

# Cache-timing attacks on AES

Daniel J. Bernstein \*

Department of Mathematics, Statistics, and Computer Science (M/C 249)  
The University of Illinois at Chicago  
Chicago, IL 60607–7045  
djb@cr.y.p.to

**Abstract.** This paper warns against the use of S-boxes in cryptography. In particular, this paper shows that a simple cache-timing attack against AES software reveals some key bits; this paper also discusses some of the obstacles to constant-time array access on modern CPUs.

**Keywords:** side channels, timing attacks, software timing attacks, cache timing, load timing, S-boxes, AES, Salsa10

## 1 Introduction

The point of this paper is that it is difficult for software to perform an array lookup in time independent of the array index. Consequently, typical high-speed AES software quickly reveals an unacceptable number of key bits to the simplest conceivable timing attack; see Section 3. Using secret data as an array index is a recipe for disaster.

One potential response is to load all necessary arrays from DRAM into L1 cache before using them. Unfortunately, this (1) is difficult to do correctly and (2) does *not* guarantee constant-time array access. See Section 4.

Another potential response is to try to randomize the information being leaked. This is a common response to hardware side channels; see, e.g., [4]. Unfortunately, randomization is expensive.

A third potential response, and the solution I recommend, is to stop using S-boxes: when a cryptographic design relies on S-boxes to achieve reasonable software speed (as AES does), throw it away. See Section 5.

---

\* The author was supported by the National Science Foundation under grant CCR–9983950, and by the Alfred P. Sloan Foundation. Date of this document: 2004.11.11. Permanent ID of this document: [cd9faae9bd5308c440df50fc26a517b4](https://arxiv.org/abs/cd9faae9bd5308c440df50fc26a517b4). This is a preliminary version meant to announce ideas; it will be replaced by a final version meant to record the ideas for posterity. There may be big changes before the final version. Future readers should not be forced to look at preliminary versions, unless they want to check historical credits; if you cite a preliminary version, please repeat all ideas that you are using from it, so that the reader can skip it.

## Previous work

Koeune and Quisquater in [10] presented details of a fast timing attack on a “careless implementation” of AES in which doublings used branches rather than table lookups. “The result presented here is not an attack against Rijndael, but against bad implementations of it,” Koeune and Quisquater wrote in [10, Section 7]. In contrast, variable-time S-box lookups are not a sign of a “careless implementation” of AES; they are difficult to avoid.

Kelsey, Schneier, Wagner, and Hall in [9, Section 7] expressed the belief that timing attacks could exploit “cache hit ratio in large S-box ciphers like Blowfish, CAST, and Khufu.” Page in [12], and then Tsunoo, Saito, Suzuki, Shigeri, and Miyauchi in [14], presented fast timing attacks on DES when all S-boxes were out of cache before each DES call.

Tsunoo et al. mentioned a strategy for attacking out-of-cache AES if many (unusual) “plaintexts with long encryption time” could be found. See [14, Section 4.1]. This paper demonstrates a much simpler timing attack on AES. This paper also points out that the “guaranteed” countermeasure suggested in [12, Section 5.1] and [14, Section 5] fails to achieve constant timings on real computers.

## 2 Summary of AES

AES scrambles a 16-byte input  $n$  using a 16-byte key  $k$ , a constant 256-byte “S-box”  $S = (99, 124, 119, 123, 242, \dots)$ , and a constant 256-byte “doubled S-box”  $S' = (198, 248, 238, 246, 255, \dots)$ . These S-boxes are expanded into tables  $T_0, T_1, T_2, T_3$  defined by

$$\begin{aligned}T_0[b] &= (S[b] \oplus S'[b], S[b], S[b], S'[b]), \\T_1[b] &= (S[b], S[b], S'[b], S[b] \oplus S'[b]), \\T_2[b] &= (S[b], S'[b], S[b] \oplus S'[b], S[b]), \\T_3[b] &= (S'[b], S[b] \oplus S'[b], S[b], S[b]).\end{aligned}$$

Here  $\oplus$  means xor, i.e., addition of vectors modulo 2.

AES works with two 16-byte auxiliary arrays,  $x$  and  $y$ . The first array is initialized to  $k$ . The second array is initialized to  $n \oplus k$ .

AES first modifies  $x$  as follows. View  $x$  as four 4-byte arrays  $x_0, x_1, x_2, x_3$ . Compute the 4-byte array  $e = (S[x_3[1]] \oplus 1, S[x_3[2]], S[x_3[3]], S[x_3[0]])$ . Replace  $(x_0, x_1, x_2, x_3)$  with  $(e \oplus x_0, e \oplus x_0 \oplus x_1, e \oplus x_0 \oplus x_1 \oplus x_2, e \oplus x_0 \oplus x_1 \oplus x_2 \oplus x_3)$ .

AES then modifies  $y$  as follows. View  $y$  as four 4-byte arrays  $y_0, y_1, y_2, y_3$ . Replace  $(y_0, y_1, y_2, y_3)$  with

$$\begin{aligned}(T_0[y_0[0]] \oplus T_1[y_1[1]] \oplus T_2[y_2[2]] \oplus T_3[y_3[3]] \oplus x_0, \\T_0[y_1[0]] \oplus T_1[y_2[1]] \oplus T_2[y_3[2]] \oplus T_3[y_0[3]] \oplus x_1, \\T_0[y_2[0]] \oplus T_1[y_3[1]] \oplus T_2[y_0[2]] \oplus T_3[y_1[3]] \oplus x_2, \\T_0[y_3[0]] \oplus T_1[y_0[1]] \oplus T_2[y_1[2]] \oplus T_3[y_2[3]] \oplus x_3).\end{aligned}$$

I learned this view of AES from software published by Barreto in [3].

AES modifies  $x$  again, using  $\oplus 2$  instead of  $\oplus 1$ ; modifies  $y$  again; modifies  $x$  again, using  $\oplus 4$ ; modifies  $y$  again; and so on for a total of ten rounds. The constants for the  $x$  modifications are 1, 2, 4, 8, 16, 32, 64, 128, 27, 54. On the tenth round, the modification of  $y$  uses  $(S[], S[], S[], S[])$  rather than  $T_0[] \oplus T_1[] \oplus T_2[] \oplus T_3[]$ . The final value of  $y$  is the output  $\text{AES}_k(n)$ .

Note that the evolution of  $x$  is independent of  $n$ . One can expand  $k$  into the eleven values of  $x$ , together occupying 176 bytes, and then reuse those values for each  $n$ . Beware, however, that the time saved may be outweighed by the time to load the precomputed values of  $x$  from memory, when many keys are handled simultaneously.

### 3 A simple cache-timing attack on AES

I wrote reasonably straightforward AES code for the x86 architecture, and timed it on a 900MHz Athlon as follows:

```
int aescycles(void)
{
    timing tstart;
    timing tend;
    int t;
    do {
        timing_now(&tstart);
        aes(out,key,in);
        timing_now(&tend);
        t = timing_diff(&tend,&tstart);
    } while (t <= 0 || t >= 1000);
    return t;
}
```

Here `timing_now` uses the CPU's RDTSC instruction to report the number of clock cycles since boot; `timing_diff` subtracts two such numbers. Timings outside a reasonable range are ignored: occasionally the AES computation is interrupted by a network packet or by another program. (Some CPUs, notably the Pentium M, have also been reported to erroneously add  $2^{32}$  to occasional cycle counts.)

I did this repeatedly with pseudorandom inputs under a single key, identifying the byte at input position `b` that made AES slowest:

```
for (x = 0;x < 256;++x) {
    xt = 0;
    for (i = 0;i < loops;++i) {
        for (j = 0;j < 16;++j) in[j] = random() >> 16;
        in[b] = x;
        xt += aescycles(); xt += aescycles(); xt += aescycles();
        xt += aescycles(); xt += aescycles();
    }
}
```

```

    }
    if (xt > bestxt) { bestx = x; bestxt = xt; }
}
return bestx;

```

Unlike [14], I did not go to any extra effort to knock table entries out of cache.

I then tried, for increasingly many repetitions and for various input positions, seeing how well that byte predicted the corresponding key byte, for ten random keys:

```

for (loops = 4; loops <= 65536; loops *= 2) {
    for (b = 0; b < 16; ++b) {
        printf("%d, %d loops:", b, loops);
        for (k = 0; k < 10; ++k) {
            for (j = 0; j < 16; ++j) key[j] = getchar();
            printf(" %d", bestx() ^ key[b]);
            fflush(stdout);
        }
        printf("\n");
    }
}

```

Here `getchar()` was reading bytes from `/dev/urandom`, a source of hard-to-predict data. The operation `^` is xor; recall that AES begins by xor'ing the key with the input and then using each byte of the result as a table index.

The program's output showed within a few minutes that this simple timing analysis reliably predicted 24 bits of the key—specifically, key byte 1 was  $135 \oplus x$  where  $x$  was the input byte 1 that maximized the AES time, and similarly for key bytes 5 and 9:

```

0, 1024 loops: 228 220 220 90 146 12 192 217 130 137
1, 1024 loops: 135 135 135 135 135 135 135 135 135 135
2, 1024 loops: 234 39 14 151 153 53 8 70 6 155
3, 1024 loops: 218 149 169 94 204 41 145 72 202 0
4, 1024 loops: 177 38 46 15 146 81 119 59 25 187
5, 1024 loops: 135 135 135 135 135 135 135 135 135 135
6, 1024 loops: 0 175 64 9 3 67 78 201 215 185
7, 1024 loops: 204 98 59 128 253 186 21 57 209 43
8, 1024 loops: 239 56 110 255 254 255 30 189 94 63
9, 1024 loops: 135 135 135 178 135 173 135 135 135 135
10, 1024 loops: 20 184 24 83 142 64 227 182 48 230
11, 1024 loops: 165 155 170 245 149 246 97 245 172 64
12, 1024 loops: 238 226 91 138 21 140 222 118 132 174
13, 1024 loops: 85 83 115 24 67 6 142 90 151 243
14, 1024 loops: 54 40 227 43 181 228 82 110 67 170
15, 1024 loops: 133 33 254 70 115 59 217 120 140 204

```

I stopped the analysis at this point; 24 key bits is already an unacceptable level of information leakage.

To check that these results were not specific to my AES implementation, I substituted the (slower, even without key expansion) implementation from [3], and again found that this simple timing analysis reliably predicted some key bytes:

```
8, 512 loops: 253 253 253 253 253 253 10 253 253 253
12, 512 loops: 253 253 253 253 253 253 110 253 253 253
```

Of course, the same attacks can be carried out on AES implementations exposed to the network. See [5] for an introduction to remote timing attacks.

## 4 Load-timing variability

This section discusses several ways in which the time taken to load an array element depends on the array index. I do not guarantee the completeness of this list; my goal in this section is to demonstrate that constant S-box timing is difficult to achieve, not to demonstrate that it *can* be achieved.

### Cache is faster than DRAM

Here is the simplest model of timing variability; see, e.g., [12, Section 2] and [14, Section 2.1]. All recently used **lines** of memory are automatically saved in a limited-size **cache**. Reading from a cached line takes less time than reading from an uncached line.

The Athlon, for example, divides memory into 64-byte lines. There are 327680 bytes of cache containing 5120 recently used lines.

Similarly, a typical Pentium III divides memory into 32-byte lines, and has 524288 bytes of cache containing 16384 recently used lines.

### L1 cache is faster than L2 cache

Modern CPUs actually have two (and occasionally more) levels of cache. All recently used lines of **level-2 cache**, up to a certain limit, are automatically saved in (or moved to) the **level-1 cache**. Reading from a line in L1 cache takes less time than reading from a line in L2 cache.

On the Athlon, for example, there are 65536 bytes of L1 cache containing 1024 recently used lines, and 262144 bytes of L2 cache containing 4096 other recently used lines. The result of a load from L1 cache is available after 3 cycles; the result of a load from L2 cache is available after 11 cycles.

Similarly, a typical Pentium III has 16384 bytes of L1 cache containing 512 recently used lines, and 524288 bytes of L2 cache containing 16384 recently used lines, including all the lines in L1 cache.

## Cache associativity is limited

The Athlon's L1 cache is **two-way associative**. (The Pentium III's L1 cache is four-way associative.) This means that each line of memory has only two locations in L1 cache where it can be placed; all lines with the same address modulo 32768 are competing for those two locations. If three lines with the same address modulo 32768 are read, the first line is booted out of the L1 cache.

Typical AES software uses several different arrays: the AES input; the key (perhaps expanded); the AES output; the stack, for temporary variables; and the AES tables. The positions of these arrays can easily overlap modulo 32768. Consequently, even if the AES software starts by making sure that the AES tables are in L1 cache, it might lose some table entries from cache by accessing (e.g.) the key and the stack.

A solution is to *first* copy all relevant data onto the stack and *then* make sure that the AES tables are in L1 cache. This takes some time. Copying 208 bytes of input, key, and output takes at least 26 cycles; loading 64 cache lines, to guarantee that those lines are in L1 cache, takes at least 32 cycles, plus several cycles waiting in case the last few lines are still being copied from L2 cache.

(Note that similar penalties apply to many other functions. Consider, for example, the 4096-byte expanded key arrays used in [11, Section 3.2]—“GCM”—for moderate-speed multiplication in a field of size  $2^{128}$ . The array indices depend on secret data. The message and the stack can overlap the arrays modulo 32768, forcing various portions of the arrays out of L1 cache and leaking a great deal of secret information.)

Another difficulty is that the operating system will interrupt the program during the AES computation if (for example) a network packet arrives, or a disk has found a block of data, or another program has been waiting for a while to use the CPU. When the program starts again, the AES tables are unlikely to be in L1 cache. A solution is for the program to call `cache_level1(table,size)` to tell the operating system to put the tables back into L1 cache before restarting the program. Unfortunately, I am not aware of any current operating-system support for this `cache_level1()` function.

## Cache-bank throughput is limited

The time taken to load an array element depends on the array index *even if all array entries are in L1 cache*.

On the Athlon, for example, each 64-byte cache line is divided into 8 **banks**. Usually the L1 cache can perform two loads from L1 cache every cycle. However, if the two loads are from bank 0 of two lines with different addresses modulo 32768, or from bank 1 of two such lines, etc., then the second load will wait for a cycle.

(I am not aware of any official Athlon documentation saying this. Thanks to Andreas Kaiser for drawing my attention to [2, Section 5.8], which documents the analogous fact for the Athlon 64. See [7, Section 14.7] for a similar observation about the Pentium Pro/II/III.)

One can try to avoid this problem by using 64-byte spacing between elements of an array, so that all elements are in the same bank. This, however, slows down the computation in two ways. First, multiplying an array index by 64 takes an extra instruction compared to multiplying it by 1, 2, 4, or 8; there are 160 array lookups in AES (aside from key expansion), and the 160 extra instructions consume more than 53 Athlon cycles. Second, it is now necessary to load 256 lines into L1 cache rather than 64, taking at least 128 Athlon cycles rather than 32.

A more subtle approach is to use 32-byte spacing or 16-byte spacing. This reduces the penalty from  $53 + 128$  cycles to  $53 + 64$  or  $53 + 32$  cycles respectively. However, it means that each array occupies 2 or 4 banks respectively. One must carefully schedule loads so that loads from the same bank cannot bump into each other.

## 5 Cryptography without S-boxes

When I explained cache-timing attacks in a course on cryptography in 1997, I categorically recommended against all use of S-boxes in new cryptographic designs. I have seen nothing since then to change my opinion. *S-boxes are unsafe and unnecessary.*

The simple operations  $+$  (addition modulo  $2^{32}$ ),  $\ll c$  (shift of a 32-bit word  $c$  bits to the left, where  $c$  is constant),  $\gg c$  (shift of a 32-bit word  $c$  bits to the right),  $\lll c$  (rotation of a 32-bit word  $c$  bits to the left),  $\&$  (32-bit and),  $|$  (32-bit or), and  $\wedge$  (32-bit xor) take constant time on typical CPUs and appear to provide high security at high speed. Three examples of this approach are TEA, published by Wheeler and Needham in [16]; SHA-256, published in [1]; and Helix, published by Ferguson, Whiting, Schneier, Kelsey, Lucks, and Kohno in [6].

(One can implement AES using these operations, of course, but the result is painfully slow. AES cannot achieve reasonable software speed without S-boxes.)

Here is a new example of a cryptographic function with input-independent timings on typical CPUs. The function, denoted Salsa10, produces a 64-byte output (representing 16 little-endian 4-byte integers) from a 64-byte input as follows:

```
#define R(a,b) (((a) << (b)) | ((a) >> (32 - (b))))
void salsa10_specification(uint32 x[16],uint32 in[16])
{
    int i;
    for (i = 0;i < 16;++i) x[i] = in[i];
    for (i = 10;i > 0;--i) {
        x[ 4] ^= R(x[ 0]+x[12], 6);  x[ 8] ^= R(x[ 4]+x[ 0],17);
        x[12] += R(x[ 8]|x[ 4],16);  x[ 0] += R(x[12]^x[ 8], 5);
        x[ 9] += R(x[ 5]|x[ 1], 8);  x[13] += R(x[ 9]|x[ 5], 7);
        x[ 1] ^= R(x[13]+x[ 9],17);  x[ 5] += R(x[ 1]^x[13],12);
    }
}
```

```

x[14] ^= R(x[10]+x[ 6], 7); x[ 2] += R(x[14]^x[10],15);
x[ 6] ^= R(x[ 2]+x[14],13); x[10] ^= R(x[ 6]+x[ 2],15);
x[ 3] += R(x[15]|x[11],20); x[ 7] ^= R(x[ 3]+x[15],16);
x[11] += R(x[ 7]^x[ 3], 7); x[15] += R(x[11]^x[ 7], 8);
x[ 1] += R(x[ 0]|x[ 3], 8)^i;x[ 2] ^= R(x[ 1]+x[ 0],14);
x[ 3] ^= R(x[ 2]+x[ 1], 6); x[ 0] += R(x[ 3]^x[ 2],18);
x[ 6] += R(x[ 5]^x[ 4], 8); x[ 7] += R(x[ 6]^x[ 5],12);
x[ 4] += R(x[ 7]|x[ 6],13); x[ 5] ^= R(x[ 4]+x[ 7],15);
x[11] ^= R(x[10]+x[ 9],18); x[ 8] += R(x[11]^x[10],11);
x[ 9] ^= R(x[ 8]+x[11], 8); x[10] += R(x[ 9]|x[ 8], 6);
x[12] += R(x[15]^x[14],17); x[13] ^= R(x[12]+x[15],15);
x[14] += R(x[13]|x[12], 9); x[15] += R(x[14]^x[13], 7);
}
for (i = 0;i < 16;++i) x[i] += in[i];
}

```

A casual Athlon implementation of Salsa10 takes under 12 cycles per output byte—less time than AES.

I designed Salsa10 to be a cryptographically strong pseudorandom number generator. I conjecture that the following two oracles are indistinguishable:

- An oracle for a uniform random function from 16-byte strings to 64-byte strings.
- An oracle for  $n \mapsto \text{Salsa10}(n, 0, k, 0)$ . Here  $n$  is a 16-byte string and  $k$  is a uniform random 16-byte string.

I offer \$1000 for a proof that one can distinguish these oracles with better price-performance ratio than searching through  $2^{128}$  keys.

I simultaneously designed Salsa10 to be a strong hash function. I offer \$1000 for a 64-byte string  $x$  such that  $\text{Salsa10}(x)$  begins with 16 bytes of 0s, or for two different 64-byte strings  $x, x'$  such that  $\text{Salsa10}(x)$  and  $\text{Salsa10}(x')$  begin with the same 32 bytes.

The Salsa10 data flow follows a row-column structure as in AES, for fast diffusion and high parallelizability. The three basic operations in Salsa10 were mixed randomly (except that I decided to have exactly eight  $|$  operations), and the rotation amounts were chosen randomly between 5 and 20; I conjecture that random choices within this framework are overwhelmingly likely to be secure.

I am not aware of any patents or patent applications relevant to Salsa10.

## References

1. —, *Secure hash standard*, Federal Information Processing Standard 180-2, National Institute of Standards and Technology, Washington, 2002. URL: <http://csrc.nist.gov/publications/fips/>.
2. —, *Software optimization guide for AMD Athlon 64 and AMD Opteron processors*, Advanced Micro Devices, 2004. URL: [http://www.amd.com/us-en/assets/content\\_type/white\\_papers\\_and\\_tech\\_docs/25112.PDF](http://www.amd.com/us-en/assets/content_type/white_papers_and_tech_docs/25112.PDF).



3. Paulo S. L. M. Barreto, *The AES block cipher in C++* (2003). URL: <http://planeta.terra.com.br/informatica/paulobarreto/EAX++.zip>.
4. Johannes Bloemer, Jorge Guajardo Merchan, Volker Krummel, *Provably secure masking of AES* (2004). URL: <http://eprint.iacr.org/2004/101/>.
5. David Brumley, Dan Boneh, *Remote timing attacks are practical* (2003). URL: <http://crypto.stanford.edu/~dabo/papers/ssl-timing.pdf>.
6. Niels Ferguson, Doug Whiting, Bruce Schneier, John Kelsey, Stefan Lucks, Tadayoshi Kohno, *Helix: fast encryption and authentication in a single cryptographic primitive*, in [8] (2003), 330–346. URL: <http://www.macfergus.com/helix/>.
7. Agner Fog, *How to optimize for the Pentium family of microprocessors*, 2004. URL: <http://www.agner.org/assem/>.
8. Thomas Johansson (editor), *Fast software encryption: 10th international workshop, FSE 2003, Lund, Sweden, February 24–26, 2003, revised papers*, Lecture Notes in Computer Science, 2887, Springer-Verlag, Berlin, 2003. ISBN 3–540–20449–0.
9. John Kelsey, Bruce Schneier, David Wagner, Chris Hall, *Side channel cryptanalysis of product ciphers*, *Journal of Computer Security* **8** (2000), 141–158. ISSN 0926–227X.
10. François Koeune, Jean-Jacques Quisquater, *A timing attack against Rijndael* (1999). URL: <http://www.dice.ucl.ac.be/crypto/techreports.html>.
11. David A. McGrew, John Viega, *The security and performance of the Galois/counter mode of operation* (2004). URL: <http://eprint.iacr.org/2004/193/>.
12. D. Page, *Theoretical use of cache memory as a cryptanalytic side-channel* (2002). URL: <http://eprint.iacr.org/2002/169/>.
13. Bart Preneel (editor), *Fast software encryption: second international workshop, Leuven, Belgium, 14–16 December 1994, proceedings*, Lecture Notes in Computer Science, 1008, Springer-Verlag, Berlin, 1995. ISBN 3–540–60590–8.
14. Yukiyasu Tsunoo, Teruo Saito, Tomoyasu Suzaki, Maki Shigeri, Hiroshi Miyauchi, *Cryptanalysis of DES implemented on computers with cache*, in [15] (2003), 62–76.
15. Colin D. Walter, Cetin K. Koc, Christof Paar (editors), *Cryptographic hardware and embedded systems—CHES 2003*, Springer-Verlag, Berlin, 2003. ISBN 3–540–40833–9.
16. David J. Wheeler, Roger M. Needham, *TEA, a tiny encryption algorithm*, in [13] (1995), 363–366.