

Microarchitectural Security of AWS Firecracker VMM for Serverless Cloud Platforms

Zane Weissman
zweissman@wpi.edu
Worcester Polytechnic Institute
Worcester, MA, USA

Thomas Eisenbarth
thomas.eisenbarth@uni-luebeck.de
University of Lübeck
Lübeck, S-H, Germany

Thore Tiemann
t.tiemann@uni-luebeck.de
University of Lübeck
Lübeck, S-H, Germany

Berk Sunar
sunar@wpi.edu
Worcester Polytechnic Institute
Worcester, MA, USA

ABSTRACT

Firecracker is a virtual machine manager (VMM) purpose-built by Amazon Web Services (AWS) for serverless cloud platforms—services that run code for end users on a per-task basis, automatically managing server infrastructure. Firecracker provides fast and lightweight VMs and promises a combination of the speed of containers, typically used to isolate small tasks, and the security of VMs, which tend to provide greater isolation at the cost of performance. This combination of security and efficiency, AWS claims, makes it not only possible but safe to run thousands of user tasks from different users on the same hardware, with the host system rapidly and frequently switching between active tasks. Though AWS states that microarchitectural attacks are included in their threat model, this class of attacks directly relies on shared hardware, just as the scalability of serverless computing relies on sharing hardware between unprecedented numbers of users.

In this work, we investigate just how secure Firecracker is against microarchitectural attacks. First, we review Firecracker’s stated isolation model and recommended best practices for deployment, identify potential threat models for serverless platforms, and analyze potential weak points. Then, we use microarchitectural attack proof-of-concepts to test the isolation provided by Firecracker and find that it offers little protection against Spectre or MDS attacks. We discover two particularly concerning cases: 1) a Medusa variant that threatens Firecracker VMs but *not* processes running outside them, and is not mitigated by defenses recommended by AWS, and 2) a Spectre-PHT variant that remains exploitable even if recommended countermeasures are in place and SMT is disabled in the system. In summary, we show that AWS overstates the security inherent to the Firecracker VMM and provides incomplete guidance for properly securing cloud systems that use Firecracker.

CCS CONCEPTS

• Security and privacy → Virtualization and security; Side-channel analysis and countermeasures.

KEYWORDS

system security, microarchitectural security, virtual machines, hypervisor, serverless, cloud systems

1 INTRODUCTION

Serverless computing is an emerging trend in cloud computing where cloud service providers (CSPs) serve runtime environments to their customers. This way, customers can focus on maintaining their function code while leaving the administrative work related to hardware, operating system (OS), and sometimes runtime to the CSPs. Common serverless platform models include function-as-a-service (FaaS) and container-as-a-service (CaaS). Since individual functions are typically small, but customers’ applications can each be running anywhere from one to thousands of functions, CSPs aim for fitting as many functions on a single server as possible to minimize idle times and, in turn, maximize profit. A rather light-weight approach to serving runtime environments is to run containers, which encapsulate a process with its dependencies so that only the necessary files for each process are loaded in virtual filesystems top of a shared kernel. This reduces a switch between containers to little more than a context switch between processes. On the other hand, full virtualization provides good isolation between virtual machines (VMs) and therefore security between tenants, while being rather heavy-weight as each VM comes with its own kernel.

Neither of these approaches, container or VM, is ideal for use in serverless environments, where ideally many short-lived functions owned by many users will run simultaneously and switch often, so new mechanisms of isolation have been developed for this use case. For example, mechanisms for in-process isolation [38, 45, 49] set out to improve the security of containers by reducing the attack surface of the runtime and underlying kernel. Protecting the kernel is important, as a compromised kernel directly leads to a fully compromised system in the container case. However, certain powerful protections, like limiting syscalls, also limit the functionality that is available to the container and even break compatibility with some applications. In VM research, developers created ever smaller and more efficient VMs, eventually leading to so-called microVMs. MicroVMs provide the same isolation guarantees as usual virtual machines, but are very limited in their capabilities when it comes to device or OS support, which makes them more light-weight compared to usual VMs and therefore better suited for serverless computing.

Firecracker [1] is a virtual machine manager (VMM) designed to run microVMs while providing memory overhead and start times comparable to those of common container systems. Firecracker is actively developed by Amazon Web Services (AWS) and has

been used in production for AWS Lambda [5] and AWS Fargate [4] serverless compute services since 2018 [1]. AWS’s design paper [1] describes the features of Firecracker, how it diverges from more traditional virtual machines, and the intended isolation model that it provides: safety for “multiple functions run[ning] on the same hardware, protected against privilege escalation, information disclosure, covert channels, and other risks” [1]. Furthermore, AWS provides production host setup recommendations [8] for securing parts of the CPU and kernel that a Firecracker VM interacts with. *In this paper, we challenge the claim that Firecracker protects functions from covert and side-channels across microVMs. We show that Firecracker itself does not add to the microarchitectural attack countermeasures but fully relies on the host and guest Linux kernels and CPU firmware/microcode updates.*

Microarchitectural attacks like the various Spectre [10, 13, 22, 30, 31, 33, 52] and microarchitectural data sampling (MDS) [14, 37, 46, 50] variants pose a threat to multi-tenant systems as they are often able to bypass both software and architectural isolation boundaries, including those of VMs. Spectre and MDS threaten tenants that share CPU core resources like the branch prediction unit (BPU) or the line-fill buffer (LFB). CSPs providing more traditional services can mitigate the problem of shared hardware resources by pinning the long-lived VMs tenants to separate CPU cores, which effectively partitions the resources between the tenants and ensures that the microarchitectural state is only effected by a single tenant at a time.

In serverless environments, however, the threat of microarchitectural attacks is greater. The reason for this is the short-livedness of the functions that are run by the different tenants. Server resources in serverless environments are expected to be over-committed, which leads to tenant functions competing for compute resources on the same hardware. Disabling simultaneous multi-threading (SMT), which would disable the concurrent use of CPU resources by two sibling threads, reduces the compute power of a CPU by up to 30% [34]. If customers rent specific CPU cores, this performance penalty may be acceptable, or both threads on a CPU core might be rented together. But for serverless services, the performance penalty directly translates to 30% fewer customers that can be served in a given amount of time. This is why it has to be assumed that most serverless CSPs keep SMT enabled in their systems unless they state otherwise. The microarchitectural attack surface is largest if SMT is enabled and the malicious thread has concurrent access to a shared core. But there are also attack variants that perform just as well if the attacker thread prepares the microarchitecture before it yields the CPU core to the victim thread or executes right after the victim thread has paused execution. And even if SMT is disabled by the CSP (as is the case for AWS Lambda), tenants still share CPUs with multiple others in this time-sliced fashion.

AWS claims that Firecracker running on a system with up-to-date microarchitectural defenses will provide sufficient hardening against microarchitectural attacks [1]. The Firecracker documentation also contains specific recommendations for microarchitectural security measures that should be enabled. *In this work, we examine Firecracker’s security claims and recommendations and reveal oversights in its guidance as well as wholly unmitigated threats.*

In summary, our main contributions are:

- We provide a comprehensive security analysis of the cross-tenant and tenant-hypervisor isolation of serverless compute when based on Firecracker VM.
- We test Firecracker’s defense capabilities against microarchitectural attack proof-of-concepts (PoCs), employing protections available through microcode updates and the Linux kernel. We show that the virtual machine itself provides *negligible protection* against major classes of microarchitectural attacks.
- We identify a variant of the Medusa MDS attack that becomes exploitable from within Firecracker VMs *even though it is not present on the host*. The kernel mitigation that protects against this exploit, and most known Medusa variants, is not mentioned by AWS’s Firecracker host setup recommendations. Additionally, we show that disabling SMT provides insufficient protection against the identified Medusa variant which urges the need of the kernel mitigation.
- We identify Spectre-PHT and Spectre-BTB variants which leak data even with recommended countermeasures in place. The Spectre-PHT variants even remain a problem when SMT is disabled if the attacker and victim share a CPU core in a time-sliced fashion.

1.1 Responsible Disclosure

We informed the AWS security team about our findings and discussed technical details. The AWS security team claims that the AWS services are not affected by our findings due to additional security measurements. AWS agreed that Firecracker does not provide micro-architectural security on its own but only in combination with microcode updates and secure host and guest operating systems and plans to update its host setup recommendations for Firecracker installations.

2 BACKGROUND

2.1 KVM

The Linux kernel-based virtual machine (KVM) [29] provides an abstraction of the hardware-assisted virtualization features like Intel VT-x or AMD-V that are available in modern CPUs. To support near-native execution, a *guest mode* is added to the Linux kernel in addition to the existing *kernel mode* and *user mode*. If in Linux guest mode, KVM causes the hardware to enter the hardware virtualization mode which replicates ring 0 and ring 3 privileges.¹

With KVM, I/O virtualization is performed mostly in user space by the process that created the VM, referred to as the VMM or hypervisor, in contrast to earlier hypervisors which typically required a separate hypervisor process [41]. A KVM hypervisor provides each VM guest with its own memory region that is separate from the memory region of the process that created the guest. This is true for guests created from kernel space as well as from user space. Each VM maps to a process on the Linux host and each virtual CPU assigned to the guest is a thread in that host process. The VM’s userspace hypervisor process makes system calls to KVM only when privileged execution is required, minimizing context switching and reducing the VM to kernel attack surface. Besides

¹The virtualized ring 0 and ring 3 are one of the core reasons why near-native code execution is achieved.

driving performance improvements across all sorts of applications, this design has allowed for the development of lightweight hypervisors that are especially useful for sandboxing individual programs and supporting cloud environments where many VMs are running at the same time.

2.2 Serverless Cloud Computing

An increasingly popular model for cloud computing is serverless computing, in which the CSP manages scalability and availability of the servers that run the user’s code. One implementation of serverless computing is called function-as-a-service (FaaS). In this model, a cloud user defines functions that are called as necessary through the service provider’s application programming interface (API) (hence the name “function-as-a-service”) and the CSP manages resource allocation on the server that executes the user’s function (hence the name “serverless computing”—the user does no server management). Similarly, container-as-a-service (CaaS) computing runs containers, portable runtime packages, on demand. The centralized server management of FaaS and CaaS is economically attractive to both CSPs and users. The CSP can manage its users’ workloads however it pleases, optimize for minimal operating cost, and implement flexible pricing where users pay for the execution time that they use. The user does not need to worry about server infrastructure design or management, and so reduces development costs and outsources maintenance cost to the CSP at a relatively small and predictable rate.

2.3 MicroVMs and AWS Firecracker

FaaS and CaaS providers use a variety of systems to manage running functions and containers. Container systems like Docker, Podman, and LXDE provide a convenient and lightweight way to package and run sandboxed applications in any environment. However, compared to the virtual machines used for many more traditional forms of cloud computing, containers offer less isolation and therefore less security. In recent years, major CSPs have introduced microVMs that back traditional containers with lightweight virtualization for extra security. [1, 55] The efficiency of hardware virtualization with KVM and lightweight design of microVMs means that code in virtualized, containerized or container-like systems can run nearly as fast as unvirtualized code and with comparable overhead to a traditional container.

Firecracker [1] is a microVM developed by AWS to isolate each of the AWS Lambda FaaS and AWS Fargate CaaS workloads in a separate VM. It only supports Linux guests on x86 or ARM Linux-KVM hosts and provides a limited number of devices that are available to guest systems. These limitations allow Firecracker to be very light-weight in the size of its code base and in memory overhead for a running VM, as well as very quick to boot or shut down. Additionally, the use of KVM lightens the requirements of Firecracker, since some virtualization functions are handled by kernel system calls and the host OS manages VMs as standard processes. Because of its small code base written in Rust, Firecracker is assumed to be very secure, even though security flaws have been identified in earlier versions (see CVE-2019-18960). Interestingly, the Firecracker white paper declares microarchitectural attacks to be in-scope of

its attacker model [1] but lacks a detailed security analysis or special countermeasures against microarchitectural attacks beyond common secure system configuration recommendations for the guest and host kernel. The Firecracker documentation does provide system security recommendations [8] that include a specific list of countermeasures, which we cover in section 2.6.1.

2.4 Meltdown and MDS

In 2018, the Meltdown [32] attack showed that speculatively accessed data could be exfiltrated across security boundaries by encoding it into a cache side-channel. This soon led to a whole class of similar attacks, known as microarchitectural data sampling (MDS), including Fallout [14], Rogue In-flight Data Load (RIDL) [50], TSX Asynchronous Abort (TAA) [50], and Zombieload [46]. These attacks all follow the same general pattern to exploit speculative execution:

- (1) The victim runs a program that handles secret data, and the secret data passes through a cache or CPU buffer.
- (2) The attacker runs a specifically chosen instruction that will cause the CPU to mistakenly predict that the secret data will be needed. The CPU forwards the secret data to the attacker’s instruction.
- (3) The forwarded secret data is used as the index for a memory read to an array that the attacker is authorized to access, causing a particular line of that array to be cached.
- (4) The CPU finishes checking the data and decides that the secret data was forwarded incorrectly, and reverts the execution state to before it was forwarded, but the state of the cache is not reverted.
- (5) The attacker probes all of the array to see which line was cached; the index of that line is the value of the secret data.

The original Meltdown vulnerability targeted cache forwarding and allowed data extraction in this manner from *any* memory address that was present in the cache. MDS attacks target smaller and more specific buffers in the on-core microarchitecture, and so make up a related but distinct class of attacks that are mitigated in a significantly different way. While Meltdown targets the main memory that is updated relatively infrequently and shared across all cores, threads, and processes, MDS attacks tend to target buffers that are local to cores (though sometimes shared across threads) and updated more frequently during execution.

2.4.1 Basic MDS Variants. Figure 1 charts the major known MDS attack pathways on Intel CPUs and the names given to different variants by Intel and by the researchers who reported them. Most broadly, Intel categorizes MDS vulnerabilities in their CPUs by the specific buffer from which data is speculatively forwarded, since these buffers tend to be used for a number of different operations. RIDL MDS vulnerabilities can be categorized as Microarchitectural Load Port Data Sampling (MLPDS), for variants that leak from the CPU’s load port, and Microarchitectural Fill Buffer Data Sampling (MFBDS), for variants that leak from the CPU’s LFB. Along the same lines, Intel calls the Fallout vulnerability Microarchitectural Store Buffer Data Sampling (MSBDS), as it involves a leakage from the store buffer. Vector Register Sampling (VRS) is a variant of MSBDS that targets data that is handled by vector operations as it passes through the store buffer. VERW bypass exploits a bug in the

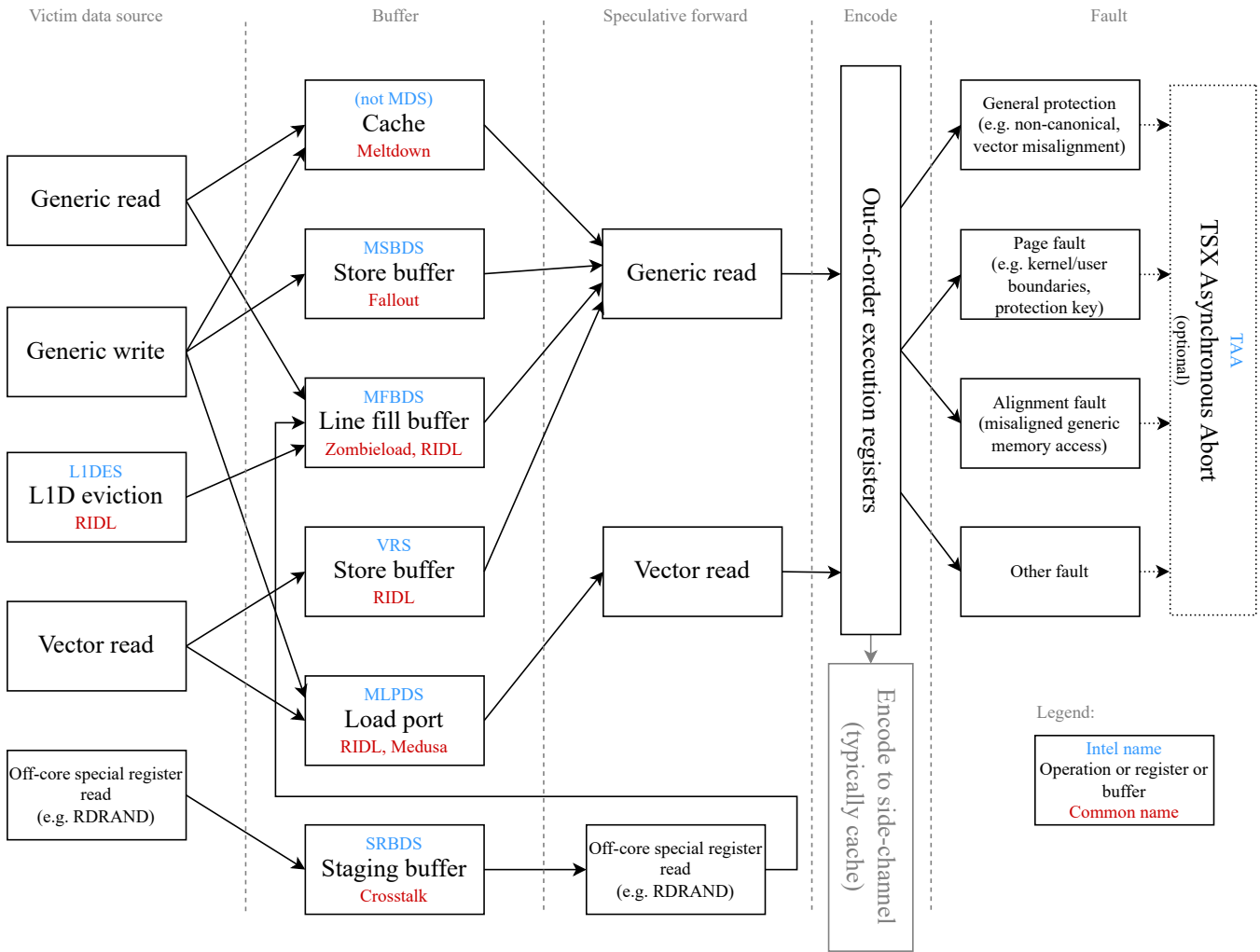


Figure 1: Major MDS attack pathways and variant names on Intel CPUs. The blue names at the top are names of vulnerabilities given by Intel; the red names at the bottom are names given by researchers or the names of the papers in which the vulnerabilities were reported. Not all fault types work with all vulnerabilities on all systems—successful forwarding and encoding during speculative execution is dependent on the exact microarchitecture including any microarchitectural countermeasures that are in place, so cataloging every known combination would be beyond the scope of this paper.

microcode fixes for MFBDS that loads stale and potentially secret data into the LFB. The basic mechanism of leakage is the same, and VERW bypass can be considered a special case of MFBDS. L1 Data Eviction Sampling (L1DES) is another special case of MFBDS, where data that is evicted from the L1 data cache passes through the LFB and becomes vulnerable to an MDS attack. Notably, L1DES is a case where the attacker can actually trigger the secret data’s presence in the CPU (by evicting it), whereas other MDS attacks rely directly on the victim process accessing the secret data to bring it into the right CPU buffer.

2.4.2 Medusa. Medusa [37] is a category of MDS attacks classified by Intel as MLPDS variants [25]. The Medusa vulnerabilities exploit the imperfect pattern-matching algorithms used to *speculatively*

combine stores in the write-combine (WC) buffer of Intel processors. Intel considers the WC buffer to be part of the load port, so Intel categorizes this vulnerability as a case of MLPDS. There are three known Medusa variants which each exploit a different feature of the write-combine buffer to cause a speculative leakage:

- Cache Indexing:** a faulting load is speculatively combined with an earlier load with a matching cache line offset.
- Unaligned Store-to-Load Forwarding:** a valid store followed by a dependent load that triggers a misaligned memory fault causes random data from the WC to be forwarded.
- Shadow REP MOV:** a faulting REP MOV instruction followed by a dependent load leaks the data of a different REP MOV.

2.4.3 TSX Asynchronous Abort. The hardware vulnerability TSX Asynchronous Abort (TAA) [24] provides a different speculation mechanism for carrying out an MDS attack. While standard MDS attacks access restricted data with a standard speculated execution, TAA uses an atomic memory transaction as implemented by TSX. When an atomic memory transaction encounters an asynchronous abort, for example because another process reads a cache line marked for use by the transaction or because the transaction encounters a fault, all operations in the transaction are rolled back to the architectural state before the transaction started. However, during this rollback, instructions inside the transaction that have already started execution can continue speculative execution, as in steps (2) and (3) of other MDS attacks. TAA impacts all Intel processors that support TSX, and the case of certain newer processors that are not affected by other MDS attacks, MDS mitigations or TAA-specific mitigations (such as disabling TSX) must be implemented in software to protect against TAA [24].

2.4.4 Mitigations. Though Meltdown and MDS-class vulnerabilities exploit low level microarchitectural operations, they can be mitigated with microcode firmware patches on most vulnerable CPUs.

Page table isolation. Historically, kernel page tables have been included in user-level process page tables so that a user-level process can make a system call to the kernel with minimal overhead. Page table isolation (first proposed by Gruss et al. as KAISER [19]) maps only the bare minimum necessary kernel memory into the user page table and introduces a second page table only accessible by the kernel. With the user process unable to access the kernel page table, accesses to all but a small and specifically chosen fraction of kernel memory are stopped before they reach the lower level caches where a Meltdown attack begins.

Buffer overwrite. MDS attacks that target on-core CPU buffers require a lower-level and more targeted defense. Intel introduced a microcode update that overwrites vulnerable buffers when the first-level data (L1d) cache (a common target of cache timing side-channel attacks) is flushed or the VERW instruction is run [25]. The kernel can then protect against MDS attacks by triggering a buffer overwrite when switching to an untrusted process.

The buffer overwrite mitigation targets MDS attacks at their source, but is imperfect to say the least. Processes remain vulnerable to attacks from concurrently running threads on the same core when SMT is enabled (since both threads share vulnerable buffers without the active process actually changing on either thread). Furthermore, shortly after the original buffer overwrite microcode update, the RIDL team found that on some Skylake CPUs, buffers were overwritten with stale and potentially sensitive data [50], and remained vulnerable even with mitigations enabled and SMT disabled. Still other processors are vulnerable to TAA but not non-TAA MDS attacks, and did not receive a buffer overwrite microcode update and as such require that TSX be disabled completely to prevent MDS attacks [20, 24].

2.5 Spectre

In 2018, Jan Horn and Paul Kocher [30] independently reported the first Spectre variants. Since then, many different Spectre variants [22, 30, 31, 33] and sub-variants [10, 13, 16, 28, 52] have been discovered. Spectre attacks make the CPU speculatively access memory that is architecturally inaccessible and leak the data into the architectural state. Therefore, all Spectre variants consist of three components [27]:

The first component is the Spectre gadget that is speculatively executed. Spectre variants are usually separated by the source of the misprediction they exploit. The outcome of a conditional direct branch, e.g., is predicted by the Pattern History Table (PHT). Mispredictions of the PHT can lead to a speculative bounds check bypass for load and store instructions [13, 28, 30]. The branch target of an indirect jump is predicted by the Branch Target Buffer (BTB). If an attacker can influence the result of a misprediction of the BTB, then speculative return-oriented programming attacks are possible [10, 13, 16, 30]. The same is true for predictions served by the Return Stack Buffer (RSB) that predicts return addresses during the execution of return instructions [13, 31, 33]. Recent results showed that some modern CPUs use the BTB for their return address predictions if the RSB underflows [52]. Another source of Spectre attacks is the prediction of store-to-load dependencies. If a load is mispredicted to *not* depend of a previous store, it speculatively executes on stale data which may lead to a speculative store bypass [22]. All of these gadgets are not exploitable by default but depend on the other two components discussed now.

The second component is how an attacker controls inputs to the aforementioned gadgets. Attackers may be able to define gadget input values directly through user input, file contents, network packets or other architectural mechanisms. On the other hand attackers may be able to inject data into the gadget transiently through load value injection [12] or floating point value injection [42]. Attackers are able to successfully control gadget inputs if they can influence which data or instructions are accessed or executed during the speculation window.

The third component is the covert channel that is used to transfer the speculative microarchitectural state into an architectural state and therefore exfiltrate the speculatively accessed data into a persistent environment. Cache covert channels [39, 40, 54] are applicable if the victim code performs a transient memory access depending on speculatively accessed secret data [30]. If a secret is accessed speculatively and loaded into an on-core buffer, an attacker can rely on an MDS-based channel [14, 46, 50] to transiently transfer the exfiltrated data to the attacker thread where the data is transferred to the architectural state through, e. g., a cache covert channel. Last but not least, if the victim executes code depending on secret data, the attacker can learn the secret by observing port contention [3, 11, 18, 43, 44].

2.5.1 Mitigations. Many countermeasures were developed to mitigate the various Spectre variants. A specific Spectre variant is effectively disabled if one of the three required components is removed. An attacker without control over inputs to Spectre gadgets is unlikely to successfully launch an attack. The same is true if a covert channel for transforming the speculative state into an architectural state is unavailable. But since this is usually hard

to guarantee, Spectre countermeasures mainly focus on stopping mispredictions. Inserting `lfence` instructions before critical code sections disable speculative execution beyond this point and can therefore be used as a generic countermeasure. But because of its high performance overhead, more specific countermeasures were developed. Spectre-BTB countermeasures include Retpoline [48] and microcode updates like IBRS, STIBP, or IBPB [23]. Spectre-RSB and Spectre-BTB-via-RSB can be mitigated by filling the RSB with values to overwrite malicious entries and prevent the RSB from underflowing or by installing IBRS microcode updates. Spectre-STL can be mitigated by the SSBD microcode update [23]. Another drastic option to stop an attacker from tampering with shared branch prediction buffers is to disable SMT. Disabling SMT effectively partitions branch prediction hardware resources between concurrent tenants at the cost of a significant performance loss.

2.6 AWS’s isolation model

Firecracker is specifically built for serverless and container applications [1] and is currently used by AWS’ Fargate CaaS and Lambda FaaS. In both of these service models, Firecracker is the primary isolation system that supports every individual Fargate task or Lambda event. Both of these service models are also designed for running very high numbers of relatively small and short-lived tasks. AWS itemizes the design requirements for the isolation system that eventually became Firecracker as follows:

Isolation: It must be safe for multiple functions to run on the same hardware, protected against privilege escalation, information disclosure, covert channels, and other risks.

Overhead and Density: It must be possible to run thousands of functions on a single machine, with minimal waste.

Performance: Functions must perform similarly to running natively. Performance must also be consistent, and isolated from the behavior of neighbors on the same hardware.

Compatibility: Lambda allows functions to contain arbitrary Linux binaries and libraries. These must be supported without code changes or recompilation.

Fast Switching: It must be possible to start new functions and clean up old functions quickly.

Soft Allocation: It must be possible to over commit CPU, memory and other resources, with each function consuming only the resources it needs, not the resources it is entitled to. [1]

We are particularly interested in the **isolation** requirement and stress that microarchitectural attacks are declared *in-scope* for the Firecracker threat model. The “design” page in AWS’s public Firecracker Git repository elaborates on the isolation model and provides a useful diagram which we reproduce in Figure 2. This diagram pertains mostly to protection against privilege escalation. The outermost layer of protection is the jailer, which uses container isolation techniques to limit the Firecracker’s access to the host kernel while running the VMM and other management components

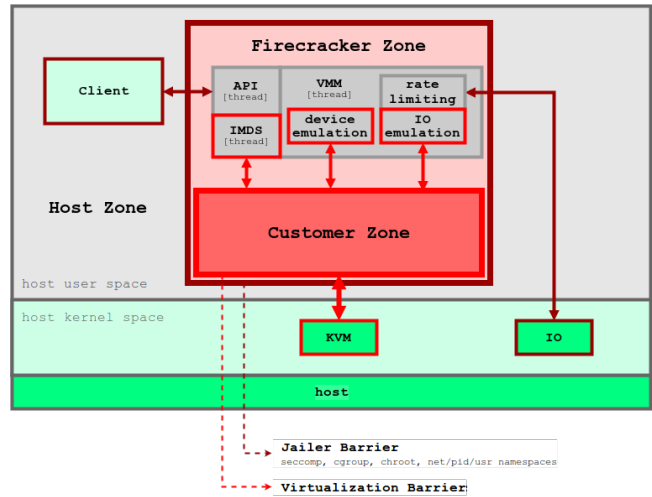


Figure 2: AWS provides this threat containment diagram in a design document in the Firecracker GitHub repository [6]. In this model, the jailer provides container-like protections around Firecracker’s VMM, API, instance metadata service (IMDS), all of which run in the host user space, and the customer’s workload, which runs inside the virtual machine. The VM isolates the customer’s workload in the guest, ensuring that it only directly interacts with predetermined elements of the host (in both user and kernel space).

of Firecracker as threads of a single process in the host userspace. Within the the Firecracker process, the user’s workload is run on other threads. The workload threads execute the guest operating system of the virtual machine and any programs running in the guest. Running the user’s code in the virtual machine guest restricts its direct interaction with the host to prearranged interactions with KVM and certain portions of the Firecracker management threads. So from the perspective of the host kernel, the VMM and the VM including the user’s code are run in the same process. This is the reason why AWS states that each VM resides in a single process. But, since the VM is isolated via hardware virtualization techniques, the user’s code, the guest kernel, and the VMM operate in separate address spaces. Therefore, the guest’s code cannot architecturally or transiently access VMM or guest kernel memory addresses as they are not mapped in the guest’s address space. The remaining microarchitectural attack surface is limited to MDS attacks that leak information from CPU internal buffers ignoring address space boundaries and Spectre attacks where an attacker manipulates the branch prediction of other processes to self-leak information.

Not shown in Figure 2, but equally important to AWS’s threat model, is the **isolation** of functions from each other when hardware is shared, especially in light of the **soft allocation** requirement. Besides the fact that compromising the host kernel could compromise the security of any guests, microarchitectural attacks that target the host hardware can also threaten user code directly. Since a single Firecracker process contains all the necessary threads to run a virtual machine with a user’s function, soft allocation can simply be performed by the host operating system [1]. This means

that standard Linux process isolation systems are in place on top of virtual machine isolation.

2.6.1 *Firecracker security recommendations.* The Firecracker documentation also recommends the following precautions for protecting against microarchitectural side-channels [8]:

- Disable SMT
- Enable kernel page-table isolation
- Disable kernel kame-page merging
- Use a kernel compiled with Spectre-BTB mitigation (e.g., IBRS and IBPB on x86)
- Verify Spectre-PHT mitigation
- Enable L1TF mitigation
- Enable Spectre-STL mitigation
- Use memory with Rowhammer mitigation
- Disable swap or use secure swap

3 THREAT MODELS

We propose two threat models applicable to Firecracker-based serverless cloud systems:

- (1) The *user-to-user* model (Figure 3): a malicious user runs arbitrary code sandboxed within a Firecracker VM and attempts to leak data, inject data, or otherwise gain information about or control over another user’s sandboxed application. In this model, we consider
 - (a) the time-sliced sharing of hardware, where the instances of the two users execute in turns on the CPU core, and
 - (b) physical co-location, where the two users’ code runs concurrently on hardware that is shared in one way or another (for example, two cores on the same CPU or two threads in the same core if SMT is enabled).
- (2) The *user-to-host* model (Figure 4): a malicious user targets some component of the host system: the Firecracker VMM, KVM, or another part of the host system kernel. For this scenario, we only consider time-sliced sharing of hardware resources. This is because the host only executes code if the guest user’s VM exits, e.g. due to a page fault that has to be handled by the host kernel or VMM.

For both models, we assume that a malicious user is able to control the runtime environment of its application. In our models, malicious users do *not* possess guest kernel privileges. Therefore, both models grant the attacker slightly less privileges than the model assumed by [1] where the guest kernel is chosen and configured by the VMM but assumed to be compromised at runtime. Rather, the attacker’s capabilities in our models match the capabilities granted to users in deployments of Firecracker in AWS Lambda and Fargate.

4 ANALYSIS OF FIRECRACKER’S CONTAINMENT SYSTEMS

Figure 2 shows the containment offered by Firecracker, as presented by AWS. In this section, we analyze each depicted component and their defenses against and vulnerabilities to microarchitectural attacks.

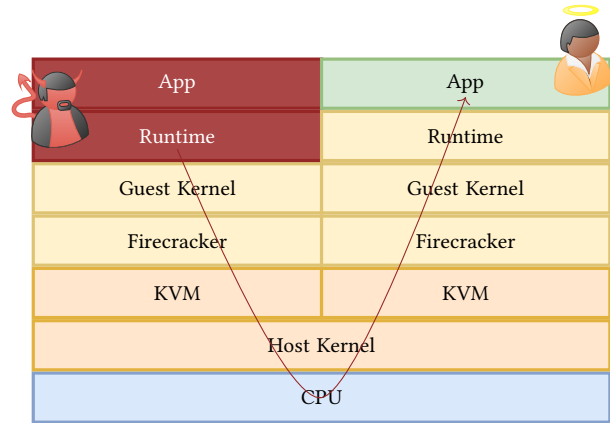


Figure 3: In the *user-to-user* threat model, we assume that a malicious cloud service tenant attempts to ex-filtrate information from another tenant. We assume the attacker to have control over the app and runtime of its VM while the guest kernel is provided by the CSP.

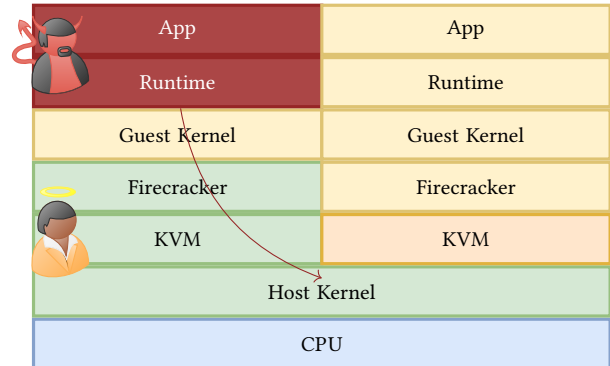


Figure 4: In the *user-to-host* threat model, the malicious tenant aims for information ex-filtration from the host system, e.g. the virtual machine manager or the host kernel. The attacker has control over the runtime and app in its virtual machine while the guest kernel is provided by the CSP.

KVM. Linux kernel-based virtual machine (KVM) is the hypervisor implemented in modern Linux kernels and therefore part of the Linux code base. It virtualizes the supervisor and user modes of the underlying hardware, manages context switches between VMs, and handles most VM-exit reasons unless they are related to I/O operations. Besides these architectural isolation mechanisms, KVM also implements mitigations against Spectre attacks on a VM-exit to protect the host OS or hypervisor from malicious guests. Firecracker heavily relies on KVM as its hypervisor. However, since KVM is part of the Linux source code and developed by the Linux community, we define KVM to not be a part of Firecracker. Therefore,

countermeasures against microarchitectural attacks that are implemented in KVM cannot be attributed to Firecracker’s containment system.

Metadata, device, and I/O services. The metadata, device, and I/O services are the parts of the Firecracker VMM and API that interact directly with a VM, collecting and managing metrics and providing connectivity. AWS touts the simplicity of these interfaces (for a reduced attack surface) and that they are written from scratch for Firecracker in Rust, a programming language known for its security features [9]. However, Rust most notably provides in-process protection against invalid and out-of-bounds memory accesses, but microarchitectural attacks like cache attacks, Spectre, and MDS can leak information between processes rather than directly hijacking a victim’s process.

Another notable difference between Firecracker and many other VMMs is that all of these services are run within the same host process as the VM itself, albeit in another thread. While the virtualization of memory addresses within the VM provides some obfuscation between the guest’s code and the I/O services, some Spectre attacks work specifically within a single process. Intra-process attacks may pose less of a threat to real world systems, however, since two guests running on the same hardware each have their own copy of these essential services.

Jailer barrier. In the event that the API or VMM are compromised, the jailer provides one last barrier of defense around a Firecracker instance. It protects the host system’s files and resources with namespaces and control groups (cgroups), respectively [7]. Microarchitectural attacks do not directly threaten files, which are by definition outside of the microarchitectural state. Cgroups allow a system administrator to assign processes to groups and then allocate, constrain, and monitor system resource usage on a per-group basis [17]. It is plausible that limitations applied with cgroups could impede an attacker’s ability to carry out certain microarchitectural attacks. For example, memory limitations might make it difficult to carry out eviction-based cache attacks, or CPU time limitations could prevent an attacker from making effective use of a CPU denial-of-service tool like HyperDegrade [2] which can slow down a victim process, simplifying the timing of a microarchitectural side-channel exfiltration or injection. In practice, Firecracker is not distributed with any particular cgroup rules [7]; in fact, it is specifically designed for the efficient operation of many Firecracker VMs under the default Linux resource allocation [6].

None of the isolation and containment systems in Firecracker seem to directly protect against user-to-user or user-to-host attacks. Therefore, we proceeded to test various microarchitectural attack proof of concepts inside and outside of Firecracker VMs.

5 ANALYSIS OF MICROARCHITECTURAL ATTACKS AND DEFENSES IN FIRECRACKER MICROVMS

In this section we present our analysis of a number of microarchitectural side-channel and speculative attack PoCs on Firecracker microVMS. We test these PoCs on bare metal and in Firecracker instances, and test relevant microcode defenses in the various scenarios. We run our tests on a server with an Intel Skylake 4114 CPU

Table 1: Presence of Medusa side-channels with all microarchitectural defense kernel options disabled. Note that the combination of cache indexing leak and block write secret (highlighted in yellow) works in Firecracker VMs but not on bare metal.

Leak	Secret	Bare Metal	Firecracker
Cache Indexing	Block Write	⇒	⇒
	REP MOV	⇒	⇒
Unaligned Store-to-Load	Block Write	⇒	⇒
	REP MOV	⇒	⇒
Shadow REP MOV	Block Write	⇒	⇒
	REP MOV	⇒	⇒

⇒ – Side-channel leakage is observable with all mitigations disabled.

⇒ – Side-channel leakage is not observable.

which has virtualization hardware extensions and SMT enabled. The CPU runs on microcode version 0x2006b06². The host OS is Ubuntu 20.04 with a Linux 5.10 kernel. We used Firecracker v1.0.0 and v1.4.0, the latest version as of July 2023, to run an Ubuntu 18.04 guest with Linux kernel 5.4 which is provided by Amazon when following the quick-start guide³.

In summary, the recommended production host setup provided with AWS Firecracker is insufficient when it comes to protecting tenants from malicious neighbors. Firecracker therefore fails in providing its claimed isolation guarantees. This is because

- (1) we identify a Medusa variant that only becomes exploitable when it is run across microVMS. In addition, the recommended countermeasures do not contain the necessary steps to mitigate the side-channel, or most other Medusa variants.
- (2) we show that tenants are not properly protected from information leaks induced through Spectre-PHT or Spectre-BTB when applying the recommended countermeasures. The Spectre-PHT variants remain a problem even when disabling SMT.
- (3) we observed no differences in PoC performance between Firecracker v1.0.0 and v1.4.0.

We conclude that the virtualization layer provided by Firecracker has little effect on microarchitectural attacks, and Firecracker’s system security recommendations are insufficient.

5.1 Medusa

We evaluated Moghimi’s PoCs [35] for the Medusa [37] side-channels (classified by Intel as MLPDS variants of MDS [25]) on the bare metal of our test system and in Firecracker VMs. There is one leaking PoC for each of the three known variants described in section 2.4.2. We used two victim programs from the PoC library:

²Updating the microcode to a newer version would disable TSX on our system which would make tests with TSX-based MDS variants impossible.

³<https://github.com/firecracker-microvm/firecracker/blob/dbd9a84b11a63b5e5bf201e244fe83f0bc76792a/docs/getting-started.md>

Table 2: Mitigations necessary to protect bare metal vs. Firecracker victims from Medusa attacks. Note that AWS’s recommended mitigation, `nosmt`, does not prevent the highlighted cache indexing/block write variant that is enabled by Firecracker (cf. Table 1), or any other variant except shadow REP MOV.

Leak	Secret	Bare metal		Firecracker		
		mds	nosmt	mds(VM)	mds(H)	nosmt
Cache Indexing	Block Write	N/A	N/A	✗	✓	✗
	REP MOV	✓	✓	✓	✓	✓
Unaligned Store	Block Write	✓	✗	✗	✓	✗
	REP MOV	+	+	✗	+	+
Shadow REP MOV	Block Write	N/A	N/A	N/A	N/A	N/A
	REP MOV	✓	✓	+	✓	+

✓ – mitigation prevents side-channel attack.

+ – mitigation prevents attack only in combination with other mitigation(s) marked with +.

✗ – mitigation has no effect on this attack.

- The “Block Write” program writes a large amount of consecutive data in a loop (so that the processor will identify repeated stores and combine them).
- The “REP MOV” program performs a similar operation, but with the REP MOV instruction instead of many instructions moving smaller blocks of data in a loop.

5.1.1 Results. Table 1 shows the cases in which data is successfully leaked with all microarchitectural protections in the kernel disabled. The left two columns show the possible combinations of the three Medusa PoCs and the two included victim programs. The right columns indicate which configurations work on bare metal and with the secret and leaking program running in parallel Firecracker instances. Most notably, with the Cache Indexing variant, the Block Write secret only works with Firecracker. This is likely because of the memory address virtualization that the virtual machine provides: the guest only sees virtual memory regions mapped by KVM, and KVM traps memory access instructions and handles the transactions on behalf of the guest.

We tested the effectiveness of `mds` and `nosmt` defenses against each combination of attacker and victim PoC on bare metal and in Firecracker VMs. Table 2 lists the protections necessary to prevent Medusa attacks in bare metal and Firecracker scenarios. Across the four vulnerabilities in Firecracker, only one is mitigated by `nosmt` alone, and AWS does *not explicitly recommend* enabling the `mds` protection, though most Linux distributions ship with it enabled by default. That is to say, a multi-tenant cloud platform could be using Firecracker in a way that is compliant with AWS’s recommended security measures and still be vulnerable to the majority of Medusa variants, including one where the Firecracker VMM itself leaks the user’s data that would not otherwise be leaked.

5.2 RIDL and More

In this section, we present an evaluation of the RIDL PoC programs [51] provided alongside van Schaik et al.’s 2019 paper [50]. RIDL is a class of MDS attacks that exploits speculative loads from buffers inside the CPU (not from cache or memory). The RIDL PoC repository also includes examples of attacks released in later addenda to the paper as well as one variant of the Fallout MDS attack.

5.2.1 Results. Table 3 shows some basic information about the RIDL PoCs that we tested and the efficacy of relevant countermeasures at preventing the attacks. We compared attacks on bare metal and in Firecracker to evaluate Amazon’s claims of the heightened hardware security of the Firecracker microVM system. For tests on the Firecracker system, we distinguish between countermeasure flags enabled on the host system (H) and the Firecracker guest kernel (VM). Besides the `nosmt` and `mds` kernel flags, we tested other relevant flags (cf. section 2.4.4, [21]), including `kaslr`, `pti`, and `l1tf`, but did not find that they had an effect on any of these programs. We excluded the `tsx_async_abort` mitigation since the CPU we tested on includes `mds` mitigation which makes the `tsx_async_abort` kernel flag redundant [20].

In general, we found that the `mds` protection does not adequately protect against the majority of RIDL attacks. However, disabling SMT does mitigate the majority of these exploits. This is consistent with Intel’s [25] and the Linux developers’ [21] statements that SMT must be disabled to prevent MDS attacks across hyperthreads. The two outliers among these PoCs are `alignment_write`, which requires both `nosmt` and `mds` on the host, and `pgtable_leak_notsox`, which is mitigated only by `mds` countermeasures. The leak relying on `alignment_write` uses an alignment fault rather than a page table fault leak to trigger speculation [50]. The RIDL paper [50] and Intel’s documentation of the related VRS exploit [26] are unclear about what exactly differentiates this attack from the page-fault-based MFBDS attacks found in other PoCs, but our experimental findings indicate that the microarchitectural mechanism of the leakage is distinct. There is a simple and reasonable explanation for the behavior of `pgtable_leak_notsox`, which is unique among these PoCs for one key reason: it is the only exploit that crosses security boundaries (leaking page table values from the kernel) within a single thread rather than leaking from another thread. It is self-evident that disabling multi-threading would have little effect on this single-threaded exploit. However, the `mds` countermeasure flushes microarchitectural buffers before switching from kernel-privilege execution to user-privilege execution within the same thread, wiping the page table data accessed by kernel code from the LFB before the attacking user code can leak it.

In contrast to Medusa, most of these PoCs are mitigated by AWS’s recommendation of disabling `smt`. However, as with Medusa, the Firecracker VMM itself provides no microarchitectural protection against these attacks.

5.3 Spectre

Next we focus on Spectre vulnerabilities. While there have been many countermeasures developed since Spectre attacks were first discovered, many of them either come with a (significant) performance penalty or only partially mitigate the attack. Therefore,

Table 3: Mitigations necessary to protect bare metal vs. Firecracker victims from RIDL and other MDS attacks. The recommended nosmt mitigation protects against most but not all of these variants. All proof of concepts were tested on Firecracker v1.0.0 and v1.4.0 with identical results.

Exploit	Details			Bare Metal		Firecracker			TSX required?
	Common Name	Target Buffer	Fault Type	nosmt	mds	nosmt (H)	mds (H)	mds (VM)	
alignment_write		Fill Buffer	Alignment	+	+	+	+	×	no
pgtable_leak_notsx	RIDL	Fill Buffer	Page	× ^b	✓	× ^b	✓	✓	no
ridl_basic		Fill Buffer	Page	✓	×	N/A ^c	N/A ^c	N/A ^c	no
ridl_invalidpage		Fill Buffer	Page	✓	×	N/A ^c	N/A ^c	N/A ^c	no
pgtable_leak	RIDL/TAA	Fill Buffer	TSX abort	✓	×	✓	×	×	yes
taa_read		Fill Buffer	TSX abort	✓	×	✓	×	×	yes
taa_basic		Fill Buffer	TSX abort	✓	×	N/A ^c	N/A ^c	N/A ^c	yes
verw_bypass_l1des	RIDL/TAA	Fill Buffer	TSX abort	✓	×	✓	×	×	yes
loadport	RIDL	Load Port	Page	✓	×	✓	×	×	no
vrs	RIDL/VRS	Store buffer	Page	✓	×	✓	×	×	no
cpuid_leak	Crosstalk	Fill Buffer	Page	✓	×	N/A ^a	N/A ^a	N/A ^a	no

✓ – mitigation prevents side-channel attack.

+ – mitigation prevents attack only in combination with other mitigation(s) marked with +.

× – mitigation has no effect on this attack.

^a CPUID instruction is emulated by virtual machine and has no microarchitectural effect.

^b This attack leaks information about pages used in its own thread.

^c PoCs had to be modified slightly to run in two processes before they could be tested in the virtual machine. These PoCs did not work on bare metal or in virtual machines when split into two processes.

system operators often have to decide for a performance vs. security trade-off. In this section we evaluate the Spectre attack surface available to Firecracker tenants in both threat models described earlier. To evaluate the wide range of Spectre attacks, we rely on the PoCs provided in [15]. For Spectre-PHT, Spectre-BTB, and Spectre-RSB, the repository contains four PoCs each. They differ in the way the attacker mistrains the BPU. The four possibilities are (1) *same-process*—the attacker has control over the victim process or its inputs to mistrain the BPU—(2) *cross-process*—the attacker runs its own code in a separate process to influence the branch predictions of the victim process—(a) *in-place*—the attacker mistrains the the BPU with branch instruction that reside at the same virtual address as the target branch that the attacker wants to be misspredicted in the victim process—(b) *out-of-place*—the attacker mistrains the BPU with branch instructions that reside at addresses that are congruent to the target branches in the victim process.

- (1) *same-process*: the attacker has control over the victim process or its inputs to mistrain the BPU,
- (2) *cross-process*: the attacker runs its own code in a separate process to influence the branch predictions of the victim process,
- (3) *in-place*: the attacker mistrains the the BPU with branch instruction that reside at the same virtual address as the target branch that the attacker wants to be misspredicted in the victim process
- (4) *out-of-place*: the attacker mistrains the BPU with branch instructions that reside at addresses that are congruent to the target branches in the victim process.

The first two and latter two situations are orthogonal, so each PoC combines two of them. For Spectre-STL, only same-process variants are known, which is why the repository only provides two PoCs in this case. For cross-VM experiments, disabled address space layout randomization for the host and guest kernels as well as for the host and guest user level to ease finding congruent addresses that are used for mistraining.

5.3.1 Results. With *AWS recommended countermeasures* [8] (the default for the Linux kernels in use, cf. Figure 5) enabled on the host system and inside Firecracker VMs, we see that Spectre-RSB is successfully mitigated both on the host and inside and across VMs (cf. Table 4). On the other hand, *Spectre-STL*, *Spectre-BTB*, and *Spectre-PHT* allowed information leakage in particular situations.

The PoCs for Spectre-STL show leakage. However, the leakage only occurs within the same process and the same privilege level. Since no cross-process variants are known, we didn’t test the cross-VM scenario for Spectre-STL. In our user-to-user threat model, Spectre-STL is not a possible attack vector, as no cross-process variants are known. If two tenant workloads would be isolated by in-process isolation within the same VM, Spectre-STL could still be a viable attack vector. In the user-to-host model, Spectre-STL is mitigated by countermeasures that are included in current Linux kernels and enabled by default.

For Spectre-PHT, the kernel countermeasures include the sanitization of user-pointers and the utilization of barriers (lfence) on privilege level switches. We therefore conclude that Spectre-PHT poses little to no threat to the host system. However, these

Table 4: Spectre PoCs run with AWS Firecracker recommended countermeasures (cf. Figure 5 and [8]). These countermeasures—which are the default for the Linux kernels in use—are insufficient when it comes to protecting tenants from Spectre attacks. Experiments with Firecracker v1.0.0 and v1.4.0 yielded the same results.

Variant	Platform	Thread Configuration	same-process		cross-process	
			in-place	out-of-place	in-place	out-of-place
PHT	☰ v1.0.0	any	✗	✗	✗	✗
		any	✗	✗	✗	✗
		any	✗	✗	✗	✗
	↔ v1.0.0	⚡(1PT), 1vCPU each	N/A	N/A	✗	✗
		⚡(2PT), 1vCPU each	N/A	N/A	✗	✓
	↔ v1.4.0	⚡(1PT), 1vCPU each	N/A	N/A	✗	✗
⚡(2PT), 1vCPU each		N/A	N/A	✗	✓	
BTB	☰	⚡(1PT)	✗	✓	✗	✓
		⚡(2PT)	✗	✓	✗	✓
	v1.0.0	⚡(1PT), 1vCPU	✗	✓	✗	✓
		⚡(1PT), 2vCPU	✗	✓	✓	✓
		⚡(2PT), 2vCPU	✗	✓	✗	✓
	v1.4.0	⚡(1PT), 1vCPU	✗	✓	✗	✓
		⚡(1PT), 2vCPU	✗	✓	✓	✓
		⚡(2PT), 2vCPU	✗	✓	✗	✓
	↔ v1.0.0	⚡(1PT), 1vCPU each	N/A	N/A	✓	✓
		⚡(2PT), 1vCPU each	N/A	N/A	✗	✓
	↔ v1.4.0	⚡(1PT), 1vCPU each	N/A	N/A	✓	✓
		⚡(2PT), 1vCPU each	N/A	N/A	✗	✓
RSB	☰	any	✓	✓	✓	✓
		any	✓	✓	✓	✓
	↔ v1.0.0	any	✓	✓	✓	✓
		any	N/A	N/A	✓	✓
		any	N/A	N/A	✓	✓
STL	☰	any	✗	N/A	N/A	N/A
	v1.0.0	any	✗	N/A	N/A	N/A
	v1.4.0	any	✗	N/A	N/A	N/A

✓ – mitigation prevents side-channel attack.

✚ – mitigation prevents attack only in combination with other mitigation(s) marked with ✚.

✗ – mitigation has no effect on this attack.

☰ – Experiments were run on the host system

🔥 v1.x.0 – Experiments were run inside a single Firecracker v1.x.0 VM

↔ v1.x.0 – Experiments were run across two Firecracker v1.x.0 VMs

⚡(XPT) – bare metal PoC processes/Firecracker VM vCPUs were pinned to X physical sibling threads.

XvCPU – the Firecracker VM has X virtual CPUs.

^a the PoC does not support running without being assigned to a core.

mitigations do not protect two hyperthreads from each other if they execute on the same physical core in parallel. This is why all four Spectre-PHT mistraining variants are fully functional on the host system as well as inside Firecracker VMs. As can be seen in Table 4, this remains true *even if SMT is disabled*⁴. In fact, pinning both VMs to the same physical thread *enables* the cross-process

out-of-place version of Spectre-PHT whereas we did not observe leaks in the SMT case. This makes Spectre-PHT a significant threat for user-to-user.

Spectre-BTB PoCs are partially functional when AWS recommended countermeasures are enabled. The original variant that mistrains the BTB in the same process and at the same address is fully functional while same-process out-of-place mistraining is

⁴This is simulated by pinning attacker and victim process to the same core (⚡(1PT))

```

1 Vulnerability Spec store bypass: Mitigation; Speculative Store Bypass disabled via prctl and seccomp
2 Vulnerability Spectre v1:      Mitigation; usercopy/swapgs barriers and __user pointer sanitization
3 Vulnerability Spectre v2:      Mitigation; Full generic retpoline, IBPB conditional, IBRS_FW, STIBP conditional, RSB
    filling

```

Figure 5: Spectre mitigations enabled in the host and guest kernel during the Spectre tests. This setup is recommended by AWS for host production systems [8].

successfully mitigated. Also, all attempts to leak information across process boundaries via out-of-place mistraining did not show any leakage. With cross-process in-place mistraining, however, we observed leakage. On the host system, the leakage occurred independent of SMT. Inside a VM, the leakage only occurred if all virtual CPU cores were assigned to a separate physical thread. Across VMs, disabling SMT removed the leakage.

Besides the countermeasures listed in Figure 5, the host kernel has Spectre countermeasures compiled into the VM entry and exit point⁵ to fully disable malicious guests from attacking the host kernel while the kernel handles a VM exit.

In summary, we can say that the Linux default countermeasures—which are recommended by AWS Firecracker—only partially mitigate Spectre. Precisely, we show:

- Spectre-PHT and Spectre-BTB can still leak information between tenants in the guest-to-guest scenario with the AWS recommended countermeasures—which includes disabling SMT—in place.
- The host kernel is likely sufficiently protected by the additional precautions that are compiled into the Linux kernel to shield VM entries and exits. This, however, is orthogonal to security measures provided by Firecracker.

All leakage observed was independent of the Firecracker version in use.

5.3.2 Evaluation. We find that Firecracker does not add to the mitigations against Spectre but solely relies on general protection recommendations, which include mitigations provided by the host and guest kernels and optional microcode updates. Even worse, the recommended countermeasures insufficiently protect serverless applications from leaking information to other tenants. We therefore claim that Firecracker does *not* achieve its isolation goal on a microarchitectural level, even though microarchitectural attacks are considered in-scope of the Firecracker threat model.

The alert reader might wonder why Spectre-BTB remains an issue with the STIBP countermeasure in place (cf. Figure 5) as this microcode patch was designed to stop the branch prediction from using prediction information that originates from another thread. This also puzzled us for a while until recently Google published a security advisory⁶ that identifies a flaw in Linux 6.2 that kept disabling the STIBP mitigation when IBRS is enabled. We verified that the code section that was identified as being responsible for the issue is also present in the Linux 5.10 source code. Our assumption therefore is that the same problem identified by Google also occurs on our system.

6 CONCLUSIONS

Cloud technologies constantly shift to meet the needs of their customers. At the same time, CSPs aim for maximizing efficiency and profit, which incentivizes serverless CSPs to over-commit available compute resources. While this is reasonable from an economic perspective, the resulting system behavior can be disastrous in the context of microarchitectural attacks that exploit shared hardware resources. In the past few years, the microarchitectural threat landscape changed frequently and rapidly. There mitigations that work reasonably well to prevent many attacks, but they often lead to significant performance costs, which forces CSPs to find a tradeoff between economic value and security. Furthermore, some microarchitectural attacks simply are not hindered by existing mitigations. The CSP customers have little control over the microarchitectural defenses deployed and must trust their providers to keep up with the pace of microarchitectural attack and mitigation development. Defense-in-depth requires security at every level, from the microcode to VMM to container. Each system must be considered as a whole, as some protections at one system level may open a vulnerabilities at another.

We showed that default countermeasures as they are recommended for the Firecracker VMM are insufficient to meet its isolation goals. In fact, many of the tested attack vectors showed leakage while countermeasures were in place. We identified the Medusa cache indexing/block write variant as an attack vector that only works across VMs, i. e. with additional isolation mechanisms in place. Additionally, we showed that disabling SMT—an expensive mitigation technique recommended and performed by AWS—does not provide full protection against Medusa variants. The aforementioned Medusa variant, and Spectre-PHT are still capable of leaking information between cloud tenants even if SMT is disabled, as long as the attacker and target threads keep competing for hardware resources of the same physical CPU core. Unfortunately this is inevitably the case in high-density serverless environments. In the present, serverless CSPs must remain vigilant in keeping firmware up-to-date and employing all possible defenses against microarchitectural attacks. Users must not only trust their CSPs of choice to keep their systems up-to-date and properly configured, but also be aware that some microarchitectural vulnerabilities, particularly certain Spectre variants, are still able to cross containment boundaries. Furthermore, processor designs continue to evolve and speculative and out-of-order execution remain important factors in improving performance from generation to generation. So, it is unlikely that we have seen the last of new microarchitectural vulnerabilities, as the recent wave of newly discovered attacks [36, 47, 53] shows.

⁵<https://elixir.bootlin.com/linux/v5.10/source/arch/x86/kvm/vmx/vmenter.S#L191>

⁶<https://github.com/google/security-research/security/advisories/GHSA-mj4w-6495-6cxc>

ACKNOWLEDGMENTS

This work was supported by the German Research Foundation (DFG) under Grants No. 439797619 and 456967092, by the German Federal Ministry for Education and Research (BMBF) under Grants SASVI and SILGENTAS, by the National Science Foundation (NSF) under Grant CNS-2026913, and in part by a grant from the Qatar National Research Fund.

REFERENCES

- [1] Alexandru Agache, Marc Brooker, Alexandra Iordache, Anthony Liguori, Wolf Neugebauer, Phil Piwonka, and Diana-Maria Popa. 2020. Firecracker: Lightweight Virtualization for Serverless Applications. In *NSDI*. USENIX Association, 419–434.
- [2] Alejandro Cabrera Aldaya and Billy Bob Brumley. 2022. HyperDegradate: From GHz to MHz Effective CPU Frequencies. In *USENIX Security Symposium*. USENIX Association, 2801–2818.
- [3] Alejandro Cabrera Aldaya, Billy Bob Brumley, Sohaib ul Hassan, Cesar Pereira García, and Nicola Tuveri. 2019. Port Contention for Fun and Profit. In *IEEE Symposium on Security and Privacy*. IEEE, 870–887.
- [4] Amazon Web Services. 2023. AWS Fargate. <https://docs.aws.amazon.com/eks/latest/userguide/fargate.html> accessed: Aug 17, 2023.
- [5] Amazon Web Services. 2023. AWS Lambda Features. <https://aws.amazon.com/lambda/features/> accessed: Aug 17, 2023.
- [6] Amazon Web Services. 2023. Firecracker Design. <https://github.com/firecracker-microvm/firecracker/blob/9c51dc6852d68d0f6982a4017a63645fa75460c0/docs/design.md>.
- [7] Amazon Web Services. 2023. The Firecracker Jailer. <https://github.com/firecracker-microvm/firecracker/blob/main/docs/jailer.md> accessed: August 14, 2023.
- [8] Amazon Web Services. 2023. Production Host Setup Recommendations. <https://github.com/firecracker-microvm/firecracker/blob/9ddea322a74c20cfb6b5af745112c95b7cecb75/docs/prod-host-setup.md> accessed: May 22, 2023.
- [9] Abhiram Balasubramanian, Marek S. Baranowski, Anton Burtsev, Aurojit Panda, Zvonimir Rakamaric, and Leonid Ryzhyk. 2017. System Programming in Rust: Beyond Safety. In *HotOS*. ACM, 156–161.
- [10] Enrico Barberis, Pietro Frigo, Marius Muench, Herbert Bos, and Cristiano Giuffrida. 2022. Branch History Injection: On the Effectiveness of Hardware Mitigations Against Cross-Privilege Spectre-v2 Attacks. In *USENIX Security Symposium*. USENIX Association, 971–988.
- [11] Atri Bhattacharyya, Alexandra Sandulescu, Matthias Neuschwandtner, Alessandro Sorniotti, Babak Falsafi, Mathias Payer, and Anil Kurmus. 2019. SMOthersSpectre: Exploiting Speculative Execution through Port Contention. In *CCS*. ACM, 785–800.
- [12] Jo Van Bulck, Daniel Moghimi, Michael Schwarz, Moritz Lipp, Marina Minkin, Daniel Genkin, Yuval Yarom, Berk Sunar, Daniel Gruss, and Frank Piessens. 2020. LVI: Hijacking Transient Execution through Microarchitectural Load Value Injection. In *IEEE Symposium on Security and Privacy*. IEEE, 54–72.
- [13] Claudio Canella, Jo Van Bulck, Michael Schwarz, Moritz Lipp, Benjamin von Berg, Philipp Ortner, Frank Piessens, Dmitry Evtushkin, and Daniel Gruss. 2019. A Systematic Evaluation of Transient Execution Attacks and Defenses. In *USENIX Security Symposium*. USENIX Association, 249–266.
- [14] Claudio Canella, Daniel Genkin, Lukas Giner, Daniel Gruss, Moritz Lipp, Marina Minkin, Daniel Moghimi, Frank Piessens, Michael Schwarz, Berk Sunar, Jo Van Bulck, and Yuval Yarom. 2019. Fallout: Leaking Data on Meltdown-resistant CPUs. In *CCS*. ACM, 769–784.
- [15] Claudio Canella, Jo Van Bulck, Michael Schwarz, Daniel Gruss, Catherine Easdon, and Saagar Jha. 2019. Transient Fail [Source Code]. <https://github.com/LAIK/transientfail>
- [16] Guoxing Chen, Sanchuan Chen, Yuan Xiao, Yinqian Zhang, Zhiqiang Lin, and Ten-Hwang Lai. 2019. SgxPectre: Stealing Intel Secrets from SGX Enclaves Via Speculative Execution. In *EuroS&P*. IEEE, 142–157.
- [17] Marie Doležalová, Milan Navrátil, Eva Majoršínová, Peter Ondrejka, Douglas Silas, Martin Prpič, and Rüdiger Landmann. 2020. *Red Hat Enterprise Linux 7 Resource Management Guide—Using cgroups to manage system resources on RHEL*. Red Hat, Inc. https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux/7/pdf/resource_management_guide/red_hat_enterprise_linux-7-resource_management_guide-en-us.pdf accessed: Aug 17, 2023.
- [18] Jacob Fustos, Michael Garrett Bechtel, and Heechul Yun. 2020. SpectreRewind: Leaking Secrets to Past Instructions. In *ASHES@CCS*. ACM, 117–126.
- [19] Daniel Gruss, Moritz Lipp, Michael Schwarz, Richard Fellner, Clémentine Maurice, and Stefan Mangard. 2017. KASLR is Dead: Long Live KASLR. In *ESSoS (Lecture Notes in Computer Science, Vol. 10379)*. Springer, 161–176.
- [20] Pawan Gupta. 2020. *TAA - TSX Asynchronous Abort*. The Linux Kernel Organization. https://www.kernel.org/doc/html/latest/admin-guide/hw-vuln/tsx_async_abort.html accessed: Aug 17, 2023.
- [21] Tyler Hicks. 2019. *MDS - Microarchitectural Data Sampling*. The Linux Kernel Organization. <https://www.kernel.org/doc/html/latest/admin-guide/hw-vuln/mds.html> accessed: Aug 17, 2023.
- [22] Jann Horn. 2018. Speculative execution, variant 4: speculative store bypass. <https://bugs.chromium.org/p/project-zero/issues/detail?id=1528> accessed: Aug 17, 2023.
- [23] Intel. 2018. Speculative Execution Side Channel Mitigations. <https://www.intel.com/content/dam/develop/external/us/en/documents/336996-speculative-execution-side-channel-mitigations.pdf>. rev. 3.0 accessed: Mar 22, 2023.
- [24] Intel. 2019. *Intel Transactional Synchronization Extensions (Intel TSX) Asynchronous Abort*. Technical Report. Intel Corp. <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/technical-documentation/intel-tsx-asynchronous-abort.html> accessed: Aug 17, 2023.
- [25] Intel. 2019. *Microarchitectural Data Sampling*. Technical Report. Intel Corp. <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/technical-documentation/intel-analysis-microarchitectural-data-sampling.html> ver. 3.0, accessed: Aug 17, 2023.
- [26] Intel. 2020. *Vector Register Sampling*. Technical Report. Intel Corp. <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/advisory-guidance/vector-register-sampling.html> accessed: Aug 17, 2023.
- [27] Brian Johannsmeyer, Jakob Koschel, Kaveh Razavi, Herbert Bos, and Cristiano Giuffrida. 2022. Kasper: Scanning for Generalized Transient Execution Gadgets in the Linux Kernel. In *NDSS*. The Internet Society.
- [28] Vladimir Kiriansky and Carl A. Waldspurger. 2018. Speculative Buffer Overflows: Attacks and Defenses. *CoRR* abs/1807.03757 (2018).
- [29] Avi Kivity, Yaniv Kamay, Dor Laor, Uri Lubin, and Anthony Liguori. 2007. kvm: the Linux Virtual Machine Monitor. In *Linux Symposium*, Vol. 1. kernel.org, 225–230.
- [30] Paul Kocher, Jann Horn, Anders Fogh, Daniel Genkin, Daniel Gruss, Werner Haas, Mike Hamburg, Moritz Lipp, Stefan Mangard, Thomas Prescher, Michael Schwarz, and Yuval Yarom. 2019. Spectre Attacks: Exploiting Speculative Execution. In *IEEE Symposium on Security and Privacy*. IEEE, 1–19.
- [31] Esmail Mohammadian Koruyeh, Khaled N. Khasawneh, Chengyu Song, and Nael B. Abu-Ghazaleh. 2018. Spectre Returns! Speculation Attacks using the Return Stack Buffer. In *WOOT @ USENIX Security Symposium*. USENIX Association.
- [32] Moritz Lipp, Michael Schwarz, Daniel Gruss, Thomas Prescher, Werner Haas, Anders Fogh, Jann Horn, Stefan Mangard, Paul Kocher, Daniel Genkin, Yuval Yarom, and Mike Hamburg. 2018. Meltdown: Reading Kernel Memory from User Space. In *USENIX Security Symposium*. USENIX Association, 973–990.
- [33] Giorgi Maisuradze and Christian Rossow. 2018. ret2spec: Speculative Execution Using Return Stack Buffers. In *CCS*. ACM, 2109–2122.
- [34] Debora T. Marr, Frank Binns, David L. Hill, Glenn Hinton, David A. Koufaty, J. Alan Miller, and Michael Upton. 2002. Hyper-Threading Technology Architecture and Microarchitecture. *Intel Technology Journal* 6, 1 (2002), 4–15.
- [35] Daniel Moghimi. 2020. Medusa Code Repository [Source Code]. <https://github.com/flowyroll/medusa>
- [36] Daniel Moghimi. 2023. Downfall: Exploiting Speculative Data Gathering. In *USENIX Security Symposium*. USENIX Association, 7179–7193.
- [37] Daniel Moghimi, Moritz Lipp, Berk Sunar, and Michael Schwarz. 2020. Medusa: Microarchitectural Data Leakage via Automated Attack Synthesis. In *USENIX Security Symposium*. USENIX Association, 1427–1444.
- [38] Shraavan Narayan, Craig Disselkoben, Daniel Moghimi, Sunjay Cauligi, Evan Johnson, Zhao Gang, Anjo Vahldiek-Oberwagner, Ravi Sahita, Hovav Shacham, Dean M. Tullsen, and Deian Stefan. 2021. Swivel: Hardening WebAssembly against Spectre. In *USENIX Security Symposium*. USENIX Association, 1433–1450.
- [39] Dag Arne Osvik, Adi Shamir, and Eran Tromer. 2006. Cache Attacks and Countermeasures: The Case of AES. In *CT-RSA (Lecture Notes in Computer Science, Vol. 3860)*. Springer, 1–20.
- [40] Antoon Purnal, Furkan Turan, and Ingrid Verbauwhede. 2021. Prime+Scope: Overcoming the Observer Effect for High-Precision Cache Contention Attacks. In *CCS*. ACM, 2906–2920.
- [41] Qumranet Inc. 2006. *KVM: Kernel-based Virtualization Driver, White Paper*. Technical Report. Qumranet Inc. <https://docs.huihoo.com/kvm/kvm-white-paper.pdf> accessed: Aug 17, 2023.
- [42] Hany Ragab, Enrico Barberis, Herbert Bos, and Cristiano Giuffrida. 2021. Rage Against the Machine Clear: A Systematic Analysis of Machine Clears and Their Implications for Transient Execution Attacks. In *USENIX Security Symposium*. USENIX Association, 1451–1468.
- [43] Thomas Rokicki, Clémentine Maurice, Marina Botvinnik, and Yossi Oren. 2022. Port Contention Goes Portable: Port Contention Side Channels in Web Browsers. In *AsiaCCS*. ACM, 1182–1194.

- [44] Thomas Rokicki, Clémentine Maurice, and Michael Schwarz. 2022. CPU Port Contention Without SMT. In *ESORICS (3) (Lecture Notes in Computer Science, Vol. 13556)*. Springer, 209–228.
- [45] David Schrammel, Samuel Weiser, Stefan Steinegger, Martin Schwarzl, Michael Schwarz, Stefan Mangard, and Daniel Gruss. 2020. Donky: Domain Keys - Efficient In-Process Isolation for RISC-V and x86. In *USENIX Security Symposium*. USENIX Association, 1677–1694.
- [46] Michael Schwarz, Moritz Lipp, Daniel Moghimi, Jo Van Bulck, Julian Stecklina, Thomas Prescher, and Daniel Gruss. 2019. ZombieLoad: Cross-Privilege-Boundary Data Sampling. In *CCS*. ACM, 753–768.
- [47] Daniël Trujillo, Johannes Wikner, and Kaveh Razavi. 2023. Inception: Exposing New Attack Surfaces with Training in Transient Execution. In *USENIX Security Symposium*. USENIX Association, 7303–7320.
- [48] Paul Turner. 2018. Retpoline: a software construct for preventing branch-target-injection. <https://support.google.com/faqs/answer/7625886>. accessed: Mar 22, 2023.
- [49] Anjo Vahldiek-Oberwagner, Eslam Elnikety, Nuno O. Duarte, Michael Sammler, Peter Druschel, and Deepak Garg. 2019. ERIM: Secure, Efficient In-process Isolation with Protection Keys (MPK). In *USENIX Security Symposium*. USENIX Association, 1221–1238.
- [50] Stephan van Schaik, Alyssa Millburn, Sebastian Österlund, Pietro Frigo, Giorgi Maisuradze, Kaveh Razavi, Herbert Bos, and Cristiano Giuffrida. 2019. RIDL: Rogue In-Flight Data Load. In *IEEE Symposium on Security and Privacy*. IEEE, 88–105.
- [51] Stephan van Schaik, Alyssa Millburn, genBTC, Paul Menzel, jun1x, Stephen Kitt, pit fr, Sebastian Österlund, and Cristiano Giuffrida. 2020. RIDL [Source Code]. <https://github.com/vusec/ridl>
- [52] Johannes Wikner and Kaveh Razavi. 2022. RETBLEED: Arbitrary Speculative Code Execution with Return Instructions. In *USENIX Security Symposium*. USENIX Association, 3825–3842.
- [53] Johannes Wikner, Daniël Trujillo, and Kaveh Razavi. 2023. Phantom: Exploiting Decoder-detectable Mispredictions. In *MICRO (to appear)*. IEEE.
- [54] Yuval Yarom and Katrina Falkner. 2014. FLUSH+RELOAD: A High Resolution, Low Noise, L3 Cache Side-Channel Attack. In *USENIX Security Symposium*. USENIX Association, 719–732.
- [55] Ethan G. Young, Pengfei Zhu, Tyler Caraza-Harter, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. 2019. The True Cost of Containing: A gVisor Case Study. In *HotCloud*. USENIX Association.