

Meltdown: Reading Kernel Memory from User Space

Moritz Lipp¹, Michael Schwarz¹, Daniel Gruss¹, Thomas Prescher²,
Werner Haas², Anders Fogh³, Jann Horn⁴, Stefan Mangard¹,
Paul Kocher⁵, Daniel Genkin^{6,9}, Yuval Yarom⁷, Mike Hamburg⁸

¹*Graz University of Technology*, ²*Cyberus Technology GmbH*,

³*G-Data Advanced Analytics*, ⁴*Google Project Zero*,

⁵*Independent (www.paulkocher.com)*, ⁶*University of Michigan*,

⁷*University of Adelaide & Data61*, ⁸*Rambus, Cryptography Research Division*

Abstract

The security of computer systems fundamentally relies on memory isolation, e.g., kernel address ranges are marked as non-accessible and are protected from user access. In this paper, we present Meltdown. Meltdown exploits side effects of out-of-order execution on modern processors to read arbitrary kernel-memory locations including personal data and passwords. Out-of-order execution is an indispensable performance feature and present in a wide range of modern processors. The attack is independent of the operating system, and it does not rely on any software vulnerabilities. Meltdown breaks all security guarantees provided by address space isolation as well as paravirtualized environments and, thus, every security mechanism building upon this foundation. On affected systems, Meltdown enables an adversary to read memory of other processes or virtual machines in the cloud without any permissions or privileges, affecting millions of customers and virtually every user of a personal computer. We show that the KAISER defense mechanism for KASLR has the important (but inadvertent) side effect of impeding Meltdown. We stress that KAISER must be deployed immediately to prevent large-scale exploitation of this severe information leakage.

1 Introduction

A central security feature of today’s operating systems is memory isolation. Operating systems ensure that user programs cannot access each other’s memory or kernel memory. This isolation is a cornerstone of our computing environments and allows running multiple applications at the same time on personal devices or executing processes of multiple users on a single machine in the cloud.

On modern processors, the isolation between the kernel and user processes is typically realized by a supervi-

sor bit of the processor that defines whether a memory page of the kernel can be accessed or not. The basic idea is that this bit can only be set when entering kernel code and it is cleared when switching to user processes. This hardware feature allows operating systems to map the kernel into the address space of every process and to have very efficient transitions from the user process to the kernel, e.g., for interrupt handling. Consequently, in practice, there is no change of the memory mapping when switching from a user process to the kernel.

In this work, we present Meltdown¹⁰. Meltdown is a novel attack that allows overcoming memory isolation completely by providing a simple way for any user process to read the entire kernel memory of the machine it executes on, including all physical memory mapped in the kernel region. Meltdown does not exploit any software vulnerability, *i.e.*, it works on all major operating systems. Instead, Meltdown exploits side-channel information available on most modern processors, e.g., modern Intel microarchitectures since 2010 and potentially on other CPUs of other vendors.

While side-channel attacks typically require very specific knowledge about the target application and are tailored to only leak information about its secrets, Meltdown allows an adversary who can run code on the vulnerable processor to obtain a dump of the entire kernel address space, including any mapped physical memory. The root cause of the simplicity and strength of Meltdown are side effects caused by *out-of-order execution*.

Out-of-order execution is an important performance feature of today’s processors in order to overcome latencies of busy execution units, e.g., a memory fetch unit needs to wait for data arrival from memory. Instead of stalling the execution, modern processors run operations

⁹Work was partially done while the author was affiliated to University of Pennsylvania and University of Maryland.

¹⁰Using the practice of responsible disclosure, disjoint groups of authors of this paper provided preliminary versions of our results to partially overlapping groups of CPU vendors and other affected companies. In coordination with industry, the authors participated in an embargo of the results. Meltdown is documented under CVE-2017-5754.

out-of-order i.e., they look ahead and schedule subsequent operations to idle execution units of the core. However, such operations often have unwanted side-effects, e.g., timing differences [55, 63, 23] can leak information from both sequential and out-of-order execution.

From a security perspective, one observation is particularly significant: vulnerable out-of-order CPUs allow an unprivileged process to load data from a privileged (kernel or physical) address into a temporary CPU register. Moreover, the CPU even performs further computations based on this register value, e.g., access to an array based on the register value. By simply discarding the results of the memory lookups (e.g., the modified register states), if it turns out that an instruction should not have been executed, the processor ensures correct program execution. Hence, on the architectural level (e.g., the abstract definition of how the processor should perform computations) no security problem arises.

However, we observed that out-of-order memory lookups influence the cache, which in turn can be detected through the cache side channel. As a result, an attacker can dump the entire kernel memory by reading privileged memory in an out-of-order execution stream, and transmit the data from this elusive state via a microarchitectural covert channel (e.g., Flush+Reload) to the outside world. On the receiving end of the covert channel, the register value is reconstructed. Hence, on the microarchitectural level (e.g., the actual hardware implementation), there is an exploitable security problem.

Meltdown breaks all security guarantees provided by the CPU’s memory isolation capabilities. We evaluated the attack on modern desktop machines and laptops, as well as servers in the cloud. Meltdown allows an unprivileged process to read data mapped in the kernel address space, including the entire physical memory on Linux, Android and OS X, and a large fraction of the physical memory on Windows. This may include the physical memory of other processes, the kernel, and in the case of kernel-sharing sandbox solutions (e.g., Docker, LXC) or Xen in paravirtualization mode, the memory of the kernel (or hypervisor), and other co-located instances. While the performance heavily depends on the specific machine, e.g., processor speed, TLB and cache sizes, and DRAM speed, we can dump arbitrary kernel and physical memory with 3.2 KB/s to 503 KB/s. Hence, an enormous number of systems are affected.

The countermeasure KAISER [20], developed initially to prevent side-channel attacks targeting KASLR, inadvertently protects against Meltdown as well. Our evaluation shows that KAISER prevents Meltdown to a large extent. Consequently, we stress that it is of utmost importance to deploy KAISER on all operating systems immediately. Fortunately, during a responsible disclosure window, the three major operating systems (Windows,

Linux, and OS X) implemented variants of KAISER and recently rolled out these patches.

Meltdown is distinct from the Spectre Attacks [40] in several ways, notably that Spectre requires tailoring to the victim process’s software environment, but applies more broadly to CPUs and is not mitigated by KAISER.

Contributions. The contributions of this work are:

1. We describe out-of-order execution as a new, extremely powerful, software-based side channel.
2. We show how out-of-order execution can be combined with a microarchitectural covert channel to transfer the data from an elusive state to a receiver on the outside.
3. We present an end-to-end attack combining out-of-order execution with exception handlers or TSX, to read arbitrary physical memory without any permissions or privileges, on laptops, desktop machines, mobile phones and on public cloud machines.
4. We evaluate the performance of Meltdown and the effects of KAISER on it.

Outline. The remainder of this paper is structured as follows: In Section 2, we describe the fundamental problem which is introduced with out-of-order execution. In Section 3, we provide a toy example illustrating the side channel Meltdown exploits. In Section 4, we describe the building blocks of Meltdown. We present the full attack in Section 5. In Section 6, we evaluate the performance of the Meltdown attack on several different systems and discuss its limitations. In Section 7, we discuss the effects of the software-based KAISER countermeasure and propose solutions in hardware. In Section 8, we discuss related work and conclude our work in Section 9.

2 Background

In this section, we provide background on out-of-order execution, address translation, and cache attacks.

2.1 Out-of-order execution

Out-of-order execution is an optimization technique that allows maximizing the utilization of all execution units of a CPU core as exhaustive as possible. Instead of processing instructions strictly in the sequential program order, the CPU executes them as soon as all required resources are available. While the execution unit of the current operation is occupied, other execution units can run ahead. Hence, instructions can be run in parallel as long as their results follow the architectural definition.

In practice, CPUs supporting out-of-order execution allow running operations *speculatively* to the extent that

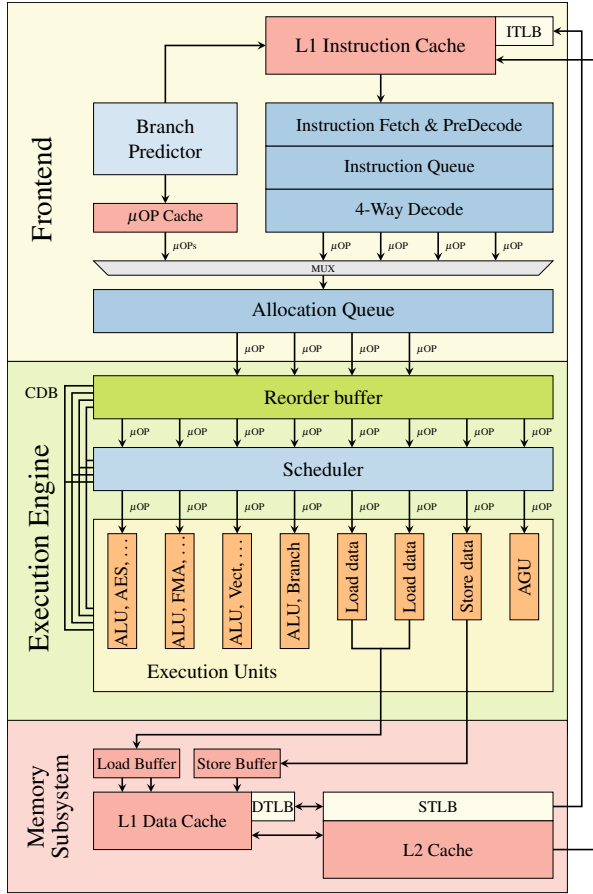


Figure 1: Simplified illustration of a single core of the Intel's Skylake microarchitecture. Instructions are decoded into μ OPs and executed out-of-order in the execution engine by individual execution units.

the processor's out-of-order logic processes instructions before the CPU is certain that the instruction will be needed and committed. In this paper, we refer to speculative execution in a more restricted meaning, where it refers to an instruction sequence following a branch, and use the term out-of-order execution to refer to any way of getting an operation executed before the processor has committed the results of all prior instructions.

In 1967, Tomasulo [61] developed an algorithm that enabled dynamic scheduling of instructions to allow out-of-order execution. Tomasulo [61] introduced a unified reservation station that allows a CPU to use a data value as it has been computed instead of storing it in a register and re-reading it. The reservation station renames registers to allow instructions that operate on the same physical registers to use the last logical one to solve read-after-write (RAW), write-after-read (WAR) and write-after-write (WAW) hazards. Furthermore, the reservation unit connects all execution units via a common data

bus (CDB). If an operand is not available, the reservation unit can listen on the CDB until it is available and then directly begin the execution of the instruction.

On the Intel architecture, the pipeline consists of the front-end, the execution engine (back-end) and the memory subsystem [32]. x86 instructions are fetched by the front-end from memory and decoded to micro-operations (μ OPs) which are continuously sent to the execution engine. Out-of-order execution is implemented within the execution engine as illustrated in Figure 1. The *Reorder Buffer* is responsible for register allocation, register renaming and retiring. Additionally, other optimizations like move elimination or the recognition of zeroing idioms are directly handled by the reorder buffer. The μ OPs are forwarded to the *Unified Reservation Station* (Scheduler) that queues the operations on exit ports that are connected to *Execution Units*. Each execution unit can perform different tasks like ALU operations, AES operations, address generation units (AGU) or memory loads and stores. AGUs, as well as load and store execution units, are directly connected to the memory subsystem to process its requests.

Since CPUs usually do not run linear instruction streams, they have branch prediction units that are used to obtain an educated guess of which instruction is executed next. Branch predictors try to determine which direction of a branch is taken before its condition is actually evaluated. Instructions that lie on that path and do not have any dependencies can be executed in advance and their results immediately used if the prediction was correct. If the prediction was incorrect, the reorder buffer allows to rollback to a sane state by clearing the reorder buffer and re-initializing the unified reservation station.

There are various approaches to predict a branch: With static branch prediction [28], the outcome is predicted solely based on the instruction itself. Dynamic branch prediction [8] gathers statistics at run-time to predict the outcome. One-level branch prediction uses a 1-bit or 2-bit counter to record the last outcome of a branch [45]. Modern processors often use two-level adaptive predictors [64] with a history of the last n outcomes, allowing to predict regularly recurring patterns. More recently, ideas to use neural branch prediction [62, 38, 60] have been picked up and integrated into CPU architectures [9].

2.2 Address Spaces

To isolate processes from each other, CPUs support virtual address spaces where virtual addresses are translated to physical addresses. A virtual address space is divided into a set of pages that can be individually mapped to physical memory through a multi-level page translation table. The translation tables define the actual virtual to physical mapping and also protection properties that

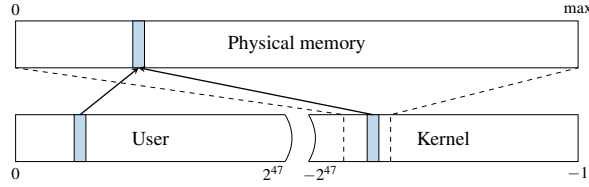


Figure 2: The physical memory is directly mapped in the kernel at a certain offset. A physical address (blue) which is mapped accessible to the user space is also mapped in the kernel space through the direct mapping.

are used to enforce privilege checks, such as readable, writable, executable and user-accessible. The currently used translation table is held in a special CPU register. On each context switch, the operating system updates this register with the next process’ translation table address in order to implement per-process virtual address spaces. Because of that, each process can only reference data that belongs to its virtual address space. Each virtual address space itself is split into a user and a kernel part. While the user address space can be accessed by the running application, the kernel address space can only be accessed if the CPU is running in privileged mode. This is enforced by the operating system disabling the user-accessible property of the corresponding translation tables. The kernel address space does not only have memory mapped for the kernel’s own usage, but it also needs to perform operations on user pages, e.g., filling them with data. Consequently, the entire physical memory is typically mapped in the kernel. On Linux and OS X, this is done via a direct-physical map, *i.e.*, the entire physical memory is directly mapped to a pre-defined virtual address (cf. Figure 2).

Instead of a direct-physical map, Windows maintains a multiple so-called *paged pools*, *non-paged pools*, and the *system cache*. These pools are virtual memory regions in the kernel address space mapping physical pages to virtual addresses which are either required to remain in the memory (*non-paged pool*) or can be removed from the memory because a copy is already stored on the disk (*paged pool*). The *system cache* further contains mappings of all file-backed pages. Combined, these memory pools will typically map a large fraction of the physical memory into the kernel address space of every process.

The exploitation of memory corruption bugs often requires knowledge of addresses of specific data. In order to impede such attacks, address space layout randomization (ASLR) has been introduced as well as non-executable stacks and stack canaries. To protect the kernel, kernel ASLR (KASLR) randomizes the offsets where drivers are located on every boot, making attacks harder as they now require to guess the location of kernel

data structures. However, side-channel attacks allow to detect the exact location of kernel data structures [21, 29, 37] or derandomize ASLR in JavaScript [16]. A combination of a software bug and the knowledge of these addresses can lead to privileged code execution.

2.3 Cache Attacks

In order to speed-up memory accesses and address translation, the CPU contains small memory buffers, called caches, that store frequently used data. CPU caches hide slow memory access latencies by buffering frequently used data in smaller and faster internal memory. Modern CPUs have multiple levels of caches that are either private per core or shared among them. Address space translation tables are also stored in memory and, thus, also cached in the regular caches.

Cache side-channel attacks exploit timing differences that are introduced by the caches. Different cache attack techniques have been proposed and demonstrated in the past, including Evict+Time [55], Prime+Probe [55, 56], and Flush+Reload [63]. Flush+Reload attacks work on a single cache line granularity. These attacks exploit the shared, inclusive last-level cache. An attacker frequently flushes a targeted memory location using the `clflush` instruction. By measuring the time it takes to reload the data, the attacker determines whether data was loaded into the cache by another process in the meantime. The Flush+Reload attack has been used for attacks on various computations, e.g., cryptographic algorithms [63, 36, 4], web server function calls [65], user input [23, 47, 58], and kernel addressing information [21].

A special use case of a side-channel attack is a covert channel. Here the attacker controls both, the part that induces the side effect, and the part that measures the side effect. This can be used to leak information from one security domain to another, while bypassing any boundaries existing on the architectural level or above. Both Prime+Probe and Flush+Reload have been used in high-performance covert channels [48, 52, 22].

3 A Toy Example

In this section, we start with a toy example, *i.e.*, a simple code snippet, to illustrate that out-of-order execution can change the microarchitectural state in a way that leaks information. However, despite its simplicity, it is used as a basis for Section 4 and Section 5, where we show how this change in state can be exploited for an attack.

Listing 1 shows a simple code snippet first raising an (unhandled) exception and then accessing an array. The property of an exception is that the control flow does not continue with the code after the exception, but jumps to an exception handler in the operating system. Regardless

```

1 raise_exception();
2 // the line below is never reached
3 access(probe_array[data * 4096]);

```

Listing 1: A toy example to illustrate side-effects of out-of-order execution.

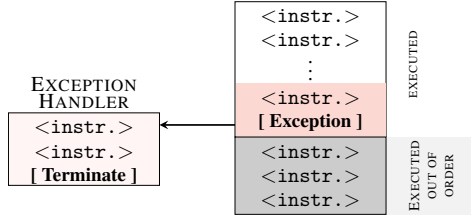


Figure 3: If an executed instruction causes an exception, diverting the control flow to an exception handler, the subsequent instruction must not be executed. Due to out-of-order execution, the subsequent instructions may already have been partially executed, but not retired. However, architectural effects of the execution are discarded.

of whether this exception is raised due to a memory access, e.g., by accessing an invalid address, or due to any other CPU exception, e.g., a division by zero, the control flow continues in the kernel and not with the next user space instruction.

Thus, our toy example cannot access the array in theory, as the exception immediately traps to the kernel and terminates the application. However, due to the out-of-order execution, the CPU might have already executed the following instructions as there is no dependency on the instruction triggering the exception. This is illustrated in Figure 3. Due to the exception, the instructions executed out of order are not retired and, thus, never have architectural effects.

Although the instructions executed out of order do not have any visible architectural effect on registers or memory, they have microarchitectural side effects. During the out-of-order execution, the referenced memory is fetched into a register and also stored in the cache. If the out-of-order execution has to be discarded, the register and memory contents are never committed. Nevertheless, the cached memory contents are kept in the cache. We can leverage a microarchitectural side-channel attack such as Flush+Reload [63], which detects whether a specific memory location is cached, to make this microarchitectural state visible. Other side channels can also detect whether a specific memory location is cached, including Prime+Probe [55, 48, 52], Evict+Reload [47], or Flush+Flush [22]. As Flush+Reload is the most accurate known cache side channel and is simple to implement, we do not consider any other side channel for this example.

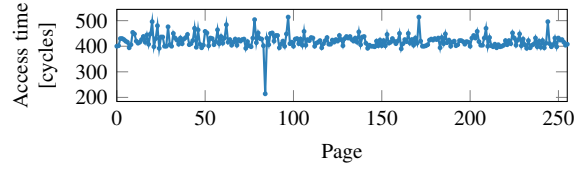


Figure 4: Even if a memory location is only accessed during out-of-order execution, it remains cached. Iterating over the 256 pages of probe_array shows one cache hit, exactly on the page that was accessed during the out-of-order execution.

Based on the value of data in this example, a different part of the cache is accessed when executing the memory access out of order. As data is multiplied by 4096, data accesses to probe_array are scattered over the array with a distance of 4 KB (assuming an 1 B data type for probe_array). Thus, there is an injective mapping from the value of data to a memory page, *i.e.*, different values for data never result in an access to the same page. Consequently, if a cache line of a page is cached, we know the value of data. The spreading over pages eliminates false positives due to the prefetcher, as the prefetcher cannot access data across page boundaries [32].

Figure 4 shows the result of a Flush+Reload measurement iterating over all pages, after executing the out-of-order snippet with data = 84. Although the array access should not have happened due to the exception, we can clearly see that the index which would have been accessed is cached. Iterating over all pages (e.g., in the exception handler) shows only a cache hit for page 84. This shows that even instructions which are never actually executed, change the microarchitectural state of the CPU. Section 4 modifies this toy example not to read a value but to leak an inaccessible secret.

4 Building Blocks of the Attack

The toy example in Section 3 illustrated that side-effects of out-of-order execution can modify the microarchitectural state to leak information. While the code snippet reveals the data value passed to a cache-side channel, we want to show how this technique can be leveraged to leak otherwise inaccessible secrets. In this section, we want to generalize and discuss the necessary building blocks to exploit out-of-order execution for an attack.

The adversary targets a secret value that is kept somewhere in physical memory. Note that register contents are also stored in memory upon context switches, *i.e.*, they are also stored in physical memory. As described in Section 2.2, the address space of every process typically includes the entire user space, as well as the entire kernel

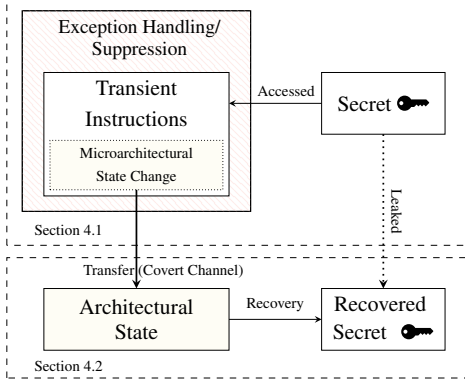


Figure 5: The Meltdown attack uses exception handling or suppression, e.g., TSX, to run a series of transient instructions. These transient instructions obtain a (persistent) secret value and change the microarchitectural state of the processor based on this secret value. This forms the sending part of a microarchitectural covert channel. The receiving side reads the microarchitectural state, making it architectural and recovers the secret value.

space, which typically also has all physical memory (in-use) mapped. However, these memory regions are only accessible in privileged mode (cf. Section 2.2).

In this work, we demonstrate leaking secrets by bypassing the privileged-mode isolation, giving an attacker full read access to the entire kernel space, including any physical memory mapped and, thus, the physical memory of any other process and the kernel. Note that Kocher et al. [40] pursue an orthogonal approach, called Spectre Attacks, which trick speculatively executed instructions into leaking information that the victim process is authorized to access. As a result, Spectre Attacks lack the privilege escalation aspect of Meltdown and require tailoring to the victim process’s software environment, but apply more broadly to CPUs that support speculative execution and are not prevented by KAISER.

The full Meltdown attack consists of two building blocks, as illustrated in Figure 5. The first building block of Meltdown is to make the CPU execute one or more instructions that would never occur in the executed path. In the toy example (cf. Section 3), this is an access to an array, which would normally never be executed, as the previous instruction always raises an exception. We call such an instruction, which is executed out of order and leaving measurable side effects, a *transient instruction*. Furthermore, we call any sequence of instructions containing at least one transient instruction a transient instruction sequence.

In order to leverage transient instructions for an attack, the transient instruction sequence must utilize a secret value that an attacker wants to leak. Section 4.1 describes

building blocks to run a transient instruction sequence with a dependency on a secret value.

The second building block of Meltdown is to transfer the microarchitectural side effect of the transient instruction sequence to an architectural state to further process the leaked secret. Thus, the second building described in Section 4.2 describes building blocks to transfer a microarchitectural side effect to an architectural state using a covert channel.

4.1 Executing Transient Instructions

The first building block of Meltdown is the execution of transient instructions. Transient instructions occur all the time, as the CPU continuously runs ahead of the current instruction to minimize the experienced latency and, thus, to maximize the performance (cf. Section 2.1). Transient instructions introduce an exploitable side channel if their operation depends on a secret value. We focus on addresses that are mapped within the attacker’s process, *i.e.*, the user-accessible user space addresses as well as the user-inaccessible kernel space addresses. Note that attacks targeting code that is executed within the context (*i.e.*, address space) of another process are possible [40], but out of scope in this work, since all physical memory (including the memory of other processes) can be read through the kernel address space regardless.

Accessing user-inaccessible pages, such as kernel pages, triggers an exception which generally terminates the application. If the attacker targets a secret at a user-inaccessible address, the attacker has to cope with this exception. We propose two approaches: With *exception handling*, we catch the exception effectively occurring after executing the transient instruction sequence, and with *exception suppression*, we prevent the exception from occurring at all and instead redirect the control flow after executing the transient instruction sequence. We discuss these approaches in detail in the following.

Exception handling. A trivial approach is to fork the attacking application before accessing the invalid memory location that terminates the process and only access the invalid memory location in the child process. The CPU executes the transient instruction sequence in the child process before crashing. The parent process can then recover the secret by observing the microarchitectural state, e.g., through a side-channel.

It is also possible to install a signal handler that is executed when a certain exception occurs, e.g., a segmentation fault. This allows the attacker to issue the instruction sequence and prevent the application from crashing, reducing the overhead as no new process has to be created.

Exception suppression. A different approach to deal with exceptions is to prevent them from being raised in the first place. Transactional memory allows to group memory accesses into one seemingly atomic operation, giving the option to roll-back to a previous state if an error occurs. If an exception occurs within the transaction, the architectural state is reset, and the program execution continues without disruption.

Furthermore, speculative execution issues instructions that might not occur on the executed code path due to a branch misprediction. Such instructions depending on a preceding conditional branch can be speculatively executed. Thus, the invalid memory access is put within a speculative instruction sequence that is only executed if a prior branch condition evaluates to true. By making sure that the condition never evaluates to true in the executed code path, we can suppress the occurring exception as the memory access is only executed speculatively. This technique may require sophisticated training of the branch predictor. Kocher et al. [40] pursue this approach in orthogonal work, since this construct can frequently be found in code of other processes.

4.2 Building a Covert Channel

The second building block of Meltdown is the transfer of the microarchitectural state, which was changed by the transient instruction sequence, into an architectural state (cf. Figure 5). The transient instruction sequence can be seen as the sending end of a microarchitectural covert channel. The receiving end of the covert channel receives the microarchitectural state change and deduces the secret from the state. Note that the receiver is not part of the transient instruction sequence and can be a different thread or even a different process e.g., the parent process in the fork-and-crash approach.

We leverage techniques from cache attacks, as the cache state is a microarchitectural state which can be reliably transferred into an architectural state using various techniques [55, 63, 22]. Specifically, we use Flush+Reload [63], as it allows to build a fast and low-noise covert channel. Thus, depending on the secret value, the transient instruction sequence (cf. Section 4.1) performs a regular memory access, e.g., as it does in the toy example (cf. Section 3).

After the transient instruction sequence accessed an accessible address, *i.e.*, this is the sender of the covert channel; the address is cached for subsequent accesses. The receiver can then monitor whether the address has been loaded into the cache by measuring the access time to the address. Thus, the sender can transmit a ‘1’-bit by accessing an address which is loaded into the monitored cache, and a ‘0’-bit by not accessing such an address.

Using multiple different cache lines, as in our toy example in Section 3, allows to transmit multiple bits at once. For every of the 256 different byte values, the sender accesses a different cache line. By performing a Flush+Reload attack on all of the 256 possible cache lines, the receiver can recover a full byte instead of just one bit. However, since the Flush+Reload attack takes much longer (typically several hundred cycles) than the transient instruction sequence, transmitting only a single bit at once is more efficient. The attacker can simply do that by shifting and masking the secret value accordingly.

Note that the covert channel is not limited to microarchitectural states which rely on the cache. Any microarchitectural state which can be influenced by an instruction (sequence) and is observable through a side channel can be used to build the sending end of a covert channel. The sender could, for example, issue an instruction (sequence) which occupies a certain execution port such as the ALU to send a ‘1’-bit. The receiver measures the latency when executing an instruction (sequence) on the same execution port. A high latency implies that the sender sends a ‘1’-bit, whereas a low latency implies that sender sends a ‘0’-bit. The advantage of the Flush+Reload cache covert channel is the noise resistance and the high transmission rate [22]. Furthermore, the leakage can be observed from any CPU core [63], *i.e.*, rescheduling events do not significantly affect the covert channel.

5 Meltdown

In this section, we present Meltdown, a powerful attack allowing to read arbitrary physical memory from an unprivileged user program, comprised of the building blocks presented in Section 4. First, we discuss the attack setting to emphasize the wide applicability of this attack. Second, we present an attack overview, showing how Meltdown can be mounted on both Windows and Linux on personal computers, on Android on mobile phones as well as in the cloud. Finally, we discuss a concrete implementation of Meltdown allowing to dump arbitrary kernel memory with 3.2 KB/s to 503 KB/s.

Attack setting. In our attack, we consider personal computers and virtual machines in the cloud. In the attack scenario, the attacker has arbitrary unprivileged code execution on the attacked system, *i.e.*, the attacker can run any code with the privileges of a normal user. However, the attacker has no physical access to the machine. Furthermore, we assume that the system is fully protected with state-of-the-art software-based defenses such as ASLR and KASLR as well as CPU features like SMAP, SMEP, NX, and PXN. Most importantly, we assume a completely bug-free operating system, thus, no

```

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

```

Listing 2: The core of Meltdown. An inaccessible kernel address is moved to a register, raising an exception. Subsequent instructions are executed out of order before the exception is raised, leaking the data from the kernel address through the indirect memory access.

software vulnerability exists that can be exploited to gain kernel privileges or leak information. The attacker targets secret user data, e.g., passwords and private keys, or any other valuable information.

5.1 Attack Description

Meltdown combines the two building blocks discussed in Section 4. First, an attacker makes the CPU execute a transient instruction sequence which uses an inaccessible secret value stored somewhere in physical memory (cf. Section 4.1). The transient instruction sequence acts as the transmitter of a covert channel (cf. Section 4.2), ultimately leaking the secret value to the attacker.

Meltdown consists of 3 steps:

- Step 1** The content of an attacker-chosen memory location, which is inaccessible to the attacker, is loaded into a register.
- Step 2** A transient instruction accesses a cache line based on the secret content of the register.
- Step 3** The attacker uses Flush+Reload to determine the accessed cache line and hence the secret stored at the chosen memory location.

By repeating these steps for different memory locations, the attacker can dump the kernel memory, including the entire physical memory.

Listing 2 shows the basic implementation of the transient instruction sequence and the sending part of the covert channel, using x86 assembly instructions. Note that this part of the attack could also be implemented entirely in higher level languages like C. In the following, we will discuss each step of Meltdown and the corresponding code line in Listing 2.

Step 1: Reading the secret. To load data from the main memory into a register, the data in the main memory is referenced using a virtual address. In parallel to translating a virtual address into a physical address, the CPU also checks the permission bits of the virtual ad-

dress, *i.e.*, whether this virtual address is user accessible or only accessible by the kernel. As already discussed in Section 2.2, this hardware-based isolation through a permission bit is considered secure and recommended by the hardware vendors. Hence, modern operating systems always map the entire kernel into the virtual address space of every user process.

As a consequence, all kernel addresses lead to a valid physical address when translating them, and the CPU can access the content of such addresses. The only difference to accessing a user space address is that the CPU raises an exception as the current permission level does not allow to access such an address. Hence, the user space cannot simply read the contents of such an address. However, Meltdown exploits the out-of-order execution of modern CPUs, which still executes instructions in the small time window between the illegal memory access and the raising of the exception.

In line 4 of Listing 2, we load the byte value located at the target kernel address, stored in the RCX register, into the least significant byte of the RAX register represented by AL. As explained in more detail in Section 2.1, the MOV instruction is fetched by the core, decoded into μ OPs, allocated, and sent to the reorder buffer. There, architectural registers (e.g., RAX and RCX in Listing 2) are mapped to underlying physical registers enabling out-of-order execution. Trying to utilize the pipeline as much as possible, subsequent instructions (lines 5-7) are already decoded and allocated as μ OPs as well. The μ OPs are further sent to the reservation station holding the μ OPs while they wait to be executed by the corresponding execution unit. The execution of a μ OP can be delayed if execution units are already used to their corresponding capacity, or operand values have not been computed yet.

When the kernel address is loaded in line 4, it is likely that the CPU already issued the subsequent instructions as part of the out-of-order execution, and that their corresponding μ OPs wait in the reservation station for the content of the kernel address to arrive. As soon as the fetched data is observed on the common data bus, the μ OPs can begin their execution. Furthermore, processor interconnects [31, 3] and cache coherence protocols [59] guarantee that the most recent value of a memory address is read, regardless of the storage location in a multi-core or multi-CPU system.

When the μ OPs finish their execution, they retire in-order, and, thus, their results are committed to the architectural state. During the retirement, any interrupts and exceptions that occurred during the execution of the instruction are handled. Thus, if the MOV instruction that loads the kernel address is retired, the exception is registered, and the pipeline is flushed to eliminate all results of subsequent instructions which were executed out of

order. However, there is a race condition between raising this exception and our attack step 2 as described below.

As reported by Gruss et al. [21], prefetching kernel addresses sometimes succeeds. We found that prefetching the kernel address can slightly improve the performance of the attack on some systems.

Step 2: Transmitting the secret. The instruction sequence from step 1 which is executed out of order has to be chosen in a way that it becomes a transient instruction sequence. If this transient instruction sequence is executed before the MOV instruction is retired (*i.e.*, raises the exception), and the transient instruction sequence performed computations based on the secret, it can be utilized to transmit the secret to the attacker.

As already discussed, we utilize cache attacks that allow building fast and low-noise covert channels using the CPU’s cache. Thus, the transient instruction sequence has to encode the secret into the microarchitectural cache state, similar to the toy example in Section 3.

We allocate a probe array in memory and ensure that no part of this array is cached. To transmit the secret, the transient instruction sequence contains an indirect memory access to an address which is computed based on the secret (inaccessible) value. In line 5 of Listing 2, the secret value from step 1 is multiplied by the page size, *i.e.*, 4 KB. The multiplication of the secret ensures that accesses to the array have a large spatial distance to each other. This prevents the hardware prefetcher from loading adjacent memory locations into the cache as well. Here, we read a single byte at once. Hence, our probe array is 256×4096 bytes, assuming 4 KB pages.

Note that in the out-of-order execution we have a noise-bias towards register value ‘0’. We discuss the reasons for this in Section 5.2. However, for this reason, we introduce a retry-logic into the transient instruction sequence. In case we read a ‘0’, we try to reread the secret (step 1). In line 7, the multiplied secret is added to the base address of the probe array, forming the target address of the covert channel. This address is read to cache the corresponding cache line. The address will be loaded into the L1 data cache of the requesting core and, due to the inclusiveness, also the L3 cache where it can be read from other cores. Consequently, our transient instruction sequence affects the cache state based on the secret value that was read in step 1.

Since the transient instruction sequence in step 2 races against raising the exception, reducing the runtime of step 2 can significantly improve the performance of the attack. For instance, taking care that the address translation for the probe array is cached in the translation-lookaside buffer (TLB) increases the attack performance on some systems.

Step 3: Receiving the secret. In step 3, the attacker recovers the secret value (step 1) by leveraging a microarchitectural side-channel attack (*i.e.*, the receiving end of a microarchitectural covert channel) that transfers the cache state (step 2) back into an architectural state. As discussed in Section 4.2, our implementation of Meltdown relies on Flush+Reload for this purpose.

When the transient instruction sequence of step 2 is executed, exactly one cache line of the probe array is cached. The position of the cached cache line within the probe array depends only on the secret which is read in step 1. Thus, the attacker iterates over all 256 pages of the probe array and measures the access time for every first cache line (*i.e.*, offset) on the page. The number of the page containing the cached cache line corresponds directly to the secret value.

Dumping the entire physical memory. Repeating all 3 steps of Meltdown, an attacker can dump the entire memory by iterating over all addresses. However, as the memory access to the kernel address raises an exception that terminates the program, we use one of the methods from Section 4.1 to handle or suppress the exception.

As all major operating systems also typically map the entire physical memory into the kernel address space (cf. Section 2.2) in every user process, Meltdown can also read the entire physical memory of the target machine.

5.2 Optimizations and Limitations

Inherent bias towards 0. While CPUs generally stall if a value is not available during an out-of-order load operation [28], CPUs might continue with the out-of-order execution by assuming a value for the load [12]. We observed that the illegal memory load in our Meltdown implementation (line 4 in Listing 2) often returns a ‘0’, which can be clearly observed when implemented using an add instruction instead of the mov. The reason for this bias to ‘0’ may either be that the memory load is masked out by a failed permission check, or a speculated value because the data of the stalled load is not available yet.

This inherent bias results from the race condition in the out-of-order execution, which may be won (*i.e.*, reads the correct value), but is often lost (*i.e.*, reads a value of ‘0’). This bias varies between different machines as well as hardware and software configurations and the specific implementation of Meltdown. In an unoptimized version, the probability that a value of ‘0’ is erroneously returned is high. Consequently, our Meltdown implementation performs a certain number of retries when the code in Listing 2 results in reading a value of ‘0’ from the Flush+Reload attack. The maximum number of retries is an optimization parameter influencing the attack performance and the error rate. On the Intel Core i5-6200U

using exception handling, we read a '0' on average in 5.25 % ($\sigma = 4.15$) with our unoptimized version. With a simple retry loop, we reduced the probability to 0.67 % ($\sigma = 1.47$). On the Core i7-8700K, we read on average a '0' in 1.78 % ($\sigma = 3.07$). Using Intel TSX, the probability is further reduced to 0.008 %.

Optimizing the case of 0. Due to the inherent bias of Meltdown, a cache hit on cache line '0' in the Flush+Reload measurement, does not provide the attacker with any information. Hence, measuring cache line '0' can be omitted and in case there is no cache hit on any other cache line, the value can be assumed to be '0'. To minimize the number of cases where no cache hit on a non-zero line occurs, we retry reading the address in the transient instruction sequence until it encounters a value different from '0' (line 6). This loop is terminated either by reading a non-zero value or by the raised exception of the invalid memory access. In either case, the time until exception handling or exception suppression returns the control flow is independent of the loop after the invalid memory access, *i.e.*, the loop does not slow down the attack measurably. Hence, these optimizations may increase the attack performance.

Single-bit transmission. In the attack description in Section 5.1, the attacker transmitted 8 bits through the covert channel at once and performed $2^8 = 256$ Flush+Reload measurements to recover the secret. However, there is a trade-off between running more transient instruction sequences and performing more Flush+Reload measurements. The attacker could transmit an arbitrary number of bits in a single transmission through the covert channel, by reading more bits using a MOV instruction for a larger data value. Furthermore, the attacker could mask bits using additional instructions in the transient instruction sequence. We found the number of additional instructions in the transient instruction sequence to have a negligible influence on the performance of the attack.

The performance bottleneck in the generic attack described above is indeed, the time spent on Flush+Reload measurements. In fact, with this implementation, almost the entire time is spent on Flush+Reload measurements. By transmitting only a single bit, we can omit all but one Flush+Reload measurement, *i.e.*, the measurement on cache line 1. If the transmitted bit was a '1', then we observe a cache hit on cache line 1. Otherwise, we observe no cache hit on cache line 1.

Transmitting only a single bit at once also has drawbacks. As described above, our side channel has a bias towards a secret value of '0'. If we read and transmit multiple bits at once, the likelihood that all bits are '0' may be quite small for actual user data. The likelihood that a single bit is '0' is typically close to 50 %. Hence,

the number of bits read and transmitted at once is a trade-off between some implicit error-reduction and the overall transmission rate of the covert channel.

However, since the error rates are quite small in either case, our evaluation (cf. Section 6) is based on the single-bit transmission mechanics.

Exception Suppression using Intel TSX. In Section 4.1, we discussed the option to prevent that an exception is raised due an invalid memory access. Using Intel TSX, a hardware transactional memory implementation, we can completely suppress the exception [37].

With Intel TSX, multiple instructions can be grouped to a transaction, which appears to be an atomic operation, *i.e.*, either all or no instruction is executed. If one instruction within the transaction fails, already executed instructions are reverted, but no exception is raised.

If we wrap the code from Listing 2 with such a TSX instruction, any exception is suppressed. However, the microarchitectural effects are still visible, *i.e.*, the cache state is persistently manipulated from within the hardware transaction [19]. This results in higher channel capacity, as suppressing the exception is significantly faster than trapping into the kernel for handling the exception, and continuing afterward.

Dealing with KASLR. In 2013, kernel address space layout randomization (KASLR) was introduced to the Linux kernel (starting from version 3.14 [11]) allowing to randomize the location of kernel code at boot time. However, only as recently as May 2017, KASLR was enabled by default in version 4.12 [54]. With KASLR also the direct-physical map is randomized and not fixed at a certain address such that the attacker is required to obtain the randomized offset before mounting the Meltdown attack. However, the randomization is limited to 40 bit.

Thus, if we assume a setup of the target machine with 8 GB of RAM, it is sufficient to test the address space for addresses in 8 GB steps. This allows covering the search space of 40 bit with only 128 tests in the worst case. If the attacker can successfully obtain a value from a tested address, the attacker can proceed to dump the entire memory from that location. This allows mounting Meltdown on a system despite being protected by KASLR within seconds.

6 Evaluation

In this section, we evaluate Meltdown and the performance of our proof-of-concept implementation.¹¹ Section 6.1 discusses the information which Meltdown can

¹¹<https://github.com/IAIK/meltdown>

Table 1: Experimental setups.

Environment	CPU Model	Cores
Lab	Celeron G540	2
Lab	Core i5-3230M	2
Lab	Core i5-3320M	2
Lab	Core i7-4790	4
Lab	Core i5-6200U	2
Lab	Core i7-6600U	2
Lab	Core i7-6700K	4
Lab	Core i7-8700K	12
Lab	Xeon E5-1630 v3	8
Cloud	Xeon E5-2676 v3	12
Cloud	Xeon E5-2650 v4	12
Phone	Exynos 8890	8

leak, and Section 6.2 evaluates the performance of Meltdown, including countermeasures. Finally, we discuss limitations for AMD and ARM in Section 6.3.

Table 1 shows a list of configurations on which we successfully reproduced Meltdown. For the evaluation of Meltdown, we used both laptops as well as desktop PCs with Intel Core CPUs and an ARM-based mobile phone. For the cloud setup, we tested Meltdown in virtual machines running on Intel Xeon CPUs hosted in the Amazon Elastic Compute Cloud as well as on DigitalOcean. Note that for ethical reasons we did not use Meltdown on addresses referring to physical memory of other tenants.

6.1 Leakage and Environments

We evaluated Meltdown on both Linux (cf. Section 6.1.1), Windows 10 (cf. Section 6.1.3) and Android (cf. Section 6.1.4), without the patches introducing the KAISER mechanism. On these operating systems, Meltdown can successfully leak kernel memory. We also evaluated the effect of the KAISER patches on Meltdown on Linux, to show that KAISER prevents the leakage of kernel memory (cf. Section 6.1.2). Furthermore, we discuss the information leakage when running inside containers such as Docker (cf. Section 6.1.5). Finally, we evaluate Meltdown on uncached and uncacheable memory (cf. Section 6.1.6).

6.1.1 Linux

We successfully evaluated Meltdown on multiple versions of the Linux kernel, from 2.6.32 to 4.13.0, without the patches introducing the KAISER mechanism. On all these versions of the Linux kernel, the kernel address space is also mapped into the user address space. Thus, all kernel addresses are also mapped into the address space of user space applications, but any access is prevented due to the permission settings for these addresses.

As Meltdown bypasses these permission settings, an attacker can leak the complete kernel memory if the virtual address of the kernel base is known. Since all major operating systems also map the entire physical memory into the kernel address space (cf. Section 2.2), all physical memory can also be read.

Before kernel 4.12, kernel address space layout randomization (KASLR) was not active by default [57]. If KASLR is active, Meltdown can still be used to find the kernel by searching through the address space (cf. Section 5.2). An attacker can also simply de-randomize the direct-physical map by iterating through the virtual address space. Without KASLR, the direct-physical map starts at address `0xffff 8800 0000 0000` and linearly maps the entire physical memory. On such systems, an attacker can use Meltdown to dump the entire physical memory, simply by reading from virtual addresses starting at `0xffff 8800 0000 0000`.

On newer systems, where KASLR is active by default, the randomization of the direct-physical map is limited to 40 bit. It is even further limited due to the linearity of the mapping. Assuming that the target system has at least 8 GB of physical memory, the attacker can test addresses in steps of 8 GB, resulting in a maximum of 128 memory locations to test. Starting from one discovered location, the attacker can again dump the entire physical memory.

Hence, for the evaluation, we can assume that the randomization is either disabled, or the offset was already retrieved in a pre-computation step.

6.1.2 Linux with KAISER Patch

The KAISER patch by Gruss et al. [20] implements a stronger isolation between kernel and user space. KAISER does not map any kernel memory in the user space, except for some parts required by the x86 architecture (e.g., interrupt handlers). Thus, there is no valid mapping to either kernel memory or physical memory (via the direct-physical map) in the user space, and such addresses can therefore not be resolved. Consequently, Meltdown cannot leak any kernel or physical memory except for the few memory locations which have to be mapped in user space.

We verified that KAISER indeed prevents Meltdown, and there is no leakage of any kernel or physical memory.

Furthermore, if KASLR is active, and the few remaining memory locations are randomized, finding these memory locations is not trivial due to their small size of several kilobytes. Section 7.2 discusses the security implications of these mapped memory locations.

6.1.3 Microsoft Windows

We successfully evaluated Meltdown on a recent Microsoft Windows 10 operating system, last updated just before patches against Meltdown were rolled out. In line with the results on Linux (cf. Section 6.1.1), Meltdown also can leak arbitrary kernel memory on Windows. This is not surprising, since Meltdown does not exploit any software issues, but is caused by a hardware issue.

In contrast to Linux, Windows does not have the concept of an identity mapping, which linearly maps the physical memory into the virtual address space. Instead, a large fraction of the physical memory is mapped in the paged pools, non-paged pools, and the system cache. Furthermore, Windows maps the kernel into the address space of every application too. Thus, Meltdown can read kernel memory which is mapped in the kernel address space, *i.e.*, any part of the kernel which is not swapped out, and any page mapped in the paged and non-paged pool, and the system cache.

Note that there are physical pages which are mapped in one process but not in the (kernel) address space of another process, *i.e.*, physical pages which cannot be attacked using Meltdown. However, most of the physical memory will still be accessible through Meltdown.

We were successfully able to read the binary of the Windows kernel using Meltdown. To verify that the leaked data is actual kernel memory, we first used the Windows kernel debugger to obtain kernel addresses containing actual data. After leaking the data, we again used the Windows kernel debugger to compare the leaked data with the actual memory content, confirming that Meltdown can successfully leak kernel memory.

6.1.4 Android

We successfully evaluated Meltdown on a Samsung Galaxy S7 mobile phone running LineageOS Android 14.1 with a Linux kernel 3.18.14. The device is equipped with a Samsung Exynos 8 Octa 8890 SoC consisting of a ARM Cortex-A53 CPU with 4 cores as well as an Exynos M1 "Mongoose" CPU with 4 cores [6]. While we were not able to mount the attack on the Cortex-A53 CPU, we successfully mounted Meltdown on Samsung's custom cores. Using *exception suppression* described in Section 4.1, we successfully leaked a predefined string using the direct-physical map located at the virtual address `0xffff fbf c000 0000`.

6.1.5 Containers

We evaluated Meltdown in containers sharing a kernel, including Docker, LXC, and OpenVZ and found that the attack can be mounted without any restrictions. Running Meltdown inside a container allows to leak information

not only from the underlying kernel but also from all other containers running on the same physical host.

The commonality of most container solutions is that every container uses the same kernel, *i.e.*, the kernel is shared among all containers. Thus, every container has a valid mapping of the entire physical memory through the direct-physical map of the shared kernel. Furthermore, Meltdown cannot be blocked in containers, as it uses only memory accesses. Especially with Intel TSX, only unprivileged instructions are executed without even trapping into the kernel.

Thus, the isolation of containers sharing a kernel can be entirely broken using Meltdown. This is especially critical for cheaper hosting providers where users are not separated through fully virtualized machines, but only through containers. We verified that our attack works in such a setup, by successfully leaking memory contents from a container of a different user under our control.

6.1.6 Uncached and Uncacheable Memory

In this section, we evaluate whether it is a requirement for data to be leaked by Meltdown to reside in the L1 data cache [33]. Therefore, we constructed a setup with two processes pinned to different physical cores. By flushing the value, using the `clflush` instruction, and only reloading it on the other core, we create a situation where the target data is not in the L1 data cache of the attacker core. As described in Section 6.2, we can still leak the data at a lower reading rate. This clearly shows that data presence in the attacker's L1 data cache is not a requirement for Meltdown. Furthermore, this observation has also been confirmed by other researchers [7, 35, 5].

The reason why Meltdown can leak uncached memory may be that Meltdown implicitly caches the data. We devise a second experiment, where we mark pages as *uncacheable* and try to leak data from them. This has the consequence that every read or write operation to one of those pages will directly go to the main memory, thus, bypassing the cache. In practice, only a negligible amount of system memory is marked uncacheable. We observed that if the attacker is able to trigger a legitimate load of the target address, *e.g.*, by issuing a system call (regular or in speculative execution [40]), on the same CPU core as the Meltdown attack, the attacker can leak the content of the uncacheable pages. We suspect that Meltdown reads the value from the line fill buffers. As the fill buffers are shared between threads running on the same core, the read to the same address within the Meltdown attack could be served from one of the fill buffers allowing the attack to succeed. However, we leave further investigations on this matter open for future work.

A similar observation on uncacheable memory was also made with Spectre attacks on the System Manage-

ment Mode [10]. While the attack works on memory set uncacheable over Memory-Type Range Registers, it does not work on memory-mapped I/O regions, which is the expected behavior as accesses to memory-mapped I/O can always have architectural effects.

6.2 Meltdown Performance

To evaluate the performance of Meltdown, we leaked known values from kernel memory. This allows us to not only determine how fast an attacker can leak memory, but also the error rate, *i.e.*, how many byte errors to expect. The race condition in Meltdown (cf. Section 5.2) has a significant influence on the performance of the attack, however, the race condition can always be won. If the targeted data resides close to the core, *e.g.*, in the L1 data cache, the race condition is won with a high probability. In this scenario, we achieved average reading rates of up to 582 KB/s ($\mu = 552.4, \sigma = 10.2$) with an error rate as low as 0.003 % ($\mu = 0.009, \sigma = 0.014$) using exception suppression on the Core i7-8700K over 10 runs over 10 seconds. With the Core i7-6700K we achieved 569 KB/s ($\mu = 515.5, \sigma = 5.99$) with an minimum error rate of 0.002 % ($\mu = 0.003, \sigma = 0.001$) and 491 KB/s ($\mu = 466.3, \sigma = 16.75$) with a minimum error rate of 10.7 % ($\mu = 11.59, \sigma = 0.62$) on the Xeon E5-1630. However, with a slower version with an average reading speed of 137 KB/s, we were able to reduce the error rate to 0. Furthermore, on the Intel Core i7-6700K if the data resides in the L3 data cache but not in L1, the race condition can still be won often, but the average reading rate decreases to 12.4 KB/s with an error rate as low as 0.02 % using exception suppression. However, if the data is uncached, winning the race condition is more difficult and, thus, we have observed reading rates of less than 10 B/s on most systems. Nevertheless, there are two optimizations to improve the reading rate: First, by simultaneously letting other threads prefetch the memory locations [21] of and around the target value and access the target memory location (with exception suppression or handling). This increases the probability that the spying thread sees the secret data value in the right moment during the data race. Second, by triggering the hardware prefetcher through speculative accesses to memory locations of and around the target value. With these two optimizations, we can improve the reading rate for uncached data to 3.2 KB/s.

For all tests, we used Flush+Reload as a covert channel to leak the memory as described in Section 5, and Intel TSX to suppress the exception. An extensive evaluation of exception suppression using conditional branches was done by Kocher et al. [40] and is thus omitted in this paper for the sake of brevity.

6.3 Limitations on ARM and AMD

We also tried to reproduce the Meltdown bug on several ARM and AMD CPUs. While we were able to successfully leak kernel memory with the attack described in Section 5 on different Intel CPUs and a Samsung Exynos M1 processor, we did not manage to mount Meltdown on other ARM cores nor on AMD. In the case of ARM, the only affected processor is the Cortex-A75 [17] which has not been available and, thus, was not among our devices under test. However, appropriate kernel patches have already been provided [2]. Furthermore, an altered attack of Meltdown targeting system registers instead of inaccessible memory locations is applicable on several ARM processors [17]. Meanwhile, AMD publicly stated that none of their CPUs are not affected by Meltdown due to architectural differences [1].

The major part of a microarchitecture is usually not publicly documented. Thus, it is virtually impossible to know the differences in the implementations that allow or prevent Meltdown without proprietary knowledge and, thus, the intellectual property of the individual CPU manufacturers. The key point is that on a microarchitectural level the load to the unprivileged address and the subsequent instructions are executed while the fault is only handled when the faulting instruction is retired. It can be assumed that the execution units for the load and the TLB are designed differently on ARM, AMD and Intel and, thus, the privileges for the load are checked differently and occurring faults are handled differently, *e.g.*, issuing a load only after the permission bit in the page table entry has been checked. However, from a performance perspective, issuing the load in parallel or only checking permissions while retiring an instruction is a reasonable decision. As trying to load kernel addresses from user space is not what programs usually do and by guaranteeing that the state does not become architecturally visible, not squashing the load is legitimate. However, as the state becomes visible on the microarchitectural level, such implementations are vulnerable.

However, for both ARM and AMD, the toy example as described in Section 3 works reliably, indicating that out-of-order execution generally occurs and instructions past illegal memory accesses are also performed.

7 Countermeasures

In this section, we discuss countermeasures against the Meltdown attack. At first, as the issue is rooted in the hardware itself, we discuss possible microcode updates and general changes in the hardware design. Second, we discuss the KAISER countermeasure that has been developed to mitigate side-channel attacks against KASLR which inadvertently also protects against Meltdown.

7.1 Hardware

Meltdown bypasses the hardware-enforced isolation of security domains. There is no software vulnerability involved in Meltdown. Any software patch (e.g., KAISER [20]) will leave small amounts of memory exposed (cf. Section 7.2). There is no documentation whether a fix requires the development of completely new hardware, or can be fixed using a microcode update.

As Meltdown exploits out-of-order execution, a trivial countermeasure is to disable out-of-order execution completely. However, performance impacts would be devastating, as the parallelism of modern CPUs could not be leveraged anymore. Thus, this is not a viable solution.

Meltdown is some form of race condition between the fetch of a memory address and the corresponding permission check for this address. Serializing the permission check and the register fetch can prevent Meltdown, as the memory address is never fetched if the permission check fails. However, this involves a significant overhead to every memory fetch, as the memory fetch has to stall until the permission check is completed.

A more realistic solution would be to introduce a hard split of user space and kernel space. This could be enabled optionally by modern kernels using a new hard-split bit in a CPU control register, e.g., CR4. If the hard-split bit is set, the kernel has to reside in the upper half of the address space, and the user space has to reside in the lower half of the address space. With this hard split, a memory fetch can immediately identify whether such a fetch of the destination would violate a security boundary, as the privilege level can be directly derived from the virtual address without any further lookups. We expect the performance impacts of such a solution to be minimal. Furthermore, the backwards compatibility is ensured, since the hard-split bit is not set by default and the kernel only sets it if it supports the hard-split feature.

Note that these countermeasures only prevent Meltdown, and not the class of Spectre attacks described by Kocher et al. [40]. Likewise, their presented countermeasures [40] do not affect Meltdown. We stress that it is important to deploy countermeasures against both attacks.

7.2 KAISER

As existing hardware is not as easy to patch, there is a need for software workarounds until new hardware can be deployed. Gruss et al. [20] proposed KAISER, a kernel modification to not have the kernel mapped in the user space. This modification was intended to prevent side-channel attacks breaking KASLR [29, 21, 37]. However, it also prevents Meltdown, as it ensures that there is no valid mapping to kernel space or physical memory available in user space. In concurrent work

to KAISER, Gens et al. [14] proposed LAZARUS as a modification to the Linux kernel to thwart side-channel attacks breaking KASLR by separating address spaces similar to KAISER. As the Linux kernel continued the development of the original KAISER patch and Windows [53] and macOS [34] based their implementation on the concept of KAISER to defeat Meltdown, we will discuss KAISER in more depth.

Although KAISER provides basic protection against Meltdown, it still has some limitations. Due to the design of the x86 architecture, several privileged memory locations are still required to be mapped in user space [20], leaving a residual attack surface for Meltdown, *i.e.*, these memory locations can still be read from user space. Even though these memory locations do not contain any secrets, e.g., credentials, they might still contain pointers. Leaking one pointer can suffice to break KASLR, as the randomization can be computed from the pointer value.

Still, KAISER is the best short-time solution currently available and should therefore be deployed on all systems immediately. Even with Meltdown, KAISER can avoid having any kernel pointers on memory locations that are mapped in the user space which would leak information about the randomized offsets. This would require trampoline locations for every kernel pointer, *i.e.*, the interrupt handler would not call into kernel code directly, but through a trampoline function. The trampoline function must only be mapped in the kernel. It must be randomized with a different offset than the remaining kernel. Consequently, an attacker can only leak pointers to the trampoline code, but not the randomized offsets of the remaining kernel. Such trampoline code is required for every kernel memory that still has to be mapped in user space and contains kernel addresses. This approach is a trade-off between performance and security which has to be assessed in future work.

The original KAISER patch [18] for the Linux kernel has been improved [24, 25, 26, 27] with various optimizations, e.g., support for PCIDs. Afterwards, before merging it into the mainline kernel, it has been renamed to kernel page-table isolation (KPTI) [49, 15]. KPTI is active in recent releases of the Linux kernel and has been backported to older versions as well [30, 43, 44, 42].

Microsoft implemented a similar patch inspired by KAISER [53] named KVA Shadow [39]. While KVA Shadow only maps a minimum of kernel transition code and data pages required to switch between address spaces, it does not protect against side-channel attacks against KASLR [39].

Apple released updates in iOS 11.2, macOS 10.13.2 and tvOS 11.2 to mitigate Meltdown. Similar to Linux and Windows, macOS shared the kernel and user address spaces in 64-bit mode unless the `-no-shared-cr3` boot option was set [46]. This option unmaps the user space

while running in kernel mode but does not unmap the kernel while running in user mode [51]. Hence, it has no effect on Meltdown. Consequently, Apple introduced *Double Map* [34] following the principles of KAISER to mitigate Meltdown.

8 Discussion

Meltdown fundamentally changes our perspective on the security of hardware optimizations that manipulate the state of microarchitectural elements. The fact that hardware optimizations can change the state of microarchitectural elements, and thereby imperil secure software implementations, is known since more than 20 years [41]. Both industry and the scientific community so far accepted this as a necessary evil for efficient computing. Today it is considered a bug when a cryptographic algorithm is not protected against the microarchitectural leakage introduced by the hardware optimizations. Meltdown changes the situation entirely. Meltdown shifts the granularity from a comparably low spatial and temporal granularity, e.g., 64-bytes every few hundred cycles for cache attacks, to an arbitrary granularity, allowing an attacker to read every single bit. This is nothing any (cryptographic) algorithm can protect itself against. KAISER is a short-term software fix, but the problem we have uncovered is much more significant.

We expect several more performance optimizations in modern CPUs which affect the microarchitectural state in some way, not even necessarily through the cache. Thus, hardware which is designed to provide certain security guarantees, e.g., CPUs running untrusted code, requires a redesign to avoid Meltdown- and Spectre-like attacks. Meltdown also shows that even error-free software, which is explicitly written to thwart side-channel attacks, is not secure if the design of the underlying hardware is not taken into account.

With the integration of KAISER into all major operating systems, an important step has already been done to prevent Meltdown. KAISER is a fundamental change in operating system design. Instead of always mapping everything into the address space, mapping only the minimally required memory locations appears to be a first step in reducing the attack surface. However, it might not be enough, and even stronger isolation may be required. In this case, we can trade flexibility for performance and security, by e.g., enforcing a certain virtual memory layout for every operating system. As most modern operating systems already use a similar memory layout, this might be a promising approach.

Meltdown also heavily affects cloud providers, especially if the guests are not fully virtualized. For performance reasons, many hosting or cloud providers do not have an abstraction layer for virtual memory. In

such environments, which typically use containers, such as Docker or OpenVZ, the kernel is shared among all guests. Thus, the isolation between guests can simply be circumvented with Meltdown, fully exposing the data of all other guests on the same host. For these providers, changing their infrastructure to full virtualization or using software workarounds such as KAISER would both increase the costs significantly.

Concurrent work has investigated the possibility to read kernel memory via out-of-order or speculative execution, but has not succeeded [13, 50]. We are the first to demonstrate that it is possible. Even if Meltdown is fixed, Spectre [40] will remain an issue, requiring different defenses. Mitigating only one of them will leave the security of the entire system at risk. Meltdown and Spectre open a new field of research to investigate to what extent performance optimizations change the microarchitectural state, how this state can be translated into an architectural state, and how such attacks can be prevented.

9 Conclusion

In this paper, we presented Meltdown, a novel software-based attack exploiting out-of-order execution and side channels on modern processors to read arbitrary kernel memory from an unprivileged user space program. Without requiring any software vulnerability and independent of the operating system, Meltdown enables an adversary to read sensitive data of other processes or virtual machines in the cloud with up to 503 KB/s, affecting millions of devices. We showed that the countermeasure KAISER, originally proposed to protect from side-channel attacks against KASLR, inadvertently impedes Meltdown as well. We stress that KAISER needs to be deployed on every operating system as a short-term workaround, until Meltdown is fixed in hardware, to prevent large-scale exploitation of Meltdown.

Acknowledgments

Several authors of this paper found Meltdown independently, ultimately leading to this collaboration. We want to thank everyone who helped us in making this collaboration possible, especially Intel who handled our responsible disclosure professionally, communicated a clear timeline and connected all involved researchers. We thank Mark Brand from Google Project Zero for contributing ideas and Peter Cordes and Henry Wong for valuable feedback. We would like to thank our anonymous reviewers for their valuable feedback. Furthermore, we would like to thank Intel, ARM, Qualcomm, and Microsoft for feedback on an early draft.

Daniel Gruss, Moritz Lipp, Stefan Mangard and Michael Schwarz were supported by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No 681402).

Daniel Genkin was supported by NSF awards #1514261 and #1652259, financial assistance award 70NANB15H328 from the U.S. Department of Commerce, National Institute of Standards and Technology, the 2017-2018 Rothschild Postdoctoral Fellowship, and the Defense Advanced Research Project Agency (DARPA) under Contract #FA8650-16-C-7622.

References

- [1] AMD. Software techniques for managing speculation on AMD processors, 2018.
- [2] ARM. AArch64 Linux kernel port (KPTI base), <https://git.kernel.org/pub/scm/linux/kernel/git/arm64/linux.git/log/?h=kpti> 2018.
- [3] ARM LIMITED. *ARM CoreLink CCI-400 Cache Coherent Interconnect Technical Reference Manual*, r1p5 ed. ARM Limited, 2015.
- [4] BENDER, N., VAN DE POL, J., SMART, N. P., AND YAROM, Y. "Ooh Aah... Just a Little Bit": A small amount of side channel can go a long way. In *CHES'14* (2014).
- [5] BOLDIN, P. Meltdown reading other process's memory, <https://www.youtube.com/watch?v=EMBGXswJC4s> Jan 2018.
- [6] BURGESS, B. Samsung Exynos M1 Processor. In *IEEE Hot Chips* (2016).
- [7] CARVALHO, R. Twitter: Meltdown with uncached memory, https://twitter.com/raphael_scarv/status/952078140028964864 Jan 2018.
- [8] CHENG, C.-C. The schemes and performances of dynamic branch predictors. *Berkeley Wireless Research Center, Tech. Rep* (2000).
- [9] DEVIES, A. M. AMD Takes Computing to a New Horizon with Ryzen™ Processors, <https://www.amd.com/en-us/press-releases/Pages/amd-takes-computing-2016dec13.aspx> 2016.
- [10] ECLYPSIUM. System Management Mode Speculative Execution Attacks, <https://blog.eclipsium.com/2018/05/17/system-management-mode-speculative-execution-attacks/> May 2018.
- [11] EDGE, J. Kernel address space layout randomization, <https://lwn.net/Articles/569635/> 2013.
- [12] EICKEMEYER, R., LE, H., NGUYEN, D., STOLT, B., AND THOMPSON, B. Load lookahead prefetch for microprocessors, 2006. <https://encrypted.google.com/patents/US20060149935> US Patent App. 11/016,236.
- [13] FOGH, A. Negative Result: Reading Kernel Memory From User Mode, <https://cyber.wtf/2017/07/28/negative-result-reading-kernel-memory-from-user-mode/> 2017.
- [14] GENS, D., ARIAS, O., SULLIVAN, D., LIEBCHEN, C., JIN, Y., AND SADEGHI, A.-R. Lazarus: Practical side-channel resilient kernel-space randomization. In *International Symposium on Research in Attacks, Intrusions, and Defenses* (2017).
- [15] GLEIXNER, T. x86/kpti: Kernel Page Table Isolation (was KAISER), <https://lkml.org/lkml/2017/12/4/709> Dec 2017.
- [16] GRAS, B., RAZAVI, K., BOSMAN, E., BOS, H., AND GIUFFRIDA, C. ASLR on the Line: Practical Cache Attacks on the MMU. In *NDSS* (2017).
- [17] GRIENTHWAITE, R. Cache Speculation Side-channels, 2018.
- [18] GRUSS, D. [RFC, PATCH] x86_64: KAISER - do not map kernel in user mode, <https://lkml.org/lkml/2017/5/4/220> May 2017.
- [19] GRUSS, D., LETTNER, J., SCHUSTER, F., OHRIMENKO, O., HALLER, I., AND COSTA, M. Strong and Efficient Cache Side-Channel Protection using Hardware Transactional Memory. In *USENIX Security Symposium* (2017).
- [20] GRUSS, D., LIPP, M., SCHWARZ, M., FELLNER, R., MAURICE, C., AND MANGARD, S. KASLR is Dead: Long Live KASLR. In *International Symposium on Engineering Secure Software and Systems* (2017), Springer, pp. 161–176.
- [21] GRUSS, D., MAURICE, C., FOGH, A., LIPP, M., AND MANGARD, S. Prefetch Side-Channel Attacks: Bypassing SMAP and Kernel ASLR. In *CCS* (2016).
- [22] GRUSS, D., MAURICE, C., WAGNER, K., AND MANGARD, S. Flush+Flush: A Fast and Stealthy Cache Attack. In *DIMVA* (2016).
- [23] GRUSS, D., SPREITZER, R., AND MANGARD, S. Cache Template Attacks: Automating Attacks on Inclusive Last-Level Caches. In *USENIX Security Symposium* (2015).
- [24] HANSEN, D. [PATCH 00/23] KAISER: unmap most of the kernel from userspace page tables, <https://lkml.org/lkml/2017/10/31/884> Oct 2017.
- [25] HANSEN, D. [v2] KAISER: unmap most of the kernel from userspace page tables, <https://lkml.org/lkml/2017/11/8/752> Nov 2017.
- [26] HANSEN, D. [v3] KAISER: unmap most of the kernel from userspace page tables, <https://lkml.org/lkml/2017/11/10/433> Nov 2017.
- [27] HANSEN, D. [v4] KAISER: unmap most of the kernel from userspace page tables, <https://lkml.org/lkml/2017/11/22/956> Nov 2017.
- [28] HENNESSY, J. L., AND PATTERSON, D. A. *Computer Architecture: A Quantitative Approach*, 6 ed. Morgan Kaufmann, 2017.
- [29] HUND, R., WILLEMS, C., AND HOLZ, T. Practical Timing Side Channel Attacks against Kernel Space ASLR. In *S&P* (2013).
- [30] HUTCHINGS, B. Linux 3.16.53, <https://cdn.kernel.org/pub/linux/kernel/v3.x/ChangeLog-3.16.53> 2018.
- [31] INTEL. An introduction to the intel quickpath interconnect, Jan 2009.
- [32] INTEL. Intel® 64 and IA-32 Architectures Optimization Reference Manual, 2017.
- [33] INTEL. Intel analysis of speculative execution side channels, <https://newsroom.intel.com/wp-content/uploads/sites/11/2018/01/Intel-Analysis-of-Speculative-Execution-Side-Channels.pdf> Jan 2018.
- [34] IONESCU, A. Twitter: Apple Double Map, <https://twitter.com/aionescu/status/948609809540046849> 2017.
- [35] IONESCU, A. Twitter: Meltdown with uncached memory, <https://twitter.com/aionescu/status/950994906759143425> Jan 2018.
- [36] IRAZOQUI, G., INCI, M. S., EISENBARTH, T., AND SUNAR, B. Wait a minute! A fast, Cross-VM attack on AES. In *RAID'14* (2014).

- [37] JANG, Y., LEE, S., AND KIM, T. Breaking Kernel Address Space Layout Randomization with Intel TSX. In *CCS* (2016).
- [38] JIMÉNEZ, D. A., AND LIN, C. Dynamic branch prediction with perceptrons. In *High-Performance Computer Architecture, 2001. HPCA. The Seventh International Symposium on* (2001), IEEE, pp. 197–206.
- [39] JOHNSON, K. KVA Shadow: Mitigating Meltdown on Windows, <https://blogs.technet.microsoft.com/srd/2018/03/23/kva-shadow-mitigating-meltdown-on-windows/> Mar 2018.
- [40] KOCHER, P., HORN, J., FOGH, A., GENKIN, D., GRUSS, G., HAAS, W., HAMBURG, M., LIPP, M., MANGARD, S., PRESCHER, T., SCHWARZ, M., AND YAROM, Y. Spectre attacks: Exploiting speculative execution. In *S&P* (2019). A preprint was published in 2018 as arXiv:1801.01203.
- [41] KOCHER, P. C. Timing Attacks on Implementations of Diffie-Hellman, RSA, DSS, and Other Systems. In *CRYPTO* (1996).
- [42] KROAH-HARTMAN, G. Linux 4.14.11, <https://cdn.kernel.org/pub/linux/kernel/v4.x/ChangeLog-4.14.11> 2018.
- [43] KROAH-HARTMAN, G. Linux 4.4.110, <https://cdn.kernel.org/pub/linux/kernel/v4.x/ChangeLog-4.4.110> 2018.
- [44] KROAH-HARTMAN, G. Linux 4.9.75, <https://cdn.kernel.org/pub/linux/kernel/v4.x/ChangeLog-4.9.75> 2018.
- [45] LEE, B., MALISHEVSKY, A., BECK, D., SCHMID, A., AND LANDRY, E. Dynamic branch prediction. *Oregon State University*.
- [46] LEVIN, J. *Mac OS X and IOS Internals: To the Apple's Core*. John Wiley & Sons, 2012.
- [47] LIPP, M., GRUSS, D., SPREITZER, R., MAURICE, C., AND MANGARD, S. ARMageddon: Cache Attacks on Mobile Devices. In *USENIX Security Symposium* (2016).
- [48] LIU, F., YAROM, Y., GE, Q., HEISER, G., AND LEE, R. B. Last-Level Cache Side-Channel Attacks are Practical. In *IEEE Symposium on Security and Privacy – SP* (2015), IEEE Computer Society, pp. 605–622.
- [49] LWN. The current state of kernel page-table isolation, <https://lwn.net/SubscriberLink/741878/eb6c9d3913d7cb2b/> Dec. 2017.
- [50] MAISURADZE, G., AND ROSSOW, C. Speculose: Analyzing the Security Implications of Speculative Execution in CPUs. *arXiv:1801.04084* (2018).
- [51] MANDT, T. Attacking the iOS Kernel: A Look at 'evasi0n', www.nislabs.no/content/download/38610/481190/file/NISlecture201303.pdf 2013.
- [52] MAURICE, C., WEBER, M., SCHWARZ, M., GINER, L., GRUSS, D., ALBERTO BOANO, C., MANGARD, S., AND RÖMER, K. Hello from the Other Side: SSH over Robust Cache Covert Channels in the Cloud. In *NDSS* (2017).
- [53] MILLER, M. Mitigating speculative execution side channel hardware vulnerabilities, <https://blogs.technet.microsoft.com/srd/2018/03/15/mitigating-speculative-execution-side-channel-hardware-vulnerabilities/> Mar 2018.
- [54] MOLNAR, I. x86: Enable KASLR by default, <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=6807c84652b0b7e2e198e50a9ad47ef41b236e59> 2017.
- [55] OSVIK, D. A., SHAMIR, A., AND TROMER, E. Cache Attacks and Countermeasures: the Case of AES. In *CT-RSA* (2006).
- [56] PERCIVAL, C. Cache missing for fun and profit. In *Proceedings of BSDCan* (2005).
- [57] PHORONIX. Linux 4.12 To Enable KASLR By Default, https://www.phoronix.com/scan.php?page=news_item&px=KASLR-Default-Linux-4.12 2017.
- [58] SCHWARZ, M., LIPP, M., GRUSS, D., WEISER, S., MAURICE, C., SPREITZER, R., AND MANGARD, S. KeyDrown: Eliminating Software-Based Keystroke Timing Side-Channel Attacks. In *NDSS'18* (2018).
- [59] SORIN, D. J., HILL, M. D., AND WOOD, D. A. *A Primer on Memory Consistency and Cache Coherence*. 2011.
- [60] TERAN, E., WANG, Z., AND JIMÉNEZ, D. A. Perceptron learning for reuse prediction. In *Microarchitecture (MICRO), 2016 49th Annual IEEE/ACM International Symposium on* (2016), IEEE, pp. 1–12.
- [61] TOMASULO, R. M. An efficient algorithm for exploiting multiple arithmetic units. *IBM Journal of research and Development* 11, 1 (1967), 25–33.
- [62] VINTAN, L. N., AND IRIDON, M. Towards a high performance neural branch predictor. In *Neural Networks, 1999. IJCNN'99. International Joint Conference on* (1999), vol. 2, IEEE, pp. 868–873.
- [63] YAROM, Y., AND FALKNER, K. Flush+Reload: a High Resolution, Low Noise, L3 Cache Side-Channel Attack. In *USENIX Security Symposium* (2014).
- [64] YEH, T.-Y., AND PATT, Y. N. Two-level adaptive training branch prediction. In *Proceedings of the 24th annual international symposium on Microarchitecture* (1991), ACM, pp. 51–61.
- [65] ZHANG, Y., JUELS, A., REITER, M. K., AND RISTENPART, T. Cross-Tenant Side-Channel Attacks in PaaS Clouds. In *CCS* (2014).

A Meltdown in Practice

In this section, we show how Meltdown can be used in practice. In Appendix A.1, we show physical memory dumps obtained via Meltdown, including passwords of the Firefox password manager. In Appendix A.2, we demonstrate a real-world exploit.

A.1 Physical-memory Dump using Meltdown

Listing 3 shows a memory dump using Meltdown on an Intel Core i7-6700K running Ubuntu 16.10 with the Linux kernel 4.8.0. In this example, we can identify HTTP headers of a request to a web server running on the machine. The XX cases represent bytes where the side channel did not yield any results, *i.e.*, no Flush+Reload hit. Additional repetitions of the attack may still be able to read these bytes.

Listing 4 shows a memory dump of Firefox 56 using Meltdown on the same machine. We can clearly identify some of the passwords that are stored in the internal password manager, *i.e.*, Dolphin18, insta_0203, and secretpwd0. The attack also recovered a URL which appears to be related to a Firefox add-on.

```

79cbb80: 6c4c 48 32 5a 78 66 56 44 73 4b 57 39 34 68 6d |1LH2ZxfVDeKW94hm|
79cbb90: 3364 2f 41 4d 41 45 44 41 41 41 41 51 45 42 |3d/AMAEAAAAAAAAQEB|
79cbb9a: 4141 41 41 41 41 3d 3d XX XX XX XX XX XX XX |AAAAAA==.....|
79cbb9b: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb9c: XXXX XX 65 2d 68 65 61 64 XX XX XX XX XX XX |...e-head.....|
79cbb9d: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb9e: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb9f: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb00: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb01: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb02: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb03: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb04: XXXX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
79cbb05: XXXX XX 0d 0a XX 6f 72 69 67 69 6e 61 6c 2d |.....original-|
79cbb06: 7265 73 70 6f 6e 73 65 2d 68 65 61 64 65 72 73 |response-headers|
79cbb07: XX44 61 74 65 3a 20 53 61 74 2c 20 30 39 20 44 |.Date: Sat, 09 D|
79cbb08: 6563 20 32 30 31 37 20 32 32 3a 32 39 3a 32 35 |ec 2017 22:29:25|
79cbb09: 2047 4d 54 0d 0a 43 6f 6e 74 65 6e 74 2d 4c 65 | GMT..Content-Le|
79cbb0a: 6e67 74 68 3a 20 31 0d 0a 43 6f 6e 74 65 6e 74 |ngth: 1..Content|
79cbb0b: 2d54 79 70 65 3a 20 74 65 78 74 2f 68 74 6d 6c |-Type: text/html|
79cbb0c: 3b20 63 68 61 72 73 65 74 3d 75 74 66 2d 38 0d |; charset=utf-8.|

```

Listing (3) Memory dump showing HTTP Headers on Ubuntu 16.10 on a Intel Core i7-6700K

A.2 Real-world Meltdown Exploit

In this section, we present a real-world exploit showing the applicability of Meltdown in practice, implemented by Pavel Boldin in collaboration with Raphael Carvalho. The exploit dumps the memory of a specific process, provided either the process id (PID) or the process name.

First, the exploit de-randomizes the kernel address space layout to be able to access internal kernel structures. Second, the kernel's task list is traversed until the victim process is found. Finally, the root of the victim's multilevel page table is extracted from the task structure and traversed to dump any of the victim's pages.

The three steps of the exploit are combined to an end-to-end exploit which targets a specific kernel build and a specific victim. The exploit can easily be adapted to work on any kernel build. The only requirement is access to either the binary or the symbol table of the kernel, which is true for all public kernels which are distributed as packages, *i.e.*, not self-compiled. In the remainder of this section, we provide a detailed explanation of the three steps.

A.2.1 Breaking KASLR

The first step is to de-randomize KASLR to access internal kernel structures. The exploit locates a known value inside the kernel, specifically the Linux banner string, as the content is known and it is large enough to rule out false positives. It starts looking for the banner string at the (non-randomized) default address according to the symbol table of the running kernel. If the string is not found, the next attempt is made at the next possible randomized address until the target is found. As the Linux KASLR implementation only has an entropy of 6 bits [37], there are only 64 possible randomization offsets, making this approach practical.

The difference between the found address and the non-randomized base address is then the randomization offset

```

f94b76f0: 12 XX e0 81 19 XX e0 81 44 6f 6c 70 68 69 6e 31 |.....Dolphin|
f94b7700: 38 e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 |8.....|
f94b7710: 70 52 b8 6b 96 7f XX XX XX XX XX XX XX XX |pR.k.....|
f94b7720: XX XX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
f94b7730: XX XX XX XX 4a XX XX XX XX XX XX XX XX XX |.....J.....|
f94b7740: XX XX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
f94b7750: XX XX XX XX XX XX XX XX XX XX e0 81 69 6e 73 74 |.....inst|
f94b7760: 61 5f 30 32 30 33 e5 e5 e5 e5 e5 e5 e5 e5 |a_0203.....|
f94b7770: 70 52 18 7d 28 7f XX XX XX XX XX XX XX XX |pR.}.....|
f94b7780: XX XX XX XX XX XX XX XX XX XX XX XX XX XX |.....T.....|
f94b7790: XX XX XX XX 54 XX XX XX XX XX XX XX XX XX |.....|
f94b77a0: XX XX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
f94b77b0: XX XX XX XX XX XX XX XX XX XX XX XX 73 65 63 72 |.....secl|
f94b77c0: 65 74 70 77 64 30 e5 e5 e5 e5 e5 e5 e5 e5 |etpud0.....|
f94b77d0: 30 b4 18 7d 28 7f XX XX XX XX XX XX XX XX |0.}.....|
f94b77e0: XX XX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
f94b77f0: XX XX XX XX XX XX XX XX XX XX XX XX XX XX |.....|
f94b7800: e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 e5 |.....|
f94b7810: 68 74 74 70 73 3a 2f 2f 61 64 64 6f 6e 73 2e 63 |https://addons.c|
f94b7820: 64 6e 2e 6d 6f 7a 69 6c 6c 61 2e 6e 65 74 2f 75 |dn.mozilla.net/u|
f94b7830: 73 65 72 2d 6d 65 64 69 61 2f 61 64 64 6f 6e 5f |ser-media/addon_|

```

Listing (4) Memory dump of Firefox 56 on Ubuntu 16.10 on a Intel Core i7-6700K disclosing saved passwords.

of the kernel address space. The remainder of this section assumes that addresses are already de-randomized using the detected offset.

A.2.2 Locating the Victim 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, which is at a fixed offset that only varies among different kernel builds. Thus, knowledge of the kernel build is sufficient to locate the task list.

Among other members, each task list structure contains a pointer to the next element in the task list as well as a task's PID, name, and the root of the multilevel page table. Thus, the exploit traverses the task list until the victim process is found.

A.2.3 Dumping the Victim Process

The root of the multilevel page table is extracted from the victim's task list entry. The page table entries on all levels are physical page addresses. Meltdown can read these addresses via the direct-physical map, *i.e.*, by adding the base address of the direct-physical map to the physical addresses. This base address is `0xffff 8800 0000 0000` if the direct-physical map is not randomized. If the direct-physical map is randomized, it can be extracted from the kernel's `page_offset_base` variable.

Starting at the root of the victim's multilevel page table, the exploit can simply traverse the levels down to the lowest level. For a specific address of the victim, the exploit uses the paging structures to resolve the respective physical address and read the content of this physical address via the direct-physical map. The exploit can also be easily extended to enumerate all pages belonging to the victim process, and then dump any (or all) of these pages.