

Diffusion and confusion of chaotic iteration based hash functions

Zhuosheng Lin*, Christophe Guyeux†, Qianxue Wang*, and Simin Yu*

* College of Automation, Guangdong University of Technology, Guangzhou, China

Email: zhuoshenglin@163.com, wangqianxue@gdut.edu.cn, siminyu@163.com

† Femto-st Institute, University of Bourgogne Franche-Comté, Besançon, France

Email: christophe.guyeux@univ-fcomte.fr

Abstract—To guarantee the integrity and security of data transmitted through the Internet, hash functions are fundamental tools. But recent researches have shown that security flaws exist in the most widely used hash functions. So a new way to improve their security performance is urgently demanded. In this article, we propose new hash functions based on chaotic iterations, which have chaotic properties as defined by Devaney. The corresponding diffusion and confusion analyzes are provided and a comparative study between the proposed hash functions is carried out, to make their use more applicable in any security context.

Keywords—hash function, security flaws, chaotic iterations, diffusion and confusion

I. INTRODUCTION

Hash functions, as one of the key technologies in information security and cryptographic application domain, are widely used in digital signatures, file integrity checking, authentication, password protection, and so on. At the same time, the analysis of hash functions has recently made some breakthroughs. Xiaoyu Wang and her team presented new collision search attacks on SHA0 and SHA1 [1]–[3]. These research results not only shocked people, but also encouraged researchers to construct more secure hash functions.

Chaos, with its high sensitiveness to small changes and initial conditions and long-term unpredictable characteristics, has become an important branch of modern nonlinear science and applications. For instance, a lot of one-way hash functions that are based on chaotic characteristics have been recently proposed [4], [5]. However, through research, we found that most of these chaotic systems are on real domain. Due to the limited-length when realized in computer or digital devices, this will inevitably lead to finite precision effects and result in dynamical degradation of chaotic systems [6]. Such flaws will make the security performance of hash function declines.

Chaotic iterations (CIs), defined on integer domains, have been proven to achieve a real chaotic system under the definition of Devaney topological chaos [7], which solves degradation of chaotic dynamic properties fundamentally. CIs have been applied to pseudorandom number generation, information hiding, symmetric cryptography, and so on [8], [9]. In this article, we intend to construct a one-way keyed hash function with CIs. Then the diffusion and confusion are analyzed.

The remainder of this article is organized as follows. The basic recalls of CIs and hash function are given in Section II. Our CI-based hash function is proposed and reformulated

in Section III. Section IV shows its experimental evaluation. This research work ends by a conclusion section in which our article is summarized and intended future work is outlined.

II. BASIC RECALLS

This section gives some recalls on topological chaotic iterations and hash functions.

A. Chaotic iterations

Let us first define some notations that are used in the remainder of this article. \mathbb{N} is the set of natural (non-negative) numbers. The domain $\mathbb{N}^* = \{1, 2, 3, \dots\}$ is the set of positive integers and $\mathbb{B} = \{0, 1\}$. $\llbracket 1; N \rrbracket = \{1, 2, 3, \dots, N\}$. A sequence which elements belong in $\llbracket 1; N \rrbracket$ is called a strategy. The set of all strategies is denoted by \mathcal{S} . S^n denotes the n^{th} term of a sequence S , X_i stands for the i^{th} components of a vector X .

Definition 1. Let $f : \mathbb{B}^N \rightarrow \mathbb{B}^N$ be a function and $S \in \mathcal{S}$ be a strategy. The so-called chaotic iterations are defined by:

$$x^0 \in \mathbb{B}^N, \\ \forall n \in \mathbb{N}^*, \forall i \in \llbracket 1; N \rrbracket, x_i^n = \begin{cases} x_i^{n-1}, & \text{if } S^n \neq i \\ (f(x^{n-1}))_{S^n}, & \text{if } S^n = i \end{cases} \quad (1)$$

In other words, at n^{th} iteration, only the S^n -th component of vector x^n is updated.

For a given function f , let us define a function $F_f : \llbracket 1; N \rrbracket \times \mathbb{B}^N \rightarrow \mathbb{B}^N$ by:

$$F_f(k, x) = \left(x_j \cdot \delta(k, j) + (f(x))_k \cdot \overline{\delta(k, j)} \right)_{j=1,2,3,\dots,N} \quad (2)$$

where $\delta(k, j) = 0 \Leftrightarrow k = j$. Consider the phase space: $\mathcal{X} = \llbracket 1; N \rrbracket \times \mathbb{B}^N$, and the map defined on \mathcal{X} by:

$$G_f(S, E) = (\sigma(S), F_f(i(S), E)), \quad (3)$$

where σ is the shift function that removes the first term of the strategy. So the chaotic iterations defined in Equ.1 can be described by the following iterations:

$$\begin{cases} X^0 \in \mathcal{X} \\ X^{k+1} = G_f(X^k) \end{cases} \quad (4)$$

For given two points $X = (S, E), Y = (\check{S}, \check{E}) \in \mathcal{X}$, we define the distance between these two points by:

$$d(X, Y) = d_e(E, \check{E}) + d_s(S, \check{S}), \text{ where} \quad (5)$$

$$\begin{cases} d_e(E, \check{E}) = \sum_{k=1}^N \delta(E_k, \check{E}_k) \\ d_s(S, \check{S}) = \frac{9}{N} \sum_{k=1}^{\infty} \frac{\|S^k - \check{S}^k\|}{10^k} \end{cases}$$

in which $\lfloor d(X, Y) \rfloor = d_e(E, \check{E})$ is the Hamming distance between E and \check{E}^k . So $d(X, Y) - \lfloor d(X, Y) \rfloor = d_s(S, \check{S})$ measures the difference between strategies S and \check{S} . More precisely, this floating part is lower than 10^{-k} if and only if the first k terms of the two strategies are equal. Moreover, if the k^{th} digit is nonzero, then $S_k \neq \check{S}_k$.

Considering the distance between d on \mathcal{X} , it has already been proven that [10]:

- G_f is continuous.
- Iterations defined in Equ.4 are regular.
- G_f is topologically transitive.
- G_f has sensitive dependence on initial conditions.

Thus, according to the Devaney's definition [7], [11], G_f is chaotic.

B. Hash functions

Let $k \in K$ be a key in a given key space K . So a function $h_k(\cdot)$ that maps a key k and a binary bit string x to a string of a fixed length l is a Secure Keyed One-Way Hash Function (SKOWHF) [12], if it satisfies the following properties:

- Given k and x , it is easy to compute $h(k, x)$.
- Without knowledge of k , it is hard to compute $h(k, x)$.
- For any x or given (possibly many) pairs x and $h(k, x)$, it is hard to compute k .
- For a given k , it is hard to find two values x and y such that $h(k, x) = h(k, y)$, but $x \neq y$.
- Length l has to be larger than 128 bits in order to counter birthday attack.
- Key space K has to be sufficiently large in order to counter exhaustive key search.

III. CI-BASED HASH FUNCTIONS

Let us now present our hash function $H_h : K \times \mathbb{B}^* \rightarrow \mathbb{B}^N$ which is based on chaotic iterations recalled before. The key $k = \{k_1, k_2, prng_type\}$ is in key space $K = \mathbb{B}^{k_1} \times \mathbb{B}^{k_2} \times \mathbb{N}$. All the steps are described in the following paragraphs.

The first step of the algorithm is to choose the traditional hash function h that we will use in our own hash function. For our implementations, we have chosen MD5, SHA-256, and SHA-512. And the selective traditional hash function determines the length (N) of the output hash value. For MD5, $N = 128$, for SHA-256, $N = 256$, and for SHA-512, $N = 512$.

Then for the input message x , we need to transform the it into a multiple normalized N bits sequence. This pre-treatment is similar to the SHA-1 case. After that, the length of the treated sequence X is L .

In the third step, we use k_2 as a seed to generate L bits pseudorandom numbers m . In our implementation, the Pseudo-Random Number Generator can be Mersenne Twister (MT), Blum Blum Shub (B.B.S.), XORshift, or Linear Congruential Generator (LCG). This is the *prng_type* in key space K . The generated pseudorandom numbers are used to construct the strategies. As $2^n = N$, we split its sequence to be $m = S^0 S^1 \dots$, where the length of S^i is n . Then the strategy is $S = \{S^0 S^1 \dots\}$, where S^i is transformed to the decimal value.

In the fourth step, we first transform the input k_1 to binary value which length is N . Here we split X into $X = \{X^0 X^1 \dots\}$. Each X^i will be combined with k_1 using exclusive-or operation. Then we combine the result with pseudorandom numbers m using exclusive-or operation too. After that, we use this result as the input of traditional hash function h .

Lastly, to construct the digest, chaotic iteration of G_f are realized with the traditional hash function output $h(k_1, X, m)$ and strategies S as defined above. The result of these iterations is a N bits vector. It is translated into hexadecimal numbers to finally obtain the hash value.

So we define the keyed hash function $H_h : K \times \mathbb{B}^* \rightarrow \mathbb{B}^N$ by the following procedure

Algorithm 1 The proposed hash function H_h

Input:

The key, $k = (k_1, k_2, prng_type) \in K$;
The input message $x \in \mathbb{B}^*$;

Output:

Hash value H ;

- 1: Transforming x to sequence X which length is L ;
 - 2: Use PRNG to generate m which using k_2 as a seed and construct strategy $S = \{S^0 S^1 \dots\}$ with m ;
 - 3: Use standard hash function to generate hash value $H = h(k_1, X, m)$;
 - 4: **for** $i = 1 \dots \mathbf{do}$
 - 5: Chaotic iterations, to generate hash value: $H = G_f(S^i, H)$;
 - 6: **end for**
 - 7: **return** H ;
-

Thus H_h is a chaotic iterations based post-treatment on the inputted hash function. If h satisfies the collision resistance property, then it is the case too for H_h . Moreover, if h satisfies the second-preimage resistance property, then it is the case too for H_h , as proven in [8].

IV. EXPERIMENTAL EVALUATION

Before discussing diffusion and confusion, we will give some examples of hash values.

A. Hash Value

Let us now consider that the input message is the poem Ulalume (E.A.Poe), which is constituted by 104 lines and 3582 characters. The traditional hash function used here will be the

MD5. So $N = 128$. To give illustration of the keys properties, we will use this hash function H_h to generate hash values in the following cases:

- Case 1. $k_1 = 50, k_2 = 50, prng_type$ is B.B.S..
- Case 2. $k_1 = 51, k_2 = 50, prng_type$ is B.B.S..
- Case 3. $k_1 = 50, k_2 = 51, prng_type$ is B.B.S..
- Case 4. $k_1 = 50, k_2 = 50, prng_type$ is LCG.
- Case 5. $k_1 = 50, k_2 = 50, prng_type$ is MT.
- Case 6. $k_1 = 50, k_2 = 50, prng_type$ is XORshift.

The corresponding hash values in hexadecimal format are:

- Case 1. D8ED0DDD1A611C1AEDE0915BE2CA91D3.
- Case 2. 54B7B1E2C2239CF0FBC327D55CFA7BF2.
- Case 3. 8453BA95FB088DA84219F1AFCD14E9EE.
- Case 4. 663F90CB4ECD5E8AF53D2760E01491C8.
- Case 5. 01C142B339413DEF49E7A65FF43A50DF.
- Case 6. FD2B8ABC6BE956718669D92367E1680A.

From simulation results, we can see that any change in key space K seems to cause a substantial modification in the final hash value, which is coherent with the topological properties of chaos.

For a security hash function, the repartition of its hash values should be uniform. In other words, the algorithm should make full use of cryptogram space to make that the hash values are evenly distributed across the cryptogram space. The parameter we use here is the same as in Case 1. In Figure 1a, the ASCII codes are localized within a small area, whereas in Figure 1b the hexadecimal numbers of the hash values are uniformly distributed in the area of cryptogram space.

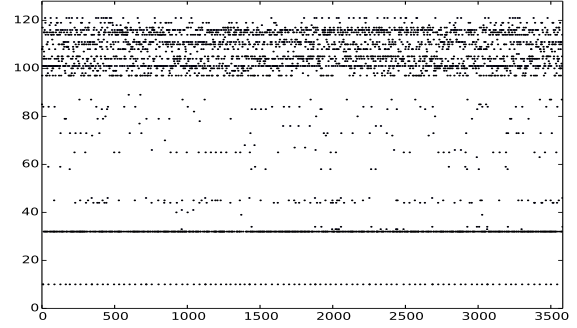
We will now test our hash function with some changes in the input message, and observe the distribution of hash values. The hash function is set with $k_1 = 50, k_2 = 50$, and $prng_type$ is B.B.S.. The hash function used here will be the MD5.

- Case 1. The input message is the poem Ulalume (E.A.Poe).
- Case 2. We replace the last point ‘.’ with a coma ‘,’.
- Case 3. In “The skies they were ashen and sober”, ‘The’ become ‘the’.
- Case 4. In “The skies they were ashen and sober”, ‘The’ become ‘Th’.
- Case 5. We add a space at the end of the poem.

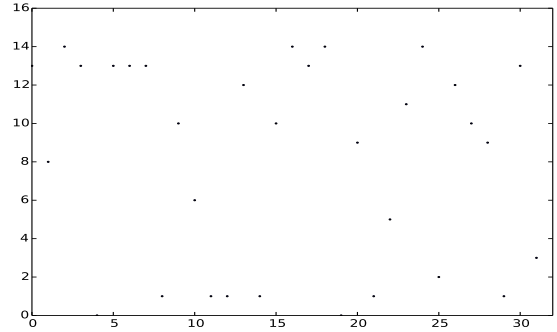
The corresponding hash values in binary format are shown in Figure 2. Through this experiment, we can check that the propose hash function is sensitive to any alteration in the input message, which will cause the modification of the hash value.

B. Diffusion and Confusion

In cryptography, diffusion and confusion are two important properties of a secure cipher that has been identified by Claude Shannon in his 1945 classified report “A Mathematical Theory of Cryptography”. Diffusion means that the redundancy of the plain text must be dispersed into the space of cryptogram space so as to hide the statistics of plain text. Confusion refers to the desire to make the statistical relationship between plain text, ciphertext, and keys as complex as possible, which makes



(a) Plain text sequence (ASCII)



(b) Hash value (Hexadecimal)

Fig. 1: Distribution of Ulalume poem

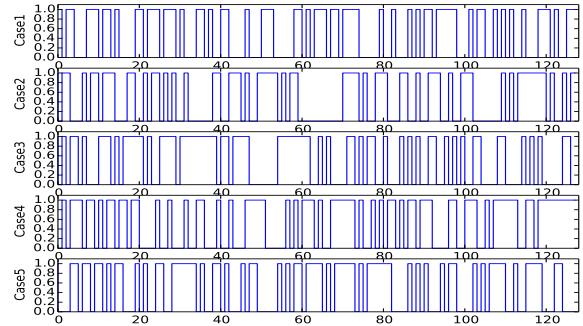


Fig. 2: 128 bit hash values in various cases

attackers difficult to get relation about keys from ciphertext. These concepts are important too in the design of robust hash functions. We now focus on the illustration of diffusion and confusion properties.

To analyze the statistic of diffusion and confusion, the following common statistics are used:

- Mean changed bit number: $\bar{B} = \frac{1}{N} \sum_{i=1}^N B_i$.
- Mean changed probability: $P = \frac{\bar{B}}{L} \times 100\%$.
- Mean square error of B: $\Delta B = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (B_i - \bar{B})^2}$.

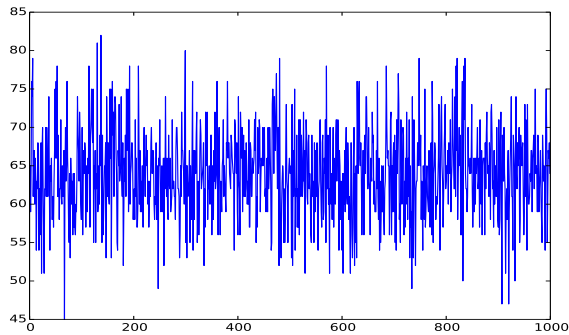


Fig. 3: Distribution of changed bit numbers B_i

- Mean square error of P:

$$\Delta P = \sqrt{\frac{1}{N-1} \sum_{i=1}^N \left(\frac{B_i}{L} - \bar{P}\right)^2} \times 100\%,$$

where N denotes the statistical times, and B_i denotes the changed bits of hash value in i^{th} test, while L denotes the bits of hash value in binary format.

We use again the poem Ulalume (E.A.Poe) as input message. Using our hash function H_h , we will get the original hash value. For this sequence, we toggle only one bit each time. Then we will obtain another hash values. Let $k_1 = 50, k_2 = 50$, and $prng_type$ is B.B.S. The hash function h used is MD5 while test times $N = 1000$. The distribution of B_i is shown in Figure 3. From the figure, we can see that a one bit change in the plain text will modify about 64 bits in the 128 bits hash value. In other words, the proposed hash function achieves desired value for such properties.

TABLE I: Statical performance of the proposed hash function

$prng_type$	$hash_type$	\bar{B}	$P(\%)$	ΔB	$\Delta P(\%)$
B.B.S	MD5	64.008	50.006	5.788	4.522
	SHA-256	128.085	50.033	7.880	3.078
	SHA-512	256.353	50.069	10.911	2.131
Mersenne Twister	MD5	63.977	49.982	5.452	4.260
	SHA-256	128.316	50.123	7.858	3.070
	SHA-512	255.534	49.909	11.691	2.283
LCG	MD5	64.355	50.277	5.795	4.528
	SHA-256	128.056	50.022	7.842	3.063
	SHA-512	256.106	50.021	11.539	2.254
XORshift	MD5	63.963	49.971	5.648	4.412
	SHA-256	127.596	49.842	8.036	3.139
	SHA-512	255.955	49.991	11.573	2.260

TABLE II: Statical performance of the standard hash function

standard hash function	\bar{B}	$P(\%)$	ΔB	$\Delta P(\%)$
MD5	63.893	49.916	5.437	4.248
SHA-256	127.746	49.901	8.405	3.283
SHA-512	256.084	50.016	11.232	2.194

The desired distribution of hash algorithm should be that small toggle in plain text causes 50% change of hash value. ΔB and ΔP show the stability of diffusion and confusion properties. The hash algorithm is more stable if these two values are close to 0. Observing Table I, both the mean changed bit number \bar{B} and the mean changed probability P

are close to the desired value. ΔB and ΔP are quite small. Both of them illustrates the diffusion and confusion of our hash function H_h and these capabilities are quite stable. From Table I, we can also know that when $prng_type = LCG$ or $prng_type = B.B.S$, all P are larger than 50%. But when $prng_type = LCG$, ΔB and ΔP are smaller. To sum up, in our proposed hash function, it is better to choose B.B.S. as pseudorandom number generator. Furthermore, compared with the performance of standard hash functions which is shown in Table II, the proposed one in some situations shows better results.

V. CONCLUSION AND FUTURE WORK

In this article, a new hash function based on chaotic iterations has been presented. We used pseudorandom number generator to construct a strategy S . Then we simulated the proposed hash function's sensitivity to keys and plain text. At last, the performance of diffusion and confusion is discussed. The experimental results show that this hash function is a secure keyed one-way hash function. Through the statical performance of the proposed hash function, we found that B.B.S is a better pseudorandom number generator to construct strategies.

In future work, we will try to apply chaotic iteration to construct pseudorandom number generators. Then we will use this kind of PRNG to construct strategies for hash functions. At the meantime, other properties induced by CIs will be explored.

REFERENCES

- [1] X. Wang, D. Feng, X. Lai, and H. Yu, "Collisions for hash functions md4, md5, haval-128 and ripemd." *IACR Cryptology ePrint Archive*, vol. 2004, p. 199, 2004.
- [2] X. Wang, Y. L. Yin, and H. Yu, "Finding collisions in the full sha-1," in *Advances in Cryptology—CRYPTO 2005*. Springer, 2005, pp. 17–36.
- [3] X. Wang, H. Yu, and Y. L. Yin, "Efficient collision search attacks on sha-0," in *Advances in Cryptology—CRYPTO 2005*. Springer, 2005, pp. 1–16.
- [4] Y. Wang, X. Liao, D. Xiao, and K.-W. Wong, "One-way hash function construction based on 2d coupled map lattices," *Information Sciences*, vol. 178, no. 5, pp. 1391–1406, 2008.
- [5] W. Guo, X. Wang, D. He, and Y. Cao, "Cryptanalysis on a parallel keyed hash function based on chaotic maps," *Physics Letters A*, vol. 373, no. 36, pp. 3201–3206, 2009.
- [6] S. Li, G. Chen, and X. Mou, "On the dynamical degradation of digital piecewise linear chaotic maps," *International Journal of Bifurcation and Chaos*, vol. 15, no. 10, pp. 3119–3151, 2005.
- [7] R. L. Devaney, L. Devaney, and L. Devaney, *An introduction to chaotic dynamical systems*. Addison-Wesley Reading, 1989, vol. 13046.
- [8] C. Guyeux, Q. Wang, X. Fang, and J. M. Bahi, "Introducing the truly chaotic finite state machines and theirs applications in security field," in *Nolta 2014, International Symposium on Nonlinear Theory and ITS Applications*, 2014.
- [9] C. Guyeux, Q. Wang, and J. M. Bahi, "Improving random number generators by chaotic iterations. application in data hiding," *Computer Science*, vol. 13, pp. V13–643 – V13–647, 2010.
- [10] C. Guyeux and J. M. Bahi, "Topological chaos and chaotic iterations application to hash functions," in *Neural Networks (IJCNN), The 2010 International Joint Conference on*. IEEE, 2010, pp. 1–7.
- [11] J. Banks, J. Brooks, G. Cairns, G. Davis, and P. Stacey, "On devaney's definition of chaos," *The American mathematical monthly*, vol. 99, no. 4, pp. 332–334, 1992.

- [12] S. Bakhtiari, R. Safavi-Naini, and J. Pieprzyk, "Keyed hash functions." in *Cryptography: Policy and Algorithms, International Conference, Brisbane, Queensland, Australia, July 3-5, 1995, Proceedings*, 1995, pp. 201–214.