

Self Encrypting Data

David Irvine, email: david.irvine@maidsafe.net, LifeStuff David

maidsafe.net limited (registered in Scotland Sc 297540)

September, 2010 *Abstract*—This paper presents a system of encryption that requires no user intervention or passwords. The resultant data item then has to be saved or stored somewhere as in all methods. The encryption here is aimed at creating cipher-text (encrypted) objects that are extremely strong and closer to perfect in terms of reversibility, as opposed to known encryption ciphers available today. This paper focuses on symmetric encryption, but does not introduce a new cipher. Instead the paper describes a method of enhancing the use of this technology to produce highly secure data and, to do so in many situations and implementations.

Index Terms—security, freedom, privacy, encryption

CONTENTS

I	Introduction	1
I-A	The Issues Addressed by this Paper . . .	1
I-B	Conventions Used	1
I-C	Symmetric Encryption	1
I-D	Cryptographically Secure Hash	1
II	Implementation	2
II-A	Overview	2
II-B	File Chunking	2
II-C	Obfuscation Step	2
II-D	Encryption Step	2
II-E	Data Map	2
II-F	Data Atlas or Recursive Data Maps . .	3
III	Conclusions	3
	References	3
	Biographies	3
	David Irvine	3

LIST OF TABLES

LIST OF ALGORITHMS

I. INTRODUCTION

ENCRYPTION has been a goal of man since before the times of the Romans, and Caesar Ciphers (simple replacement ciphers) through Enigma machine ciphers to modern day complex matrix manipulations are present in abundance in computer enhanced algorithms. This paper describes a way to use such algorithms in addition to direct encryption that clearly shows significant improvements in our use of encryption.

A. The Issues Addressed by this Paper

The issue with today's encryption of data is in just that; we encrypt data, as a whole. This reduces the potential set of possible inputs, i.e. if we are chasing somebody's bank balance, we may expect the output to be roughly the size of a bank statement, and guess what? It is! Furthermore, the security of a whole piece of data encrypted with a single algorithm depends upon just that single algorithm not getting broken. Almost all encryption ciphers appear to reduce in effectiveness as we understand the mathematics better and create more powerful computers. One may believe then, that the answer is to encrypt bits of files. However, this would require many passwords or algorithms; like putting more locks on a door to make it more secure, it gives people a headache to think that they may have to remember multiple passwords for each file. This is unlikely to be a successful manoeuvre.

B. Conventions Used

This paper does not require all the operations listed below, but lists these for example implementations, which are outlined later.

Hash = Hash function such as SHA, MD5, etc.
 Symm = Symmetrical encryption such as AES, 3DES, etc.
 PBKDF2 = Password-Based Key Derivation Function or similar

C. Symmetric Encryption

This paper will use symmetric encryption as a term to cover all algorithms (to some extent) and therefore will use a key and initialisation vector and plain-text input data. There will be no mention of MAC or similar additions to the algorithms in the hope that the reader would not attempt to implement a poorly stated or incorrect algorithm.

D. Cryptographically Secure Hash

A hash function can be thought of as a digital fingerprint. Just as a fingerprint of a person is supposed to be unique, then a digital hash function is also supposedly unique. We have all heard of two people with identical fingerprints (but perhaps have never met any!) and in the digital world it can be possible to get two pieces of data with the same hash result. This is referred to as a collision and reduces the security of the hash algorithm. The more secure the algorithm, then the likelihood of a collision is reduced. It is very similar to taking more points of reference on an actual fingerprint to reduce collisions in that area of science also. This is an area where both systems share a similarity in the increasing complexity of measurement and recording of data points of reference.

In cryptographically secure hashing, the data is analysed and a fixed length key called the hash of the data is produced. Again similarly to human fingerprinting a hash cannot reveal data just as a fingerprint cannot reveal a person (i.e. you cannot recreate the person from the print) and you cannot recreate the data from the hash.

Early hash algorithms such as MD4, MD5 and even early SHA are considered broken, in the sense that they simply allow too many collisions to occur. Hence larger descriptors (keylengths) and more efficient algorithms are almost always required.¹

II. IMPLEMENTATION

A. Overview

- 1) Get hash of whole file.
- 2) Split into several chunks of various sizes.
- 3) Take hash of each chunk.
- 4) Create another chunk by concatenating the hashes of other chunks.
- 5) XOR (or an equivalent logical operation) the chunks together as in a one time pad.
- 6) Take the hash from (3) and use the hash of chunk + 1 as the key, use the hash of chunk + 2 as the IV.
- 7) Rename each with the hash of the new content and save these hashes.

B. File Chunking

Definition 1.

$f_c \equiv \text{file content}$, $f_m \equiv \text{file metadata}$, $hf \equiv H(f_c)$

- 1) Take hf and create a structure, which we will refer to as a data map.
- 2) Take the size of the file ($f.size()$) and calculate number n of chunks²
- 3) Create chunks by taking the first n bytes of hf and randomise the chunk sizes to remove the possibility of gathering a set of similar chunks and associating these with a piece of data. This process must allow chunks to be of varying sizes in a manner not calculable without the original file (this is why the hash is used). Many algorithms can be employed at this stage.

The chunks are created with random size to ensure the set required to recreate the file is as large as the number of available chunks in any data store.

C. Obfuscation Step

In the obfuscation step, we pollute each chunk with data from other chunks.

Example 2. For c_n , create an identically-sized data chunk by repeatedly rehashing the hash of chunk $n + 2$ and appending the result, i.e. $H(c_{n+2}) + H(H(c_{n+2})) + H(H(H(c_{n+2}))) + \dots$

¹you can see where the problem exists as the logical conclusion is a key longer than any known piece of uncompress-able data, to ensure no collisions

²Number of chunks is a setting and depends on implementation, you may wish a max number of chunks, or maximum chunk size, this decision and code is left to the reader.

This is called the XOR chunk n (c_{XORn}) and is unsurprisingly XORed (\oplus) with chunk n .³

Example 3. $c_{XOR1} \oplus c_1 \equiv c_{x1}$, $c_{XOR2} \oplus c_2 \equiv c_{x2}$, etc.

D. Encryption Step

In the encryption stage, we require two separate non deterministic pieces of data, the encryption key (or password) and the Initialisation Vector (IV). To ensure all data encrypts to the same end result we determine the IV from what can be considered non deterministic data⁴, that being the hash of one of the chunks.

Definition 4. Encrypt with key and IV is shown as $Enc_{[key][IV]}(data)$ in the following example. It is assumed the key and the IV for chunk n are derived from separate portions of the hash of chunk $n + 1$.

Example 5. $Enc_{[H(c_{n+1})][H(c_{n+1})]}(c_{xn}) \equiv c_{xen}$

Therefore we have the final stage of the data represented as chunks of highly obfuscated chunks. We then take the hash of each chunk again $H(c_{xen})$ and rename each chunk with the hash of its content.

E. Data Map

In the previous sections, we described the process of self encrypting data. However, it did leave an important question unanswered. How do we reverse this process to retrieve the plain-text from the cipher-text chunks? The answer is data maps.

In the II-A steps 1, 3 & 7 we collected important data. This data alone is enough to reverse the encryption process and this is stored in a structure we refer to as a data map. This is described in the following table.

fh		<i>Data Map structure</i>
$H(c_1)$	$H(c_{xe1})$	
$H(c_2)$	$H(c_{xe2})$	
\dots	\dots	
$H(c_n)$	$H(c_{xen})$	

With this structure the names of all the chunks are in the right hand column and all passwords and IV's (which are derived from the original chunk hashes) are stored in the left hand column. The file hash in the top row identifies the data element and acts as the unique key for this file. Reversing the process is now obvious.

- 1) Retrieve the chunks listed in right hand column.
- 2) Decrypt each.
- 3) Create each XOR chunk again.
- 4) Reverse the obfuscation stage.
- 5) Concatenate the results.

This is the complete encrypt / decrypt process for each file.

³In this case we have selected XOR to represent a logical operation to obfuscate the data, this is not restrictive in any way and may be replaced by other obfuscation methods.

⁴This is an area of debate as to whether this is non deterministic data, in this case the argument is that the only way to determine the data is to have the original data in the first place, therefore there is no need to determine keys as it would be fruitless. This is somewhat of a philosophical debate and likely to be the topic of a few furled eyebrows over a few drams in a few bars for a few years to come.

F. Data Atlas or Recursive Data Maps

The data maps (dm) from multiple files can be concatenated into a new structure we call the Data Atlas (da). Therefore $dm_1 + dm_2 + \dots dm_c \equiv da$. This data atlas is itself now a large piece of data and is fed into the self-encryption process once more. This produces a single data map and more chunks. These chunks are stored somewhere and the single remaining data map is the key to all the data.

III. CONCLUSIONS

This process allows for multiple data elements to be encrypted in a very powerful fashion. Indeed there may be some debate as to whether the encryption or obfuscation stages cannot be left out (well, at least one of them). It is decided this is not a bottleneck in such a system, as data can be processed at speeds in excess of current networking capabilities in many cases. This is open to further research for differing situations though.

An important issue here is that all data is encrypted using no user information or input. This means that if the container for all the chunks is a single container then duplicate files will produce the exact same chunks and the storage system can automatically remove duplicate information. It is estimated the savings in such a system would be greater than 95%.

Also interesting is the fact that the encryption may be seen as a “step too far”; nevertheless it does indicate that any break in an encryption cipher will not reveal any data to an attacker. This is a valuable and important point.

It is hoped the research in this field will continue and measures of number of chunks versus data map size, etc. would reveal interesting scope for optimisations and improvements.

REFERENCES

- [1] David Irvine, maidsafe - a new network paradigm. david.irvine@maidsafe.net

David Irvine is a Scottish Engineer and innovator who has spent the last 12 years researching ways to make computers function in a more efficient manner.

He is an Inventor listed on more than 20 patent submissions and was Designer / Project Manager of one of the World's largest private networks (Saudi Aramco, over \$300M). He is an experienced Project Manager and has been involved in start up businesses since 1995 and has provided business consultancy to corporates and SMEs in many sectors.

He has presented technology at Google (Seattle), British Computer Society (Christmas Lecture) and many others.

He has spent many years as a lifeboat Helmsman and is a keen sailor when time permits.