Moritz Lipp (@mlqxyz)

Michael Schwarz (@misc0110)

Daniel Gruss (@lavados)

# Meltdown

## Basics, Details, Consequences

**Moritz Lipp**

PhD student @ Graz University of Technology

@mlqxyz

moritz.lipp@iaik.tugraz.at

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology
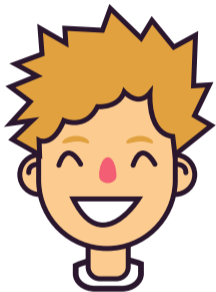
**Michael Schwarz**

PhD student @ Graz University of Technology

🐦 @misc0110

✉ michael.schwarz@iaik.tugraz.at

**Daniel Gruss**

PostDoc @ Graz University of Technology
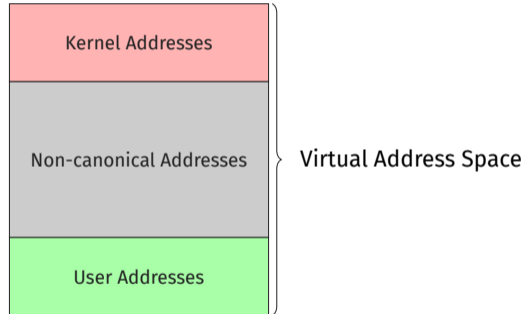
@lavados

daniel.gruss@iaik.tugraz.at

- Anders Fogh
- Daniel Genkin
- Werner Haas
- Mike Hamburg
- Jann Horn
- Paul Kocher
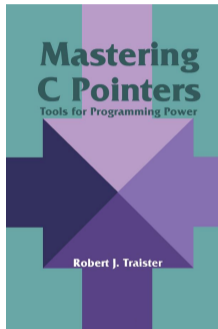- Stefan Mangard
- Thomas Prescher
- Yuval Yarom

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**Let's Read Kernel Memory from User Space!**

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Find something human readable, e.g., the Linux version

```
# sudo grep linux_banner /proc/kallsyms
ffffffff81a000e0 R linux_banner
```

**Mastering C Pointers**
Tools for Programming Power

Robert J. Traister

```c
char data = *(char*) 0xffffffff81a000e0;
printf("%c\n", data);
```

- Compile and run

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

• Compile and run

```
segfault at ffffffff81a000e0 ip 0000000000400535
      sp 00007ffce4a80610 error 5 in reader
```

- Compile and run

```
segfault at ffffffff81a000e0 ip 0000000000400535
      sp 00007ffce4a80610 error 5 in reader
```

- Kernel addresses are of course not accessible

- Compile and run

```
segfault at ffffffff81a000e0 ip 0000000000400535
      sp 00007ffce4a80610 error 5 in reader
```

- Kernel addresses are of course not accessible
- Any invalid access throws an exception → segmentation fault

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Just catch the segmentation fault!

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Just catch the segmentation fault!
- We can simply install a signal handler

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Just catch the segmentation fault!
- We can simply install a signal handler
- And if an exception occurs, just jump back and continue

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Just catch the segmentation fault!
- We can simply install a signal handler
- And if an exception occurs, just jump back and continue
- Then we can read the value

- Still no kernel memory

- Still no kernel memory
- Privilege checks seem to work

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Still no kernel memory
- Privilege checks seem to work
- Maybe it is not that straight forward
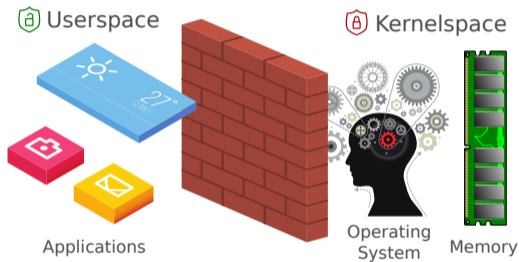
Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Still no kernel memory
- Privilege checks seem to work
- Maybe it is not that straight forward
- Back to the drawing board

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

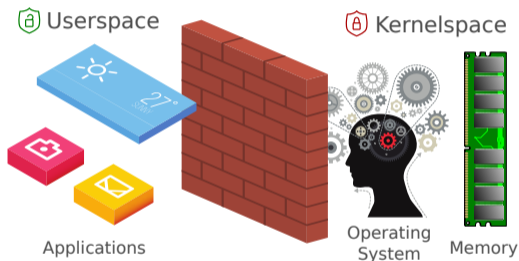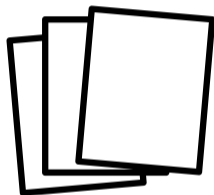# Operating Systems 101

- Kernel is isolated from user space

- Kernel is isolated from user space
- This isolation is a combination of hardware and software

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

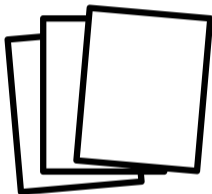Userspace

Kernelspace

Applications

Operating System

Memory

- Kernel is isolated from user space
- This isolation is a combination of hardware and software
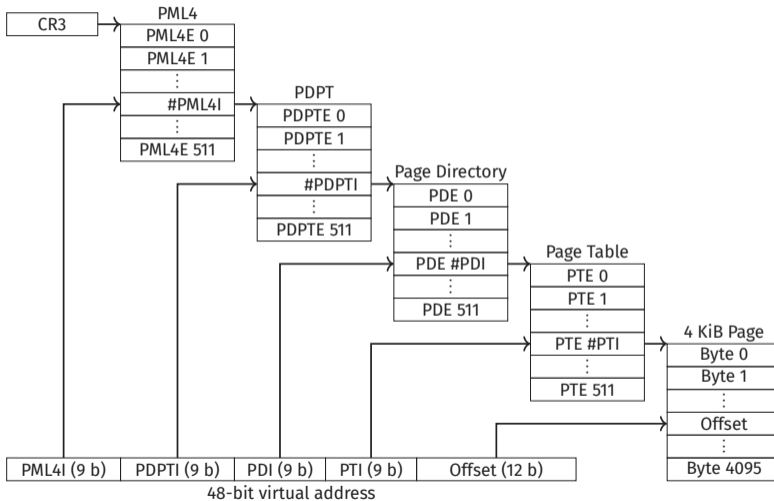- User applications cannot access anything from the kernel
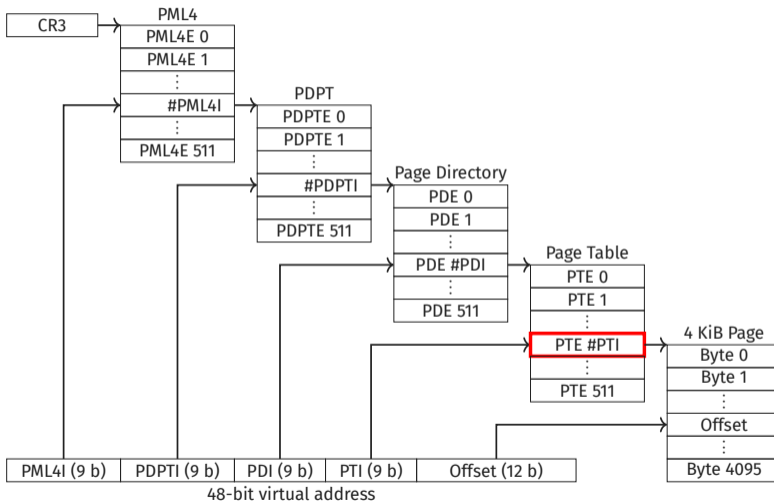
- CPU support virtual address spaces to isolate processes

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- CPU support virtual address spaces to isolate processes
- Physical memory is organized in page frames

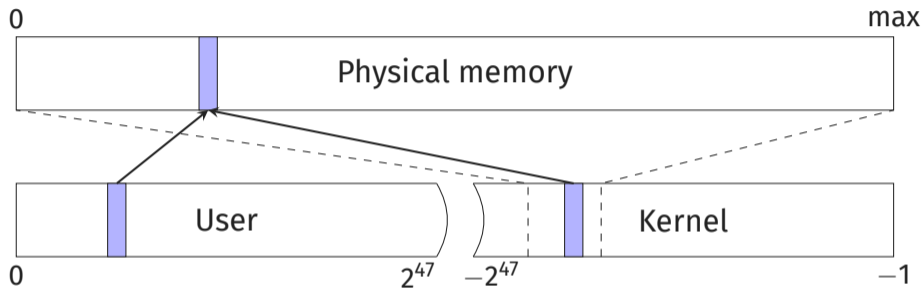Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- CPU support virtual address spaces to isolate processes
- Physical memory is organized in page frames
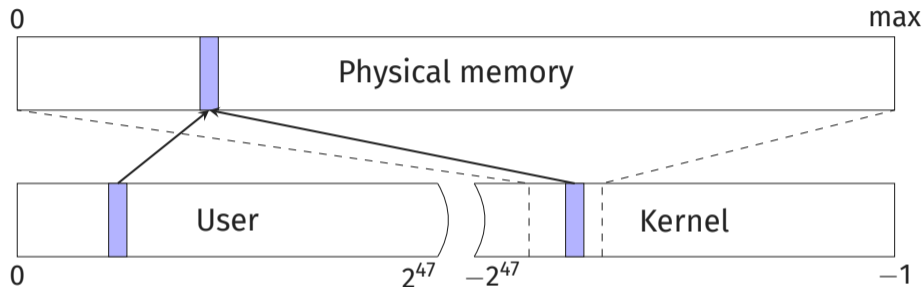- Virtual memory pages are mapped to page frames using page tables

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

| P | RW | US | WT | UC | R | D | S | G | Ignored | |
|---|----|----|----|----|----|----|----|----|---------|---|
| Physical Page Number | | | | | | | | | | |
| | | | | | | | | | | |
| | | Ignored | | | | | | | | X |

- User/Supervisor bit defines in which privilege level the page can be accessed

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Kernel is typically mapped into every address space

0                                                                          max

Physical memory

0                                              $2^{47}$    $-2^{47}$                    $-1$

User                                                    Kernel

- Kernel is typically mapped into every address space
- Entire physical memory is mapped in the kernel

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

# Side-channel Attacks

- Safe software infrastructure does not mean safe execution

- Safe software infrastructure does not mean safe execution
- Information leaks because of the <span style="color:red">underlying hardware</span>

- Safe software infrastructure does not mean safe execution
- Information leaks because of the underlying hardware
- Exploit unintentional information leakage by side-effects

- Safe software infrastructure does not mean safe execution
- Information leaks because of the underlying hardware
- Exploit unintentional information leakage by side-effects
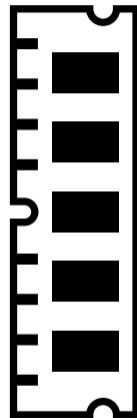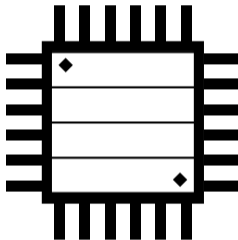
Power consumption

Execution time

CPU caches

• • •

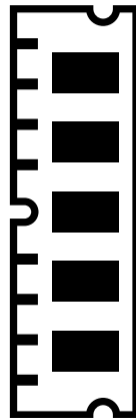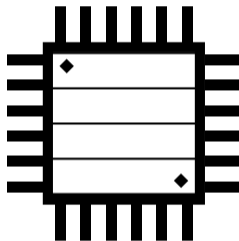Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology
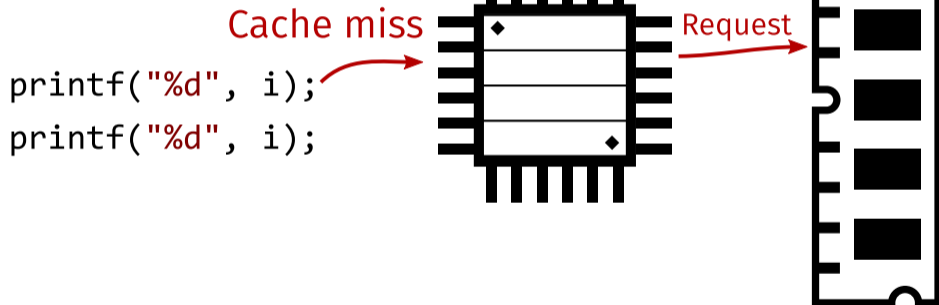
# Caches and Cache Attacks
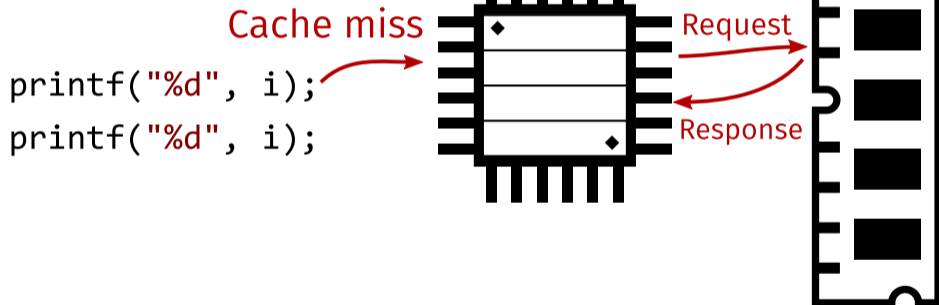
```
printf("%d", i);
printf("%d", i);
```

Cache miss
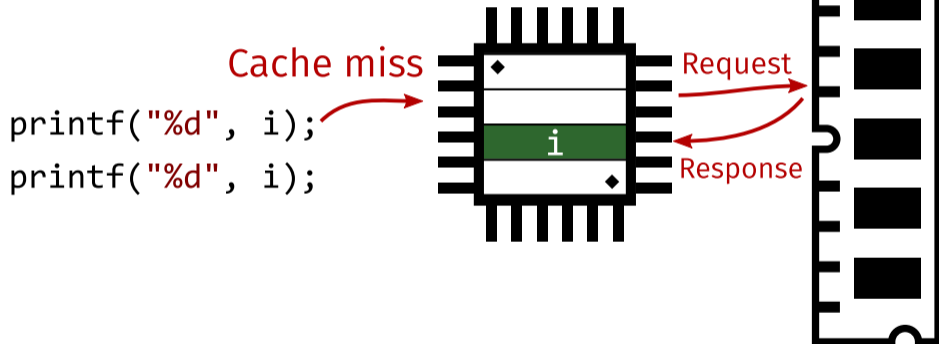
```
printf("%d", i);
printf("%d", i);
```

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Cache miss

Request

```
printf("%d", i);
printf("%d", i);
```

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Cache miss

```
printf("%d", i);
printf("%d", i);
```

Request

Response

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Cache miss

`printf("%d", i);`
`printf("%d", i);`

Cache hit

Request

i

Response

No DRAM access,
much faster

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

DRAM access, slow

Cache miss

Request

printf("%d", i);

i

printf("%d", i);

Response

Cache hit

No DRAM access, much faster

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Shared Memory

ATTACKER

flush
access

VICTIM

access

Shared Memory

Victim accessed
(fast)

vs

Victim did not access
(slow)

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

# Microarchitecture

- Instruction Set Architecture (ISA) is an abstract model of a computer (x86, ARMv8, SPARC, …)

- Instruction Set Architecture (ISA) is an abstract model of a computer (x86, ARMv8, SPARC, …)
- Serves as the interface between hardware and software

- Instruction Set Architecture (ISA) is an abstract model of a computer (x86, ARMv8, SPARC, ...)
- Serves as the interface between hardware and software
- Microarchitecture is an actual implementation of the ISA

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Instruction Set Architecture (ISA) is an abstract model of a computer (x86, ARMv8, SPARC, …)
- Serves as the interface between hardware and software
- Microarchitecture is an actual implementation of the ISA

- Instructions are…
  - fetched (IF) from the L1 Instruction Cache

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Instructions are...
  - fetched (IF) from the L1 Instruction Cache
  - decoded (ID)

| IF | ID | EX | MEM | WB | | |
|----|----|----|-----|----|----|----|
| | IF | ID | EX | MEM | WB | |
| | | IF | ID | EX | MEM | WB |
| | | | IF | ID | EX | MEM | WB |
| | | | | IF | ID | EX | MEM | WB |

- Instructions are…
  - fetched (`IF`) from the L1 Instruction Cache
  - decoded (`ID`)
  - executed (`EX`) by execution units

| IF | ID | EX | MEM | WB |
|----|----|----|-----|----|

(pipeline diagram)

- Instructions are…
  - fetched (`IF`) from the L1 Instruction Cache
  - decoded (`ID`)
  - executed (`EX`) by execution units
- Memory access is performed (`MEM`)

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Instructions are…
  - fetched (`IF`) from the L1 Instruction Cache
  - decoded (`ID`)
  - executed (`EX`) by execution units
- Memory access is performed (`MEM`)
- Architectural register file is updated (`WB`)

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Instructions are executed in-order

- Instructions are executed in-order
- Pipeline stalls when stages are not ready

- Instructions are executed in-order
- Pipeline stalls when stages are not ready
- If data is not cached, we need to wait

```c
int width = 10, height = 5;

float diagonal = sqrt(width * width
                    + height * height);
int area = width * height;

printf("Area %d x %d = %d\n", width, height, area);
```

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**Parallelize**

**Dependency**

```c
int width = 10, height = 5;

float diagonal = sqrt(width * width
                    + height * height);
int area = width * height;

printf("Area %d x %d = %d\n", width, height, area);
```

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Instructions are

- fetched and decoded in the front-end

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Instructions are

- fetched and decoded in the front-end
- dispatched to the backend

Instructions are

- fetched and decoded in the front-end
- dispatched to the backend
- processed by individual execution units

Instructions

- are executed out-of-order

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Instructions

- are executed out-of-order
- wait until their dependencies are ready

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Instructions

- are executed out-of-order
- wait until their dependencies are ready
  - Later instructions might execute prior earlier instructions

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Instructions

- are executed out-of-order
- wait until their dependencies are ready
  - Later instructions might execute prior earlier instructions
- retire in-order

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Instructions

- are executed out-of-order
- wait until their dependencies are ready
  - Later instructions might execute prior earlier instructions
- retire in-order
  - State becomes architecturally visible

Instructions

- are executed out-of-order
- wait until their dependencies are ready
  - Later instructions might execute prior earlier instructions
- retire in-order
  - State becomes architecturally visible
- Exceptions are checked during retirement

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Instructions

- are executed out-of-order
- wait until their dependencies are ready
  - Later instructions might execute prior earlier instructions
- retire in-order
  - State becomes architecturally visible
- Exceptions are checked during retirement
  - Flush pipeline and recover state

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**The state does not become architecturally visible but ...**

# The state does not become architecturally visible but ...

- New code

```
*(volatile char*) 0;
array[84 * 4096] = 0;
```

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- New code

```
*(volatile char*) 0;
array[84 * 4096] = 0;
```

- volatile because compiler was not happy

```
warning: statement with no effect [−Wunused−value]
          *(char*)0;
```

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- New code

```
*(volatile char*) 0;
array[84 * 4096] = 0;
```

- volatile because compiler was not happy

warning: statement with no effect [−Wunused−value]
            *(char*)0;

- Static code analyzer is still not happy

warning: Dereference of null pointer
            *(volatile char*)0;

- Flush+Reload over all pages of the array

- Flush+Reload over all pages of the array



- "Unreachable" code line was actually executed

- Flush+Reload over all pages of the array



- "Unreachable" code line was actually executed
- Exception was only thrown afterwards

- Out-of-order instructions leave microarchitectural traces

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Out-of-order instructions leave microarchitectural traces
  - We can see them for example in the cache

- Out-of-order instructions leave microarchitectural traces
    - We can see them for example in the cache
- Give such instructions a name: transient instructions

- Out-of-order instructions leave microarchitectural traces
  - We can see them for example in the cache
- Give such instructions a name: transient instructions
- We can indirectly observe the execution of transient instructions

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**December**

**3**

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Add another layer of indirection to test

```
char data = *(char*) 0xffffffff81a000e0;
array[data * 4096] = 0;
```

- Add another layer of indirection to test

```
char data = *(char*) 0xffffffff81a000e0;
array[data * 4096] = 0;
```

- Then check whether any part of `array` is cached

- Flush+Reload over all pages of the array



- Index of cache hit reveals data

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Flush+Reload over all pages of the array



- Index of cache hit reveals data
- Permission check is in some cases not fast enough

MELTDOWN

- Using out-of-order execution, we can read data at any address

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**MELTDOWN**

- Using out-of-order execution, we can read data at any address
- Index of cache hit reveals data

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**MELTDOWN**

- Using out-of-order execution, we can read data at any address
- Index of cache hit reveals data
- Permission check is in some cases not fast enough

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**MELTDOWN**

- Using out-of-order execution, we can read data at any address
- Index of cache hit reveals data
- Permission check is in some cases not fast enough
- Entire physical memory is typically accessible through kernel space

I SHIT YOU NOT

THERE WAS KERNEL MEMORY ALL OVER THE TERMINAL

**Demo**

- Basic Meltdown code leads to a crash (segfault)

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Basic Meltdown code leads to a crash (segfault)
- How to prevent the crash?

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Basic Meltdown code leads to a crash (segfault)
- How to prevent the crash?

Fault
Handling

Fault
Suppression

Fault
Prevention

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Intel TSX to suppress exceptions instead of signal handler

```
if(xbegin() == XBEGIN_STARTED) {
  char secret = *(char*) 0xffffffff81a000e0;
  array[secret * 4096] = 0;
  xend();
}

for (size_t i = 0; i < 256; i++) {
  if (flush_and_reload(array + i * 4096) == CACHE_HIT) {
    printf("%c\n", i);
  }
}
```

- Speculative execution to prevent exceptions

```
int speculate = rand() % 2;
size_t address = (0xffffffff81a000e0 * speculate) +
                 ((size_t)&zero * (1 - speculate));
if(!speculate) {
  char secret = *(char*) address;
  array[secret * 4096] = 0;
}

for (size_t i = 0; i < 256; i++) {
  if (flush_and_reload(array + i * 4096) == CACHE_HIT) {
    printf("%c\n", i);
  }
}
```

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Mastering C Pointers — Second Edition, Robert J. Traister

- Improve the performance with a NULL pointer dereference

- Improve the performance with a NULL pointer dereference

```
if(xbegin() == XBEGIN_STARTED) {
  *(volatile char*) 0;
  char secret = *(char*) 0xffffffff81a000e0;
  array[secret * 4096] = 0;
  xend();
}
```

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

SO YOU ARE TELLING ME

YOU CAN DUMP THE MEMORY STORED IN L1?

- Assumed that one can only read data stored in the L1 with Meltdown

- Assumed that one can only read data stored in the L1 with Meltdown
- Experiment where a thread flushes the value constantly and a thread on a different core reloads the value

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Assumed that one can only read data stored in the L1 with Meltdown
- Experiment where a thread flushes the value constantly and a thread on a different core reloads the value
  - Target data is not in the L1 cache of the attacking core

- Assumed that one can only read data stored in the L1 with Meltdown
- Experiment where a thread flushes the value constantly and a thread on a different core reloads the value
  - Target data is not in the L1 cache of the attacking core
- We can still leak the data at a lower reading rate

- Assumed that one can only read data stored in the L1 with Meltdown
- Experiment where a thread flushes the value constantly and a thread on a different core reloads the value
  - Target data is not in the L1 cache of the attacking core
- We can still leak the data at a lower reading rate
- Meltdown might implicitly cache the data

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

• Mark pages in page tables as UC (uncachable)

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Mark pages in page tables as UC (uncachable)
  - Every read or write operation will go to main memory

- Mark pages in page tables as UC (uncachable)
  - Every read or write operation will go to main memory
- If the attacker can trigger a legitimate load (system call, …) on the same CPU core, the data still can be leaked

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Mark pages in page tables as UC (uncachable)
  - Every read or write operation will go to main memory
- If the attacker can trigger a legitimate load (system call, …) on the same CPU core, the data still can be leaked
- Meltdown might read the data from one of the fill buffers

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Mark pages in page tables as UC (uncachable)
    - Every read or write operation will go to main memory
- If the attacker can trigger a legitimate load (system call, …) on the same CPU core, the data still can be leaked
- Meltdown might read the data from one of the fill buffers
    - as they are shared between threads running on the same core

**So you can dump the entire memory.**

So you can dump the entire memory. But it takes ages?

- Dumping the entire physical memory takes some time

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Dumping the entire physical memory takes some time
  - Not very practical in most scenarios

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Dumping the entire physical memory takes some time
  - Not very practical in most scenarios
- Can we mount more targeted attacks?

- Open-source utility for disk encryption

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

VeraCrypt

- Open-source utility for disk encryption
- Fork of TrueCrypt

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Open-source utility for disk encryption
- Fork of TrueCrypt
- Cryptographic keys are stored in RAM

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

VeraCrypt

- Open-source utility for disk encryption
- Fork of TrueCrypt
- Cryptographic keys are stored in RAM
  - With Meltdown, we can extract the keys from DRAM

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Demo

- De-randomize KASLR to access internal kernel structures

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- De-randomize KASLR to access internal kernel structures
- Locate a known value inside the kernel, e.g., Linux banner

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- De-randomize KASLR to access internal kernel structures
- Locate a known value inside the kernel, e.g., Linux banner
  - Start at the default address according to the symbol table of the running kernel

- De-randomize KASLR to access internal kernel structures
- Locate a known value inside the kernel, e.g., Linux banner
  - Start at the default address according to the symbol table of the running kernel
  - Linux KASLR has an entropy of 6 bits $\Rightarrow$ only 64 possible randomization offsets

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- De-randomize KASLR to access internal kernel structures
- Locate a known value inside the kernel, e.g., Linux banner
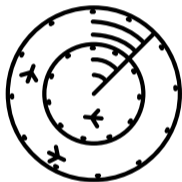    - Start at the default address according to the symbol table of the running kernel
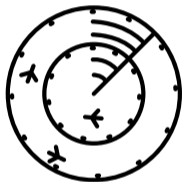    - Linux KASLR has an entropy of 6 bits $\Rightarrow$ only 64 possible randomization offsets
- Difference between the found address and the non-randomized base address is the randomization offset

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Linux manages all processes in a linked list

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Linux manages all processes in a linked list
- Head of the list is stored at `init_task` structure

- Linux manages all processes in a linked list
- Head of the list is stored at `init_task` structure
    - At a fixed offset that varies only among kernel builds

- Linux manages all processes in a linked list
- Head of the list is stored at `init_task` structure
  - At a fixed offset that varies only among kernel builds
- Each task list structure contains a pointer to the next task and
  - PID of the task
  - name of the task
  - Root of the multi-level page table

- Resolve physical address using paging structures

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Resolve physical address using paging structures
- Read the content using the direct-physical map

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Resolve physical address using paging structures
- Read the content using the direct-physical map
- Enumerate all mapped pages and dump entire process memory

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Resolve physical address using paging structures
- Read the content using the direct-physical map
- Enumerate all mapped pages and dump entire process memory
- Location of the key known, we can just dump the key directly

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

• `aeskeyfind` to extract AES keys from the memory dump

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- `aeskeyfind` to extract AES keys from the memory dump
- `pytruecrypt` to decrypt disk image using the extracted key

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology
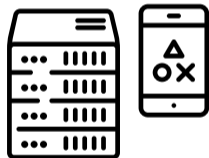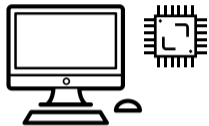
- `aeskeyfind` to extract AES keys from the memory dump
- `pytruecrypt` to decrypt disk image using the extracted key

- Affects every application that stores its secret in DRAM

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

# Who is affected?

• **Intel**: Almost every CPU

- **Intel**: Almost every CPU
- **AMD**: Seems not to be affected

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- **Intel**: Almost every CPU
- **AMD**: Seems not to be affected
- **ARM**: Only the Cortex-A75

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- **Intel**: Almost every CPU
- **AMD**: Seems not to be affected
- **ARM**: Only the Cortex-A75
- **IBM**: System Z, Power Architecture, POWER8 and POWER9

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- **Intel**: Almost every CPU
- **AMD**: Seems not to be affected
- **ARM**: Only the Cortex-A75
- **IBM**: System Z, Power Architecture, POWER8 and POWER9
- **Apple**: All Mac and iOS devices

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- But there are other CPU manufacturers as well ...

Samsung Galaxy S7

- Exynos Mongoose M1 CPU Architecture

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Samsung Galaxy S7

- Exynos Mongoose M1 CPU Architecture
  - Custom CPU core in the Exynos 8 Octa (8890)

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

SAMSUNG GALAXY S7

- Exynos Mongoose M1 CPU Architecture
  - Custom CPU core in the Exynos 8 Octa (8890)
  - Perceptron Branch Prediction

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Samsung Galaxy S7

- Exynos Mongoose M1 CPU Architecture
  - Custom CPU core in the Exynos 8 Octa (8890)
  - Perceptron Branch Prediction
  - Full Out-of-Order Instruction Execution

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Samsung Galaxy S7
- Exynos Mongoose M1 CPU Architecture
    - Custom CPU core in the Exynos 8 Octa (8890)
    - Perceptron Branch Prediction
    - Full Out-of-Order Instruction Execution
        - Full Out-of-Order loads and stores

**Demo**

Samsung Galaxy S7

- Luckily they already fixed it

Samsung Galaxy S7
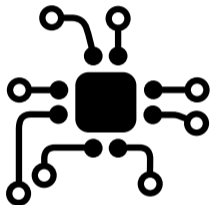
- Luckily they already fixed it
- With their latest update

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

SAMSUNG GALAXY S7

- Luckily they already fixed it
- With their latest update on July 10, 2018

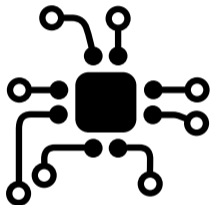Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- But there are other CPU manufacturers as well …
- …which are affected
- Need to evaluate the attack on other CPUs as well
- Notify the users …
- …and custom ROM developers, e.g., LineageOS

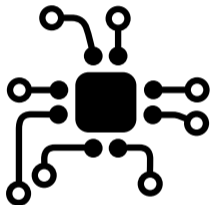Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology
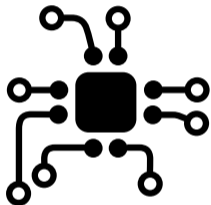
But wait, what about **privileged registers**?

- ARM found a closely related Meltdown variant

- ARM found a closely related Meltdown variant
- Read of system registers that are not accessible from current exception level

- ARM found a closely related Meltdown variant
- Read of system registers that are not accessible from current exception level
- ARM Cortex-A15, Cortex-A57 and Cortex-A72 are vulnerable

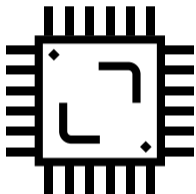Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- ARM found a closely related Meltdown variant
- Read of system registers that are not accessible from current exception level
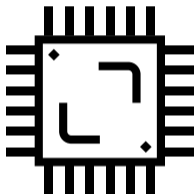- ARM Cortex-A15, Cortex-A57 and Cortex-A72 are vulnerable
- Impact: breaking KASLR and pointer authentication

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**Demo**

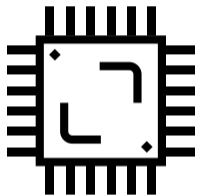- Intel is affected too (May 21, 2018)

- Intel is affected too (May 21, 2018)
  - Almost every CPU (Core i3/i5/i7, 2nd-8th Intel Core, Xeon, Atom, Pentium, ...)

- Intel is affected too (May 21, 2018)
  - Almost every CPU (Core i3/i5/i7, 2nd-8th Intel Core, Xeon, Atom, Pentium, ...)
- Rogue System Register Read (RSRE) (CVE-2018-3640)

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

No.

No.

- We read the data directly

No.
- We read the data directly
- We use a side channel internally for transmission

No.
- We read the data directly
- We use a side channel internally for transmission
$\rightarrow$ does not make the entire thing a side-channel attack

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**Is Meltdown a variant of Spectre? Is it speculative execution?**

black hat

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

No.

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

No.

- Often heard: "Meltdown is speculating beyond faulting instructions"

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

No.

- Often heard: "Meltdown is speculating beyond faulting instructions"
- $\rightarrow$ That's not speculative execution

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

No.

- Often heard: "Meltdown is speculating beyond faulting instructions"
- → That's not speculative execution
- "Speculating beyond faulting instructions" - not even the actual problem

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

No.

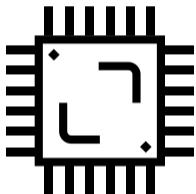- Often heard: "Meltdown is speculating beyond faulting instructions"
- → That's not speculative execution
- "Speculating beyond faulting instructions" - not even the actual problem
- AMD does that - but is not affected!

**Is Meltdown a variant of Spectre? Is it speculative execution?**

black hat

No.

- Often heard: "Meltdown is speculating beyond faulting instructions"
- → That's not speculative execution
- "Speculating beyond faulting instructions" - not even the actual problem
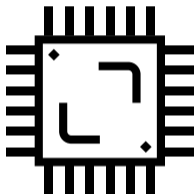- AMD does that - but is not affected!
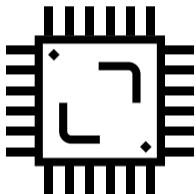- → Actual problem: fetching & using real values for instructions after faulting ones

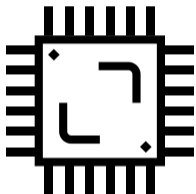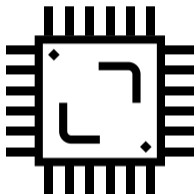How can this all be fixed?

- Problem is rooted in hardware

- Problem is rooted in hardware
- Race condition between the memory fetch and corresponding permission check
  - Serialize both of them

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Problem is rooted in hardware
- Race condition between the memory fetch and corresponding permission check
  - Serialize both of them
- Hard split of user space and kernel space
  - New bit in control register

- Problem is rooted in hardware
- Race condition between the memory fetch and corresponding permission check
  - Serialize both of them
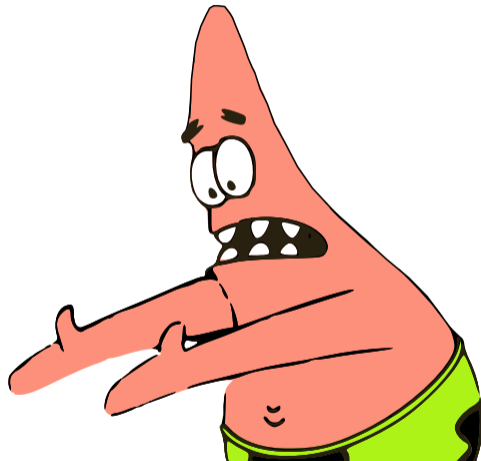- Hard split of user space and kernel space
  - New bit in control register

- Fix the hardware → long-term solution

- Problem is rooted in hardware
- Race condition between the memory fetch and corresponding permission check
  - Serialize both of them
- Hard split of user space and kernel space
  - New bit in control register

- Fix the hardware → long-term solution
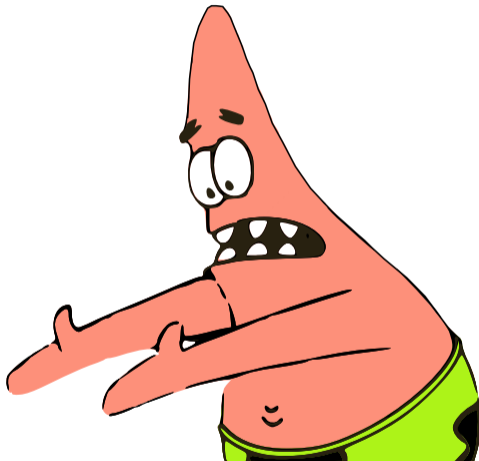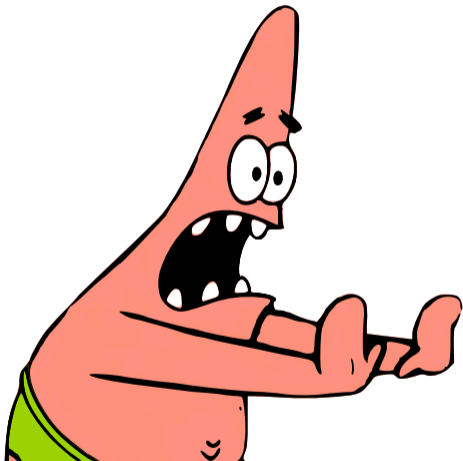- Can we fix it in software?

- Kernel addresses in user space are a problem

- Kernel addresses in user space are a problem
- Why don't we take the kernel addresses…

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- …and remove them if not needed?

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology
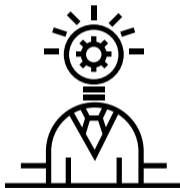
- …and remove them if not needed?
- User accessible check in hardware is not reliable

- Unmap the kernel in user space

- Unmap the kernel in user space
- Kernel addresses are then no longer present

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Unmap the kernel in user space
- Kernel addresses are then no longer present
- Memory which is not mapped cannot be accessed at all

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Userspace

Kernelspace

Applications

Operating System

Memory

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Kernel View

User View

- We published KAISER in May 2017 ...

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- We published KAISER in May 2017 …
- …as a countermeasure against other side-channel attacks

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- We published KAISER in May 2017 …
- …as a countermeasure against other side-channel attacks
- Inadvertently defeats Meltdown as well

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- We published KAISER in May 2017 …
- …as a countermeasure against other side-channel attacks
- Inadvertently defeats Meltdown as well
- PoC implementation for the Linux kernel

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Hardware interrupt while running in user mode
  - Kernel needs to deal with interrupt but does not exist anymore in address space
  - Traps, NMI, system calls, …
- Must map some kernel code in user space

- Need to update CR3 in order to switch to other address space

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Need to update CR3 in order to switch to other address space
- How can we do this efficiently?

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Need to update CR3 in order to switch to other address space
- How can we do this efficiently?
- Instead of one PGD, two PGDs are allocated
  - 8k in size and 8k aligned

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Need to update CR3 in order to switch to other address space
- How can we do this efficiently?
- Instead of one PGD, two PGDs are allocated
  - 8k in size and 8k aligned
- **Trick**: Just flip bit 12 in the pointer to swap between both halves

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Need to update CR3 in order to switch to other address space
- How can we do this efficiently?
- Instead of one PGD, two PGDs are allocated
  - 8k in size and 8k aligned
- **Trick**: Just flip bit 12 in the pointer to swap between both halves

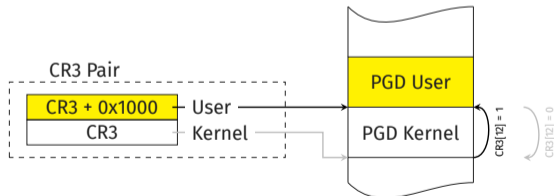Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Need to update CR3 in order to switch to other address space
- How can we do this efficiently?
- Instead of one PGD, two PGDs are allocated
    - 8k in size and 8k aligned
- **Trick**: Just flip bit 12 in the pointer to swap between both halves

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Intel and others improved KAISER

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Intel and others improved KAISER
- Merged it into upstream as KPTI (Kernel Page-table Isolation)

- Intel and others improved KAISER
- Merged it into upstream as KPTI (Kernel Page-table Isolation)
- Kernel patches are available for arm64 as well

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Apple released updates in iOS 11.2, macOS 10.13.2 and tvOS 11.2 to mitigate Meltdown
- Boot option: `-no-shared-cr3`
  - Unmaps the user space while running in kernel mode
  - But not vice versa

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Kernel Virtual Address (KVA) Shadow
- Meltdown Mitigation for Microsoft Windows

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Introducing such a fundamental change to the operating system is extremely challenging
- Our PoC implementation contained many bugs as well

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Discovered by Ulf Frisk (@ulffrisk) in the 2018-02 security update
  - CVE-2018-1038
- Modified the PML4 entry of `0x1ed` to allow to access page from user-mode
- On Windows 7 and Server 2018 R2: Self-Referencing Entry
- Allows to read and modify entire physical memory

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

**What now?**

- More attacks exploiting performance optimizations in hardware
  - New variants are disclosed frequently

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

A unique chance to

- rethink processor design
- grow up, like other fields (car industry, construction industry)

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

You can find our proof-of-concept implementation on:

- `https://github.com/IAIK/meltdown`

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

- Underestimated microarchitectural attacks for a long time
- Meltdown exploit performance optimizations
  - Allow to leak arbitrary memory
- Countermeasures come with a performance impact
- Find trade-offs between security and performance

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

Moritz Lipp (@mlqxyz)

Michael Schwarz (@misc0110)

Daniel Gruss (@lavados)

# Meltdown
## Basics, Details, Consequences

# References

R. Grisenthwaite. Cache Speculation Side-channels. 2018.

D. Gruss, M. Lipp, M. Schwarz, R. Fellner, C. Maurice, and S. Mangard. KASLR is Dead: Long Live KASLR. In: ESSoS. 2017.

Intel. Intel Analysis of Speculative Execution Side Channels. Jan. 2018. URL:
`https://newsroom.intel.com/wp-content/uploads/sites/11/2018/01/Intel-Analysis-of-Speculative-Execution-Side-Channels.pdf`.

K. Johnson. KVA Shadow: Mitigating Meltdown on Windows. Mar. 2018. URL:
`https://blogs.technet.microsoft.com/srd/2018/03/23/kva-shadow-mitigating-meltdown-on-windows/`.

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology

M. Lipp, M. Schwarz, D. Gruss, T. Prescher, W. Haas, A. Fogh, J. Horn, S. Mangard, P. Kocher, D. Genkin, Y. Yarom, and M. Hamburg. Meltdown: Reading Kernel Memory from User Space. In: USENIX Security Symposium. 2018.

LWN. The current state of kernel page-table isolation. 2017. URL: https://lwn.net/SubscriberLink/741878/eb6c9d3913d7cb2b/.

Y. Yarom and K. Falkner. Flush+Reload: a High Resolution, Low Noise, L3 Cache Side-Channel Attack. In: USENIX Security Symposium. 2014.

Moritz Lipp, Michael Schwarz, Daniel Gruss | Graz University of Technology