

DAWG: A Defense Against Cache Timing Attacks in Speculative Execution Processors*

Vladimir Kiriansky*, Iliia Lebedev*, Saman Amarasinghe*, Srinivas Devadas*, Joel Emer‡
*MIT CSAIL, ‡NVIDIA / MIT CSAIL
{vlk, ilebedev, saman, devadas, emer}@csail.mit.edu

ABSTRACT

Software side channel attacks have become a serious concern with the recent rash of attacks on speculative processor architectures. Most attacks that have been demonstrated exploit the cache tag state as their exfiltration channel. While many existing defense mechanisms that can be implemented solely in software have been proposed, these mechanisms appear to patch specific attacks, and can be circumvented. In this paper, we propose minimal modifications to hardware to defend against a broad class of attacks, including those based on speculation, with the goal of eliminating the entire attack surface associated with the cache state covert channel.

We propose DAWG, Dynamically Allocated Way Guard, a generic mechanism for secure way partitioning of set associative structures including memory caches. DAWG endows a set associative structure with a notion of protection domains to provide strong isolation. When applied to a cache, unlike existing quality of service mechanisms such as Intel’s Cache Allocation Technology (CAT), DAWG isolates hits and metadata updates across protection domains. We describe how DAWG can be implemented on a processor with minimal modifications to modern operating systems. We argue a non-interference property that is orthogonal to speculative execution and therefore argue that existing attacks such as Spectre Variant 1 and 2 will not work on a system equipped with DAWG. Finally, we evaluate the performance impact of DAWG on the cache subsystem.

1. INTRODUCTION

1.1 Context

For decades, processors have been architected for performance or power-performance. While it is generally assumed by computer architects that performance and security are orthogonal concerns, there are a slew of examples, including the recent Google Project Zero attacks [22] (and Spectre [31] and Meltdown [35]) that show that performance and security are not independent, and micro-architectural optimizations that preserve architectural correctness can affect the security of the system.

The objective of the attacker is to create some software that can steal some secret that another piece of code should have

*Student and faculty authors listed in alphabetical order.

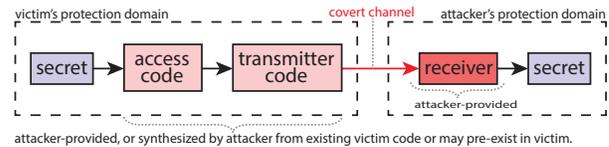


Figure 1: Attack Schema: an adversary 1) accesses a victim’s secret, 2) transmits it via a covert channel, 3) receives it in their own protection domain.

exclusive access to. The access to the secret may be made directly, e.g., by reading the value of a memory location, or indirectly, e.g., inferred from the execution flow a program takes. In either case, this is referred to as violating *isolation*, which is different from violating *integrity* (corrupting the results obtained through program execution).¹

In a well-designed system the attacker cannot architecturally observe this secret, as the secret should be confined to some protection domain that prevents other programs from observing it. Even without physical access to the machine, an attacker may observe side effects of execution, i.e., via software means alone.

The mechanism by which such observations are made are referred to as *software side channels*. Currently, the most widely explored of these channels is the ability to ascertain something about the state of a shared cache. The attacker measures the time it takes for its program to make specific references, which vary due to cache hits and misses, and infers something about others’, i.e., victims’ use of the cache. If the attacker observes a hit on an address, the address must be cached already, meaning some party had recently accessed it, and it had not yet been displaced.

When a side channel conveys a “secret” to an attacker it is referred to as a *covert communication channel*, and an attack would include code for a transmitter inside the victim’s protection domain that has access to the secret and sends it to a receiver controlled by the attacker, and outside the victim’s protection domain. This is pictorially illustrated in Figure 1.

A classic attack on RSA relied on such a scenario [10]. Specifically, existing RSA code followed a conditional execution sequence that was a function of the secret, and inad-

¹Violating isolation and obtaining a secret may result in the attacker being able to violate integrity as well, since it may now have the capability to modify memory, but in this paper we focus on the initial attack that would violate isolation.

vertently transmitted private information via the pattern of its instruction cache access. This resulted in a covert communication that let an observing adversary determine bits of the secret. In this case, the code that accessed the secret and the transmitter that conveyed the secret were pre-existing in the RSA code and the attacker provided only a receiver that was able to receive the secret via cache tag state. Recent work has shown that a broad space of viable attacks exfiltrate information via shared caches.

1.2 Generalized Attack Schema

Recently, multiple security researchers (e.g., [22, 31, 35]) have found ways for an attacker to create a *new* transmitter in the victim. In addition to the well-studied scenarios of side channel (unwitting transmitter code) and covert channel (cooperating transmitter) transmission, here an attacker is able to *create* a transmitter in the victim’s domain and/or influences the transmitter to access a chosen secret. Spectre and Meltdown have shown that code executing *speculatively* has full access to any secret.

While speculative execution is broadly defined, we focus on control flow speculation in this paper. Modern processors execute instructions *out of order*, allowing downstream instructions to execute prior to upstream instructions as long as dependencies are preserved. Most instructions on modern out-of-order processors are also *speculative*, i.e., they create checkpoints and execute along a predicted front while one or more prior conditional branches are pending resolution. When a processor does not know the future instruction stream because a conditional branch instruction is reached whose direction depends on upstream instructions that have not executed yet, the processor predicts the path that the program will follow, and speculates along that path after making a checkpoint. A prediction resolved to be correct discards a checkpoint state, while an incorrect one forces the processor to roll back to the checkpoint and resume along the correct path. Incorrectly predicted instructions are executed, for a time, but do not modify architectural state. However, micro-architectural state such as cache tag state *is* modified as a result of (incorrect) speculative execution, allowing information to leak.

Note that information leakage as a result of correct speculation requires that the transmitter already existed in the code, meaning it is a pre-existing transmitter, as shown in Figure 1. By exploiting mis-speculated execution, an attacker can exercise code paths that are normally not reachable, circumventing software invariants. One example has the attacker illegally referencing the secret from within the attacker’s domain, causing micro-architectural side effects before an exception is raised [35]. Another example has the attacker coercing branch predictor state to encourage mis-speculation along an attacker-selected code path, which implements a transmitter. There are therefore three ways of creating the transmitter:

1. Transmitter pre-exists in victim’s code, which we described in the RSA attack [10].
2. Attacker explicitly programs the transmitter. Meltdown [35] is an example of this.
3. Attacker synthesizes it out of existing code in the victim. This is exemplified by the Spectre-style attacks [22, 31].

This framework can be applied for side channels other than the cache state, describing exfiltration via branch predictor logic or TLB state, for example. Given the intensified research interest in variants of this new attack class, we also imagine that there will be new ways that transmitters can be constructed. We therefore wish to design a defense against a broad class of current and future attacks.

1.3 Our approach to defense

Defense mechanisms that can be implemented solely in software have been proposed (e.g., [12, 43]). Unfortunately, these mechanisms appear very attack specific: e.g., a compiler analysis [43] identifies some instances of code vulnerable to Spectre Variant 1; microcode updates or compiler and linker fixes reduce exposure to Spectre Variant 2 [12]. Instructions to turn off speculation in vulnerable regions have been introduced (e.g., [2]) for future compilers to use. In this paper, we target minimal modifications to hardware that defend against a broad class of side channel attacks, including those based on speculation, with the goal of eliminating the entire attack surface associated with exfiltration via cache timing.

To prevent exfiltration, we require strong isolation between protection domains. Cache partitioning is an appealing mechanism to achieve isolation. Unfortunately, set (e.g., page coloring [30, 51]) and way (e.g., Intel’s Cache Allocation Technology (CAT) [21, 23]) partitioning mechanisms available in today’s processors are either low-performing or do not provide isolation.

We propose DAWG, *Dynamically Allocated Way Guard*, a generic mechanism for secure way partitioning of set associative structures including caches. DAWG endows a set associative structure with a notion of *protection domains* to provide strong isolation. Unlike existing mechanisms such as CAT, DAWG disallows hits across protection domains. This affects hit paths and cache coherence [42], and DAWG handles these issues with minimal modification to modern operating systems, while reducing the attack surface of operating systems to a small set of annotated sections where data moves across protection domains, or where domains are resized/reallocated. Only in these handful of routines, DAWG protection is relaxed, and other defensive mechanisms such as speculation fences are applied as needed. We evaluate the performance implications of DAWG using a combination of architectural simulation and real hardware and compare to conventional and quality-of-service partitioned caches. We conclude that DAWG provides strong isolation with reasonable performance overhead.

In Section 5, we argue a non-interference property that is orthogonal to speculative execution and therefore argue that existing attacks such as Spectre Variant 1 and 2 will not work on a system equipped with DAWG.

1.4 Contributions and organization

The contributions of our paper are:

1. We motivate strong isolation of replacement metadata by demonstrating that the replacement policy can leak information (cf. Section 2.2.2) in a way-partitioned cache.
2. We design a cache way partitioning scheme, DAWG, with strong isolation properties that blocks old and new

timing attacks based on the cache state exfiltration channel (cf. Section 3). DAWG does not require invasive changes to modern operating systems, and preserves the semantics of copy-on-write resource management.

3. We analyze the security of DAWG and argue its security against recent cache timing attacks that exploit speculative execution (cf. Section 5.1).
4. We illustrate the limitations of cache partitioning for isolation by discussing a hypothetical leak framed by our attack schema (cf. Figure 1) that circumvents a partitioned cache. For completeness, we briefly describe a defense against this type of attack (cf. Section 5.3).
5. We evaluate the performance impact of DAWG in comparison to CAT [21, 23] and non-partitioned caches with a variety of workloads, detailing the overhead of DAWG’s protection domains, which limit data sharing in the system (cf. Section 6).

The paper is organized as follows. We provide background and discuss related work in Section 2. The hardware modifications implied by DAWG are presented in Section 3, and software support is detailed in Section 4. Security analysis and evaluation are the subjects of Section 5 and Section 6, respectively. Section 7 concludes.

2. BACKGROUND AND RELATED WORK

We focus on the cache state exfiltration channel and isolating domains through cache hardware modifications. We state our threat model in Section 2.1, describe relevant attacks in Section 2.2, and existing defenses in Section 2.3.

2.1 Threat model

Our focus is on blocking attacks that utilize the cache state exfiltration channel. We do not claim to protect other channels, such as L3 cache slice contention, L2 cache bank contention, network-on-chip bandwidth, DRAM bandwidth, branch data structures, TLBs or shared functional units in a physical core. In the case of the branch data structures channel, or any other set associative structure, however, we believe that a DAWG-like technique can be used to block the channel. Given the large number of attacks against simultaneous-multithreading (SMT) we assume that different protection domains are never simultaneously allowed to share physical cores. The victim process can be privileged (kernel) or an unprivileged peer. We assume an unprivileged attacker.

2.2 Attacks

A conventional set-associative cache is shown in Figure 2. In the figure, we identify the data structures, which when observed by an adversary can leak information: cache tag state and cache metadata corresponding to coherence information and replacement information.

The replacement policy selects a cache line to evict (a “victim”) when no invalid line is available on a fill.

A review [3] of recent cache replacement and protection policies includes valuable commentary on less-cited excellent papers which have left little room for improvement [26, 27, 28, 45, 46, 59].

The most common attack strategy corresponds to the adversary presetting the cache tag state to a particular value, and then after the victim runs, observing a difference in the cache

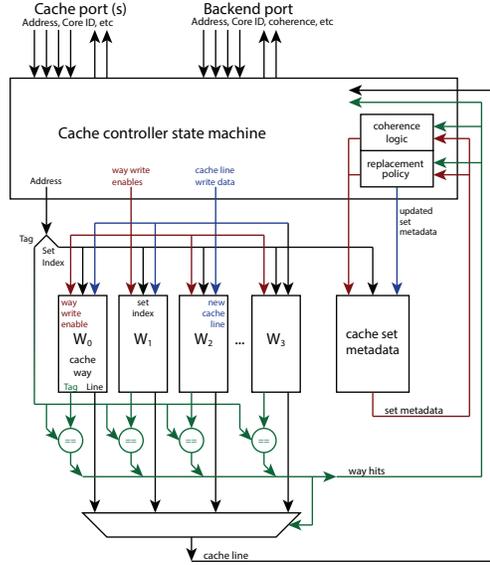


Figure 2: Conventional Set-Associative Cache.

tag state to learn something about the victim process. A less common yet viable attack strategy corresponds to observing changes in coherence [57] or replacement metadata.

2.2.1 Cache tag state based attacks

Attacks observing cache tag state are known to retrieve cryptographic keys from a growing body of cryptographic implementations: AES [8, 40], RSA [10], Diffie-Hellman [32], and elliptic-curve cryptography [9], to name a few. Cache timing attacks can be mounted by unprivileged software sharing a computer with the victim software [4]. While early attacks required access to the victim’s CPU core, more recent sophisticated attacks such as flush+reload [58] and variants of prime+probe [38] target the last-level cache (LLC), which is shared by all cores in a socket. The evict+reload attack variant of flush+reload uses cache contention rather than flushing [38]. A cache timing attack in JavaScript was demonstrated [39] to automatically exfiltrate private information upon a web page visit.

These attacks focus on various levels of the memory cache hierarchy and exploit cache lines shared between an attacker’s program and the victim process. Regardless of the specific mechanism for inspecting shared tags, the underlying concepts are the same: two entities separated by a trust boundary share the computer system’s resources, specifically sets in the memory hierarchy. Given that any given cache set is a limited resource, the entities can communicate (or “unwittingly” communicate, in the case of an attack) across the trust boundary by modulating the presence of a cache tag in a set. The receiver can detect the transmitter’s fills of tag T_A either directly, by observing whether it had fetched a shared line, or indirectly, by observing capacity misses on the receiver’s own data caused by the transmitter’s accesses, as shown in Figure 3.

In a situation where both entities share an address with which the communication happens, the receiver (or attacker)

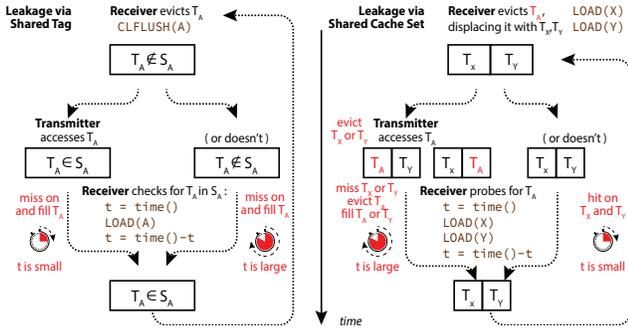


Figure 3: Leakage via a shared cache set, implemented via a shared tag T_A directly, or indirectly via $T_X, T_Y \cong T_A$.

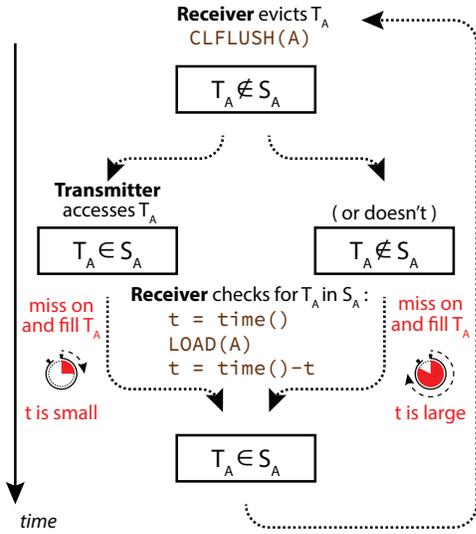


Figure 4: Covert channel via a shared tag T_A .

can both *evict* and *fill* the address in the cache, making communication straightforward. As illustrated in Figure 4, the receiver forces the cache set to a known state by flushing the address, yields for a time, and then measures the latency of loading the address (in order to detect whether the transmitter re-filled it in the cache). In a typical system, shared addresses are commonplace, due to common libraries and aggressive de-duplication by the system software to improve efficiency.

Communication via cache tag state is also possible without a shared address. Although the receiver (attacker) can neither directly load, nor directly flush a transmitter’s (victim’s) cache line, indirect communication via conflict misses is possible. Given that cache sets have limited room, a receiver can force the transmitter’s address out of the cache by filling the set with their own cache tags, and detect the transmitter’s re-fill by observing an eviction on one of the receiver’s own tags, as shown in Figure 5.

2.2.2 A cache metadata based attack

Even without shared cache lines (as is the case in a way-partitioned cache), the replacement metadata associated with each set may leak information across partitions. Most replacement policies employ a replacement state bit vector

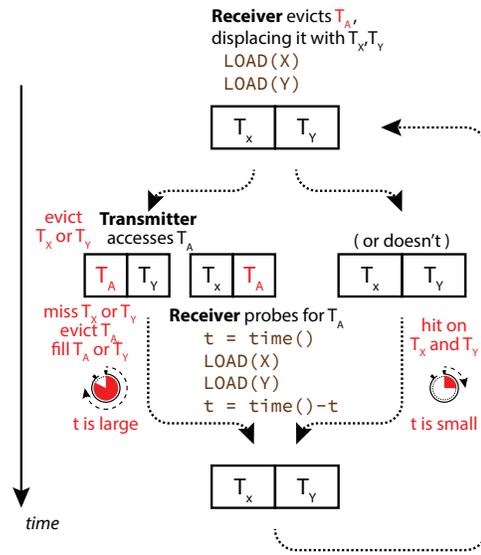


Figure 5: Covert channel via a shared cache set, as exemplified by a 2-way set-associative cache with an LRU, via tags $T_A \cong T_X \cong T_Y$.

that encodes access history to the cache set in order to track (or approximate) the ways least costly to evict in case of a miss. If the cache does not explicitly partition the replacement state metadata across protection domains, some policies may violate isolation in the cache by allowing one protection domain’s accesses to affect victim selection in another partition. Figure 6 exemplifies this with Tree-PLRU replacement (Section 3.9.2): a metadata update after an access to a small partition overwrites metadata bits used to select the victim in a larger partition. A securely way-partitioned cache must ensure that replacement metadata does not allow information flow across the cache partition(s).

This means cache attack defenses have to take into account the cache replacement policy and potentially modify the policy to ensure isolation.

2.3 Defenses

Broadly speaking, there are five classes of defenses, with each class corresponding to blocking one of the steps of the attack described in Figure 1.

1. *Prevent access to the secret.* For example, KAISER [14], which removes virtual address mappings of kernel memory when executing in user mode, is effective against Meltdown [35].
2. *Make it difficult to construct the transmitter.* For example, randomizing virtual addresses of code, flushing the Branch Table Buffer (BTB) when entering victim’s domain [47].
3. *Make it difficult to launch the transmitter.* For example, not speculatively executing through permission checks, keeping predictor state partitioned between domains, and preventing user arguments from influencing code with access to secrets. The Retpoline [53] defense against Spectre Variant 2 [12] makes it hard to launch (or construct) a transmitter via an indirect branch.

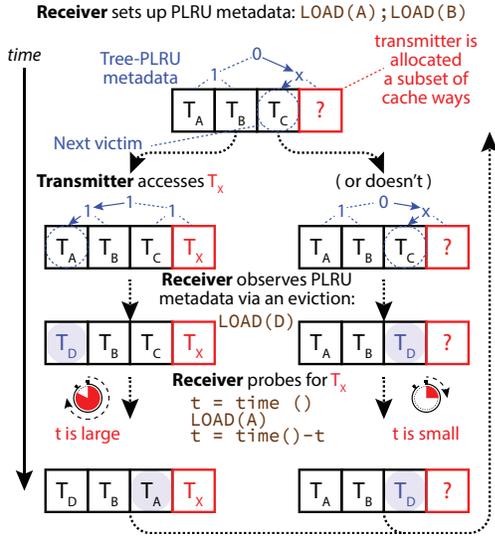


Figure 6: Covert channel via shared replacement metadata, exemplified by a 4-way set-associative cache with a Tree-PLRU policy, cache allocation boundary. Tags $T_A \cong T_B \cong T_C \cong T_D \cong T_x$.

4. *Reduce the bandwidth of side channels.* For example, removing the APIs for high resolution timestamps in JavaScript, as well as support for shared memory buffers to prevent attackers from creating timers.
5. *Close the side channels.* For example, partitioning of cache state or predictor state. This is the strategy of choice in our paper.

2.3.1 Set partitioning via page coloring

Set partitioning can provide isolation against tag and metadata attacks. It has the advantage of working with existing hardware when allocating sets at page granularity [34, 60] via page coloring [30, 51]. Linux currently does not support page coloring, since most early OS coloring was driven by the needs of low-associativity data caches [52].

Set partitioning allows communication between protection domains without destroying cache coherence. The downsides are that it requires some privileged entity, or collaboration, to move large regions of data around in memory when allocating cache sets, as set partitioning via page coloring binds cache set allocation to physical address allocation. For example, in order to give a protection domain 1/8 of the cache space, the same 12.5% of the system’s physical address space must be given to the process. In an ideal situation, the amount of allocated DRAM and the amount of allocated cache space should be decoupled.

Furthermore, cache coloring at page granularity is not straightforwardly compatible with large pages, drastically reducing the TLB reach, and therefore performance, of processes. On current processors, the index bits placement requires that small (4KB) pages are used, and coloring is not possible for large (2MB) pages. Large pages provide critical performance benefits for virtualization platforms used in the public cloud [44], and reverting to small pages would be deleterious.

2.3.2 Insecure way partitioning

Intel’s Cache Allocation Technology (CAT) [21, 23] provides a mechanism to configure each logical process with a *class of service*, and allocates LLC cache ways to logical processes. The CAT manual explicitly states that a cache access will hit if the line is cached in *any* of the cache’s ways, regardless of CAT settings. However, one CAT domain will not cause evictions in another domain. To achieve CAT’s properties, no critical path changes in the cache are required: CAT’s behavior on a cache hit is identical to a generic cache. Victim selection (replacement policy), however, must be made aware of the CAT configuration in order to constrain ways on an eviction.

Via this quality of service (QoS) mechanism, CAT improves system performance because an inefficient, cache-hungry process can be reined in and made to only cause evictions in a subset of the LLC, instead of trashing the entire cache. The fact that the cache checks all ways for cache hits is also good for performance: shared data need not be duplicated, and overhead due to internal fragmentation of cache ways is reduced. The number of ways for each domain can also be dynamically adjusted. For example, DynaWay [17] uses CAT with online performance monitoring to adjust the ways per domain.

CAT-style partitioning is unfortunately insufficient for stopping privacy leaks through the cache: even though a process can be made to only evict in their own ways, access patterns leak through metadata updates on hitting loads. Furthermore, an attacker can spy on the victim’s private ways if the address of a transmitting line is mapped by the attacker (flush&reload, or flush&flush [20]). This is viable in the case of Spectre-style attacks, since the victim can be made to speculatively prefetch arbitrary addresses, including those in shared pages (OpenSSL, kernel, etc.).

2.3.3 Reducing privacy leakage from caches

PLcache [33, 55] and the Random Fill Cache Architecture (RFill, [37]) were designed and analyzed in the context of a small region of sensitive data. RPlcache [33, 55] trusts the OS to assign different hardware process IDs to mutually mistrusting entities, and its mechanism does not directly scale to large LLCs. The non-monopolizable cache [15] uses a well-principled partitioning scheme, but does not completely stop leakage, and relies on the OS to assign hardware process IDs. CATalyst [36] trusts the Xen hypervisor to correctly tame Intel’s Cache Allocation Technology into providing cache pinning, which can only secure software whose code and data fits into a fraction of the LLC, e.g., each virtual machine is given 8 “secure” pages. [49] similarly depends on CAT for the KVM (Kernel-based Virtual Machine) hypervisor.

3. DYNAMICALLY ALLOCATED WAY GUARD (DAWG) HARDWARE

3.1 High-level design

Consider a conventional set-associative cache, a structure comprised of several *ways*, each of which is essentially a direct-mapped cache, as well as a controller mechanism. In order to implement Dynamically Allocated Way Guard

state across protection domains, considering both the cache data and the metadata. This simple scenario will be generalized to dynamic allocation in Section 3.8 and we discuss the handling of cache replacement metadata in Section 3.9 for a variety of replacement policies.

Consider a read access with address $A \implies (T_A, S_A)$, where T_A and S_A are the set index and tag derived from the bits of A , respectively. A conventional set associative cache scans set S_A in all ways W_i and compares the stored tags T_{W_i} against T_A . If no match is found, this access is a cache *miss*, and will require a cache *fill* (whereby the desired cache line is brought in), and possibly an *eviction*, where an existing cache line is *flushed* out of the cache to make room for the incoming fill. Without DAWG, one domain’s accesses may cause evictions in another, as the cache does not respect domain boundaries.

Straightforwardly, cache misses in a DAWG cache must not cause fills or evictions outside the requesting privacy domain’s ways in order to enforce DAWG’s isolation. Like Intel’s CAT (Section 2.3.2), our design ensures that only the ways that a process has been allocated (via its protection domain’s `policy_fillmap` policy MSRs) are candidates for eviction; but we also restrict CLFLUSH instructions. Hardware instrumentation needed to accomplish this is highlighted in Figure 7.

3.5 DAWG’s cache metadata isolation

The *cache set metadata* structure in Figure 7 stores per-line helper data including replacement policy and cache coherence state. The metadata update logic uses tag comparisons (hit information) from all ways to modify set replacement state. DAWG does not leak via the coherence metadata, as coherence traffic is tagged with the requestors’s protection domain and does not modify lines in other domains (with a sole exception described in Section 3.7).

DAWG’s replacement metadata isolation requirement, at a high level, is a non-interference property: victim selection in a protection domain should not be affected by the accesses performed against any other protection domain(s). Furthermore, the cache’s replacement policy must allow system software to sanitize the replacement data of a way in order to implement safe protection domain resizing. Details of implementing DAWG-friendly partitionable cache replacement policies are explored in Section 3.9,

3.6 DAWG’s cache hit isolation

Cache hits in a DAWG cache must also be isolated, requiring a change to the critical path of the cache structure: a cache access must not hit in ways it was not allocated – a possibility if physical tags are shared across protection domains.

Consider a read access with address $A \implies (T_A, S_A)$ (tag and set, respectively) in a conventional set associative cache. A match on any of the way comparisons indicates a cache *hit* ($\exists i \mid T_{W_i} == T_A \implies hit$); the associated cache line data is returned to the requesting core, and the replacement policy metadata is updated to make note of the access. This allows a receiver (attacker) to communicate via the cache state by probing the cache tag or metadata state as described in Section 2.2.

In DAWG, tag comparisons must be masked with a policy (`policy_hitmap`) that white-lists ways allocated to the re-

quester’s protection domain ($\exists i \mid \text{policy_hitmap}[i] \ \& \ (T_{W_i} == T_A) \implies hit$). By configuring `policy_hitmap`, system software can ensure cache hits are not visible across protection domains. While the additional required hardware in DAWG caches’ hit path adds a gate delay to each cache access, we note that modern L1 caches are usually pipelined. We expect hardware designers will be able to manage an additional low-fanout gate without affecting clock frequency.

In addition to masking hits, DAWG’s metadata update must use this policy-masked hit information to modify any replacement policy state safely, preventing information leakage across protection domains via the replacement policy state, as described in Section 3.5.

3.7 Cache lines shared across domains

DAWG effectively hides cache hits outside the white-listed ways as per `policy_hitmap`. While this prevents information leakage via adversarial observation of cached lines, it also complicates the case where addresses are shared across two or more protection domains by allowing ways belonging to different protection domains to have copies of the same line. Read-only data and instruction misses acquire lines in the Shared state of the MESI protocol [42] and its variants.

Neither a conventional set associative cache nor Intel’s CAT permit duplicating a cache line within a cache: their hardware enforces a simple invariant that *a given tag can only exist in a single way of a cache at any time*. In the case of a DAWG cache, the hardware does not strictly enforce this invariant across protection domains; we allow read-only cache lines (in Shared state) to be replicated across ways in different protection domains.

As a concrete example where this may require attention, consider again a DAWG LLC configured with two protection domains, with half of the ways allocated to the kernel, and the remaining half to unprivileged software. Suppose the kernel has address A cached in Shared state, and the same address is read by a user-mode application, in this case permitted by the kernel. The domain’s `policy_hitmap` hides the cache hit in kernel’s ways, reports a miss, and triggers a cache fill, causing A to be cached in two ways of set S_A , both lines being Shared.

Cache line sharing may leak information via the cache coherence protocol (whereby one domain can invalidate lines in another), or entirely violate invariants expected by the cache coherence protocol (by creating a situation where multiple copies of a line exist when one is Modified).

In order to maintain isolation, cache coherence traffic must respect DAWG’s protection domain boundaries. Requests on the same line from different domains are therefore non-matching, and are filled by the memory controller. Cache flush instructions (CLFLUSH, CLWB) affect only the ways allocated to the requesting `domain_id`. Cross-socket invalidation requests must likewise communicate their originating protection domain. DAWG caches are *not*, however, expected to handle a replicated Modified line, meaning system software must not allow shared writable pages across protection domains via a TLB invariant, as described in Section 4.2.2.

A further optimization for single-socket systems is to allow Modified lines to remain in the LLC without writeback to DRAM when a page is marked Copy-on-Write. The con-

tents of Modified cache lines overrides misses on behalf of other protection domains. Therefore, data reflects the latest contents, while timing matches that of DRAM access.

Stale Shared lines of de-allocated pages may linger in the cache; DAWG must invalidate these before zeroing a page to be granted to a process (see Section 4.2.2). To this end, DAWG requires a new privileged MSR, with which to *invalidate all copies* of a Shared line, given an address, regardless of protection domain. DAWG relies on system software to prevent the case of a replicated Modified line.

3.8 Dynamic allocation of ways

It is unreasonable to implement a static protection domain policy, as it would make inefficient use of the cache resources due to internal fragmentation of ways. Instead, DAWG caches can be provisioned with updated security policies *dynamically*, as the system’s workload changes.

In order to maintain its security properties, system software must manage protection domains by manipulating the domains’ `policy_hitmap` and `policy_fillmap` MSRs in the DAWG cache. These MSRs are normally equal, but diverge to enable concurrent use of shared caches.

In order to re-assign a DAWG cache way, when creating or modifying the system’s protection domains, the way must be *invalidated*, destroying any private information in form of the cache tags and metadata for the way(s) in question. In the case of write-back caches, dirty cache lines in the affected ways must be written-back, or swapped within the set. A privileged software routine flushes one or more ways via a hardware affordance to perform fine-grained cache flushes by set&way, e.g., available on ARM [1].

We require hardware mechanisms to flush a line and/or perform write-back (if M), of a specified way in a DAWG memory cache, allowing privileged software to orchestrate way-flushing as part of its software management of protection domains. This functionality is exposed for each cache, and therefore accommodates systems with diverse hierarchies of DAWG caches. We discuss the software mechanism to accommodate dynamic protection domains in Section 4.1.2.

While this manuscript does describe the mechanism to adjust the DAWG policies in order to create, grow, or shrink protection domains, we leave as future work resource management support to securely determine the efficient sizes of protection domains for a given workload.

3.9 Replacement Policies

In this section, we will exemplify the implementation of several common replacement policies compatible with DAWG’s isolation requirement. We focus here on several commonplace replacement policies, given that cache replacement policies are diverse. The optimal policy for a workload depends on the effective associativity and may even be software-selected, e.g., ARM A72 [1] allows pseudo-random or pseudo-LRU cache-replacement.

Because victim selection among the ways of a cache can be seen as a decision tree, a convenient means to partition a replacement policy of a DAWG cache is to allocate ways at sub-tree granularity. Without indirection between the policy bitmap and the ways of the cache, this imposes some restrictions on way partitioning: processes can be given one

Consider a cache access that misses in its protection domain:

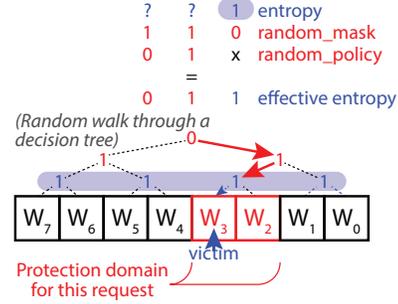


Figure 8: Victim selection with a DAWG-partitioned random policy.

or more ways as long as these ways are *consecutive*, the number of allocated ways is a power-of-two, and the allocation is aligned to the same power of two. An example: in the case of an 8-way cache, when allocating 2 ways to a process, the process may receive [0,1],[2,3],[4,5],[6,7], but *not* [1,2]. When partitioning a 4-way cache, [0:1], or [2:3] are the only allowed allocations. We assume that a system where this is unacceptable would add a level of indirection between the policy bitmap and the cache ways, as this is not on the critical path.

3.9.1 Random replacement policy

In a typical cache with a random replacement policy, a word of $\log_2(\text{ways})$ bits of entropy (for example, from an LFSR) decide a walk through a decision tree to select a victim way (where “1” signals “select left subtree”). A partitioned replacement policy forces a prefix of these decisions, as directed by `policy_fillmap`.

From the selected `policy_fillmap`, a DAWG cache with a random replacement policy derives two new quantities: `random_mask` and `random_policy`. To simplify this discussion, assume protection domains are allocated at subtree granularity. If the system employs an indirection to virtualize the mapping of policy bits to cache ways, as motivated in Sections 3.2, 3.4, this constraint can be relaxed.

The cache derives `random_mask` from `policy_fillmap`, by *clearing* the i^{th} bit of the policy if and only if any entire left subtree of $2^{(i+1)}$ leaves is allocated by the `policy_fillmap`. For example, if a protection domain is allocated ways W_2, W_3 of 8 ways, as shown in Figure 8, then `random_mask=0b110`.

Similarly, the i^{th} bit of `random_policy` is set if any *left* subtree of 2^i leaves is *entirely* allocated by `policy_fillmap`. In the example above, `random_policy=0b011x` (0b011, with the low-order bit ignored as per `random_mask`).

This logic is concurrent with metadata and way lookups on cache accesses, and is not critical. To generalize:

$$\text{random_mask}[i] = \bigvee_{j=\text{multiples of } (2^{i+1})}^{\text{ways}} \left(\bigwedge_{k=0}^{2^{i+1}-1} \text{policy_fillmap}[j+k] \right)$$

$$\text{random_policy}[i] = \bigvee_{j=\text{multiples of } (2^i)}^{\text{ways}} \left(\bigwedge_{k=2^i}^{2^{i+1}-1} \text{policy_fillmap}[j+k] \right)$$

When changing the policy, no special consideration is

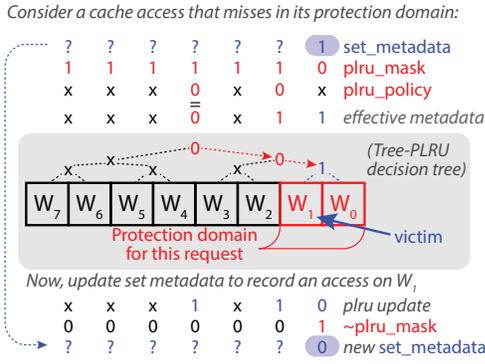


Figure 9: Victim selection and metadata update with a DAWG-partitioned Tree-PLRU policy.

required in a DAWG cache with a random replacement policy, as no metadata persists across accesses, and therefore across policy changes.

3.9.2 Tree-PLRU: pseudo least recently used

Tree-PLRU “approximates” LRU with a small set of bits stored per cache line. The victim selection is a (complete) decision tree informed by metadata bits. 1 signals “go left”, whereas 0 signals “go right” to reach the PLRU element, as shown in Figure 9.

The cache derives a `plru_mask` and `plru_policy` from `policy_fillmap`. These fields augment a decision tree over the ways of the cache; a bit of `plru_mask` is 0 if and only if its corresponding subtree in `policy_fillmap` has no zeroes (if the subtree of the decision tree is entirely allocated to the protection domain). Similarly, `plru_policy` bits are set if their corresponding *left* subtrees contain one or more ways allocated to the protection domain. For example, if a protection domain is allocated ways W_0, W_1 of 8 ways, then `plru_mask=0b11111110`, and `plru_policy=0bxxx0x0x` (0b0000001, to be precise, with x marking masked and unused bits).

The mapping of `policy_fillmap` to `plru_mask` and `plru_policy` for an 8-way cache is given below:

i	$\neg \bigwedge_j \text{policy_fillmap}[j]$, where j is:	$\bigvee_j \text{policy_fillmap}[j]$, where j is:
0	0,1	1
1	0,1,2,3	2,3
2	2,3	3
3	0,1,2,3,4,5,6,7	4,5,6,7
4	4,5	5
5	4,5,6,7	6,7
6	6,7	6

At each access, `set_metadata` is updated by changing each bit on the branch leading to the hitting way to be the *opposite* of the direction taken, i.e., “away” from the most recently used way. For example, when accessing W_5 , metadata bits are updated by $b_0 \rightarrow 0$, $b_2 \rightarrow 1$, $b_5 \rightarrow 0$. These updates are masked to avoid modifying PLRU bits *above* the allocated subtree. For example, when $\{W_2, W_3\}$ are allocated to the process, and it hits W_3 , b_0 and b_1 remain unmodified to avoid leaking information via the metadata updates.

Furthermore, we must mask `set_metadata` bits that are made irrelevant by the allocation. For example, when $\{W_2, W_3\}$

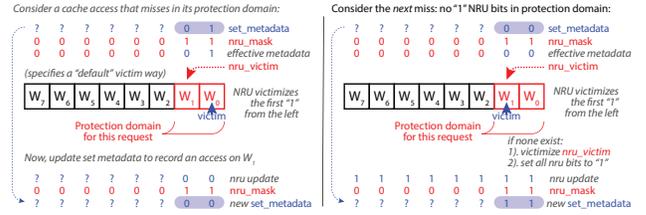


Figure 10: Victim selection and metadata update with a DAWG-partitioned NRU policy.

are allocated to the process, the victim selection should always reach the b_4 node when searching for the pseudo-LRU way. To do this, ignore $\{b_0, b_1\}$ in the metadata table, and use values 0 and 1, respectively.

Observe that both are straightforwardly implemented via `plru_mask` and `plru_policy`. This forces a subset of decision tree bits, as specified by the policy: victim selection logic uses $(\text{set_metadata} \& \sim \text{plru_mask}) \mid (\text{plru_mask} \& \text{plru_policy})$. This ensures that system software is able to restrict victim selection to a subtree over the cache ways. Metadata updates are partitioned also, by constraining updates to `set_metadata & ~plru_mask`. When system software alters the cache’s policies, and re-assigns a way to a different protection domain, it must take care to force the way’s metadata to a known value in order to avoid private information leakage.

3.9.3 SRRIP and NRU: Not recently used

An NRU policy requires one bit of metadata per way be stored with each set. On a cache hit, the accessed way’s NRU bit is set to “0”. On a cache miss, the victim is the first (according to some pre-determined order, such as left-to-right) line with a “1” NRU bit. If none exists, the first line is victimized, and all NRU bits of the set are set to “1”.

Enforcing DAWG’s isolation across protection domains for an NRU policy is a simple matter, as shown in Figure 10. As before, restrict metadata updates to the ways white-listed by `nru_mask = policy_fillmap`. In order to victimize only among ways white-listed by the policy, mask the NRU bits of all other ways via `set_metadata & nru_mask` at the input to the NRU replacement logic.

Instead of victimizing the *first* cache line if no “1” bits are found, the victim way must fall into the current protection domain. To implement this, specify the default victim via `nru_victim`, which selects the leftmost way with a corresponding “1” bit of `nru_mask`, whereas the unmodified NRU is hard-wired to evict a specific way.

The SRRIP [28] replacement policy is similar, but expands the state space of each line from two to four (or more) states by adding a *counter* to track ways less favored to be victimized. Much like NRU, SRRIP victimizes the first (up to some pre-determined order) line with the largest counter during a fill that requires eviction. To partition SRRIP, use the same `nru_mask = policy_fillmap`, mask each line’s metadata with the way’s bit of `nru_mask` to other domains’ lines to be stuck at “recently used” and not be candidates for eviction.

3.9.4 Other replacement policies

A true LRU (least recently used) policy has an interesting property: it does not leak private information across pro-

tection domains, despite requiring metadata updates across statically defined protection domains (the LRU rank of all lines in a set must be modified as part of each update). This is because software lacks a means to detect the LRU rank of the ways in its protection domains, as this information is only exposed via victim selection on a cache miss. Therefore, only the *relative* rank of lines in a protection domain may be observed.

By way of an analogy, consider a deck of cards. If a subset of cards is *marked* (granted to a protection domain), and an *unmarked* card is drawn and placed at the top of the deck, the relative order of the *marked* cards has not changed. Equivalently, updating the LRU metadata of a line in a given protection domain cannot cause the relative order of lines in another protection domain to change. The LRU policy is therefore DAWG-friendly by design: one would not modify its implementation in a DAWG cache (but would need to restrict victim selection to the allocated sets only, and sanitize LRU rank of any ways re-allocated to another domain).

Complex cache replacement policies such as DRRIP [28] are more challenging to partition securely. DRRIP employs dueling sets, dedicating a small group of cache sets to each of its candidate replacement policies, occasionally switching the policy of the remaining sets (follower sets) to employ the “winning” policy. The performance counters associated with dueling sets must be securely partitioned among protection domains, which is made challenging by the dynamic nature of DAWG’s policies. A detailed explanation of a partitioned DRRIP policy is out of the scope of this manuscript.

4. SOFTWARE MODIFICATIONS

We describe software provisions for modifying DAWG’s protection domains, and also describe small, required modifications to several well-annotated sections of kernel software to implement cross-domain communication primitives robust against speculative execution attacks.

4.1 Software management of DAWG policies

Protection domains are a software abstraction implemented by system software via DAWG’s policy MSRs. The policy MSRs themselves (a table mapping protection `domain_id` to a `policy_hitmap` and `policy_fillmap` at each cache, as described in Section 3.2) reside in the DAWG cache hardware, and are atomically modified.

4.1.1 DAWG Resource Allocation

Protection domains for a process tree should be specified using the same `cgroup`-like interface as Intel’s CAT. In order to orchestrate DAWG’s protection domains and policies, the operating system must track the mapping of process IDs to protection domains. In a system with 16 ways in the most associative cache, no more than 16 protection domains can be concurrently scheduled, meaning if the OS has need for more mutually distrusting entities to schedule, it needs to virtualize protection domains by time-multiplexing protection domain IDs, and flushing the ways of the multiplexed domain whenever it is re-allocated.

We augment the process data structure (`task_struct`) with a `protection_domain` byte in order to track the allocation of DAWG domains to processes. Processes with no

protection domain are mapped to an invalid protection domain (0xFF), and must not be scheduled until a protection domain can be granted. There is no requirement to grant different protection domains to each process.

Another data structure, `dawg_policy`, tracks the resources (cache ways) allocated to each protection domain. This is a table mapping `domain_id` to pairs (`policy_hitmap`, `policy_fillmap`) for each DAWG cache. The kernel uses this table when resizing, creating, or destroying protection domains in order to maintain an exclusive allocation of ways to each protection domain. Whenever one or more ways are re-allocated, the supervisor must look up the current `domain_id` of the owner, accomplished via either a search or a persistent inverse map cache way to `domain_id`.

4.1.2 Secure Dynamic Way Reassignment

When modifying an existing allocation of ways in a DAWG cache (writing policy MSRs), as necessary to create or modify protection domains, system software must sanitize (including any replacement metadata, as discussed in Section 3.5) the re-allocated way(s) before they may be granted to a new protection domain. The process for re-assigning cache way(s) proceeds as follows:

1. Update the `policy_fillmap` MSRs to disallow fills in the way(s) being transferred out of the shrinking domain.
2. A software loop iterates through the cache’s set indexes and flushes all sets of the re-allocated way(s). The shrinking domain may hit on lines yet to be flushed, as `policy_hitmap` is not yet updated.
3. Update the `policy_hitmap` MSRs to exclude ways to be removed from the shrinking protection domain.
4. Update the `policy_hitmap` and `policy_fillmap` MSRs to grant the ways to the growing protection domain.

Higher level policies can be built on this dynamic way-reassignment mechanism.

4.1.3 Code Prioritization

Programming the domain selectors for code and data separately allows ways to be dedicated to code without data interference. Commercial studies of code cache sensitivity of production server workloads [25, 29, 41] show large instruction miss rates in L2, but even the largest code working sets fit within 1–2 L3 ways. Code prioritization will also reduce the performance impact of disallowing code sharing across domains, especially when context switching between untrusted domains sharing code. DAWG domains cannot reuse shared code in L2 in the worst case scenario of context switching to untrusted domains. While partitioning protects against context switching to other code, adding dynamic way reassignment and code prioritization can offset loss of sharing.

4.2 Kernel changes required by DAWG

Consider a likely configuration where a user-mode application and the OS kernel are in different protection domains. In order to perform a system call, communication must occur across the protection domains: the supervisor extracts the (possibly cached) data from the caller by copying into its own memory. In DAWG, this presents a challenge due to strong

isolation in the cache.

4.2.1 DAWG augments SMAP-annotated sections

We take advantage of modern OS support for the Supervisor Mode Access Prevention (SMAP) feature available in recent x86 architectures, which allows supervisor mode programs to raise a trap on accesses to user-space memory. The intent is to harden the kernel against malicious programs attempting to influence privileged execution via untrusted user-space memory. At each routine where supervisor code intends to access user-space memory, SMAP must be temporarily disabled and subsequently re-enabled via `stac` (Set AC Flag) and `clac` (Clear AC Flag) instructions, respectively. We observe that a modern kernel’s interactions with user-space memory are diligently annotated with these instructions, and will refer to these sections as *annotated sections*.

Currently Linux kernels use seven such sections for simple memory copy or clearing routines: `copy_from_user`, `copy_to_user`, `clear_user`, etc. We propose extending these annotated sections with short instruction sequences to correctly handle DAWG’s communication requirements on system calls and inter-process communication, in addition to the existing handling of the SMAP mechanism. Specifically, sections implementing data movement *from* user to kernel memory are annotated with an MSR write to `domain_id`: *ifetch* and *store* accesses proceed on behalf of the kernel, as before, but *load* accesses use the caller’s (user) protection domain. This allows the kernel to efficiently copy from warm cache lines, but preserves isolation. After copying from the user, the `domain_id` MSR is restored to perform all accesses on behalf of the kernel’s protection domain. Likewise, sections implementing data movement *to* user memory *ifetch* and *load* on behalf of the kernel’s domain, but *store* in the user’s cache ways. While the annotated sections may be interrupted by asynchronous events, interrupt handlers are expected to explicitly set `domain_id` to the kernel’s protection domain, and restore the MSR to its *prior* state afterwards.

As described in Section 3.3, DAWG’s `domain_id` MSR writes are a fence, preventing speculative disabling of DAWG’s protection mechanism. Current Linux distributions diligently pair a `stac` instruction with an `lfence` instruction to prevent speculative execution within regions that access user-mode memory, meaning DAWG does not significantly serialize annotated sections over its insecure baseline.

Finally, to guarantee isolation, we require the annotated sections to contain only code that obeys certain properties: to protect against known and future speculative attacks, indirect jumps or calls, and potentially unsafe branches are not to be used. Further, we cannot guarantee that these sections will not require patching as new attacks are discovered, although this is reasonable given the small number and size of the annotated sections.

4.2.2 Read-only and CoW sharing across domains

For memory efficiency, DAWG allows securely mapping read-only pages across protection domains, e.g., for shared libraries, requiring hardware cache coherence protocol changes (see Section 3.7), and OS/hypervisor support.

This enables conventional system optimizations via page sharing, such as read-only `mmap` from page caches, Copy-on-

Write (CoW) conventionally used for `fork`, or for page deduplication across VMs (e.g., Transparent Page Sharing [54]; VM page sharing is typically disabled due to concerns raised by shared cache tag attacks [24]). DAWG maintains security with read-only mappings across protection domains to maintain memory efficiency.

Dirty pages can be prepared for CoW sharing eagerly, or lazily (but cautiously [56]) by installing non-present pages in the consumer domain mapping. Preparing a dirty page for sharing *requires* a write-back of any dirty cache lines on behalf of the producer’s domain (via `CLWB` instructions and an appropriate `load domain_id`). The writeback guarantees that read-only pages appear only as `Shared` lines in DAWG caches, and can be replicated across protection domains as described in Section 3.7.

A write to a page read-only shared across protection domains signals the OS to create a new, private copy using the original producer’s `domain_id` for reads, and the consumer’s `domain_id` for writes.

During copying, cache lines of the original page must also get invalidated. The last step prepares the consumer domain for possible future uses of the original page.

4.2.3 Reclamation of shared physical pages

Before cache lines may be filled in a new protection domain, pages reclaimed from a protection domain must be removed from DAWG caches as part of normal OS page cleansing. Prior to zeroing (or preparing for DMA) a page previously shared across protection domains, the OS must *invalidate all cache lines* belonging to the page, as described in Section 3.7. The same is required between `unmap` and `mmap` operations over the same physical addresses. For most applications, therefore cache line invalidation can be deferred to wholesale destruction of protection domains at `exit`, given ample physical memory.

5. SECURITY ANALYSIS

We explain why DAWG protects against attacks realized thus far on speculative execution processors in Section 5.1 by stating and arguing a non-interference property. We then argue in Section 5.2 that system calls and other cross-domain communication are safe. Finally, we show a generalization of our attack schema and point out the limitations of cache partitioning in Section 5.3.

5.1 DAWG Isolation Property

DAWG enforces isolation of *exclusive* protection domains among cache tags and replacement metadata, as long as:

1. victim selection is restricted to the ways allocated to the protection domain (an invariant maintained by system software), and
2. metadata updates as a result of an access in one domain do not affect victim selection in another domain (a requirement on DAWG’s cache replacement policy).

Together, this guarantees non-interference – the hits and misses of a program running in one protection domain are unaffected by program behavior in different protection domains. As a result, DAWG blocks the cache tag and metadata channels of non-communicating processes separated by DAWG’s protection domains.

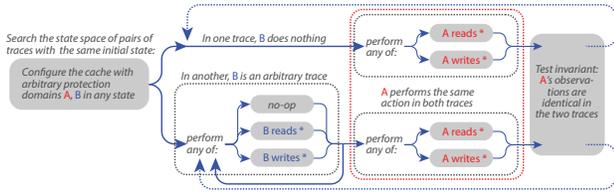


Figure 11: Formal model verification for a DAWG cache’s non-interference property on statically configured protection domains.

Figure 11 illustrates DAWG’s non-interference property, which can be proven in a manner similar to the proof for the Sanctum processor’s cache partitioning scheme [50].

5.2 No leaks from system calls

Consider a case where the kernel (victim) and a user program (attacker) reside in different protection domains. While both use the same cache hierarchy, they share neither cache *lines* nor *metadata* (Section 3.2), effectively closing the cache exfiltration channel. In few, well-defined instances where data is passed between them (such as `copy_to_user`), the kernel accesses the attacker’s ways to read/write user memory, leaking the (public) access pattern associated with the bulk copying of the syscall inputs and outputs (Section 4.2). Writes to DAWG’s MSRs are *fences*, and the annotated sections must not offer opportunity to maliciously mis-speculate control flow (see Section 4.2.1), thwarting speculative disabling or misuse of DAWG’s protection domains. DAWG also blocks leaks via coherence Metadata, as coherence traffic is restricted to its protection domain (Section 3.7), with the sole exception of cross-domain *invalidation*, where physical pages are reclaimed and sanitized.

When re-allocating cache ways, as part of resizing or multiplexing protection domains, no private information is transferred to the receiving protection domain: the kernel sanitizes ways before they are granted, as described in Section 4.1.2. Physical pages re-allocated across protection domains are likewise sanitized (Section 4.2.3).

When an application makes a system call, the necessary communication (data copying) between kernel and user program must not leak information beyond what is communicated. The OS’s correct handling of `domain_id` MSR within annotated sections, as described in Section 4.2 ensures user space cache side effects reflect the section’s explicit memory accesses.

5.3 Limitations of cache partitioning

DAWG’s cache isolation goals are meant to approach the isolation guarantees of separate machines, yet, even remote network services can fall victim to leaks employing cache tag state for communication. Consider the example of the attacker and victim residing in different protection domains, sharing no data, but communicating via some API, such as system calls. As in a remote network timing leak [11], where network latency is used to communicate some hidden state in the victim, the *completion time* of API calls can communicate insights about the cache state [31] within a protection domain.

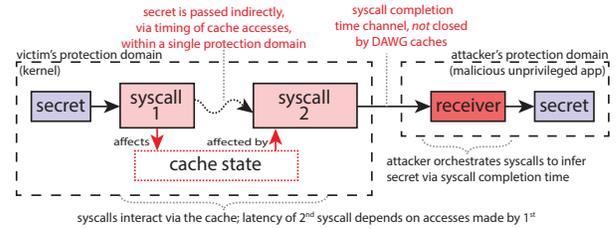


Figure 12: Generalized Attack Schema: an adversary 1) accesses a victim’s secret, 2) reflects it to a transmitter 3) transmits it via a covert channel, 4) receives it in their own protection domain.

Leakage via *reflection* through the cache is thus possible: the receiver invokes an API call that accesses private information, which affects the state of its private cache ways. The receiver then exfiltrates this information via the latency of another API call. Figure 12 shows a cache timing leak which relies on cache reflection entirely within the victim’s protection domain. The syscall completion time channel is used for exfiltration, meaning no private information crosses DAWG’s domain boundaries in the caches, rendering DAWG, and cache partitioning in general, ineffective at closing a leak of this type.

The transmitter is instructed via an API call to access $a[b[i]]$, where i is provided by the receiver (via `syscall1`), while a, b reside in the victim’s protection domain. The cache tag state of the transmitter now reflects $b[i]$, affecting the latency of subsequent syscalls in a way dependent on the secret $b[i]$. The receiver now exfiltrates information about $b[i]$ by selecting a j from the space of possible values of $b[i]$ and measuring the completion time of `syscall2`, which accesses $a[j]$. The syscall completion time communicates whether the transmitter hits on $a[j]$, which implies $a[j] \cong a[b[i]]$, and for a compact a , that $b[i] = j$ – a leak. This leak can be amplified by initializing cache state via a bulk memory operation, and, for a machine-local receiver by malicious mis-speculation.

While not the focus of this paper, for completeness, we outline a few countermeasures for this type of leak. Observe that the completion time of a public API is used here to exfiltrate private information. The execution time of a syscall can be padded to guarantee constant (and worst-case) latency, no matter the input or internal state. This can be relaxed to bound the leak to a known number of bits per access [18].

A zero leak countermeasure requires destroying the transmitting domain’s cache state across syscalls/API invocations, preventing reflection via the cache. DAWG can make this less inefficient: in addition to dynamic resizing, setting the replacement mask `policy_fillmap` to a subset of the `policy_hitmap` allows locking cache ways to preserve the hot working set. This ensures that all unique cache lines accessed during one request have constant observable time.

6. EVALUATION

To evaluate DAWG, we use the `zsim` [48] execution-driven x86-64 simulator and Haswell hardware [16] for our experiments.

6.1 Configuration of insecure baseline

Table 1 summarizes the characteristics of the simulated

Cores		DRAM		Bandwidth	
Count	Frequency	Controllers		<i>Peak</i>	
8 OoO	3 GHz	4 x DDR3-1333		42 GB/s	
Private Caches			Shared Cache		
L1	L2	Organization		L3	Organization
2 × 32 KB	256 KB	8-way PLRU		8 × 2 MB	16-way NRU

Table 1: Simulated system specifications.

environment. The out-of-order model implemented by `zsim` is calibrated against Intel Westmere, informing our choice of cache and network-on-chip latencies. The DRAM configuration is typical for contemporary servers at ~ 5 GB/s theoretical DRAM bandwidth per core. Our baseline uses the Tree-PLRU (Section 3.9.2) replacement policy for private caches, and a 2-bit NRU for the shared LLC. The simulated model implements inclusive caches, although DAWG domains with reduced associativity would benefit from relaxed inclusion [25]. We simulate CAT partitioning at *all* levels of the cache, while modern hardware only offers this at the LLC. We do this by restricting the replacement mask `policy_fillmap`, while white-listing all ways via the `policy_hitmap`.

6.2 DAWG Policy Scenarios

We evaluate several protection domain configurations for different resource sharing and isolation scenarios.

6.2.1 VM or container isolation on dedicated cores

Isolating peer protection domains from one another requires equitable LLC partitioning, e.g., 50% of ways allocated to two active domains. In the case of cores dedicated to each workload (no context switches), each scheduled domain is assigned the entirety of its L1 and L2.

6.2.2 VM or container isolation on time-shared cores

To allow the OS to overcommit cores across protection domains (thus requiring frequent context switches between domains), we also evaluate a partitioned L2 cache.

6.2.3 OS isolation

Only two DAWG domains are needed to isolate an OS from applications. For processes with few OS interventions in the steady state, e.g., SPECCPU workloads, the OS can reserve a single way in the LLC, and flush L1 and L2 ways to service the rare system calls. Processes utilizing more OS services would benefit from more ways allocated to OS’s domain.

6.3 DAWG versus insecure baseline

Way partitioning mechanisms reduce cache capacity and associativity, which increases conflict misses, but improves fairness and reduces contention. We refer to CAT [21] for analysis of the performance impact of way partitioning on a subset of SPEC CPU2006. Here, we evaluate CAT and DAWG on parallel applications from PARSEC [7], and parallel graph applications from the GAP Benchmark Suite (GAPBS) [5], which allows a sweep of workload sizes.

Figure 13 shows DAWG partitioning of private L1 and L2 (Section 6.2.2) caches in addition to the L3. We explore DAWG configurations on a subset of PARSEC benchmarks

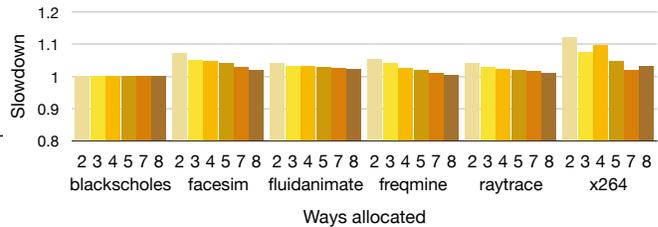


Figure 13: Way partitioning performance at low associativity in all caches (8-way L1, 8-way L2, and 16-way L3).

on `simlarge` workloads. The cache insensitive `blackscholes` (or omitted swaptions with 0.001 L2 MPKI (Misses Per 1000 Instructions)) are unaffected at any way allocation. For a VM isolation policy (Section 6.2.1) with $8/16$ of the L3, even workloads with higher MPKI such as `facesim` show at most 2% slowdown. The $\langle 2/8$ L2, $2/16$ L3) configuration is affected by both capacity and associativity reductions, yet most benchmarks have 4–7% slowdown, up to 12% for `x264`. Such an extreme configuration can accommodate 4 very frequently context switched protection domains.

Figure 14 shows performance overhead of variously sized protection domains at the L3 cache for one 4-thread instance. We use GAPBS’s synthetic *power law* graphs [13, 19] that match the structure of real-world social and web graphs and therefore exhibit cache locality [6]. The power law structure, however, implies that there is diminishing return from each additional L3 way. As shown, at half cache capacity ($8/16$ L3, Section 6.2.1), there is at most 15% slowdown (`bc` and `tc` benchmarks) at the largest simulated size (2^{20} vertices). A characteristic *eye* is formed when the performance curves of different configurations cross over the working set boundary (e.g., graph size of 2^{17}). Performance with working sets smaller or larger than the effective cache capacity are unaffected — at the largest size `cc`, `pr`, and `sssp` show 1–4% slowdown.

Reserving for the OS (Section 6.2.3), one way (6% of LLC capacity) adds no performance overhead to most workloads. The only exception would be a workload caught in the *eye*, e.g., PageRank at 2^{17} has 30% overhead (Figure 14), while at 2^{16} or 2^{18} — 0% difference.

6.4 CAT versus DAWG

We analyze and evaluate scenarios based on the degree of code and data sharing across domains.

6.4.1 No Sharing

There is virtually no performance difference between secure DAWG partitioning, and insecure CAT partitioning in the absence of read-sharing across domains.

DAWG reduces interference in replacement metadata updates and enforces the intended replacement strategy within a domain, while CAT may lose block history effectively exhibiting random replacement — a minor, workload-dependent perturbation. In simulations (not shown), we replicate a known observation that random replacement occasionally performs better than LRU near cache capacity. We did not observe this effect with NRU replacement.

6.4.2 Read-only Sharing

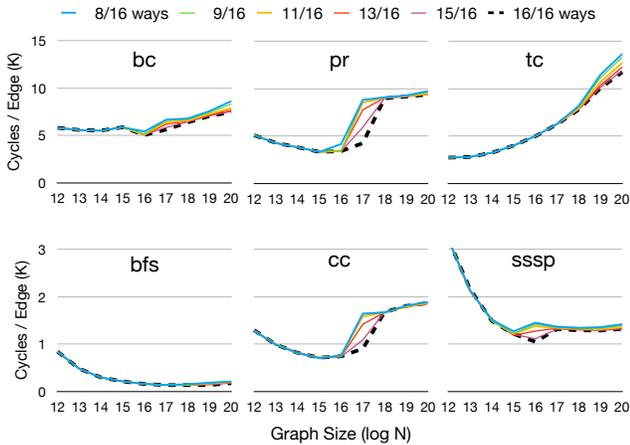


Figure 14: Way partitioning performance with varying working set on graph applications. Simulated 16-way L3.

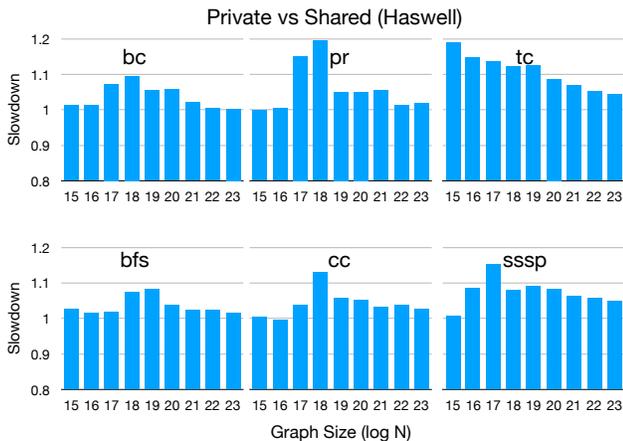


Figure 15: Read-only sharing effects of two instances using Shared vs Private data of varying scale. Two 1-thread instances. Actual Haswell 20-way 30 MB L3.

CAT QoS guarantees a lower bound on a workload’s effective cache capacity, while DAWG isolation forces a tight upper bound. DAWG’s isolation reduces cache capacity compared to CAT when cache lines are read-only shared across mutually untrusting protection domains. CAT permits hits across partitions where code or read-only data are unsafely shared. We focus on read-only data in our evaluation, as benchmarks with low L1i MPKI like GAPBS, PARSEC, or SPEC CPU are poorly suited to study code cache sensitivity.

We analyze real applications using one line modifications to GAPBS to fork (a single-thread process) either before or after creating in-memory graph representations. The first results in a private graph for each process, while the latter simulates mmap of a shared graph. The shared graphs access read-only data across domains in the baseline and CAT, while DAWG has to replicate data in domain-private ways. Since zsim does not simulate TLBs, we ensure different virtual addresses are used to avoid false sharing. We first verified

in simulation that DAWG, with memory shared across protection domains, behaves identically to CAT and the baseline with private data.

Next, we demonstrate (in Figure 15) that these benchmarks show little performance difference on real hardware [16] for most data sizes; Shared baseline models Shared CAT, while Private baseline models Shared DAWG. The majority of cycles are spent on random accesses to read-write data, while read-only data is streamed sequentially. Although read-only data is much larger than read-write data (e.g., 16 times more edges than vertices), prefetching and scan- and thrash-resistant policies [28, 46] further reduce the need for cache resident read-only data. Note that even at 2^{23} vertices these effects are immaterial; real-world graphs have billions of people or pages.

6.5 Domain copy microbenchmark

We simulated a privilege level change at simulated system calls for user-mode TCP/IP. Since `copy_from_user` and `copy_to_user` permit hits in the producer’s ways, there is no performance difference against the baseline (not shown).

7. CONCLUSION

DAWG protects against cache-timing attacks on speculative execution processors with reasonable overheads. The same policies can be applied to any set-associative structure, e.g., TLB or branch history tables. DAWG has its limitations and additional techniques are required to block exfiltration channels different from the cache channel. We believe that techniques like DAWG are needed to restore our confidence in public cloud infrastructure, and hardware and software co-design will help minimize performance overheads.

A good proxy for the performance overheads of secure DAWG is Intel’s existing, though insecure, CAT hardware. Traditional QoS uses of CAT, however, differ from desired DAWG protection domains’ configurations. Research on software resource management strategies can therefore commence with evaluation of large scale workloads on CAT. CPU vendors can similarly analyze the cost-benefits of increasing cache capacity and associativity to accommodate larger numbers of active protection domains.

8. ACKNOWLEDGMENTS

Funding for this research was partially provided by the National Science Foundation under contract number CNS-1413920, Delta Electronics, and DARPA & SPAWAR under contract N66001-15-C-4066. We are grateful to Carl Waldspurger for his valuable feedback on the initial design as well as the final presentation of this paper.

9. REFERENCES

- [1] ARM, “ARM Cortex-A72 MPCore processor technical reference manual,” 2015.
- [2] ARM, “ARM Software Speculation Barrier,” <https://github.com/ARM-software/speculation-barrier>, January 2018.
- [3] R. Balasubramonian, N. Jouppi, and N. Muralimanohar, *Multi-Core Cache Hierarchies*, 1st ed. Morgan & Claypool Publishers, 2011.
- [4] S. Banescu, “Cache timing attacks,” 2011, [Online; accessed 26-January-2014].

- [5] S. Beamer, K. Asanović, and D. A. Patterson, "The GAP benchmark suite," *CoRR*, vol. abs/1508.03619, 2015. [Online]. Available: <http://arxiv.org/abs/1508.03619>
- [6] S. Beamer, K. Asanović, and D. A. Patterson, "Locality exists in graph processing: Workload characterization on an Ivy Bridge server," in *2015 IEEE International Symposium on Workload Characterization, IISWC 2015, Atlanta, GA, USA, October 4-6, 2015*, 2015, pp. 56–65.
- [7] C. Bienia, "Benchmarking modern multiprocessors," Ph.D. dissertation, Princeton University, January 2011.
- [8] J. Bonneau and I. Mironov, "Cache-collision timing attacks against AES," in *Cryptographic Hardware and Embedded Systems-CHES 2006*. Springer, 2006, pp. 201–215.
- [9] B. B. Brumley and N. Tuveri, "Remote timing attacks are still practical," in *Computer Security-ESORICS*. Springer, 2011.
- [10] D. Brumley and D. Boneh, "Remote timing attacks are practical," *Computer Networks*, 2005.
- [11] D. Brumley and D. Boneh, "Remote timing attacks are practical," *Computer Networks*, 2005.
- [12] C. Carruth, "Introduce the "retpoline" x86 mitigation technique for variant #2 of the speculative execution vulnerabilities," <http://lists.lvm.org/pipermail/llvm-commits/Week-of-Mon-20180101/513630.html>, January 2018.
- [13] D. Chakrabarti, Y. Zhan, and C. Faloutsos, "R-MAT: A recursive model for graph mining," in *Proceedings of the Fourth SIAM International Conference on Data Mining, Lake Buena Vista, Florida, USA, April 22-24, 2004*, 2004, pp. 442–446.
- [14] J. Corbet, "KAISER: hiding the kernel from user space," <https://lwn.net/Articles/738975/>, November 2017.
- [15] L. Domnitscher, A. Jaleel, J. Loew, N. Abu-Ghazaleh, and D. Ponomarev, "Non-monopolizable caches: Low-complexity mitigation of cache side channel attacks," *Transactions on Architecture and Code Optimization (TACO)*, 2012.
- [16] E5v3, "Intel Xeon Processor E5-2680 v3(30M Cache, 2.50 GHz)," http://ark.intel.com/products/81908/Intel-Xeon-Processor-E5-2680-v3-30M-Cache-2_50-GHz.
- [17] N. El-Sayed, A. Mukkara, P.-A. Tsai, H. Kasture, X. Ma, and D. Sanchez, "KPart: A hybrid cache partitioning-sharing technique for commodity multicores," in *Proceedings of the 24th international symposium on High Performance Computer Architecture (HPCA-24)*, February 2018.
- [18] C. W. Fletcher, L. Ren, X. Yu, M. V. Dijk, O. Khan, and S. Devadas, "Suppressing the oblivious RAM timing channel while making information leakage and program efficiency trade-offs," in *2014 IEEE 20th International Symposium on High Performance Computer Architecture (HPCA)*, Feb 2014, pp. 213–224.
- [19] Graph500, "Graph 500 benchmark," <http://www.graph500.org/specifications>.
- [20] D. Gruss, C. Maurice, K. Wagner, and S. Mangard, "Flush+Flush: a fast and stealthy cache attack," in *International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer, 2016, pp. 279–299.
- [21] A. Herdrich, E. Verplanke, P. Autee, R. Illikkal, C. Gianos, R. Singhal, and R. Iyer, "Cache QoS: From concept to reality in the Intel Xeon processor E5-2600 v3 product family," in *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, March 2016, pp. 657–668.
- [22] J. Horn, "Reading privileged memory with a side-channel," <https://googleprojectzero.blogspot.com/2018/01/>, January 2018.
- [23] Intel Corp., "Improving real-time performance by utilizing Cache Allocation Technology," April 2015.
- [24] G. Irazoqui, M. S. Inci, T. Eisenbarth, and B. Sunar, "Wait a minute! a fast, cross-VM attack on AES," in *International Workshop on Recent Advances in Intrusion Detection*. Springer, 2014, pp. 299–319.
- [25] A. Jaleel, J. Nuzman, A. Moga, S. C. Steely, and J. Emer, "High performing cache hierarchies for server workloads: Relaxing inclusion to capture the latency benefits of exclusive caches," in *2015 IEEE 21st International Symposium on High Performance Computer Architecture (HPCA)*, Feb 2015, pp. 343–353.
- [26] A. Jaleel, E. Borch, M. Bhandaru, S. C. Steely Jr, and J. S. Emer, "Achieving non-inclusive cache performance with inclusive caches: Temporal locality aware TLA cache management policies," in *Proceedings of the 2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE Computer Society, 2010, pp. 151–162.
- [27] A. Jaleel, W. Hasenplaugh, M. Qureshi, J. Sebot, S. Steely, Jr., and J. Emer, "Adaptive insertion policies for managing shared caches," in *Proceedings of the 17th International Conference on Parallel Architectures and Compilation Techniques*, ser. PACT '08. New York, NY, USA: ACM, 2008, pp. 208–219. [Online]. Available: <http://doi.acm.org/10.1145/1454115.1454145>
- [28] A. Jaleel, K. B. Theobald, S. C. S. Jr., and J. S. Emer, "High performance cache replacement using re-reference interval prediction (RRIP)," in *37th International Symposium on Computer Architecture (ISCA 2010), June 19-23, 2010, Saint-Malo, France*, 2010, pp. 60–71. [Online]. Available: <http://doi.acm.org/10.1145/1815961.1815971>
- [29] S. Kanev, J. P. Darago, K. Hazelwood, P. Ranganathan, T. Moseley, G. Y. Wei, and D. Brooks, "Profiling a warehouse-scale computer," in *2015 ACM/IEEE 42nd Annual International Symposium on Computer Architecture (ISCA)*, June 2015, pp. 158–169.
- [30] R. E. Kessler and M. D. Hill, "Page placement algorithms for large real-indexed caches," *Transactions on Computer Systems (TOCS)*, 1992.
- [31] P. Kocher, D. Genkin, D. Gruss, W. Haas, M. Hamburg, M. Lipp, S. Mangard, T. Prescher, M. Schwarz, and Y. Yarom, "Spectre attacks: Exploiting speculative execution," *ArXiv e-prints*, Jan. 2018.
- [32] P. C. Kocher, "Timing attacks on implementations of Diffie-Hellman, RSA, DSS, and other systems," in *Advances in Cryptology (CRYPTO)*. Springer, 1996.
- [33] J. Kong, O. Acicmez, J.-P. Seifert, and H. Zhou, "Deconstructing new cache designs for thwarting software cache-based side channel attacks," in *workshop on Computer security architectures*. ACM, 2008.
- [34] J. Lin, Q. Lu, X. Ding, Z. Zhang, X. Zhang, and P. Sadayappan, "Gaining insights into multicore cache partitioning: Bridging the gap between simulation and real systems," in *HPCA*. IEEE, 2008.
- [35] M. Lipp, M. Schwarz, D. Gruss, T. Prescher, W. Haas, S. Mangard, P. Kocher, D. Genkin, Y. Yarom, and M. Hamburg, "Meltdown," *ArXiv e-prints*, Jan. 2018.
- [36] F. Liu, Q. Ge, Y. Yarom, F. Mckeen, C. Rozas, G. Heiser, and R. B. Lee, "CATalyst: Defeating last-level cache side channel attacks in cloud computing," in *HPCA*, Mar 2016.
- [37] F. Liu and R. B. Lee, "Random fill cache architecture," in *Microarchitecture (MICRO)*. IEEE, 2014.
- [38] F. Liu, Y. Yarom, Q. Ge, G. Heiser, and R. B. Lee, "Last-level cache side-channel attacks are practical," in *Security and Privacy*. IEEE, 2015.
- [39] Y. Oren, V. P. Kemerlis, S. Sethumadhavan, and A. D. Keromytis, "The spy in the sandbox – practical cache attacks in javascript," *arXiv preprint arXiv:1502.07373*, 2015.
- [40] D. A. Osvik, A. Shamir, and E. Tromer, "Cache attacks and countermeasures: the case of AES," in *Topics in Cryptology-CT-RSA 2006*. Springer, 2006, pp. 1–20.
- [41] G. Ottoni and B. Maher, "Optimizing function placement for large-scale data-center applications," in *2017 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, Feb 2017, pp. 233–244.
- [42] M. S. Papamarcos and J. H. Patel, "A low-overhead coherence solution for multiprocessors with private cache memories," *SIGARCH Comput. Archit. News*, vol. 12, no. 3, pp. 348–354, Jan. 1984.
- [43] A. Pardoe, "Spectre mitigations in msvc," <https://blogs.msdn.microsoft.com/vcblog/2018/01/15/spectre-mitigations-in-msvc/>, January 2018.
- [44] B. Pham, J. Vesely, G. H. Loh, and A. Bhattacharjee, "Large pages and lightweight memory management in virtualized environments: Can you have it both ways?" in *Proceedings of the 48th International Symposium on Microarchitecture*, ser. MICRO-48. New York, NY, USA: ACM, 2015, pp. 1–12. [Online]. Available: <http://doi.acm.org/10.1145/2830772.2830773>
- [45] M. K. Qureshi, D. Thompson, and Y. N. Patt, "The v-way cache: demand-based associativity via global replacement," in *32nd International Symposium on Computer Architecture (ISCA'05)*, June

2005, pp. 544–555.

- [46] M. K. Qureshi, A. Jaleel, Y. N. Patt, S. C. Steely, and J. S. Emer, “Adaptive insertion policies for high performance caching,” in *Proceedings of the 34th Annual International Symposium on Computer Architecture*, ser. ISCA ’07. New York, NY, USA: ACM, 2007, pp. 381–391. [Online]. Available: <http://doi.acm.org/10.1145/1250662.1250709>
- [47] Richard Grisenthwaite, “Cache Speculation Side-channels,” January 2018.
- [48] D. Sanchez and C. Kozyrakis, “ZSim: Fast and accurate microarchitectural simulation of thousand-core systems,” in *Proceedings of the 40th Annual International Symposium on Computer Architecture-ISCA*, vol. 13. Association for Computing Machinery, 2013, pp. 23–27.
- [49] R. Sprabery, K. Evchenko, A. Raj, R. B. Bobba, S. Mohan, and R. H. Campbell, “A novel scheduling framework leveraging hardware cache partitioning for cache-side-channel elimination in clouds,” *CoRR*, vol. abs/1708.09538, 2017. [Online]. Available: <http://arxiv.org/abs/1708.09538>
- [50] P. Subramanyan, R. Sinha, I. Lebedev, S. Devadas, and S. A. Seshia, “A formal foundation for secure remote execution of enclaves,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS ’17, 2017, pp. 2435–2450.
- [51] G. Taylor, P. Davies, and M. Farmwald, “The TLB slice - a low-cost high-speed address translation mechanism,” *SIGARCH Computer Architecture News*, 1990.
- [52] L. Torvalds, “Re: Page colouring,” 2003. [Online]. Available: http://yarchive.net/comp/linux/cache_coloring.html
- [53] P. Turner, “Retpoline: a software construct for preventing branch-target-injection,” <https://support.google.com/faqs/answer/7625886>, January 2018.
- [54] C. A. Waldspurger, “Memory resource management in VMware ESX server,” in *Proceedings of the 5th Symposium on Operating Systems Design and Implementation*, ser. OSDI ’02. Berkeley, CA, USA: USENIX Association, 2002, pp. 181–194. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1060289.1060307>
- [55] Z. Wang and R. B. Lee, “New cache designs for thwarting software cache-based side channel attacks,” in *International Symposium on Computer Architecture (ISCA)*, 2007.
- [56] Y. Xu, W. Cui, and M. Peinado, “Controlled-channel attacks: Deterministic side channels for untrusted operating systems,” in *2015 IEEE Symposium on Security and Privacy*, May 2015, pp. 640–656.
- [57] F. Yao, M. Doroslovacki, and G. Venkataramani, “Are coherence protocol states vulnerable to information leakage?” in *2018 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, Feb 2018, pp. 168–179.
- [58] Y. Yarom and K. Falkner, “FLUSH+RELOAD: A high resolution, low noise, L3 cache side-channel attack.” in *USENIX Security Symposium*, 2014.
- [59] M. Zhang and K. Asanovic, “Victim replication: maximizing capacity while hiding wire delay in tiled chip multiprocessors,” in *32nd International Symposium on Computer Architecture (ISCA’05)*, June 2005, pp. 336–345.
- [60] X. Zhang, S. Dwarkadas, and K. Shen, “Towards practical page coloring-based multicore cache management,” in *Proceedings of the 4th ACM European Conference on Computer Systems*, ser. EuroSys ’09. New York, NY, USA: ACM, 2009, pp. 89–102. [Online]. Available: <http://doi.acm.org/10.1145/1519065.1519076>