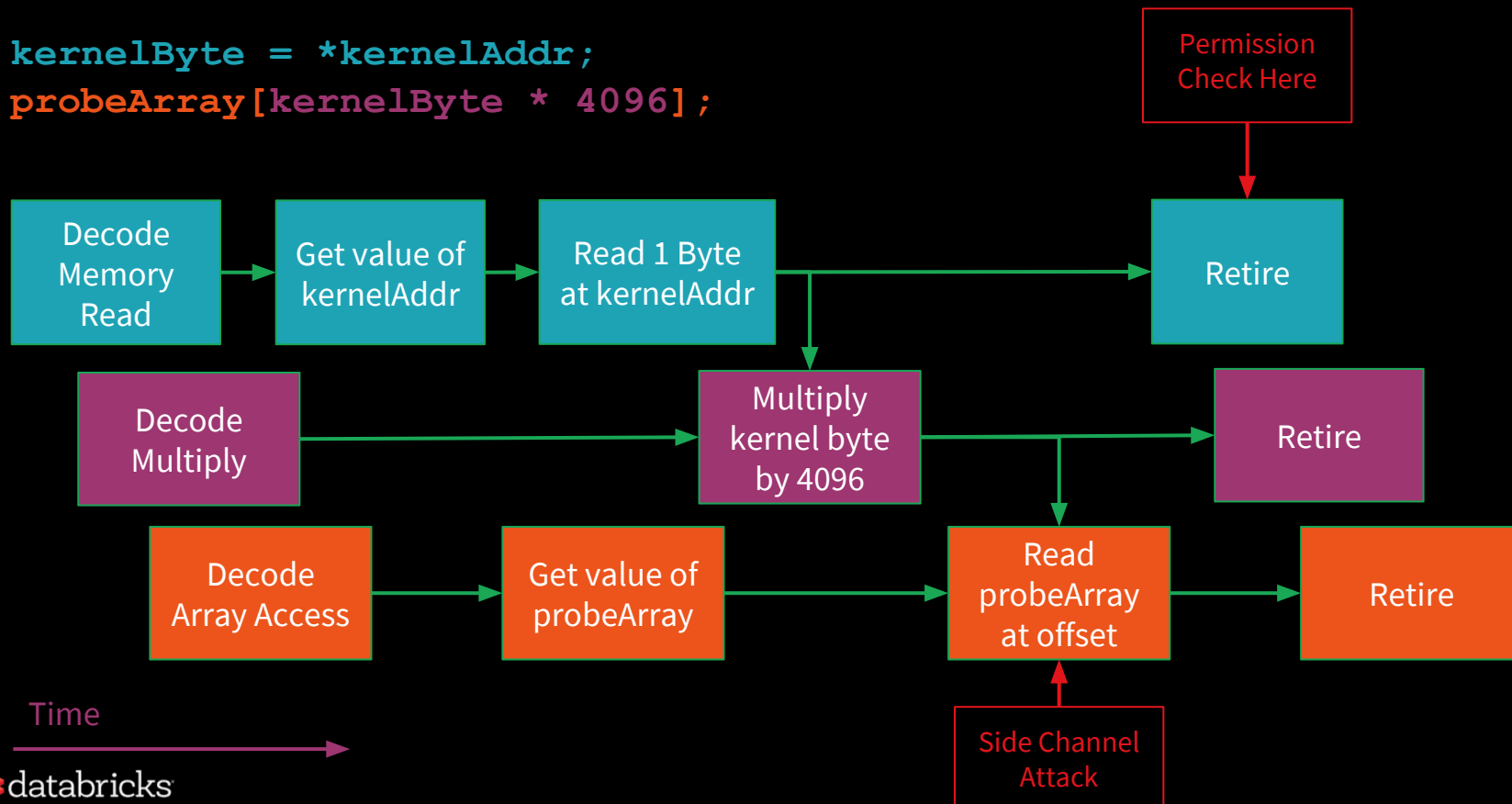
```
kernelByte = *kernelAddr;  
probeArray[kernelByte * 4096];
```



# Meltdown



```
kernelByte = *kernelAddr;  
probeArray[kernelByte * 4096];
```





# Side-Channel Attack



## FLUSH

0*4096	INVALID
1*4096	INVALID
2*4096	INVALID
3*4096	INVALID

Flush each page in  
probeArray

## ATTACK

0*4096	INVALID
1*4096	INVALID
2*4096	0xEFEF
3*4096	INVALID

Speculative execution  
reads probeArray at  
kernelByte \* 4096

## +RELOAD

0*4096	0xABAB
1*4096	0xCDCD
2*4096	0xEFEF
3*4096	0x1212

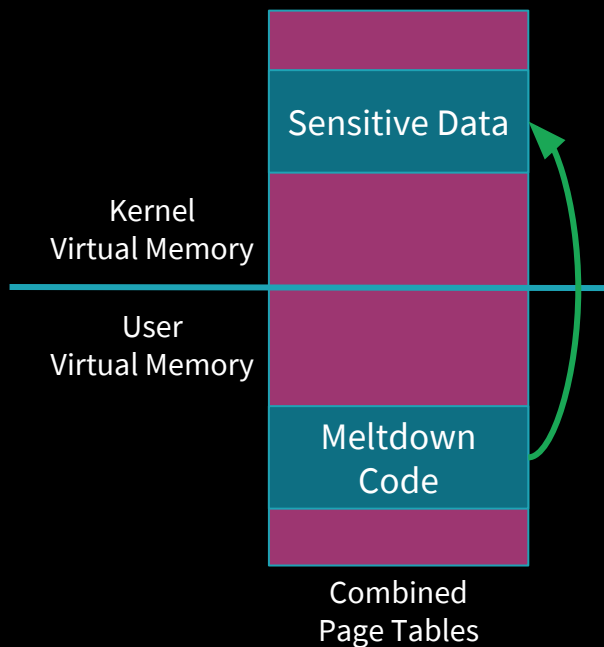
Reload measures the time  
to read each page in  
probeArray

Sees that page 2  
was the fastest =>  
kernelByte = 0x2

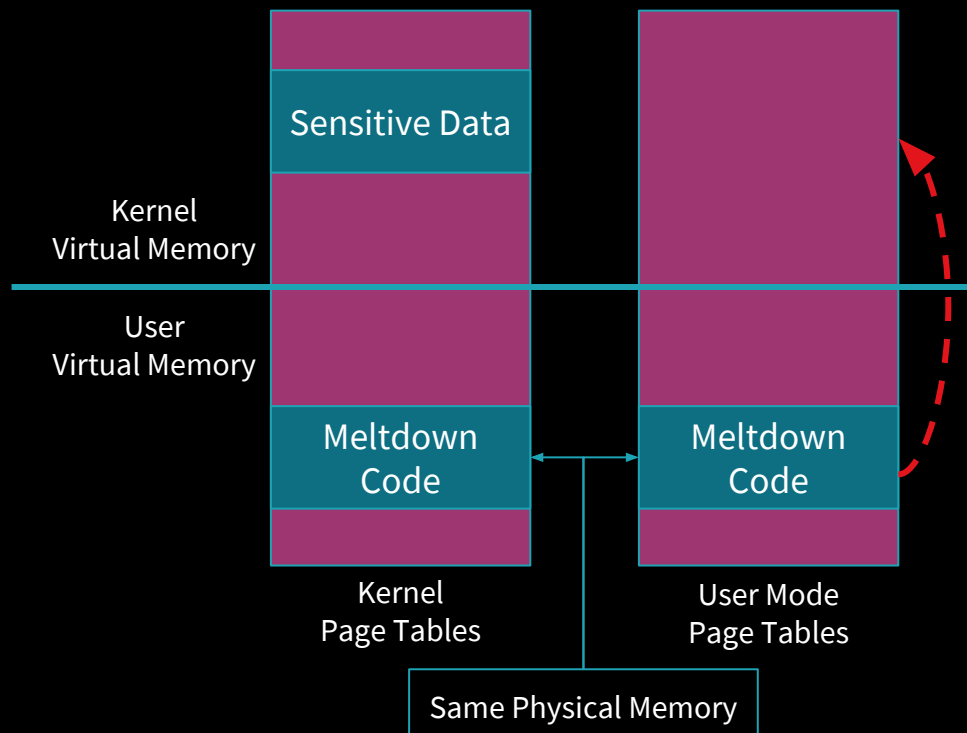
# Kernel Page-Table Isolation



Process Page Tables without KPTI



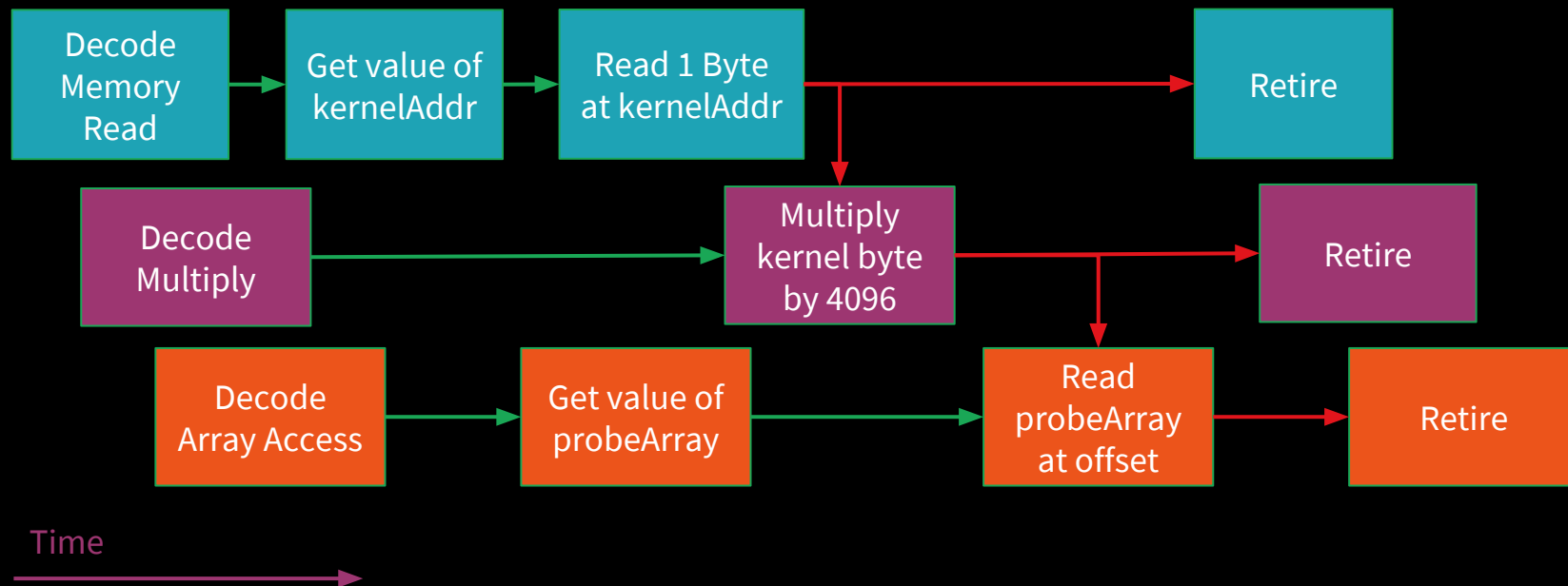
Process Page Tables with KPTI



# Meltdown



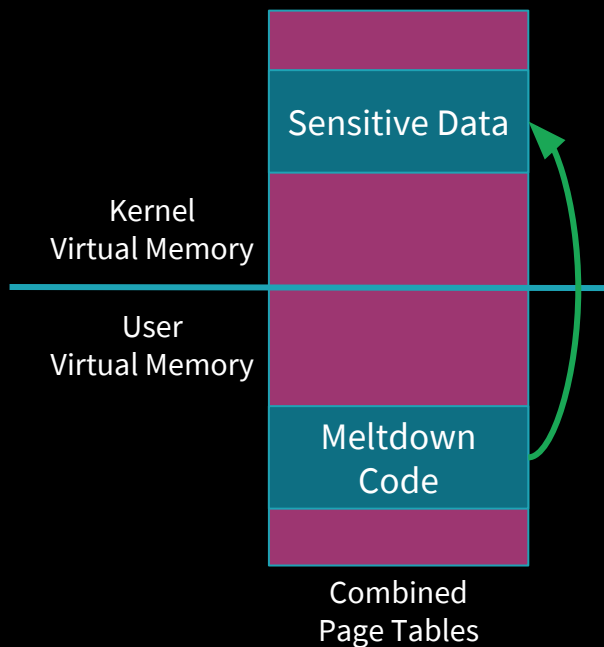
```
kernelByte = *kernelAddr;  
probeArray[kernelByte * 4096];
```



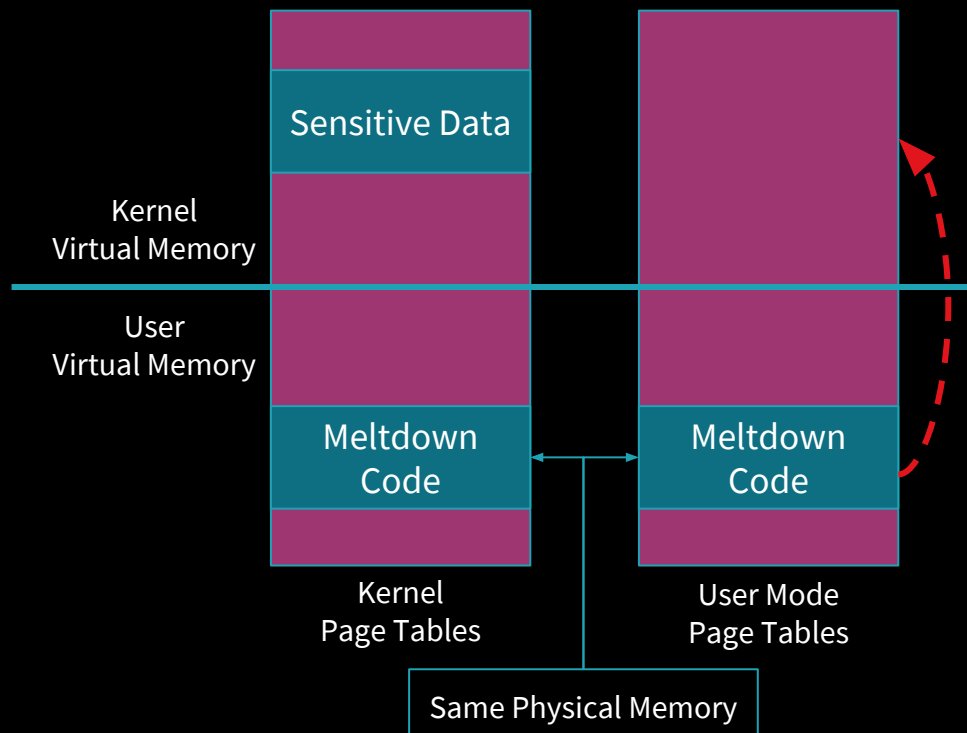
# Kernel Page-Table Isolation



Process Page Tables without KPTI



Process Page Tables with KPTI



# TLB before KPTI



```
while True:
    print *UserAddress
    sys_clock_gettime
```

## TLB MISS

Virtual Address	Physical Address
-	-
-	-

print \*UserAddress

## TLB MISS

Virtual Address	Physical Address
UserAddress	Page 1
-	-

sys\_clock\_gettime

## TLB HIT

Virtual Address	Physical Address
UserAddress	Page 1
KernelTime	Page 2

print \*UserAddress

# TLB with KPTI



```
while True:
    print *UserAddress
    sys_clock_gettime
```

## TLB MISS

Virtual Address	Physical Address
-	-
-	-

print \*UserAddress

## TLB MISS

Virtual Address	Physical Address
-	-
-	-

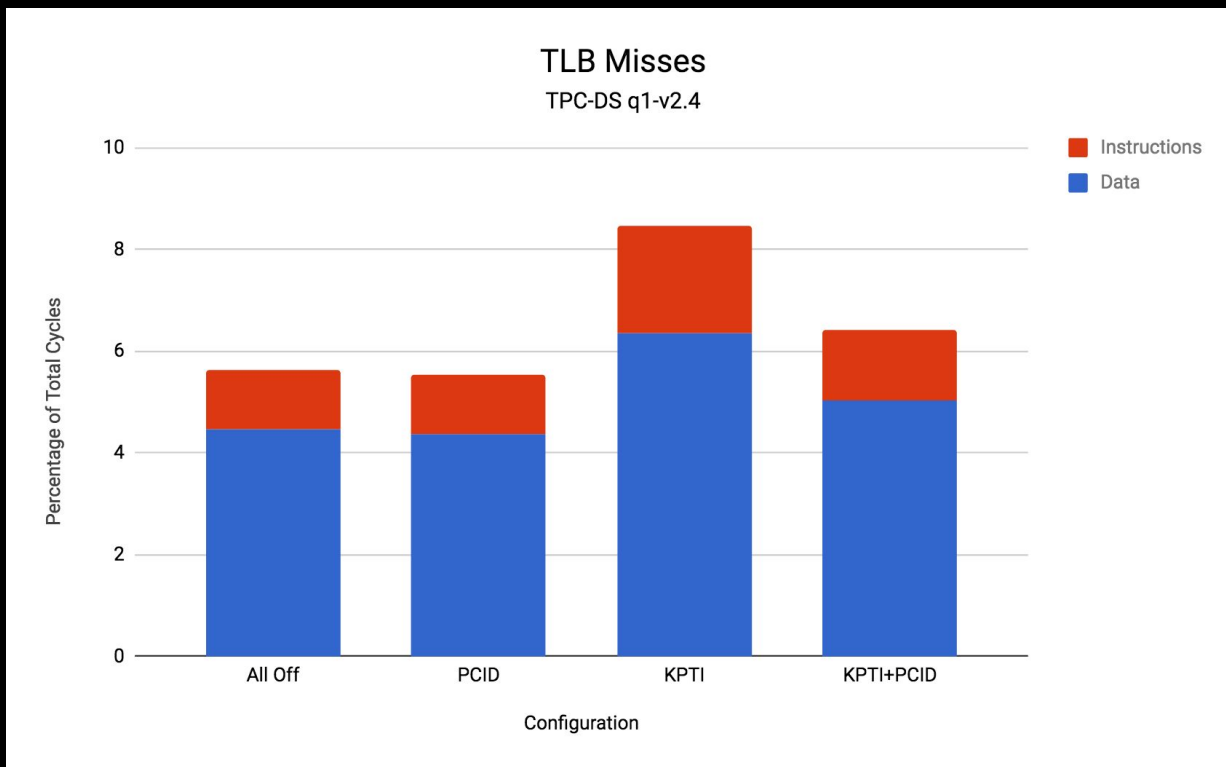
sys\_clock\_gettime

## TLB MISS

Virtual Address	Physical Address
-	-
-	-

print \*UserAddress

# Meltdown TLB Misses



# TLB with KPTI and PCID



```
while True:
    print *UserAddress
    sys_clock_gettime
```

## TLB MISS

Virtual Address	PCID	Physical Address
-	-	-
-	-	-

print \*UserAddress

## TLB MISS

Virtual Address	PCID	Physical Address
UserAddress	1	Page 1
-	-	-

sys\_clock\_gettime

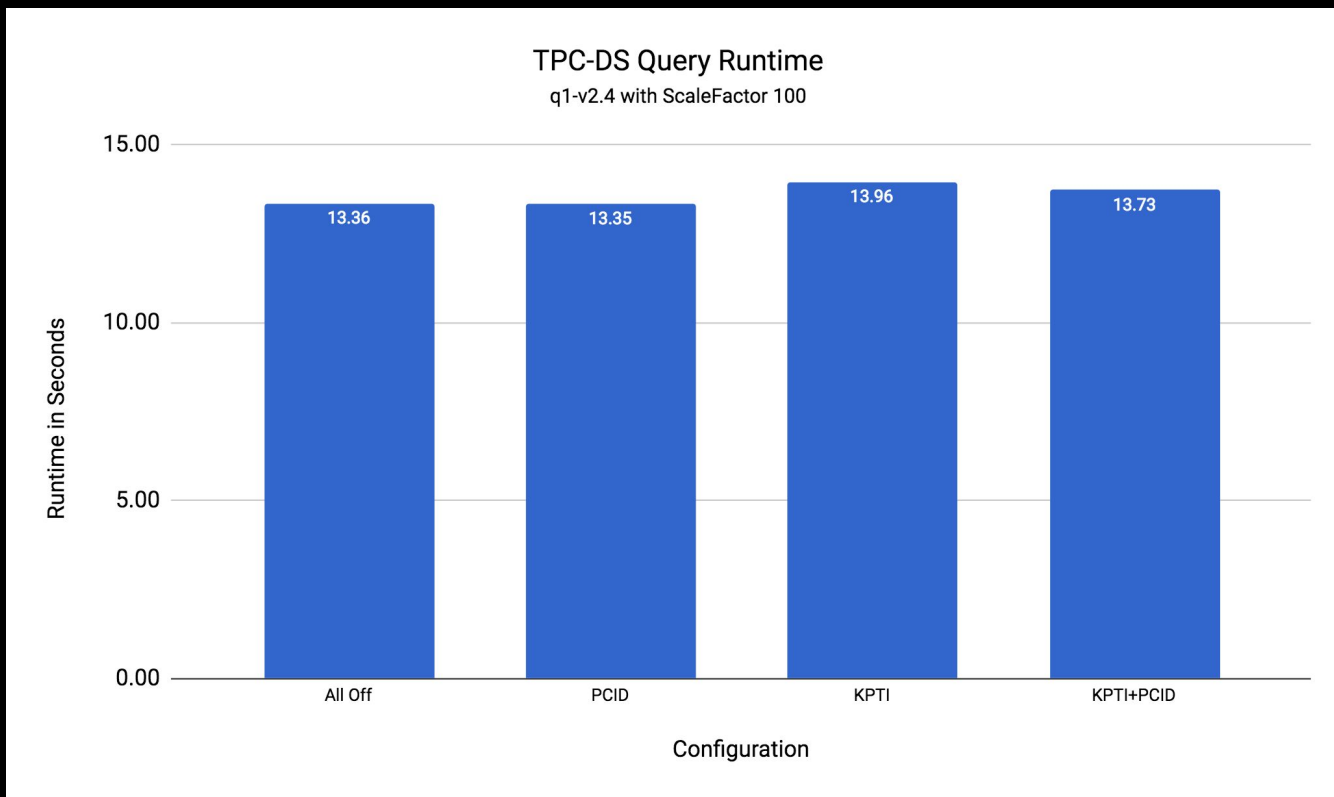
## TLB HIT

Virtual Address	PCID	Physical Address
UserAddress	1	Page 1
KernelTime	0	Page 2

print \*UserAddress



# Meltdown Runtime





# Spectre V1

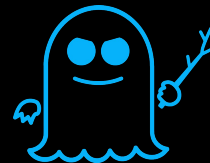
Bounds Check Bypass

# Bounds Check Bypass

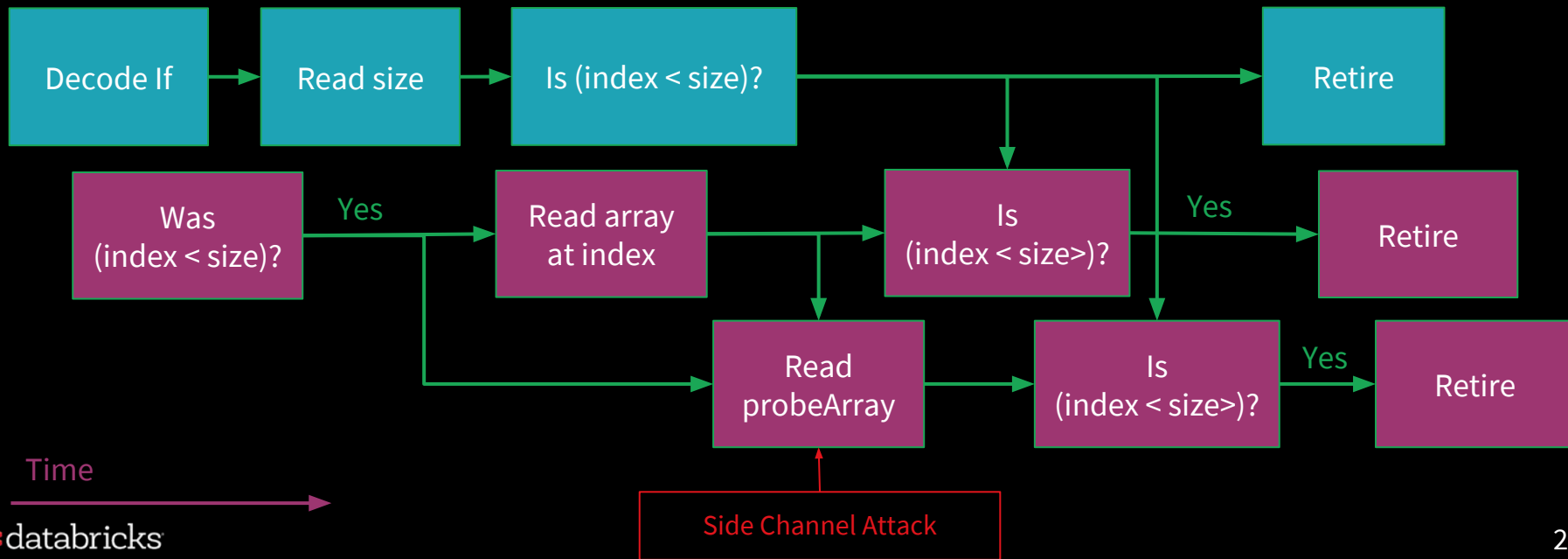


```
if (index < size) {  
    val = array[index];  
    probeArray[val];  
}
```

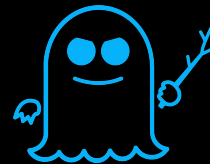
# Bounds Check Bypass



```
if (index < size) {  
    val = array[index];  
    probeArray[val];  
}
```

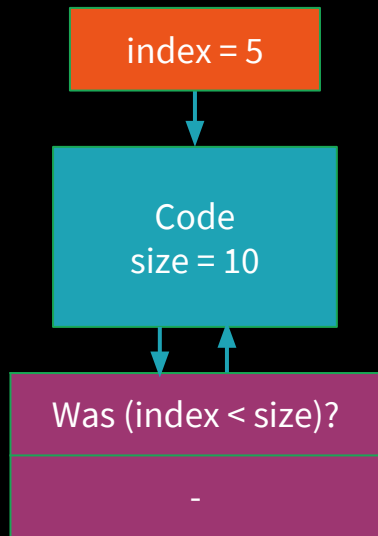


# Bounds Check Bypass

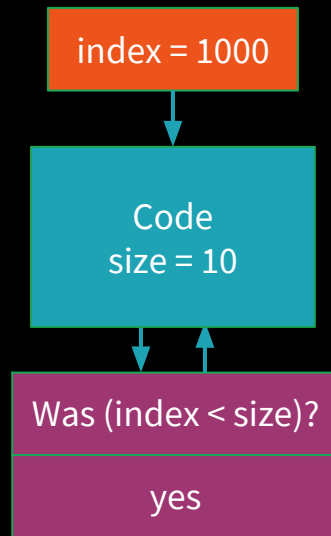


```
if (index < size) {  
    val = array[index];  
    probeArray[val];  
}
```

## TRAIN



## ATTACK



# Observable Speculation Barrier



- Protects against Spectre V1 - Bounds Check Bypass
- ~14 in the Ubuntu kernel and drivers
- Stops speculative array access with the LFENCE barrier

Before

```
if (index < size) {  
    val = array[index];  
    probeArray[val];  
}
```

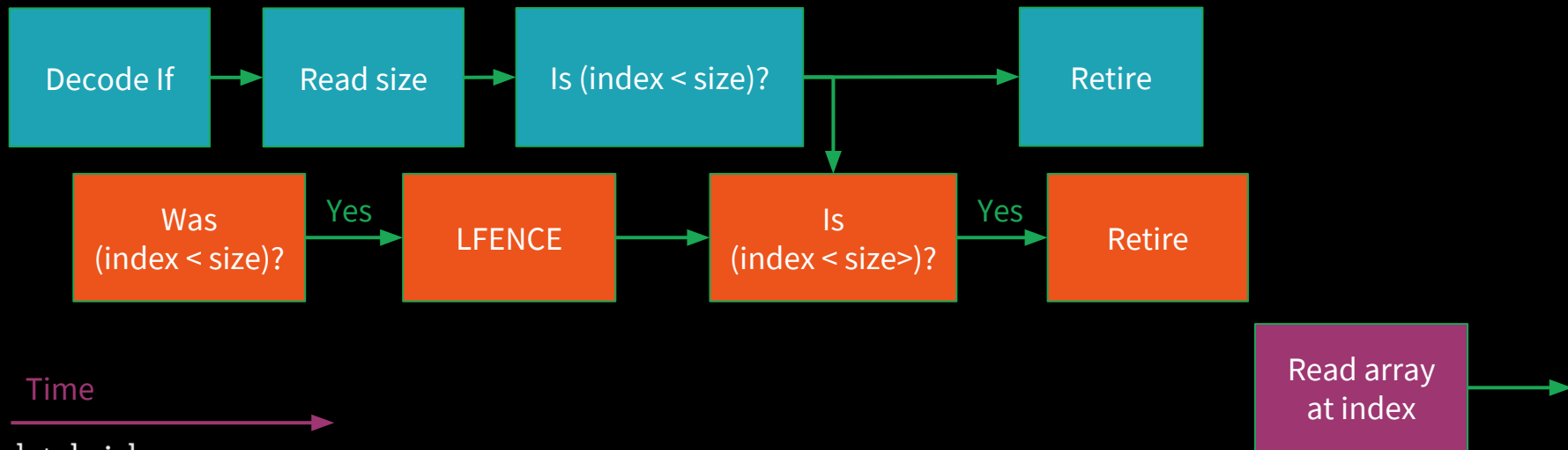
After

```
if (index < size) {  
    osb();  
    val = array[index];  
    probeArray[val];  
}
```

# Observable Speculation Barrier



```
if (index < size) {  
  osb()  
  val = array[index];  
  probeArray[val];  
}
```





# Spectre V2

Branch Target Injection



# Branch Target Injection

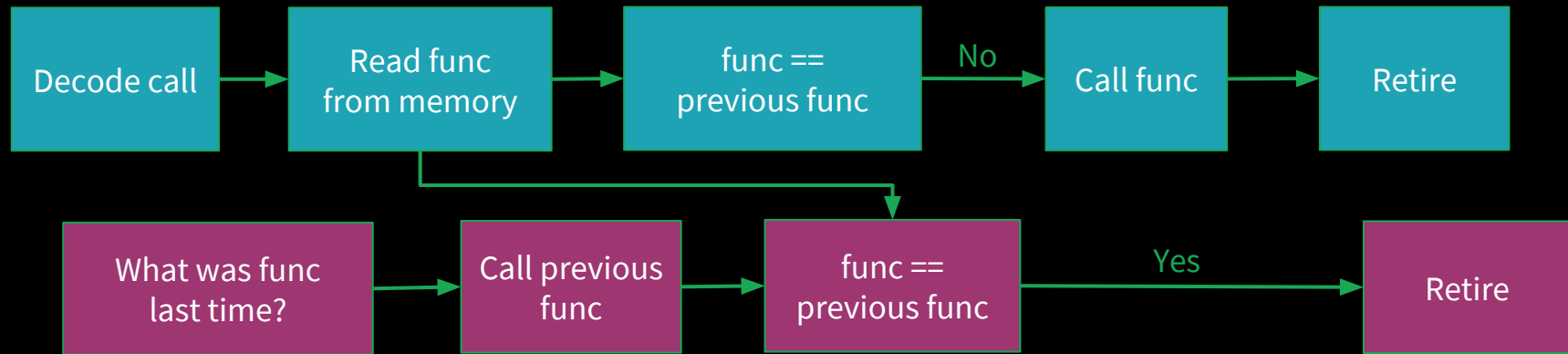


```
def call_func(func, arg):  
    func(arg)
```

# Branch Target Injection



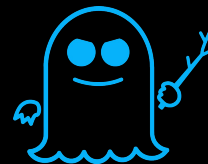
```
def call_func(func, arg):  
    func(arg)
```



Time

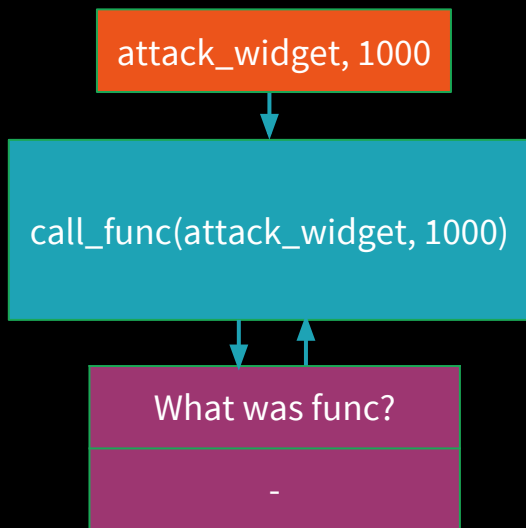


# Branch Target Injection

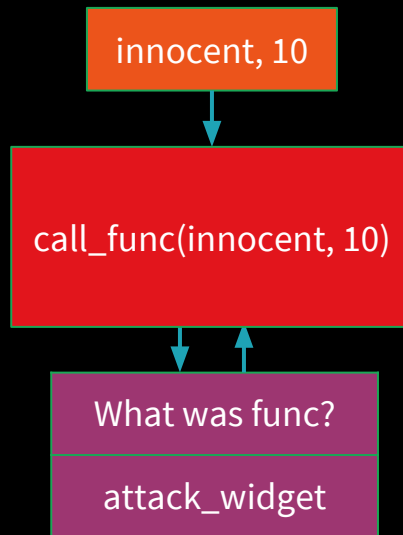


```
def call_func(func, arg):  
    func(arg)
```

## TRAIN



## ATTACK



# Intel Microcode Updates

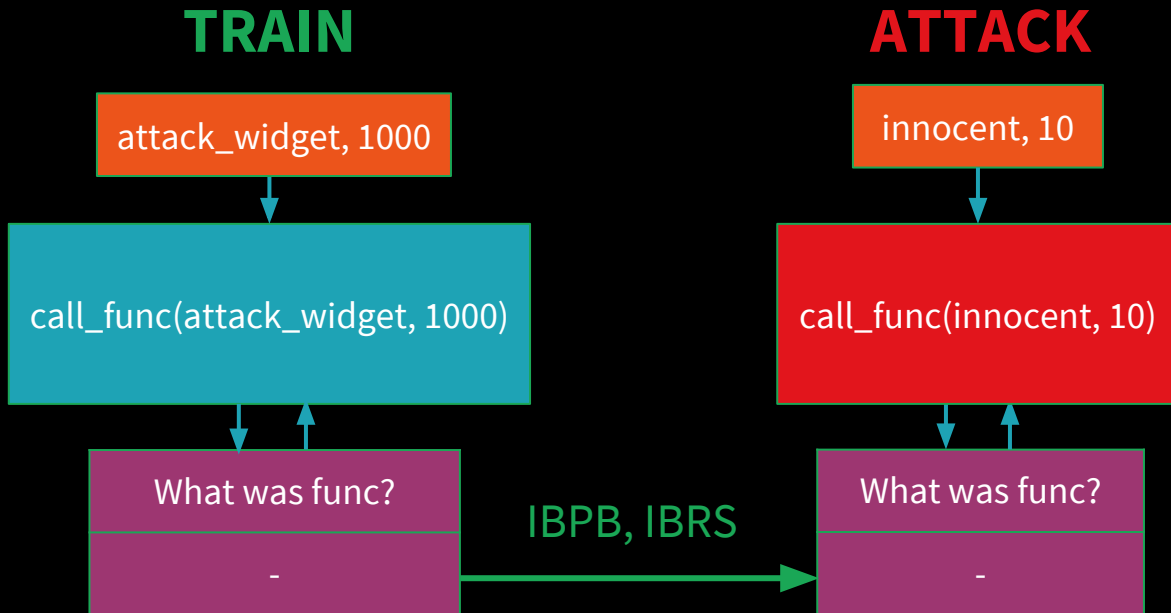


- Protect against Spectre V2 - Branch Target Injection
- Indirect Branch Restricted Speculation (IBRS)
  - Stops attacks from code running at lower privilege levels
  - Stops attacks from code running on the sibling hyperthread (STIBP)
- Indirect Branch Prediction Barrier (IBPB)
  - Stops attacks from code running at the same privilege level
  - Inserted at User-to-User and Guest-to-Guest transitions

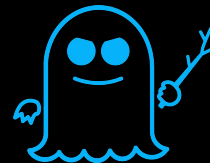
# IBPB and IBRS Patches



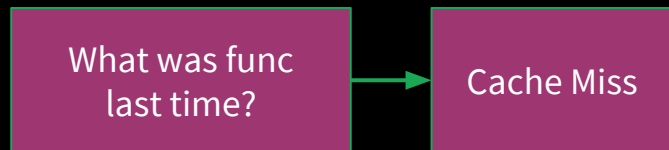
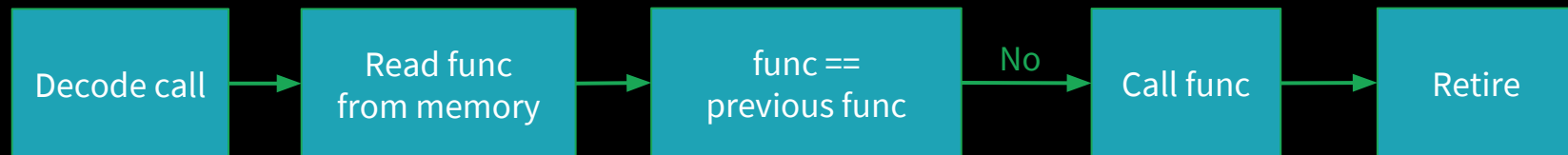
```
def call_func(func, arg):  
    func(arg)
```



# Branch Target Injection

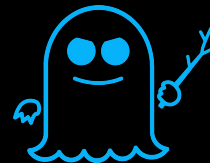


```
def call_func(func, arg):  
    func(arg)
```

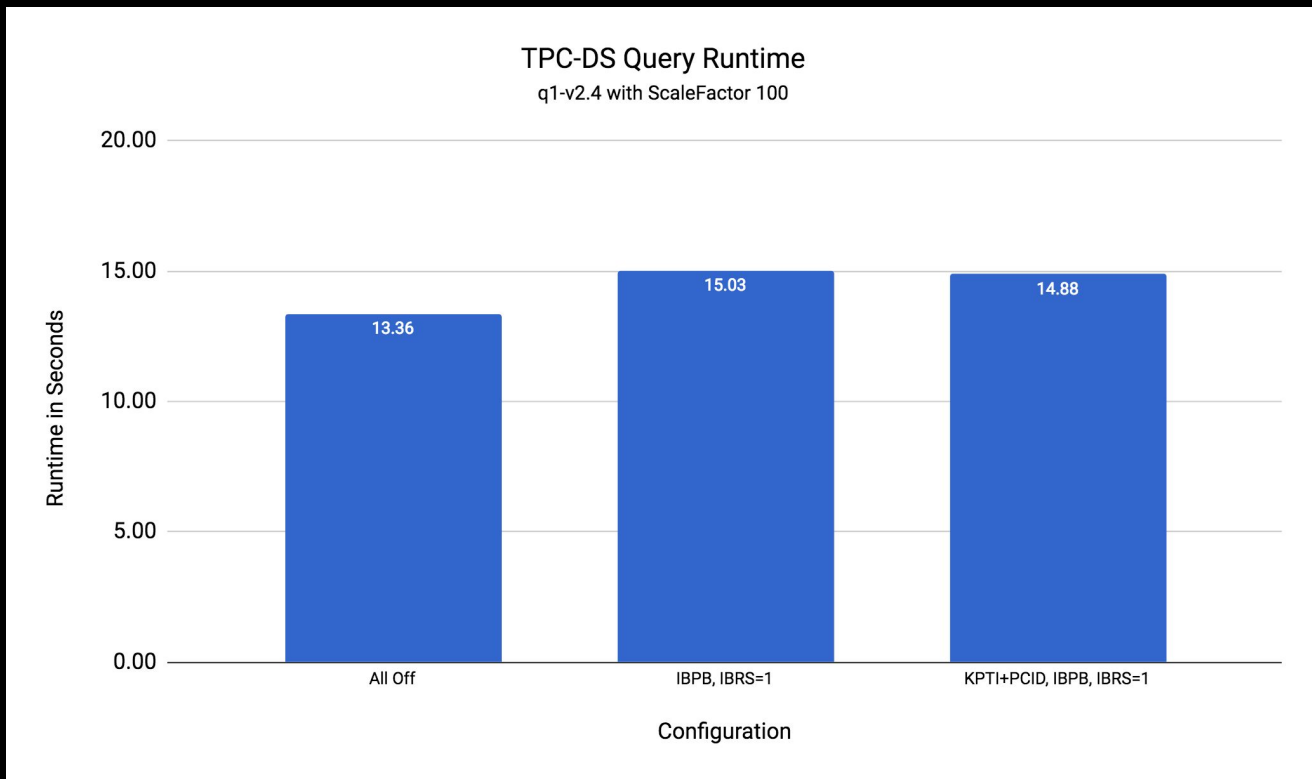
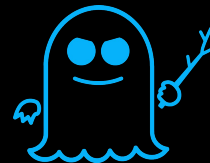


Time  
→

# Spectre V2 Branch Misses



# Spectre V2 Runtime





# Retpoline



- Protects indirect branches/calls from speculative execution
- Uses an infinite loop to do it

Before

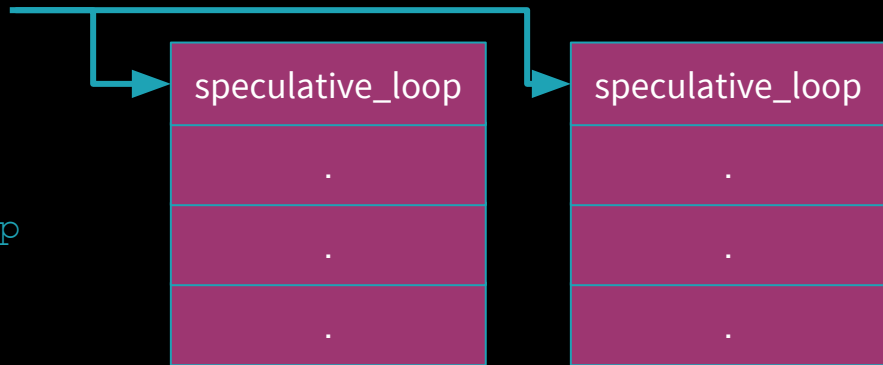
```
jmp rax
```

After

```
call set_up_target  
  
speculative_loop:  
  pause  
  lfence  
  jmp speculative_loop  
  
set_up_target:  
  mov rax, [rsp]  
  ret
```

Return Stack Buffer

Stack



# Retpoline



- Protects indirect branches/calls from speculative execution
- Uses an infinite loop to do it

Before

```
jmp rax
```

After

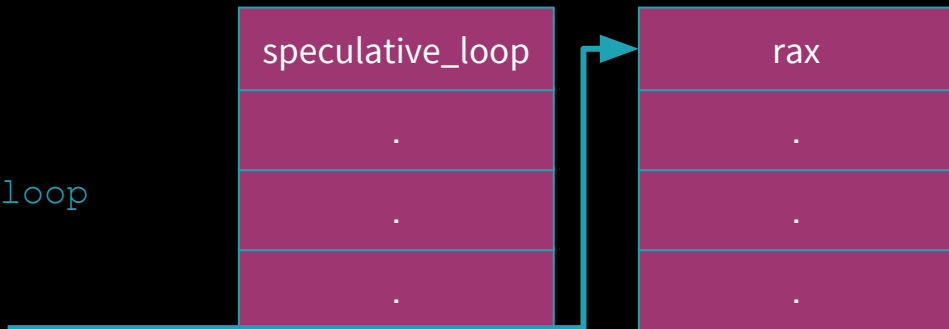
```
    call set_up_target

speculative_loop:
    pause
    lfence
    jmp speculative_loop

set_up_target:
    mov rax, [rsp]
    ret
```

Return Stack Buffer

Stack



# Retpoline



- Protects indirect branches/calls from speculative execution
- Uses an infinite loop to do it

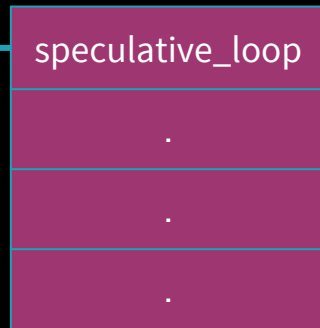
Before

```
jmp rax
```

After

```
    call set_up_target  
  
speculative_loop:  
    pause  
    lfence  
    jmp speculative_loop  
  
set_up_target:  
    mov rax, [rsp]  
    ret
```

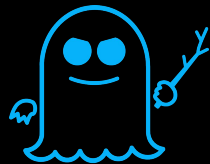
Return Stack Buffer



Stack



# Retpoline



- Protects indirect branches/calls from speculative execution
- Uses an infinite loop to do it

Before

```
jmp rax
```

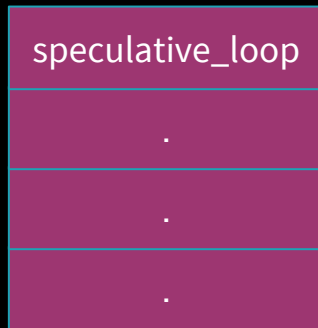
After

```
    call set_up_target

speculative_loop:
    pause
    lfence
    jmp speculative_loop

set_up_target:
    mov rax, [rsp]
    ret
```

Return Stack Buffer



Stack



# Retpoline



- Protects indirect branches/calls from speculative execution
- Uses an infinite loop to do it

Before

```
jmp rax
```

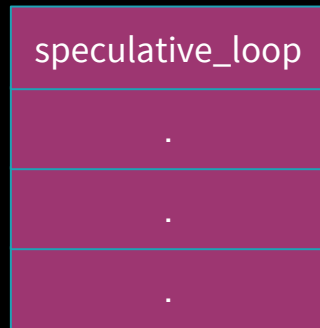
After

```
    call set_up_target

speculative_loop:
    pause
    lfence
    jmp speculative_loop

set_up_target:
    mov rax, [rsp]
    ret
```

Return Stack Buffer



Stack



# Retpoline (Improved)



- Check against a known target and direct branch

Before

```
jmp rax
```

After

```
cmp rax, known_target
jne retpoline
jmp known_target
```



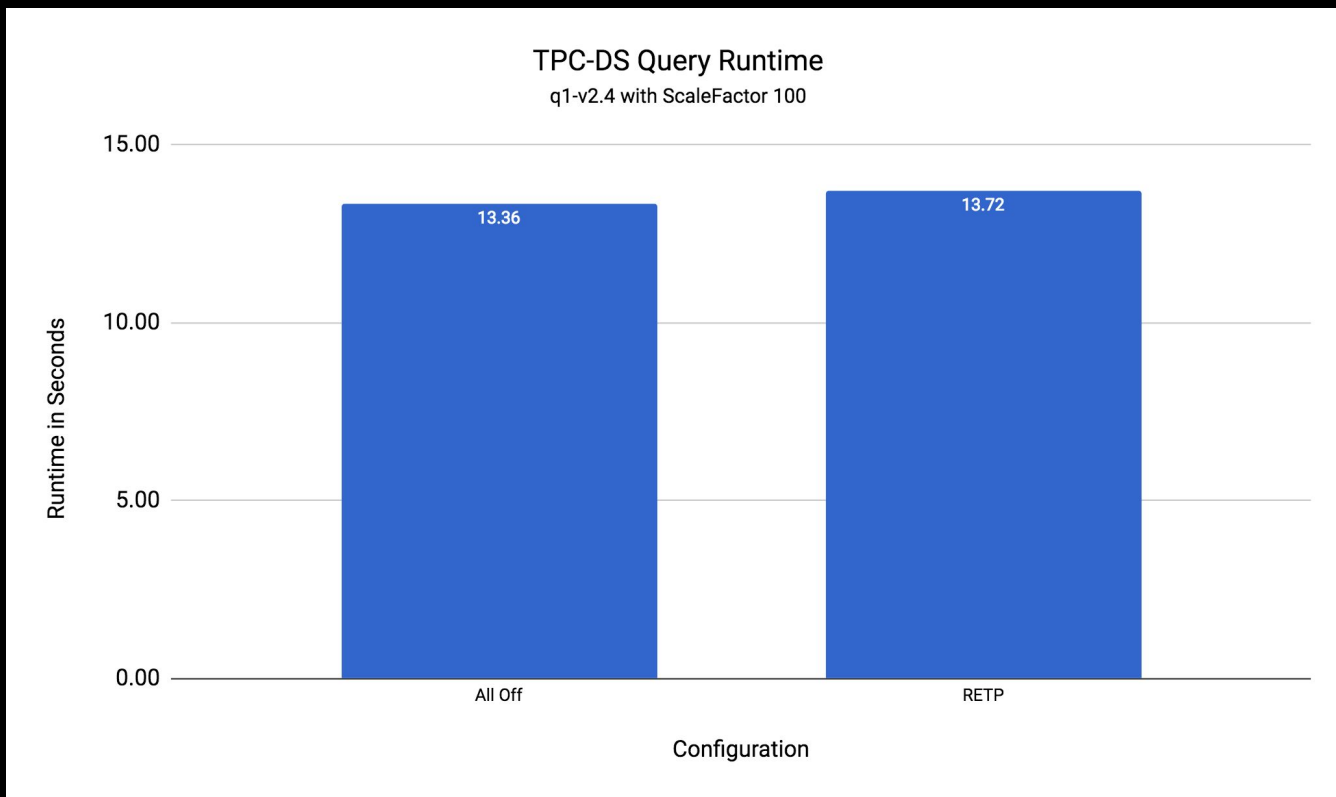
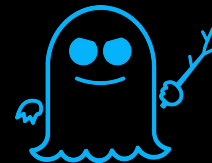
Direct Branch (fast)

```
retpoline:
  call set_up_target
```

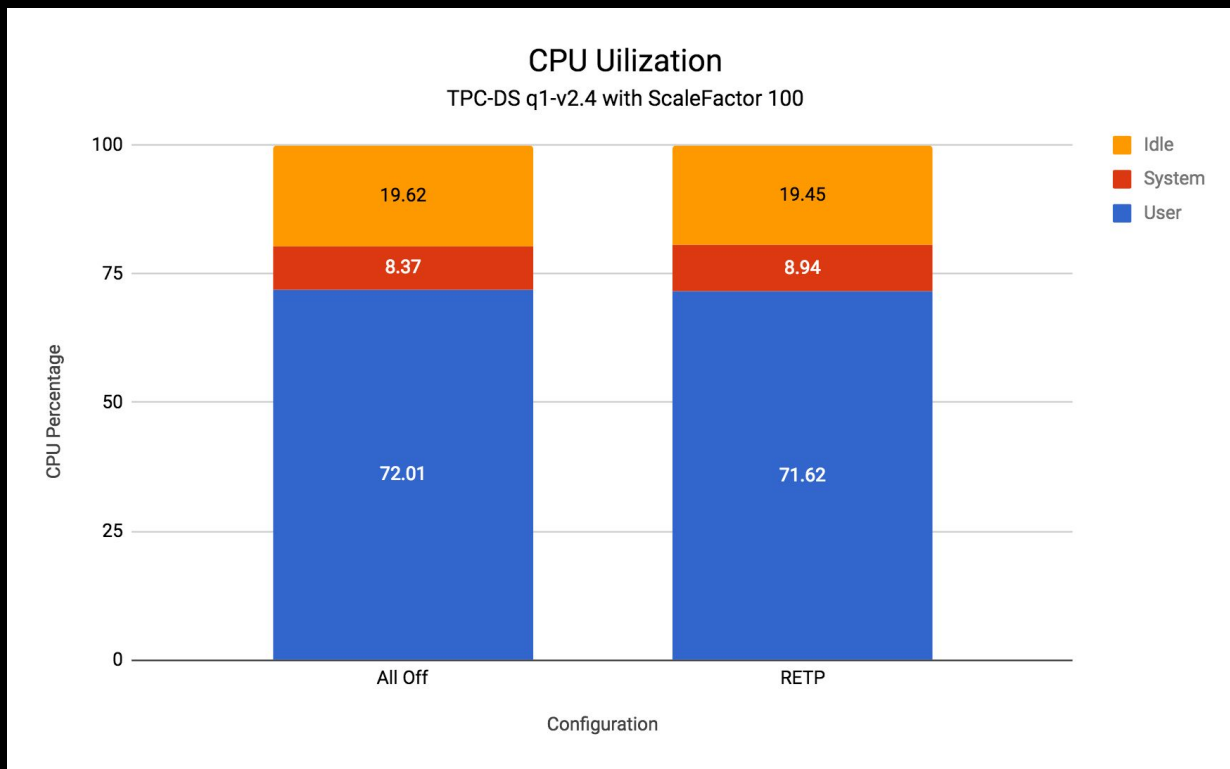
```
speculative_loop:
  pause
  jmp speculative_loop
```

```
set_up_target:
  mov rax, [rsp]
  ret
```

# Retpoline Runtime

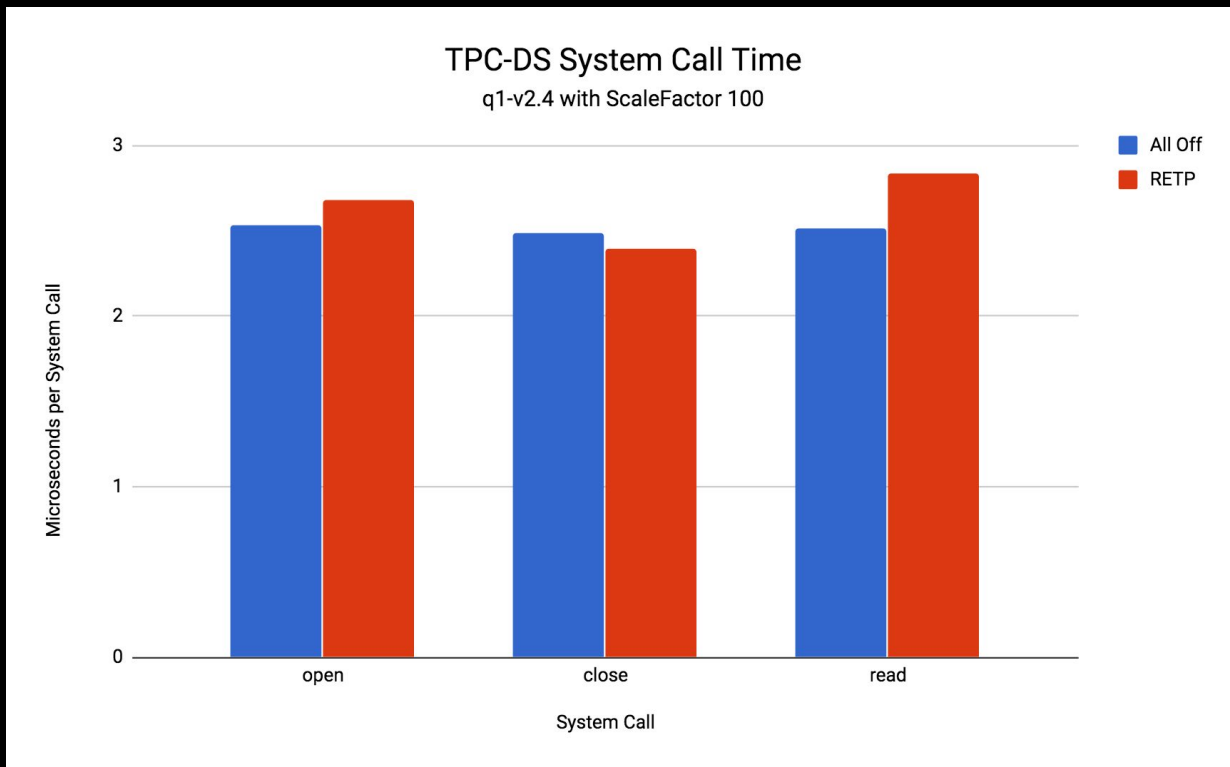


# Retpoline CPU Utilization





# Retpoline System Calls



# A Tale of Two Queries

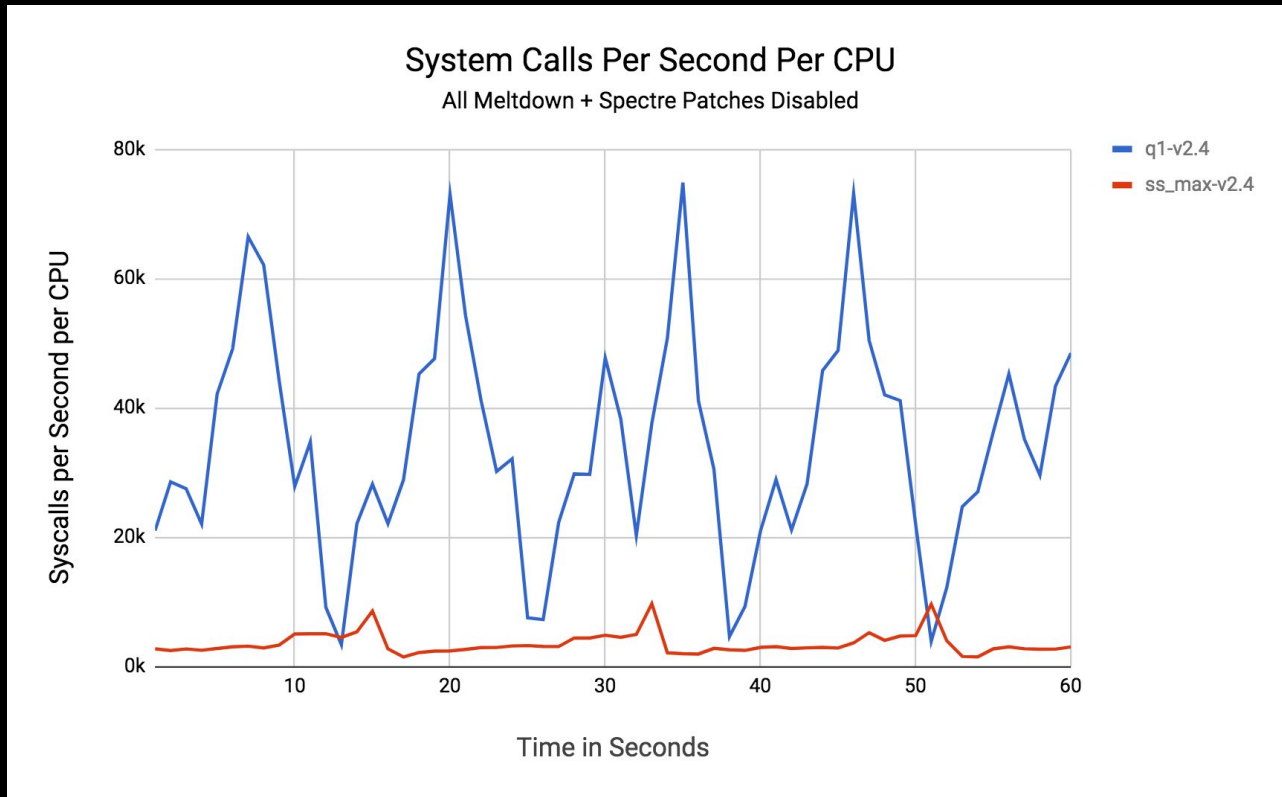
## q1-v2.4

```
WITH customer_total_return AS
  (SELECT sr_customer_sk AS ctr_customer_sk,
          sr_store_sk AS ctr_store_sk,
          sum(sr_return_amt) AS ctr_total_return
   FROM store_returns, date_dim
   WHERE sr_returned_date_sk = d_date_sk
   AND d_year = 2000
   GROUP BY sr_customer_sk, sr_store_sk)
SELECT c_customer_id
FROM customer_total_return ctr1, store, customer
WHERE ctr1.ctr_total_return >
  (SELECT avg(ctr_total_return)*1.2
   FROM customer_total_return ctr2
   WHERE ctr1.ctr_store_sk = ctr2.ctr_store_sk)
AND s_store_sk = ctr1.ctr_store_sk
AND s_state = 'TN'
AND ctr1.ctr_customer_sk = c_customer_sk
ORDER BY c_customer_id LIMIT 100
```

## ss\_max-v2.4

```
SELECT
  count(*) AS total,
  count(ss_sold_date_sk) AS not_null_total,
  count(DISTINCT ss_sold_date_sk) AS unique_days,
  max(ss_sold_date_sk) AS max_ss_sold_date_sk,
  max(ss_sold_time_sk) AS max_ss_sold_time_sk,
  max(ss_item_sk) AS max_ss_item_sk,
  max(ss_customer_sk) AS max_ss_customer_sk,
  max(ss_cdemo_sk) AS max_ss_cdemo_sk,
  max(ss_hdemo_sk) AS max_ss_hdemo_sk,
  max(ss_addr_sk) AS max_ss_addr_sk,
  max(ss_store_sk) AS max_ss_store_sk,
  max(ss_promo_sk) AS max_ss_promo_sk
FROM store_sales
```

# q1-v2.4 vs ss\_max-v2.4





# Thank You

[chris.stevens@databricks.com](mailto:chris.stevens@databricks.com)

<https://www.linkedin.com/in/chrisrstevens/>



# Backup Slides

# array\_index\_nospec()



- Protects against Spectre V1 - Bounds Check Bypass
- ~50 in the Linux kernel and drivers (not in Ubuntu)
- Stops speculative array access by clamping the index

Before

```
if (index < size) {  
    val = array[index];  
    probeArray[val];  
}
```

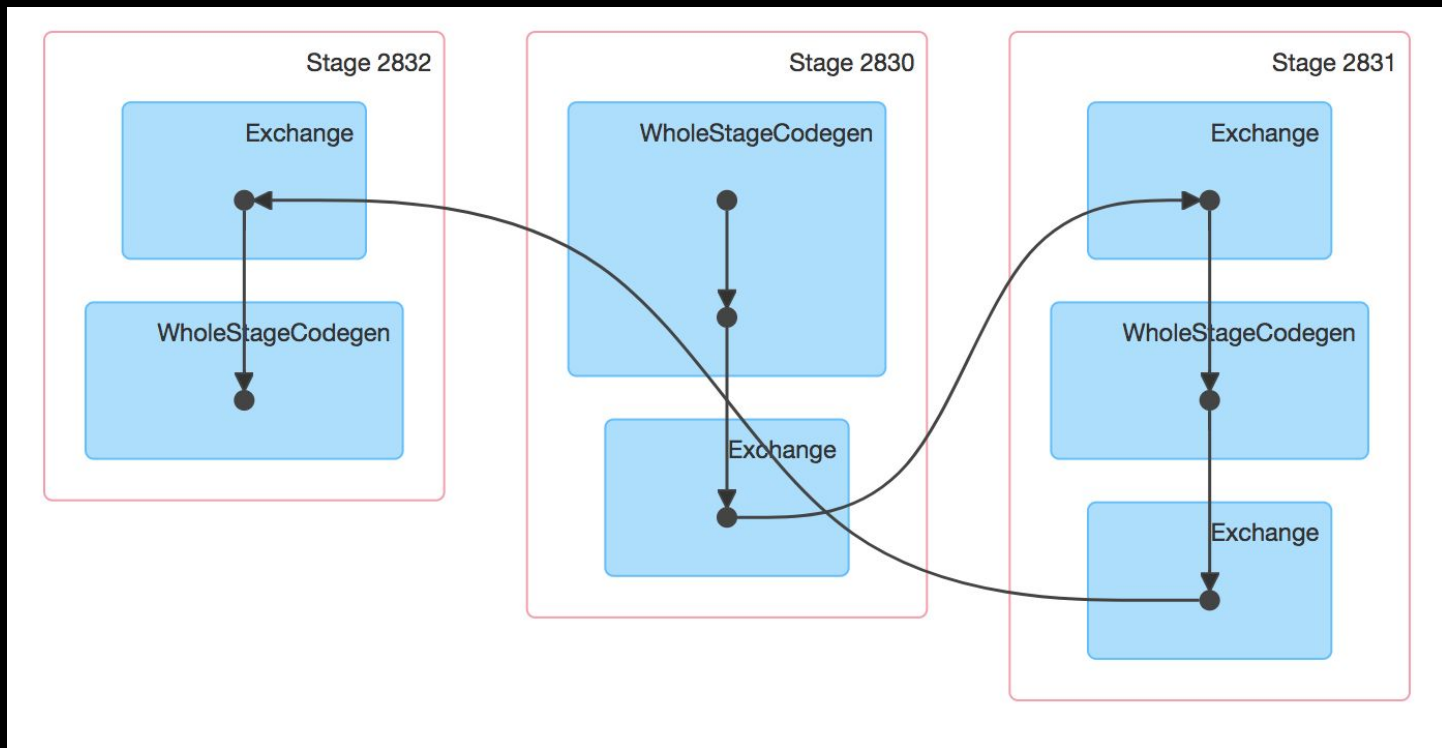
After

```
if (index < size) {  
    index = array_index_nospec(index, size);  
    val = array[index];  
    probeArray[val];  
}
```

# Speculative Store Bypass

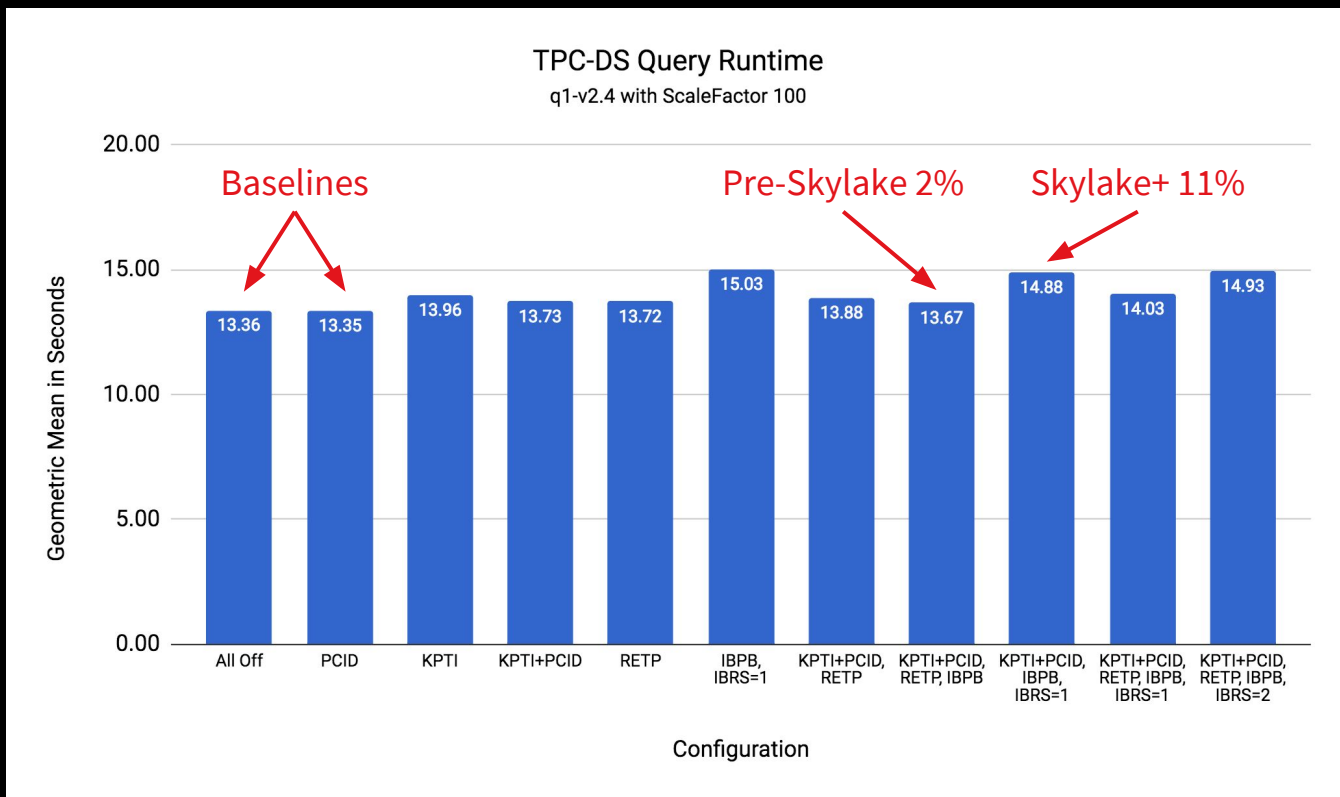
- CVE-2018-3639
- “Variant 4”
- Public Date: May 21, 2018
- Details:
  - <https://access.redhat.com/security/vulnerabilities/ssbd>
- Ubuntu Command Line Parameter:
  - `spec_store_bypass_disable=on`

# TPC-DS: ss\_max-v2.4

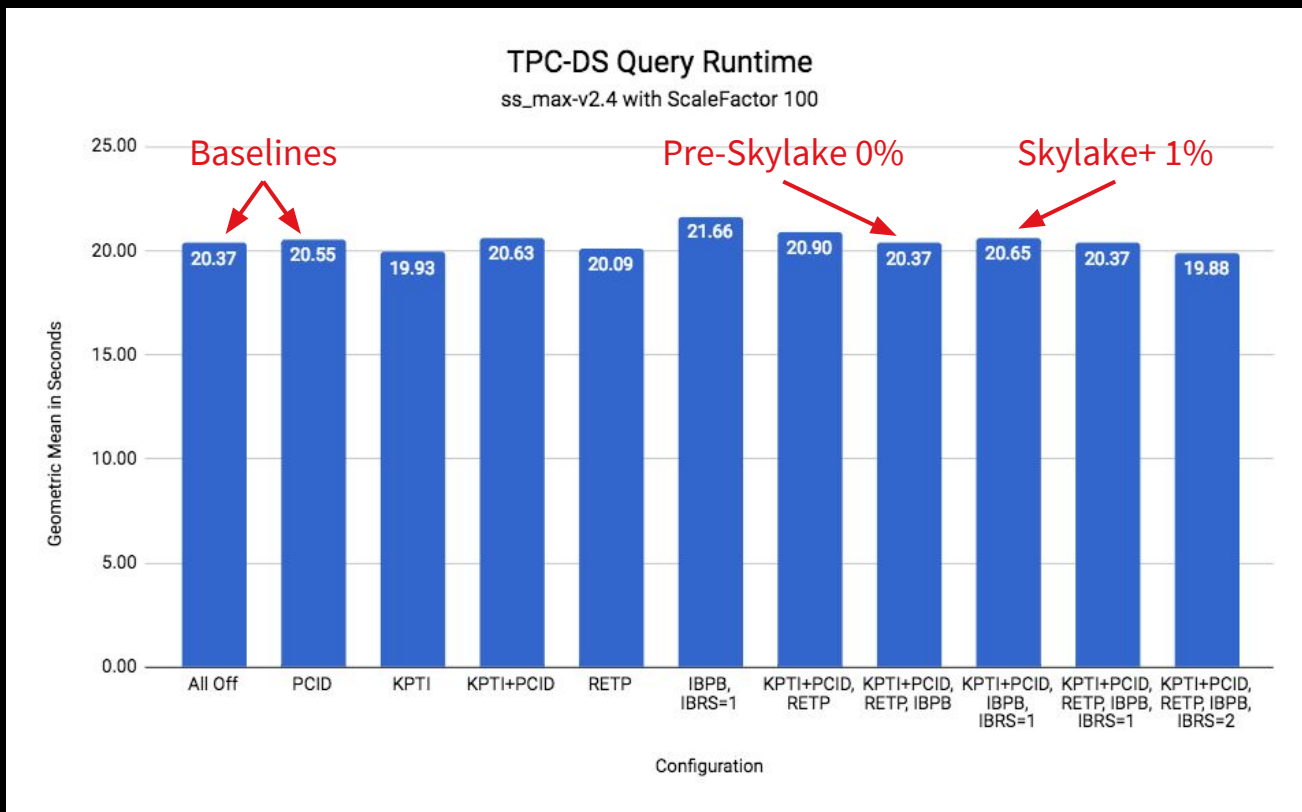




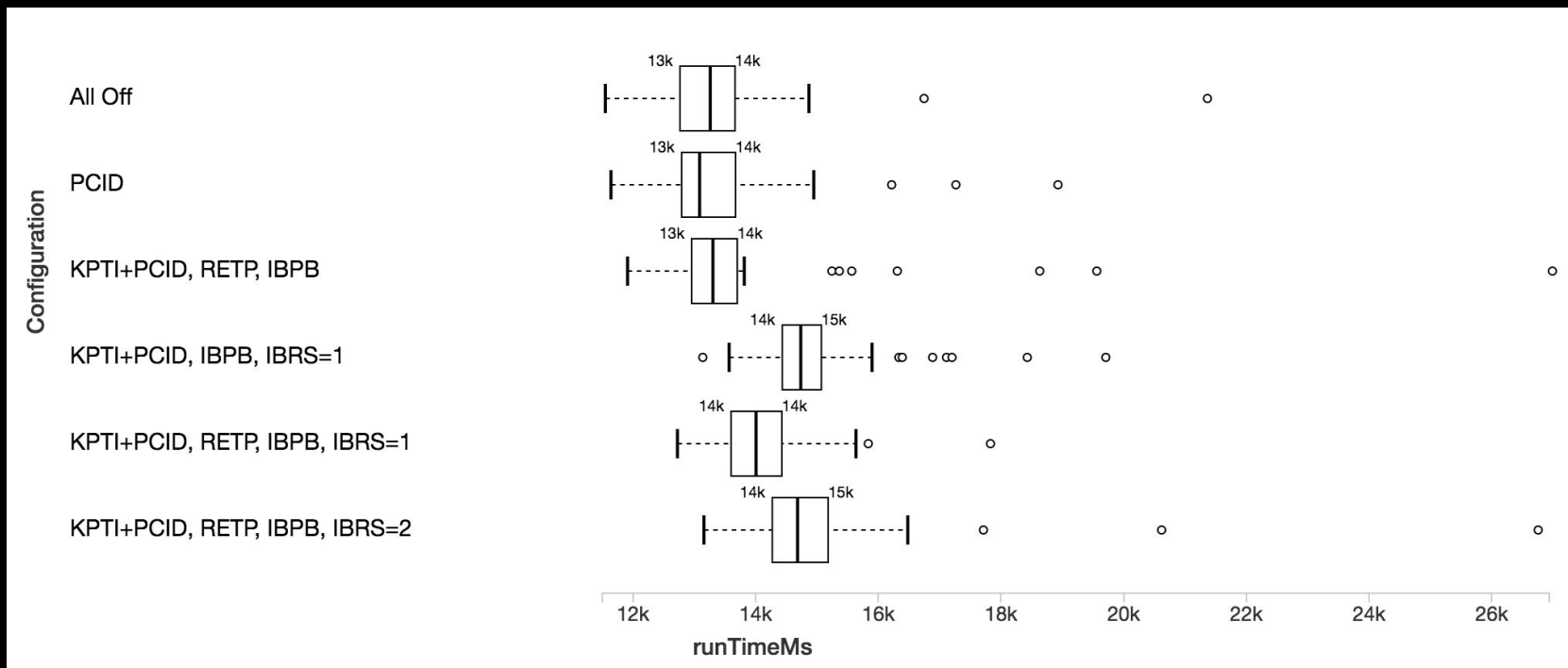
# TPC-DS: q1-v2.4



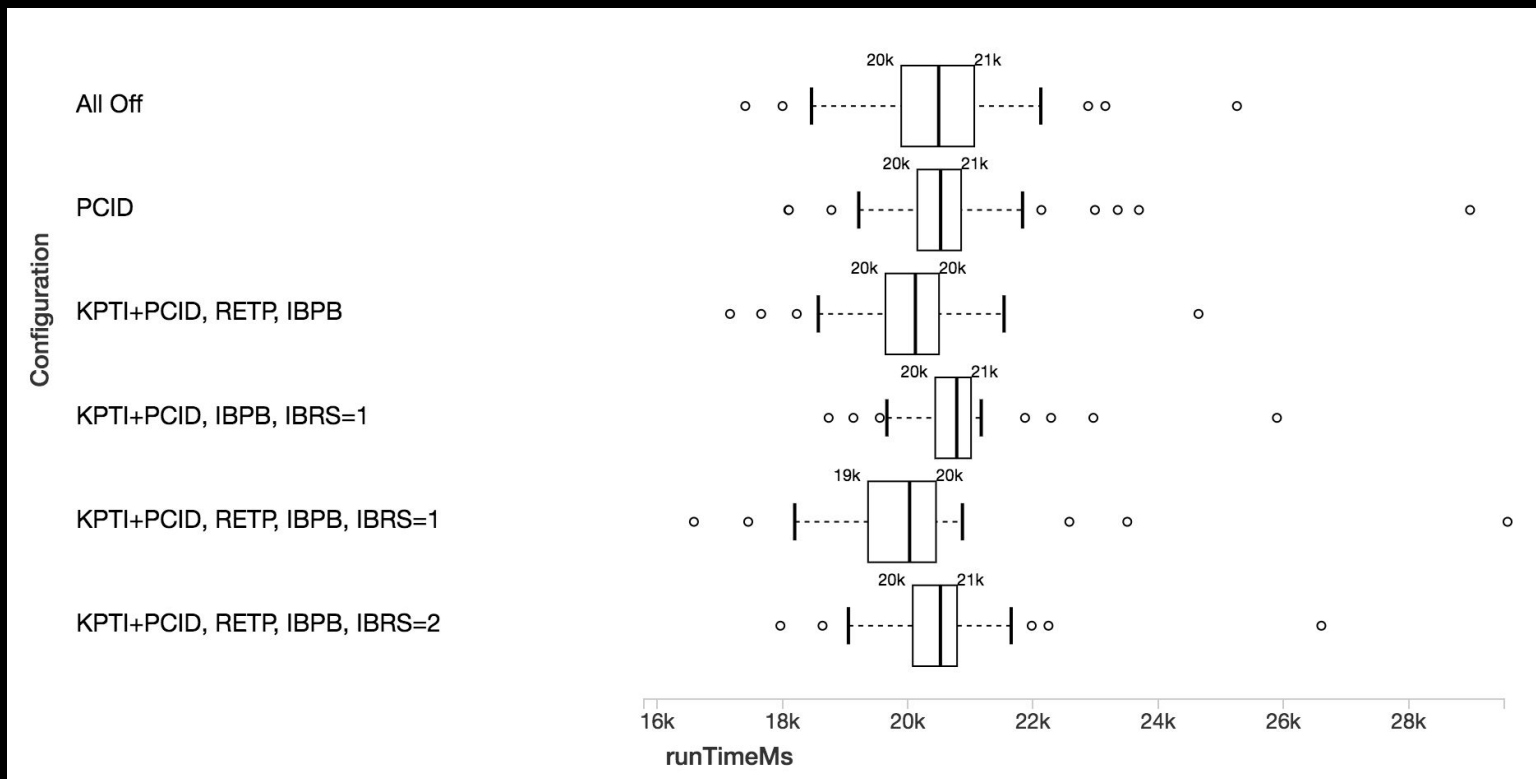
# TPC-DS: ss\_max-v2.4



# q1-v2.4 Runtime Box Plots



# ss\_max-v2.4 Runtime Box Plots



# Single Node TPC-DS Setup

- Intel Core i7-4820k @ 3.70GHz (Ivy Bridge - ca. 2012)
- 8 CPUs
- 8GB of RAM
- Disk: Western Digital Blue
  - WDC10EZEX-08M2NA0
  - 1TB
  - 7200 rpm
  - 146MB/s sequential read (<http://hdd.userbenchmark.com/>)
- Ubuntu 16.04.4 LTS (Xenial Xerus)
  - 64-bit server image
  - 4.4.0-116-generic Linux kernel

# Repositories

- <https://github.com/apache/spark.git>
- <https://github.com/databricks/spark-sql-perf.git>
- <https://github.com/databricks/tpcds-kit>
- <https://github.com/speed47/spectre-meltdown-checker>
- <https://github.com/brendangregg/pmc-cloud-tools>
- <https://github.com/brendangregg/perf-tools>