# Generating short-output digest functions

L.H. Nguyen and A.W. Roscoe

Oxford University Computing Laboratory

E-mail addresses: {Long.Nguyen, Bill.Roscoe@comlab.ox.ac.uk} This paper forms the sixth chapter of my thesis, and reports my recent research result towards the first goal of my proposal: "Efficient and secure digest functions".

This paper has also been submitted to AFRICACRYPT 2010 in January 2010.

#### Abstract

This paper introduces two related methods of generating a new cryptographic primitive termed *digest* which has similarities to  $\epsilon$ -balanced and almost universal hash functions. Digest functions, however, typically have a very short output, e.g. 16-64 bits, and hence they are not required to resist collision and inversion attacks. They also have the potential to be very fast to compute relative to long-output hash functions. The first construction uses Toeplitz matrix multiplication, which is similar to a Toeplitz based universal hashing algorithm of Krawczyk, whose security requirements can be reduced to the underlying  $\epsilon$ -biased sequences of random variables. The second is based on integer multiplications which have, perhaps surprisingly, a similar structure to Toeplitz matrix multiplication. However, due to the complication of carry bits, a rigorous mathematical proof of the second construction cannot be provided. We instead exploit the short output of digest functions to carry out statistical analysis, including chi-square tests, quantilequantile plots and maximum median calculation, of digest collision and distribution test results to argue for the security of the second construction.

## 1 Introduction

In this paper, we investigate the design, constructions and security of a new cryptographic primitive termed ( $\epsilon$ -balanced) digest function, which was recently introduced in [?]. Digest functions have similarities to cryptographic hash and universal hash functions, the second of which was extensively studied in the early nineties by many researchers [?, ?, ?, ?, ?]. The main difference between digest functions and nearly all other families of hash functions is that digest functions typically have a very short output in the range between 16 and 64 bits. For this reason, it is not possible to generate digest functions. This feature however opens the way for efficiently computed digest functions arising from both their short output and the potential of parallel computation. We observe that parallelism is unfortunately not available with many current cryptographic hash functions, such as SHA-I and MD5, whose design strictly follows the sequential Merkle-Damgård construction [?].<sup>1</sup> The main goal of this

<sup>&</sup>lt;sup>1</sup>There have been a number of newly proposed cryptographic hash functions, such as MD6 [?] and SANDstorm [?] (both of which were entered into the hash competition [?] organised by NIST recently), which make use of the hash-tree structure introduced by Merkle [?] to benefit from parallel computation. We note that there are two weaknesses of a hash-tree construction: (1) the speed-up factor of parallel computation is

paper is therefore to study efficient digest constructions which are provably secure regarding their security properties, including distribution and some other imposed restrictions in a digest function, e.g. the likelihood of digest collision with respect to distinct keys are taken into consideration.

The specification of a short-output digest function arose from a new (and non-standard) family of authentication protocols [?, ?, ?, ?] which can create secure communication from human trust and human interactions to overcome the needs for passwords as in Bluetooth and security certificates as in PKI. In this new authentication technology, to complete a protocol run, devices' human owners have to manually compare a short digest of the information they want to authenticate. It is therefore desirable to design efficient and secure short-output digest functions which can be implemented in small devices and lightweight cryptographic applications. Another application of digest functions is a new family of digital signature schemes where the use of a cryptographic hash function is replaced by digest to potentially increase efficiency. A typical example of this new type of digital signature is Flexi-MAC [?] whose single certificate can be checked to various degrees, depending on the perceived security threat, the time and computing power available.

Although short-output digest functions can be computed by just truncating the first small number of bits of a cryptographic or universal hash function, as discussed in Section ??, we will point out that this strategy does not exploit the short output of digest and parallel computation to increase computational efficiency. We subsequently prove that a well-studied universal hashing algorithm based on Toeplitz matrix multiplication and  $\epsilon$ -biased sequences of random variables can generate digests with the properties we require in Definition ??.

The second main contribution of this paper is the discovery of a structural similarity between Toeplitz matrix multiplication and the "school book" algorithm for integer long multiplication. Even though these are not exactly the same due to the effect of carry bits in integer multiplications, this work suggests another method using word multiplications which can be computed efficiently on any processor. We note that carry bits in integer multiplication make it not feasible to give a mathematical proof of digest properties, we will instead exploit the short output of digest functions to facilitate our digest collision and distribution tests with statistical analysis. This approach includes computing the *p*-values in chi-square tests, quantile-quantile plots and the maximum median calculation of collision and distribution test results on word multiplication based digest algorithm. This will suggests carry bits in integer long multiplication do not introduce much bias into the digest computation regime relative to the Toeplitz based method. To the best of our knowledge, this is the first time when statistical approach has been used to analyse the security properties of hash, universal hash and digest functions. Without carry bits our word multiplication based digest construction is identical to the Toeplitz one, and indeed such a carry-less instruction for word multiplication has been introduced into the next generation of Intel processor [?].

Both our digest constructions in this paper can benefit from parallel computation should they be computed in multi-core microprocessors, such as a 48-core CPU recently announced by Intel [?].

bounded above by the depth of the hash tree, i.e.  $\log n$  where n is the length of input messages, relative to the sequential Merkle-Damgård construction; and (2) to the best of our knowledge, there might be security weaknesses due to the variable depth of the tree arising from variable length input message, and thus this structure has not been widely adopted to construct efficient hash functions for commercial purposes.

### 2 Notation

Every operator used in this paper is bitwise, which is equivalent to working over the finite field  $\mathbb{F}_2$ . Therefore, we will replace bitwise exclusive-OR and bitwise AND with addition (or summation) and multiplication in  $\mathbb{F}_2$ , respectively. We note that addition is equivalent to subtraction in  $\mathbb{F}_2$ , and so we will stick to the use of addition to simplify the notation. Also, to avoid confusion, '×' and '\*' are used to indicate finite field multiplication and respectively integer multiplication.

In this paper, we only consider non-zero messages representable by M bits, and hence the set of all messages is denoted  $X = \{1...(2^M - 1)\}$ . K, M and b are the bitlengths of digest or hash key, input message, and digest or hash output.

**Definition 1** [?] An  $\epsilon$ -balanced universal hash function is a set of  $2^K$  hash functions  $h_k()$  where  $k \in \{0...(2^K-1)\}$  such that  $h_k : X \to Y$ ; here  $X = \{1...(2^M-1)\}$  and  $Y = \{0...(2^b-1)\}$ . Moreover, for every non-zero message m and every hash output  $y \in Y$ , we have:

$$\Pr_{\{k \in R\}}[digest(k,m) = y] \le \epsilon$$

Extending the definition of an  $\epsilon$ -balanced universal hash function, we define an  $\epsilon$ -balanced digest function as follows.<sup>2</sup>

**Definition 2** [?] A  $\epsilon$ -balanced digest function: digest :  $R \times X \to Y$  where  $R = \{0...(2^K - 1)\}$ ,  $X = \{1...(2^M - 1)\}$  and  $Y = \{0...(2^b - 1)\}$  are the set of all keys, input messages and digest outputs, and moreover:

- for every  $m \in X$  and  $y \in Y$ ,  $\Pr_{\{k \in R\}}[digest(k,m) = y] \le \epsilon$
- for every  $m, m' \in X$   $(m \neq m')$  and any  $\theta \in R$ :  $\Pr_{\{k \in R\}}[digest(k, m) = digest(k + \theta, m')] \leq \epsilon$

Note that any function satisfying Definition ?? also satisfies both Definition ?? and the requirements of an  $\epsilon$ -almost universal hash function defined in [?, ?].

**Definition 3** [?, ?] Let S be a distribution of sequences of length l. Let  $(\alpha, s)$  denote the scalar (or inner) product modulo 2 of  $\alpha \in \{0, 1\}^l$  and  $s \in \{0, 1\}^l$ .

- S is said to pass the linear test  $\alpha \in \{0,1\}^l$  with bias  $\epsilon$  if  $|\Pr_{\{s \in S\}}[(\alpha, s) = 1] 1/2| \le \epsilon$ .
- S is said to be an  $\epsilon$ -biased distribution if it passes all linear tests  $\alpha \neq 0$  with bias  $\epsilon$ .

**Definition 4** A Toeplitz or diagonal-constant matrix A is a (not necessary square) matrix where each descending diagonal from left to right is constant, i.e. for all pairs of indexes (i, j):  $A_{i,j} = A_{i+1,j+1}$ .

<sup>&</sup>lt;sup>2</sup>In the second requirement of a digest function, we can replace  $+\theta$  with  $\oplus\theta$  since the bitwise exclusive-or operator is equivalent to addition over the finite field  $\mathbb{F}_2$ .

## **3** Background information

It would be a great mistake to compute digest(k,m) as some function f() of k and some similar length digest'(m), i.e. digest'() is not collision-resistant or the lengths of digest' and digest are similar, which it might be tempting to do.

$$digest(k,m) = f(k, digest'(m))$$

This is because an intruder could search for a set of different input messages all digesting to the same value digest'(m) irrespective of the value of k, i.e. any pair m and m' of these would violate the second requirement of a digest function. We therefore need to embed the key k fundamentally into the calculation of digest(k,m): it cannot be an afterthought.

Although digest output is short, it is challenging to construct digest algorithms that have both efficiency and provable security relative to the digest's requirements. Thus the majority of short-output digest constructions invented to date [?, ?, ?] make use of cryptographic hash functions, as seen in the following examples.

Pasini and Vaudenay [?], Gehrmann et al.[?], and the Bluetooth white-paper [?] suggested the following scheme:

$$digest(k,m) = trunc_b (hash(k \parallel m))$$

where  $trunc_b()$  function truncates to the leading b bits. We make two observations about this.

- 1. The definition of a cryptographic hash function does not normally specify the distribution of individual groups of bits, and so whether or not this is a good idea will very much depend on the properties we require of the hash function and the properties we want from the digest function. It follows that a specific analysis would be required for any particular standard cryptographic hash function proposed.
- 2. Computing a longhash operating on long words as in the sequential Merkle-Damgård construction can be potentially expensive, i.e. neither does the Merkle-Damgård structure permit parallelism nor exploit the short output of a digest function to increase efficiency.

Taking a different approach, Gehrmann and Nyberg [?] proposed using error-correcting codes, such as the extended Reed-Solomon codes, to construct the digest function. This has the advantage of having a coherent mathematical structure. On the other hand, to have unconditional security, the algorithm limits the bitlength of the input message to some fixed number, such as 128 or 256, which is the length of a dataword input by the algorithm. To digest any significant amount of data, the algorithm first compresses the input message into that number of bits by using a cryptographic hash which can be inefficient as discussed above. This feature also makes the scheme be potentially vulnerable to attacks should the intruder find (off line) a collision on the cryptographic hash, i.e. this is equivalent to finding a pair of distinct messages which violate the second requirement of a digest function. The reason for this weakness is because the input message is not entirely linked to the key in digest computation, as discussed at the beginning of this section.

# 4 Toeplitz matrix product construction of a digest function

We now describe a universal hash function construction using Toeplitz matrixes whose binary elements are drawn from  $\epsilon$ -biased sequences of (pseudo)random bits. This construction was

proved by Krawczyk in [?] to satisfy the requirements of an  $\epsilon$ -balanced universal hash function. We will then show the *same* algorithm can be used to generate a digest function with the properties we require.

Suppose we want to compute a *b*-bit universal hash of a (non-zero) *M*-bit message *m*, then the key *k* must be drawn from an  $\epsilon$ -biased distribution *R* of length K = M + b - 1. Using the key *k*, we can generate a Toeplitz matrix A(k) of *M* rows and *b* columns. Using matrix product, we define:

$$h_k(m) = m \times A(k) \tag{1}$$

The symbol '×' represents the product of the vector m and the matrix A(k) in  $\mathbb{F}_2$ .

**Theorem 1** [?]  $h_k(m)$  defined in Equation (??) above satisfies the definition of an  $\epsilon$ -balanced universal hash function.

Readers who are interested in the proof of this theorem can find it in the paper of Krawczyk [?] (Theorem 9, Lemma 10 and Theorem 5). The proof makes use of a known relation between Discrete Fourier Transform of matrix multiplication and  $\epsilon$ -biased distributions.

If we replace a universal hash function  $h_k(m)$  with a digest function  $digest_T(k, m)$ , where index T indicates a Toeplitz based digest construction, we have:

$$digest_T(k,m) = m \times A(k)$$
 (2)

**Theorem 2**  $digest_T(k,m)$  defined in Equation (??) above satisfies the definition of an  $\epsilon$ -balanced digest function.

**Proof** The first requirement of an  $\epsilon$ -balanced digest function is satisfied thanks to Theorem ??.

An observation on the digest construction using Toeplitz matrix multiplication is that this construction is linear in terms of both input message m and key k. The latter is true because for any k and  $\theta$  of length K = M + b - 1 bits, we always have  $A(k + \theta) = A(k) + A(\theta)$ and hence  $digest_T(k + \theta, m) = digest_T(k, m) + digest_T(\theta, m)$ .

For any pair of distinct messages (m, m') and any  $\theta$  of length K bit, a digest collision  $digest_T(k, m) = digest_T(k + \theta, m')$  is equivalent to:

$$digest_{T}(k,m) = digest_{T}(k+\theta,m')$$

$$m \times A(k) = m' \times A(k+\theta)$$

$$m \times A(k) = m' \times A(k) + m' \times A(\theta)$$

$$(m-m') \times A(k) = m' \times A(\theta)$$

$$digest_{T}(k,m-m') = m' \times A(\theta)$$

Since we have fixed m, m' and  $\theta$ , let us denote m'' = m - m' and  $\alpha = m' \times A(\theta)$ . And the second requirement of a digest function is satisfied because:

$$\Pr_{\{k \in R\}}[digest_T(k,m) = digest_T(k+\theta,m')] = \Pr_{\{k \in R\}}[digest_T(k,m'') = \alpha] \le \epsilon$$

The inequality in the above equation is true thanks to the first requirement of a digest function (or Theorem ??).

In the above Toeplitz based digest computation, the key k is assumed to have a bitlength of K = M + b - 1 which can be long if the input message is large. In practice, however, a key is normally of the size of a typical cryptographic hash function, say r = 160 or 256 bits, and hence we need to derive the required number K of bits pseudorandomly from the initial r bits.

One possible way is to use a Linear Feedback Shift Register (LFSR) as suggested by Alon et al. [?] (Proposition 1 of Section 3). In this LFSR construction, the first r/2 bits of key k will be used to select the linear structure of the underlying r/2-bit LFSR used, i.e. this is equivalent to selecting an irreducible polynomial of degree r/2. This LFSR is then seeded by the other r/2 bits of key k. This construction has been shown by Alon et al. [?] to produce an  $\epsilon$ -biased distribution on length K = M + b - 1 with  $\epsilon = \frac{K}{2^{r/2}}$ , where each sequence is generated out of r initial bits.

There are a number of other well-studied constructions of  $\epsilon$ -biased sequences due to Alon et al. [?], including Legendre symbol construction and scalar product construction, both of which can be found in their paper.

#### 4.1 Efficiency of Toeplitz matrix based digest functions

Given an  $\epsilon$ -biased sequence of length K = M + b - 1 bits,<sup>3</sup> it is straightforward to design a customised hardware implementation of  $digest_T(k, m)$ , which can also benefit from parallelism.

We first look at the complexity of a sequential implementation of this algorithm. The cost of multiplying a M-bit vector by a  $M \times b$  matrix is 2Mb bit operations, which are either AND or XOR in this case. On a w-bit microprocessor, we can expect this to take at least 2Mb/w bitwise word operations. This figure clearly shows that the cost of computing digest function is at least proportional to its output length. To simplify the calculation, we assume that w equals b, and so this is equivalent to 2M word operations. Since the latency of most logical instructions (AND, OR, NOT, and XOR) of modern microprocessors is 1 clock cycle, executing 2M bitwise word operations takes 2M arithmetic clock cycles.

As regards parallel computation, we observe that a matrix multiplication permits parallelism in two distinct ways. First, each digest output bit can be computed on its own, and hence the complexity can be improved by a factor of b. Secondly, a  $M \times b$  matrix multiplication can be split into t = M/w parallel  $w \times b$  matrices multiplications with their resulting b-bit vectors being XORed to produce the digest value. The second method is thus t = M/wtimes more efficient than the sequential computation.

# 5 Comparing Toeplitz matrix and integer long multiplication

In this section, we will show a structural similarity between the "school book" algorithm for integer long multiplication and Toeplitz matrix multiplication. This similarity, in the next section, will lead to a more efficient algorithm for digest function.

Suppose as in the Toeplitz-based method of Section ??, key k is drawn from an  $\epsilon$ -biased distribution R of length K = M + b - 1. Consider a "school book" algorithm for integer long multiplication of two numbers k and m represented in binary, i.e.  $k = (k_1 \dots k_K)$  and

<sup>&</sup>lt;sup>3</sup>The K-bit  $\epsilon$ -biased sequence can be derived (off-line) pseudorandomly from a shorter key and stored in both devices for subsequent digest computation.

	$k_1$	$k_2$	• • •	$k_{b-1}$	$k_b$	•••	$k_K$
					$m_1$	• • •	$m_M$
	$m_M \times k_1$	$m_M \times k_2$	•••	$m_M \times k_{b-1}$	$m_M \times k_b$	•••	
•••	$m_{M-1} \times k_2$	$m_{M-1} \times k_3$	• • •	$m_{M-1} \times k_b$	$m_{M-1} \times k_{b+1}$	•••	
•••	$m_{M-2} \times k_3$	$m_{M-2} \times k_4$	• • •	$m_{M-2} \times k_{b+1}$	$m_{M-2} \times k_{b+2}$	•••	
:	:	:	÷	:	:	÷	
•••	$m_1 \times k_M$	$m_1 \times k_{M+1}$	• • •	$m_1 \times k_{K-1}$	$m_1 \times k_K$		
•••	$d_1$	$d_2$	• • •	$d_{b-1}$	$d_b$	• • •	•••

 $m = (m_1 \dots m_M)$ . The following is, however, only a part of the table that is used to carry out the integer long multiplication.

We are only interested in the overlap of the expanded multiplication, where each bit of the resulting product  $(d_1 \ldots d_b)$  is influenced *directly* by every bit of the message m. In other words, each carry bit  $c_i$  which is instrumental in the computation of  $d_i$  is considered to have an *indirect* impact.

We note that the overlapping part of the expanded multiplication can be interpreted as a matrix of b columns and M rows. Moreover, the bits  $k_i$ s in each row of this matrix are shifted by 1 position to the right hand side as we move from one to the upper row. This means that if we reverse the order of these rows in the expanded multiplication, we will have the property of a Toeplitz matrix. As a result, the b-bit digest value is equivalent to:

$$digest_M(k,m) = m \times A(k) + (c_1 \dots c_b)$$
(3)

Note that the index "M" in  $digest_M(k,m)$  indicates an integer long multiplication based digest construction. Here A(k) is a Toeplitz matrix of b columns and M rows, which is constructed by the same K = M + b - 1 bits of (or derived from) key k, and  $(c_1 \dots c_b)$  is a bit vector whose elements are carry bits which are instrumental in the computation of the corresponding output or digest bits.

The above description demonstrates a structural similarity between integer long multiplication and Toeplitz matrix multiplication, however they are not equivalent due to the impact of carry bits which are accumulated from  $d_b$  all the way to  $d_1$  in an integer long multiplication. Nevertheless, it opens the way for devising new digest constructions using word multiplications which can be computed fast on any microprocessor. This will be addressed in the next section.

### 6 Word multiplication based digest functions

Although the Toeplitz matrix based digest function can be computed efficiently on customised hardware, many applications of digest functions will need to carry out this operation in software on standard and sometimes basic microprocessors. In this section, we will show that another digest function  $digest_{WM}(k,m)$  with strong structural similarities to both  $digest_T(k,m)$  and  $digest_M(k,m)$  can be calculated using standard integer multiplication for half or whole word blocks that are implemented efficiently in just about all microprocessors. This method is inspired by the relation between the "school book" algorithm for integer long multiplication and Toeplitz matrix multiplication, as demonstrated in Section ??.



Figure 1: Word multiplication model  $digest_{WM}(k, m)$ . Each parallelogram equals the expansion of a word multiplication between a b-bit key block and a b-bit message block.

Let us divide message m into b-bit blocks  $[m_1, \ldots, m_{t=\frac{M}{b}}]$ . We generate (t+1) pseudorandom b-bit blocks  $r_i$  derived from key k.<sup>4</sup> In the following equation, the index "WM" of  $digest_{WM}(k,m)$  indicates a b-bit word multiplication based digest function.

$$digest_{WM}(k,m) = \sum_{i=1}^{t} [Low(m_i * r_i) + Up(m_i * r_{i+1})]$$
 (4)

Here, \* refers to a word multiplication of two *b*-bit blocks which produces a 2*b*-bit output. The output of each word multiplication can be partitioned into the lower and upper (*b*-bit) halves by applying Low() and Up() functions respectively. Both '+' and  $\sum$  used in Equation (??) above are additions in  $\mathbb{F}_2$ .

To see why  $digest_{WM}(k, m)$  resembles  $digest_M(k, m)$ , we give a simple example when the bitlength of the input message is M = tb = 3b in Figure ??.

As was pointed out in Section ??, there may be some asymmetry due to carry bits in integer or word multiplications. However, the effect of carry bits in this construction is reduced significantly relative to an integer long multiplication of  $digest_M(k,m)$  as described in Section ??, where carry bits are accumulated from the least to the most significant bits of the product. In contrast, by partitioning an integer long multiplication into multiple (*b*-bit) word multiplications, carry bits are only accumulated locally, i.e. up to 2b - 1 bits to the left hand side of each position as opposed to M(M + b - 1) bits in an integer long multiplication of Section ??.

<sup>&</sup>lt;sup>4</sup>Any pseudorandom number generator (PRNG) which produces  $\epsilon$ -biased sequences can be used here. A possible candidate is a LFSR whose linear structure and initial seed are derived from key k as described in Section ??.

The above construction produces a b-bit digest, in practice, we might want to construct longer digests whose output is a multiple of b, i.e. xb-bit digest function. Using the same idea, it is not hard to generalise the above algorithm to multiple-word digest construction as follows.

We still divide m into b-bit blocks  $[m_1, \ldots, m_{t=\frac{M}{b}}]$ . However, we will need to generate t + x pseudorandom b-bit blocks  $r_i$  derived from key k to compute a xb-bit digest. For all  $j \in \{1...x\}$ , we then define:

$$d_j = \sum_{i=1}^{t} [\text{Low}(m_i * r_{i+j-1}) + \text{Up}(m_i * r_{i+j})]$$

And

$$digest_{WM}(k,m) = (d_1 \cdots d_x)$$

Note that the negative and biased impact of carry bits in integer multiplications can be eliminated entirely if we make use of the carry-less multiplication instruction, termed PCLMULQDQ in the next generation of Intel processors [?], that computes the carry-less multiplication of two 64-bit operands without the generation or propagation of carry bits.<sup>5</sup> This will make  $digest_{WM}(k,m)$  identical to  $digest_T(k,m)$  as defined in Equation (??) of Section ??.

#### 6.1 Efficiency of word multiplication based digest functions

Given an  $\epsilon$ -biased sequence of length K = M + b - 1, the computation of a (b = 64)-bit  $digest_{WM}(k,m)$  as defined in Equation (??) consists of 2t = 2M/b word multiplications and 2t = 2M/b bitwise XOR operations.<sup>6</sup> In the worst case, each word multiplication takes w = b clock cycles to complete, but it is always significantly faster in many modern microprocessors which have RISC pipeline instructions. For example, AMD (K8 and K10) and Intel CPUs [?] can dispatch a  $64 \times 64 = 128$ -bit MUL once every 2-4 cycles with a latency of 4-7 cycles. To simplify the calculation, we only consider non-pipeline execution, and thus the number of arithmetic clock cycles of computing a (b=64)-bit  $digest_{WM}(k,m)$  in an AMD CPU is  $\frac{2M}{64}4 + \frac{2M}{64} = 5M/32$  compared to 2M arithmetic clock cycles of computing  $digest_T(k,m)$  of Section ??.

Another observation we want to make is that the computation of  $digest_T(k, m)$  is done via many (2M) logical instructions (AND, OR and XOR) as opposed to a significantly smaller number of MUL and logical instructions in  $digest_{WM}(k,m)$ . As a consequence, the number of *control* clock cycles (e.g. loop indexing, fetching and writing) involved in  $digest_{WM}(k,m)$ is potentially reduced by a factor of  $\frac{2M}{2M/b+2M/b} = b/2$  relative to  $digest_T(k,m)$ .

The word multiplication based method also permits a high level of parallelism in multicore processors [?] thanks to its modular structure, i.e. the digest computation can be easily split into many groups of equal or different sizes of independent word multiplications.

<sup>&</sup>lt;sup>5</sup>The primary purpose of this new instruction is for computing the Galois Hash, which is the underlying computation of the Galois Counter Mode (GCM) and AES-GCM.

<sup>&</sup>lt;sup>6</sup>There is no cost of applying Low() and Up() functions to the product of 2 word operands because in a 64/32-bit processor the multiplication result of two 64/32-bit words is typically stored in two 64/32-bit registers which are the upper and lower halves of the product.

#### 7 Statistical tests of word and Toeplitz matrix multiplications

Since if digest keys are derived from a standard pseudorandom number generator (PRNG), which produces  $\epsilon$ -biased sequences of pseudorandom bits, clearly we should expect good statistical results on Toeplitz based digest construction  $digest_T(k,m)$  thanks to its mathematical structure (Theorem ??).<sup>7</sup> In this section, we concentrate on verifying the quality of word multiplication based digest function  $digest_{WM}(k,m)$  with respect to its distribution and collision properties.

In our digest collision tests, we consider N pairs of distinct messages in combination with a large set  $\mathcal{K}$  of keys. For each *i*th pair of messages  $(m_i, m'_i)$ ,  $x_i$  denotes the number of distinct keys from set  $\mathcal{K}$  of all keys which result in digest collisions. What we want to verify is the following null hypothesis  $H_0$ :

 $H_0$ : The observed numbers of collision keys of all N trials  $\{x_1, \ldots, x_N\}$  should fall into a binomial distribution:  $x_i \sim Binom(||\mathcal{K}||, \epsilon)$ , where  $\epsilon$  is the pairwise-collision probability of a digest function and  $||\mathcal{K}||$  is the cardinality of set  $\mathcal{K}$ .

The motivation behind this null hypothesis comes from the fact that in an ideal digest function the chance  $\epsilon$  that a pair of distinct messages digesting to the same value under key  $k_1$  is the same as and independent from under another randomly chosen (different) key  $k_2$ . To check the accuracy of this hypothesis, we first carry out the *chi-square test* to compute the corresponding *p*-value [?]:

$$p$$
-value =  $\Pr(\mathcal{X}_V^2 > \text{Observed } X^2)$ 

Secondly we compare the distribution of  $\{x_1 \dots x_N\}$  against an ideal binomial distribution  $Binom(||\mathcal{K}||, \epsilon)$  by plotting their quantiles against each other. This is called the *Quantile-Quantile* or Q-Q plot [?]. If the two distributions being compared are similar, then the points in the Q-Q plot will approximately lie on the line y = x, which is also drawn in each Q-Q plot to illustrate any difference between the two.

To calculate the collision probability of a digest construction, we will look at the maximum x' of the set  $\{x_1 \cdots x_N\}$ . And suppose that the null hypothesis  $H_0$  is correct then the value of x' should be very near to the maximum median x at which all N trials from the binomial distribution  $Binom(||\mathcal{K}||, \epsilon)$  are 50% likely to be less than x. This means that

$$(\Pr\{x_i \le x\})^N = \left(\sum_{t=0}^x \binom{\|\mathcal{K}\|}{t} \epsilon^t (1-\epsilon)^{\|\mathcal{K}\|-t}\right)^N \approx \left(\frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \epsilon^{-\frac{(t-\mu)^2}{2\sigma^2}} dt\right)^N = 1/2$$

Here  $\Pr\{x_i \leq x\}$  is the probability that the *i*th trial (the number of collision keys for the *i*th pair of messages) is smaller than x, and  $\mu = \|\mathcal{K}\|\epsilon$  and  $\sigma^2 = \|\mathcal{K}\|\epsilon(1-\epsilon)$ . Since all N trials are themselves independent, we need to raise  $\Pr\{x_i \leq x\}$  to the power of N. The approximation between the cumulative binomial distribution function and the cumulative normal distribution function is due to the DeMoive-Laplace limit theorem [?]. The point x, which is used to check the accuracy of x' (the observed maximum number of collision keys across N pairs of messages), can then be found by using tabulations of the probability mass functions of the normal distribution [?]. For example, when  $N = 2^{16}$ ,  $\epsilon = 2^{-10}$  and  $\|\mathcal{K}\| =$ 

<sup>&</sup>lt;sup>7</sup>Statistical test results of Toeplitz based digest functions can be found in additional supporting information section at the end of this paper.

102400, 1024000 and 10240000, then the expected values for observed x's are respectively 142, 1134 and 10424, as given in Tables ??, ??, ?? and ??.

Regarding a digest distribution test, the main difference between this and collision one is that we will fix N distinct pairs of message and digest output, and then for each *i*th pair  $(m_i, d_i)$  we denote  $y_i$  the number of keys from set  $\mathcal{K}$  satisfying the equation  $digest_{WM}(k, m_i) = d_i$ . Since for any pair  $(m_i, d_i)$ , the probability that key  $k_1$  satisfying  $digest(k_1, m_i) = d_i$  is the same as and independent from another randomly chosen key  $k_2$ , the set  $\{y_1 \dots y_N\}$  should also fall into a binomial distribution, i.e.  $y_i \sim Binom(||\mathcal{K}||, \epsilon)$ . As a result, the same methods of analysing collision tests (chi-square test, Q-Q plots and max median) described above apply to distribution tests.

In our tests, we generate 32-bit digest functions, and so it is sufficient to consider 32-bit input messages. The resulting length for digest keys that are pseudorandomly generated in our test is therefore 64 bits. Since digest keys are usually not influenced by an attacker at the point when they are invented in many applications of digest functions [?, ?, ?, ?, ?], we use the Mersenne Twister pseudorandom number generator [?] to derive large sets  $\mathcal{K}$  of 64-bit keys.

A problem in carrying out collision (and distribution) tests is that it is very expensive when we want to acquire meaningful and reliable results about collision (and distribution) on the whole of a 32-bit digest output.<sup>8</sup> For this reason, we will only consider collisions on groups of 10 adjacent bits of a digest output, which perhaps unexpectedly gives us evidence on any difference in distribution between these group of bits potentially caused by the impact of carry bits in integer multiplication. This also makes the expected probability of a pairwise 10-bit digest collision be  $\epsilon = 2^{-10}$  in all of our tests. In particular, we will look at the following three groups of bits: 10LSB, 10MIB and 10MSB which denote the 10 least, the 10 middle and respectively the 10 most significant bits of a 32-bit digest output of  $digest_{WM}(k, m)$ .

As opposed to digest keys, input messages are public and can be under the control of an attacker, we therefore consider the following four types of input messages:

- Series 1: Pseudorandom messages derived from the Mersenne Twister [?];
- Series 2: Sequential (or cluster) messages where the initial message or pair of messages is derived from the Mersenne Twister;
- Series 3: Sparse and cluster messages where the initial pair of messages in collision tests is  $(1,2^8+1)$ ; and
- Series 4: Dense and cluster messages where the initial pair of messages in collision tests is  $(2^{32} 1, 2^{32} 257)$ .

#### 7.1 Analysis of statistical test results

We have implemented our tests relative to three different numbers of keys, i.e. capitals A, B and C denote 102400, 1024000 and 10240000 which are the cardinalities of set  $\mathcal{K}$ . However, due to lack of space, we will only give *p*-values of chi-square tests and their Q-Q plots for the

<sup>&</sup>lt;sup>8</sup>For example, to have around  $1024 = 2^{10}$  32-bit digest collisions for each pair of messages, we need to consider or generate pseudorandomly  $2^{32} \times 2^{10} = 2^{42}$  distinct keys. Since there are  $2^{16}$  pairs of messages in our tests, the number of digest computations is  $2 \times 2^{16} \times 2^{42} = 2^{59}$  or  $2^{59} \times 2^{11} = 2^{70}$  bit operations, i.e. each 32-bit digest computation of a 64-bit key and 32-bit messages requires two 32-bit word multiplication adding up to  $2 \times 2^5 \times 2^5 = 2^{11}$  bit operations. This therefore goes beyond the computation capability of our computers, and it will be much more than this for longer digest functions and longer messages.

case of 10240000 keys as presented in Tables ??, ??, ??, ??, ?? and ??. We note that to ensure uniformity in our computation of *p*-values in chi-square tests for collision and distribution, we always use the same set of bins (or intervals) to count the observed occurrences in each interval.<sup>9</sup>

Since a large number of tables reporting the statistical results of our collision and distribution tests considered here do not always make it easy to see their purposes and importance, we summarise the contents of all of these tables in the following figure. Another digest construction termed  $digest_{XWM}(k,m)$  will be described at the end of this section.

Table	Digest	Digest	$\ \mathcal{K}\ $	N	Tests
	method	property			
??	$digest_{WM}()$	distribution	С	$2^{18}$	Chi-square and Q-Q plot
??	$digest_{WM}()$	distribution	A,B,C	$2^{18}$	mean, max (no of keys) and variance
??	$digest_{WM}()$	collision	С	$2^{16}$	Chi-square and Q-Q plot
??	$digest_{WM}()$	collision	A,B,C	$2^{16}$	mean, max (no of keys) and variance
??	$digest_{XWM}()$	collision	С	$2^{16}$	Chi-square and Q-Q plot
??	$digest_{XWM}()$	collision	A,B,C	$2^{16}$	mean, max (no of keys) and variance
??	$digest_{XWM}()$	distribution	С	$2^{18}$	Chi-square and Q-Q plot
??	$digest_{XWM}()$	distribution	A,B,C	$2^{18}$	mean, max (no of keys) and variance
??	$digest_T()$	distribution	С	$2^{17}$	Chi-square and Q-Q plot
??	$digest_T()$	distribution	A,B,C	$2^{17}$	mean, max (no of keys) and variance
??	$digest_T()$	collision	С	$2^{16}$	Chi-square and Q-Q plot
??	$digest_T()$	collision	A,B,C	$2^{16}$	mean, max (no of keys) and variance

It is clear from Tables ?? that the distribution property of  $digest_{WM}(k,m)$  approaches optimality (i.e. as in an ideal binomial distribution) because of the right range of the observed *p*-values and nearly all the points in the Q-Q plots approximately lie on the line y = x in Series 1-4. This observation is strengthened by results reported in Table ?? which indicates that the ratio between max and mean numbers of keys converges to 1 as the cardinality of set  $\mathcal{K}$  increases. Please note that to exploit the 10-bit digest output and to discover if any correlation exits among all distinct digest values with respect to the same message, in our distribution tests we always pair each message with each of every possible digest output from 0 to  $1023 = 2^{10} - 1$ , i.e. in Table ??,  $N = 2^{18}$  pairs are made up from  $2^8$  distinct messages and  $2^{10}$  digest values. As a result, the mean values of numbers of keys in all distribution tests equal  $||\mathcal{K}||\epsilon$  exactly as seen in Tables ??, ?? and ??. What we are interested in are the distribution properties of  $\{y_1 \dots y_N\}$  represented by *p*-value, Q-Q plot, max and variance.<sup>10</sup>

As regards collision property, Table ?? suggests when input messages are pseudorandomly generated (Series 1-2), the numbers of collision keys in  $digest_{WM}(k,m)$  fall into a binomial distribution (i.e. the null hypothesis  $H_0$  is accurate), because both *p*-values are in the right range and nearly all points in the Q-Q plots lie approximately on the line y = x. Moreover, the maximum number of observed collision keys across all N pair of distinct messages, their

<sup>&</sup>lt;sup>9</sup>The sets  $\{x_1, \ldots, x_N\}$  and  $\{y_1, \ldots, y_N\}$  are counted into the following 16 intervals: [0,9650), [9650,9700), [9700,9750), [9750,9800), [9800,9850), [9850,9900), [9900,9950), [9950,10000), [10000,10050), [10050,10100), [10100,10150), [10150,10200), [10200,10250), [10250,10300), [10300,10350), [10350,+ $\infty$ ), and therefore the degree of freedom in all chi-square tests is 15. Here  $10000 = 10240000 * 2^{-10} = N\epsilon$  is the expected value of mean of number of keys across N trials in our chi-square tests.

<sup>&</sup>lt;sup>10</sup>When each message is not paired up with each of every possible digest output, it is highly unlikely that the observed value of mean equals  $\|\mathcal{K}\|\epsilon$  exactly. This has been tested and yielded similar statistical results as the above case.

mean and variance values (Series 1-2 of Table ??) are all near to the expected values in an ideal scenarios. As in distribution tests, the ratio between the maximum number of observed collision keys and Mean converges to 1 with each additional order of magnitude in the size of the set  $\mathcal{K}$ , i.e. 10-time increase in the number of keys.

We also obtain similar results regarding 10MSB, 10 MIB and 10MSB in both distribution and collision tests of  $digest_{WM}(k, m)$  when messages are pseudorandomly generated as seen in Series 1-2 of Tables ??-??, and Series 1-4 of Table ??-??. This suggests that the digest output bits are equivalent in their quality, and thus carry bits in integer multiplication does not introduce much bias into the digest computation regime.

However, when it comes to either sparse or dense (and cluster) messages as seen in Series 3-4 of Table ??, we detect unsatisfactory behaviours. For example, there are many points in Q-Q plots moving far away (both above and below) from the line y = x, and their *p*-values are extremely small compared to Series 1-2. To confirm the negative impacts of dense and sparse messages, we look at the ratio between Max and Mean numbers of collision keys in Series 3-4 of Table ??, and realise that it is around 2 and does not seem to improve even when we increase the number of keys. We therefore conclude that the null hypothesis  $H_0$  stated above is not correct in this case.

To fix this weakness, we need to destroy known patterns within input messages. In other words, we want to randomise input messages in an unpredictable, but also deterministic, way prior to them being processed by our digest schemes. One possible and efficient way is to exclusive-or the message with the output of some function f() applying to the key k, resulting in the following improved digest construction:

$$digest_{XWM}(k,m) = digest_{WM}(k,m \oplus f(k))$$

Since k is random and fresh in each protocol session, and so is the modified version of m. In particular, when we try  $f(k) = k_1 \oplus k_2$ , where  $k_1$  and  $k_2$  are the two (32-bit) halves of key k, then we seem to be able to get around the negative impact caused by very dense or sparse messages (except a bias in the 10MIB relative to 10MSB and 10LSB of a 32-bit digest output) as reported in Series 3-4 of Tables ??-??. Devising good ways to randomise input messages is therefore a subject for future research.

Testing results regarding distribution and collision properties of both  $digest_T(k,m)$  and  $digest_{XWM}(k,m)$  can be found in the additional supporting information at the end of this paper.

### 8 Conclusions and future research for short-output functions

There are two main contributions about a new short-output digest function presented in this paper. First, we prove that Toeplitz matrix multiplication coupled with  $\epsilon$ -biased sequences of (pseudo)random variables can be used to generate digest functions with the properties we require. Secondly, a structural similarity between Toeplitz matrix and the "school book" algorithm for integer long multiplications exists as pointed out in Section ??. This similarity leads us to introduce more efficient digest computation based on word multiplications which can be calculated fast in any processors. Although there can be bias due to carry bits in an integer multiplication, we have conducted statistical and collision tests to show that no significant difference seems to exist between our two digest constructions.

It is clear that both of our digest constructions  $(digest_T(k,m) \text{ and } digest_{WM}(k,m))$ permit parallel computation in multi-core processors to increase efficiency. This is however only possible if the key k is either really long, i.e. the same order as message bitlength K = M + b - 1, or derived (off-line) from a short number of bits in advance. We note that there is a theoretical bound on the digest (or universal hash) key length which says that the message length can grow exponentially with the key length:  $K \ge \log(M/b)$  [?]. This suggests that other digest constructions which require shorter key length might exit.

# References

- [1] http://www.intel.com/pressroom/archive/releases/20091202comp\_sm.htm
- [2] Simple Pairing White Paper. See: www.bluetooth.com/NR/rdonlyres/ 0A0B3F36-D15F-4470-85A6-F2CCFA26F70F/0/SimplePairing\_WP\_V10r00.pdf
- [3] http://software.intel.com/en-us/articles/ carry-less-multiplication-and-its-usage-for-computing-the-gcm-mode/
- [4] http://csrc.nist.gov/groups/ST/hash/sha-3/index.html
- [5] http://www.sandia.gov/scada/documents/SANDstorm\_Submission\_2008\_10\_30.pdf
- [6] M. Abramowitz and I.A. Stegun. Handbook of Mathematical Functions: with Formulas, Graphs, and Mathematical Tables. ISBN 0-486-61272-4.
- [7] N. Alon, O. Goldreich, J. Hastad and R. Peralta. Simple Constructions of Almost k-wise Independent Random Variables. In Proceedings of the IEEE Symposium on Foundations of Computer Science, 1990, vol. 2, pp. 544-553.
- [8] D.J. Bernstein, H.-C. Chen, M.-S. Chen, C.-M. Cheng, C.-H. Hsiao, T. Lange, Z.-C. Lin, and B.-Y. Yang. A Billion-mulmods-per-second PC. Eurocrypt 2009, Rump Session.
- [9] J.L. Carter and M.N. Wegman. Universal Classes of Hash Functions. Journal of Computer and System Sciences, vol. 18 (1979), pp. 143-154.
- [10] I. Damgård. A Design Principle for Hash Functions. CRYPTO, LNCS vol. 435 (1989), 416-427.
- [11] C. Gehrmann, C. Mitchell and K. Nyberg. Manual Authentication for Wireless Devices. RSA Cryptobytes, vol. 7, no. 1, pp. 29-37, 2004.
- [12] T. Johansson, G.A. Kabatianskii and B. Smeets. On the relation between A-Codes and Codes correcting independent errors. Advances in Cryptology - Eurocrypt 1993, LNCS vol. 765, pp. 1-11.
- [13] H. Krawczyk. New Hash Functions For Message Authentication. Advances in Cryptology - Eurocrypt 1995, LNCS vol. 921, pp. 301-310.
- [14] S. Laur and K. Nyberg. Efficient Mutual Data Authentication Using Manually Authenticated Strings. LNCS vol. 4301, pp. 90-107, 2006.
- [15] M. Matsumoto. Mersenne twister: a 623-dimensionally equidistributed uniform pseudorandom number generator. ACM Transactions on Modeling and Computer Simulation 8: 3-1 (1998).



Table 1: Chi-square and Q-Q plots of distribution test results of 32-bit  $digest_{WM}(k,m)$ :  $N = 2^{18}$  pairs of message and digest output (m, d) over  $\|\mathcal{K}\| = 10240000$  keys. In each Q-Q plot, the vertical (Keys) and horizontal (Binom) axes denote quantiles on observed keys and respectively binomial distribution.

Test	No of	Map	Expected Mean & Max	Observed	Observed	Observed
Series	keys	range	numbers of keys,	Mean	Max	Variance
	$\ \mathcal{K}\ $		and Variance	$\mu'$	y'	$\sigma^{\prime 2}$
			$\mu = \ \mathcal{K}\ \epsilon$	$\sum_{i=1}^{N} y_i / N$		
			$\sigma^2 = \ \mathcal{K}\ \epsilon(1-\epsilon)$		max of	$\sum_{i=1}^{N} \frac{(y_i - \mu')^2}{N}$
			$\int_{-\infty}^{y} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt = 2^{-1/N}$		$\{y_1 \dots y_N\}$	
			Pseudorandom 32-bit 1	messages.		
		10LSB		100.00	148	100.45
1A	102400	10MIB	100 - 145 - 99.90	100.00	149	99.86
		10MSB		100.00	148	99.80
		10LSB		1000.00	1139	1001.43
1B	1024000	10MIB	1000 - 1144 - 999.0	1000.00	1149	998.39
		10MSB		1000.00	1146	1002.97
		10LSB		10000.00	10442	10049.08
1C	10240000	10MIB	10000 - 10445 - 9990	10000.00	10481	9941.39
		10MSB		10000.00	10401	10027.96
	Cluster 3	32-bit me	essages: the initial messa	ge is derived	pseudorando	omly.
		10LSB		100.00	148	99.63
2A	102400	10MIB	100 - 145 - 99.90	100.00	149	99.94
		10MSB		100.00	152	99.99
		10LSB		1000.00	1143	998.92
2B	1024000	10MIB	1000 - 1144 - 999.0	1000.00	1141	1002.83
		10MSB		1000.00	1149	1001.46
		10LSB		10000.00	10437	10036.47
2C	10240000	10MIB	10000 - 10445 - 9990	10000.00	10449	9999.81
		10MSB		10000.00	10419	9944.26
	Spar	se and C	luster 32-bit messages	: the initial	message is 1.	I
		10LSB		100.00	155	99.38
3A	102400	10MIB	100 - 145 - 99.90	100.00	149	99.10
		10MSB		100.00	150	99.90
		10LSB		1000.00	1147	996.47
3B	1024000	10MIB	1000 - 1144 - 999.0	1000.00	1142	991.83
		10MSB		1000.00	1139	987.38
		10LSB		10000.00	10490	10022.99
3C	10240000	10MIB	10000 - 10445 - 9990	10000.00	10443	9914.02
		10MSB		10000.00	10430	10063.45
	Dense	and Clus	ster 32-bit messages: th	he initial me	ssage is $2^{32}$ –	1.
		10LSB	<u> </u>	100.00	155	99.54
4A	102400	10MIB	100 - 145 - 99.90	100.00	148	100.31
		10MSB		100.00	152	100.08
		10LSB		1000.00	1146	1002.80
4B	1024000	10MIB	1000 - 1144 - 999.0	1000.00	1137	1001.23
		10MSB		1000.00	1143	998.89
		10LSB		10000.00	10486	9998.82
4C	10240000	10MIB	10000 - 10445 - 9990	10000.00	10477	9986.94
		10MSB		10000.00	10473	10032.88

Table 2: Comparing the mean, max (no of keys) and variance of distribution test results of 32-bit  $digest_{WM}(k,m)$  against their expected values in an ideal binomial distribution:  $N = 2^{18}$  distinct pairs of message and digest output (m, d). The set  $\{y_1 \dots y_N\}$  denotes the number of keys observed and counted for all N pairs (m, d).



Table 3: Chi-square and Q-Q plots of collision test results of 32-bit  $digest_{WM}(k,m)$ :  $N = 2^{16}$  pairs of 32-bit messages over  $\|\mathcal{K}\| = 10240000$  keys. In each Q-Q plot, the vertical (Collision) and horizontal (Binom) axes denote quantiles on observed collision keys and respectively binomial distribution.

Test	No of	Map	Expected Mean & Max	Observed	Observed	Observed
Series	keys	range	numbers of collision	Mean	Max	Variance
	$\ \mathcal{K}\ $		keys, and Variance	$\mu'$	x'	$\sigma^{\prime 2}$
			$\mu = \ \mathcal{K}\ \epsilon$	$\sum_{i=1}^{N} x_i / N$		
			$\sigma^2 = \ \mathcal{K}\ \epsilon(1-\epsilon)$		max of	$\sum_{i=1}^{N} \frac{(x_i - \mu')^2}{N}$
			$\int_{-\infty}^{x} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt = 2^{-1/N}$		$\{x_1 \dots x_N\}$	
	·		Pseudorandom 32-bit	messages.	·	
		10LSB		100.10	156	99.5
1A	102400	10MIB	100 - 142 - 99.90	99.61	162	98.1
		10MSB		99.99	140	99.8
		10LSB		999.91	1153	993.9
1B	1024000	10MIB	1000 - 1134 - 999.0	999.98	1128	994.3
		10MSB		999.97	1134	1000.5
		10LSB		10000.73	10412	10039.8
1C	10240000	10MIB	10000 - 10424 - 9990	10000.36	10433	10010.3
		10MSB		9999.99	10404	10102.2
C	luster 32-k	oit messa	<b>ges</b> : the initial pair of me	essages is der	rived pseudora	andomly.
		10LSB		100.05	153	99.0
2A	102400	10MIB	100 - 142 - 99.90	99.97	143	100.4
		10MSB		100.04	145	99.1
		10LSB		1000.07	1135	1004.7
2B	1024000	10MIB	1000 - 1134 - 999.0	1000.03	1133	1000.6
		10MSB		999.68	1146	1004.4
		10LSB		9999.74	10406	10037.5
2C	10240000	10MIB	10000 - 10424 - 9990	10000.10	10431	10026.4
		10MSB		10000.41	10415	9998.5
S	parse and	Cluster	<b>32-bit messages</b> : the ini	tial pair of r	nessages is (1	$(2^8 + 1).$
		10LSB		100.32	214	114.2
3A	102400	10MIB	100 - 142 - 99.90	100.02	195	101.8
		10MSB		100.12	224	99.99
		10LSB		1002.80	2059	2434.8
3B	1024000	10MIB	1000 - 1134 - 999.0	1000.27	1976	1355.1
		10MSB		1000.01	1998	1061.6
		10LSB		10027.32	20123	155437.4
3C	10240000	10MIB	10000 - 10424 - 9990	10003.16	19348	44738.7
		10MSB		10005.45	20047	17346.7
Dens	e and Clus	ster 32-b	it messages: the initial p	pair of messa	ges is $(2^{32} - 1)$	$1, 2^{32} - 257).$
		10LSB		100.06	202	99.8
4A	102400	10MIB	100 - 142 - 99.90	100.02	163	100.3
		10MSB		100.18	224	100.2
		10LSB		1000.29	2059	1069.7
4B	1024000	10MIB	1000 - 1134 - 999.0	1000.45	1692	1077.6
		10MSB		999.97	2003	1062.7
		10LSB		10003.03	19905	17230.5
$4\mathrm{C}$	10240000	10MIB	10000 - 10424 - 9990	10003.72	17478	19184.4
		10MSB		10005.41	20090	17413.0

Table 4: Comparing the mean, max (no of keys) and variance of collision test results of 32bit  $digest_{WM}(k,m)$  against their expected values in an ideal binomial distribution:  $N = 2^{16}$ pairs of 32-bit messages. The set  $\{x_1 \dots x_N\}$  denotes the number of collision keys observed and counted for all N trials.

- [16] R.C. Merkle. A digital signature based on a conventional encryption function. Crypto 1987, LNCS vol. 293, pp. 369 - 378.
- [17] L.H. Nguyen and A.W. Roscoe. Authenticating ad hoc networks by comparison of short digests. Information and Computation, vol. 206 (2008), pp. 250-271. Proceedings of Workshop on Foundation of Computer Security and Automated Reasoning Protocol Security Analysis, pp. 9-31, 2006.
- [18] L.H. Nguyen and A.W. Roscoe. Separating two roles of hashing in one-way message authentication. Proceedings of FCS-ARSPA-WITS 2008, pp. 195-210.
- [19] S. Pasini and S. Vaudenay. SAS-based Authenticated Key Agreement. Public Key Cryptography - PKC 2006: The 9th international workshop on theory and practice in public key cryptography, LNCS vol. 3958, pp. 395-409.
- [20] J.A. Rice Mathematical Statistics and data Analysis. ISBN 0-534-20934-3
- [21] R.L. Rivest. The MD6 Hash Function. Crypto 2008.
- [22] S. Ross. A First Course in Probability. ISBN 0-13-033851-6.
- [23] D.R. Stinson. Universal Hashing and Authentication Codes. Advances in Cryptology -Crypto 1991, LNCS vol. 576, pp. 74-85, 1992.
- [24] M.N. Wegman and J.L. Carter. New Hash Functions and Their Use in Authentication and Set Equality. Journal of Computer and System Sciences, vol. 22, pp. 265-279, 1981.

# A Additional supporting information about statistical test results of Toeplitz based and improved word multiplication based digest functions

Please see Tables ??-?? for testing results of  $digest_T(k,m)$ , and Tables ??-?? for testing results of  $digest_{XWM}(k,m)$ .



Table 5: Chi-square and Q-Q plots of collision test results of 32-bit  $digest_{XWM}(k,m)$  (input messages are XORed with two 32-bit halves of keys prior to multiplication):  $N = 2^{16}$  pairs of 32-bit messages over  $\|\mathcal{K}\| = 10240000$  distint keys. In each Q-Q plot, the vertical (Collision) and horizontal (Binom) axes denote quantiles on observed collision keys and respectively binomial distribution.

Test	No of	Map	Expected Mean & Max	Observed	Observed	Observed
Series	keys	range	numbers of collision	Mean	Max	Variance
	$\ \mathcal{K}\ $		keys, and Variance	$\mu'$	x'	$\sigma^{\prime 2}$
			$\mu = \ \mathcal{K}\ \epsilon$	$\Sigma_{i=1}^N x_i/N$		
			$\sigma^2 = \ \mathcal{K}\ \epsilon(1-\epsilon)$		max of	$\sum_{i=1}^{N} \frac{(x_i - \mu')^2}{N}$
			$\int_{-\infty}^{x} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt = 2^{-1/N}$		$\{x_1 \dots x_N\}$	
	1		Pseudorandom 32-bit	messages.	1	1
		10LSB		99.96	141	99.5
1A	102400	10MIB	100 - 142 - 99.90	100.00	147	98.1
		10MSB		100.10	143	100.1
		10LSB		1000.01	1136	997.2
1B	1024000	10MIB	1000 - 1134 - 999.0	999.89	1148	1010.0
		10MSB		1000.05	1143	996.6
		10LSB		9999.82	10389	10007.3
1C	10240000	10MIB	10000 - 10424 - 9990	10000.52	10474	10033.5
		10MSB		9999.73	10439	9996.1
C	luster 32-b	oit messa	ges: the initial pair of me	essages is der	rived pseudora	andomly.
		10LSB		99.98	150	100.1
2A	102400	10MIB	100 - 142 - 99.90	100.01	144	99.7
		10MSB		100.02	148	99.1
		10LSB		999.91	1159	997.6
2B	1024000	10MIB	1000 - 1134 - 999.0	999.98	1132	999.6
		10MSB		1000.00	1142	994.5
		10LSB		9999.89	10485	9937.5
2C	10240000	10MIB	10000 - 10424 - 9990	10000.17	10417	10037.8
		10MSB		9999.82	10430	10073.5
S	parse and	Cluster	32-bit messages: the ini	tial pair of r	messages is $(1)$	$(2^8 + 1).$
		10LSB		100.01	147	99.6
3A	102400	10MIB	100 - 142 - 99.90	100.02	143	98.4
		10MSB		100.14	147	99.2
		10LSB		999.84	1129	1002.1
3B	1024000	10MIB	1000 - 1134 - 999.0	1000.04	1359	1010.3
		10MSB		999.79	1133	1002.1
		10LSB		9999.58	10427	10058.8
3C	10240000	10MIB	10000 - 10424 - 9990	10000.69	13734	10834.1
		10MSB		10002.48	10425	9905.4
Dens	e and Clus	ster 32-b	it messages: the initial p	pair of messa	ges is $(2^{32} - 1)$	$1, 2^{32} - 257).$
		10LSB		99.99	145	98.7
4A	102400	10MIB	100 - 142 - 99.90	99.95	144	98.8
		10 MSB		100.13	145	99.4
		10LSB		999.70	1142	990.1
4B	1024000	10MIB	1000 - 1134 - 999.0	1000.01	1374	1004.7
		10MSB		999.95	1143	994.4
		10LSB		9999.48	10398	10044.7
4C	10240000	10MIB	10000 - 10424 - 9990	9999.25	13707	11166.1
		10MSB		10001.97	10503	9926.0

Table 6: Comparing the mean, max (no of keys) and variance of collision test results of 32-bit  $digest_{XWM}(k,m)$  (input messages are 21 ORed with two 32-bit halves of the key prior to multiplication) against their expected values in an ideal binomial distribution:  $N = 2^{16}$  pairs of 32-bit messages. The set  $\{x_1 \dots x_N\}$  denotes the number of collision keys observed and counted for all N trials.



Table 7: Chi-square and Q-Q plots of distribution test results of 32-bit  $digest_{XWM}(k,m)$  (input messages are XORed with two 32-bit halves of keys prior to multiplication):  $N = 2^{18}$  pairs of 32-bit messages over  $||\mathcal{K}|| = 10240000$  distint keys. In each Q-Q plot, the vertical and horizontal axes denote quantiles on observed (collision) keys and respectively binomial distribution.

Test	No of	Map	Expected Mean & Max	Observed	Observed	Observed
Series	keys	range	numbers of keys,	Mean	Max	Variance
	$\ \mathcal{K}\ $		and Variance	$\mu'$	y'	$\sigma^{\prime 2}$
			$\mu = \ \mathcal{K}\ \epsilon$	$\sum_{i=1}^{N} y_i / N$		
			$\sigma^2 = \ \mathcal{K}\ \epsilon(1-\epsilon)$		max of	$\sum_{i=1}^{N} \frac{(y_i - \mu')^2}{N}$
			$\int_{-\infty}^{y} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt = 2^{-1/N}$		$\{y_1 \dots y_N\}$	
			Pseudorandom 32-bit	messages.		
		10LSB		100.00	151	100.34
1A	102400	10MIB	100 - 145 - 99.90	100.00	151	99.76
		10MSB		100.00	158	99.71
		10LSB		1000.00	1146	995.60
1B	1024000	10MIB	1000 - 1144 - 999.0	1000.00	1160	999.70
		10MSB		1000.00	1165	1002.67
		10LSB		10000.00	10456	10031.05
1C	10240000	10MIB	10000 - 10445 - 9990	10000.00	10428	9968.10
		10MSB		10000.00	10439	9907.97
	Cluster 3	32-bit me	essages: the initial messa	ge is derived	pseudorando	omly.
		10LSB		100.00	146	99.80
2A	102400	10MIB	100 - 145 - 99.90	100.00	149	99.67
		10MSB		100.00	148	100.06
		10LSB		1000.00	1142	997.07
2B	1024000	10MIB	1000 - 1144 - 999.0	1000.00	1141	997.58
		10MSB		1000.00	1150	1000.52
		10LSB		10000.00	10425	10038.38
2C	10240000	10MIB	10000 - 10445 - 9990	10000.00	10490	9983.42
		10MSB		10000.00	10463	10009.36
	Spar	se and C	Uluster 32-bit messages	: the initial	message is 1.	<u> </u>
		10LSB		100.00	149	99.67
3A	102400	10MIB	100 - 145 - 99.90	100.00	151	100.12
		10MSB		100.00	153	99.94
		10LSB		1000.00	1149	996.76
3B	1024000	10MIB	1000 - 1144 - 999.0	1000.00	1150	999.42
		10MSB		1000.00	1137	997.70
		10LSB		10000.00	10436	9989.51
3C	10240000	10MIB	10000 - 10445 - 9990	10000.00	10440	9960.43
		10MSB		10000.00	10440	9989.53
	Dense	and Clus	ster 32-bit messages: th	ne initial me	ssage is $2^{32}$ –	1.
		10LSB		100.00	153	100.15
4A	102400	10MIB	100 - 145 - 99.90	100.00	147	99.79
		10MSB		100.00	147	99.78
		10LSB		1000.00	1144	1001.88
4B	1024000	10MIR	1000 - 1144 - 999 0	1000.00	1145	1002.80
	1021000	10MSR	1000 1111 000.0	1000.00	1150	995 98
		10LSR		10000.00	10447	9982.14
4C	10240000	10 MIR	10000 - 10445 - 9990	10000.00	10418	9987 62
	10210000	10MSR	10000 10110 0000	10000.00	10450	10011 10
	1	TOWIDD		10000.00	10100	10011.10

Table 8: Comparing the mean, max (no of keys) and variance of distribution test results of 32-bit  $digest_{XWM}(k,m)$  against their expected values in an ideal binomial distribution:  $N = 2^{18}$  distinct pairs of message and digest output (m, d). The set  $\{y_1 \dots y_N\}$  denotes the number of keys observed and counted for all N pairs (m, d).



Table 9: Chi-square and Q-Q plots of distribution test results of 32-bit  $digest_T(k, m)$ :  $N = 2^{17}$  pairs of 32-bit messages over  $\|\mathcal{K}\| = 10240000$  keys. In each Q-Q plot, the vertical and horizontal axes denote quantiles on observed keys and respectively binomial distribution.

Test	No of 64-bit	Expected Mean & Max	Observed	Observed	Observed
Series	random keys	numbers of keys,	Mean	Max	Variance
	$\ \mathcal{K}\ $	and Variance	$\mu'$	x'	$\sigma^{\prime 2}$
		$\mu = \ \mathcal{K}\ \epsilon$	$\Sigma_{i=1}^N y_i/N$		
		$\sigma^2 = \ \mathcal{K}\ \epsilon(1-\epsilon)$			$\sum_{i=1}^{N} \frac{(y_i - \mu')^2}{N}$
		$\int_{-\infty}^{x} e^{-\frac{(t-\mu)^{2}}{2\sigma^{2}}} dt = 2^{-1/N}$		maximum of	
				$\{y_1 \dots y_N\}$	
	1	1			
		Pseudorandom 32-l	oit message	es.	
1A	102400	100 - 144 - 99.90	100.00	145	100.19
1B	1024000	1000 - 1139 - 999.0	1000.00	1147	1002.38
1C	10240000	10000 - 10439 - 9990	10000.00	10436	9997.77
			•	•	
	Cluster 32-bit	t <b>messages</b> : the initial m	essage is der	ived pseudoran	domly.
2A	102400	100 - 144 - 99.90	100.00	148	99.95
2B	1024000	1000 - 1139 - 999.0	1000.00	1152	1003.82
2C	10240000	10000 - 10439 - 9990	10000.00	10427	9920.94
	Sparse an	nd Cluster 32-bit messa	<b>ages</b> : the ini	tial message is	1.
3A	102400	100 - 144 - 99.90	100.00	145	99.53
3B	1024000	1000 - 1139 - 999.0	1000.00	1130	982.61
3C	10240000	10000 - 10439 - 9990	10000.00	10445	9949.84
	Dense and	Cluster 32-bit message	$\mathbf{s}$ : the initial	l message is $2^{32}$	-1.
4A	102400	100 - 144 - 99.90	100.00	143	100.16
4B	1024000	1000 - 1139 - 999.0	1000.00	1143	998.97
4C	10240000	10000 - 10439 - 9990	10000.00	10409	9910.73

Table 10: Comparing the mean, max (no of keys) and variance of distribution test results of 32-bit  $digest_T(k,m)$  (Toeplitz method) against their expected values in an ideal binomial distribution:  $N = 2^{17}$  pairs of 32-bit messages. The set  $\{y_1 \dots y_N\}$  denotes the number of collision keys observed and counted for all N trials.



Table 11: Chi-square and Q-Q plots of collision test results of 32-bit  $digest_T(k, m)$ :  $N = 2^{16}$  pairs of 32-bit messages over  $||\mathcal{K}|| = 10240000$  keys. In each Q-Q plot, the vertical and horizontal axes denote quantiles on observed keys and respectively binomial distribution.

Test	No of 64-bit	Expected Mean & Max	Observed	Observed	Observed
Series	random keys	number of collision	Mean	Max	Variance
	$\ \mathcal{K}\ $	keys and Variance	$\mid \mu'$	x'	$\sigma'^2$
		$\mu = \ \mathcal{K}\ \epsilon$	$\sum_{i=1}^{N} x_i / N$		
		$\sigma^2 = \ \mathcal{K}\ \epsilon(1-\epsilon)$			$\sum_{i=1}^{N} \frac{(x_i - \mu')^2}{N}$
		$\int_{-\infty}^{x} \epsilon^{-\frac{(t-\mu)^2}{2\sigma^2}} dt = 2^{-1/N}$		maximum of	
				$\{x_1 \dots x_N\}$	
	1	1	1		
		Pseudorandom 32	-bit messag	ges.	
1A	102400	100 - 142 - 99.90	99.99	146	99.1
1B	1024000	1000 - 1134 - 999.0	1000.12	1153	997.6
1C	10240000	10000 - 10424 - 9990	10000.69	10425	9958.3
	1				
Cl	uster 32-bit n	nessages: the initial pair	of messages	is derived pseu	dorandomly.
2A	102400	100 - 142 - 99.90	100.04	146	99.4
2B	1024000	1000 - 1134 - 999.0	1000.09	1144	993.3
2C	10240000	10000 - 10424 - 9990	10000.31	10409	10016.5
Sp	barse and Clu	ster 32-bit messages: t	he initial pai	r of messages is	$s (1, 2^8 + 1).$
3A	102400	100 - 142 - 99.90	100.20	142	99.6
3B	1024000	1000 - 1134 - 999.0	999.85	1137	1003.7
3C	10240000	10000 - 10424 - 9990	9998.26	10444	10005.6
		·	•		
Dense	and Cluster	32-bit messages: the init	itial pair of r	messages is $(2^{32})$	$(2^2 - 1, 2^{32} - 257).$
4A	102400	100 - 142 - 99.90	100.06	142	99.7
4B	1024000	1000 - 1134 - 999.0	1000.16	1127	1005.6
4C	10240000	10000 - 10424 - 9990	9999.42	10461	9940.6

Table 12: Comparing the mean, max (no of keys) and variance of collision test results of 32-bit  $digest_T(k,m)$  (Toeplitz method) against their expected values in an ideal binomial distribution:  $N = 2^{16}$  pairs of 32-bit messages. The set  $\{x_1 \dots x_N\}$  denotes the number of collision keys observed and counted for all N trials.