

securitypitfalls

A blog about cryptography and security by Hubert Kario

“Constant time” compare in Python

You may be familiar with the following piece of code to implement the constant time comparison function for strings:

```
01 def constant_time_compare(val1, val2):
02     if len(val1) != len(val2):
03         return False
04     result = 0
05     for x, y in zip(val1, val2):
06         result |= x ^ y
07     return result == 0
```

The idea behind this code is to compare all bytes of input using a flag value that will be flipped in any of the comparisons fail. Only when all the bytes were compared, is the ultimate result of the method returned. This is used to thwart attacks that use the time of processing queries to guess secret values.

Unfortunately, because of CPython specifics, this code **doesn't work** for its intended purpose.

Sensitive code should always use `hmac.compare_digest()` (https://docs.python.org/2/library/hmac.html#hmac.compare_digest) method and you should not write code that needs to be side channel secure in Python.

With the tl;dr version out, let's investigate why.

Timing side channel

Many attacks against cryptographic implementations don't actually use maths to compromise the systems. The [Bleichenbacher Million Messages attack](http://archiv.infsec.ethz.ch/education/fs08/secsem/bleichenbacher98.pdf) (<http://archiv.infsec.ethz.ch/education/fs08/secsem/bleichenbacher98.pdf>), [POODLE](https://www.openssl.org/~bodo/ssl-poodle.pdf) (<https://www.openssl.org/~bodo/ssl-poodle.pdf>) and [Lucky 13](http://www.isg.rhul.ac.uk/tls/TLStiming.pdf) (<http://www.isg.rhul.ac.uk/tls/TLStiming.pdf>) attacks use some kind of a side channel to guess the contents of encrypted messages, either the timing of responses or contents of responses (different TLS Alert description field values).

Side channel attacks don't impact only cryptographic protocols, other places where secret values need to be compared to values that are controlled by attacker, like checking password equality, API tokens and HMAC value validation need to be performed in constant time too.

Already back in 00's, differences in timing as low as 100ns could be distinguished over LAN environment. See research by Crosby et al. in [Opportunities And Limits Of Remote Timing Attacks](https://www.cs.rice.edu/~dwallach/pub/crosby-timing2009.pdf) (<https://www.cs.rice.edu/~dwallach/pub/crosby-timing2009.pdf>) and Brumley and Boneh in [Remote Timing Attacks are Practical](http://crypto.stanford.edu/~dabo/papers/ssl-timing.pdf) (<http://crypto.stanford.edu/~dabo/papers/ssl-timing.pdf>). Currently we also have to worry about cross VM or cross-process attacks where ability to distinguish between single cycles may be possible.

Measuring timing differences

Let's see what happens if we use the simple way to compare two strings in python, the `==` operator.

Benchmarking code:

```
timing-eq_cmp-perf.py
01 import perf
02
03 setup = """
04 str_a = b'secret API key'
05
06 alt_s = b'XXXXXXXXXXXXX'
07
08 str_b = str_a[:{0}] + alt_s[{0}:]
09
10 assert len(str_a) == len(str_b)
11 """
12
13 fun = """str_a == str_b"""
14
15 if __name__ == "__main__":
16     total_runs = 128
17     runs_per_process = 4
18     runner = perf.Runner(values=runs_per_process,
19                          warmups=16,
20                          processes=total_runs//runs_per_process)
21     vals = list(range(14)) # length of str_a
22     for delta in vals:
23         runner.timeit("eq_cmp delta={0:#04x}".format(delta),
24                      fun,
25                      setup=setup.format(delta))
```

Running it will simulate what timings does the attacker see when the difference from the expected value is at different positions in the attacker provided string.

```
command
01 PYTHONHASHSEED=1 python3 timing-eq_cmp-perf.py \
02 -o timing-eq_cmp-perf-1.py --fast
```

output

```

01 .....
02 eq_cmp delta=0x00: Mean +- std dev: 18.4 ns +- 0.0 ns
03 .....
04 eq_cmp delta=0x01: Mean +- std dev: 20.8 ns +- 0.0 ns
05 .....
06 eq_cmp delta=0x02: Mean +- std dev: 20.8 ns +- 0.0 ns
07 .....
08 eq_cmp delta=0x03: Mean +- std dev: 20.8 ns +- 0.0 ns
09 .....
10 eq_cmp delta=0x04: Mean +- std dev: 20.8 ns +- 0.0 ns
11 .....
12 eq_cmp delta=0x05: Mean +- std dev: 20.8 ns +- 0.0 ns
13 .....
14 eq_cmp delta=0x06: Mean +- std dev: 20.8 ns +- 0.0 ns
15 .....
16 eq_cmp delta=0x07: Mean +- std dev: 20.8 ns +- 0.0 ns
17 .....
18 eq_cmp delta=0x08: Mean +- std dev: 21.3 ns +- 0.0 ns
19 .....
20 eq_cmp delta=0x09: Mean +- std dev: 21.3 ns +- 0.0 ns
21 .....
22 eq_cmp delta=0x0a: Mean +- std dev: 21.3 ns +- 0.0 ns
23 .....
24 eq_cmp delta=0x0b: Mean +- std dev: 21.3 ns +- 0.0 ns
25 .....
26 eq_cmp delta=0x0c: Mean +- std dev: 21.3 ns +- 0.0 ns
27 .....
28 eq_cmp delta=0x0d: Mean +- std dev: 21.3 ns +- 0.0 ns

```

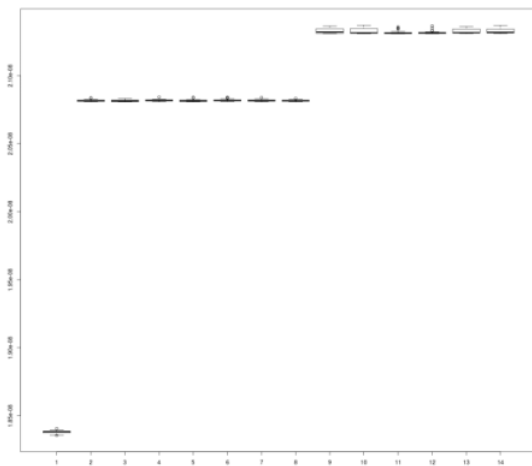
Already we can see that the difference in timing when the first different byte is at fist position or the second position is quite huge, looking at a [box plot](https://en.wikipedia.org/wiki/Box_plot) (https://en.wikipedia.org/wiki/Box_plot) of the specific values makes it quite obvious:

R script

```

01 a = read.csv(file="timing-eq_cmp-perf-1.csv", header=FALSE)
02 data = as.matrix(a)
03 boxplot(t(data))

```



https://securitypitfalls.files.wordpress.com/2018/07/timing-eq_cmp-perf-1-boxplot1.png

Let's see how does it compare to the "constant_time" _compare. First the code:

```

timing-ct_eq_cmp-perf.py
01 import perf
02
03 setup = """
04 def constant_time_compare(val1, val2):
05     if len(val1) != len(val2):
06         return False
07     result = 0
08     for x, y in zip(val1, val2):
09         result |= x ^ y
10     return result == 0
11
12 str_a = b'secret API key'
13
14 alt_s = b'XXXXXXXXXXXXXX'
15
16 str_b = str_a[:{0}] + alt_s[{0}:]
17
18 assert len(str_a) == len(str_b)
19 """
20
21 fun = """constant_time_compare(str_a, str_b)"""
22
23 if __name__ == "__main__":
24     total_runs = 128
25     runs_per_process = 4
26     runner = perf.Runner(values=runs_per_process,
27                           warmups=16,
28                           processes=total_runs//runs_per_process)
29     vals = list(range(14)) # length of str_a
30     for delta in vals:
31         runner.timeit("ct_eq_cmp delta={0:#04x}".format(delta),
32                       fun,
33                       setup=setup.format(delta))

```

The test run:

command

```
01 PYTHONHASHSEED=1 python3 timing-ct_eq_cmp-perf.py \  
02 -o timing-ct_eq_cmp-perf-1.json --fast
```

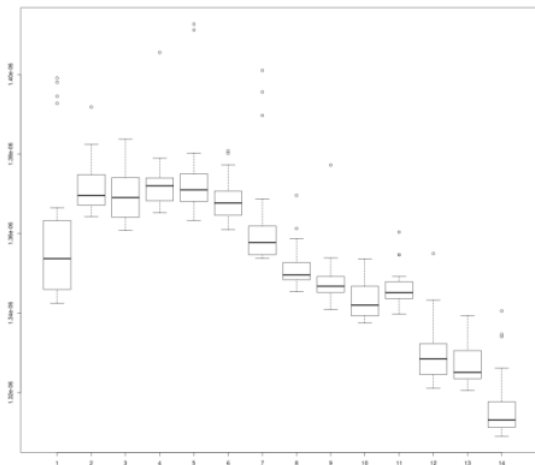
output

```
01 .....  
02 ct_eq_cmp delta=0x00: Mean +- std dev: 1.36 us +- 0.02 us  
03 .....  
04 ct_eq_cmp delta=0x01: Mean +- std dev: 1.37 us +- 0.01 us  
05 .....  
06 ct_eq_cmp delta=0x02: Mean +- std dev: 1.37 us +- 0.01 us  
07 .....  
08 ct_eq_cmp delta=0x03: Mean +- std dev: 1.37 us +- 0.01 us  
09 .....  
10 ct_eq_cmp delta=0x04: Mean +- std dev: 1.37 us +- 0.01 us  
11 .....  
12 ct_eq_cmp delta=0x05: Mean +- std dev: 1.37 us +- 0.00 us  
13 .....  
14 ct_eq_cmp delta=0x06: Mean +- std dev: 1.36 us +- 0.01 us  
15 .....  
16 ct_eq_cmp delta=0x07: Mean +- std dev: 1.35 us +- 0.01 us  
17 .....  
18 ct_eq_cmp delta=0x08: Mean +- std dev: 1.35 us +- 0.01 us  
19 .....  
20 ct_eq_cmp delta=0x09: Mean +- std dev: 1.34 us +- 0.00 us  
21 .....  
22 ct_eq_cmp delta=0x0a: Mean +- std dev: 1.35 us +- 0.00 us  
23 .....  
24 ct_eq_cmp delta=0x0b: Mean +- std dev: 1.33 us +- 0.01 us  
25 .....  
26 ct_eq_cmp delta=0x0c: Mean +- std dev: 1.33 us +- 0.01 us  
27 .....  
28 ct_eq_cmp delta=0x0d: Mean +- std dev: 1.32 us +- 0.01 us
```

The results don't look too bad, but there's definitely a difference between the first and last one, even accounting for one standard deviation between them. Let's see the box plot:

R script

```
01 a = read.csv(file="timing-ct_eq_cmp-perf-1.csv", header=FALSE)  
02 data = as.matrix(a)  
03 boxplot(t(data))
```



(https://securitypitfalls.files.wordpress.com/2018/07/timing-ct_eq_cmp-perf-1.png)

That doesn't look good. Indeed, if we compare the distributions for the different delta values using the [Kolmogorov–Smirnov test](https://en.wikipedia.org/wiki/Kolmogorov%E2%80%93Smirnov_test) (https://en.wikipedia.org/wiki/Kolmogorov%E2%80%93Smirnov_test), we'll see that results for all deltas are statistically different:

R code

```
01 a = read.csv(file="timing-ct_eq_cmp-perf-1.csv", header=FALSE)  
02 data = as.matrix(a)  
03 r = c()  
04 for (i in c(1:length(data[,1]))) {  
05   r[i] = ks.test(data[1,i], data[i,i])$p.value  
06   which(unlist(r) < 0.05/(length(data[,1])-1))
```

"output"

```
01 [1] 2 3 4 5 6 7 9 10 11 12 13 14
```

Which means that the distributions are statistically distinguishable. (The 0.05 p-value is divided by the amount of performed tests because we're applying the [Bonferroni correction](https://en.wikipedia.org/wiki/Bonferroni_correction) (https://en.wikipedia.org/wiki/Bonferroni_correction))

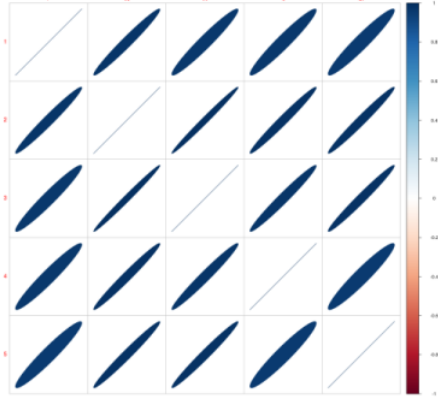
To make sure, we re-run the test 4 more times and check for correlation between medians (as the distributions are unimodal, median is [robust statistic](https://en.wikipedia.org/wiki/Robust_statistics) (https://en.wikipedia.org/wiki/Robust_statistics)).

R code

```

01 require(corrplot)
02
03 a = read.csv(file="timing-ct_eq_cmp-perf-1.csv", header=FALSE)
04 data = as.matrix(a)
05 vals = cbind(apply(data, 1, median))
06
07 for (i in 2:5) {
08   name = paste("timing-ct_eq_cmp-perf-", i, ".csv", sep="")
09   a = read.csv(file=name, header=FALSE)
10   data = as.matrix(a)
11   vals = cbind(vals, apply(data, 1, median))
12 }
13 corrplot(cor(vals, method="spearman"), method="ellipse")

```



https://securitypitfalls.files.wordpress.com/2018/07/timing-ct_eq_cmp-perf-corrplot.png

There is a very strong correlation between all the different runs, so indeed, it does look like the function is leaking timing information.

Note, we're using the "spearman" correlation statistic (https://en.wikipedia.org/wiki/Spearman%27s_rank_correlation_coefficient) as the values are not normally distributed.

Let's compare it to the `hmac.compare_digest()` method:

```

01 import perf
02
03 setup = """
04 from hmac import compare_digest
05
06 str_a = b'secret API key'
07
08 str_b = b''.join((str_a[:{0}], b'X' * (14 - {0})))
09
10 assert len(str_a) == len(str_b)
11 assert len(str_a) == 14
12 """
13
14 fun = """compare_digest(str_a, str_b)"""
15
16 if __name__ == "__main__":
17     total_runs = 128
18     runs_per_process = 4
19     runner = perf.Runner(values=runs_per_process,
20                          warmups=64,
21                          processes=total_runs//runs_per_process)
22     vals = list(range(14)) # length of str_a
23     for delta in vals:
24         runner.timeit("compare_digest delta={0:#04x}".format(delta),
25                      fun,
26                      setup=setup.format(delta))

```

command

```

01 PYTHONHASHSEED=1 python3 timing-compare_digest-perf.py \
02 -o timing-compare_digest-perf-1.json --rigorous

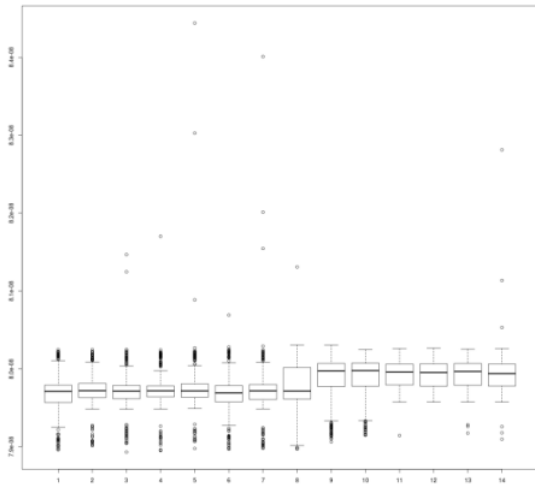
```

R code

```

01 a = read.csv(file="timing-compare_digest-perf-1.csv", header=FALSE)
02 data = as.matrix(a)
03 boxplot(t(data))

```



https://securitypitfalls.files.wordpress.com/2018/07/timing-compare_digest-perf-1.png

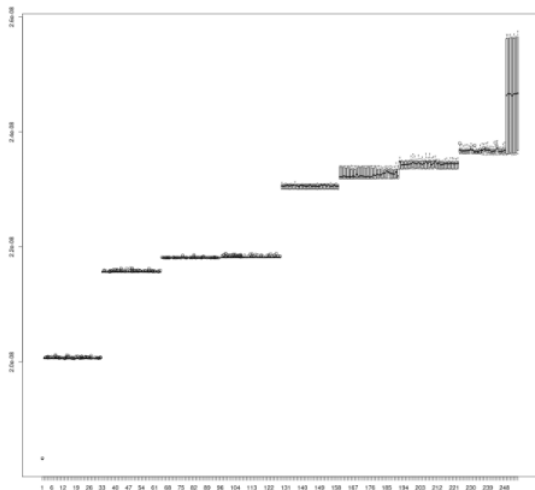
While there is some difference that depends on where the first differing byte is, there is no difference between first and second byte, and the “step” around 8th byte is only around it (when comparing longer strings, I still see just one step at the beginning and one at the end). I have no good explanation for it. That being said, the difference between medians of the 2nd byte and 11th byte is 0.240 ns, for comparison, one cycle of the CPU (4Ghz) on which the test is running takes 0.250 ns. So I’m assuming that it is not detectable over the network, but may be detectable in cross-VM attacks.

To confirm the results I’ve run the test with simple == for 255 byte long strings and with using the `hmac.compare_digest()`.

Results for ==:

R code

```
01 a = read.csv(file="timing-eq_cmp-2-perf-1.csv", header=FALSE)
02 data = as.matrix(a)
03 boxplot(t(data))
```



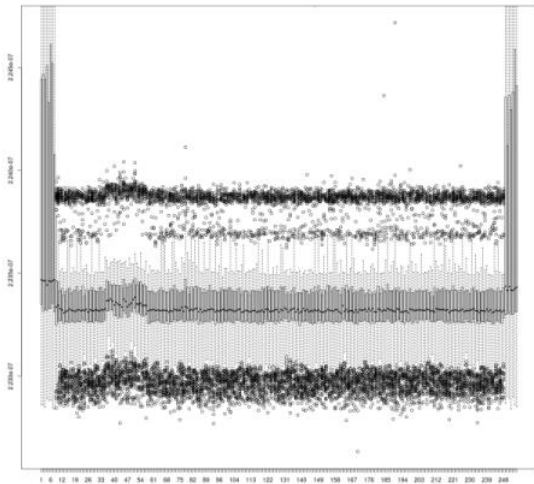
https://securitypitfalls.files.wordpress.com/2018/07/timing-eq_cmp-2-perf-1.png

As expected, obvious steps that are directly dependant on the amount of matching data between the two parameters to the operator.

Results for `compare_digest()`:

R code

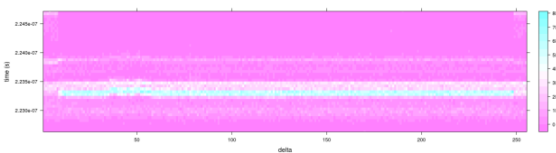
```
1 a = read.csv(file="timing-compare_digest-8-perf-1.csv", header=FALSE)
2 data = as.matrix(a)
3 boxplot(t(data), ylim=c(min(data), quantile(data, 0.99)))
```



(https://securitypitfalls.files.wordpress.com/2018/07/timing-compare_digest-8-perf-1.png).

They are quite noisy, but what the grouping around $2.239e-7$ hints at (the thick horizontal line comprised of circles), is that the distribution is not unimodal (otherwise the outliers would look like the ones below the boxes). Let's see what are the counts for different time bins, as in a histogram, in detail:

```
R code
1 require("lattice")
2 a = read.csv(file="timing-compare_digest-8-perf-1.csv", header=FALSE)
3 data = as.matrix(a)
4 h <- hist(data, breaks=200, plot=FALSE)
5 breaks = c(h$breaks)
6 mids = c(h$mids)
7 hm <- rbind(hist(data[,1], breaks=breaks, plot=FALSE)$counts)
8 for (i in c(2:length(data[,1]))) {
9   hm <- rbind(hm, hist(data[,i], breaks=breaks, plot=FALSE)$counts)}
10
11 d = data.frame(x=rep(seq(1, nrow(hm), length=nrow(hm)), ncol(hm)),
12               y=rep(mids, each=nrow(hm)),
13               z=c(hm))
14 levelplot(z~x*y, data=d, xlab="delta", ylab="time (s)",
15           ylim=c(min(data), quantile(data, 0.99)))
```



(https://securitypitfalls.files.wordpress.com/2018/07/timing-compare_digest-8-perf-1-levelplot.png).

We can see now, that even though the measurements with delta between 0 and 8 and 249 and 255 look very different on the box plot, it's more because a third mode was added to them rather than one of the other two was removed. Statistical test confirms this:

```
R code
1 a = read.csv(file="timing-compare_digest-8-perf-1.csv", header=FALSE)
2 data = as.matrix(a)
3 r = c()
4 for (i in c(1:length(data[,1]))){
5   r[i] = ks.test(data[,19], data[,i])$p.value}
6 which(unlist(r) < 0.05/nrow(data))

"result"
1 [1] 1 2 3 4 5 6 7 8 36 37 38 39 41 45 46 49 50 51 52
2 [20] 53 54 249 250 251 252 253 254 255
```

(the deltas between 36 and 54 are a fluke that subsequent quick runs didn't show).

Note about benchmarking

You may have noticed that the data we have collected, has very low amounts of noise. While it is partially the result of use of the `perf` (<https://py-pi.org/project/perf/>), module instead of the `timeit` (<https://docs.python.org/3/library/timeit.html>), library module, it mostly is the result of careful system configuration.

On the benchmarking system, the following tasks were performed:

- 3rd and 4th core were isolated
- kernel RCU was disabled on the isolated cores
- HyperThreading was disabled in BIOS
- Intel TurboBoost was disabled
- Intel power management was disabled (no C-states or P-states other than C0 were allowed)
- CPU frequency was locked in place to 4Ghz (the nominal for the i7 4970K of the workstation used)
- Decreasing maximum perf probe query rate to 1 per second
- Disabling irqbalance and setting default IRQ affinity to un-isolated cores
- ASRL disabled
- Python hash table seed fixed

Those operations can be performed by:

1. Adding `isolcpus=2,3 rcu_nocbs=2,3 processor.max_cstate=1 idle=poll` to the kernel command line
2. Disabling HyperThreading in BIOS
3. Running `python3 -m perf system tune`
4. Disabling ASLR by running `echo 0 > /proc/sys/kernel/randomize_va_space`
5. exporting the PYTHONSEED environment variable

[Documentation \(https://perf.readthedocs.io/en/latest/system.html\)](https://perf.readthedocs.io/en/latest/system.html), of the perf module provides most of the explanations of the particular options, but we diverge in two places: ASLR and Python hash seed. The purpose of the perf module is to test the overall performance of a piece of Python code (and compare it to either compilation or different implementation). Because Python is a language that answers the question "what if everything was a hash table" ;), that means the names of variables, memory positions of variables or [code \(https://vstinner.github.io/analysis-python-performance-issue.html\)](https://vstinner.github.io/analysis-python-performance-issue.html), number of variables, and particular hash table key have significant impact on performance. But, because we are interested if an attacker is able to tell behaviour of code between two different inputs, and those two inputs will likely be processed by the same process, both the ASLR seed and the Python hash table seeds will be constant from the point of view of the attacker. To speed up finding the expected value for particular inputs I thus opted out of those randomisation mechanisms.

Expectations of behaviour

You may wonder, why is the Python code so unstable, so data dependant, if the implementation of `hmac.compare_digest()` is [doing exactly the same thing \(https://github.com/python/cpython/blob/80b762f010097ab8137782e5fbd89c5c620ed4e/Modules/_operator.c#L727-L762\)](https://github.com/python/cpython/blob/80b762f010097ab8137782e5fbd89c5c620ed4e/Modules/_operator.c#L727-L762) (xor-ing the values together and then or-ing result with a guard variable)? The problem stems from the fact that the Python `int` and C `unsigned char` are vastly different data types – one is used for arbitrary precision arithmetic while the other can store just 256 unique values. Thus, even such simple operations like xor or or with two small integers are data dependant in Python.

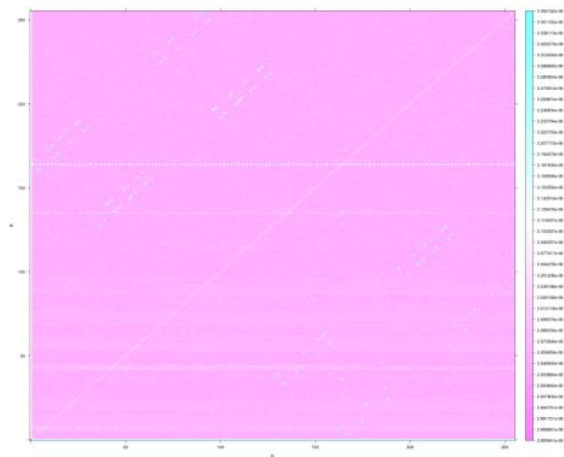
Let's see how much time does the Python VM need for those two small integers. (Unfortunately, it looks like perf uses the slow json module, and because it exports results after every loop iteration, after few hundred results, the export takes more time than benchmarking. To make it fast enough, and not waste few days on exporting the same data over and over again, we will use `timeit` module.)

Script:

```
timing-xor-2-timeit.py
1 import timeit
2 import sys
3 import math
4
5 setup = """
6 val_a = {0}
7
8 val_b = {1}
9 """
10
11 fun = """val_a ^ val_b"""
12
13 def std_dev(vals):
14     avg = sum(vals)/len(vals)
15     sum_sq = sum((i - avg)**2 for i in vals)
16     return math.sqrt(sum_sq / (len(vals) - 1))
17
18 if __name__ == "__main__":
19     total_runs = 20
20     runs_per_process = 3
21     warmups = 16
22
23     runner = timeit.Timer(fun, setup=setup.format(0, 0))
24     number, delay = runner.autorange()
25     number //= 2
26     delay /= 2
27
28     print(("will do {0} iterations per process, "
29           "expecting {1:7.2} s per process")
30           .format(number, delay), file=sys.stderr)
31     print("warmups:", file=sys.stderr, end='')
32     sys.stderr.flush()
33     for _ in range(warmups):
34         timeit.repeat(fun, setup=setup.format(0, 0), repeat=1,
35                       number=number)
36         print(".", file=sys.stderr, end='')
37         sys.stderr.flush()
38     print(file=sys.stderr)
39
40     for a in range(256):
41         for b in range(256):
42             res = []
43             for _ in range(total_runs // runs_per_process):
44                 # drop the first result as a local warmup
45                 res.extend(i / number for i in
46                           timeit.repeat(fun,
47                                         setup=setup.format(a, b),
48                                         repeat=runs_per_process + 1,
49                                         number=number)[1:])
50             print(".", file=sys.stderr, end='')
51             sys.stderr.flush()
52             if std_dev(res) > 0.001:
53                 with open('timing-xor-2-timeit-1.csv', 'a') as f:
54                     f.write(f'{a},{b},{std_dev(res)}\n')
```

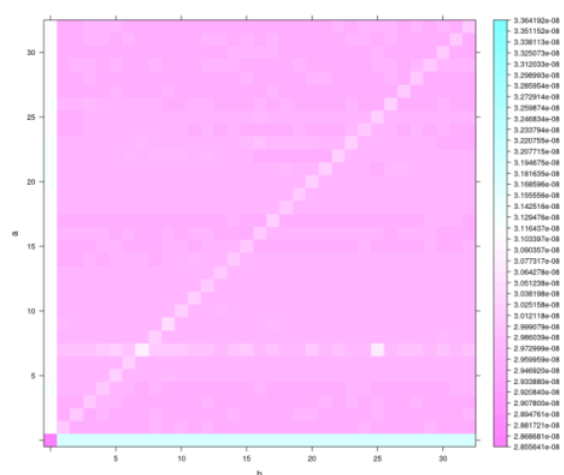
R code

```
01 require("lattice")
02 a = read.csv(file="timing-xor-2-timeit-1.csv",
03             header=FALSE, col.names=seq(1, 20),
04             fill=TRUE)
05 data = as.matrix(a)
06 med = apply(data, 1, median, na.rm=TRUE)
07 # full lines
08 len = length(med)
09 columns = ceiling(length(med) / 256)
10 d = data.frame(x=rep(seq(0, 255), length.out=len, 256),
11               y=rep(seq(0, 255), length.out=len, each=256),
12               z=med)
13 my.at = seq(min(med), max(med), length=40)
14 levelplot(z~x*y, data=d, xlab="b", ylab="a",
15           at=my.at, aspects="iso",
16           colorkey=list(at=my.at, labels=list(at=my.at)))
```



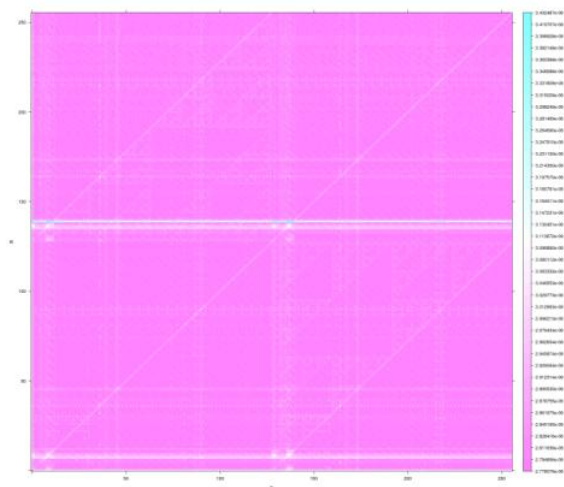
(<https://securitypitfalls.files.wordpress.com/2018/07/timing-xor-2-timeit-1.png>),

While there are few repeating patterns, there are 4 things that are of particular importance – behaviour when the two numbers are equal (the lighter diagonal), when both are zero, or when one of the operands is zero. The difference between the background and the diagonal is small, just 0.555 ns, but that translates to about 2 cycles at 4GHz. The difference between the 0, 0 and the backgrounds is even smaller, just 0.114 ns, so half a cycle. The difference between the background and the situations when the second variable is non-zero is about 2.24 ns which translates to about 9 cycles. When the first variable is non-zero and the second is, the difference is about 1.39 ns which is about 6 cycles. Here's the zoomed-in part of the graph for the small numbers:



(<https://securitypitfalls.files.wordpress.com/2018/07/timing-xor-2-timeit-1-zoom.png>),

The binary or operator is similarly dependant on values of parameters:



(<https://securitypitfalls.files.wordpress.com/2018/08/timing-or-timeit-1.png>)

Both of those things put together mean that using the supposedly constant time compare doesn't actually protect against timing attacks, but rather makes them *easier*. The strength of the signal for different inputs is about 100 times stronger, likely allowing them even over Internet, not only over LAN (as is the case for == operator).

Anything else?

Because I started looking into those microbenchmarks to verify the “constant” time CBC MAC and pad check from `tlsite-ng` (<https://github.com/tomato42/tlsite-ng/blob/02852484e5ca9ea5e0c69d525ff45a3a999b386b/tlsite/utls/constanttime.py#L96>), needed to protect against Lucky 13 (see the [very extensive article](https://www.imperialviolet.org/2013/02/04/luckythirteen.html) (<https://www.imperialviolet.org/2013/02/04/luckythirteen.html>)) by Adam Langley on the topic, I've also checked if it is possible to speed up the process of hashing data. Because on the Python level we don't have the luxury of access to lower level hash APIs, as the developers of OpenSSL have, to implement the CBC check, I wrote code that in fact calculates [256 different hmacs](https://github.com/tomato42/tlsite-ng/blob/02852484e5ca9ea5e0c69d525ff45a3a999b386b/tlsite/utls/constanttime.py#L178-L186) (<https://github.com/tomato42/tlsite-ng/blob/02852484e5ca9ea5e0c69d525ff45a3a999b386b/tlsite/utls/constanttime.py#L178-L186>) for every record that contains at least 256 bytes of data + padding. That means that for every record processed, the

client and server actually process 64 KiB of additional data. In theory (that is, if the hmac itself is constant time), we could speed the process of checking the mac and de-padding in TLS dramatically, if we could hash the data just once, as OpenSSL is doing in its TLS implementation. You may say, “but hashes are implemented in C, surely they are constant time!”. To which I’ll answer, “what did we say about trusting assumptions?”.

Let’s see how our assumptions hold. First code that hashes all of provided data, but returns also a hash from the “middle” of data (in a TLS implementation that would be the real HMAC that we need to compare to the one from record):

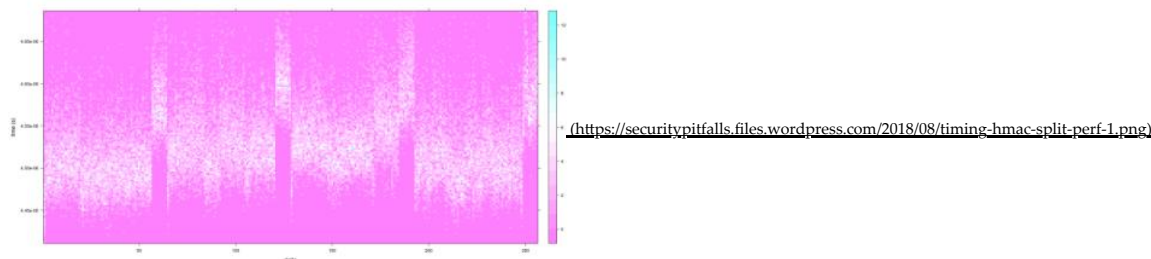
```
timing-hmac-split-perf.py
01 import perf
02
03 setup = ""
04 import hmac
05 from hashlib import sha1
06
07 def fun(digest, data, split):
08     digest.update(data[:split])
09     ret = digest.copy().digest()
10     digest.update(data[split:])
11     return ret, digest.digest()
12
13 str_a = memoryview(b'X'*256)
14 key = b'a' * 32
15
16 val_b = {}
17
18 mac = hmac.new(key, digestmod=sha1)
19 """
20
21 fun = """fun(mac.copy(), str_a, val_b)"""
22
23 if __name__ == "__main__":
24     total_runs = 128
25     runs_per_process = 4
26     runner = perf.Runner(values=runs_per_process,
27                          warmups=16,
28                          processes=total_runs//runs_per_process)
29     vals = list(range(256)) # length of str_a
30     for delta in vals:
31         runner.timeit("hmac split delta={0:#04x}".format(delta),
32                      fun,
33                      setup=setup.format(delta))
```

Command to gather the statistics:

```
command
01 PYTHONHASHSEED=1 python3 timing-hmac-split-perf.py \
02 -o timing-hmac-split-perf-1.json
```

And a way to visualise them:

```
R code
01 require("lattice")
02 a = read.csv(file="timing-hmac-split-perf-1.csv", header=FALSE)
03 data = as.matrix(a)
04 h <- hist(data, breaks=200, plot=FALSE)
05 breaks = c(h$breaks)
06 mids = c(h$mids)
07 hm <- rbind(hist(data[,1], breaks=breaks, plot=FALSE)$counts)
08 for (i in c(2:length(data[,1]))) {
09     hm <- rbind(hm, hist(data[,i], breaks=breaks, plot=FALSE)$counts)}
10
11 d = data.frame(x=rep(seq(1, nrow(hm), length=nrow(hm)), ncol(hm)),
12               y=rep(mids, each=nrow(hm)),
13               z=c(hm))
14 levelplot(z~x*y, data=d, xlab="delta", ylab="time (s)",
15           ylim=c(min(data), quantile(data, 0.99)))
```



Besides the obvious peaks since 56th to 64th byte every 64 bytes (caused by an additional hash block that had to be padded (<https://tools.ietf.org/html/rfc4634#section-4>) to calculate the intermediate HMAC), there is also a dip for the first byte of the 64 byte block and a second dip for bytes between 20 and 55 of every block. Finally, when the split is about even (in that the intermediate hash is calculated over the first 120 bytes), the whole operation takes measurably longer. In short, if the position of the intermediate hash comes from the last byte of encrypted data (as it does in TLS), calculating HMAC like this has a definite sidechannel leak.

To confirm, let’s perform Kolomogorov-Smirnov test:

```
R code
01 a = read.csv(file="timing-hmac-split-perf-1.csv", header=FALSE)
02 data = as.matrix(a)
03 r=c()
04 for (i in c(1:nrow(data))) {
05     r[i] = ks.test(data[2,], data[i,])$p.value}
06 which(unlist(r) < 0.05/(nrow(normalised)-1))
```

(we’re testing against second row as the first row (for delta of 0) is obviously different from the others so all tests failing wouldn’t be unexpected)

result

```

01 | [1] 1 21 23 26 44 57 58 59 60 61 62 63 64 69 71 75 77 78
02 | [19] 80 81 82 83 92 93 94 95 96 98 99 100 103 104 109 114 118 121
03 | [37] 122 123 124 125 126 127 128 130 131 132 133 134 135 136 137 138 139 140
04 | [55] 141 142 143 144 145 146 147 160 161 163 165 169 170 171 172 173 174 175
05 | [73] 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 215
06 | [91] 217 218 249 250 251 252 253 254 255 256

```

Quite obviously different, even with just 128 samples per delta value.

Summary

Moral of the story is, don't use something without testing if it behaves as it claims to. If it does have tests, verify that they check *your* expectations, not only the programmers that wrote it in the first place.

State your assumptions and test them. If values look similar, measure them multiple times, and use statistical methods to compare them.

Test setup

Tests were performed, as previously mentioned, on an Intel i7 4790K CPU. The system was running Linux 4.17.5-1-ARCH with Python 3.6.6-1 and perf 1.5.1 from Archlinux.

Conversion from json files to csv files was performed using `json-to-csv.py` script available at the testing repo, together with raw results, at [github \(https://github.com/tomato42/python-timing-tests\)](https://github.com/tomato42/python-timing-tests).

Post scriptum

Other operations on integers, including equality are also not constant time:

```

timing-eq-timeit.py
01 | import timeit
02 | import sys
03 | import math
04 |
05 |
06 | setup = """
07 | val_a = {}
08 |
09 | val_b = {1}
10 | """
11 |
12 |
13 | fun = """val_a == val_b"""
14 |
15 |
16 | def std_dev(vals):
17 |     avg = sum(vals)/len(vals)
18 |     sum_sq = sum((i - avg)**2 for i in vals)
19 |     return math.sqrt(sum_sq / (len(vals) - 1))
20 |
21 |
22 | if __name__ == "__main__":
23 |     total_runs = 3
24 |     runs_per_process = 3
25 |     warmups = 16
26 |
27 |     runner = timeit.Timer(fun, setup=setup.format(0, 0))
28 |     number, delay = runner.autorange()
29 |     number //= 100
30 |     delay /= 100
31 |
32 |     print("will do {0} iterations per process, "
33 |           "expecting {1:7.2} s per process"
34 |           .format(number, delay), file=sys.stderr)
35 |     print("warmups:", file=sys.stderr, end='')
36 |     sys.stderr.flush()
37 |     for _ in range(warmups):
38 |         timeit.repeat(fun, setup=setup.format(0, 0), repeat=1,
39 |                        number=number)
40 |         print(".", file=sys.stderr, end='')
41 |         sys.stderr.flush()
42 |     print(file=sys.stderr)
43 |
44 |     for a in range(256):
45 |         for b in range(256):
46 |             res = []
47 |             for _ in range(total_runs // runs_per_process):
48 |                 # drop the first result as a local warmup
49 |                 res.extend(i / number for i in
50 |                             timeit.repeat(fun,
51 |                                           setup=setup.format(a, b),
52 |                                           repeat=runs_per_process+1,
53 |                                           number=number)[1:])
54 |             print(".", file=sys.stderr, end='')
55 |             sys.stderr.flush()
56 |             if std_dev(res) > 0.001:
57 |                 print("timing-eq-timeit-1.csv")

```

Execution:

```

command
01 | PYTHONHASHSEED=1 taskset -c 2 python3 \
02 | -u timing-eq-timeit.py > timing-eq-timeit-1.csv

```

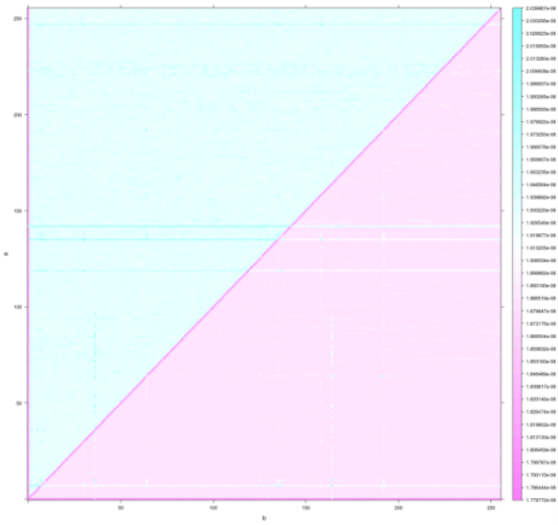
Code to create the graph:

R code

```

01 require("lattice")
02 a = read.csv(file="timing-eq-timeit-1.csv", header=FALSE,
03             col.names=seq(1, 20), fill=TRUE)
04 data = as.matrix(a)
05 med = apply(data, 1, median, na.rm=TRUE)
06 # full lines
07 len = length(med)
08 columns = ceiling(length(med) / 256)
09 d = data.frame(x=rep(seq(0, 255), length.out=len, 256),
10             y=rep(seq(0, 255), length.out=len, each=256),
11             z=med)
12 my.at = seq(min(med), max(med), length=40)
13 levelplot(z~x*y, data=d, xlab="b", ylab="a",
14          at=my.at, colorkey=list(at=my.at, labels=list(at=my.at)))

```



(<https://securitypitfalls.files.wordpress.com/2018/08/timing-eq-timeit-1.png>)

So it looks to me like the xor operation is actually one of the more constant time primitives...

Posted in [Articles](#) and tagged [constant time](#), [lucky 13](#), [python](#), [side channel](#), [testing](#), [timing attacks](#) on [2018/08/03](#) by [Hubert Kario](#). [5 Comments](#)

5 comments

1. [Alexander Kostadinov](#) says:

[2018/08/03 at 17:45](#)

Wow, thank you for the highly educational post. I think original `constant_time_compare` has another flaw as well. It will exit early when string length does not match.

I'm not a python dev but I'm wondering if this flaw can be removed as well the type conversion to be eliminated like this:

```

def constant_time_compare(password, guess):
    len_guess = len(guess)
    len_pass = len(password)
    result = len_pass == len_guess
    generated_password = ""
    for index in range(len_guess):
        generated_password += (password if index < len_pass else guess)[index]
    for index in range(len_guess):
        result = (guess[index] == generated_password[index]) & result
    return result

```

Basically timing should be based on the `guess` length always. The only thing I see that could expose password length data is that with a longer guess, we generate `generated_password` using another string. My "hope" though is that time for this would be comparable to getting from the other.

I wonder if one can get any better without adding artificial noise.

[REPLY](#)

1. [Hubert Kario](#) says:

[2018/08/03 at 18:06](#)

> It will exit early when string length does not match.

This is expected, and rather hard to exploit – it only leaks if the secret is as long as the attacker provided data. Given that you should always compare the attacker provided data to some hash (HMAC, crypt, PBKDF2) the length of that secret value will be constant. Because of Kerckhoffs's principle, we must assume that the attacker knows what hash and algorithm we are using.

> (password, guess)

To compare attacker controlled data to a password, you need to have plaintext password stored somewhere. But I'm sure you don't have them and that was just an example... 🙄

> I wonder if one can get any better without adding artificial noise.

adding artificial noise in general does not work: <https://eprint.iacr.org/2015/1129.pdf>

[REPLY](#)

1. [Alexander Kostadinov](#) says:

[2018/08/03 at 18:39](#)

Yeah, forgot we have passwords hashed and thus fixed length. In such case though we have one less thing to worry about. Now what if we avoid this data conversion with comparing characters for equality instead of xor-ing (as in my example above)? Shouldn't this eliminate the value based timing differences?

wrt noise, I remember this doc. I was thinking more about using high precision timer to sleep such that artificially make operation always take around X milliseconds. Much preferable if language allows making sure something takes constant time. But if it truly doesn't, then what are the options? And what if next version of language introduces some optimizations that invalidate previous measurements for constant time? I guess you'd say automated tests 😊

2. [Hubert Kario](#) says:

[2018/08/03 at 19:34](#)

> Now what if we avoid this data conversion with comparing characters for equality instead of xor-ing (as in my example above)

there is no "character" class in Python, and hash outputs are binary strings, which means that single "characters" (bytes, really) are integers

> Shouldn't this eliminate the value based timing differences?

I don't think so. But you have the tools, run the benchmarks 😊

> I was thinking more about using high precision timer to sleep such that artificially make operation always take around X milliseconds.

For low level operations, the differences are measured in single cycles, not doable from python level. For high level: what is the baseline? How do you ensure it doesn't change? What if thermal throttling kicks in? What if the server is under heavier load than when we measured the baseline?

> But if it truly doesn't, then what are the options?

use language that does – most compiled languages should meet this, I wouldn't be surprised if Cython fixed it (if the compiler wasn't too clever)

but in general, minimise the critical code, delegate it to languages that operate on machine intrinsics (are not interpreted) and then test those modules

> I guess you'd say automated tests 😊

to quote the Indiana Jones and the Last Crusade: "You have chosen... wisely"

3. [Hubert Kario](#) says:

[2018/08/04 at 11:39](#)

I've updated the article to include the test for the "==" operator: it's quite obviously not constant time too.

[CREATE A FREE WEBSITE OR BLOG AT WORDPRESS.COM.](#)

