

Hardware-Assisted Fault Isolation: Going Beyond the Limits of Software-Based Sandboxing

Shravan Narayan , University of Texas at Austin, Austin, TX, 78712, USA

Tal Garfinkel , University of California, San Diego, San Diego, CA, 92093, USA

Mohammadkazem Taram , Purdue University, West Lafayette, IN, 47907, USA

Joey Rudek , University of California, San Diego, San Diego, CA, 92093, USA

Daniel Moghimi , Google, Mountain View, CA, 94043, USA

Evan Johnson , University of California, San Diego, San Diego, CA, 92093, USA

Chris Fallin , Fastly, San Francisco, CA, 94107, USA

Anjo Vahldiek-Oberwagner , Intel Labs, Berlin, Berlin, 14167, Germany

Michael LeMay , Intel Labs, Hillsboro, OR, 97124, USA

Ravi Sahita , Rivos, Santa Clara, CA, 95054, USA

Dean Tullsen  and Deian Stefan , University of California, San Diego, San Diego, CA, 92093, USA

Hardware-assisted fault isolation (HFI) is a minimal extension to current processors that supports secure, flexible, and efficient in-process isolation. HFI addresses the limitations of existing software-based fault isolation (SFI) systems, including runtime overheads, limited scalability, vulnerability to Spectre attacks, and limited compatibility with existing code and binaries. HFI can be seamlessly integrated into existing SFI systems (e.g., WebAssembly) or directly sandbox unmodified native binaries. To ease adoption, HFI relies only on incremental changes to existing high-performance processors.

Isolation primitives have a critical impact on how we organize software systems for security, reliability, flexibility, and performance.

Today, familiar isolation primitives, such as virtual machines (VMs), processes, and containers, all rely on a common set of underlying hardware primitives—i.e., page tables and protection rings. However, in recent years, software developers have been reaching for another technique—software-based fault isolation (SFI)^{1,2}—to overcome the limitations of existing primitives.

SFI enforces isolation in process; i.e., it allows the creation of isolated “sandboxes” in a process (or kernel) address space through a combination of compiler

instrumentation and virtual memory tricks. SFI is attractive, as it avoids the limitations of existing hardware-based isolation primitives,^{1,3} such as high context-switch overheads, slow cold starts, and scaling limitations.

With the widespread adoption of WebAssembly (Wasm), SFI has become ubiquitous. In the browser, Wasm powers applications used by billions of people daily (e.g., Zoom, Figma, and Photoshop). Beyond the browser, developers are using Wasm to support use cases that existing isolation primitives cannot—thanks to the novel capabilities SFI offers.

To start, SFI supports context switches for roughly the same cost as a function call—orders of magnitude cheaper than existing isolation primitives. This enables isolation that is tightly integrated with high-performance applications. For example, Wasm is used to provide extensibility in microservice data planes (Istio), real-time

0272-1732 © 2024 IEEE

Digital Object Identifier 10.1109/MM.2024.3422977

Date of publication 9 July 2024; date of current version 14 August 2024.

databases (SingleStore), data streaming platforms (RedPanda), and software-as-a-service applications (Shopify). This tight integration also enables sandboxing to be retrofitted into existing applications—to contain the impact of memory safety vulnerabilities—without costly re-engineering. For example, Firefox⁴ now relies on Wasm to sandbox the third-party C libraries it relies on for audio decoding, XML parsing, spell-checking, and many other tasks.

SFI also enables fast context creation. For example, production function-as-a-service (FaaS) systems can spin up a new Wasm instance in 5 μ s⁵ instead of the tens to hundreds of milliseconds it takes to spin up a container or virtual machine. This has enabled a new class of high-concurrency, low-latency edge computing platforms from Fastly, Cloudflare, Akamai, etc.

While SFI's unique capabilities have opened the door to many new use cases, it also has a variety of limitations, including performance overheads, scaling limitations, limited compatibility with existing code and binaries, and Spectre safety. These are the natural result of the fact that SFI is a hack—an ad hoc software technique trying to bridge the gap between past models of hardware protection and the current needs of software systems.

To overcome these limitations, we developed *hardware-assisted fault isolation (HFI)*—a minimal set of non-intrusive architecture extensions that bring first-class support for in-process isolation to modern processors.

HFI offers primitives that systematically eliminate typical software (and hardware) isolation overheads by design: it imposes near-zero overhead on sandbox setup, teardown, and resizing; it can support an arbitrary number of concurrent sandboxes; it offers context-switch overheads on the same order as a function call; it can share memory between sandboxes at near-zero cost; and it provides flexible low-cost mitigations for Spectre and near-zero-cost system call interposition (for native binaries).

HFI provides first-class assistance for Wasm and similar systems by offering secure, scalable, and efficient hardware primitives that can be used as a drop-in replacement for SFI; it also provides first-class support for backward-compatible in-process isolation, allowing it to sandbox existing native binaries and dynamically generated code. HFI achieves this with minimal additional hardware and minor changes to the control and data paths of existing processors, making it easy to adopt.

SFI AND ITS LIMITATIONS

SFI is the only widely deployed technology for fine-grain, in-process isolation today. Its unique capabilities

are enabling the growth of in-process isolation in many domains, in particular, in the form of Wasm. As discussed, SFI achieves this through an ad hoc software technique that allows it to skirt some of the key limitations³ of existing hardware isolation—however, this ad hoc approach brings limitations of its own.

Conceptually, SFI enforces isolation by interposing on all memory access through compiler instrumentation. To implement this, memory is viewed as a set of contiguous memory regions with a base (starting address) and a size—a compiler adds the base address of a sandbox to the operand of any memory operation, e.g., the load, and then checks that the result is in bounds, rather like a poor man's version of segmentation.

Naively, we could imagine implementing this by adding explicit bounds checks to each memory load/store and instruction fetch; however, this can easily slow down code by a factor of $2\times$ or more.² Instead, Wasm and other modern SFI systems (e.g., Native Client) rely on a faster technique that relies on the memory management unit (MMU) to enforce bounds implicitly using a system of large address spaces and guard regions.

Wasm runtimes implement this by allocating a 4-GiB address space (called a *linear memory* in Wasm), followed by a 4-GiB guard region (unmapped address space) for each sandbox. When accessing linear memory, each Wasm load/store instruction takes two 32-bit unsigned operands that are added, resulting in a 33-bit address—and the result is added to a 64-bit *base address*—the starting address of a linear memory. By construction, any 33-bit unsigned offset plus a base address will be within 8 GB of the base; thus, access beyond the first 32-bit (4-GB) address space will trap. To isolate control flow, Wasm also relies on software control flow integrity.⁶

Despite this clever design, Wasm still has many limitations, some fundamental to SFI and others specific to Wasm's design:

- ▶ *32-bit address spaces*: Guard-region-based SFI only works for 32-bit address spaces on 64-bit architectures—supporting larger Wasm sandboxes or smaller processors requires falling back to explicit bounds checks, with the high overheads (up to $2\times$) this implies.
- ▶ *Performance overheads*: Even with these tricks, Wasm/SFI can still easily impose performance overheads of 40%—sometimes less and sometimes a lot more.⁴
- ▶ *Spectre*: Wasm cannot protect itself against Spectre attacks without performance penalties—to wit, software-based mitigations add an additional 62%–100% of overhead.⁷

- › *Compatibility*: A compiler must explicitly target Wasm to use it. Thus, assembly language, platform-specific compiler intrinsics, dynamically generated code, and existing binaries (e.g., precompiled libraries) are not supported.
- › *Scaling (virtual memory consumption)*: As previously noted, every Wasm instance consumes 8 GiB of virtual address space, even if it only uses only a few hundred megabytes or fewer, as most serverless edge workloads do today.

This limits scaling, as virtual address space is finite—typical x86-64 CPUs provide 2^{47} bytes (128 TiB) worth of user-accessible virtual address space.^a Thus, at 8 GiB (2^{33} bytes) per instance, we can run, at most, 16,000 (2^{14}) instances concurrently per process. High-performance edge computing platforms are already running up against this limit today. These systems spin up a new instance in microseconds for every incoming network request, and requests often block for input-output; thus, massive concurrency is the norm.

At present, their only recourse is to spin up more processes. Unfortunately, this leads to load imbalances and expensive context-switch overheads as processes contend for physical cores. Also, applications that use FaaS platforms do not always consist of just one function; they can be multiple functions that want to communicate (function chaining). In a single address space, this communication is as fast as a function call; however, this is easily 1000–10,000× slower across process boundaries due to use of inter-process communication (IPC).

The main reason FaaS providers use Wasm is to avoid these overheads in the first place. FaaS providers would rather schedule more instances in fewer processes—ideally, one. If used efficiently, 128 TiB can support quite a lot of serverless instances.

After struggling against these limitations for years in the context of building and deploying production library sandboxing and edge computing systems, we began to explore what would be possible with some additional help from hardware. Our resulting hardware extension (HFI) offers all of the benefits of SFI—and more—while eliminating these limitations.

HFI Overview

HFI is an instruction set architecture (ISA) extension for modern processors that offers secure, flexible, and efficient in-process isolation with minimal additional hardware complexity.

^aIntel supports 52/57-bit address spaces in certain high-end server CPUs.

Several key design choices enable these unique properties:

- › HFI does not rely on the MMU for in-process isolation—instead, sandboxing is enforced via a new mechanism called *regions*; regions enable coarse-grain isolation (e.g., heaps) and fine-grain sharing (e.g., objects) within a process's address space.
- › HFI is fully available in the user space; thus, there are no overheads from ring transitions or system calls when changing memory restrictions or entering and leaving a sandbox.
- › HFI only keeps the on-chip state for the currently executing sandbox; thus, it can scale to an arbitrary number of concurrent sandboxes. In contrast, many other systems hit a hard limit as they keep the on-chip state for all active sandboxes.^{1,8}

Functionally, HFI allows the creation of one or more in-process sandboxes. These can be either *hybrid* sandboxes, which integrate HFI's primitives into an existing SFI system, like Wasm, for increased performance, scalability, security, etc., or they can be *native* sandboxes, which can confine arbitrary untrusted binaries without any additional software support—thus shedding some of the compatibility and complexity baggage that accompanies SFI-based technologies, such as Wasm, which rely on significant modifications to existing language toolchains.

To use HFI, a sandboxing runtime uses the `hfi_enter` instruction, resulting in sandboxing constraints (memory and control isolation) being enforced until the `hfi_exit` instruction is invoked, at which time HFI returns control back to the runtime. The runtime is responsible for saving and restoring the sandbox's context (e.g., general-purpose registers) and can multiplex many sandboxes across cores, scheduling them as it sees fit.

A sandbox's semantics are dictated by a set of per-core HFI registers, including 1) region registers that grant access to memory; 2) a register with the sandbox exit handler—where system calls and sandbox exits are redirected—supporting full control of privileged instructions and control flow; and 3) a register with sandbox option flags, e.g., whether the sandbox is a hybrid or native sandbox.

Efficient Memory Access Control With Regions

HFI's memory and control isolation is configured using a finite set of regions. Regions in HFI are *base* and *bound* pairs that identify a range of memory that can be accessed (e.g., stack, heap, and code), and permissions (read, write, and execute) that apply to that range.

These checks are located in the heart of the processor's control and data path, where every additional gate matters. To optimize this, HFI utilizes multiple types of *specialized regions*.

Concretely, an approach that uses two 64-bit comparators (per region) would be the obvious design choice for ensuring that memory accesses are restricted to the configured regions; however, these are large circuits that would add unacceptable delay and power consumption to a processor's critical path. Instead, HFI offers multiple region types, each of which reduces hardware complexity by being specialized to a particular task. Next, we describe these different region types and how they are used.

Implicit Data Regions

HFI uses implicit regions to apply checks to every memory access and grant access on a first-match basis. For example, if sandboxed code executes an instruction—*load address X into register Y*—HFI will check if *any* region register has a range that includes X in parallel and then apply the permissions from the first matching region. If the first match has read permission, the operation will proceed; otherwise, HFI will trap.

Implicit regions are essential for situations when every memory operation in an application must be checked. In exchange for this power, they assume some constraints on a region's size and alignment; in particular, implicit regions require a region's address to occupy an aligned location and have a size that is a power of two.

This restriction typically requires developers to only make slight modifications to how memory allocation occurs; however, it allows hardware to implement the region checks as simple, fast masking operations. Again, for hardware simplicity, implicit regions are further specialized into code and data regions. In total, HFI provides four implicit data regions and two implicit code regions for a sandbox.

Explicit Data Regions

HFI also provides four explicit data regions that trade the generality of implicit regions for precision. Specifically, all loads and stores to explicit regions are region relative; i.e., they name the region they apply to using the new `hmov1,2,3,4` instruction, a variant of the x86-64 `mov` instruction, with the region number encoded in the instruction.

In exchange for this constraint, the size and alignment constraints of implicit regions are relaxed, allowing more precise control over memory layout. There are two types of explicit regions: small, which support byte granular size alignment and are limited to a

maximum of 4G B, and large, which are 64,000-aligned and can be any multiple of 64,000 in size. This added specialization again supports simpler hardware; explicit regions can be supported with just a single 32-bit comparator. In total, HFI provides four explicit data regions for a sandbox.

Using Regions

By default, a sandbox cannot access any memory—region registers must be set before sandbox entry (`hfi_enter`) to grant memory access to a sandbox using a handful of new instructions: `hfi_set_region`, `hfi_get_region`, and `hfi_clear_region`.

Using these instructions to configure HFI's explicit and implicit regions meets a diverse set of needs. HFI's implicit regions are ideal for sandboxing unmodified native binaries, where large regions of memory (e.g., the stack, heap, and code) are not as particular about their size or alignment. The added precision of explicit regions makes them well suited to supporting heaps for SFI systems like Wasm—where heaps grow in fixed size increments, and all operations are already explicitly relative to a region of memory—as well as for fine-grain, zero-copy sharing of objects between two different sandboxes.

System Call Interposition

System call interposition is another important feature of HFI. Architecture support makes this very efficient (`syscall` instructions are optionally converted to jumps to a specified location in the processor's decode stage), simple, and accessible at the user level. With this, HFI can sandbox unmodified native binaries while ensuring that sandboxing is not bypassed or disabled.⁹ System call interposition is used by setting a flag in the HFI option register prior to sandbox entry as well as providing an exit handler—the location where system calls are redirected to.

HFI offers a variety of other features to enable secure and efficient sandboxing—for example, to mitigate certain classes of Spectre attacks. Please consult our full paper for more details.¹⁰

EVALUATION

We implemented HFI in the Gem5 simulator for detailed performance analysis. We also used compiler-based emulation (that emits instructions, such as `cpuid`, `lfence`, etc., to insert appropriate slowdowns) to efficiently evaluate longer running workloads and validated that our emulator has similar performance characteristics as our Gem5 simulator. Please consult our full paper for more details.¹⁰

Performance Evaluation

We evaluated the performance and scalability of HFI in three use cases: library sandboxing in a browser, an FaaS workload, and native sandboxing in a server workload (NGINX).

- › *Wasm sandboxing in Firefox:* To understand the end-to-end performance impact of HFI, we evaluate the performance of HFI in sandboxing font rendering in Firefox. This benchmark stresses HFI's transitions as the code rapidly jumps from the sandboxed code to Firefox code—one jump per glyph/letter. The benchmark reflows the text on a page 10 times via the sandboxed `libgraphite` and takes 1823 ms when using Wasm with guard pages, 2022 ms when using Wasm with bounds checking, and 1677 ms when Wasm is powered by HFI—a 17% and 8% speedup, respectively.
- › *HFI and scalability:* As HFI eliminates the need for guard regions, it can support far more concurrent instances in a process's address space, a highly desirable trait in high-scale serverless settings. To exercise this, we modified Wasmtime, a Wasm runtime popular in server-side settings, to use HFI instead of guard regions, and spun up as many 1-GiB sandboxes as it would support. As expected, it could create up to 256,000 1-GiB sandboxes in a single process, making full use of the process's address space.^b
- › *HFI for sandboxing native binaries:* We used HFI to sandbox a native binary and compared it to ERIM,⁸ a state-of-the-art system for sandboxing. ERIM uses Intel memory protection keys (MPKs) for sandboxing memory accesses and Linux's `seccomp-BPF` to filter system calls.

We first compare the overhead of ERIM's `seccomp-BPF` system call filter to HFI's microarchitectural support for the same. For this, we ran a custom `syscall` benchmark that opens a file, reads it, and closes it 100,000 times. We found that using the `seccomp-BPF` version imposes an overhead of 2.1% over HFI.

Next, we compare the throughput of the NGINX Web server delivering content when protecting session keys with HFI and MPK, respectively (similar to ERIM⁸). HFI's overhead ranges from 2.9% to 6.1%, compared to MPK's range from 1.9% to 5.3%—a slight increase due to the cost of HFI's metadata

^bWasmtime normally supports up to 21,504 sandboxes. It is able to slightly exceed the 16,000 limit through a sophisticated combination of guard regions and bounds checks.

manipulation as well as serialization on `hfi_enter` and `hfi_exit` for Spectre protection.

Spectre Evaluation

We use proof-of-concept attacks from the Transient-Fail (the in-place Spectre pattern history table attack) and Google Safe-Side (the in-place Spectre branch target buffer attack) test suites to ensure that HFI is resistant to Spectre attacks. HFI prevents both of these attacks, ensuring that sandboxed code cannot speculatively access secret data.

We compared HFI's Spectre protection overheads to the performance of Swivel,⁷ the fastest known compiler-based approach to mitigating Spectre in Wasm. We evaluated this by measuring the costs of Spectre protection on several common FaaS Wasm workloads running in the Rocket Web server. We compiled our Wasm workloads with a Lucet Wasm compiler without Spectre protections, with Lucet plus Swivel-SFI protections, and with Lucet plus HFI using the native sandbox. Finally, we also recorded the binary sizes of all three builds.

Table 1 shows that HFI guards against Spectre with a very low drop-in tail latency and no noticeable binary bloat, while Swivel incurs noticeable overheads for the same. In fact, the only overheads imposed by native sandbox HFI are due to region construction and sandbox state transitions (two per connection), and these costs are amortized by the cost of the workload.

HFI ACROSS DIFFERENT ARCHITECTURES

Our initial HFI prototype¹⁰ was implemented in the context of the x86-64 architecture (see Figure 1). However, its general design and key principles, such as the region specialization to minimize circuit area, avoiding extra pipeline stages, restricting the saved state to only the currently executing sandbox, etc., are equally applicable to other ISAs, with some modifications to address the particular constraints of other ISAs.

For example, on x86s, we added four new `hmov1`, `hmov2`, `hmov3`, and `hmov4` instruction prefixes that can be added to any instruction accessing memory, indicating it should operate on one of the explicit data regions. However, when adapting HFI to reduced-instruction-set computing V (RISC-V), we found that this approach was incompatible with RISC-V, as it does not support instruction prefixes. Rather, each RISC-V instruction accessing memory—`load`, `store`, `load_fp`, and `store_fp`—would each need four new opcodes to create versions of these instructions accessing explicit regions. This would put unacceptable pressure on RISC-V's already very constrained

TABLE 1. Impact of HFI Spectre protection on tail latency.*

HFI Protection	XML to JSON				Image classification				Check SHA-256				Templated HTML			
	Average Latency	Tail Latency	Throughput (requests/s)	Bin size	Average Latency	Tail Latency	Throughput (requests/s)	Bin size	Average Latency	Tail Latency	Throughput (requests/s)	Bin size	Average Latency	Tail Latency	Throughput (requests/s)	Bin size
Lucet (unsafe)	421 ms	466 ms	231	3.5 MiB	12.2 s	14.7 s	1.62	34.3 MiB	589 ms	667 ms	161	3.9 MiB	45.6 ms	61.8 ms	2190	3.6 MiB
Lucet + HFI	431 ms	480 ms	227	3.5 MiB	12.2 s	14.7 s	1.62	34.3 MiB	602 ms	647 ms	165	3.9 MiB	45.7 ms	61.2 ms	2180	3.6 MiB
Lucet + Swivel	559 ms	616 ms	174	4.1 MiB	11.5 s	12.8 s	1.72	34.5 MiB	645 ms	709 ms	150	4.6 MiB	78.9 ms	97.9 ms	1260	4.2 MiB

*We compared HFI and Swivel—the fastest software-based Spectre mitigation—on several Wasm FaaS workloads. Swivel increased tail latency by 9%–42%. HFI increased tail latency by 0%–2%.

opcode space. Instead, for RISC-V, we add a new instruction `h_set_region` to set the active explicit region and limit `hload`, `hstore`, etc. to only modify the active region. While this leads to a less compact instruction encoding in the case where multiple regions are being accessed concurrently, it alleviates this pressure on the opcode space and also keeps the ISA simpler, keeping with the RISC design philosophy.

Another aspect of our design that differs across ISAs (and processors more generally) is how it is instantiated at the microarchitectural level. In particular, HFI is informed by a “living off the land” design philosophy; i.e., we aimed to leverage existing hardware already present in the pipeline as much as possible.

Again, going back to our x86 example, HFI’s explicit regions can be implemented by repurposing the circuits for 32-bit segmented memory and 64-bit segment relative addressing. Together, these features are implemented using a 32-bit comparator and a 64-bit adder, respectively—precisely the circuits needed for HFI’s `hmov` instruction—and are both circuits that would be otherwise unused for this instruction. On other architectures, the circuits available for reuse will differ. For example, future RISC-V standards will incorporate support for pointer masking, which could be reused to support implicit regions, while architectures such as the ARM Cortex-M feature hardware for data watchpoint and trace range checks, which could be repurposed to implement data region bounds checks.

APPLICATIONS AND IMPLICATIONS

By providing first-class support for in-process isolation, HFI opens the door to greater capabilities and new applications in a variety of domains. Here, we explore just a few.

Easier Library Sandboxing

Billions of lines of vulnerable C/C++ underlie today’s software ecosystem.¹¹ A cursory inspection of any open source ecosystem of a “safe” language, such as JavaScript, Java, Rust, Python, etc., reveals thousands of packages that are essentially wrappers around existing C/C++ libraries.

Sandboxing is a potentially powerful tool for mitigating memory safety vulnerabilities in these libraries as well as other defending against other attack vectors, e.g., supply chain attacks. Unfortunately, sandboxing is rarely utilized because of the high cost of context switches (and, consequently, IPC), and sandbox creation with process sandboxing constrains the generality

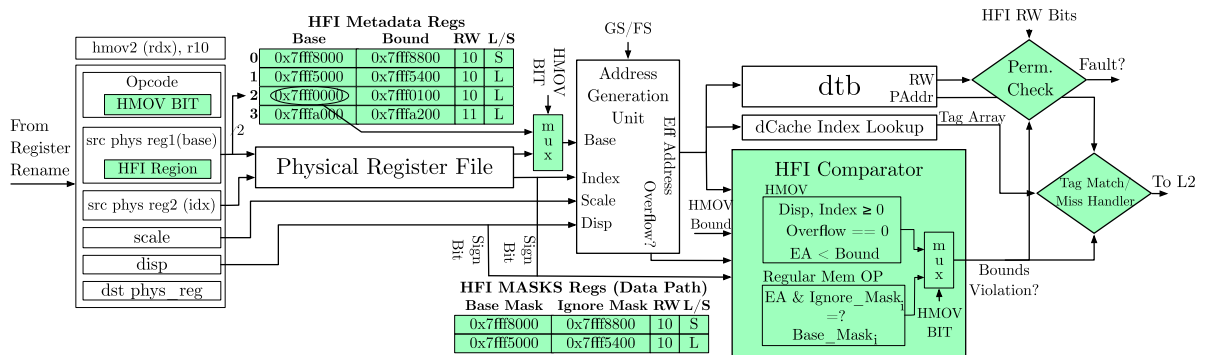


FIGURE 1. HFI impact on the x86 data pipeline. The x86 data pipeline is shown with added any HFI components in green. HFI adds no overhead to the data path—no new pipeline stages are added—and all new operations take place in parallel with the data translation lookaside buffer lookup or instruction decode stages. src: source; phys reg: physical register; idx: index; disp: displacement; dst: destination; RW: read-write; L/S: large or small; Eff: effective; Regs: registers; Perm: permission; PAddr: physical address; dCache: data cache; EA: effective address; Mem: memory; OP: operation.

of this technique and often requires significant re-engineering of applications to use it.

While SFI toolchains, such as in Wasm, have opened the door to retrofitting sandboxing⁴—i.e., sandboxing existing libraries in place without expensive re-engineering—Wasm’s many limitations still restrict the usefulness of this technique. In practice, we have seen that many libraries have handwritten assembly code; dynamically generated code; single-instruction, multiple data code that was poorly supported/unsupported; etc.

HFI renders these practical obstacles moot, eliminating many of the performance and engineering overheads that have limited the usefulness of this technique.¹¹

Secure, Compatible, and Efficient Serverless

Currently, serverless infrastructure exists in two radically different forms. On one hand, data center serverless architectures, like AWS Lambda, Azure Functions, Google Cloudfunctions, etc., rely on heavy weight isolation mechanisms, such as VMs and containers, to isolate individual instances. While this offers excellent compatibility with existing languages and libraries and a strong foundation for multitenant isolation, start-up overheads and latency are often very high (in the range of tens to hundreds of milliseconds).

On the other hand, edge computing platforms, such as Cloudflare Workers and Fastly’s Compute@Edge, can offer orders-of-magnitude improvements in throughput and latency by relying on software-based isolation

(utilizing Wasm and JavaScript)—but, in turn, sacrifice compatibility. For example, Java and .NET both have mature just-in-time compilers (JITs) and runtimes, which have been optimized and hardened over decades, that cannot be used in these environments.

HFI offers the potential to eliminate this disconnect, opening the door to future serverless platforms that could offer edge-level scaling, efficiency, and responsiveness while retaining the security and compatibility of existing data center serverless platforms.

Secure and Efficient Browser JIT Hardening

One of the hard challenges in browser security today is memory safety bugs in JITs. Modern JavaScript JITs are marvels of engineering, offering astonishingly good performance for a dynamic language like JavaScript. Unfortunately, the combination of complex optimizations required to achieve this and the complexity of JavaScript itself result in a very large attack surface.

Further, some of the best techniques for mitigating these bugs in most applications are not applicable in a JIT context. For example, as bugs are in code optimization and generation, safe languages, such as Rust, are often of little benefit since the bugs are in machine code generated by the JIT. Normal write-XOR-execute ($W \oplus X$) restrictions that prevent code injection in normal applications also break in a JIT environment. JIT compilers often have to modify generated executable code; requiring JIT compilers to make expensive system calls to change permissions of memory prior to such changes imposes both a direct slowdown of the

JIT compiler and a systemwide slowdown due to frequent modifications of the translation look-aside buffer (TLB).

One compelling approach to limiting the impact of these bugs is sandboxing JIT's code, for example, Chrome's Ubercage.¹² Unfortunately, such efforts currently have to trade off security for performance; for example, Ubercage it does not incorporate software protections for control flow hijacking, let alone Spectre, offering multiple potential ways to bypass it.

HFI can advance the state of the art in JIT hardening in a number of ways.

First, HFI offers per-thread memory permissions that can be set without system calls; thus, a JIT compiler running on one thread can enforce a write permission on memory, while code compiled by a JIT executing on a separate thread can enforce an executable permission, restoring the benefits of write-XOR-execute in a JIT setting.

Next, HFI can eliminate the limitations of JIT sandboxes, such as Ubercage, by both eliminating the SFI overheads it suffers from today and allowing it to efficiently add new restrictions, such as efficiently enforcing Control flow and Spectre isolation—thus eliminating the paths to bypass sandboxing that currently exist.¹²

PATH TO ADOPTION

HFI offers a unique solution for hardware-assisted isolation that balances novel acceleration capabilities with practical adoption.

Previous approaches to sandboxing that repurpose existing hardware,⁸ such as MPKs, consistently fall short with respect to scaling, performance, security,⁹ etc., while ambitious novel architectures, such as CHERI, impose high costs in terms of hardware complexity and modifying existing software stacks, posing significant barriers to adoption. HFI illustrates a middle path between these two extremes, offering hardware support for in-process isolation that is both capable and easy to adopt.

In recent years, SFI, in the form of Wasm, has seen broad adoption on the Web, on edge computing platforms, for the sandboxing of buggy C libraries,⁴ and for providing safe extensibility in numerous applications. Beyond its use in Wasm, SFI is also used by Chrome to isolate its JavaScript JIT and by the Linux kernel (eBPF) to provide extensibility, performance analysis, and security monitoring.

The widespread use of SFI allows HFI to avoid the classic chicken-and-egg problem faced by new hardware security features—there is already a huge base of

software using SFI that can (with minor changes) immediately benefit from the added security, performance, scalability, etc. that HFI offers—while also opening the door to use cases that are simply not possible with SFI today.

To further ease adoption, HFI was refined through numerous iterations of feedback from architects at multiple hardware companies including Intel. As a result, HFI offers precisely what is needed to support in-process isolation with a minimal hardware bill of materials.

Finally, HFI requires very little support from the operating system (OS) (on the order of tens of lines of code), again informed by the reality that getting OS kernel support is a significant barrier to rolling out new hardware features.

Our work with HFI is just beginning. Engineers who develop Chrome and Firefox are excited about the benefits that HFI can bring to their hardening efforts¹²; engineers at edge computing companies are excited about its potential to improve scalability, security, and performance on their platforms; and processor architects are excited about enabling these benefits with minimal hardware costs.

ACKNOWLEDGMENTS

The authors thank Dan Gohman and Luke Wagner from Fastly and the architects from Intel for their insightful discussions as well as the anonymous reviewers and shepherd for their valuable comments for improving the quality of the original published paper. This work was supported in part by a Sloan Research Fellowship; by the National Science Foundation under Grant CNS-2155235, Grant CNS-2120642, and Grant CAREER CNS-2048262; by gifts from Intel, Google, Cisco, and Mozilla; and by the Defense Advanced Research Projects Agency HARDEN under NIWC Contract N66001-23-9-4004.

REFERENCES

1. G. Tan, "Principles and implementation techniques of software-based fault isolation," in *Foundations and Trends in Privacy and Security*. Camp Hill, PA, USA: The Pennsylvania State University, 2017, pp. 137–198.
2. R. Wahbe, S. Lucco, T. E. Anderson, and S. L. Graham, "Efficient software-based fault isolation," in *Proc. 14th ACM Symp. Operating Syst. Princ. (SOSP)*, 1993, pp. 203–216, doi: [10.1145/168619.168635](https://doi.org/10.1145/168619.168635).
3. R. M. Norton, "Hardware support for compartmentalisation," Univ. of Cambridge,

- Cambridge, England, Tech. Rep., UCAM-CL-TR-887, May 2016.
4. S. Narayan et al., "Retrofitting fine grain isolation in the Firefox renderer," in *Proc. 29th USENIX Secur. Symp.*, 2020, pp. 699–716.
 5. L. Clark. "Wasmtime reaches 1.0: Fast, safe and production ready!" Bytecode Alliance. Accessed: Mar. 1, 2024. [Online]. Available: <https://bytecodealliance.org/articles/wasmtime-1-0-fast-safe-and-production-ready>
 6. A. Haas et al., "Bringing the web up to speed with WebAssembly," in *Proc. 38th ACM SIGPLAN Conf. Program. Lang. Design Implementation (PLDI)*, 2017, pp. 185–200, doi: [10.1145/3062341.3062363](https://doi.org/10.1145/3062341.3062363).
 7. S. Narayan et al., "Swivel: Hardening WebAssembly against Spectre," in *Proc. USENIX Secur. Symp.*, 2021, pp. 1433–1450.
 8. A. Vahldiek-Oberwagner, E. Elnikety, N. O. Duarte, M. Sammler, P. Druschel, and D. Garg, "ERIM: Secure, efficient in-process isolation with protection keys (MPK)," in *Proc. 28th USENIX Secur. Symp.*, 2019, pp. 1221–1238.
 9. R. J. Connor, T. McDaniel, J. M. Smith, and M. Schuchard, "PKU pitfalls: Attacks on PKU-based memory isolation systems," in *Proc. 29th USENIX Secur. Symp.*, 2020, pp. 1409–1426.
 10. S. Narayan et al., "Going beyond the limits of SFI: Flexible and secure hardware-assisted in-process isolation with HFI," in *Proc. 28th ACM Int. Conf. Archit. Support Program. Lang. Operating Syst. (ASPLOS)*, 2023, pp. 266–281, doi: [10.1145/3582016.3582023](https://doi.org/10.1145/3582016.3582023).
 11. B. Lord. "The urgent need for memory safety in software products." CISA (.gov). Accessed: Mar. 1, 2024. [Online]. Available: <https://www.cisa.gov/news-events/news/urgent-need-memory-safety-software-products>
 12. S. Groß (Saelo). "V8 sandbox - hardware support." Accessed: Mar. 1, 2024. [Online]. Available: <https://docs.google.com/document/d/12MsaG6BYRB-jQWNkZiuM3bY8X2B2cAsCMLL dgErvK4c/>

SHRAVAN NARAYAN is an assistant professor at University of Texas at Austin, Austin, TX, 78712, USA. His research interests include building secure systems, program verification, and hardware-based security. Narayan received his doctorate degree from the University of California, San Diego. Contact him at shr@cs.utexas.edu.

TAL GARFINKEL is a research scientist at the University of California, San Diego, San Diego, CA, 92093, USA. His

research interests include hardware/software systems from programming languages to computer architecture, focusing on security, reliability, and performance. Garfinkel received his doctorate degree from Stanford University. Contact him at talg@cs.ucsd.edu.

MOHAMMADKAZEM TARAM is an assistant professor at Purdue University, West Lafayette, IN, 47907, USA. His research interests include computer architecture and computer security, focusing on microarchitectural attacks, high-performance mitigations, and architecture support for security and privacy. Taram received his doctorate degree from the University of California, San Diego. Contact him at kazem@purdue.edu.

JOEY RUDEK is a Ph.D. student at University of California, San Diego, San Diego, CA, 92093, USA. His research interests include secure hardware design and the security implications of novel hardware optimizations. Contact him at jrudek@ucsd.edu.

DANIEL MOGHIMI is a senior research scientist at Google, Mountain View, CA, 94043, USA. His research interests include computer and hardware security, including microarchitectural vulnerabilities, side-channel cryptanalysis, and security architecture. Moghimi received his doctorate degree from Worcester Polytechnic Institute. Contact him at danielmm@google.com.

EVAN JOHNSON is a Ph.D. student at University of California, San Diego, San Diego, CA, 92093, USA. His research focuses on interests include software security and software correctness, particularly securing and verifying software sandboxing systems. Contact him at e5johnso@ucsd.edu.

CHRIS FALLIN is a principal software engineer on the WebAssembly team at Fastly, San Francisco, CA, 94107, USA. His research interests include compilers, runtimes, and system software with an emphasis on correctness. Fallin received his doctorate degree from Carnegie Mellon University. Contact him at chris@cfallin.org.

ANJO VAHLDIK-OBERWAGNER is a research scientist with the Datacenter Security Group, Intel Labs, Berlin, Berlin, 14167, Germany. His research interests include analyzing and building secure software and hardware computing systems, focusing on the usability of in-process isolation in data center workloads. Vahldiek-Oberwagner received his doctorate

degree from the Max Planck Institute for Software Systems and Saarland University. Contact him at anjo.lucas.vahldiek-oberwagner@intel.com.

MICHAEL LEMAY is a senior staff research scientist in Intel Labs, Hillsboro, OR, 97124, USA. His research interests include memory management architectures for security. Lemay received his doctorate degree from the University of Illinois at Urbana-Champaign. Contact him at michael.lemay@intel.com.

RAVI SAHITA is a principal member of the technical staff at Rivos Inc., Santa Clara, CA, 95054, USA, where he leads the platform security architecture and ISAs, and vice chair of the Security Horizontal Committee at RISC-V International. His research interests include exploit prevention instruction set architectures and confidential computing. Sahita received

his M.S. degree in computer science from Iowa State University. Contact him at ravi@rivosinc.com.

DEAN TULLSEN is a professor at the University of California, San Diego, San Diego, CA, 92093, USA. His research interests include computer architecture and architectural security. Tullsen received his doctorate degree from the University of Washington. He is a fellow of the Association for Computing Machinery and a Fellow of IEEE. Contact him at tullsen@ucsd.edu.

DEIAN STEFAN is an associate professor with the University of California, San Diego, San Diego, CA, 92093, USA. His research interests include principled and practical secure systems using techniques that span security, programming languages, and systems. Stefan received his doctorate degree from Stanford University. Contact him at deian@cs.ucsd.edu.



IEEE Security & Privacy magazine provides articles with both a practical and research bent by the top thinkers in the field.

- stay current on the latest security tools and theories and gain invaluable practical and research knowledge,
- learn more about the latest techniques and cutting-edge technology, and
- discover case studies, tutorials, columns, and in-depth interviews and podcasts for the information security industry.



computer.org/security



IEEE
COMPUTER
SOCIETY

