

On rearrangement hashing with Haskell

Roman Maksimovich

Abstract

In this paper I introduce and develop a mathematical method of producing a cryptographic hash of adjustable length, given a public key and a private key. The hashing is done through encoding selections and permutations with natural numbers, and then composing the hash from a set of source strings with respect to the permutations encoded by the keys. The attempts to construct a suitable integer-to-selection mapping leads to interesting mathematical definitions and statements, which are discussed in this paper and applied to give bounds on the reliability of the hashing algorithm. An implementation is provided in the Haskell programming language (source available at <https://github.com/thornoar/password-hash>) and applied in the setting of password creation. In the paper, the details of the implementation are discussed, as well as the connections between it and the corresponding mathematical model.



March 23, 2024

Contents

1	Introduction	3
2	The theory	3
2.1	Preliminary terminology and notation	3
2.2	Enumerating list selections	4
2.3	Elevating the choice function	6
2.4	The hash function	7

1 Introduction

The motivation behind the topic lies in the management of personal passwords. Nowadays, the average person requires tens of different passwords for different websites and services. Overall, one can distinguish between two ways of managing this set of passwords:

- **Keeping everything in one's head.** This is a method employed by many, yet it inevitably leads to certain risks. First of all, in order to fit the passwords in memory, one will probably make them similar to each other, or at least have them follow a simple pattern like "[shortened name of website]+[fixed phrase]". As a result, if even one password is guessed or leaked, it will be almost trivial to retrieve most of the others, following the pattern. Furthermore, the passwords themselves will tend to be memorable and connected to one's personal life, which will make them easier to guess. There is, after all, a limit to one's imagination.
- **Storing the passwords in a secure location.** Arguably, this is a better method, but there is a natural risk of this location being revealed, or of the passwords being lost, especially if they are stored physically on a piece of paper. Currently, various "password managers" are available, which are software programs that will create and store your passwords for you. It is usually unclear, however, how this software works and whether it can be trusted with one's potentially very sensitive passwords. After all, guessing the password to the password manager is enough to have all the other passwords exposed.

In this paper I suggest a way of doing neither of these things. The user will not know the passwords or have any connection to them whatsoever, and at the same time the passwords will not be stored anywhere, physically or digitally. In this system, every password is a cryptographic hash produced by a fixed hashing algorithm. The algorithm requires two inputs: the public key, i.e. the name of the website or service, and the private key, which is an arbitrary positive integer known only to the user. Every time when retrieving a password, the user will use the keys to re-create it from scratch. Therefore, in order to be reliable, the algorithm must be "pure", i.e. must always return the same output given the same input. Additionally, the algorithm must be robust enough so that, even if a hacker had full access to it and its working, they would still not be able to guess the user's private key or the passwords that it produces. These considerations naturally lead to exploring pure mathematical functions as hashing algorithms and implementing them in a functional programming language such as Haskell.

2 The theory

There are many ways to generate hash strings. In our case, these strings are potential passwords, meaning they should contain lower-case and upper-case letters, as well as numbers and special characters. Instead of somehow deriving such symbol sequences directly from the public and private keys, we will be creating the strings by selecting them from a pre-defined set of distinct elements (i.e. the English alphabet or the digits from 0 to 9) and rearranging them. The keys will play a role in determining the rearrangement scheme. With regard to this strategy, some preliminary definitions are in order.

2.1 Preliminary terminology and notation

Symbols A , B , C will denote arbitrary sets (unless specified otherwise). \mathbb{N}_0 is the set of all non-negative integers.

By E we will commonly understand a finite set of distinct elements, called a *source*. When multiple sources E_0, E_1, \dots, E_{N-1} are considered, we take none of them to share any elements between each other. In other words, their pair-wise intersections will be assumed to be empty. By $|E|$ we will denote the cardinality of a source E , and $E:i$ will represent its i -th element, with the numeration starting from $i = 0$. On the opposite, the expression $E!i$ will denote the set difference $E \setminus \{E:i\}$.

The symbol " $\#$ " will be used to describe the number of ways to make a combinatorial selection. For example, $\#^m(E)$ is the number of ways to choose m elements from a source E with significant order.

The expression $[A]$ will denote the set of all ordered lists composed from elements of the set A . The subset $[A]_m \subset [A]$ will include only the lists of length m . Extending the notation, we will define $[A_0, A_1, \dots, A_{N-1}]$ as the set of lists $\alpha = [a_0, a_1, \dots, a_{N-1}]$ of length N where the first element is from A_0 , the second from A_1 , and so on, until the last one from A_{N-1} . Finally, if $\alpha \in [A]$ and $\beta \in [B]$, the list $\alpha \mathbin{++} \beta \in [A \cup B]$ will be the concatenation of lists α and β .

Let $k \in \mathbb{N}_0$, $n \in \mathbb{N}$. The numbers ${}^Nk, {}_Nk \in \mathbb{N}_0$ are defined to be such that $0 \leq {}^Nk < N$ and ${}_Nk \cdot N + {}^Nk = k$. The number Nk is the remainder after division by N , and ${}_Nk$ is the result of division.

For a number $N \in \mathbb{N}$, the expression (N) will represent the semi-open integer interval from 0 to N : $(N) = \{0, 1, \dots, N-1\}$.

Let $n, m \in \mathbb{N}$, $m \leq n$. The quantity $n!/(n-m)!$ will be called a *relative factorial* and denoted by $(n \mid m)!$.

If f is a function of many arguments a_0, \dots, a_{n-1} , the expression $f(a_0, \dots, a_{i-1}, -, a_{i+1}, \dots, a_{n-1})$ will represent the function of one argument a_i where all others are held constant.

2.2 Enumerating list selections

The defining feature of the public key is that it is either publicly known or at least very easy to guess. Therefore, it should play little role in actually encrypting the information stored in the private key. It exists solely for the purpose of producing different passwords with the same private key. So for now we will forget about it. In this and the following subsection we will focus on the method of mapping a private key $k \in \mathbb{N}_0$ to an ordered selection from a set of sources in an effective and reliable way.

Definition 2.1. Let E be a source, $k \in \mathbb{N}_0$. The *choice function of order 1* is defined as the following one-element list:

$$C^1(E, k) = [E: {}^{|E|}k].$$

It corresponds to picking one element from the source according to the key. For a fixed source E , the choice function is periodic with a period of $|E|$ and is injective on the interval $(|E|)$ with respect to k . Injectivity is a very important property for a hash function, since it determines the number of keys that produce different outputs. When describing injectivity on intervals, the following definition proves useful:

Definition 2.2. Let A be a finite set and let $f: \mathbb{N}_0 \rightarrow A$ be a function. The *spread* of f is defined to be the largest number n such that, for all $k_1, k_2 \in \mathbb{N}_0$, $k_1 \neq k_2$, the following implication holds:

$$f(k_1) = f(k_2) \implies |k_1 - k_2| \geq n.$$

This number exists due to A being finite. We will denote this number by $\text{spr}(f)$.

Trivially, if $\text{spr}(f) \geq n$, then f is injective on (n) , but the inverse is not always true. Therefore, a lower bound on the spread of a function serves as a guarantee of its injectivity. Furthermore, if $\text{spr}(f) \geq n$ and f is bijective on (n) , then f is periodic with period n and therefore has a spread of exactly n . We leave this as a simple exercise for the reader.

Proposition 2.3. Let $f: \mathbb{N}_0 \rightarrow A$, $g: \mathbb{N}_0 \rightarrow B$ be functions such that $\text{spr}(f) \geq n$ and $\text{spr}(g) \geq m$. Define the function $h: \mathbb{N}_0 \rightarrow [A, B]$ as follows:

$$h(k) = [f({}^nk), g({}_nk + T({}^nk))],$$

where $T: \mathbb{N}_0 \rightarrow \mathbb{N}_0$ is a fixed function, referred to as the argument shift function. It is then stated that $\text{spr}(h) \geq nm$.

Proof. Assume that $k_1 \neq k_2$ and $h(k_1) = h(k_2)$. Since h returns an ordered list, the equality of lists is equivalent to the equality of all their corresponding elements:

$$f({}^nk_1) = f({}^nk_2), \quad (1)$$

$$g({}_nk_1 + T({}^nk_1)) = g({}_nk_2 + T({}^nk_2)). \quad (2)$$

Since f is injective on (n) , we see that ${}^nk_1 = {}^nk_2$. Consequently, it follows from $k_1 \neq k_2$ that ${}_nk_1 \neq {}_nk_2$ and ${}_nk_1 + T({}^nk_1) \neq {}_nk_2 + T({}^nk_2)$. We can then proceed to utilize the definition of spread for the function g :

$$\begin{aligned} |{}_nk_1 + T({}^nk_1) - {}_nk_2 - T({}^nk_2)| &\geq m, \\ |{}_nk_1 - {}_nk_2| &\geq m, \\ \left| \frac{k_1 - {}^nk_1}{n} - \frac{k_2 - {}^nk_2}{n} \right| &\geq m, \\ \left| \frac{k_1 - k_2}{n} \right| &\geq m, \\ |k_1 - k_2| &\geq nm. \end{aligned}$$

■

With this proposition at hand, we have a natural way of extending the definition of the choice function:

Definition 2.4. Let E be a source with cardinality $|E| = n$, $k \in \mathbb{N}_0$, $2 \leq m \leq n$. The choice function of order m is defined recursively as

$$\mathcal{C}^m(E, k) = [E: {}^nk] ++ \mathcal{C}^{m-1}(E', k'),$$

where $E' = E! {}^nk$ and $k' = {}_nk + T({}^nk)$, while $T: \mathbb{N}_0 \rightarrow \mathbb{N}_0$ is a fixed argument shift function.

Proposition 2.5. Let E be a source with cardinality n . Then the choice function $\mathcal{C}^m(E, k)$ of order $m \leq n$, as a function of k , has a spread of at least $(n | m)!$.

Proof. We will conduct a proof by induction over m . In the base case, $m = 1$, we notice that $(n | m)! = n$, and the statement trivially follows from the definition of $\mathcal{C}^1(E, k)$.

Let us assume that the statement is proven for all choice functions of order $m-1$. Under closer inspection it is clear that the definition of $\mathcal{C}^m(E, k)$ follows the scheme given in proposition 2.3, with $\mathcal{C}^1(E, k)$ standing for f and $\mathcal{C}^{m-1}(E', k')$ standing for g . The application of the proposition is not straightforward, and we encourage the reader to consider the caveats. Thus, we can utilize the statement of the proposition as follows:

$$\text{spr}(\mathcal{C}^m(E, -)) \geq \text{spr}(\mathcal{C}^1(E, -)) \cdot \text{spr}(\mathcal{C}^{m-1}(E', -)) \geq n \cdot ((n-1) | (m-1))! = (n | m)!,$$

q.e.d.

■

The preceding result is especially valuable considering the fact that there are exactly $(n \mid m)!$ ways to select an ordered sub-list from a list, meaning that $\mathcal{C}^m(E, k)$ is not only injective, but also surjective with respect to k on the interval $((n \mid m)!).$ This makes it a bijection

$$\mathcal{C}^m(E, -): ((n \mid m)!) \rightarrow [E]_m,$$

and therefore a periodic function with a spread of exactly $(n \mid m)! = \#^m(E).$

These properties make the choice function a fine candidate for a hash mapping. Suppose that the source E is composed from lower-case and upper-case Latin characters, as well as special symbols and digits:

$$E = \text{"qwertyuiopasdfghjklzxcvbnmQWERTYUIOPASDFGHJKLZXCVBNM0123456789!@#\$%"}.$$

The choice function gives us a way to enumerate all possible ways to select a sub-list from E . What's more, these selections can be made more "random" and unpredictable by means of complicating the argument shift function T . A reasonable practice is to set $T({}^n k)$ to the ASCII value of the character $E: {}^n k$. This way, each chosen character will influence the choice of the next, creating what is called a "chaotic system", where its behavior is fully determined, but even small changes to inputs eventually produce large changes in the output. Here is a little input-output table for the choice function of order 10 with the specified source and shift function:

123	"41BeGs9\$Dd"
124	"52NgJfZIk7"
125	"63MfHs9\$Da"
126	"740VbDo6@u"
127	"851Br469\$S"

There is, however, a serious problem. This selection method does not guarantee that the chosen 10 symbols will contain lower-case and upper-case characters, as well as digits and spacial symbols, all at the same time. Since the choice function is bijective, there is a key that produces the combination "djaktpsnei", which will not be accepted as a password in many placed, because it contains only one category of symbols. Fortunately, there is a solution.

2.3 Elevating the choice function

Definition 2.6. Let α be the list of pairs (E_i, m_i) , where E_i are sources, $|E_i| = n_i$, $m_i \leq n_i$, for $i \in (N)$. The *elevated choice function* corresponding to these data is defined for a key $k \in \mathbb{N}_0$ by means of the following recursion:

$$\bar{\mathcal{C}}(\alpha, k) = [\mathcal{C}^{m_0}(E_0, {}^{n_0}k)] ++ \bar{\mathcal{C}}(\alpha ! 0, {}_{n_0}k + T({}^{n_0}k)),$$

where T is an argument shift function. The base of the recursion is given when α is empty, in which case $\bar{\mathcal{C}}([], k) = []$. Otherwise, for every key k , its image is an element of

$$\text{cod } \bar{\mathcal{C}}(\alpha, -) = [[E_0]_{m_0}, [E_1]_{m_1}, \dots, [E_N]_{m_N}].$$

In this context, the list α will be called a *source configuration*.

In other words, the elevated choice function is a "mapping" of the choice function over a list of sources, it selects a sub-list from every source and then composes the results in a list, which we will call a *multiselection*. A trivial application of proposition 2.3 shows that the spread of $\bar{\mathcal{C}}(\alpha, -)$ is at least

$$\prod_{i=0}^{N-1} \text{spr}(\mathcal{C}^{m_i}(E_i, -)) = \prod_{i=0}^{N-1} (n_i \mid m_i)! \quad (3)$$

where E_i , n_i , and m_i compose the configuration α . In fact, due to the rule of product in combinatorics, we see that the expression in (3) directly corresponds to the number of possible multiselections from α , or $\#^{\bar{\mathcal{C}}}(\alpha)$ for convenience. Therefore, $\bar{\mathcal{C}}(\alpha, -)$ is bijective on the interval $(\#^{\bar{\mathcal{C}}}(\alpha))$ and periodic with period $\#^{\bar{\mathcal{C}}}(\alpha)$.

This solves the problem with lacking symbol categories — now we can separate upper-case letters, lower-case letters, numbers, etc. into different sources and apply the elevated choice function, specifying the number of symbols from each source. However, there are two issues arising:

- The result of the elevated choice function will be something like "amwYXT28@!", which is not a bad password, but it would be nice to be able to shuffle the individual selections between each other instead of lining them up one after another.
- Despite the fact that the argument shift function makes the password selection chaotic, the function is a bijection, which means that it can be reversed. With sufficient knowledge of the algorithm, a hacker can write an inverse algorithm that retrieves the private key from the resulting password. This is a deal breaker for our function, because it defeats the purpose — you may as well have one password for everything. The way to solve this problem is to make the choice function artificially non-injective, or overlapping, in a controlled way. In such case, many different keys will produce the same password, and it will be impossible to know which one of them is the correct one. This violates the common non-collision property of hash functions, but it is necessary given the nature of the function we are developing.

We will solve both problems at the same time.

2.4 The hash function

Definition 2.7. Let α be a source configuration consisting of pairs (E_i, m_i) for $i \in (N)$. For two numbers $k_1, k_2 \in \mathbb{N}_0$, define their *hash* by the following expression:

$$\mathcal{H}(\alpha, k_1, k_2) = \mathcal{C}^{\Sigma m_i} (++\bar{\mathcal{C}}(\alpha, k_1), k_2),$$

where $++: [[A]] \rightarrow [A]$ is list concatenation. This definition can be re-written in a more readable way by defining the *shuffle function* $\mathcal{S}(E, k)$ as $\mathcal{C}^{|E|}(E, k)$ and letting the source configuration α be the varying argument:

$$\mathcal{H}(-, k_1, k_2) = \mathcal{S}(-, k_2) \circ (++) \circ \bar{\mathcal{C}}(-, k_1).$$

The hash function makes a selection from every source in the configuration, then concatenates all these selections, and finally reshuffles the resulting source. The two keys k_1 and k_2 used to make a selection will be called the *choice key* and the *shuffle key* respectively.

Note that the term "hash" is used loosely here, as it may not adhere to the formal definition of a cryptographic hash. Still, such naming is somewhat justified, given that \mathcal{H} is designed to be a uniformly distributed encryption mapping that is very hard to invert. We will now discuss the properties of $\mathcal{H}(\alpha, k_1, k_2)$ for a given source configuration α :

- **Injectivity.** \mathcal{H} is injective with respect to the choice key k_1 on the interval from zero up to

$$\#^{\bar{\mathcal{C}}}(\alpha) = \prod_{i=0}^{N-1} (n_i \mid m_i)!$$

This is because $\bar{\mathcal{C}}(\alpha, -)$ is injective on this interval, and $\mathcal{S}(-, k_2)$ is injective as well. With respect to the shuffle key k_2 , the hash function is injective on the spread of \mathcal{S} , which is

$$\text{spr}(\mathcal{S}(E, -)) = \#^{|E|}(E) = (|E| \mid |E|)! = |E|! = \left(\sum_{i=1}^{N-1} m_i \right)!$$

Therefore, the number of relevant key pairs for the hash function, denoted by $\#^{(k_1, k_2)}(\alpha)$:

$$\#^{(k_1, k_2)}(\alpha) = \#^{\mathcal{S}}(\alpha) \cdot \#^{\bar{\mathcal{C}}}(\alpha) = \prod_{i=0}^{N-1} (n_i \mid m_i)! \cdot \left(\sum_{i=1}^{N-1} m_i \right)!$$

- **Overlapping.** However, this number does not equal the number of all possible values of \mathcal{H} . When applying \mathcal{S} after $\bar{\mathcal{C}}$, we are changing the order of elements in each source twice. That is, the information about the order of these elements, stored in the output of $\bar{\mathcal{C}}$, is lost after this output is concatenated and reshuffled with \mathcal{S} . Since there are $m_i!$ ways to reorder every sub-list chosen by $\bar{\mathcal{C}}$ from E_i , the amount of lost information accounts for a total of $\prod_{i=0}^{N-1} (m_i!)$, which will be denoted by $\#^{\cap}(\alpha)$.