

Amplifying Side Channels Through Performance Degradation

Thomas Allan
The University of Adelaide and Data61, CSIRO
tom.allan@student.adelaide.edu.au

Katrina Falkner
The University of Adelaide
katrina.falkner@adelaide.edu.au

Billy Bob Brumley
Tampere University of Technology
billy.brumley@tut.fi

Joop van de Pol
University of Bristol
joop.vandepol.2011@my.bristol.ac.uk

Yuval Yarom
The University of Adelaide and Data61, CSIRO
yval@cs.adelaide.edu.au

ABSTRACT

Interference between processes executing on shared hardware can be used to mount performance-degradation attacks. However, in most cases, such attacks offer little benefit for the adversary. In this paper, we demonstrate that software-based performance-degradation attacks can be used to amplify side-channel leaks, enabling the adversary to increase both the amount and the quality of information captured.

We identify a new information leak in the OpenSSL implementation of the ECDSA digital signature algorithm, albeit seemingly unexploitable due to the limited granularity of previous trace procurement techniques. To overcome this imposing hurdle, we combine the information leak with a microarchitectural performance-degradation attack that can slow victims down by a factor of over 150. We demonstrate how this combination enables the amplification of a side-channel sufficiently to exploit this new information leak. Using the combined attack, an adversary can break a private key of the `secp256k1` curve, used in the Bitcoin protocol, after observing only 6 signatures—a four-fold improvement over all previously described attacks.

1. INTRODUCTION

Executing multiple clients' workloads on a single hardware platform can help achieve high resource utilisation. A consequence of this resource sharing is that workloads of different clients can interfere with each other due to shared-resource contention [53, 54].

Malicious clients can exploit this interference to mount performance-degradation attacks against co-resident clients [13, 17, 33, 43]. Fortunately, such attacks have a limited usability. In most cases, the attacker does not gain any direct benefit from the attack. The main benefit an attacker gets from mounting a performance-degradation attack is harming the victim. An excep-

tion is the attack of Varadarajan et al. [48], in which the attack frees resources for the attacker's use by forcing the victim to wait on other resources.

In this paper we investigate the use of performance-degradation attacks to amplify side-channel attacks. Slower encryption provides the adversary with more opportunities for collecting side-channel information [50]. Hence, actively slowing the victim down carries the promise of better side-channel attacks. This concept has been used in power analysis attacks which reduce the clock frequency to achieve a better signal acquisition [27, 37]. Gruss et al. [14] demonstrate that performance-degradation attacks can slow encryption down and speculate on possible benefits for side-channel attack. However, they do not demonstrate that the attack indeed improves the side channel and the benefits for an adversary remain unclear. This work is the first to demonstrate that performance-degradation attacks can amplify side-channel information and provide tangible benefits to the adversary.

More specifically, we identify a new information leak in the OpenSSL implementation of ECDSA with the `secp256k1` curve and use it to extend the attack of Van de Pol et al. [42]. In a nutshell, Van de Pol et al. [42] trace the use of point addition and point doubling used throughout the scalar multiplication of ECDSA signature generation. From this sequence, they infer information on the ephemeral key used for the signature. The long term private key is then reconstructed from the information collected from multiple signatures by using a lattice attack.

We identify that tracing point inversions can increase the amount of information collected on each ephemeral key, potentially reducing the number of signatures required for breaking the key. However, tracing point inversions introduces two formidable problems: (1) due to the high resolution required, there is a high probability of missing point inversions; (2) adding the trace for point inversions increases the number of memory locations we trace, limiting the applicability of the attack.

To overcome these limitations, we apply a performance-degradation attack against the scalar multiplication code. We repeatedly evict parts of the victim's code from the cache, slowing elliptic group operations by a factor of over 40, and the scalar multiplication by a factor of 32, allowing a virtually error-free trace of the curve operations sequence. By using this technique, we can break the private key of the `secp256k1` curve used in Bitcoin after observing as few as 6 signatures.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

ACSAC '16 December 05-09, 2016, Los Angeles, CA, USA

© 2016 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-4771-6/16/12.

DOI: <http://dx.doi.org/10.1145/2991079.2991084>

To further our understanding of the attack we study the performance-degradation and the side-channel attacks independently of each other. We demonstrate that strategically evicting code from the cache can slow programs down by a factor of over 150, with a mean slow-down factor of 18 over the integer SPEC 2006 [18] benchmarks and 15 over the floating-point benchmarks.

We also analyse the FLUSH+RELOAD side-channel attack [15, 56]. We demonstrate that the attack has a maximum resolution which depends on the number of memory locations it attempts to probe. We further show a relationship between the resolution of events in the victim and the likelihood of the attacker missing an event.

The contributions of this paper are:

- We investigate the cache-eviction performance-degradation attack and demonstrate that it is about 8 times more potent than previously disclosed attacks. (Section 4)
- We analyse the FLUSH+RELOAD attack, identifying the tradeoffs between the attack resolution, the number of memory locations monitored and the probability of missing a monitored events. (Section 5)
- We identify point inversions as a new source of leaked information in the implementation of ECDSA over prime fields in OpenSSL and show how to exploit this information. (Section 6)
- We use the performance-degradation attack to amplify the side channel, allowing an attacker to observe the leaked information. The combined attack requires roughly a quarter of observed signatures compared to any prior attack. (Section 7)

2. BACKGROUND

2.1 The Memory Hierarchy

The cache is part of the memory hierarchy that exploits the spatial and temporal locality of memory access to bridge the performance gap between the fast processor and the slower memory. Modern processors feature a hierarchy of caches, with higher-level caches, which are closer to the processor core, being smaller but faster than lower-level caches, which are closer to the main memory. In recent Intel architecture, there are, typically, three levels of cache. Each core has two levels of caches, called the L1 and L2 caches. The cores share access to a larger Last-Level Cache (LLC).

To exploit spatial locality, caches are organised in fixed-size *lines*, which are the units of allocation and transfer of data in the memory hierarchy. When the processor needs to access a memory address, it first checks if the line containing the address is cached in the top-level L1 cache. In a *cache hit*, the data is served from a copy of the data in the cache. Otherwise, in a *cache miss*, the processor repeats the search for the line in the next lower level in the memory hierarchy. When the line is found, the processor stores its contents in the cache, reducing the time required for accessing it in the near future. See [39, Ch. 8] for a good overview of caching in computer architecture.

Modern caches are typically *set associative*. A set associative cache is divided into multiple *sets*, each consisting of multiple *ways*. Each memory line is mapped to a single cache set. The memory line can only be cached in the set it is mapped to, but can be cached in any of the ways of the set. Typically, the set a memory line maps to is determined by a sequence of bits in the physical address of the memory line. However, the LLC in modern Intel processor uses a more complex hash function to determine the mapping [20, 29, 57].

Several cache optimisations result in memory lines being brought to the cache without the code accessing data in these lines. In the Intel architecture, the *spatial prefetcher* pairs consecutive memory lines and attempts to fetch the pair of a missed line [21]. Another optimisation is to detect sequences of accesses to consecutive memory addresses and prefetch memory lines that the processor anticipates may be required [21]. A third optimisation is *speculative execution*, where the processor attempts to follow both paths of a conditional branch before the branch condition is evaluated [47], bringing the code of both paths into the cache.

When multiple programs share the same cache, one program's use of the cache may evict another program's data from the cache, which due to the timing difference between cache hits and cache misses may create noticeable timing variations in the sharing programs. These timing variations have been used to mount side-channel attacks [1, 5, 26, 38, 40, 44, 58].

2.2 The Flush+Reload Attack

FLUSH+RELOAD [56] is a cache-based side-channel attack technique. Unlike other techniques, which infer the memory lines the victim accesses based on activity in cache sets, FLUSH+RELOAD positively identifies access to memory lines, giving it high accuracy, high signal to noise ratio and high resolution. The attack has been used in various settings, including between non-trusting processes, between isolated containers and across virtual machines and has been shown to be effective against multiple algorithms [3, 14, 22, 23, 42, 55, 59].

FLUSH+RELOAD relies on memory sharing between the victim and the adversary. Such sharing could be achieved via the use of shared libraries or using page de-duplication [2, 49]. To identify victim access to a shared memory line, the adversary flushes or evicts the memory line from the cache, waits a bit and then measures the time it takes to reload the memory line. If the victim accesses the line during the wait, the line will be cached and the reload will retrieve it from the cache. Otherwise, the line will not be cached and reloading will have to retrieve it from the main memory. As retrieving the line from the memory takes longer than accessing a cached copy, the adversary can distinguish between the two options and identify whether the victim has accessed the line during the wait.

The FLUSH+RELOAD attack needs processor support for evicting memory lines from the cache. So far, all published reports of the attack use the `clflush` instructions of the x86 and x86-64 instruction sets. In those instruction sets, `clflush` is an unprivileged instruction, which every process can use.

Gruss et al. [14] suggest a variant of FLUSH+RELOAD, called EVICT+RELOAD, which does not require a specific instruction for evicting the memory line. Instead, they evict the victim memory line by accessing a number of memory lines that map to the same cache set as the victim line. Evicting the victim memory line using this technique takes significantly longer than using the `clflush` instruction. (325 cycles compared with 41 for `clflush`.) Furthermore, the eviction may fail, resulting in a false positive.

Both FLUSH+RELOAD and EVICT+RELOAD need to evict the victim cache line from all of the caches that the victim uses. When the victim and the adversary do not execute on the same core, they do not share the L1 and L2 caches. In this case, the attack relies on the *inclusion* property of the LLC. The contents of an inclusive cache is a superset of the contents of all higher level caches. To maintain the inclusion property, when a memory line is evicted from the LLC, the processor also evicts it from all of the L1 and L2 caches above it. All of the published attacks run on Intel processors, which use inclusive LLCs. Yarom and Falkner [56] report that the

FLUSH+RELOAD attack does not work on AMD processors due to their non-inclusive LLCs.

2.3 Related Work

Several works have investigated performance-degradation attacks by co-located adversaries. Grunwald and Ghiasi [13] implement two attacks against Intel HyperThreading (HT), a Simultaneous Multithreading (SMT) technique. The first attack uses denormalised floating point numbers [12], which flush the instruction pipeline of the Pentium 4 processor used. The second attack uses self-modifying code, which results in flushing both the pipeline and the processor’s trace cache. To test the attack, they use a compute-bound victim which repeatedly calculates the MD5 hash. The victim is slowed by about 120% with the first attack and by a factor of 20 with the second.

Heat stroke [17] is a performance-degradation attack that exploits the thermal management of the processor chip. Certain components of the chip tend to overheat when experiencing high utilisation, forcing the processor to reduce the utilisation of the hot components until they cool down. The authors use a simulated multi-threaded processor to test the attack. The adversary generates many register accesses causing overheating in the shared register file. The processor responds to overheating by slowing access to the register file. The attack achieves a mean slow down by a factor of 8 over the SPEC 2000 benchmark suite.

Matthews et al. [28] compare the performance isolation properties of virtualisation. They implement multiple adversaries, each attempting to monopolise a system resource. The main finding is that OS-level virtualisation (e.g. Solaris containers) provides less isolation than system-level hypervisors such as VMware or Xen. In particular, it performs poorly under memory or process number pressure. Other than that, all systems at most experience minor interference.

Moscibroda and Mutlu [33] note that the scheduling policy of memory banks favours requests for the currently open DRAM row. Consequently, an adversary that issues many requests to the same row can cause memory-access delays for programs that access the same DRAM bank. These delays can slow the victim down by a factor of 2.9 for one adversary and up to a factor of 4 for multiple adversaries. The suggested fix is to change the DRAM scheduling algorithm.

Woo and Lee [52] investigate attacks against a shared LLC. The attacks aim to evict entries from the LLC and rely on the LLC inclusiveness to also evict data from the victim L1. Two forms of attack are suggested and are tested using a simulator—no tests on a real processor are performed. Attacks using load instructions slow victims down by 50% on average, with a maximum slowdown of 100%. (The amount of degradation is estimated from the graphs provided due to the absence of exact figures.) The second form of attack uses atomic instructions which lock access to the bus. The mean slowdown with this attack is by a factor of 5, with a maximum slowdown factor of 10.

Another LLC monopolising attack is suggested by Weng et al. [51] which demonstrate a significant performance drop in co-resident VMs. The paper does not present exact figures, but judging from the supplied graphs, the performance seems to drop by about 30%. As a countermeasure, Weng et al. [51] suggest not scheduling non-trusting VMs concurrently on the same processor package.

Cardenas and Boppana [7] also use an adversary that tries to monopolise the LLC. The attack reduces the performance of the victim by 50% with a single attacking thread and up to 75% with multiple threads. Based on the observation that the adversary also suffers

LLC misses, the paper suggests using the performance management unit (PMU) to identify the adversary and eventually mitigate the attack. We note that because our adversary does not suffer cache misses, this mitigation does not apply to the attack we present.

Richter et al. [43] investigate multiple techniques for degrading the performance of a shared PCI bus. They show that when I/O virtualisation is used, a malicious VM can cause a drop of 27% in TCP throughput. With multiple attackers the drop reaches 35%.

Swiper [8] generates adversarial I/O workload to slow a target application down, achieving a reduction of up to 31% in the throughput of Web and media servers.

In all the systems described above, the only motivation for adversarial behaviour is the damage it causes to the victim. Performance degradation attacks are therefore a form of vandalism, whose only benefit is harming the victim. Varadarajan et al. [48] is the only prior work to offer direct benefits to the adversary. The resource freeing attack suggested uses a performance-degradation attack to slow a victim down. The adversary can then benefit from the victim slow down by using resources that the victim would otherwise use. The paper demonstrates how increasing the load on the victim gives the adversary a 60% performance boost.

Gruss et al. [14] suggest using repeated cache evictions to slow encryption down. They demonstrate a slowdown of AES encryption from 320 cycles to up to 20,000 cycles, and speculate that this may be used to perform a trace-driven attack using FLUSH+RELOAD. No further analysis is provided and given our analysis of the FLUSH+RELOAD attack (Section 5) it is not clear how much information an adversary can collect using FLUSH+RELOAD during 20,000 cycles.

Walter [50] demonstrates that longer keys are more vulnerable to side channel attacks than shorter keys because the operations on longer keys are slower. He does not, however, suggest slowing encryption to amplify the side channel.

Örs et al. [37] and Mangard et al. [27] reduce the hardware clock speed to improve the results of power analysis of ASIC-based AES implementations. To the best of our knowledge, no prior work has demonstrated the use of software-based performance degradation to improve side-channel attacks.

Recently, Pereida et al. [41] published an attack on OpenSSL DSA. They use our performance-degradation technique to amplify the side-channel information leaked from the OpenSSL implementation of modular exponentiation. Their work demonstrates that the utility of our work extends beyond the examples given in this paper.

3. THREAT MODEL

In the attack scenario, the adversary executes code concurrently with victim code on the same hardware. This scenario is common in multi-user operating systems and in virtualised environments. The operating system or the hypervisor prevent the adversary from accessing the victim’s data.

We assume that the system supports a form of read-only sharing between the adversary and the victim. This sharing could be based on file mapping, e.g. shared libraries, or it can be based on coalescing identical contents through memory de-duplication. Memory de-duplication is known to be vulnerable to side-channel attacks [46], and is one of the requirements for the FLUSH+RELOAD attack [56]. We show that it also enables performance-degradation attacks. Like the FLUSH+RELOAD attack, we also assume a shared inclusive LLC and require an efficient method of evicting memory lines from the cache.

4. A PERFORMANCE DEGRADATION ATTACK

The performance-degradation attack we describe is based on the observation that programs tend to spend a significant part of their execution within a small “hot” section of the program code. Under normal execution, the frequently executed code is in the processor cache, hence access to it is fast.

If the memory that contains the hot code is shared between the adversary and the victim, the adversary can evict memory lines that contain that code from the last-level cache. This forces the victim to wait until the processor loads the code from the memory, introducing delays to the victim’s process. Repeatedly evicting the hot code would negate the performance benefits of the cache, slowing the victim down.

The amount of slowdown depends primarily on the difference between the latencies of the cache and the memory. We measure the time it takes to load data from the L1 cache and from memory on an HP Elite 8300 running CentOS 6.5. (Intel i5-3470 processor, running at 3.2 GHz, with 8 GiB of DDR3-1600 CL-11 memory.)

Figure 1 shows the distribution over 100,000 measurements.

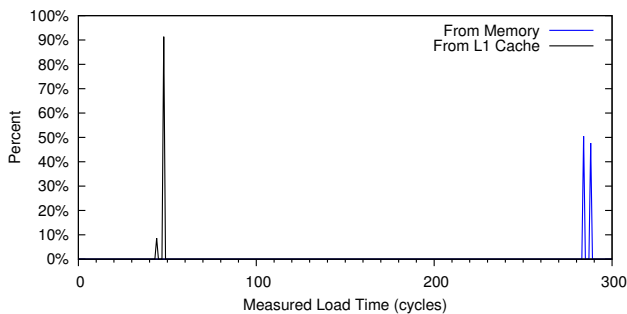


Figure 1: Distribution of L1 cache and memory access times.

As we can see, virtually all loads from the L1 cache take 48 cycles. Over 98% of the loads from the memory take between 280 and 290 cycles, with the rest spread over the interval 250–1200 cycles.

In addition to data access latency, the measurements include the overhead of the measurement code. Due to optimisations, such as instruction pipelining and reordering, parallel use of multiple functional units and data prefetching, we cannot measure this overhead. Given that the L1 cache latency is 4 cycles [21], we can conclude that the memory latency is around 240 cycles.

To measure the effects of the attack, we test it with the SPEC CPU 2006 [18] benchmark suite. To generate a baseline performance measurement, we pin the SPEC benchmarks to one core and run them on an otherwise idle machine. The measurements follow the SPEC reporting guidelines. In particular, we use the SPEC ref workload, and for each benchmark we use the median time of three runs.

We then measure the performance of the benchmark under the attack. We measure under two scenarios—with a single attacking thread and with three attacking threads running in parallel. To avoid affecting the SPEC benchmark through time sharing, we pin the SPEC benchmarks and the attacking threads, each to a separate core. As in the baseline case, the machine is otherwise idle.

To apply the attack, we need to identify the hot sections of each of the SPEC benchmarks. One possible way of doing that is to read and understand the code of each benchmark and use that understanding to identify frequently used code sections. However, due

to the size of the code base, such an approach would require significant effort and is prone to errors due to limited understanding of the code [45].

Instead, we use automatic tools for analysing the SPEC benchmarks. We build the SPEC benchmarks with instrumentation for collecting code-coverage information. We then use the program gcov to find out which source lines are the most frequently executed. Our attack targets this code.

Because the instrumentation skews the performance of the program, we do not use the instrumented binaries for the performance testing. Instead, we build optimised SPEC benchmarks with debugging symbols and use these debugging symbols to find the memory addresses corresponding to the lines identified through code coverage. The result of this process is a list of candidate memory lines for the attack. We note that debugging symbols are not loaded into memory when the program executes and do not affect its performance.

Usually, to achieve an efficient attack, we cannot use all the candidate memory lines. The reason being that evicting a line from the cache takes time. If we try to evict too many memory lines, we reduce the frequency of evicting each of the lines. Hence lines stay longer in the cache, allowing the victim to benefit from faster access to them. With cache eviction taking around 70 cycles and memory access around 240, we should be able to evict three lines from memory before the first is reloaded. Hence, in our settings, evicting more than three cache lines in a single attacking thread reduces the efficiency of the attack.

To implement an efficient attack we therefore need to select a small number of the candidate memory lines identified above. A naïve approach is to pick memory lines corresponding to the most frequently accessed source lines. Such an approach, however, does not guarantee the most efficient attack, mostly due to the effects of optimisations, both at the compiler level and at the processor level.

Instead of attempting to accurately predict the best memory lines to use for the attack, we test the efficiency of the attack with several different selections of candidate memory lines. The results we report are for the selection that produced the most effective attack. We acknowledge that other selections may produce more effective attacks. Hence the results below may understate the strength of the attack. These results are visualised in Figure 2 and Figure 3. Detailed timing for the SPEC benchmarks are shown in Table 1.

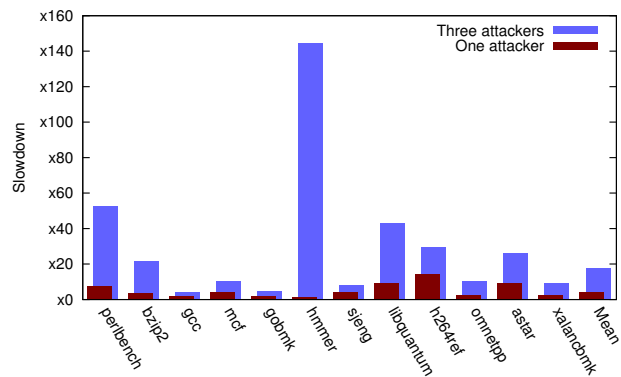


Figure 2: SPEC CPU2006 Integer Results.

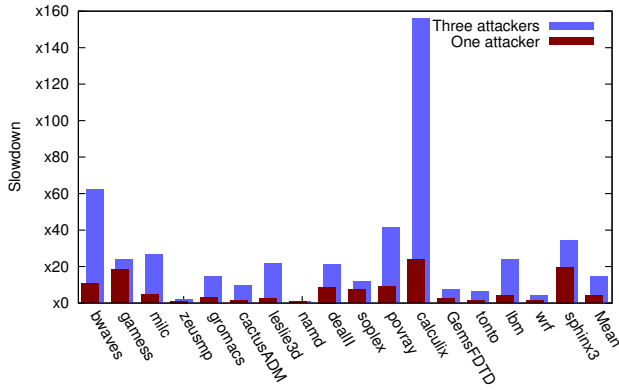


Figure 3: SPEC CPU2006 Floating Point Results.

Table 1: SPEC CPU 2006 running times (seconds)

	Baseline	One attacker	Three attackers
perlbench	396	3,052	20,922
bzip2	443	1,651	9,538
gcc	312	660	1,369
mcf	286	1,145	2,928
gobmk	446	970	2,180
hmmmer	432	514	62,507
sjeng	513	2,048	4,288
libquantum	587	5,492	25,395
h264ref	523	7,381	15,482
omnetpp	290	723	2,935
astar	375	3,364	9,792
xalancbmk	219	602	1,990
SpecINT Mean	387	1,574	6,841
bwaves	756	8,004	46,993
gamess	689	12,493	16,367
milc	405	1,846	10,737
zeusmp	387	426	823
gromacs	375	1,050	5,390
cactusADM	660	817	6,408
leslie3d	628	1,426	13,695
namd	397	414	405
dealII	317	2,761	6,723
soplex	320	2,403	3,829
povray	163	1,439	6,759
calculix	780	18,558	121,759
GemsFDTD	674	1,740	4,859
tonto	465	739	2,956
lbm	370	1,506	8,868
wrf	586	752	2,473
sphinx3	591	11,640	20,225
SpecFP Mean	469	1,989	6,885

As the results show, a single attacking thread reduces the mean execution speed to about a quarter of the normal speed, whereas three threads have a mean slowdown by a factor of 15–18. However, there is a large variance in the effectiveness of the attack. The effective slowdown with one attacker ranges from 4% (*namd*) to 2,279% (*calculix*). For three attackers the range is even bigger. *namd* is hardly affected whereas *calculix* is over 150 times slower under the attack.

The attack is less effective on *namd* and *zeusmp* because both benchmarks do not have a tight internal loop. Instead, the internal loops in these benchmarks span a relatively large amount of code. For example, the main loop in *namd* contains 256 lines of C++ code, which span over 93 memory lines. The attack only evicts a small

fraction of this code, so the overall performance hit is very small.

Considering that memory accesses are 60 times slower than cache accesses, the results for *hmmmer* and *calculix* are surprising. The observed slowdowns by factors of 145 and 156, respectively, are much larger than would be expected from the cache vs. memory speed difference. We speculate that the reason for this slowdown is the interaction of instruction fetching with the cache. Under normal circumstances the processor fetches instructions in batches of up to five instructions. While each of these fetches takes four cycles, they execute in parallel, achieving a rate of one batch per cycle [10]. If the attack is very efficient, the targeted cache line could be evicted after fetching only one batch, potentially reducing the performance by a factor of 240.

Unlike previous microarchitectural performance-degradation attacks, which affect all of the programs that use a microarchitectural component, the attack is very specific. It only targets programs that use specific code segments.

The rest of this paper describes how we exploit this property of the attack.

5. LIMITATIONS OF THE FLUSH+RELOAD ATTACK

To better understand why slowing down victims can potentially allow the adversary to improve side-channel attacks, we first study the FLUSH+RELOAD attack to see what limits its accuracy and resolution. Our focus is on asynchronous attacks, i.e. on attacks in which the adversary executes concurrently with the victim.

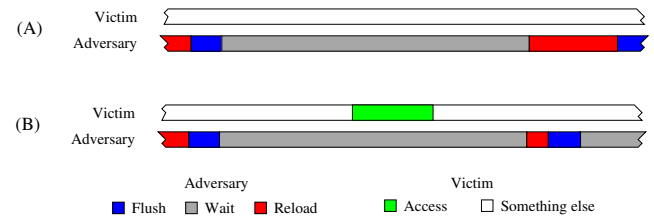


Figure 4: Timing of FLUSH+RELOAD. (A) No Victim Access (B) With Victim Access

Typically, the adversary divides time into fixed length *slots*. At the start of a time slot, the monitored memory line is flushed from the cache hierarchy. The adversary, then, waits to allow the victim time to access the memory line. At the end of the slot, the adversary reloads the memory line, measuring the time to load it. If the victim accesses the memory line during the wait, the line will be available in the cache and the reload operation will take a short time. If, on the other hand, the victim has not accessed the memory line, the line will need to be brought from memory and the reload will take significantly longer. Figure 4 (A) and (B) show the timing of the attack phases without and with victim access.

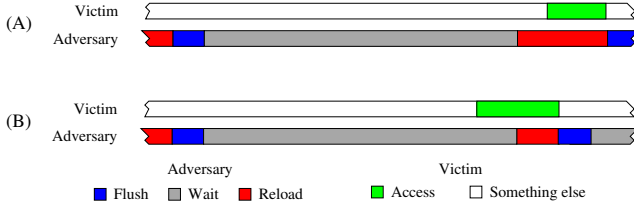
The length of the time slot determines the granularity of the attack. The adversary cannot distinguish between multiple victim accesses to the probed memory line if they all occur within the same time slot. Consequently, a shorter time slot allows for a higher attack resolution. However, because the flush and reload operations are not instantaneous, they pose a lower bound on the length of the slot. This lower bound may be more significant when the adversary needs to monitor multiple lines, in which case the slot cannot be shorter than the time required for flushing and reloading all of the probed memory lines.

Another factor that limits the slot size is the probability of missing a victim access due to overlap. In an asynchronous attack, the

Table 2: Number of missed accesses out of 10,000 tries for slot length (cycles)

Slot	Missed	Slot	Missed	Slot	Missed	Slot	Missed	Slot	Missed
1,000	5,286	7,000	807	13,000	379	19,000	281	45,000	102
2,000	2,637	8,000	660	14,000	376	20,000	241	50,000	98
3,000	1,864	9,000	607	15,000	364	25,000	226	55,000	84
4,000	1,364	10,000	531	16,000	338	30,000	150	60,000	79
5,000	1,079	11,000	589	17,000	296	35,000	155	65,000	93
6,000	913	12,000	431	18,000	209	40,000	150	70,000	93

victim operates independently of the adversary. As such, victim access to a memory location can occur at the same time the adversary reloads the location to test if it is cached, depicted in Figure 5 (A). In such a case, the victim access will not trigger a cache fill. Instead, the victim will use the cached data from the reload phase. Consequently, the adversary will miss the access.

**Figure 5: Overlap in FLUSH+RELOAD. (A) Total overlap (B) Partial overlap**

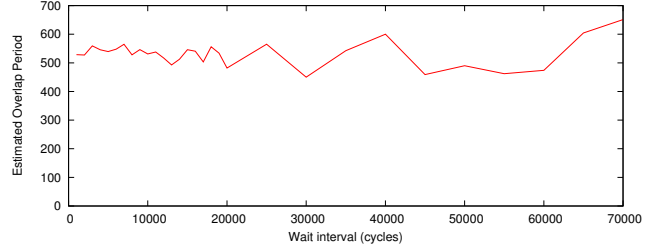
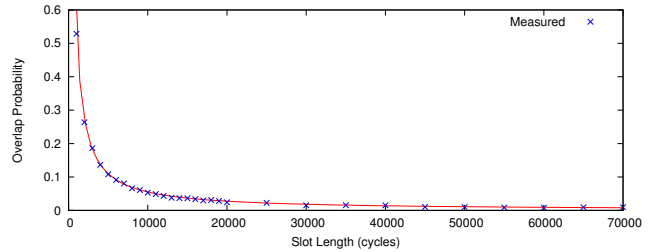
A similar scenario occurs when the reload operation partially overlaps the victim access. In this case, depicted in Figure 5 (B), the reload phase starts while the victim is waiting for the data. The reload benefits from the victim access and terminates faster than if the data has to be loaded from memory. However, the timing may still be longer than a load from the cache. Whether the adversary recognises a partial overlap as a read from the cache or from memory depends on the time difference between the start of the victim access and the start of the adversary reload.

As we can see, there is a short *overlap period* that starts a bit before the adversary probe and ends when the adversary evicts the monitored line from the cache. Victim accesses to the monitored cache line during the overlap period are missed by the adversary. Because the victim access time is independent of the adversary probe, we can expect that the probability of a miss would be the ratio between the length of the overlap period and the interval between adversary probes.

To validate this expectation we measure the miss rate with different slot sizes. We run an adversary program that monitors a memory line at a fixed rate. In parallel, we run a victim program that accesses the monitored memory line 10,000 times, and count how many of these 10,000 accesses our adversary misses. Table 2 summarises the results.

We can now multiply the miss rate by the length of the slot to estimate the length of the overlap period. Figure 6 shows the estimated overlap period for each slot length. We can see that with a few outliers, due to noise, the estimated period is fairly stable. The average estimated period is 530 cycles. Figure 7 shows the overlap probability for each slot length along with the calculated value (530/slot).

A further aspect that affects the attack accuracy is operating system activity. The operating system may suspend the adversary execution to handle some system activity, such as a network or a timer interrupt. If the interruption is short enough to be wholly contained

**Figure 6: Estimated length of the overlap period****Figure 7: Slot length and overlap probability**

within a time slot, it will not affect the attack. If, however, the adversary is interrupted for a longer period, the adversary loses the ability to distinguish between and to order multiple events occurring during the interruption.

In our experience, shorter interruptions of about 5,000 cycles are quite common, occurring, on average, about 1,000 times per second. Longer interruptions of about 30,000 cycles or 9 μ s occur at a rate of 50 per second. Significantly longer interruptions are possible when the operating system suspends the adversary in order to time-share the processor.

In summary, to achieve a high attack resolution, the adversary needs to use a short time slot. However, the length of the probe and the number of required probes present a lower limit on the slot length and, consequently, an upper limit on the attack resolution. Furthermore, the higher the attack resolution is the higher the probability of an error due to missing a victim access or being interrupted by the operating system is.

Several methods to overcome the large miss probability with short time slots have been suggested. Yarom and Falkner [56] monitor lines that are accessed in a loop. Their attack cannot distinguish between multiple consecutive accesses to the same line, but it can distinguish between periods of access and periods of no access to the line. Bengier et al. [3] and Van de Pol et al. [42] monitor memory lines that contain a call instruction. Such lines are accessed twice, once before the call and once upon return. Depending on the execution time of the called function, this approach can ensure that at most one of the two accesses is missed.

While these techniques can reduce or eliminate the probability of missing a victim access, they are not always applicable. In such scenarios, slowing the victim down can increase the interval between victim accesses and allow reducing the miss probability by using longer time slots. We describe such a scenario in the following sections.

6. ATTACKING OPENSLL

6.1 ECDSA

The ElGamal Signature Scheme [9] is the basis of the US 1994 NIST standard, Digital Signature Algorithm (DSA). The ECDSA is the adaptation of one step of the algorithm from the multiplicative group of a finite field to the group of points on an elliptic curve. The main benefit of using this group as opposed to the multiplicative group of a finite field is that smaller parameters can be used to achieve the same security level [24, 30] due to the fact that the current best known algorithms to solve the discrete logarithm problem in the finite field are sub-exponential and those used to solve the ECDLP are exponential — see Galbraith and Gaudry [11] and Koblitz and Menezes [25, Sec. 2-3] for an overview of recent ECDLP developments.

Parameters: An elliptic curve E defined over a finite field \mathbb{F}_q ; a point $G \in E$ of a large prime order n (generator of the group of points of order n). Parameters chosen as such are generally believed to offer a security level of \sqrt{n} given current knowledge and technologies. Parameters are recommended to be generated following the Digital Signature Standard [36]. The field size q is usually taken to be a large odd prime or a power of 2. The implementation of OpenSSL uses both prime fields and $q = 2^m$; the results in this paper relate to the former case.

Public-Private Key pairs: The private key is an integer d , $1 < d < n - 1$ and the public key is the point $Q = dG$. Calculating the private key from the public key requires solving the ECDLP, which is known to be hard in practice for correctly chosen parameters.

Signing: Suppose Bob, with private-public key pair $\{d_B, Q_B\}$, wishes to send a signed message m to Alice. He follows the following steps:

1. Using an approved hash algorithm, compute $e = \text{Hash}(m)$, take \bar{e} to be the leftmost ℓ bits of e (where $\ell = \min(\log_2(q), \text{bitlength of the hash})$).
2. Randomly select $1 \leq k \leq n - 1$.
3. Compute the point $(x, y) = kG \in E$.
4. Take $r = x \bmod n$; if $r = 0$ then return to **Step 2**.
5. Compute $s = k^{-1}(\bar{e} + rd_B) \bmod n$; if $s = 0$ then return to **Step 2**.
6. Bob sends (m, r, s) to Alice.

Verifying: The message m is not necessarily encrypted, the contents may not be secret, but a valid signature gives Alice strong evidence that the message was indeed sent by Bob. She verifies that the message came from Bob by:

1. Checking that all received parameters are correct, that $r, s \in [1, n - 1]$ and that Bob's public key is valid, that is $Q_B \neq 0$ and $Q_B \in E$ is of order n .
2. Using the same hash function and method as above, compute \bar{e} .
3. Compute $\bar{s} = s^{-1} \bmod n$.
4. Find the point $(x, y) = \bar{e}\bar{s}G + r\bar{s}Q_B$.
5. Verify that $r = x \bmod n$ otherwise reject the signature.

Step 2 of the signing algorithm is of vital importance — inappropriate reuse of the random integer led to the highly publicised

```

Input: Integer  $k \geq 1$ , width  $w$ 
Output:  $\text{mNAF}_w(k)$ 
 $i \leftarrow 0$ 
while  $k \geq 1$  do
  if  $k$  is odd then  $k_i \leftarrow k \bmod 2^w, k \leftarrow k - k_i$ 
  else  $k_i \leftarrow 0$ 
   $k \leftarrow k/2, i \leftarrow i + 1$ 
end
if  $k_{i-1} = 1$  and  $k_{i-1-w} < 0$  then
   $k_{i-1-w} \leftarrow k_{i-1-w} + 2^{w-1}$ 
   $k_{i-1} \leftarrow 0, k_{i-2} \leftarrow 1, i \leftarrow i - 1$ 
end
return  $(k_{i-1}, \dots, k_0)$ 

```

Figure 8: Generating modified Non-Adjacent Form for scalars. Here mods takes residues from $-(2^{w-1} - 1)$ to $2^{w-1} - 1$.

breaking of Sony PS3 implementation of ECDSA¹. Knowledge of the random value k , a.k.a. the *ephemeral key* or the *nonce*, leads to knowledge of the secret key. All values (m, r, s) can be observed by an eavesdropper, \bar{e} can be found from m , $r^{-1} \bmod n$ can be easily computed from n and r , and if k is discovered then an adversary can find Bob's secret key through the simple calculation

$$d_B = (sk - \bar{e})r^{-1}.$$

Our attack targets **Step 3** of OpenSSL's ECDSA signing algorithm.

6.2 ECC in OpenSSL

For ECDSA signing, the performance-critical component is scalar multiplication (**Step 3**) that, for an ℓ -bit integer k computes

$$kP = \sum_{i=0}^{\ell-1} k_i 2^i P$$

where k_i denotes bit i of k . Two key avenues for improving the performance of this operation are using a low-weight representation for the scalar, coupled with a scalar multiplication algorithm that interleaves elliptic curve additions and doublings, both with a goal of reducing the number of group operations. What follows is a description of how OpenSSL carries out this computation.

Scalar representation: The fact that group element inversion is cheap for elliptic curves makes signed representations for scalars a viable option: “subtraction of points on an elliptic curve is just as efficient as addition” [16, p. 98]. Generally, signed representations reduce the amount of needed precomputation by a factor of 2. A popular choice for ECC is Non-Adjacent Form (NAF) that, with a window width w represents k using digit set $\{0, \pm 1, \pm 3, \dots, \pm(2^{w-1} - 1)\}$ with the property that all non-zero digits are separated by at least $w - 1$ zeros, leading to lower average weight than other representations (e.g. binary). The modified version mNAF_w is otherwise the same but allows the most significant digit to violate the non-adjacent property if doing so decreases the length but keeps the same weight [32, Sec. 4.1]. It does so by applying the map $10^{w-1}\delta \mapsto 010^{w-2}\hat{\delta}$ if $\delta < 0$ where $\hat{\delta} = \delta + 2^{w-1}$ in the most significant digits. **Figure 8** illustrates the mNAF_w algorithm. See function `bn_compute_wNAF` in `crypto/bn/bn_intern.c` for OpenSSL's implementation of this

¹<http://arstechnica.com/gaming/2010/12/ps3-hacked-through-poor-implementation-of-cryptography/>

```

Input: Integer  $k \geq 1$ ,  $P \in E(\mathbb{F}_q)$ , width  $w$ 
Output:  $kP$ 
 $(k_{\ell-1} \dots k_0) \leftarrow \text{mNAF}_w(k)$ 
Precompute  $jP$  for all odd  $0 < j < 2^{w-1}$ 
 $Q \leftarrow k_{\ell-1}P$ 
for  $i \leftarrow \ell - 2$  to  $0$  do
   $Q \leftarrow 2Q$ 
  if  $k_i \neq 0$  then  $Q \leftarrow Q + k_i P$ 
end
return  $Q$ 

```

Figure 9: Left-to-right double-and-add scalar multiplication with mNAF_w signed representation

procedure. Lastly, it is worth noting that the most significant digit in NAF and mNAF_w for $k \geq 1$ is guaranteed to be positive.

Scalar multiplication: In the absence of any curve-specific routines, for curves over \mathbb{F}_p OpenSSL implements interleaved scalar multiplication by Möller [31, Sec. 3.2] — see scalar multiplication function `ec_wNAF_mul` in `crypto/ec/ec_mult.c` for OpenSSL’s implementation. While there are many paths through the code depending on inputs [4, Sec. 2.2], this work assumes the case of a single scalar input where no a priori precomputation structure is available. For this case, the function execution simplifies to a textbook left-to-right, double-and-add scalar multiplication routine — see e.g. Hankerson et al. [16, p. 100]. Figure 9 illustrates the algorithm that will perform ℓ point doublings and a number of point additions equaling the number of non-zero digits (minus the first point addition and plus the $2^{w-2} - 1$ point additions for ad hoc precomputation). Since point Q accumulates the partial scalar multiple, Q is termed the *accumulator*.

6.3 Attacking ECDSA

As mentioned, an attacker who knows the ephemeral key k used for a single signature (m, r, s) can obtain the secret key d_B from a simple calculation. It turns out that knowing a few bits of the nonces for *sufficiently many* signatures allows an attacker to obtain the secret key as well. One option is to embed the information for various signatures into a *lattice* such that the solution to a geometric lattice problem corresponds to the secret key [19, 34, 35].

But how does the attacker obtain any information on the ephemeral keys? As these keys are only used during the computation, a natural approach is to obtain this information through a side-channel attack. Unfortunately, using the side-channel described above to attack the wNAF implementation of the scalar multiplication does not directly reveal a fixed number of bits of every ephemeral key. This is due to the fact that the side-channel only reveals when the relevant operations take place, but in the case of an addition it does not show which value is being added. Previous works obtain information on the ephemeral key k from the double and add chains in different ways.

L1 cache targeting fixed lower bits: Brumley and Hakala [5] use the fact that the number of doubles after the last addition in the trace reveals an equal number of least significant bits of k : “From the side channel perspective, consecutive doublings allow inference of zero coefficients, and more than w point doublings reveals non-trivial zero coefficients” [5, Sec. 3.2]. They target signatures and traces that indicate a minimum of six zeros in the LSBs, in total requiring 2600 signatures and corresponding traces to recover the private key for curve `secp160r1` with a lattice attack [5, Sec. 6].

LLC targeting variable lower bits: The numerous drawbacks of the

previous attack include (1) discarding on average $1 - 2^{-6}$ percent of the traces; (2) limiting to SMT architectures like Intel’s HT; (3) rather noisy traces from the L1 data cache. Benger et al. [3] tackle all of these issues, while at the same time targeting the substantially larger and relevant curve `secp256k1`: “Prior work fixes a minimum value of [LSBs] and utilizes this single value in all equations . . . If we do this we would need to throw away [the majority] of the executions obtained. By maintaining full generality . . . we are able to utilize all information at our disposal” [3, Sec. 4]. As each trace reveals a different number of LSBs of the ephemeral key, they adjust the lattice problem accordingly and recover the private key with as little as 200 signatures and corresponding traces. However, to recover the private key with probability greater than 0.5, they require approximately 300 signatures.

LLC targeting full traces: Subsequently, Van de Pol et al. [42] show how to use roughly half of the double and add chain for group orders of a special form, i.e., $q = 2^n + \varepsilon$ where $|\varepsilon| < 2^p$ for $p \ll n$. It relies on the fact that the positions of adds in the chain reveal the positions of non-zero wNAF digits in the representation of k . Two adds are separated by at least w doubles, and every additional double reveals that the corresponding bits of k are repeating. However, a single bit of information is lost for every pair of consecutive non-zero wNAF digits, because these repeating bits of k are either zero or one depending on whether the second wNAF digit was positive or negative. Note that this method requires perfect traces, because each double is required to determine the bit position of the various additions. Therefore, whenever a double is missed, the whole trace preceding the missed double will produce inaccurate information, causing the subsequent lattice attack to fail.

6.4 Point Inversion: A New Leak

An implementation of scalar multiplication in Figure 9 requires accompanying control logic — in particular, to handle negative k_i digits. We observe the following trends in open source elliptic curve libraries for inverting points.

Invert on-the-fly: While the cost of elliptic curve point inversion can vary depending on the coordinate system choice, for many systems \mathbb{F}_p curves require only a finite field negation, i.e. flipping the sign of the y -coordinate. In these cases, since point inversion is so light many implementations opt for on-the-fly inversion. That is, when $k_i < 0$ compute $Q := Q + -(k_i P)$ inverting $k_i P$ to a temporary variable immediately preceding the point addition function call. For example, this is the approach taken by Bitcoin’s `libsecp256k1`². Scalar multiplication function `secp256k1_ecmult` in `src/ecmult_impl.h` calls macro `ECMULT_TABLE_GET_GE` which, in the case of a negative digit, calls `secp256k1_ge_neg` in `src/group_impl.h` to negate the point operand. The advantage to this approach is that it requires marginal additional storage overhead, and the disadvantage is that the algorithm will eventually end up inverting the same point more than once — duplicating a previously computed value.

Precompute inversions: As written, the precomputation table in Figure 9 requires storing 2^{w-2} points. Another strategy is to double the size of the table and additionally store the inverses of the required points. Then for negative digits, compute $Q := Q + (\hat{k}_i)P$ where \hat{k}_i is the table index for $-k_i$. Normally this will be handled in the NAF coding itself by yielding e.g. indices $(0, 1, 2, 3)$ corresponding to digits $(1, 3, -3, -1)$. For example, this is the approach taken by NSS³ — see scalar multiplication function

²<https://github.com/bitcoin/secp256k1>

³<https://developer.mozilla.org/en-US/docs/Mozilla/Projects/NSS>

`ec_GFp_pt_mul_jm_wNAF` in `lib/freebl/ec1/ecp_jm.c`. The advantage of this approach is that each point is only inverted a single time, and the disadvantage is that the required storage for the precomputation table doubles.

Invert the accumulator: Similar to the invert on-the-fly approach but without requiring a temporary point, another strategy is to track the sign of the accumulator in a variable and invert the accumulator as needed preceding point additions. That is, if the sign of the accumulator matches the sign of the digit $k_i \neq 0$, compute $Q := Q + |k_i|P$; otherwise $Q := -Q + |k_i|P$, inverting the accumulator before the point addition function call. Finally, after all digits are processed set $Q := -Q$ if the accumulator is in the inverted state. For example, this is the approach taken by OpenSSL — see scalar multiplication function `ec_wNAF_mul` in `crypto/ec/ec_mult.c` that calls `EC_POINT_invert` if *precisely* one of the following statements is true:

- the current (non-zero) digit is negative (variable `is_neg`);
- the accumulator is inverted (variable `r_is_inverted`).

All of the above approaches have potential side-channel issues — e.g. point inversion function calls dependent on secret sign bits or using secret digits as indices in memory-resident tables. We focus on the *invert the accumulator* approach since OpenSSL implements it. From the side-channel perspective, if we capture the sequence of elliptic curve point doublings, additions, and inversions for a particular k we can recover the signs of all non-zero k_i digits as follows. Denote the n inversions by $I_1 \dots I_n$ for $n > 1$. Note that n is always even since the accumulator (Q) always starts and ends in the non-inverted state, and (for completeness) that if $n = 0$ then all digits are positive. The accumulator toggles to the inverted state at I_1 , then back to the non-inverted state at I_2 , and so on — i.e. the accumulator enters the inverted state at I_j for odd j and non-inverted state for even j . Hence:

1. All the additions before I_1 correspond to positive digits.
2. For odd j , all the additions between I_j and I_{j+1} correspond to negative digits. This is due to the fact that the sign of the accumulator agrees with the sign of the current digit for such additions.
3. Similarly for even j , all the additions between I_j and I_{j+1} correspond to positive digits.

6.5 Exploiting the new leak

The improved side-channel described in this work allows us to determine whether the wNAF digits are positive or negative. This immediately gives one extra bit of information for each pair of consecutive adds in the top half of the double and add chain. Using the notation of Van de Pol et al. [42], if there are two consecutive adds at positions m and $m + l$ for $p < m < n - l$ we can write

$$k = a \cdot 2^{m+l+1} + b \cdot 2^{m+w} + c,$$

where $0 \leq a < 2^{n-m-l}$, $2^{m-1} < c < 2^{m+w} - 2^{m-1}$, and $b = 2^{l-w}$ if the second wNAF digit is positive or $b = 2^{l-w} - 1$ if it is negative.

Now, if we define $t = (r/s) \cdot 2^{n-m-l-1} \bmod q$ and $u = (b + 1/2) \cdot 2^{n+w-l-1} - (h/s) \cdot 2^{n-m-l-1} \bmod q$, it follows that $|d_B \cdot t - u|_q < q/2^{l-w+2}$, where $|\cdot|_q$ is reduction mod q into the range $[-q/2, q/2)$. Writing $z = l - w + 1$, each such triple (t, u, z) provides z bits of information about the secret key d_B , but conversely increases the dimension of the closest vector problem in the lattice by one.

Given d triples (t_i, u_i, z_i) such that $v_i = |d_B \cdot t_i - u_i|_q < q/2^{z_i+1}$, the secret key d_B can now be obtained by constructing

the following lattice problem: define a basis

$$\mathbf{B} = \begin{pmatrix} 2^{z_1+1} \cdot q & & & 2^{z_1+1} \cdot t_1 \\ & \ddots & & \vdots \\ & & 2^{z_d+1} \cdot q & 2^{z_d+1} \cdot t_d \\ & & & 1 \end{pmatrix},$$

where the columns give rise to a lattice L . Now consider the vector $\mathbf{u} = (2^{z_1+1}u_1, \dots, 2^{z_d+1}u_d, 0)$ and how close it lies to this lattice L . One lattice vector that is potentially close to \mathbf{u} is $\mathbf{B}\mathbf{x}$, where $\mathbf{x} = (\lambda_1, \dots, \lambda_d, d_B)$, which means the difference vector is given by

$$\begin{aligned} \mathbf{v} &= \mathbf{B}\mathbf{x} - \mathbf{u} \\ &= (2^{z_1+1}(d_B \cdot t_1 - u_1 + \lambda_1 q), \\ &\quad \dots, 2^{z_d+1}(d_B \cdot t_d - u_d + \lambda_d q), d_B). \end{aligned}$$

It is now possible to choose the coefficients λ_i such that the i th element of the vector \mathbf{v} has its elements within the range $[-2^{z_i}q, 2^{z_i}q)$, which implies that the difference vector is then given by $\mathbf{v} = (2^{z_1+1}v_1, \dots, 2^{z_d+1}v_d, d_B)$. Due to the bounds on the v_i coefficients, the norm of this difference vector can be bounded by $q\sqrt{d+1}$, whereas the volume of the lattice is given by $2^{(d+\sum_i z_i)q^d}$.

From the above argument we expect that, as $\sum_i z_i$ increases, the vector $\mathbf{B}\mathbf{x}$ becomes more and more likely to be the closest lattice vector to our target vector \mathbf{u} . As a result, solving the closest vector problem on L with target vector \mathbf{u} reveals the vectors \mathbf{v} and \mathbf{x} , both of which contain the key d_B as its final coefficient. This can be either solved through the embedding technique using a Shortest Vector Problem (SVP) approximation algorithm or by solving the Closest Vector Problem (CVP) directly with a CVP solver. For more information, see Van de Pol et al. [42].

To balance the hardness of the lattice problem with the information provided by the triples, we only used the 75 triples with the highest z -values which results in a lattice dimension of 76.

Table 3 shows the result of this attack for a varying number of signatures. It was implemented using the `fp111` library⁴ and executed on a single core of an Intel E5620 processor. Thus, given six error-free traces on different signatures allows an attacker to obtain the secret key in more than half the cases.

Table 3: Attack results for a given number of signatures

Signatures	Time (s)	Prob
4	15.08	0.005
5	13.94	0.165
6	12.51	0.545
7	11.50	0.735
8	9.69	0.840

7. AMPLIFICATION ATTACK

In the previous section, we identified a new leak in the OpenSSL implementation of ECDSA signing and analysed the leak under the assumption that the adversary can obtain a perfect trace of the victim's operations. In this section, we investigate the practical issues of obtaining perfect traces. We first look at why error-free traces are required for the attack. We proceed with describing how past research used the FLUSH+RELOAD attack to achieve a high probability of obtaining perfect traces. We then explain why these techniques are not sufficient when we want to capture inversions. We

⁴<http://perso.ens-lyon.fr/damien.stehle>

show that amplification allows us to overcome the limitations and demonstrate how to use it to obtain perfect traces.

7.1 The need for perfect traces

Suppose that the adversary manages to obtain an almost-perfect trace. That is, she knows the sequence of operations taken by the victim with the exception of a single error that causes it to either miss a double operation or add a spurious one. When inferring the positions of the non-zero wNAF digits, the error will propagate through the representation of the scalar, changing the position of all digits to the left of the error, which are the positions we use for the lattice attack. Consequently, the lattice attack will receive an erroneous input and will fail to find the key. Similarly, if the trace misses an inversion or contains a spurious one, the signs of any digit above the error locations are incorrect.

Even if the adversary knows or suspects that an error has occurred, correcting the error poses problems. If, for example, the adversary notes that the time between two operations in the trace is longer than expected, the adversary can suspect an error. However, because the victim may have been suspended while the processor executed some system function, the adversary cannot be certain that an error occurred.

The adversary could try to use known properties of perfect traces to identify and possibly correct errors in captured traces. However, there is very little information that the adversary can use. In particular, the adversary does not know for certain the number and position of point addition operations. She can detect, but not correct, errors like: (1) the number of point inversions must be even; (2) at least w point doublings must separate point additions. Finally, even though all the scalars used in the multiplication have the same bit length (due to a timing attack [6] resulting in CVE-2011-1945), the length of the wNAF representation may vary. For example, we look at the numbers 228 and 229. The binary representations of these is 11100100 and 11100101, i.e. both are 8 bit numbers. The 4-NAF representation of 228 is 1, 0, 0, 0, 0, 0, -7, 0, 0. That is, $228 = 1 \cdot 2^8 - 7 \cdot 2^2$. The representation of 229 is 7, 0, 0, 0, 0, 5 - 229 = $7 \cdot 2^5 + 5$. Hence, while the bit length of both numbers is 8, the length of their 4-NAF representations are 9 and 6. Consequently, the number of double operations in the trace is not fixed.

As we can see, the effects of errors in the trace are not localised, errors are hard to detect, and are almost impossible to correct. Combined with the sensitivity of the lattice attack to errors, every small error in the captured trace significantly reduces the probability of attack success. In particular, unless the adversary can get enough error free traces, she will not be able to apply the attack.

7.2 Obtaining perfect traces

Van de Pol et al. [42] attack the same implementation that we target. Unlike us, they do not try to trace the accumulator inversions, focusing instead on add and double operations. They divide time into slots of 1,200 cycles and probe memory lines within the functions that implement the group add and group double operations. As Section 5 demonstrates, with slots of 1,200 cycles the expected miss rate is around 44%.

To reduce the miss probability, they choose memory lines that contain a call to a field multiplication operation. As discussed above, the victim accesses memory lines that contain a call twice; once when executing the call and once when the call returns. Because these two accesses are related, their times are not independent and the probability of missing each is not independent of each other. Consequently, they manage to reduce the number of capture errors to 1 in 1,000 group operations. With around 300 operations

in trace, the probability of capturing an error-free trace is 58%.

For our attack, we need to further trace accumulator inversions along with group addition and double. While the group inversion code contains a call instruction, we cannot probe the memory line that contains it. The reason is that due to speculative execution, all of the code up to the call instruction is prefetched into the cache, even if the execution does not take the path. As a result, monitoring these memory lines would result in a large number of false positives. Therefore, we have to monitor memory lines that follow the call to the field negation, which do not contain additional call instructions.

In our environment (OpenSSL 1.0.2a running on an HP Elite 8300 running CentOS 6.5 64 bit), add operations take on average 3,223 cycles and double operations take 3,427 cycles. As Bengier et al. [3] discuss, the maximum slot length we can use is about half the length of the operations, or 1,600 cycles. With these time slots, the probability of missing the victim access to the memory line in the inversion code is about 33%. With such an error probability, and an expected number of 25 inversions in each scalar multiplication, the probability of capturing a perfect trace is less than 1/25,000, which is way too low for a practical attack.

One possibility of reducing the miss probability when tracing accumulator inversions is to monitor two memory lines within the code. The scalar multiplication code in OpenSSL invokes the generic elliptic curve point inversion function `EC_POINT_invert`. The function invokes the curve-specific point inversion function, which in the case of the `secp256k1` curve is `ec_GFp_simple_invert`. Said function invokes field subtraction (`BN_usub`) to negate the y component of the point. By probing the memory lines following the return of `BN_usub` and the return of `ec_GFp_simple_invert` we get the same effect as probing a memory line that contains a `call` instruction, with the adversary missing at most one of these accesses.

While this approach guarantees that the adversary does not miss accumulator inversions, it requires the adversary to monitor four memory lines: one in each of the double and add functions and two in the inversion functions. Each probe takes about 450 cycles, so probing four memory lines takes 1,800 cycles. When we set the slot size to 1,800 cycles, the traces loses accuracy because we can no longer determine the order of some of the operations in the sequence.

Increasing the slot length would allow us to consistently trace all accumulator inversions, however the speed of calculating the group addition and doubling limits the maximum slot length. To increase the limit, we can try slowing the group operations down.

7.3 A performance-degradation attack against OpenSSL

We use the performance-degradation attack to slow the group operations down. We target the `bn_mul_mont` function, which implements the field multiplication and square. We use one attacking thread and check the effect of repeatedly evicting three memory lines in the main loop of the function. Table 4 summarises the run time of the add and the double operations under the attack. As we can see, the attack slows the group add operation by a factor of 53 and the double operation by a factor of 41.

Table 4: `secp256k1` group operation times (cycles)

	Add	Double
No attack	2,892	3,086
Under attack	153,709	126,282

With group operations taking over 100,000 cycles, we can safely increase the slot size and monitor the four memory lines required for obtaining the trace. We set the slot size to 17,000 cycles and captured 1,000 traces. Comparing the traces to the ground truth we find that only five of them show errors. Hence, our attack captures error-free traces almost every time. We can now use these traces with the lattice attack of Section 6, to break the long-term ECDSA key of the victim after observing as few as six signatures.

Table 5 compares the results of this work with previous cache-based attacks on OpenSSL ECDSA. As we can see, the attack requires less than a quarter of the previous best attack. About half of this improvement is due to exploiting the leak of point inversion and the other half comes from the increased accuracy of observing the side-channel. Employing the performance-degradation attack to amplify the side-channel underpins both these improvements.

Table 5: OpenSSL ECC cache-timing attack results compared

Curve	Source	Perfect traces	Signatures
secp160r1	Brumley and Hakala [5]	-	2600
secp256k1	Benger et al. [3]	-	300
secp256k1	Van de Pol et al. [42]	13	25
secp256k1	This work	6	6

8. OTHER CIPHERS

We have demonstrated the utility of the performance-degradation attack for amplifying the side-channel information leakage from the OpenSSL implementation of elliptic curves over prime fields. To demonstrate that the technique is likely to have more general applications, we measure the slowdown we achieve with several public-key ciphers.

Table 6 provides a short summary of the results. More detailed results are available in Appendix A. As we can see, the attack is effective in slowing down all of the tested implementations. We can therefore use the attack to amplify side-channel leaks from these implementations.

We should note that the performance-degradation attack results do not imply that the cipher implementations are vulnerable to side-channel attacks. While we believe that most of the tested implementations may leak information under certain circumstances, we do not claim that all of these leaks are exploitable. Pereira et al. [41] use our performance-degradation technique for their attack on OpenSSL DSA. We leave exploiting other leaks to future work.

Table 6: Slowdown of cipher implementations

Library	Algorithm	Slowdown
OpenSSL	Exponentiation	7–158
OpenSSL	EC Scalar Multiplication	7–63
libcrypt	Exponentiation	10–21
libcrypt	EC Scalar Multiplication	8–13

9. CONCLUSION

Typical performance-degradation attacks usually do not provide any direct benefit to the attacker. Their main benefit is derived indirectly, through the damage they cause to the victim. In this paper we demonstrate that these attacks can offer tangible benefits to the attacker—they can be used to amplify a side-channel, allowing the attacker to receive more information through the channel than was otherwise possible.

To demonstrate side-channel amplification, we first investigate using cache evictions as a microarchitectural attack vector, showing that is over 8 times more potent than previously published attacks. We further identify a new information leak in the OpenSSL implementation of the ECDSA signature scheme. Lastly, we show how using the performance-degradation attack to amplify a cache side channel allows the attacker to exploit the information leak. Our combined attack allows the adversary to completely cryptanalyse the secp256k1 elliptic curve used in Bitcoin after observing the side channel over only 6 signatures, less than a quarter of any prior result.

Acknowledgements

We would like to thank Dr Naomi Benger for the useful discussions, advice and support. We would also like to thank Camilla Beck and Diclehan Erdal for performing some of the experiments for this work.

Parts of this research were performed under contract to the Defence Science and Technology Group, Maritime Division, Australia.

This research was supported in part by COST Action IC1306.

The second author was supported in part by TEKES grant 4681/31/2014 INKA EAKR Hardware Rooted Security.

The fourth author was supported in part by EPSRC via grant EP/I03126X.

References

- [1] Onur Aciçmez, Billy Bob Brumley, and Philipp Grabher. New results on instruction cache attacks. In *CHES*, Santa Barbara, CA, US, 2010.
- [2] Andrea Arcangeli, Izik Eidus, and Chris Wright. Increasing memory density by using KSM. In *2009 Ottawa Linux Symp.*, pages 19–28, Montreal, Quebec, Canada, Jul 2009.
- [3] Naomi Benger, Joop van de Pol, Nigel P. Smart, and Yuval Yarom. “Ooh aah... just a little bit”: A small amount of side channel can go a long way. In *CHES*, pages 75–92, Busan, KR, Sep 2014.
- [4] Billy Bob Brumley. Faster software for fast endomorphisms. In *6th COSADE*, pages 127–140, Berlin, DE, Apr 2015.
- [5] Billy Bob Brumley and Risto M. Hakala. Cache-timing template attacks. In *15th ASIACRYPT*, pages 667–684, Tokyo, JP, Dec 2009.
- [6] Billy Bob Brumley and Nicola Tuveri. Remote timing attacks are still practical. In *16th ESORICS*, Leuven, BE, 2011.
- [7] Carlos Cardenas and Rajendra V Boppana. Detection and mitigation of performance attacks in multi-tenant cloud computing. In *1st International IBM Cloud Academy Conference*, Research Triangle Park, NC, US, 2012.
- [8] Ron C. Chiang, Sundaresan Rajasekaran, Nan Zhang, and H. Howie Huang. Swiper: Exploiting virtual machine vulnerability in third-party clouds with competition for I/O resources. *Trans. Parall. & Distr. Syst.*, 26(6):1732–1742, Jun 2015.
- [9] Taher ElGamal. A public key cryptosystem and a signature scheme based on discrete logarithms. In *Advances in Cryptology*, Santa Barbara, CA, US, 1985.

- [10] Agner Fog. The microarchitecture of Intel, AMD and VIA CPUs: An optimization guide for assembly programmers and compiler makers. <http://www.agner.org/optimize/microarchitecture.pdf>, Aug 2014.
- [11] Steven D. Galbraith and Pierrick Gaudry. Recent progress on the elliptic curve discrete logarithm problem. *DCC*, 78(1), Jan 2016.
- [12] David Goldberg. What every computer scientist should know about floating-point arithmetic. *Comput. Surveys*, 23(1):6–48, Mar 1991.
- [13] Dirk Grunwald and Soraya Ghiasi. Microarchitectural denial of service: Insuring microarchitectural fairness. In *35th MICRO*, pages 409–418, Istanbul, TR, Nov 2002.
- [14] Daniel Gruss, Raphael Spreitzer, and Stefan Mangard. Cache template attacks: Automating attacks on inclusive last-level caches. In *24th USENIX Security*, pages 897–912, Washington, DC, US, Aug 2015.
- [15] David Gullasch, Endre Bangerter, and Stephan Krenn. Cache games – bringing access-based cache attacks on AES to practice. In *S&P*, pages 490–505, Oakland, CA, US, 2011.
- [16] Darrel Hankerson, Alfred Menezes, and Scott Vanstone. *Guide to elliptic curve cryptography*. Springer Professional Computing, 2004.
- [17] Jahangir Hasan, Ankit Jalote, T. N. Vijaykumar, and Carla E. Brodley. Heat stroke: Power-density-based denial of service in SMT. In *11th HPCA*, pages 166–177, San Francisco, CA, US, Feb 2005.
- [18] J.L. Henning. SPEC CPU2006 benchmark descriptions. *Comp. Arch. News*, 34(4), Sep 2006.
- [19] Nick Howgrave-Graham and Nigel P. Smart. Lattice attacks on digital signature schemes. *DCC*, 23(3):283–290, Aug 2001.
- [20] Mehmet Sinan İnci, Berk Gülmezoğlu, Gorka Irazoqui, Thomas Eisenbarth, and Berk Sunar. Seriously, get off my cloud! Cross-VM RSA key recovery in a public cloud. IACR Cryptology ePrint Archive, Report 2015/898, Sep 2015.
- [21] Intel 64 & IA-32 AORM. *Intel 64 and IA-32 Architectures Optimization Reference Manual*. Intel Corporation, Apr 2012.
- [22] Gorka Irazoqui, Mehmet Sinan İnci, Thomas Eisenbarth, and Berk Sunar. Wait a minute! a fast, cross-VM attack on AES. In *RAID*, pages 299–319, Gothenburg, Sweden, Sep 2014.
- [23] Gorka Irazoqui, Mehmet Sinan İnci, Thomas Eisenbarth, and Berk Sunar. Lucky 13 strikes back. In *ASIA CCS*, pages 85–96, Singapore, Apr 2015.
- [24] Neal Koblitz. Elliptic curve cryptosystems. *Mathematics Comput.*, 48(177):203–209, Jan 1987.
- [25] Neal Koblitz and Alfred Menezes. A riddle wrapped in an enigma. IACR Cryptology ePrint Archive, Report 2015/1018, Nov 2015.
- [26] Fangfei Liu, Yuval Yarom, Qian Ge, Gernot Heiser, and Ruby B Lee. Last-level cache side-channel attacks are practical. In *S&P*, pages 605–622, San Jose, CA, US, May 2015.
- [27] Stefan Mangard, Norbert Pramstaller, and Elisabeth Oswald. Successfully attacking masked AES hardware implementations. In *CHES*, pages 157–171, Edinburgh, UK, Aug 2005.
- [28] Jeanna Neefe Matthews, Wenjin Hu, Madhujith Hapuarachchi, Todd Deshane, Demetrius Dimatos, Gary Hamilton, Michael McCabe, and James Owens. Quantifying the performance isolation of virtualization systems. In *WS Experimental Comp. Sci.*, San Diego, CA, US, Jun 2007.
- [29] Clémentine Maurice, Nicolas Le Scouarnec, Christoph Neumann, Olivier Heen, and Aurélien Francillon. Reverse engineering Intel last-level cache complex addressing using performance counters. In *RAID*, Kyoto, Japan, Nov 2015.
- [30] Victor S. Miller. Use of elliptic curves in cryptography. In *CRYPTO’85*, pages 417–426, Santa Barbara, CA, US, Aug 1985.
- [31] Bodo Möller. Algorithms for multi-exponentiation. In *SAC*, pages 165–180, Toronto, ON, CA, Aug 2001.
- [32] Bodo Möller. Improved techniques for fast exponentiation. In *Inform. Security & Cryptology*, pages 298–302, Seoul, KR, Nov 2002.
- [33] Thomas Moscibroda and Onur Mutlu. Memory performance attacks: Denial of memory service in multi-core systems. In *16th USENIX Security*, Boston, MA, US, 2007.
- [34] Phong Q. Nguyen and Igor E. Shparlinski. The insecurity of the digital signature algorithm with partially known nonces. *J. Cryptology*, 15(2):151–176, Jun 2002.
- [35] Phong Q. Nguyen and Igor E. Shparlinski. The insecurity of the elliptic curve digital signature algorithm with partially known nonces. *DCC*, 30(2):201–217, Sep 2003.
- [36] NIST FIPS PUB 186-4. *Digital Signature Standard (DSS)*. NIST, 2013.
- [37] Siddika Berna Örs, Frank Gürkaynak, Elisabeth Oswald, and Bart Preneel. Power-analysis attack on an ASIC AES implementation. In *ITCC 2004*, volume 2, pages 546–552, Las Vegas, NV, US, Apr 2004.
- [38] Dag Arne Osvik, Adi Shamir, and Eran Tromer. Cache attacks and countermeasures: the case of AES. <http://www.cs.tau.ac.il/~tromer/papers/cache.pdf>, Nov 2005.
- [39] Daniel Page. *Practical Introduction to Computer Architecture*. Texts in Computer Science, 2009.
- [40] Colin Percival. Cache missing for fun and profit. In *BSDCan 2005*, Ottawa, CA, 2005.
- [41] César Pereida, Billy Bob Brumley, and Yuval Yarom. “Make sure DSA signing exponentiations really are constant-time”. In *23rd CCS*, Vienna, AT, Oct 2016.
- [42] Joop van de Pol, Nigel P. Smart, and Yuval Yarom. Just a little bit more. In *2015 CT-RSA*, pages 3–21, San Francisco, CA, USA, Apr 2015.
- [43] Andre Richter, Christian Herber, Holm Rauchfuss, Thomas Wild, and Andreas Herkersdorf. Performance isolation exposure in virtualized platforms with PCI passthrough I/O sharing. In *Architecture of Computing Systems*, pages 171–182, 2014.

- [44] Thomas Ristenpart, Eran Tromer, Hovav Shacham, and Stefan Savage. Hey, you, get off of my cloud: Exploring information leakage in third-party compute clouds. In *16th CCS*, pages 199–212, Chicago, IL, US, 2009.
- [45] Vineet Sinha, David Karger, and Rob Miller. Relo: Helping users manage context during interactive exploratory visualization of large codebases. In *OOPSLA Workshop on Eclipse Technology eXchange (ETX)*, pages 21–25, San Diego, CA, US, Oct 2005.
- [46] Kuniyasu Suzaki, Kengo Iijima, Toshiki Yagi, and Cyrille Artho. Memory deduplication as a threat to the guest OS. In *4th European Workshop on System Security*, Salzburg, AT, 2011.
- [47] Augustus K. Uht and Vijay Sindagi. Disjoint eager execution: An optimal form of speculative execution. In *28th MICRO*, pages 313–325, Nov 1995.
- [48] Venkatanathan Varadarajan, Thawan Kooburat, Benjamin Farley, Thomas Ristenpart, and Michael M Swift. Resource-freeing attacks: improve your cloud performance (at your neighbor’s expense). In *19th CCS*, Raleigh, NC, US, 2012.
- [49] Carl A. Waldspurger. Memory resource management in VMware ESX server. In *5th OSDI*, Boston, MA, US, 2002.
- [50] Colin D. Walter. Longer keys may facilitate side channel attacks. In *SAC*, pages 42–57, Waterloo, ON, Canada, Aug 2004.
- [51] Chuliang Weng, Jianfeng Zhan, and Yuan Luo. TSAC: Enforcing isolation of virtual machines in clouds. *Trans. Computers*, 64(5):1470–1482, May 2015.
- [52] Dong Hyuk Woo and Hsien-Hsin S. Lee. Analyzing performance vulnerability due to resource denial of service attack on chip multiprocessors. In *WS Chip Multiprocessor Memory Syst. & Interconnects*, Phoenix, AZ, US, 2007.
- [53] Carole-Jean Wu and Margaret Martonosi. Characterization and dynamic mitigation of intra-application cache interference. In *Int. Symp. Performance Analysis Syst. & Softw.*, ISPASS ’11, Austin, TX, US, 2011.
- [54] Tianni Xu, Xiufeng Sui, Zhicheng Yao, Jiuyue Ma, Bao Yungang, and Lixin Zhang. Rethinking virtual machine interference in the era of cloud applications. In *15th HPCC*, pages 190–197, Zhangjiajie, Hunan, China, Nov 2013.
- [55] Yuval Yarom and Naomi Benger. Recovering OpenSSL ECDSA nonces using the FLUSH+RELOAD cache side-channel attack. IACR Cryptology ePrint Archive, Report 2014/140, Feb 2014.
- [56] Yuval Yarom and Katrina Falkner. FLUSH+RELOAD: a high resolution, low noise, L3 cache side-channel attack. In *23rd USENIX Security*, pages 719–732, San Diego, CA, US, 2014.
- [57] Yuval Yarom, Qian Ge, Fangfei Liu, Ruby B. Lee, and Gernot Heiser. Mapping the Intel last-level cache. <http://eprint.iacr.org/>, Sep 2015.
- [58] Yinqian Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. Cross-VM side channels and their use to extract private keys. In *19th CCS*, pages 305–316, Raleigh, NC, US, Oct 2012.
- [59] Yinqian Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. Cross-Tenant side-channel attacks in PaaS clouds. In *21st CCS*, Scottsdale, AZ, US, 2014.

Table 7: Performance-degradation attack on OpenSSL and libcrypto

Library	Algorithm	Details	Operation	Normal	Degraded	Slowdown
OpenSSL	Constant-time exponentiation	1024 bits	power-5	7161	128650	18.0
		2048 bits	power-5	26215	653372	24.9
	Sliding-window exponentiation	512 bits	Square	391	2773	7.1
		512 bits	Multiply	437	37182	85.1
		1024 bits	Square	1181	19284	16.3
		1024 bits	Multiply	1520	186342	122.6
		2048 bits	Square	4129	129787	31.4
		2048 bits	Multiply	5849	921943	157.6
	ECC over prime fields	secp112r1	Double	1905	13917	7.3
		secp112r1	Add	1770	18338	10.4
		secp256k1	Double	3086	126282	40.9
		secp256k1	Add	2892	153709	53.1
		secp521r1	Double	6830	428443	62.7
		secp521r1	Add	8065	579372	71.8
	ECC over binary fields	sect113r1	Double	1303	26987	20.7
		sect113r1	Add	1206	22553	18.7
		sect283r1	Double	3177	84887	26.7
		sect283r1	Add	3402	71571	21.0
sect571k1		Double	5096	150070	29.4	
sect571k1		Add	7353	151913	20.7	
libcrypto	Modular exponentiation	512 bits	Multiply	1365	28996	21.2
		1024 bits	Multiply	3813	59568	15.6
		2048 bits	Multiply	11150	110420	9.9
	Elliptic curves	brainpoolP160r1	Double	13144	111155	8.5
		brainpoolP160r1	Add	18530	147401	8.0
		Ed25519	Double	10164	90945	8.9
		Ed25519	Add	14746	143850	9.8
		secp256r1	Double	9342	94603	10.1
		secp256r1	Add	16401	156736	9.6
		secp521r1	Double	22797	288595	12.7
		secp521r1	Add	39737	499954	12.3

APPENDIX

A. CIPHER SLOWDOWN RESULTS

Table 7 shows the effect of our performance-degradation attack on the core operations in multiple implementations of public-key primitives. We look at two popular cryptographic libraries: OpenSSL version 1.0.2h and libcrypto version 1.6.4.

We time the individual operations used during modular exponentiation and during elliptic curves scalar multiplications. For all of the attacks we used a single attacking thread, evicting code lines within the inner loops of the targeted operations. For each of the operations we show the average time to complete both when at-

tacked and when not attacked.

As the data demonstrates, the attack is effective at slowing the operations down, achieving a slowdown factor of between 7 and 157. With this level of amplification potential vulnerabilities in the implementations may become easier to exploit.

We note that in the majority of cases, the larger the security parameters are, the more effective the performance-degradation attack is. The reason for that is that the time complexity of the targeted code is often quadratic in the size of the input. Consequently, as the input size increases, the relative proportion of the targeted code in the operation increases and the slowdown is more noticeable.