

Here is a compression function, which is provably strongly collision-free.

However, it is rather slow and so is not practical for hashing larger data. On the other hand, its slowness could be beneficial for applications like password hashing.

**Example 199. (the discrete log hash)** Let  $p$  be a large safe prime (that is,  $q = (p - 1) / 2$  is also prime). Let  $g_1, g_2$  be two primitive roots modulo  $p$ . Define the compression function  $h$  as:

$$h: \{0, 1, \dots, q^2 - 1\} \rightarrow \{1, 2, \dots, p - 1\}, \quad h(m_1 + m_2q) = g_1^{m_1} g_2^{m_2} \pmod{p}.$$

[Note that, although not working with inputs and outputs of certain size in bits, this is a compression function, because the input space is much larger than the output space.]

Show that finding a collision of  $h(x)$  is as difficult as determining the discrete logarithm  $x$  in  $g_1^x = g_2 \pmod{p}$ .

**Solution.** Suppose we have a collision:  $g_1^{m_1} g_2^{m_2} \equiv g_1^{m'_1} g_2^{m'_2} \pmod{p}$

Hence,  $g_1^{(m_1 - m'_1) + (m_2 - m'_2)x} \equiv 1 \pmod{p}$  or, equivalently,  $(m_1 - m'_1) + (m_2 - m'_2)x \equiv 0 \pmod{p - 1}$  (because  $g_1$  is a primitive root and so has order  $p - 1$ ).

This final congruence can now be solved for  $x$ .

More precisely, if  $d = \gcd(m_2 - m'_2, p - 1)$ , there are actually  $d$  solutions for  $x$ . Since we chose  $p$  to be safe, the only factors of  $p - 1$  are  $1, 2, (p - 1) / 2, p - 1$ .

Since  $|m_2 - m'_2| < q$ , the only possibilities are  $d = 1, 2$  (unless  $m_2 = m'_2$ ; however, this cannot be the case since then also  $m_1 = m'_1$ , so that we wouldn't have a collision in the first place).

## Passwords

Let's say you design a system that users access using personal passwords. Somehow, you need to store the password information.

- The worst thing you can do is to actually store the passwords  $m$ .
  - This is an absolutely atrocious choice, even if you take severe measures to protect (e.g. encrypt) the collection of passwords.
  - Comment.** Sadly, there is still systems out there doing that. An indication that this might\* be happening is systems that require you to update passwords and then complain that your new password is too close to the original one. Any reasonably designed system should never learn about your actual password in the first place!
  - \*: On the other hand, think about how you could check for (certain kinds of) closeness of passwords without having to store the actual password.
  
- Better, but still terrible, is to instead store hashes  $H(m)$  of the passwords  $m$ .
  - Good.** An attacker getting hold of the password file, only learns about the hash of a user's password. Assuming the hash function is one-way, it is infeasible for the attacker to determine the corresponding password (if the password was random!!).
  - Still bad.** However, passwords are (usually) not random. Hence, an attacker can go through a list of common passwords (dictionary attack), compute the hashes and compare with the hashes of users (similarly, a brute-force attack can simply go through all possible passwords).
  - Even worse, it is immediately obvious if two users are using the same password (or, if the same user is using the same password for different services using the same hash function).
  - Comment.** So, storing password hashes is not OK unless all passwords are completely random.

- Better, a random value  $s$  is generated for each user, and then  $s$  and  $H(m, s)$  are stored. The value  $s$  is referred to as **salt**.

In other words, instead of storing the hash of the password  $m$ , we are storing the hash of the salted password, as well as the salt.

**Why?** Two users using the same password would have different salt and hence different hashes stored. As a consequence, an attacker can (of course) still mount a dictionary or brute-force attack but only against a single user, not all users at once.

**Comment.** Note how the concept of salt is similar to a nonce.

**Comment.** To be future-proof, the hash+salt is often stored in a single field in a format like (hash-algo, salt, salted hash).

**Comment.** There's also the concept of **pepper** (usually, sort of a secret salt). This provides extra security if the pepper is stored separately. [Sometimes pepper is used as a sort of small random salt, which is discarded; this only slows a brute-force attack down and should instead be addressed using the item below.]

[https://en.wikipedia.org/wiki/Pepper\\_\(cryptography\)](https://en.wikipedia.org/wiki/Pepper_(cryptography))

- Finally, we should not use the usual (fast!) hash functions like SHA-2.

**Why?** One of the things that makes SHA-2 a good hash function in practice is its speed. However, that actually makes SHA-2 a poor choice in this context of password hashing. An attacker can compute billions of hashes per second, which makes a dictionary or brute-force attack very efficient.

To make a dictionary or brute-force attack impractical, the hashing needs to be slowed down. See Example 200 for some scary numbers.

Hashing functions like SHA-2 are not secure password hashing algorithms.

Instead, options that are considered secure include: PBKDF2, bcrypt, scrypt.

**Comment.** For instance, WPA2 uses PBKDF2 based on SHA-1 with 4096 (fairly small!) iterations.

**Comment.** Only increasing the number of iterations increases computation time but not memory usage. scrypt is designed to also consume an arbitrarily specified amount of memory.

For a nice discussion about password hashing:

<https://security.stackexchange.com/questions/211/how-to-securely-hash-passwords>

**Example 200. (the power of brute-force)** In April 2019, the Bitcoin hashrate is about  $40E=4 \cdot 10^{19}$  hashes per second. How long would it take to brute-force a (completely random!) 8 character password, using all 94 printable ASCII characters (excluding the space)?

**Solution.** There are  $94^8 \approx 6.1 \cdot 10^{15}$  possible passwords. Hence, it would take about 0.00015 seconds!

**Comment.** Even using 10 random characters (almost no human password has that kind of entropy), there are  $94^{10} \approx 5.4 \cdot 10^{19}$  possible passwords. It would take less than 1.4 seconds to go through all of these!

**Comment.** <https://bitinfocharts.com/comparison/bitcoin-hashrate.html>

**Example 201.** Your king's webserver contains the following code to check whether the king is accessing the server.

[As is far too common, his password derives from his girlfriend's name and year of birth.]

```
def is_king(password):
    return password=="Ludmilla1310"
```

Obviously, anyone who might be able to see the code (including its binary version) learns about your king's password. With minimal change, how can this be fixed?

**Solution.** The password should be hashed. For instance, in Python, using SHA-2 (why is that actually not a good choice here?) with 256 output bits:

```
from hashlib import sha256
def is_king(password):
    phash = sha256(password).hexdigest()
    return phash == "9e4b4fe180e22bc6cdf01fe9711cf2558507e5c3ae1c3c1f6607a25741941c66"
```

**Comment.** 256 bit is 64 digits in hexadecimal.

**Comment.** Of course, a real implementation should use `digest()` instead of `hexdigest()`.

**Why is SHA not good here?** Too fast to discourage brute-force attacks.

**Example 202.** Suppose you don't like the idea of creating random salt.

- (a) How about using the same salt for all your users?
- (b) Is it a good idea to use the username as salt?

**Solution.**

- (a) This is a terrible idea and defeats the purpose of a salt. (For instance, again an attacker can immediately see if users have the same password.)

**Comment.** Essentially, this is a form of pepper (if the value is kept secret, i.e. stored elsewhere).

- (b) That is a reasonable idea. One reason against it is that, ideally, the salt should be unique (globally). However, this could be easily achieved by using the username combined with something identifying your service (like your hostname).

**Comment.** A possible practical reason against choosing the username for salt is that the username might change.

**Example 203.** You need to hash (salted) passwords for storage. Unfortunately, you only have SHA-2 available. What can you do?

**Solution.** Iterate many times! (In order to slow down the computation of the hash.) The naive way would be to simply set  $h_0 = H(m)$  and  $h_{n+1} = H(h_n)$ . Then use as hash the value  $h_N$  for large  $N$ .

In current applications, it is typical to choose  $N$  on the order of 100,000 or higher (depending on how long is reasonable to have your user wait each time she logs in and needs her password hashed for verification).