

Bachelor's Thesis

Collisions in Hash-Functions

Marco Bellmann

First Reviewer:
Prof. Dr. Coja-Oghlan

Second Reviewer:
Msc. Arnab Chatterjee

Dortmund, August 2023

Contents

| | | |
|----------|-------------------------------------|-----------|
| 1 | Introduction | 2 |
| 1.1 | May Contains | 2 |
| 1.2 | Notation Helper and Ideas | 3 |
| 2 | About MD5 | 4 |
| 2.1 | Definition | 4 |
| 2.2 | On Code | 4 |
| 3 | Collision Finding | 6 |
| 3.1 | About collisions | 6 |
| 3.2 | Differential Path | 6 |
| 3.3 | Bit Conditions | 7 |
| 3.4 | About Code | 9 |
| A | Tables | 11 |
| B | Code | 12 |
| | Abbildungsverzeichnis | 13 |
| | Algorithmenverzeichnis | 15 |
| | Bibliography | 17 |
| | Erklärung | 17 |

Chapter 1

Introduction

1.1 May Contains

1. about hashfunctions
2. Definition of md5
3. my impelmentation of md5
 - (a) padding (missing in Stevens code)
 - (b) md5compress
 - (c) potential for improvement
4. collisions for md5
 - (a) Blocks 1 and 2 finding
 - (b) 00, 01, 10, 11
 - (c) MMM
5. conclusion
6. Definition of SHA1
7. why not in SHA1
8. notation helper

Notation

| Stevens | Wang | Definition |
|------------------|------------------|---|
| $h(m) = H$ | | the hash H of some message m |
| $RL(X, Y)$ | $ROTL^Y(X)$ | cyclic left shift X by $Y \pmod{31}$ |
| $RR(X, Y)$ | - | cyclic right shift X by $Y \pmod{31}$ |
| $RC(t)$ | $S(t)$ | rotation Constant of t |
| Block 1, Block 2 | Block N, Block M | pair of first blocks for collisions finding same pair but im Code |
| Block 0, Block 1 | Block N, Block M | |

Motivation

This thesis is about a deeper look on MD5. We take a closer look at the Master thesis of M. Stevens: a fast collisions finding algorithm [1]. The goal is to work out a more clear and understandable code, which is not necessarily faster, to reevaluate the code on modern systems and the difference to SHA1.

Notes

1.2 Notation Helper and Ideas

$$RR(X, Y) \coloneqq X \circledast Y$$

$$RL(X, Y) \equiv X \otimes Y$$

$$X \oplus Y \equiv X \text{ XOR } Y$$

[illegible]

Chapter 2

About MD5

2.1 Definition

MD5 stands for *Message-Digest Algorithm 5* since it generates a digest for any given text. MD5's output length is 128 Bit, represented in hexadecimal. Since the amount of possible text is close to and the amount of MD5s is limited by 2^{128} , some text may have the MD5. If two different inputs create have the same hash value, that is what we call a collision. MD5 is usually seen as 4 steps as seen in fig md5

2.2 On Code

As mention before, we can describe the MD5 hash algorithm in four steps. The implementation of the padding is the least interesting, yet there are many mistakes one can make. To archive an adequate comparison we need build simply ignore the process of padding and pretend to expect the padding to be done, even thou we do it ourselves. The processing step is were the codes vary the most. Stevens works a

1. padding
2. procssesing
3. md5 sum
4. output

Figure 2.1: MD5 algo

$$\begin{aligned}
F_t &= f(Q'_t, Q'_{t-1}, Q'_{t-2}) \\
T_t &= F_t + Q'_{t-3} + W_t \\
R_t &= RL(T'_t, RC_t) \\
Q_{t+1} &= Q_t + R_t
\end{aligned}$$

Figure 2.2: MD5 sum

lot with global variables. All in all his code is extremely optimized for performance. The MD5 sum is build up like:

Chapter 3

Collision Finding

3.1 About collisions

Hash functions have resistance:

1. First Preimage Resistance: for a hash function $h(m) = H$ the message m is hard to find.
2. Second Preimage Resistance: for a given messages m_1 it is hard to find an m_2 with $m_1 \neq m_2$ and $h(m_1) = h(m_2)$.
3. Collision Resistance: two arbitrary messages m_1 and m_2 with $\neq m_2$ and $h(m_1) = h(m_2)$ are hard to find.

3.2 Differential Path

Stevens starts with Wang's attack, which tries to find two pairs of blocks: (B_0, B'_0) and (B_1, B'_1) that $IHV = IHV'$, with the goal to create two messages M and M' , with the same hash value:

$$\begin{array}{cccccccccccccccc} IHV_0 & \xrightarrow{M_{(1)}} & \cdots & \xrightarrow{M_k} & IHV_k & \xrightarrow{B_0} & IHV_{k+1} & \xrightarrow{B_1} & IHV_{k+2} & \xrightarrow{M_{k+1}} & \cdots & \xrightarrow{M_N} & IHV_N \\ = & & & = & & \neq & & = & & & & = & \\ IHV_0 & \xrightarrow{M_{(1)}} & \cdots & \xrightarrow{M_k} & IHV_k & \xrightarrow{B_0} & IHV'_{k+1} & \xrightarrow{B_1} & IHV'_{k+2} & \xrightarrow{M_{k+1}} & \cdots & \xrightarrow{M_N} & IHV_N \end{array}$$

The idea to manipulate a block B such that $Q_1 \dots Q_{16}$ maintain their conditions and that Q_1 to some Q_k do not change at all. We try to make k as large as possible.

$$\begin{aligned}
\delta F_t &= f(Q'_t, Q'_{t-1}, Q'_{t-2}) \\
\delta T_t &= \delta F_t + \delta Q'_{t-3} + \delta W_t \\
\delta R_t &= RL(T'_t, RC_t) - RL(T_t, RC_t) \\
\delta Q_{t+1} &= \delta Q_t + \delta R_t
\end{aligned}$$

Figure 3.1: Add difference

3.3 Bit Conditions

Bit conditions describe the differential path on bits. We need the bit conditions to avoid a carry, so a manipulation in step t stays in step t and does not propagate beyond the 31st bit. We look at conditions and restrictions. The restrictions leads to conditions, which we calculate in the following. A restriction e.g. $\Delta T_2[31] = +1$ leads to conditions $Q_1[16] Q_2[16] = Q_3[15] = 0$ and $Q_2[15] = 1$. Notice, conditions are on $\Delta T_t[i]$ a state in md5-algorithm before the rotation and restrictions are on $Q_t[i]$ states of the md5-algorithm after the rotation.

We calculate the bit conditions by using the Add-Difference for two message blocks containing tow blocks $N|M$ and $N'|M'$. The XOR-Difference is useful, too.

$$\begin{aligned}
\delta X &= X' - X \pmod{32} \text{ Add-Difference} \\
\Delta X &= X' \oplus X \text{ XOR-Difference} \\
\lambda[i] &= \text{our guess for the } i\text{th bit: } X[i]
\end{aligned}$$

if $\Delta X = \lambda \Rightarrow \delta x$ can be determined

For $\lambda[i]$ we only need to consider $i < 31$, since $X[31]$ as msb always creates a add difference of 2^{31} .

We calculate a δ for each f_t , Q_t , T_t and R_t for our add difference, to calculate Q_{t+1} . Additional we need the rotation constant RC for each t . In general we begin with the f_t since we want f_t to be in a particular state. Since we want to avoid a carries in our calculation

| t | $RC(t)$ |
|-----|---------|
| 0 | 7 |
| 1 | 12 |
| 2 | 17 |
| 3 | 22 |
| 4 | 7 |
| 5 | 12 |
| 6 | 17 |
| 7 | 22 |

1. $t \in \{0, 1, 2, 3\}$:

$Q_t = 0$ since here is no influence by an message and no calculation of f , there is nothing to change:

2. $t = 4$

$\Delta T_4 = -2^{31}$, because we must not have a carry, we *lock* the last bit. Since $RL(T_4, RC_4) = RL(-2^{31}, 7) = -2^6$ and $\delta Q_4 = 0 \Rightarrow \delta Q_5 = -2^6$

Something

3.4 About Code

There are multiple approaches to find the first block, Stevens used in his code five. Four written by himself. They all have similarities, especially in the beginnings. The first *find-first-block* algorithm though, is the algorithm by Wang. His additional algorithms are "just" improvements, if the algorithm of Wang does not success. Since the similarities we only go in Wang's algorithm with a deeper view and explain some improvements eventually.

The question: "*If there are "better" algorithms by Stevens, why start with Wang's anyway?*" may come up, yet we have to clarify what better could. The algorithm of Wang is quite efficient, but also in very minimalistic. The quote I have to check is, *if the code by Wang finds a first block, it is in the most efficient way*, but it does not always find a first block were some is. So using the algorithm by Wang as a first search algorithm, may increases our performance. It is also an more easy way to implement working code, since the algorithms of Stevens tend to be more complex, but also building up on Wang's.

The first 16 Qs can be chosen arbitrary, as long as we fulfill the conditions. Stevens does this by generating really good random values. After this he alters the random values so thy fulfill conditions. We follow Stevens approach but generate "normal" random values and alter these so they fulfill the conditions Improvements for random number generation are possible.

Example for Stevens bit manipulation for Q_t with $t = 3$:

1. the val to set the zeros (bitwise-and with *0xfe87bc3f*)
2. the val to set the ones (bitwise-or with *0x017841c0*)
3. the new bit Conds
4. the old bit Conds

| | | | | | |
|--------|----------|----------|----------|----------|-------------------|
| 1. AND | 11111110 | 10000111 | 10111100 | 00111111 | <i>0xfe87bc3f</i> |
| 2. OR | 00000001 | 01111000 | 01000001 | 11000000 | <i>0x017841c0</i> |
| 3. | | .1111... | .1....01 | 11..... | |
| 4. | |0... |0... | .0..... | |

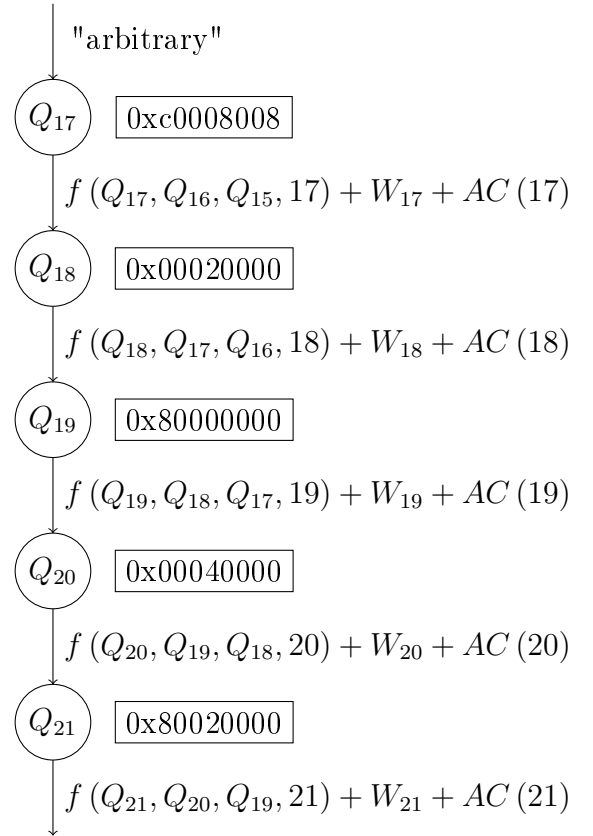
the & flips the 0 correct, the || flips the 1 correct

We now calculate the message m_t for $t \in \{0, 6, \dots, 15\}$. For this we use the reverse md5 function:

$$m_t = RR(Q_{t+1} - Q_t, RC_t) - f_t(Q_t, Q_{t-1}, Q_{t-2}) - Q_{t-3} - AC(t)$$

This function comes up by simply processing the MD5 sum 2.2 backwards. We try to pick a Q_{17} that Q_{18}, \dots, Q_{21} can be calculated with Q_{17}, \dots, Q_{15} and fulfill their conditions. We can solve this by looping about the calculation as long as the conditions are not fulfilled and then brake. Note: we may calculate the values in temp variable before pasting them into the actual Q .

For a good pick for Q_{17} we just generate a random value and force the conditions manual.



Appendix A

Tables

| |
|-----|
| +—+ |
| abc |
| +—+ |

Appendix B

Code

```
a +=a;
```

List of Figures

| | | |
|-----|--------------------|---|
| 2.1 | MD5 algo | 4 |
| 2.2 | short | 5 |
| 3.1 | short | 7 |

Algorithmenverzeichnis

Bibliography

Hiermit versichere ich, dass ich die vorliegende Arbeit selbstständig verfasst habe und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet sowie Zitate kenntlich gemacht habe.

Dortmund, den June 5, 2023

Muster Mustermann

