

CACHE MISSING FOR FUN AND PROFIT

COLIN PERCIVAL

ABSTRACT. We describe the construction of a channel between processes via the state of a shared memory cache, and its use in the cryptanalysis of RSA. Unlike earlier side-channel attacks involving memory caches, our attack has the remarkable property of only requiring that a single private key operation be observed.

We also discuss other methods in which this channel might be abused, and provide some suggestions to processor designers, operating system vendors, and the authors of cryptographic software as to how this and related attacks could be mitigated or eliminated entirely.

1. INTRODUCTION

As integrated circuit fabrication technologies have improved over the past few decades, improvements in processor performance have vastly outpaced those in memory latency; while accessing a random location in RAM might have taken a few processor cycles two decades ago, it can now easily take several hundred cycles. The answer to this processor-memory performance gap has been to add caches: By relying upon the principles of temporal and spatial locality, it has been possible to keep the *average* cost of a memory access reasonably constant relative to the cost of arithmetic operations, even though the worst case has degraded significantly.

The improvement in average performance due to caches comes at the expense of a vastly increased variability in performance. This has been known for many years to cause problems in the design of safety-critical “real time” systems where it is imperative that a series of deadlines be satisfied even as the presence of caches makes it very difficult to determine the worst-case performance [15]. More recently, it has been shown that the presence of caches and the resulting timing variability makes possible a number of cryptanalytic side channel attacks [1, 14, 17].

Key words and phrases. Cryptography, RSA, side channels, simultaneous multi-threading, caching.

The problems introduced by caches have been further exacerbated by the current trend towards increased parallelism. On recent processors implementing simultaneous multithreading [18], such as Intel’s “Hyper-Threading” processors [11], access to the L1 cache is shared between two independent instruction streams; this makes it non-trivial to even optimize the average performance, as frequently-used data needs to be located carefully to minimize the competition between threads for “popular” cache sets [7]. Of more cryptanalytic interest, the ability of one thread to influence the state of the cache encountered by another thread makes possible the step from *simple* cache-based timing attacks, where the time needed to perform cryptographic operations is measured while varying only the data used, to *differential* cache-based timing attacks, where the time needed to perform cryptographic operations is measured while varying the initial state of the cache [13].

In addition to allowing an attacker to influence the timing of a victim computation, the sharing of caches creates a channel allowing information to be transmitted from a victim to an attacker. Taking the 2.8 GHz Intel Pentium 4 with Hyper-Threading processor as an example (for reasons of availability only – we have no reason to expect it to be any more vulnerable to attack than other systems with shared caches), we will demonstrate first that a high bandwidth covert channel can be constructed using a shared cache, and second that this channel can be used cryptanalytically to attack RSA computations.

Independently of and concurrent with this work, a closely related attack against the Advanced Encryption Standard symmetric cipher (AES, [12]) has been discovered by Osvik, Shamir, and Tromer, who refer to it as the “Prime + Probe” method [13].

2. COVERT COMMUNICATION VIA PAGING

To see how shared caches can create a cryptographic side-channel, we first step back for a moment to a simpler problem — covert channels [10] — and one of the classic examples of such a channel: virtual memory paging.

Consider two processes, known as the Trojan process and the Spy process, operating at different privilege levels on a multilevel secure system, but both with access to some large reference file (naturally, on a multilevel secure system this access would necessarily be read-only). The Trojan process now reads a subset of pages in this reference file, resulting in page faults which load the selected pages from disk into memory. Once this is complete (or even in the middle of this operation) the Spy process reads *every* page of the reference file and measures the

time taken for each memory access. Attempts to read pages which have been previously read by the Trojan process will complete very quickly, while those pages which have not already been read will incur the (easily measurable) cost of a disk access. In this manner, the Trojan process can repeatedly communicate one bit of information to the Spy process in the time it takes for a page to be loaded from disk into memory, up to a total number of bits equal to the size (in pages) of the shared reference file.

If the two processes do not share any reference file, this approach will not work, but instead an opposite approach may be taken: Instead of faulting pages *into* memory, the Trojan process can fault pages *out* of memory. Assume that the Trojan and Spy processes each have an address space of more than half of the available system memory and the operating system uses a least-recently-used page eviction strategy. To transmit a “one” bit, the Trojan process reads its entire address space; to transmit a “zero” bit, the Trojan process spins for the same amount of time while only accessing a single page of memory. The Spy process now repeatedly measures the amount of time needed to read its entire address space. If the Trojan process was sending a “one” bit, then the operating system will have evicted pages owned by the Spy process from memory, and the necessary disk activity when those pages are accessed will provide an easily measurable time difference. While this covert channel has far lower bandwidth than the previous channel — it operates at a fraction of a bit per second, compared to a few hundred bits per second — it demonstrates how a shared cache can be used as a covert channel, even if the two communicating processes do not have shared access to any potentially cached data.

3. L1 CACHE MISSING

The L1 data cache in the Pentium 4 consists of 128 cache lines of 64 bytes each, organized into 32 4-way associative sets. This cache is completely shared between the two execution threads; as such, each of the 32 cache sets behaves in the same manner as the paging system discussed in the previous section: The threads cannot communicate by loading data *into* the cache, since no data is shared between the two threads, but they can communicate via the cache metadata by forcing each other’s data *out* of the cache.

A covert channel can therefore be constructed as follows: The Trojan process allocates an array of 2048 bytes, and for each 32-bit word it wishes to transmit, it accesses byte $64i$ of the array iff bit i of the word is set. The Spy process allocates an array of 8192 bytes, and repeatedly

measures the amount of time needed to read bytes $64i$, $64i + 2048$, $64i + 4096$, and $64i + 6144$ for each $0 \leq i < 32$. Each memory access performed by the Trojan will evict a cache line owned by the Spy, resulting in lines being reloaded from the L2 cache, which adds an additional latency of approximately 30 cycles if the memory accesses are dependent. This alone would not be measurable, thanks to the long latency of the RDTSC (read time stamp counter) instruction, but this difficulty is easily resolved by adding some high-latency instructions – for example, integer multiplications – into the critical path. In Figure 1 we show an example of how the Spy process could measure and record the amount of time required to access all the cache lines of each set.

Using this code, 32 bits can be reliably transmitted from the Trojan to the Spy in roughly 5000 cycles with a bit error rate of under 25%; using an appropriate error correcting code, this provides a covert channel of 400 kilobytes per second on a 2.8 GHz processor.

4. L2 CACHE MISSING

The same general approach is effective in respect of the L2 cache, with a few minor complications. The Pentium 4 L2 cache (on the particular model which we are examining) consists of 4096 cache lines of 128 bytes each, organized into 512 8-way associative sets. However, the data TLB holds only 64 entries — only enough to provide address mappings for half of the cached data. As a result, a Spy process operating in the same manner as described in the previous section will incur the cost of TLB misses on at least some of its memory accesses. To avoid allowing this to add noise to the measurements, we can resort to ensuring that *every* memory access incurs the cost of a TLB miss, by accessing each of the 128 pages (512 kB divided by 4 kB per page) before returning to the first page and accessing the second cache line it contains. (Another option would use a buffer of 16 MB, placing each potentially cached line into a separate page, but accessing the lines in a suitable order is just as effective.) Since the TLB entries have to be repeatedly reloaded, however, we also experience some additional cache misses, as the memory holding the paging tables will be repeatedly reloaded into the cache. Fortunately, this will only affect a small number of cache lines, leaving the vast majority acting as a fully operational covert channel.

Another complication is introduced by the design of the Pentium 4 as a streaming processor. The “Advanced Transfer Cache” includes a capability for hardware prefetching: If a series of cache misses occur, in arithmetic progression, within a single page, then the cache will

```
mov ecx, start_of_buffer
sub length_of_buffer, 0x2000
rdtsc
mov esi, eax
xor edi, edi

loop:
prefetcht2 [ecx + edi + 0x2800]

add cx, [ecx + edi + 0x0000]
imul ecx, 1
add cx, [ecx + edi + 0x0800]
imul ecx, 1
add cx, [ecx + edi + 0x1000]
imul ecx, 1
add cx, [ecx + edi + 0x1800]
imul ecx, 1

rdtsc
sub eax, esi
mov [ecx + edi], ax
add esi, eax
imul ecx, 1

add edi, 0x40
test edi, 0x7C0
jnz loop

sub edi, 0x7FE
test edi, 0x3E
jnz loop

add edi, 0x7C0
sub length_of_buffer, 0x800
jge loop
```

FIGURE 1. Example code for a Spy process monitoring the L1 cache on an Intel Pentium 4 with Hyper-Threading processor.

“recognize” this as a data stream and prefetch two additional cache lines. This is quite effective for reducing cache misses; but since we instead want to maximize cache misses, it becomes a disadvantage. Here we can simply trust to luck: If we access cache lines in an irregular manner (e.g., following a *de Bruijn* cycle rather than accessing the lines in increasing address order), then it is unlikely that we will exhibit three or more cache misses in arithmetic progression, and the hardware prefetcher will not activate.

Finally, since the L2 cache is used for both data and code, there will be some inevitable cache collisions (and line evictions) caused by the instruction fetching activity.

Due to the lower memory bandwidth, increased size of L2 cache sets (8 lines of 128 bytes, vs. 4 lines of 64 bytes in the L1 cache), and noise introduced by memory activity associated with TLB misses and instruction fetching, the L2 cache provides a significantly lower bandwidth covert channel than the L1 cache. In roughly 350000 cycles (on the same machine as previously used), 512 bits can be transmitted with an error rate of under 25%; this provides a channel of approximately 100 kilobytes per second.

Despite the reduced bandwidth, however, the L2-collision covert channel is potentially more interesting than the L1-collision channel: On systems without shared caches — i.e., where the Trojan and Spy processes are always separated by context switches — the contents of the L1 cache will tend to be fairly comprehensively replaced between schedulings of the Spy process; the L2 cache however, due to its larger size, is often not completely replaced, allowing it to be used as a covert channel with a bandwidth of several bits per context switch. On an otherwise quiescent system this could easily provide a covert channel of a few kilobits per second, and several times that if the kernel makes the POSIX `sched_yield(2)` system call [6] available or if there is some other mechanism allowing very frequent context switching to be obtained.

5. OPENSsl KEY THEFT

Having demonstrated the effectiveness of this cache-missing approach in the construction of a covert channel, we now examine it as a crypt-analytic side channel. Taking as a demonstration platform OpenSSL 0.9.7c [16] running on FreeBSD 5.2.1-RELEASE-p13 [4], we performed a 1024-bit private RSA operation (via the command `openssl rsautl -inkey priv.key -sign`), while running the L1 Spy process described

in Section 3. To simplify the attack, we started running the Spy process before we started OpenSSL, stopped it after the OpenSSL process had completed, and minimized the number of other processes running; without these measures, it might be necessary to make several attempts before successfully spying upon the RSA private key operation or to splice together multiple observations.

Like most implementations of RSA, OpenSSL uses the Chinese Remainder Theorem (CRT) [9] when performing private key operations: It computes a 1024-bit modular exponentiation over \mathbb{Z}_{pq} using two 512-bit modular exponentiations over the rings \mathbb{Z}_p and \mathbb{Z}_q . Further, OpenSSL utilizes a “sliding window” method of modular exponentiation, decomposing $x := a^d \bmod p$ into a series of squarings $x := x^2 \bmod p$ and multiplications $x := x \cdot a^{2k+1} \bmod p$, using a set of pre-computed multipliers $\{a, a^3, a^5 \dots a^{31}\} \bmod p$.

In Figure 2 we show a small portion of one of the two modular exponentiations, as observed by the L1 Spy process. The modular squarings and modular multiplications are easily distinguishable here; this difference results from the use of the `BN_sqr` vs. `BN_mul` functions in OpenSSL: `BN_sqr` is slightly faster, but uses a different temporary working space for performing Karatsuba multiplication [8] and consequently leaves a different “footprint” behind in the cache. From the sequence of multiplications and squarings, we can typically obtain about 200 bits out of each 512-bit exponent: For each multiplication we can infer a 1 bit, since the multipliers are all odd powers of a (this yields approximately 80 bits), and any time we have more than five squarings without an intervening multiplication, we can infer the presence of one or more 0 bits, since the multipliers are of degree at most 31 (this yields approximately 120 bits).

In addition to the exponent bits obtained from the sequence of multiplications and squarings, some information can be obtained about the individual multipliers a^{2k+1} used. These multipliers are precomputed and stored in a table; when a multiplication is performed, the appropriate multiplier is read from memory (and hence loaded into the cache), thereby allowing the spy process to determine into which cache set the multiplier is mapped. Theoretically, this might allow us to obtain all of the remaining bits of the exponent, but two factors limit the information leakage. First, there is no easy way to distinguish between a cache set being used by a copy of the multiplier a^{2k+1} and a cache set being used for temporary storage as part of the process of computing the modular multiplication; consequently, when the cache set corresponding to the multiplier used is one of the cache sets already used as part of the fixed memory-access pattern of a modular multiplication,

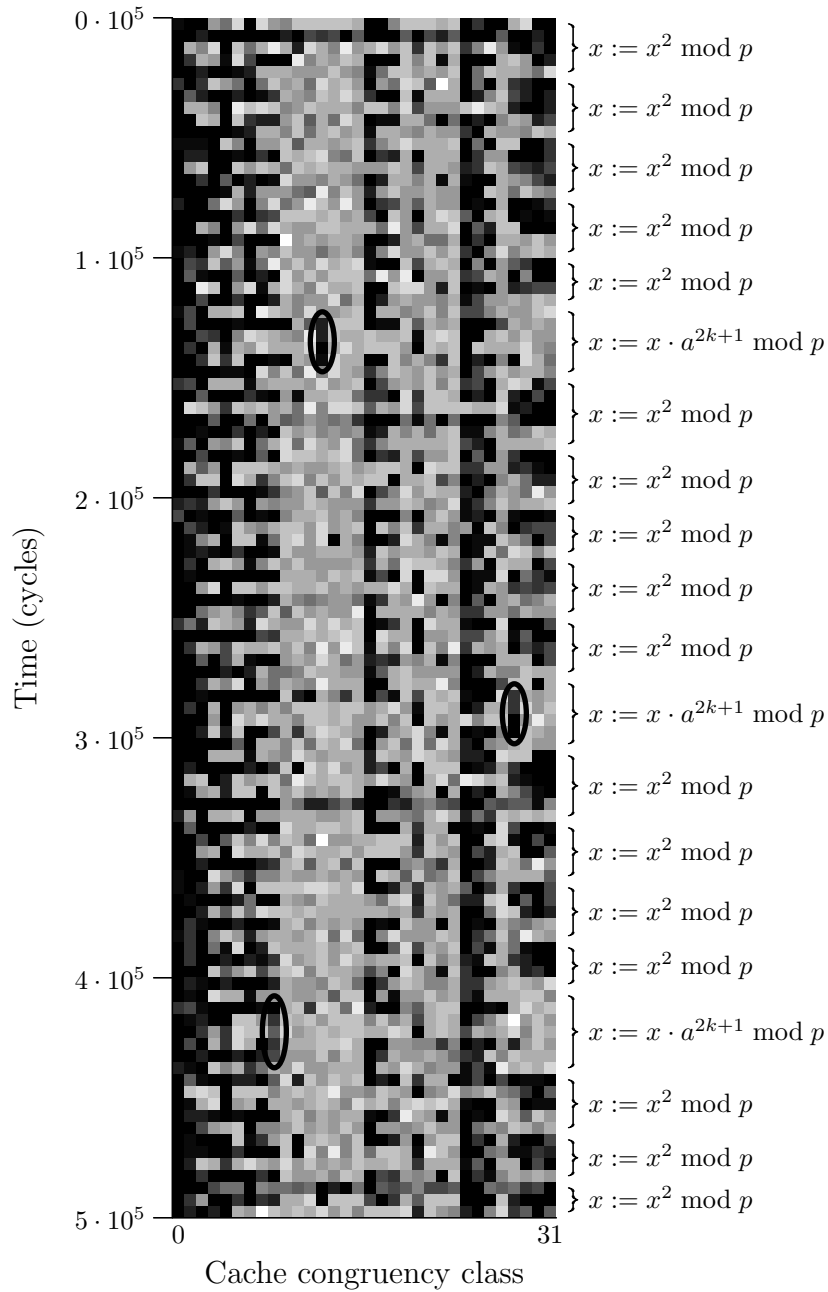


FIGURE 2. Part of a 512-bit modular exponentiation in OpenSSL 0.9.7c. The shading of each block indicates the number of cycles needed to access all the lines in a cache set, ranging from 120 cycles (white) to over 170 (black). The circled regions reveal information about the multipliers a^{2k+1} being used.

no information can be obtained about the multiplier. Second, it is possible for more than one multiplier to map to the same cache set (in our test, we found that the 16 multipliers mapped in pairs to 8 different cache sets); this reduces the information obtained about the multiplier even when the correct cache set is identified. Nevertheless, we find that roughly 110 bits from each exponent can typically be obtained in this manner.

We note again that this side channel has the unusual property of only requiring a single private key operation to be observed. Unlike earlier timing attacks which can require over 10^6 private key operations to be observed [2], the side channel which we use has high enough bandwidth and signal-to-noise ratio that repeated measurements are unnecessary.

6. OFFLINE COMPUTATIONS

Using the side channel described in the preceding section, it is possible to obtain approximately 310 bits out of each 512-bit exponent, with the bits distributed essentially randomly throughout. While it has been shown that knowing half of the bits of one of the factors is sufficient to allow N to be factored in polynomial time using lattice reduction methods [3], such methods require that all the known bits are contiguous, which we do not have in this case. Consequently, we have to construct a new method for using the data we have obtained to assist in factoring N .

Assume that we are given the public modulus N and the public exponent e , and that we have obtained some information about the private exponents d_p and d_q used in computations modulo p and q prior to the reconstruction of a value modulo N using the CRT. Then from the construction of RSA, we note that

$$ed_p \equiv 1 \pmod{p-1}$$

$$ed_q \equiv 1 \pmod{q-1}$$

and so for some $k_p, k_q \in Z_e^*$,

$$ed_p = k_p(p-1) + 1$$

$$ed_q = k_q(q-1) + 1$$

and (after some simple algebraic manipulation),

$$(1) \quad Nk_pk_q = (pk_p) \cdot (qk_q) = (ed_p + k_p - 1) \cdot (ed_q + k_q - 1)$$

Now consider sets S_n containing all ordered 4-tuples (k_p, k_q, d_p, d_q) which satisfy equation (1) modulo $e \cdot 2^n$, match our observations about specific bits of d_p and d_q , and have $0 \leq d_p, d_q < 2^n$. The set S_0 simply contains tuples $(k_p, k_q, 0, 0)$ where $Nk_pk_q \equiv (k_p - 1)(k_q - 1) \pmod{e}$, so

k_q is uniquely determined given N and k_p , and $|S_0| < e$. Further, if we have the set S_n , we can compute the set S_{n+1} by lifting each 4-tuple in S_n into two 4-tuples — the n^{th} bit of d_p can be either 0 or 1, and once it has been chosen, equation (1) will determine the n^{th} bit of d_q — and discarding any 4-tuples which do not match our observations.

If we know neither the n^{th} bit of d_p nor the n^{th} bit of d_q , then $|S_{n+1}| = 2|S_n|$. If we know one of the bits, then each 4-tuple in S_n maps to a unique 4-tuple in S_{n+1} , and we have $|S_{n+1}| = |S_n|$. If we know both bits, then the size of S_{n+1} will be approximately half the size of S_n . Consequently, the size of S_n as n increases follows a random walk starting from e ; but since there are significantly more known bits than unknown bits, the overall trend is downwards, and the sets are unlikely to ever become so large as to be unwieldy.

Once we have computed the set S_{512} , it is a simple matter to test the remaining candidate exponents and retrieve the factorization of N .

7. SOLUTIONS AND WORKAROUNDS

Both the covert channel operating through shared caches and the associated side channel can be easily blocked by processor designers. Most trivially, if shared caches (and simultaneous multithreading) are eliminated, there will be no potential for information leakage. More interestingly, the covert channel can be almost completely removed, and the side channel can be made small enough as to be cryptologically insignificant, if the cache eviction logic is changed: Rather than using a single pseudo-LRU cache eviction strategy, the cache eviction logic could be made aware of individual threads (in the case of simultaneous multithreading) or processor cores and made to only allow thread A to evict a cache line “owned” by thread B if thread B currently “owns” more than its “fair share” (e.g., one half if the cache is shared between two threads) of the cache lines in the set. As we lack expertise in the field of microarchitecture, we cannot comment upon whether such a strategy would be practical in such a performance-critical path as the L1 data cache, but it seems very likely that this or similar methods could be used on secondary caches.

These channels can also be closed by the operating system. If only one thread or processor sharing a cache is ever used — that is, if any other threads are forced to be idle — then there is effectively no cache sharing, and these channels no longer exist. This approach has been taken by FreeBSD [5] and some Linux distributions in response to the risk posed by Intel’s Hyper-Threading. A more subtle approach can also be taken via the kernel scheduler. Recognizing that a side channel

between threads is only dangerous if the threads are operating at different privileges — or, put another way, if the threads are not permitted to debug each other — a scheduler could be written in such a way as to use the credentials of threads in the process of determining which threads should be scheduled on which (virtual) processors. There are some potential dangers in this approach, however: Since the credentials of a thread can be changed during a system call, the kernel would have to re-evaluate whether a set of threads are compatible at several different points, which could lead to a loss of performance, the introduction of bugs, and quite possibly difficulties involving the locking of kernel data.

In some cases, this side channel can also be closed at the application level. If applications and libraries are written in such a manner that the code path and sequence of memory accesses are oblivious to the data and key being used, then all timing side channels are immediately and trivially closed providing that the underlying hardware does not exhibit data-dependent instruction timings. This would be a dramatic divergence from existing practice — in OpenSSL, the large integer arithmetic code alone contains over a thousand “if” statements — and would require that some existing algorithms be thrown out or reworked considerable (e.g., the “sliding window” method of modular exponentiation), which could significantly impact performance. A rather weaker approach is to require only that the sequence of *cache lines* accessed is independent of the key used; this approach has been adopted in OpenSSL’s modular exponentiation code (in response to an earlier version of this paper). We consider this weaker approach to be dangerous at best, in light of remarks from Bernstein pointing out that, even within a single cache line, different bytes may take different amounts of time to access [1]. In addition, even if all cryptographic software is rewritten to avoid information leakage, there is significant potential for leakage of sensitive non-cryptographic data, since a shared cache could allow an attacker to distinguish between `vi` and `emacs`, or even to precisely measure the timing of a sequence of keystrokes.

8. FINAL WORDS

While we have demonstrated our attack only against one very specific target, this target was selected for its high profile, not because we believed it to be in any way more vulnerable to attack. Taken in combination with the recent results of Osvik, Shamir, and Tromer [13], our work clearly demonstrates that shared caches introduce very significant dangers which should not be ignored. Sadly, in the six months

since this work was first quietly circulated within the operating system security community, and the four months since it was first publicly disclosed, some vendors have failed to provide any response. The discovery of cryptologic attacks only helps to improve security if the necessary steps are taken to guard against them; we hope that the vendor community will prove more responsive in the future.

9. ACKNOWLEDGEMENTS

The author wishes to thank Jacques Vidrine and Mike O'Connor for their invaluable aid in facilitating pre-disclosure communications within the operating system security community, and the many people who made helpful comments about earlier versions of this paper.

REFERENCES

- [1] D.J. Bernstein. Cache-timing attacks on AES, 21 Nov 2004. Document ID: cd9faae9bd5308c440df50fc26a517b4.
- [2] D. Brumley and D. Boneh. Remote timing attacks are practical. In *Proceedings of the 12th USENIX Security Symposium*, August 2003.
- [3] D. Coppersmith. Finding a small root of a bivariate integer equation; factoring with high bits known. In U. Maurer, editor, *Advances in Cryptology - EUROCRYPT '96*, LNCS 1070, pages 178–189. Springer-Verlag, 1996.
- [4] FreeBSD Project. The FreeBSD operating system.
<http://www.freebsd.org/>.
- [5] FreeBSD Project. Freebsd security advisory FreeBSD-SA-05:09.htt, May 2005.
- [6] IEEE Std 1003.1. 2004 Edition.
- [7] Intel Corporation. *IA-32 Intel Architecture Optimization Reference Manual*, June 2005.
- [8] A. Karatsuba and Y. Ofman. Multiplication of multidigit numbers on automata. *Soviet Physics - Doklady*, 7:595–596, 1963.
- [9] D.E. Knuth. *Seminumerical Algorithms*, volume 2 of *The Art of Computer Programming*. Addison Wesley, third edition, 1997.
- [10] B.W. Lampson. A note on the confinement problem. *Communications of the ACM*, 16(10):613–615, 1973.
- [11] D. Marr, F. Binns, D. Hill, G. Hinton, D. Koufaty, J. Miller, and M. Upton. Hyper-threading technology architecture and microarchitecture: A hypertext history. *Intel Technology Journal*, February 2002.
<http://developer.intel.com/technology/itj/2002/volume06issue01/>.
- [12] National Institute of Standards and Technology. Announcing the Advanced Encryption Standard (AES). NIST FIPS PUB 197, U.S. Department of Commerce, 2001.
- [13] D.A. Osvik, A. Shamir, and E. Tromer. Cache attacks and countermeasures: the case of AES. In *RSA Conference Cryptographers Track (CT-RSA) 2006*, 2005. to appear.
- [14] D. Page. Theoretical use of cache memory as a cryptanalytic side-channel. Technical Report CSTR-02-003, Department of Computer Science, University of Bristol, June 2002.

- [15] F. Sebek. Cache memories in real-time systems. Technical Report 01 / 37, Mälardalen Real-Time Research Centre, 2001.
- [16] The OpenSSL Project. OpenSSL: The open source toolkit for SSL/TLS. <http://www.openssl.org/>.
- [17] Y. Tsunoo, T. Saito, T. Suzaki, M. Shigeri, and H. Miyauchi. Cryptanalysis of DES implemented on computers with cache. In C.D. Walter, Ç.K. Koç, and C. Paar, editors, *Cryptographic Hardware and Embedded Systems - CHES 2003*, LNCS 2779, pages 62–76. Springer-Verlag, 2003.
- [18] D.M. Tullsen, S. Eggers, and H.M. Levy. Simultaneous multithreading: Maximizing on-chip parallelism. In *Proceedings of the 22nd Annual International Symposium on Computer Architecture*, pages 392–403, 1995.

IRMACS CENTRE, SIMON FRASER UNIVERSITY, BURNABY, BC, CANADA
E-mail address: cperciva@irmacs.sfu.ca