



AutoLock: Why Cache Attacks on ARM Are Harder Than You Think

Marc Green, Worcester Polytechnic Institute; Leandro Rodrigues-Lima and Andreas Zankl, Fraunhofer AISEC; Gorka Irazoqui, Worcester Polytechnic Institute; Johann Heyszl, Fraunhofer AISEC; Thomas Eisenbarth, Worcester Polytechnic Institute

<https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/green>

**This paper is included in the Proceedings of the
26th USENIX Security Symposium
August 16–18, 2017 • Vancouver, BC, Canada**

ISBN 978-1-931971-40-9

**Open access to the Proceedings of the
26th USENIX Security Symposium
is sponsored by USENIX**

AutoLock: Why Cache Attacks on ARM Are Harder Than You Think

Marc Green
Worcester Polytechnic Institute

Leandro Rodrigues-Lima
Fraunhofer AISEC

Andreas Zankl
Fraunhofer AISEC

Gorka Irazoqui
Worcester Polytechnic Institute

Johann Heyszl
Fraunhofer AISEC

Thomas Eisenbarth
Worcester Polytechnic Institute

Abstract

Attacks on the microarchitecture of modern processors have become a practical threat to security and privacy in desktop and cloud computing. Recently, cache attacks have successfully been demonstrated on ARM based mobile devices, suggesting they are as vulnerable as their desktop or server counterparts. In this work, we show that previous literature might have left an overly pessimistic conclusion of ARM's security as we unveil AutoLock: an internal performance enhancement found in inclusive cache levels of ARM processors that adversely affects Evict+Time, Prime+Probe, and Evict+Reload attacks. AutoLock's presence on system-on-chips (SoCs) is not publicly documented, yet knowing that it is implemented is vital to correctly assess the risk of cache attacks. We therefore provide a detailed description of the feature and propose three ways to detect its presence on actual SoCs. We illustrate how AutoLock impedes cross-core cache evictions, but show that its effect can also be compensated in a practical attack. Our findings highlight the intricacies of cache attacks on ARM and suggest that a fair and comprehensive vulnerability assessment requires an in-depth understanding of ARM's cache architectures and rigorous testing across a broad range of ARM based devices.

1 Introduction

The rapid growth of mobile computing illustrates the continually increasing role of digital services in our daily lives. As more and more information is processed digitally, data privacy and security are of utmost importance. One of the threats known today aims directly at the fabric of digital computing. Attacks on processors and their microarchitecture exploit the very core that handles our data. In particular, processor caches have been exploited to retrieve sensitive information across logic boundaries established by operating systems and hypervisors. As

caches speed up the access to data and instructions, timing measurements allow an adversary to infer the activity of other applications and the data processed by them. In fact, cache attacks have been demonstrated in multiple scenarios in which our personal data is processed, e.g., web browsing [41] or cloud computing [26, 58]. These attacks have severe security implications, as they can recover sensitive information such as passwords, cryptographic keys, and private user behavior. The majority of attacks have been demonstrated on classic desktop and server hardware [25, 30, 42, 51], and with Intel's market share for server processors being over 98% [31], their platforms have been targeted most frequently.

With mobile usage skyrocketing, the feasibility of cache attacks on smartphone and IoT processors – which are predominantly ARM-based – has become a relevant issue. Attacks that rely on the existence of a cache flush instruction, i.e., Flush+Reload [51] and Flush+Flush [23], work efficiently across a broad range of x86 processors, but have limited applicability on ARM devices. In general, cache flush instructions serve the legitimate purpose of manually maintaining coherency, e.g., for memory mapped input-output or self-modifying code. On any x86 processor implementing the SSE2 instruction set extension, this flush instruction is available from all privilege levels as `clflush`. A similar instruction was introduced for ARM processors only in the most recent architecture version, ARMv8. In contrast to `clflush`, it must be specifically enabled to be accessible from userspace. This leaves a significant number of ARM processors without a cache flush instruction.

For all processors with a disabled flush instruction or an earlier architecture version, e.g., ARMv7, only eviction based cache attacks can be deployed. In particular, these attacks are Evict+Time [42], Prime+Probe [42], and Evict+Reload [24]. On multi-core systems, they target the last-level cache (LLC) to succeed regardless of which core a victim process is running on. This requires the LLC to be inclusive, i.e., to always contain the con-

tents of all core-private cache levels. On Intel processors, the entire cache hierarchy fulfills the inclusiveness property and is therefore a viable target for eviction based attacks. ARM devices, on the contrary, implement inclusive and non-inclusive caches alike. Both properties can co-exist, even in the same cache hierarchy. This renders eviction based attacks to be practicable only on a limited number of devices, in particular those that implement inclusive last-level caches. Yet, our findings show that an internal performance enhancement in inclusive last-level caches, dubbed `AutoLock`, can still impede eviction based cache attacks. In short, `AutoLock` prevents a processor core from evicting a cache line from an inclusive last-level cache, if said line is allocated in any of the other cores' private cache levels. This inhibits cross-core LLC evictions, a key requirement for practical `Evict+Time`, `Prime+Probe`, and `Evict+Reload` attacks on multi-core systems, and further limits the number of ARM based attack targets in practice.

In literature, Lipp et al. [37], Zhang et al. [54], and Zhang et al. [56] confirmed the general feasibility of flush and eviction based cache attacks from unprivileged code on ARM processors. Given the lack of flush instructions on a large selection of ARM devices and the deployment of non-inclusive LLCs or inclusive LLCs implementing `AutoLock`, the authors might have left an overly pessimistic conclusion of ARM's security against cache attacks. In addition, ARM's highly flexible licensing ecosystem creates a heterogeneous market of system-on-chips (SoCs) that can exhibit significant differences in their microarchitectural implementations. Demme et al. [17] illustrate that already small changes to the cache architecture can have considerable impact on the side-channel vulnerability of the processor. Consequently, judging the true impact of cache attacks on a broad range of ARM based platforms remains to be a challenge. Our work adds another step in this process. It is a contribution to an in-depth understanding of microarchitectural features on ARM in general and an extension to our current knowledge of cache implementations in particular.

1.1 Contribution

This work unveils `AutoLock`, an internal and undocumented performance enhancement feature found in inclusive cache levels on ARM processors. It prevents cross-core evictions of cache lines from inclusive last-level caches, if said lines are allocated in any of the core-private cache levels. Consequently, it has a direct and fundamentally adverse effect on all eviction based cache attacks launched from unprivileged code on multi-core systems. Understanding `AutoLock` and determining its existence on a given system-on-chip is vital to assess the SoC's vulnerability to those attacks. Yet, neither

technical reference manuals (TRMs) nor any other public product documentation by ARM mention `AutoLock`. We therefore provide a detailed description of the feature and propose three methodologies to test for it: using a hardware debugging probe, reading the performance monitoring unit (PMU), and conducting simple cache-timing measurements. Each test strategy has different requirements and reliability; having multiple of them is vital to test for `AutoLock` under any circumstances. With the proposed test suite, we verify `AutoLock` on ARM Cortex-A7, A15, A53, and A57 processors. As `AutoLock` is likely implemented on a larger number of ARM processors, we discuss its general implications and how our results relate to previous literature.

Despite its adverse effect on eviction based cache attacks, the impact of `AutoLock` can be reduced. We discuss generic circumvention strategies and execute the attack by Irazoqui et al. [29] in a practical cross-core `Evict+Reload` scenario on a Cortex-A15 implementing `AutoLock`. We successfully recover the secret key from a table based implementation of AES and show that attacks can tolerate `AutoLock` if multiple cache lines are exploitable. Furthermore, the presented circumvention strategies implicitly facilitate cross-core eviction based attacks also on non-inclusive caches. This is because in the context of cross-core LLC evictions, inclusive last-level caches with `AutoLock` behave identically to non-inclusive ones. In summary, our main contributions are:

- the disclosure and description of `AutoLock`, an undocumented and previously unknown cache implementation feature with adverse impact on practical eviction based cache attacks on ARM devices,
- a comprehensive test suite to determine the existence of `AutoLock` on actual devices, as its presence is not documented publicly,
- a discussion of `AutoLock`'s implications and its relation to previous literature demonstrating cache attacks on ARM, and
- a set of strategies to circumvent `AutoLock` together with a practical demonstration of a cross-core `Evict+Reload` attack on a multi-core SoC implementing `AutoLock`.

The rest of this paper is organized as follows. Section 2 describes `AutoLock`. A theoretical methodology to test for it is presented in Section 3. We evaluate SoCs for `AutoLock` in Section 4 and address how recent literature relates to it in Section 5. The implications of `AutoLock` are discussed in Section 6. Circumvention strategies together with a practical cross-core attack are presented in Section 7. We conclude in Section 8.

2 AutoLock: Transparent Lockdown of Cache Lines in Inclusive Cache Levels

Processor caches can be organized in levels that build up a hierarchy. Higher levels have small capacities and fast response times. *L1* typically refers to the highest level. In contrast, lower levels have increased capacities and response times. The lowest cache level is often referred to as the last-level cache, or *LLC*. Data and instructions brought into cache reside on cache *lines*, the smallest logical storage unit. In set-associative caches, lines are grouped into *sets* of fixed size. The number of lines or *ways* per cache set is called the *associativity* of the cache level. It can be different for every level. Whether a cache level can hold copies of cache lines stored in other levels is mainly defined by the inclusiveness property. If a cache level x is *inclusive* with respect to a higher level y , then all valid cache lines contained in y must also be contained in x . If x is *exclusive* with respect to y , valid lines in y must not be contained in x . If any combination is possible, the cache is said to be *non-inclusive*.

Inclusive caches enforce two rules. If a cache line is brought into a higher cache level, a copy of the line must be stored in the inclusive lower level. Determining whether a line is stored anywhere in the hierarchy can then be achieved by simply looking into the LLC. Vice versa, if a line is evicted from the lower level, any copy in the higher levels must subsequently be evicted as well. This is an implicit consequence of the inclusiveness property that has been successfully exploited in cross-core cache attacks that target inclusive LLCs [26, 27, 39].

Evictions in higher cache levels to maintain inclusiveness can add substantial performance penalties in practice. In a patent publication by Williamson and ARM Ltd., the authors propose a mechanism that protects a given line in an inclusive cache level from eviction, if any higher cache level holds a copy of the line [50]. An *indicator storage element* that is integrated into the inclusive cache level tracks which lines are stored in higher levels. The element can be realized with a set of *indicator* or *inclusion* bits per cache line, or a tag directory. If an indicator is set, the corresponding line is protected. This mechanism therefore prevents said performance penalties, because subsequent evictions in higher cache levels are prohibited. We refer to this transparent protection of cache lines in the LLC as *Automatic Lockdown* or *AutoLock*.

The impact of AutoLock during eviction is illustrated in Figure 1. For simplicity, the illustration is based on a two-level cache hierarchy: core-private L1 caches and a shared inclusive last-level cache (L2) with one inclusion bit per cache line. S and L are placeholders for any available cache set and line in the two cache levels. The left side of the figure shows how a cache line is evicted in L1.

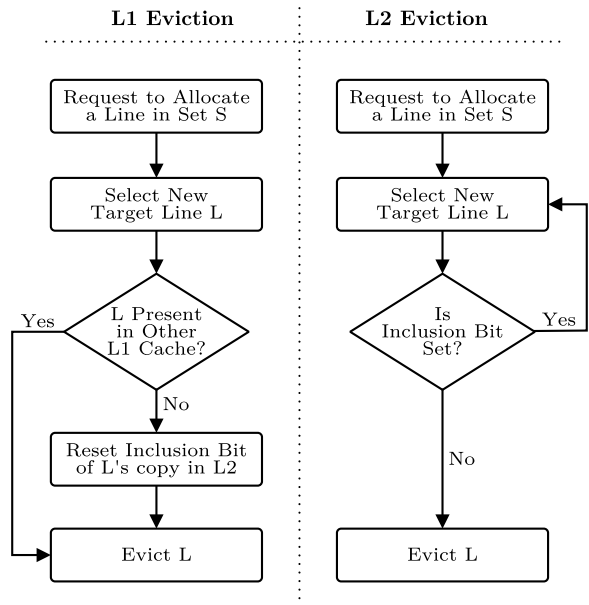


Figure 1: Simplified example of evicting a cache line from level 1 (left) and level 2 (right) cache sets. The level 2 cache is inclusive with respect to level 1 and implements AutoLock.

First, an allocation request for set S is received. Then, a target line L is selected by the replacement algorithm for eviction. If a copy of L is present in any other of the core-private L1 caches, it can immediately be evicted without updating the inclusion bit in L2. Because other L1 copies exist, the bit does not need to be changed. If no other L1 cache holds a copy of L , the inclusion bit must be reset in L2, which unlocks the copy of L in L2. After the bit is reset, L is evicted from L1.

Similarly, an allocation request in L2 triggers the replacement algorithm to select a line in set S for eviction. Before L is evicted in L2, its inclusion bit is checked. If a copy of L exists in any L1 cache, the replacement algorithm is called again to select another target line. This is repeated until one line is found whose inclusion bit is not set. This line is then evicted to allow the allocation of a new one.

If the number of ways in the inclusive lower cache level, W_l , is higher than the sum of ways in all higher cache levels, i.e., $W_{h,sum} = \sum_{i=1}^N W_{h,i}$, it can be guaranteed that at least one line is always selectable for eviction. If $W_l = W_{h,sum}$, all lines of a set in the lower cache level can be auto-locked. In this case, the patent proposes to fall back to the previous behavior, i.e., evict all copies of a line from higher level caches. This unlocks the line in the lower cache level and subsequently enables its eviction. If the number of ways in the lower cache level is further reduced, such that $W_l < W_{h,sum}$, ad-

ditional measures must be taken to implement inclusiveness and Automatic Lockdown. While not impossible per se, this case is not covered by the patent authors.

If an inclusive LLC with AutoLock is targeted in a cache attack, the adversary is not able to evict a target's data or instructions from the LLC, as long as they are contained in the target's core-private cache. In theory, the adversary can only succeed, if the scenario is reduced to a same-core attack. Then, it is possible once again to directly evict data and instructions from core-private caches. Note that the same attack limitation is encountered on systems with non-inclusive last-level caches, because cache lines are allowed to reside in higher levels without being stored in lower levels. In both cases, AutoLock and non-inclusive LLC, the attacks do not work cross-core because the attacking core cannot influence the target core's private cache. Note that it is possible, and indeed common on ARM processors, that there are separate L1 caches for instructions and data and that the LLC is inclusive with respect to one of them, but non-inclusive with respect to the other. In Section 7, we discuss possible strategies to circumvent AutoLock and re-enable cross-core cache attacks. Because of said similarities, those strategies also facilitate cross-core attacks on non-inclusive LLCs.

Distinct from Automatic Lockdown, there exists *programmable lockdown* in some ARM processors. Regardless of inclusiveness, it allows the user to explicitly lock and unlock cache lines by writing to registers of the cache controller. This has the same effect as AutoLock, i.e., the locked cache line will not be evicted until it is unlocked. In contrast, however, programmable lockdown must be actively requested by a (privileged) user. Of the four Cortex-A processors we study in this paper, the technical reference manuals do not mention programmable lockdown for any of them [5, 6, 8, 9]. AutoLock, however, is found in all of them.

3 How to Test for AutoLock

AutoLock is neither mentioned in ARM's architecture reference manuals [7, 10] nor in the technical reference manuals of the Cortex-A cores considered in this work [5, 6, 8, 9]. To the best of our knowledge, it is not publicly documented other than in patent form [50]. Based on official information, it is therefore impossible to determine which Cortex-A or thereto compliant processor cores implement AutoLock, let alone whether an actual system-on-chip features it. The presence of AutoLock, however, is crucial to assess the risk of cache attacks, in particular those that rely on cross-core evictions in the LLC. We therefore propose the following test methodology to determine the existence of AutoLock. On any device under test, two processes are spawned on

distinct cores of the processor implementing an inclusive last-level cache. The first process allocates a cache line in the private cache level of the core it is running on. This allocation is done with a simple memory access. The inclusiveness property ensures that a copy of the cache line must also be allocated in the last-level cache. The second process then tries to evict the line from the LLC by filling the corresponding cache set with dummy data. If the cache line remains in the LLC and core-private cache after the cross-core eviction, the test concludes that AutoLock is implemented. If the line is removed from both, the test concludes that AutoLock is not present.

Test Requirements and Intricacies. The proposed test strategy requires that the eviction itself works reliably, because otherwise AutoLock cannot be distinguished from a failed eviction and false positives are the consequence. We ensure a working eviction by verifying it in the same-core scenario before executing the AutoLock tests. Another requirement is that the LLC is inclusive, because AutoLock is defined only in the context of inclusive cache levels [50]. If the inclusiveness property is not documented for a processor, it must be determined experimentally. For these inclusiveness tests we recommend to refrain from using cross-core evictions, as AutoLock might interfere, which leads to wrong conclusions. Instead, hardware debugger or cache maintenance operations should be used to confidently determine whether caches are inclusive or not. If neither of those are available, cross-core evictions still allow to draw a conclusion, but only for a certain outcome. Assume the same scenario as in the AutoLock tests. The first process allocates memory on a cache line in its private cache level. The second process then fills the corresponding LLC set, after which the first process tries to re-access its memory. If the access is fulfilled from RAM, then it can be concluded that the LLC is inclusive and that AutoLock is not present. If the access is fulfilled from either private or last-level cache, no conclusion can be drawn, because either the LLC is not inclusive or AutoLock interfered. Once the inclusiveness property of the LLC is ensured, AutoLock can be tested as described in the following sections.

3.1 Cache Eviction

In order to evict a cache line from the LLC, we implement the method described by Gruss et al. [22] and Lipp et al. [37]. Assume that an address T is stored on line L in cache set S . In order to evict L from S , one has to access a number of addresses distinct from T that all map to S . These memory accesses fill up S and eventually remove L from the set. The addresses that are accessed in this process are said to be *set-congruent* to T . They are

collected in the eviction set C . Whenever C is accessed, T is forced out of the set S . The sequence of accesses to addresses in C is referred to as the *eviction strategy*. The strategy proposed by Gruss et al. [22] is shown in Algorithm 1.

Algorithm 1: Sliding window eviction of the form N-A-D [22, 37].

Input:
 C ... list of set-congruent addresses

```

1 for  $i = 0..N-1$  do           /* # of windows */
2   for  $j = 0..A-1$  do         /* # reps/window */
3     for  $k = 0..D-1$  do       /* # addrs/window */
4       |   access( $C[i+k]$ );
5       |   end
6     end
7 end
```

The idea is to always access a subset of the addresses in C for a number of repetitions, then replace one address in the subset with a new one, and repeat. This essentially yields a window that slides over all available set-congruent addresses. We therefore refer to this method as *sliding window eviction*. In the algorithm, N denotes the total number of generated windows, A defines the repetitions per window, and D denotes the number of addresses per window. The required size of C is given by the sum of N and D . The final eviction strategy is then written as the triple N-A-D. The strategy 23-4-2, for example, comprises 23 total windows, each iterated 4 consecutive times and containing 2 addresses. Lipp et al. [37] demonstrate that sliding window eviction can successfully be applied to ARM processors.

The parameters N-A-D must be determined once for each processor. This is done by creating a list of set-congruent addresses C and exhaustively iterating over multiple choices of N , A , and D . By continuously checking the success of the eviction, the strategy with the least number of memory accesses that still provides reliable eviction can be determined. Generating the list of set-congruent addresses C requires access to physical address information. This is because the last-level caches on our test devices use bits of the physical address as the index to the cache sets. If the parameter search for N , A , and D is done in a bare-metal setting, physical address information is directly available. Operating systems typically employ virtual addresses that must be translated to physical ones. Applications on Linux, for instance, can consult the file `/proc/[pid]/pagemap` to translate virtual addresses [37]. Although accessing the pagemap is efficient, access to it can be limited to privileged code or deactivated permanently. Alternatively, huge pages reveal sufficient bits of the physical address

to derive the corresponding cache set [27]. To find addresses set-congruent to T , new memory is allocated and the containing addresses are compared to T . If the least significant address bits match while the most significant bits differ, the address will map to T 's cache set but will be placed on a different line within the set. If access to physical address information is entirely prohibited, timing measurements can still be used to find set-congruent addresses [41].

Once C is filled with addresses, they are accessed according to Algorithm 1. If a processor core implements separated data and instruction caches, the manner in which a set-congruent address ought to be accessed differs. Data addresses can be accessed by loading their content to a register with the LDR assembly instruction. Instruction addresses can be accessed by executing a *branch* instruction that jumps to it. When determining N-A-D on devices that might implement AutoLock, all memory accesses to T and all evictions of it must be performed on a single core. This ensures that AutoLock is not interfering with the parameter search. Once the eviction with a triple N-A-D works reliably in the same-core setting, the AutoLock tests can be commenced.

3.2 AutoLock Tests

In the subsequent sections we propose three tests that have been designed to prove or disprove the existence of AutoLock. All of them follow the general methodology of determining the success of a cross-core eviction strategy that is known to succeed in the same-core scenario. For simplicity, all tests are explained in a dual-core setting. For a system with more processor cores, each test can either be repeated multiple times or extended in order to determine the presence of AutoLock simultaneously on all but one core. In the dual-core setting, core 0 is accessing the target address T and core 1 is trying to evict it by using the eviction set C and the processor specific eviction parameters N-A-D. Table 1 contains the parameters for the processors considered in this work. For all tests, both T and C are listed as inputs required for eviction, while the triple N-A-D is assumed to be correctly set according to Table 1. As currently nothing indicates

Table 1: List of ARM and ARM-compliant processors under test, including the number of inclusive L1 and L2 ways, as well as the eviction strategy parameters N-A-D.

Processor	L1 Ways	L2 Ways	N-A-D
Cortex-A7	2 (Instr.)	8	23-4-2
Cortex-A15	2 (Data)	16	36-6-2
Cortex-A53	2 (Instr.)	16	25-2-6
Cortex-A57	2 (Data)	16	30-4-6
Krait 450 ¹	4 (Data)	8	50-1-1

that AutoLock can be en- or disabled from software, its presence on a processor has to be determined only once.

3.2.1 Hardware Debugger

The first method to test for AutoLock is through the usage of a hardware debugger. It allows to halt a processor at will and directly monitor the cache content. A breakpoint is inserted after the eviction of T and the contents of the caches are analyzed. Through this visual inspection, it is possible to determine with very high confidence whether or not T remains in the cache after the eviction strategy is run. Given an inclusive LLC, it is sufficient to confirm that T either remains in L1 or in L2 to prove that AutoLock is present. Algorithm 2 outlines this test.

Algorithm 2: Hardware Debugger Test

Input:

T ... target address
 C ... corresponding eviction set

- 1 Core 0 brings T to L1 and L2.
 - 2 Core 1 runs eviction strategy using C .
 - 3 Halt processor and inspect caches.
 - 4 **if** T in L1 of Core 0 or L2 **then** AutoLock is present
 - 5 **else** AutoLock is not present
-

The hardware debugger test requires a debugging unit, a target platform that supports it, and physical access to the target device. In our experiments, we use the DSTREAM debugging unit [11] in combination with the ARM DS-5 development studio. Of course, the test can also be run with other debugging hardware.

3.2.2 Performance Monitoring Unit (PMU)

The second test utilizes the performance monitoring unit, which can count the occurrence of microarchitectural events in a processor. The PMU of ARMv7- and ARMv8-compliant processors can be configured to count the number of accesses (hit or miss) to the last-level cache. The corresponding event is defined under the ID 0x16 in the architectural reference manuals [7, 10]. The difference of the access counts before and after reloading the target address T indicates whether the reload fetched the address from the L1 or the L2 cache. A fetch from L1 indicates that the eviction strategy failed and suggests that AutoLock is implemented. If the eviction strategy is successful, the target address has to be fetched from external memory. Before querying the slow external memory, the L2 cache is accessed and checked for the target address. This access is counted and indicates that AutoLock is not implemented. To determine this extra access to the L2, a reference value R must be obtained before the test. This is done by reading the L2 access

counter for a reload with no previous run of the eviction strategy, which guarantees a fetch from core-private cache. The PMU test can be conducted as outlined in Algorithm 3.

Algorithm 3: PMU Test

Input:

T ... target address
 C ... corresponding eviction set
 R ... L2 access reference count

- 1 Core 0 brings T to L1 and L2.
 - 2 Core 1 runs eviction strategy using C .
 - 3 Save PMU count of L2 accesses.
 - 4 Core 0 reloads T .
 - 5 Save PMU count again and calculate difference d .
 - 6 **if** $d \approx R$ **then** AutoLock is present
 - 7 **else** AutoLock is not present
-

This test requires access to the PMU, which on ARM is typically limited to privileged code, unless otherwise configured. Some operating systems, however, allow userspace applications to access hardware performance events. On Linux, for instance, the perf subsystem of the kernel provides this access via the perf_event_open system call [36]. In general, the PMU test can be used when the target processor is not supported by a hardware debugger, or if physical access to the device is not given. Since PMU event counts can be affected by system activity unrelated to the AutoLock test, it is recommended to repeat the experiment multiple times. The best results can be obtained in a bare-metal setting, where the test code is executed without a full-scale operating system.

3.2.3 Cache-timing Measurements

The third experiment uses timing measurements to infer from where in the memory hierarchy the target address T is reloaded after the supposed eviction. If external memory access times are known, the reload time of T can indicate whether AutoLock is implemented or not. This test approach is outlined in Algorithm 4.

Algorithm 4: Cache-timing Test

Input:

T ... target address
 C ... corresponding eviction set
 M ... external memory access time

- 1 Core 0 brings T to L1 and L2.
 - 2 Core 1 runs eviction strategy using C .
 - 3 Core 0 reloads T and measures reload time t .
 - 4 **if** $t < M$ **then** AutoLock is present
 - 5 **else** AutoLock is not present
-

If after the eviction strategy the reload time is smaller than what is expected for an external memory access, the target address is likely fetched from cache, thus indicating AutoLock. If the reload time is equal to an external memory access, the eviction strategy was successful and AutoLock is likely not present.

This test has no further requirements other than running code on the system from userspace and having access to a sufficiently accurate timing source. Commonly used timing sources include hardware based time-stamp counters (PMCCNTR for ARM), the perf subsystem of Linux [36], the POSIX `clock_gettime()` function [45], and a custom thread based timer. If available, a hardware based time-stamp counter is preferred due to its high precision. Further discussions about timing sources can be found in the work by Lipp et al. [37] and Zhang et al. [56]. Similar to PMU event counts, timing measurements can significantly be affected by noise. It is therefore advisable to repeat the proposed test multiple times to get a robust conclusion about whether address T is fetched from cache or external memory. Due to the versatility of this test, we recommend its use in situations where either adequate debugging equipment is not available or the abilities to conduct the other, more robust experiments are not given (e.g., when root access on a device cannot be gained due to vendor restrictions).

4 Finding AutoLock in Existing SoCs

In this work, we evaluate the presence of AutoLock on four test devices and their corresponding system-on-chips. They are illustrated in Table 2. The Samsung Exynos 5422 and the ARM Juno r0 SoCs feature two processors with multiple cores each. They are so-called ARM big.LITTLE platforms, on which a powerful processor is paired with an energy efficient one. Together with the Samsung Exynos 5250, these SoCs are part of dedicated development boards or single-board computers. In contrast, the Qualcomm Snapdragon 805 is part of an off-the-shelf mobile phone. In total, the four test devices comprise five different processors: the ARM Cortex-A7, A15, A53, A57, and the Qualcomm Krait 450. Table 1 provides details about their cache hierarchies. It shows the number of ways in L1 and L2 caches, and the eviction strategy parameters for all of them. The illustrated processors implement separate L1 instruction and data caches. The number of L1 ways is given only for the side which the L2 cache is inclusive to. The LLCs on the Cortex-A7 and A53 are inclusive to the L1 instruction caches, while the LLCs on the Cortex-A15, A57, and the Krait 450 are inclusive to the L1 data caches. The inclusiveness properties of the A15 and A57 are explicitly stated in their respective reference manuals [5, 8] (Section *Level 2 Memory System*). The A7 and

Table 2: Platforms used for the evaluation of AutoLock. For each device, the corresponding SoC and processor cores are given.

Device	System-on-Chip	Core(s)
Arndale	Samsung Exynos 5250	2x Cortex-A15
ODROID XU4	Samsung Exynos 5422	4x Cortex-A7 4x Cortex-A15
ARM Juno	ARM Juno r0	4x Cortex-A53 2x Cortex-A57
Nexus 6	Qualcomm Snapdragon 805	4x Krait 450

A53 manuals imply inclusiveness on the instruction side, but do not explicitly state it [6, 9] (e.g. in Section *Optional integrated L2 cache*). For the A53, however, the lead architect confirmed it in an interview [47]. Public documentation of the Krait 450 is scarce and information about cache inclusiveness could only be obtained for earlier Krait generations [34]. We therefore infer its inclusiveness from successful cross-core eviction experiments that at the same time disprove the existence of AutoLock.

The tests on the Cortex-A processors are initially done in a bare-metal setting. The lack of an operating system eliminates interfering cache activity from system processes and significantly reduces noise. The experiments are then repeated on Linux for verification. On the Krait 450, the experiments are conducted on Linux only. For each processor, we verify that the eviction parameters listed in Table 1 can successfully evict cache lines in a same-core setting. More precisely, we verify successful eviction when evicting data cache lines using data addresses, and when evicting instruction cache lines using instruction addresses. We then test for AutoLock in the cross-core case with the experiments proposed in the previous section.

4.1 Test Results

The subsequent sections present the results for all test methodologies described in Section 3. Along with the conclusions about the presence of AutoLock on the test devices, details about the practical execution of the experiments are discussed.

4.1.1 Hardware Debugger

The SoCs on the ARM Juno and the Arndale development boards are the only ones among the test devices that are supported by DSTREAM². It is therefore possible to visually inspect the L1 caches of the Cortex-A15, A53, and A57 processors, and the L2 caches of the A15 and A57. A hardware limitation of the Cortex-A53 in the ARM Juno r0 SoC prevents the visual inspection of its L2 cache. To still test for AutoLock on the A53, we

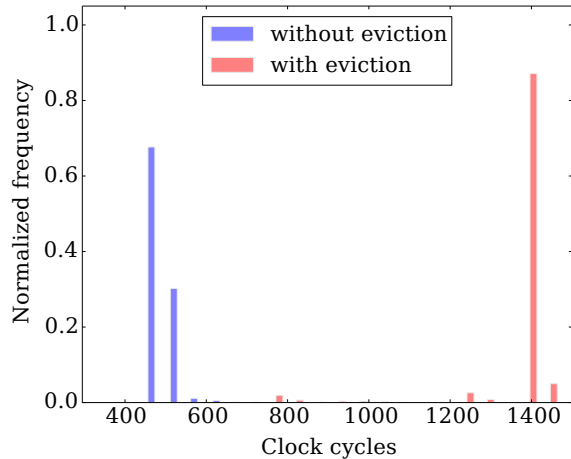


Figure 2: Memory access times with and without cross-core eviction on the Krait 450 processor. A threshold of 700 clock cycles clearly separates the two timing distributions, which indicates that AutoLock is not present.

leverage the inclusiveness property to surmise L2 contents. According to Algorithm 2, AutoLock can still be concluded, if the target address is contained in the core-private cache of core 0. This is derived from the inclusiveness property of the L2 cache.

To conduct the tests, we connect each supported board, in turn, to the DSTREAM and use breakpoints to temporarily halt program execution after the eviction algorithm is run. When halted, we use the Cache View of the DS-5 development studio to visually determine if the target cache line is present in the respective caches. For the A53, we infer the contents of the L2 based on the inclusive L1. We ran the experiments several times on the A15, A53, and A57 processors. All trials indicate each processor’s inclusive cache implements AutoLock.

4.1.2 Performance Monitoring Unit (PMU)

To verify to results of the Cortex-A53, we conduct the experiment described in Algorithm 3 with it. The PMU is configured to count accesses to the L2 cache. We then execute a target instruction on core 0 and run a 25-2-6 eviction strategy on core 1. Before and after reloading the target instruction, we insert 10 NOP instructions. This reduces the effects of pipelining, as the A53 has an 8-stage pipeline. To ensure we only measure exactly the reload of the target instruction, we execute a DSB and ISB instruction before each set of NOPs. These instructions function as *memory barriers*, guaranteeing that memory access instructions will execute sequentially. This is necessary because the ARM architecture allows memory accesses to be reordered to optimize performance.

As a result, we observe that reloading the target in-

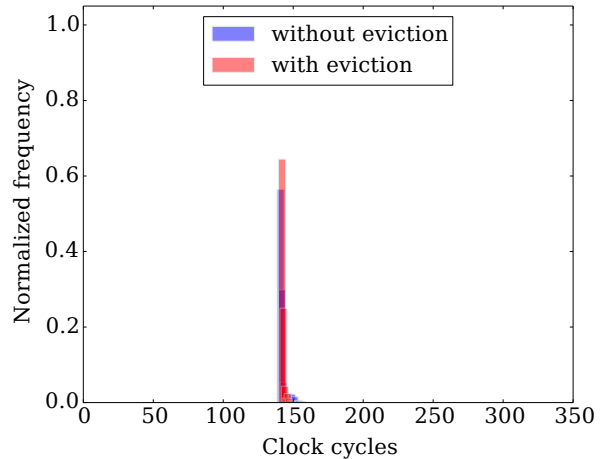


Figure 3: Memory access times with and without cross-core eviction on the ARM Cortex-A57 processor. The two timing distributions clearly overlap, which indicates that AutoLock inhibits the eviction.

struction after executing the eviction algorithm causes no additional L2 access. This indicates that the eviction failed and the reload was served from L1. If the eviction had succeeded, the event counter would have been incremented by the L2 cache miss. We ran the experiment multiple times and observed consistent results in each trial. This confirms the presence of AutoLock on the Cortex-A53.

4.1.3 Cache-timing Measurements

To determine the presence of AutoLock on the Cortex-A7 and the Krait 450, we execute the cache-timing experiment described in Algorithm 4. In addition, timing measurements are used to verify the results obtained for the Cortex-A15, A53, and A57. In all experiments, the perf subsystem of the Linux kernel, accessible from userspace, is used to measure access times with a hardware based clock cycle counter. This was taken from Lipp et al. [37]³.

Figure 2 shows the timing data collected with two separate executions of Algorithm 4 on the Krait 450. The first execution performed the eviction step as defined in the algorithm. The timings measured during the reload phase are shown in red. The second execution skipped the eviction step. These timings are shown in blue, for comparison. When the eviction step is skipped, the target address remains in the cache, and thus the access time is significantly lower, as pictured. This indicates that the Krait 450 does not implement AutoLock.

The same timing measurements are performed on the Cortex-A57. The results are shown in Figure 3. Since the timings for each execution virtually overlap in the

graph, it is clear that the target address is never evicted. This confirms the results for the A57 derived with the hardware debugger, i.e., that it implements AutoLock. For both the Krait 450 and the A57, we performed 50,000 measurements to ensure that clear trends can be seen.

Corresponding experiments on the Cortex-A7 indicate that its instruction cache side implements AutoLock. The measurements on the Cortex-A15 and A53 processors confirm the previous results and indicate once more that they implement AutoLock.

4.2 Discussion of the Test Results

The summary of the test results is shown in Table 3. All ARM Cortex-A processors on our test devices exhibit AutoLock in their inclusive last-level caches, whereas no evidence of AutoLock can be found on the Qualcomm Krait 450. The practical impact of AutoLock is that cache lines in the LLC are transparently locked during runtime. On a multi-core system, this can be triggered simultaneously in multiple cores. For each core, AutoLock essentially reduces the number of lines per LLC set that are available to store new data and instructions. Depending on the associativity of the core-private cache levels the LLC is inclusive to, a significant fraction of cache lines can be locked in an LLC set. Referring to Table 1, it can be seen that the Cortex-A7 features 2-way L1 instruction caches and an inclusive 8-way L2 cache. This means that on a quad-core A7 it is possible to lock all ways of an L2 cache set with instructions held in the core-private caches. Requests to store lines from L1 data caches in such a set subsequently fail, but do not violate inclusiveness, as the L2 is non-inclusive to L1 data caches. On the other Cortex-A processors, L2 cache sets cannot fully be locked. In quad-core Cortex-A15, A53, and A57 processors, up to 8 ways can be locked in an L2 set at once. As the Krait 450 does not implement AutoLock, the number of ways in the L2 does not need to match the sum of ways in the L1 caches. Hence, the L1 data caches contain 4 ways while the L2 contains 8.

5 Related Work

The field of cache attacks finds its origins in the early anticipation of varying memory access times compromising the security of cryptographic software [32, 35]. Ever since, this field has seen significant development in both attack and defense strategies. Tsunoo [48] and Bernstein [13] first introduced practical attacks based on the varying execution time of block ciphers. As the targeted implementations performed key-dependent memory accesses that resulted in key-dependent numbers of cache hits and misses, the execution time contained sufficient information to recover (parts of) the se-

Table 3: Evaluation results for the ARM and ARM-compliant processors of the test devices.

Processor	System-on-Chip	AutoLock
Cortex-A7	Samsung Exynos 5422	Present
Cortex-A15	Samsung Exynos 5250/5422	Present
Cortex-A53	ARM Juno r0	Present
Cortex-A57	ARM Juno r0	Present
Krait 450	Qualcomm Snapdragon 805	Not Present

cret key. The works of Tsunoo and Bernstein belong to the group of *time-driven* attacks that exploit the link between overall processing time and cache activity. Similar attacks have subsequently been demonstrated by Bonnaeu and Mironov [15], Açıçmez et al. [3], and Bogdanov et al. [14]. An alternative attack vector is introduced in the work of Page [43], where the sequence of cache hits and misses is observed during the execution of the cipher. Knowing how the key is involved in the memory accesses, this sequence (or *trace*) allows to infer bits of the key, which is the basis of so-called *trace-driven* attacks that have been studied further by Açıçmez and Koç [2], Fournier and Tunstall [19], and Gallais and Kizhvatov [20]. While AutoLock can in theory affect both time- and trace-driven attacks, we consider its practical impact to be limited. First, both attack types typically do not require active and fine-grained manipulation of the cache. Instead, attacks rely on background activity on the target system, limited sizes of cache levels, self-evictions, and the fact that unused data will typically be removed from cache after some time. These building blocks are hardly affected by AutoLock. Second, attacks often employ statistical analyses and thereby implicitly account for erroneous observations. As such, AutoLock will most likely act as an additional source of noise. The situation changes for attacks that exploit individual accesses to the cache. These so-called *access-driven* attacks have proven to be considerable threats in practice and consequently gained significant attention in literature over the past decade. What access-driven attacks have in common is that data or instructions are selectively removed from the cache hierarchy at some point during the attack. While a large number of attack papers have been published, only a few basic attack strategies exist that can be categorized depending on how this removal is implemented.

Flush-based Attacks. The first group of attacks relies on *cache flushing*, i.e., the removal of cache contents with dedicated flush instructions that are part of the processor’s instruction set. Gullasch et al. [25] introduced a flush based attack on x86 processors and used it to retrieve an AES key from a core co-resident victim by abusing the flush instruction `clflush` as well as memory deduplication and the Completely Fair Scheduler of

Linux. Yarom and Falkner [51] extended the work by Gullasch et al. and proposed the Flush+Reload attack, with which they recovered RSA keys across processor cores and virtual machines. This work was expanded by Irazoqui et al. [29, 30], who demonstrated the recovery of AES keys and TLS session messages. The Flush+Reload technique was concurrently used by Bengier et al. [12] to recover ECC secret keys, by Zhang et al. [58] to launch cross-tenant attacks on PaaS clouds, and by Gruss et al. [24] to implement template attacks. Based on the timing variations of `clflush`, Gruss et al. [23] also proposed the Flush+Flush attack. We currently believe that none of the flush based attacks are affected by AutoLock, because cache maintenance operations such as cache flushing seem to override AutoLock.

Eviction-based Attacks. The second group of attacks relies on *cache eviction*, i.e., the removal of cache contents by filling the corresponding parts of the cache with dummy or unrelated content. Osvik et al. [42] introduced two eviction based attacks, called `Evict+Time` and `Prime+Probe`, which both target software implementations of AES. As the underlying attack strategies are more generic, Aciçmez [1] and Percival [44] utilized `Prime+Probe` to steal an RSA key while Neve and Seifert used it to perform an efficient last round attack on AES [40]. Ristenpart et al. [46] used the same technique to recover keystrokes from co-resident virtual machines (VMs) in the Amazon EC2 cloud. This work was later expanded by Zhang et al. [57], proving the effectiveness of `Prime+Probe` to recover El Gamal keys across VMs. Liu et al. [39] and Irazoqui et al. [27] showed the feasibility of `Prime+Probe` via the last-level cache. Both studies opened a range of scenarios in which cache attacks could be applied. For instance, Oren et al. [41] executed `Prime+Probe` in JavaScript and Inci et al. [26] demonstrated the applicability of the technique in commercial IaaS clouds. As a cross-over to the flush based attacks, Gruss et al. [24] proposed the `Evict+Reload` attack, which removes cache contents through evictions but other than that remains identical to the Flush+Reload technique. All of the eviction based attacks are inherently affected by AutoLock, if they are executed in a cross-core scenario. While this might not hold for the same-core attacks proposed in early literature, it very much applies to the most recent attacks that pose greater threats in practice.

Cache Internals. Besides the generic structure and behavior of caches, literature has also exploited more implementation specific aspects. Irazoqui et al. [28] demonstrated the applicability of cache attacks across CPUs through the cache coherency protocol, which allows the execution of Flush+Reload style attacks by

forwarding cache flush and data request messages between two CPU clusters. Yarom et al. [52] showed that cache bank contentions introduce timing variations of accesses to different words on a single cache line. This undermines the general assumption in other work that different words on one cache line have the same timing behavior. Gruss et al. [21] introduced prefetching instructions as a way to load memory into cache without explicitly accessing it. This can be used to circumvent supervisor mode access prevention and address space layout randomization. As the exact impact of AutoLock on these attacks is unclear, we leave a more detailed assessment of AutoLock's intricacies in these cases to future work.

Countermeasures. The threat posed by cache attacks has been addressed from both hardware and software side in the form of countermeasures. Hardware based approaches often leverage programmable lockdown, i.e., the ability to actively lock cache lines. Cache contents belonging to sensitive applications are then locked and therefore protected from eviction by other, e.g., malicious applications. In literature, this programmable lockdown is for instance used by Wang and Lee [49] and by Liu et al. [38] to counter cache attacks. Although AutoLock also prevents cache lines from being evicted and therefore behaves in a similar way, there are two main differences compared to the proposed countermeasures. First, there is no means of controlling AutoLock, as it can neither be configured nor disabled. Second, AutoLock only protects lines in the LLC under certain conditions, i.e., if they are kept in core-private cache levels. Once these lines are removed from core-private levels, the protection immediately stops. As such, AutoLock cannot provide any guarantee to impede cache attacks. It is rather an additional layer of complexity that an adversary has to overcome during a cache attack. Software based countermeasures often try to separate the cache footprints of applications, such that each application gets a separate portion of the cache, which prevents interferences at the cache level. In literature, this strategy has for instance been applied by Kim et al. [33] and Zhou et al. [59]. Ultimately, applications that handle sensitive data should be designed such that both execution flow and memory accesses are independent of any sensitive input that is processed. As this is not a trivial task, several tools have been proposed that help to detect cache leaks and fix vulnerable code [4, 18, 53]. If applications cannot be fixed, other approaches try to detect and stop cache attacks in real-time, before any harm is caused [55]. For these system- and application-level countermeasures, we do not expect any particular impact from AutoLock.

Most of the previous cache attack literature is dedicated to the x86 architecture. Recent works [37, 54, 56] have

Table 4: Utility of known cache attacks in different scenarios on ARM Cortex-A processors with inclusive caches implementing AutoLock. ‘✓’ indicates the attack is unaffected by AutoLock, while ‘✗’ denotes obstruction by AutoLock. Flush+Reload and Flush+Flush are uninhibited by AutoLock but only apply to a limited number of ARMv8-A SoCs and are thus listed as ‘*’.

Attack	Same-core	Cross-core	Cross-CPU
Evict + Time [42]	✓	✗	✗
Prime + Probe [42]	✓	✗	✗
Flush + Reload [51]	*	*	*
Evict + Reload [24]	✓	✗	✗
Flush + Flush [23]	*	*	*

made several contributions to overcome the challenges of applying known userspace cache attacks from x86 to ARM processors. AutoLock is not recognized or mentioned in any of them. The following section is dedicated to discuss how AutoLock relates to these publications and why it might have stayed undetected.

5.1 AutoLock in Previous Work

Lipp et al. [37] were the first to demonstrate the feasibility of Prime+Probe, Flush+Reload, Evict+Reload, and Flush+Flush attacks on ARM devices. Despite their extensive experiments, the authors did not mention any encounter of a feature similar to AutoLock. We believe this can mainly be explained with their selection of test devices: the OnePlus One, the Alcatel One Touch Pop 2, and the Samsung Galaxy S6. Respectively, these mobile phones feature the Krait 400, the Cortex-A53, and a big.LITTLE configuration of the Cortex-A53 and Cortex-A57. We assume that the Krait 400, like the Krait 450 we experimented on, does not feature AutoLock. The Samsung Galaxy S6 features a full-hierarchy flush instruction that is available from userspace by default. Thus, the authors bypassed AutoLock on this device by flushing cache lines instead of evicting them. In contrast, the Alcatel One Touch Pop 2 features a Cortex-A53 that would potentially be affected by AutoLock according to our results. Yet, Lipp et al. successfully demonstrate a covert channel based on cross-core evictions on this chip. A possible explanation could be that the SoC manufacturers are different. While Lipp et al. experiment on a Qualcomm Snapdragon 410, we perform our tests on an ARM built Juno SoC. This could mean that the existence of AutoLock is yet more fragmented than our results might indicate. If this should be true, testing for AutoLock on a specific device is more important than ever. Another explanation is that the authors relied on evictions caused by background activity or by the measurement program itself (self-evictions).

Zhang et al. [56] implemented a variant of the Flush+Reload attack in a zero-permission Android application. The authors also experimented on processors we expect to feature AutoLock, but they did not mention any encounter with it either. In fact, they used the same processors analyzed in this work, namely, the Cortex-A7, A15, A53, A57, and the Krait 450. Since their work was focused solely on Flush+Reload, one of the two cache attacks unaffected by AutoLock, we assume they never encountered it during their experiments. One of the main contributions of their work was to implement an *instruction-side* Reload in a return-oriented manner, i.e., by executing small blocks of instructions. This contribution stemmed from using the `cacheflush` syscall in the Flush step, as it only invalidates the instruction side. As a prerequisite for their final target device selection, the authors experimentally determined the inclusiveness property of the last-level caches on all devices. Surprisingly, they concluded that all of the L2 caches in the aforementioned processors are inclusive with respect to the L1 data and instruction caches. This contradicts our experiments, which found the Cortex-A7 and A53 to only be inclusive on the instruction side, and the Cortex-A15 and A57 to only be inclusive on the data side. Further, the official ARM documentation of the Cortex-A7, for example, explains that “*Data is only allocated to the L2 cache when evicted from the L1 memory system, not when first fetched from the system.*” [6]. We understand this to mean that the L2 cache of the Cortex-A7 is *not* inclusive with respect to the L1 data caches. This complies with our observations.

In other previous work, Zhang et al. [54] implemented a Prime+Probe attack in an unprivileged Android application on an ARM Cortex-A8. Since we did not conduct experiments on this processor model, it is unclear whether AutoLock is implemented on it. However, the test system used by Zhang et al. comprised only a single Cortex-A8 processor core. As AutoLock does not affect same-core attacks, the experiments of the authors would not have been affected, even if the Cortex-A8 implemented AutoLock.

6 Implications of AutoLock

Automatic Lockdown prevents the eviction of cache lines from inclusive cache levels, if copies of that line are contained in any of the caches said level is inclusive to. Yet, the ability to evict data and instructions from a target’s cache is a key requirement for practical cross-core cache attacks. Table 4 illustrates the impact of AutoLock on state of the art cache attacks implemented on ARM Cortex-A processors. Each row shows one attack technique and the corresponding effect of AutoLock in three different scenarios: same-

Table 5: Selection of recent smartphones from manufacturers with high global market share [16]. For each listed device, the SoC, the contained processing cores, and their assumed exposure to AutoLock are given. A ‘✓’ indicates that AutoLock is currently expected to be present, while a ‘-’ states unknown exposure.

Smartphone	SoC	Core(s)	AutoLock Assumed
Apple iPhone 6	Apple A8	Typhoon	-
Apple iPhone 6s (Plus)	Apple A9	Twister	-
Apple iPhone 7 (Plus)	Apple A10	Hurricane, Zephyr	-
Huawei P9	Kirin 955	A72, A53	-(A72), ✓(A53)
Huawei P10 (Plus)	Kirin 960	A73, A53	-(A73), ✓(A53)
Huawei P10 Lite	Kirin 658	A53	✓
Huawei Y7	Snapdragon 435	A53	✓
LG Harmony	Snapdragon 425	A53	✓
LG V20	Snapdragon 820	Kryo	-
LG G6	Snapdragon 821	Kryo	-
Oppo A77	Mediatek MT6750T	A53	✓
Oppo R9s	Snapdragon 625	A53	✓
Oppo R9s Plus	Snapdragon 653	A72, A53	-(A72), ✓(A53)
Oppo R11 (Plus)	Snapdragon 660	Kryo	✓
Samsung Galaxy S6 (Edge)	Exynos 7420	A57, A53	✓(A57), ✓(A53)
Samsung Galaxy S7	Exynos 8890	M1, A53	-(M1), ✓(A53)
Samsung Galaxy S8 ^b	Exynos 8895	M2, A53	-(M2), ✓(A53)
Samsung Galaxy S7 Edge ^a	Snapdragon 820	Kryo	-
Samsung Galaxy S8+ ^b	Snapdragon 835	Kryo	-
vivo V5	Mediatek MT6750	A53	✓
vivo V5 Plus	Snapdragon 625	A53	✓
vivo Y55s	Snapdragon 425	A53	✓
Xiaomi Mi Max 2	Snapdragon 625	A53	✓
Xiaomi Mi6	Snapdragon 835	Kryo	-
Xiaomi Mi 5c	Surge S1	A53	✓

^a ... An alternative edition of the S7 Edge features an Exynos 8890.

^b ... The S8(+) can both feature either an Exynos 8895 or a Snapdragon 835.

core, cross-core, and cross-CPU. A ‘✓’ or ‘*’ signifies that an attack can be mounted, whereas a ‘✗’ indicates that AutoLock fundamentally interferes with the attack. Given the nature of AutoLock, all same-core attacks remain possible, as the adversary can evict target memory from all core-private cache levels. All attacks based on a full-hierarchy flush instruction are also not affected by AutoLock. However, said flush instruction, unlike on x86 processors, is not available on any ARMv7-A compliant processor and must be enabled in control registers on ARMv8-A compliant processors. Access to these control registers is limited to privileged, i.e., kernel or hypervisor code. These flush based attacks, namely Flush+Reload and Flush+Flush, are therefore denoted with ‘*’. In contrast to these attacks, all techniques that rely on evicting a cache line, namely Evict+Time, Prime+Probe, and Evict+Reload, are impaired by AutoLock in cross-core scenarios.

While it’s possible to state AutoLock’s fundamental impact on state of the art attack techniques, its effect and presence on actual devices is more difficult to assess. Based on our experiments, we currently assume that AutoLock is primarily implemented on Cortex-A cores designed by ARM itself. As the concept is protected by US patent [50], ARMv7-A and ARMv8-A compliant cores, such as Qualcomm’s Krait 450, would have to pay royalties to implement AutoLock. We therefore assume that compliant cores refrain from implementing

AutoLock. Based on these assumptions, Table 5 tries to illustrate the impact of AutoLock on recent smartphones. It lists devices from major manufacturers that have high global market shares [16]. For each of them, we select smartphones that have recently been introduced or announced. For every device, the corresponding SoC and processor cores are shown. In addition, the table states whether we expect AutoLock to be implemented (indicated by a ‘✓’) or whether a device’s exposure remains unknown (indicated by a ‘-’). A significant fraction of devices shown in Table 5 feature an ARM Cortex-A53, which we found to implement AutoLock. While newer cores such as the Cortex-A72 and A73 might be affected as well, it remains unclear whether this also holds for the ARM-compliant cores, such as the Kryo (the successor of the Krait), the Mongoose (M1, M2), as well as the cores integrated into the Apple SoCs.

If a device implements AutoLock, adversaries must find and employ circumvention strategies to leverage the full potential of eviction based cache attacks. In general, such strategies can also be used to target non-inclusive LLCs, where cross-core evictions are not possible, either. In the upcoming section, we discuss circumvention strategies and demonstrate that the attack proposed by Irazoqui et al. [29] can still be mounted in a cross-core Evict+Reload scenario with an inclusive LLC implementing AutoLock.

7 Circumventing AutoLock

Despite the restrictions that Automatic Lockdown poses to eviction based cache attacks, its effects can be alleviated with the following strategies:

- **Pre-select Target SoCs:** Our findings suggest that AutoLock is present on Cortex-A cores designed by ARM itself, while it is not implemented by ARM compliant cores, such as Qualcomm’s Krait 450. As previously stated, this might be due to the protection of the concept by US patent [50]. By exclusively targeting Cortex-A compliant processors not implemented by ARM, chances of not encountering AutoLock might increase. Alternatively, Flush+Reload or Flush+Flush based attacks can still be mounted on ARMv8-A SoCs that offer the cache flush instruction in userspace, i.e. the Samsung Galaxy S6 [37].
- **Achieve Same-core Scenario:** Certain attack scenarios realistically allow the adversary to execute code on the same core as the target program. Since same-core attacks are not affected by AutoLock, this entirely removes its impact. ARM TrustZone, e.g., enables the secure execution of trusted code

on an untrusted operating system. Given that the untrusted OS is compromised by the adversary, the trusted code can be scheduled to run on any given processor core. By matching the core affinity of the attacking process to the one of the respective trusted target, the attack is reduced to a same-core setting and can successfully be mounted, even across TEE boundaries [54].

- **Trigger Self-evictions:** When AutoLock is active, a cache line can only be evicted from the inclusive LLC if no higher cache level contains a copy of it. If the target program offers services to the rest of the system, the adversary can try to issue requests such that the core-private cache of the target is sufficiently polluted and the cache line under attack is evicted. The target program essentially performs a ‘self-eviction’ and thus re-enables LLC evictions and consequently cross-core attacks.
- **Increase Load and Waiting Time:** Inclusive caches with AutoLock require that the number of ways in lower levels are greater or equal than the sum of ways in all higher cache levels. This limits the associativity of core-private caches, which the LLC is inclusive to, especially on multi-core systems. If the attack allows, an adversary can take advantage of the low associativity and simply prolong the waiting time between reloads such that the target line will automatically be evicted from core-private caches by other system activity scheduled on the respective core. To amplify the effect, the adversary can also try to increase overall system load, e.g., by issuing requests to the OS aiming at increasing background activity in the targeted core.
- **Target Large Data Structures:** Self-evictions, high system load, and prolonged waiting times all increase the chances that a cache line is evicted by itself from core-private caches. The success rate of an attack is further improved, if multiple cache lines can be targeted. The more lines that are exploitable, the higher the chances that at least one of them is automatically evicted from core-private caches. For example, the transformation tables (T-tables) of AES software implementations span multiple cache lines due to their size of several kilobytes. The attack proposed by Irazoqui et al. [29] observes the first line per table to recover an entire AES key. The authors note that “*any memory line of the T table would work equally well.*” In the upcoming section, we pick up this idea and demonstrate how the attack can be extended to exploit multiple cache lines to successfully circumvent AutoLock.

Note that all of the presented strategies increase the chances of successful attacks not only on inclusive caches implementing AutoLock, but also on *non-inclusive* caches.

7.1 Attack on AES

Irazoqui et al. [29] propose an attack on table based implementations of AES using Flush+Reload. The basic strategy is to flush one cache line per table before an encryption and reload it afterwards. If any lookup table value stored on the observed cache line is used during encryption, the adversary will encounter fast reload times. If said line is not accessed, it will be fetched from external memory and reload times will be slow. With all table lookups dependent on the secret key, the adversary can infer bits of the key from the observed reload times.

In table based implementations of AES, each cache line has a certain probability with which it is *not* used during en- or decryption. This probability depends on the size and the number of the tables as well as the size of the cache lines. It is defined as

$$P_{na} = \left(1 - \frac{t}{256}\right)^n. \quad (1)$$

Variable t denotes the number of table entries stored on a cache line. For 4-byte entries and a 64-byte cache line, $t = 16$. Variable n defines the number of accesses to the table that a cache line is part of. Given an AES-128 implementation that uses four 1 kiB T-tables and performs 160 lookups per encryption, which evenly spread over the four tables, $n = 40$. With $t = 16$, this yields a no-access probability of $P_{na} = 0.0757$. Note that the attack exploits observations of not accessed cache lines. As a result, the number of required observations increases, as P_{na} gets smaller.

The attack targets the last round of AES, i.e., the 10th round. It is shown in Equation 2. The ciphertext is denoted as c_i ($i = 0..15$). The 10th round key is given as k_i^{10} , whereas the state of AES is defined as s_i . The target lookup table used in the last round is denoted as T . For each encryption, the adversary keeps track of the reload times and the ciphertext. If successful, the attack recovers the last round key. For the recovery phase, the last round is re-written as

$$c_i = k_i^{10} \oplus T[s_i] \rightarrow k_i^{10} = c_i \oplus T[s_i]. \quad (2)$$

Improvements The original attack targets one cache line per table. If the observations of said line are of poor quality, the attack is prolonged or fails. This can happen on a processor that implements AutoLock or non-inclusive caches. If LLC evictions fail, the adversary cannot determine whether the selected cache line has

been used during encryption. Irazoqui et al. [29] state that the attack works equally well with all cache lines carrying lookup table entries. As discussed in the previous section, it is likely in practice that some of them are automatically evicted from core-private caches, hence re-enabling the attack despite AutoLock. To leverage observations from all available cache lines, we propose three improvements to the original attack:

1. **Majority vote:** All available cache lines l are attacked ($l = 0..L$). This yields L recovered keys. For each key byte, a majority vote is done over all L recovered values. The value with the most frequent occurrence is assumed to be the correct one. If two or more values occur equally frequently, the lowest one is chosen. The majority vote ensures that wrong hypotheses from noisy cache lines are compensated for as long as the majority of lines yield correct results.
2. **Probability filter:** The reload times allow to calculate the actual no-access probability for each cache line, \hat{P}_{na}^l . For each table, the line closest to the expected theoretic probability, P_{na} , is taken and used in the attack. Lines showing distorted usage statistics due to noise and interference are discarded.
3. **Weighted counting:** Every cache line is assigned an individual score S_l that is counted each time a key byte hypothesis is derived from the line's reload times. The score is based on the absolute difference of the no-access probabilities, $d_l = \text{abs}(P_{na} - \hat{P}_{na}^l)$. It is defined as $S_l = 1 - \frac{d_l}{1 - P_{na}}$. After all scores have been added for all hypotheses, the recovery phase proceeds as proposed.

We implement the original attack and all improvements using Evict+Reload on a multi-core ARM Cortex-A15 processor featuring a data-inclusive LLC with AutoLock. We employ sliding window eviction with parameters 36-6-2 (see Table 1). The targeted AES implementation uses four 1 kiB T-tables. The adversary and target processes are running on top of a full-scale Linux operating system. In total, we perform five different attacks. First, we implement the original attack with adversary and target on the same core and then on separate cores. This illustrates the impact of AutoLock, which only affects cross-core attacks. The rest of the attacks are conducted with adversary and target on separate cores and each demonstrate one of the proposed improvements. The results are illustrated in Figure 4. Each attack is repeated for 100 random keys and the average number of correctly recovered key bytes is shown over an increasing number of encryptions. It

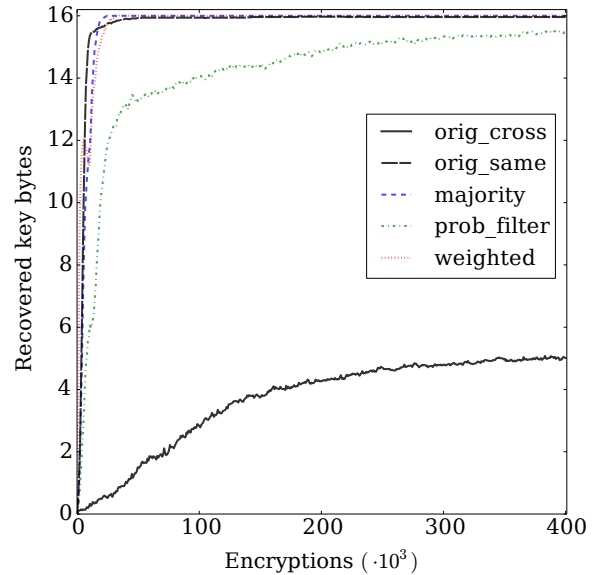


Figure 4: Evict+Reload attacks on an ARM Cortex-A15 with AutoLock; performed for 100 random keys. The number of key bytes recovered on average are displayed for an increasing number of encryptions.

can clearly be seen that AutoLock impairs the original cross-core attack (*orig_cross*). After 400,000 encryptions no more than 5 key bytes are correctly recovered. The fact that at least some key bytes are correct is owed to sporadic and automatic evictions of the observed cache lines from the target's core-private cache. These evictions can be caused by stack and heap data accesses (such as AES state and key arrays as well as their pointers), and possibly by unrelated processes scheduled on the same core. Attacks in the same-core setting (*orig_same*) are not affected by AutoLock and allow full-key recovery. Our improvements clearly demonstrate that cross-core attacks are still possible, if multiple cache lines can be observed. Both majority vote (*majority*) and weighted counting (*weighted*) recover the entire key with less than 100,000 encryptions and therefore offer similar success rates as the same-core attack. The probability filter (*prob_filter*) still allows full-key recovery within 100,000 encryptions, if a brute-force search with complexity $< 2^{32}$ is added.

The results illustrate that even on processors implementing AutoLock cache attacks can still be successful in practice, if multiple cache lines are monitored. Note that the proposed improvements are also beneficial on processors without AutoLock or on systems with non-inclusive caches. If attacks rely on observing a specific cache line, the chances of success are significantly reduced on processors implementing AutoLock.

8 Conclusion

The licensing ecosystem of ARM drives an increasingly heterogeneous market of processors with significant microarchitectural differences. Paired with a limited understanding of how ARM's cache architectures function internally, this makes assessing the practical threat of flush and eviction based cache attacks on ARM a challenging task. Although the feasibility of state of the art attacks has been demonstrated, their requirements are far from being fulfilled on all ARM processors. Flush instructions are supported only by the latest architecture version and must explicitly be enabled for userspace applications. This limits the practical utility of flush based attacks. Last-level caches can be non-inclusive, impeding practical cross-core eviction attacks that require LLCs to be inclusive. Our work shows that these attacks can still fail even on inclusive LLCs, if AutoLock is implemented. On the contrary, more sophisticated attack techniques seem to overcome both AutoLock and non-inclusive cache hierarchies. We therefore believe that a fair and comprehensive assessment of ARM's security against cache attacks requires a better understanding of the implemented cache architectures as well as rigorous testing across a broad range of ARM and thereto compliant processors.

Acknowledgments

We would like to thank our anonymous reviewers for their valuable comments and suggestions.

References

- [1] ACIİÇMEZ, O. Yet another microarchitectural attack: Exploiting i-cache. In *Proceedings of the 2007 ACM Workshop on Computer Security Architecture* (2007), CSAW '07, ACM, pp. 11–18.
- [2] ACIİÇMEZ, O., AND ÇETİN K. KOÇ. Trace-driven cache attacks on aes. In *Proceedings of the 8th International Conference on Information and Communications Security* (2006), ICICS'06, Springer-Verlag, pp. 112–121.
- [3] ACIİÇMEZ, O., SCHINDLER, W., AND ÇETİN K. KOÇ. Cache based remote timing attack on the aes. In *Topics in Cryptology – CT-RSA 2007*, M. Abe, Ed., vol. 4377 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2006, pp. 271–286.
- [4] ALMEIDA, J. B., BARBOSA, M., BARTHE, G., DUPRESSOIR, F., AND EMMI, M. Verifying constant-time implementations. In *25th USENIX Security Symposium (USENIX Security 16)* (Austin, TX, Aug 2016), USENIX Association, pp. 53–70.
- [5] ARM LIMITED. Cortex-A15 MPCore Revision: r4p0 Technical Reference Manual. http://infocenter.arm.com/help/topic/com.arm.doc.ddi0438i/DDI0438I_cortex_a15_r4p0_trm.pdf, June 2013. Accessed: 2017-06-29.
- [6] ARM LIMITED. Cortex-A7 MPCore Revision: r0p5 Technical Reference Manual. http://infocenter.arm.com/help/topic/com.arm.doc.ddi0464f/DDI0464F_cortex_a7_mpcore_r0p5_trm.pdf, April 2013. Accessed: 2017-06-29.
- [7] ARM LIMITED. ARM Architecture Reference Manual - ARMv7-A and ARMv7-R Edition, May 2014. Revision C.c.
- [8] ARM LIMITED. Cortex-A57 MPCore Processor Revision: r1p3 Technical Reference Manual. http://infocenter.arm.com/help/topic/com.arm.doc.ddi0488h/DDI0488H_cortex_a57_mpcore_trm.pdf, August 2014. Accessed: 2017-06-29.
- [9] ARM LIMITED. Cortex-A53 MPCore Processor Revision: r0p4 Technical Reference Manual. http://infocenter.arm.com/help/topic/com.arm.doc.ddi0500g/DDI0500G_cortex_a53_trm.pdf, February 2016. Accessed: 2017-06-29.
- [10] ARM LIMITED. ARM Architecture Reference Manual - ARMv8, for ARMv8-A architecture profile, March 2017. Revision B.a.
- [11] ARM LIMITED. DSTREAM. <https://developer.arm.com/products/software-development-tools/debug-probes-and-adapters/dstream>, 2017. Accessed: 2017-06-29.
- [12] BENDER, N., VAN DE POL, J., SMART, N., AND YAROM, Y. "ooh aah... just a little bit": A small amount of side channel can go a long way. In *Cryptographic Hardware and Embedded Systems – CHES 2014: 16th International Workshop, Busan, South Korea, September 23-26, 2014. Proceedings*, L. Batina and M. Robshaw, Eds., vol. 8731 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2014, pp. 75–92.
- [13] BERNSTEIN, D. J. Cache-timing attacks on AES. Tech. rep., The University of Illinois at Chicago, 2005. <http://cr.ypt.to/antiforgery/cachetiming-20050414.pdf>.
- [14] BOGDANOV, A., EISENBARTH, T., PAAR, C., AND WIENECKE, M. Differential cache-collision timing attacks on aes with applications to embedded cpus. In *The Cryptographers' Track at the RSA Conference 2010, San Francisco, CA, USA, March 1-5, 2010. Proceedings* (2010), Springer-Verlag, pp. 235–251.
- [15] BONNEAU, J., AND MIRONOV, I. Cache-collision timing attacks against aes. In *Proc. Cryptographic Hardware and Embedded Systems (CHES) 2006. Lecture Notes in Computer Science* (2006), Springer-Verlag, pp. 201–215.
- [16] COUNTERPOINT TECHNOLOGY MARKET RESEARCH. Global Smartphone Shipments Market Share Q1 2017. Infographic Q1-2017 Mobile Market Monitor, 2017.
- [17] DEMME, J., MARTIN, R., WAKSMAN, A., AND SETHUMADHAVAN, S. Side-channel vulnerability factor: A metric for measuring information leakage. In *Proceedings of the 39th Annual International Symposium on Computer Architecture* (Washington, DC, USA, 2012), ISCA '12, IEEE Computer Society, pp. 106–117.
- [18] DOYCHEV, G., FELD, D., KOPF, B., MAUBORGNE, L., AND REINEKE, J. Cacheaudit: A tool for the static analysis of cache side channels. In *Presented as part of the 22nd USENIX Security Symposium (USENIX Security 13)* (Washington, D.C., 2013), USENIX, pp. 431–446.
- [19] FOURNIER, J., AND TUNSTALL, M. Cache based power analysis attacks on aes. In *Information Security and Privacy*, L. Batten and R. Safavi-Naini, Eds., vol. 4058 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2006, pp. 17–28.
- [20] GALLAIS, J.-F., AND KIZHVATOV, I. Error-tolerance in trace-driven cache collision attacks. In *Second International Workshop on Constructive Side-Channel Analysis and Secure Design* (2011).
- [21] GRUSS, D., MAURICE, C., FOGH, A., LIPP, M., AND MANGARD, S. Prefetch side-channel attacks: Bypassing smap and

- kernel aslr. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security* (2016), CCS '16, pp. 368–379.
- [22] GRUSS, D., MAURICE, C., AND MANGARD, S. Rowhammer.js: A remote software-induced fault attack in javascript. In *13th Conference on Detection of Intrusions and Malware & Vulnerability Assessment (DIMVA)* (2016).
- [23] GRUSS, D., MAURICE, C., WAGNER, K., AND MANGARD, S. Flush+flush: A fast and stealthy cache attack. In *13th Conference on Detection of Intrusions and Malware & Vulnerability Assessment (DIMVA)* (2016).
- [24] GRUSS, D., SPREITZER, R., AND MANGARD, S. Cache template attacks: Automating attacks on inclusive last-level caches. In *24th USENIX Security Symposium (USENIX Security 15)* (Washington, D.C., Aug 2015), USENIX Association, pp. 897–912.
- [25] GULLASCH, D., BANGERTER, E., AND KRENN, S. Cache games – bringing access-based cache attacks on AES to practice. In *Proceedings of the 2011 IEEE Symposium on Security and Privacy* (Oakland, CA, USA, 2011), SP '11, IEEE Computer Society, pp. 490–505.
- [26] İNCI, M. S., GÜLMEZOĞLU, B., IRAZOQUI, G., EISENBARTH, T., AND SUNAR, B. Cache attacks enable bulk key recovery on the cloud. In *Cryptographic Hardware and Embedded Systems – CHES 2016: 18th International Conference, Santa Barbara, CA, USA, August 17–19, 2016, Proceedings* (Berlin, Heidelberg, 2016), Springer Berlin Heidelberg, pp. 368–388.
- [27] IRAZOQUI, G., EISENBARTH, T., AND SUNAR, B. S5a: A shared cache attack that works across cores and defies vm sandboxing and its application to aes. In *Proceedings of the 2015 IEEE Symposium on Security and Privacy* (San Jose, CA, USA, 2015), SP '15, IEEE Computer Society, pp. 591–604.
- [28] IRAZOQUI, G., EISENBARTH, T., AND SUNAR, B. Cross processor cache attacks. In *Proceedings of the 11th ACM Symposium on Information, Computer and Communications Security* (New York, NY, USA, 2016), ASIA CCS '16, ACM, pp. 353–364.
- [29] IRAZOQUI, G., İNCI, M., EISENBARTH, T., AND SUNAR, B. Wait a minute! a fast, cross-vm attack on aes. In *Research in Attacks, Intrusions and Defenses: 17th International Symposium, RAID 2014, Gothenburg, Sweden, September 17–19, 2014. Proceedings*, A. Stavrou, H. Bos, and G. Portokalidis, Eds., vol. 8688 of *Lecture Notes in Computer Science*. Springer International Publishing, 2014, pp. 299–319.
- [30] IRAZOQUI, G., İNCI, M. S., EISENBARTH, T., AND SUNAR, B. Lucky 13 strikes back. In *Proceedings of the 10th ACM Symposium on Information, Computer and Communications Security* (New York, NY, USA, 2015), ASIA CCS '15, ACM, pp. 85–96.
- [31] KAY, R. Intel And AMD: The Juggernaut Vs. The Squid. Forbes, Inc., <http://www.forbes.com/sites/rogerkay/2014/11/25/intel-and-amd-the-juggernaut-vs-the-squid>, 2014. Published: 2014-11-25. Accessed: 2017-06-29.
- [32] KELSEY, J., SCHNEIER, B., WAGNER, D., AND HALL, C. Side channel cryptanalysis of product ciphers. In *Computer Security - ESORICS 98*, J.-J. Quisquater, Y. Deswarte, C. Meadows, and D. Gollmann, Eds., vol. 1485 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 1998, pp. 97–110.
- [33] KIM, T., PEINADO, M., AND MAINAR-RUIZ, G. Stealthmem: System-level protection against cache-based side channel attacks in the cloud. In *Presented as part of the 21st USENIX Security Symposium (USENIX Security 12)* (Bellevue, WA, 2012), USENIX, pp. 189–204.
- [34] KLUG, B., AND SHIMPI, A. L. Qualcomm's new snapdragon s4: Msm8960 & krait architecture explored. <http://www.anandtech.com/show/4940/qualcomm-new-snapdragon-s4-msm8960-krait-architecture/2>, 2011. Published: 2011-10-07. Accessed: 2017-06-29.
- [35] KOCHER, P. C. Timing attacks on implementations of diffie-hellman, rsa, dss, and other systems. In *Proceedings of the 16th Annual International Cryptology Conference on Advances in Cryptology* (1996), CRYPTO '96, Springer-Verlag, pp. 104–113.
- [36] LINUX PROGRAMMER'S MANUAL. perf_event_open - set up performance monitoring. http://man7.org/linux/man-pages/man2/perf_event_open.2.html, 2016. Accessed: 2017-06-29.
- [37] LIPP, M., GRUSS, D., SPREITZER, R., MAURICE, C., AND MANGARD, S. Armageddon: Cache attacks on mobile devices. In *25th USENIX Security Symposium (USENIX Security 16)* (Austin, TX, Aug 2016), USENIX Association, pp. 549–564.
- [38] LIU, F., GE, Q., YAROM, Y., MCKEEN, F., ROZAS, C., HEISER, G., AND LEE, R. B. Catalyst: Defeating last-level cache side channel attacks in cloud computing. In *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)* (March 2016), pp. 406–418.
- [39] LIU, F., YAROM, Y., GE, Q., HEISER, G., AND LEE, R. B. Last-level cache side-channel attacks are practical. In *Proceedings of the 2015 IEEE Symposium on Security and Privacy* (San Jose, CA, USA, 2015), SP '15, IEEE Computer Society, pp. 605–622.
- [40] NEVE, M., AND SEIFERT, J.-P. Advances on access-driven cache attacks on aes. In *Selected Areas in Cryptography*, E. Biham and A. Youssef, Eds., vol. 4356 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2007, pp. 147–162.
- [41] OREN, Y., KEMERLIS, V. P., SETHUMADHAVAN, S., AND KEROMYTIS, A. D. The spy in the sandbox: Practical cache attacks in javascript and their implications. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security* (2015), CCS '15, pp. 1406–1418.
- [42] OSVIK, D. A., SHAMIR, A., AND TROMER, E. Cache attacks and countermeasures: The case of aes. In *Topics in Cryptology – CT-RSA 2006: The Cryptographers' Track at the RSA Conference 2006, San Jose, CA, USA, February 13–17, 2005. Proceedings* (Berlin, Heidelberg, 2006), Springer Berlin Heidelberg, pp. 1–20.
- [43] PAGE, D. Theoretical use of cache memory as a cryptanalytic side-channel. *Department of Computer Science, University of Bristol, Tech. Rep. CSTR-02-003* (2002). <http://eprint.iacr.org/2002/169>.
- [44] PERCIVAL, C. Cache missing for fun and profit. In *Proceedings of BSDCan 2005* (2005).
- [45] POSIX PROGRAMMER'S MANUAL. clock_getres, clock_gettime, clock_settime - clock and timer functions. http://man7.org/linux/man-pages/man3/clock_gettime.3p.html, 2013. Accessed: 2017-06-29.
- [46] RISTENPART, T., TROMER, E., SHACHAM, H., AND SAVAGE, S. Hey, you, get off of my cloud: Exploring information leakage in third-party compute clouds. In *Proceedings of the 16th ACM Conference on Computer and Communications Security* (New York, NY, USA, 2009), CCS '09, ACM, pp. 199–212.
- [47] SHIMPI, A. L., AND GREENHALGH, P. Answered by the Experts: ARM's Cortex A53 Lead Architect, Peter Greenhalgh. <http://www.anandtech.com/show/7591/answered-by-the-experts-arms-cortex-a53-lead-architect-peter-greenhalgh>, 2013. Published: 2013-12-17. Accessed: 2017-06-29.

- [48] TSUNOO, Y., SAITO, T., SUZAKI, T., SHIGERI, M., AND MIYAUCHI, H. Cryptanalysis of des implemented on computers with cache. In *Cryptographic Hardware and Embedded Systems - CHES 2003*, C. Walter, c. K. Koç, and C. Paar, Eds., vol. 2779 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2003, pp. 62–76.
- [49] WANG, Z., AND LEE, R. B. New cache designs for thwarting software cache-based side channel attacks. In *Proceedings of the 34th Annual International Symposium on Computer Architecture* (New York, NY, USA, 2007), ISCA '07, ACM, pp. 494–505.
- [50] WILLIAMSON, B. D. Line allocation in multi-level hierarchical data stores. Patent US8271733 B2, ARM Limited, September 2012. <https://www.google.com/patents/US8271733>.
- [51] YAROM, Y., AND FALKNER, K. Flush+reload: A high resolution, low noise, l3 cache side-channel attack. In *23rd USENIX Security Symposium (USENIX Security 14)* (San Diego, CA, 2014), USENIX Association, pp. 719–732.
- [52] YAROM, Y., GENKIN, D., AND HENINGER, N. Cachebleed: A timing attack on openssl constant time RSA. In *Cryptographic Hardware and Embedded Systems – CHES 2016: 18th International Conference, Santa Barbara, CA, USA, August 17-19, 2016, Proceedings* (Berlin, Heidelberg, 2016), B. Gierlichs and A. Y. Poschmann, Eds., Springer Berlin Heidelberg, pp. 346–367.
- [53] ZANKL, A., HEYSZL, J., AND SIGL, G. Automated detection of instruction cache leaks in modular exponentiation software. In *Smart Card Research and Advanced Applications: 15th International Conference, CARDIS 2016, Cannes, France, November 7–9, 2016, Revised Selected Papers* (Cham, 2017), K. Lemke-Rust and M. Tunstall, Eds., Springer International Publishing, pp. 228–244.
- [54] ZHANG, N., SUN, K., SHANDS, D., LOU, W., AND HOU, Y. T. Truspy: Cache side-channel information leakage from the secure world on arm devices. Cryptology ePrint Archive, Report 2016/980, 2016. <http://eprint.iacr.org/2016/980>.
- [55] ZHANG, T., ZHANG, Y., AND LEE, R. B. Cloudradar: A real-time side-channel attack detection system in clouds. In *Research in Attacks, Intrusions, and Defenses: 19th International Symposium, RAID 2016, Paris, France, September 19-21, 2016, Proceedings* (Cham, 2016), F. Monrose, M. Dacier, G. Blanc, and J. Garcia-Alfaro, Eds., Springer International Publishing, pp. 118–140.
- [56] ZHANG, X., XIAO, Y., AND ZHANG, Y. Return-oriented flush-reload side channels on arm and their implications for android devices. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security* (New York, NY, USA, 2016), CCS '16, ACM, pp. 858–870.
- [57] ZHANG, Y., JUELS, A., REITER, M. K., AND RISTENPART, T. Cross-VM Side Channels and Their Use to Extract Private Keys. In *Proceedings of the 2012 ACM Conference on Computer and Communications Security* (New York, NY, USA, 2012), CCS '12, ACM, pp. 305–316.
- [58] ZHANG, Y., JUELS, A., REITER, M. K., AND RISTENPART, T. Cross-tenant side-channel attacks in paas clouds. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security* (New York, NY, USA, 2014), CCS '14, ACM, pp. 990–1003.
- [59] ZHOU, Z., REITER, M. K., AND ZHANG, Y. A software approach to defeating side channels in last-level caches. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security* (New York, NY, USA, 2016), CCS '16, ACM, pp. 871–882.

Notes

¹ To the best of our knowledge, there is no official public documentation for the Qualcomm Krait 450. The details given in Table 1 are based on the results of our experiments as well as on statements from online articles [34].

² The list of devices supported by DSTREAM can be retrieved from ARM’s website at <https://developer.arm.com/products/software-development-tools/ds-5-development-studio/resources/supported-devices>.

³ The timing measurement code can be retrieved from the GitHub repository at <https://github.com/IAIK/armageddon>.

