

CS 1699: Privacy in the Electronic Society

Project 1 – Side-Channel Attacks

Jonathan Dyer

March 1, 2018

Contents

Task W0	1
Problem P: String Comparison	1
Task W1	1
Relevance to Privacy	1
Task W4	2
Algorithm A: Break-on-Inequality	2
Results	2
Effect on Privacy	3
Existential information	4
Secret values	4
Task W7	6
Algorithm B: Constant-Time Comparison	6
Results	6
Trade-Offs & Design Decisions	7
Task W8	8
Remaining Problems	8
References	9

Task W0

Problem P: String Comparison

The problem chosen for this assignment is *string comparison*, which involves checking one string for equality with another string. This includes comparing both the **length** and the **content** of each string for a mismatch.

Input: Two strings to be compared.

Output: Boolean indicating if the strings are *equal*.

Task W1

Relevance to Privacy

String comparison has big implications in privacy for such a seemingly innocuous problem. Comparing values is a basic, necessary task for any system involving authentication or user input, which includes most (if not all) privacy-related services. This is especially true of web apps, and string comparison has been mentioned as a potential point of weakness on a number of online resources (more on how this weakness is exploited later). [7]

For example, an application may check if a username exists or a given password is correct by comparing its

string value to the stored password for a user, whether in plaintext (unlikely for passwords) or by comparing their hash values. This typically occurs via a custom comparison for the object in question (potentially still timing sensitive!) or by using simple string comparison (both have been done historically). [2] For example, the Java `MessageDigest.isEqual` method uses a simple byte comparison to check whether two digests are equal, which is essentially the implementation of string comparison for some languages.

Task W4

Algorithm A: Break-on-Inequality

Firstly, let's examine **Algorithm A** so that we can understand how and why this naive algorithm varies in runtime (including for inputs of the same size). In the pseudocode below, first notice that if the strings are different lengths, a 'False' is returned immediately (lines 5-6). This clearly changes the runtime in the case of different-sized inputs, which can reveal private information as discussed below. More subtle is the character-by-character comparison, which breaks as soon as it finds a mismatch between the two strings (line 11). This reveals information about *how* close the two strings are (or in the case of an attacker, how correct the guess was).

```

1 // This method takes two strings, a and b, and returns True if they are equal, False otherwise
2 def algorithm_A (String a, String b)
3
4   // First check that the lengths are the same
5   if a.length != b.length
6     return False
7
8   // Now iterate through characters, comparing one by one
9   int i = 0
10  while i < a.length
11    if a[i] != b[i]
12      return False
13    i = i+1
14
15  // If we make it all the way, they match!
16  return True

```

Results

Benchmarking [12] the above code – implemented in the Ruby programming language – over 1000s of iterations gives the results displayed in the table below. Note that the "different strings of the same length" differed in the first or nearly first characters while the "close strings of the same length" differed only on the final character, for both the sentence and the hash value. The following observations (about real time taken to compare inputs) are plain from the data:

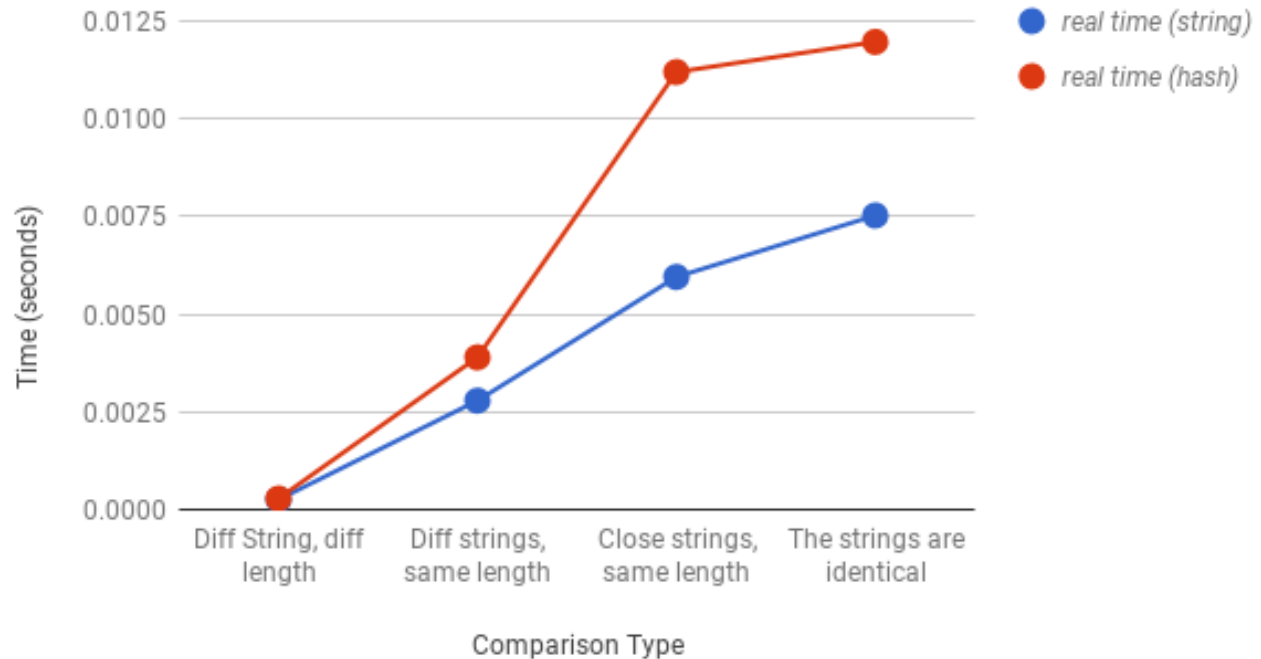
- Comparing different strings of the same length is an order of magnitude slower than strings of different length.
- Comparing nearly-identical strings is *another* order of magnitude slower.
- The difference between two close strings vs. two of the same string is very small, although it becomes more noticeable with a longer input.

Naive		This is a good dog			
Test with string	real time (string)	iterations in 5.0s	iterations / second	slower?	
Diff String, diff length	0.000255	16.386M	3.263M		
Diff strings, same length	0.002788	1.936M	386.523k	8.44x	
Close strings, same length	0.005953	860.031k	170.027k	19.19x	
The strings are identical	0.007526	863.582k	171.550k	19.02x	

		286755fad04869ca523320acce0dc6a4			
Test with hash	real time (hash)	iterations in 5.0s	i / s	slower?	
Diff String, diff length	0.000282	16.203M	3.207M		
Diff strings, same length	0.003896	1.423M	280.709k	11.43x	
Close strings, same length	0.011195	534.558k	104.985k	30.55x	
The strings are identical	0.011967	518.313k	103.680k	30.94x	

This may also be clearer from the following graph of the clock time taken by each comparison type. Notice that the difference between two different strings of the same length and two close strings is more dramatic as the length increases.

Naive Comparison



Effect on Privacy

It is clear that this measurable difference in string comparison times can be leveraged in a timing attack against some systems that make use of such an algorithm. Although my first speculation was that such an attack could be used to recover a password from a system that allows multiple login attempts, I quickly recalled my computer science class on 'Privacy in the Electronic Society', wherein we learned that it is insecure and foolish to 1) transmit passwords in cleartext, or 2) store cleartext passwords on a secure system. So the best that a timing attack could reveal in this instance (using techniques described below) is the hash of a password, or something that is normally sent in cleartext anyways. But of course due to the cryptographic properties of any good hash function, recovering the password itself from such information is all but impossible. [8] After further research and armchair rationalization, I determined that there are at least **two categories** of information that can be gained as a consequence of the above algorithm. The general setup for an attacker requires that:

1. The system being attacked makes use of Algorithm A, above.
2. The attacker be able to choose the input to the algorithm (making this a *chosen plaintext* attack).
3. The attacker has some way of **verifying** whether or not the input was accepted; for example, a login form for a web app must indicate to the user (attacker) whether or not the input was correct, rather than simply sitting without any response.

If the above conditions hold, then the system is functioning (loosely) as what is known as a *verification oracle*, meaning that it will verify as correct or incorrect whatever input you provide, specifically by way of the (vulnerable) algorithm given above. This can leak information by way of a simple attack outlined here:

1. For any given set x of known input characters and $y = y_j \dots y_n$ unknown, try every possibility for a variable a in the input $x||a||y_{j+1} \dots y_n$ and some fixed filler, say `<underscore>`, for all $y_i \neq y_j$.
2. Sample some dozens of times (or more if necessary, over a network, etc.) on each combination, and then take the mean (or median to mitigate outliers) of the time taken to return an error/False message.

3. Set a to be the character/byte with the greatest mean/median value (i.e. that takes the longest), add it to x , and then return to step 1 with $y = y_{j+1} \dots y_n$ and $a = y_{j+1}$.

By repeating this process, first to find the correct length and then to find the correct combination of characters, it is possible to discern one or both of the following two types of (private) information. This may require a more sophisticated timing software than Ruby's built-in benchmarking module (which I used), but I believe that it would not be difficult to achieve such an attack, even with a basic timing program.

Existential information

This type of information is leaked when an attacker exploits the timing differences in the above algorithm to discover whether or not a system's stored list of information contains a given value. An attacker takes note of the typical response times for nonsense inputs and then compares them to response times for potentially legitimate input values. It is even simpler than the outline given above. An example will help to clarify:

- There is a web application called Foo Services that allows users to create accounts and login with an email address and password. It turns out that Foo uses Algorithm A to check the input email against valid emails to determine if such an account exists.
- The attacker, Mallory, creates an account with Foo Services. Mallory notices that when she mistypes her email in the login form, the Foo site displays an error: **Invalid username or password**.
- Mallory wants to discover if any of her friends have an account with Foo Services, so she follows the attack format given above, but instead of testing combinations of characters she simply cycles through her list of known email contacts, testing each several times (or dozens, or hundreds) and recording the average/median time that each takes to return the error message (obviously she doesn't know the password for any of their accounts, so she leaves that field blank for every attempt, or enters in an identical filler password for each attempt).
- Finally, Mallory enters her own email address (with the wrong password) and again tests several times to find the average/median time taken for a *valid* account.
- According to the results section above, with enough samples Mallory should be able to tell which emails correspond to extant accounts. How? She can do this by noticing a (close to linear) boundary between those accounts that returned an error message more quickly and those accounts (including her own) which took longer. There may be some diffusion if some of the email addresses were similar to real accounts, but even that relationship should be distinguishable according to the results above.
- At this point Mallory has discovered which of her friends have accounts with Foo Services (with the email addresses she is aware of). She may leverage this information to perform a more targeted attack on one account (perhaps involving the password), or if Foo Services provides a sensitive or embarrassing service, then the existence information itself is enough to compromise that individual's privacy.

Some variant of this attack is applicable in many circumstances. Some examples include checking "forgot my password" email tokens, testing answers to "security questions", or login forms as above. [7] Note that an attack similar to this could be perpetrated on older Unix systems, whose `login` program had a timing variation (execution of the `crypt` function) that depended on the string being recognized as a username. Combined with a dictionary attack on passwords, this could compromise the security of a system. [11]

This attack seems somewhat trivial, and perhaps only damaging in special circumstances. An extension of it that utilizes the full power of the attack outline first described would allow discovery of values that are associated with user information but typically hidden from users or inaccessible to outside attackers. It is this variant that we turn to next.

Secret values

This type of information is revealed when an attacker follows the iterative attack outline given previously. The examples I discuss all refer to hash values being uncovered, but there may be other applications I haven't considered or listed here.

Suppose that there is a system, perhaps another web-based service, that authenticates users according to some keyed hash—a MAC or HMAC value for instance. This could be to authenticate a session cookie in a web browser [2], or to verify the author of some firmware for a system upgrade, or any other scenario wherein the hash of a

password is used to confirm an identity. Also suppose that the three prerequisites defined at the beginning of this section all hold in relation to this hash value. That is:

1. The system being attacked makes use of Algorithm A to compare the given authenticating hash value with the expected value.
2. The attacker is able to choose the hash that is passed to the system to be verified.
3. The system informs the user (attacker) of the results of that comparison in some way.

In other words, the system is functioning as a hash verification oracle for the attacker, who can pass arbitrary plaintext data along with a hash and in return find out whether or not they should "go together" (in the context of the system in question). Thus, the attacker might the more advanced character-by-character attack that we outlined before to retrieve the corresponding hash for a given input data. It is provided again for convenience:

1. For any given set x of known input characters and $y = y_j \dots y_n$ unknown, try every possibility for a variable a in the input $x||a||y_{j+1} \dots y_n$ and some fixed filler, say **<underscore>**, for all $y_i \neq y_j$.
2. Sample some dozens of times (or more if necessary, over a network, etc.) on each combination, and then take the mean (or median to mitigate outliers) of the time taken to return an error/False message.
3. Set a to be the character/byte with the greatest mean/median value (i.e. that takes the longest), add it to x , and then return to step 1 with $y = y_{j+1} \dots y_n$ and $a = y_{j+1}$.

Now, retrieval of such a hash alone may not be explicitly detrimental to privacy, since (as mentioned before) it does not enable the recovery of a preimage/password – at least for any hash function worth its "salt" (pun intended). **However**, certain attacks are facilitated specifically by having the hash of a value. Three particular examples are briefly described below, and further references are given.

1. A straightforward attack known as "hash length extension" [5] allows an attacker to falsely authenticate an "extension" of his choosing. This attack is valid for all Merkle-Damgard based hash algorithms, even for keyed hashes when the key is concatenated to the message (rather than an HMAC). A simple example is given at [9].
2. One attack, known as 'pass-the-hash', provides authentication to a system or service without the password for a given user, but just the hash. It applies on all systems using either the NTLM protocol or LanMan hash to authenticate users to the system. [10] Newer operating systems have disabled compromised versions of these protocols by default, but some of them continue to be used by third-party developers. If an attacker can acquire such a hash, then it would be possible to impersonate a user via a remote protocol such as SMB, FTP, HTTP, or others. The attack may be more sophisticated than our previous examples. Further reading: [4].
3. Finally, and perhaps most **notably**, an actual timing attack involving hash-based authentication and string comparison was used to bypass a security measure in the Xbox 360 and modify system software (thereby enabling the running of custom software on the machine). This attack worked by exploiting time differences in verifying valid MAC values for kernel updates: for each valid byte in the hash, the system would reject about 2.2 milliseconds more slowly than an invalid byte in the equivalent position, enabling *exactly our attack outlined above* to find a valid hash value. Details of this attack are presented in [1], and this and other MAC verification comparison timing attacks are mentioned in [6].

There are other examples and resources besides the above elucidating the possibilities resulting from string comparison inconsistency. For example the popular crypto library developed by Google, known as Keyczar, fell prey to this problem in their HMAC verification implementation [3], thus it is not unheard of even in well-developed libraries.

Task W7

Algorithm B: Constant-Time Comparison

Now let's examine **Algorithm B** so that we can understand how and why this revised algorithm runs in a more constant-time manner, mitigating some of the effects discussed for Algorithm A above. In the pseudocode below, first notice that before anything else we select a consistent number of iterations to run the byte-wise comparisons for (`len` in the code below, line 5)—this helps ensure that we don't reveal information about the length of the secret value (perhaps a login name), since that was an obvious and large timing difference in our naive algorithm. Next, consider the revised loop through our inputs (lines 9-13). Rather than breaking this loop once we find a mismatch, we continue for the specified length `len` (line 9). Then we check that both of the next characters are even valid to be compared (line 10). If so, we skip to the `else` (line 12) and proceed with the comparison. If either one is *not* valid, we execute a "dummy" comparison to (hopefully) take the same amount of time as a real one (line 11). Finally, we return the resulting boolean.

```
1 // This method takes two strings, a and b, and returns True if they are equal, False otherwise
2 def algorithm_B (String a, String b)
3
4     // We'll assume string a is the "system string" we're comparing the user-input against
5     int len = a.length
6
7     // Now set a boolean and then iterate through characters, comparing each (even on mismatch)
8     boolean bool = (a.length == b.length)
9     for i in 0... len // for every character in the fixed length
10         if a[i] == null || b[i] == null // if we've reached the end of one string
11             bool = (a[0] == a[0]) && bool // compare something anyways
12         else // otherwise
13             bool = (a[i] == b[i]) && bool // compare the two chars, update the boolean
14
15     return bool
```

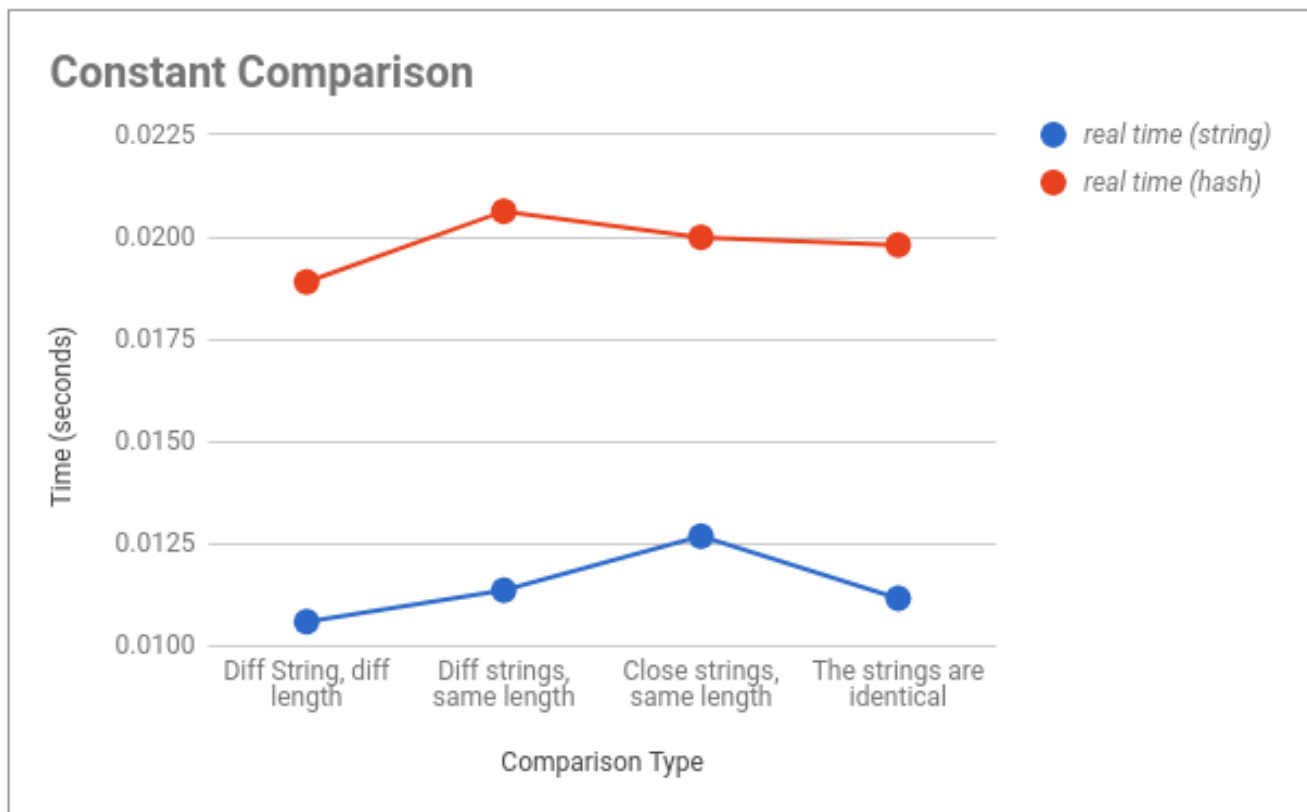
Results

Benchmarking the above code – implemented in Ruby – over 1000s of iterations gives the results displayed in the table below. Note that the timing distinctions between the different types of comparison are almost entirely eliminated for shorter strings, and are statistically insignificant for longer strings/hash values. It is worth noting that comparing strings of different length still seems to be notably faster, especially for shorter inputs. This could reveal the length of a secret, given enough samples.

Constant	This is a good dog				
Test with string	real time (string)	iterations in 5.0s	i / s	slower?	
Diff String, diff length	0.010596	495.144k	97.921k		
Diff strings, same length	0.011366	431.600k	85.590k	1.14x	
Close strings, same length	0.012686	431.400k	86.376k	1.13x	
The strings are identical	0.011168	439.416k	87.895k	1.11x	

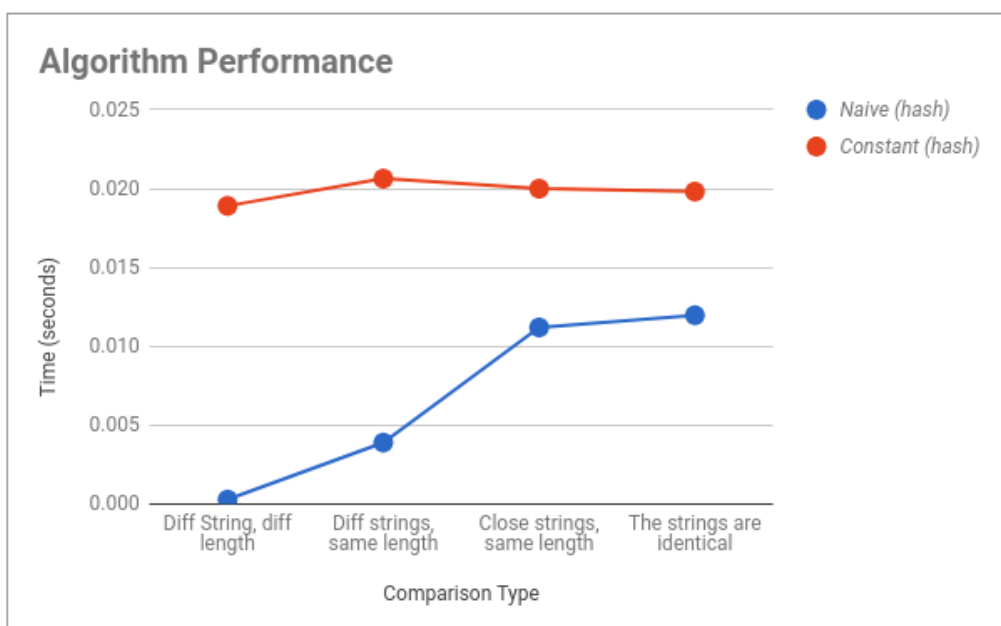
286755fad04869ca523320acce0dc6a4					
Test with hash	real time (hash)	iterations in 5.0s	i / s	slower?	
Diff String, diff length	0.018902	259.632k	52.482k		
Diff strings, same length	0.020629	255.500k	51.108k	within error	
Close strings, same length	0.019993	260.355k	51.922k	within error	
The strings are identical	0.019806	263.648k	52.152k	within error	

The improvement is especially clear given the graph below, which shows much more even timing results across different types of input. In particular, the test on longer (hash) strings shows a slight reversal in the trend across strings of the same length—which could indicate that it would lead an attacker to make incorrect guesses (a desirable trait).



Trade-Offs & Design Decisions

There is a significant performance cost associated with making constant-time comparisons, especially when trying to obscure even the length of the secret value(s). Notice in the comparison graph below for Naive vs. Constant Comparison that the constant-time algorithm takes nearly twice as long when comparing similar strings and *several times* as long when comparing very different strings. This would be a highly significant cost on a system or web-server that may receive thousands of requests a day/hour/minute.



Additionally, not all the timing difference for comparing strings of different length is mitigated. This could be because the 'else' conditional code (lines 12-13) is not being executed for every null character of the compared string, resulting in faster overall execution. Some amount of randomness could perhaps be introduced that would further obscure those lines of comparison.

Finally, because the chosen amount of time is based on the *first* string passed to the algorithm, we risk revealing length of secret values anyways if our implementation is known and the attacker can duplicate our system/working environment. An alternative to this would be to always go with the longer string, but that might open up the algorithm to DOS (denial-of-service) attacks, where the attacker simply passes excessive input. This could be checked for, but length information could be revealed by passing very short arguments and working up to the secret string length. On the other hand, if we choose to always go with the shorter string, the inverse process would also reveal lengths. One solution to this problem may be to choose a random `len` value within an acceptable range and iterate for that many bytes (but this may not mitigate all leakage).

Task W8

Remaining Problems

There is still potential to leak lengths of secret values (as mentioned in the trade-offs section). There are likely better ways to obscure those lengths that I have not implemented above, without sacrificing quite as much in performance. One option is randomizing the number of byte-comparison iterations as suggested above. However, it is plain that there will always be a performance cost associated with implementing a constant-time algorithm (since we can only come *up* to the slowest time, rather than decreasing everything else).

It is also possible that the very timing attack described above could work with my new and improved algorithm, perhaps with more samples. This would require a more detailed and granular timing software, probably written in a lower-level language (e.g. C), but would almost certainly still be possible for a dedicated attacker – especially given the comparison graph above, which still reveals some variation.

It is clear that ensuring no information is leaked by a string comparison algorithm is a nontrivial task. We've learned not only that, but that this task is in fact important to securing applications and systems that involve private information. Better implementations should be (and are more frequently) included in complete and well-tested cryptographic libraries.

References

- [1] IVC Wiki Contributors. *Xbox 360 Timing Attack*. 2007. URL: http://beta.ivc.no/wiki/index.php/Xbox_360_Timing_Attack. (accessed 28 Feb 2018).
- [2] Coda Hale. *A Lesson In Timing Attacks*. 2009. URL: <https://codahale.com/a-lesson-in-timing-attacks/>. (accessed 28 Feb 2018).
- [3] Nate Lawson. *Timing attack in Google Keyczar library*. 2009. URL: <https://rdist.root.org/2009/05/28/timing-attack-in-google-keyczar-library/>. (accessed 1 Mar 2018).
- [4] Bashar Ewaida. “Pass-the-hash attacks: Tools and Mitigation”. In: *SANS Institute Reading Room* (2010). URL: <https://www.sans.org/reading-room/whitepapers/testing/pass-the-hash-attacks-tools-mitigation-33283>.
- [5] Douglass Clem. *Hash Length Extension Attacks*. 2012. URL: <https://www.whitehatsec.com/blog/hash-length-extension-attacks/>. (accessed 27 Feb 2018).
- [6] Yehuda Lindell Jonathan Katz. *Introduction to Modern Cryptography, Second Edition*. CRC Cryptography and Network Security Series. CRC Press, 2014, pp. 115–116. ISBN: 9781466570269.
- [7] Nick Malcolm. *Timing Attacks against String Comparison*. 2016. URL: <https://thisdata.com/blog/timing-attacks-against-string-comparison/>. (accessed 1 Mar 2018).
- [8] Wikipedia Contributors. *Cryptographic hash function*. Wikipedia, The Free Encyclopedia. URL: https://en.wikipedia.org/wiki/Cryptographic_hash_function. (accessed 1 Mar 2018).
- [9] Wikipedia Contributors. *Length extension attack*. Wikipedia, The Free Encyclopedia. URL: https://en.wikipedia.org/wiki/Length_extension_attack. (accessed 27 Feb 2018).
- [10] Wikipedia Contributors. *Pass the hash*. Wikipedia, The Free Encyclopedia. URL: https://en.wikipedia.org/wiki/Pass_the_hash. (accessed 1 Mar 2018).
- [11] Wikipedia Contributors. *Timing attack*. Wikipedia, The Free Encyclopedia. URL: https://en.wikipedia.org/wiki/Timing_attack#Examples. (accessed 1 Mar 2018).
- [12] Rdoc Generator. *Module: Benchmark (Ruby 2.5.0)*. Ruby-doc.org. URL: <https://ruby-doc.org/stdlib-2.5.0/libdoc/benchmark/rdoc/Benchmark.html>. (accessed 26 Feb 2018).