

5850 Graduate Project: MD5 Collision Attack Lab

Ariannah Black

Nov 17, 2023

Introduction

The aim for SEED Labs' MD5 Collision Attack Lab is to educate students about the security of one-way hash functions and develop their understanding of how collision attacks work. Upon completion of the assignment, students should be more familiar with the one-way and collision-resistance properties of hash functions, developing and launching attacks on one-way hash functions, and the structure of the MD5 hash function specifically.

In order for students to reach these objectives, SEED Labs generated four tasks related to the MD5 hash and the function `md5collgen`, which would find MD5 collisions within seconds. Prior to the start of this lab, students are required to download the lab's source code. Since the tasks must be completed within a Linux environment, I had to set up a virtual machine for my lab environment; I used `limactl`. Following this, I was able to proceed with the lab.

Reference link to MD5 Collision Attack Lab description:

https://seedsecuritylabs.org/Labs_20.04/Crypto/Crypto_MD5_Collision/

Tasks

The MD5 Collision Attack Lab is divided into four tasks; each is listed below alongside a walkthrough of my solutions and an explanation of my results.

Task 1: Generating Two Different Files with the Same MD5 Hash

Overview: Students will generate two different files which produce the same MD5 hash value and explore how `md5collgen` works, guided by discussion questions.

Methods: I first created an arbitrary text file (`prefix.txt`) to be used in later steps. Note that the length of the text file was 1043 bytes, which is not a multiple of 64. Following this, I used the below command to run `md5collgen` on my prefix file and generate two binary files (`out1.bin`, `out2.bin`).

(Command 1)

```
./md5collgen -p prefix.txt -o out1.bin out2.bin
```

For the next portion of the task, students are to view their binary file outputs in a hex editor of their choice; the lab document suggests using `bless`, which is available within SEED Labs' VM. Since I was using my own VM for the project, and since my computer's OS is Unix-based, I used the below commands to output hex dumps of the binaries to console.

(Commands 2, 3)

```
xxd out1.bin | less  
xxd out2.bin | less
```

Below is a snippet of the end of the hex dump for `out1.bin`.

(Figure 1)

```
00000370: 206c 6965 730a 446f 6f6d 6564 2073 6176    lies.Doomed sav
00000380: 696f 722c 2074 7275 6520 6d61 6b65 720a  ior, true maker.
00000390: 4875 6e74 206d 6520 646f 776e 2027 7469 Hunt me down 'ti
000003a0: 6c20 6475 736b 0a57 6f75 6c64 2079 6f75  l dusk.Would you
000003b0: 2077 6169 7420 666f 7220 6d65 2061 6e64  wait for me and
000003c0: 0a43 686f 6f73 6520 6465 6174 6820 6f72 .Choose death or
000003d0: 206c 6f73 6520 6576 6572 7974 6869 6e67  lose everything
000003e0: 3f0a 5374 7261 6e67 6c69 6e67 2079 6f75 ?.Strangling you
000003f0: 7220 6c6f 7665 0a57 6f75 6c64 2079 6f75 r love.Would you
00000400: 2077 6169 7420 666f 7220 6d65 2074 6865  wait for me the
00000410: 7265 3f00 0000 0000 0000 0000 0000 0000 re?.....
00000420: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00000430: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00000440: 251f 942c 84ca defe 5bc8 629e 8383 7529 %.....[.b...u)
00000450: e8c6 a36e 7c48 7cdc 3fea 8082 280c 9adc ...n|H|?.?...(...
00000460: a3b9 faf2 5050 6c12 eb43 4824 1108 5994 ....PP1..CH$..Y.
00000470: d4c4 4efb 6f96 74b0 a390 6667 3527 aa7f ..N.o.t..fg5'..
00000480: 5f21 dba5 3ad0 0a14 0fc9 27ce 09fa 697a _!.....'iz
00000490: a616 9d0a fc75 4019 c197 f6c4 5162 591e .....u@....QbY.
000004a0: d130 5ed6 2d99 b234 57a9 d1ef 9502 5368 .0^...4W....Sh
000004b0: dade 050c bdःa a72c 4f9b be74 38ba 75a7 .....O..t8.u.
(END)
```

I next created another arbitrary text file, `prefix64.txt`, which was of length 64 bytes.

Below is the hex dump for the first binary file generated from calling `md5collgen` on the new prefix.

(Figure 2)

```
00000000: 446f 6f6d 6564 2073 6176 696f 722c 2074 Doomed savior, t
00000010: 7275 6520 6d61 6b65 720a 4875 6e74 206d rue maker.Hunt m
00000020: 6520 646f 776e 2027 7469 6c20 6475 736b e down 'til dusk
00000030: 0a57 6f75 6c64 2079 6f75 2077 6169 7420 .Would you wait
00000040: 943f ec38 f7cc 25c2 ada9 fb5c c235 ca23 .?.8.%....\..5.#
00000050: 392d 2121 ede3 7f26 1662 8151 b4b6 c4d6 9-!!...&.b.Q....
00000060: 77e1 a801 ee27 3e5d cb21 eafa 854b 3337 w....'>]!.K37
00000070: 9f79 1d35 539c c6de 3a9e a969 ba4a e53e .y..5S....:..i.J.>
00000080: dd35 b552 3198 428b f34e 8443 75c0 59da .5.R1.B..N.Cu.Y.
00000090: 36ad d842 2326 26b3 cb1a 6e4d ec86 02ee 6..B#&...nM...
000000a0: 109c 46d7 044f 1f80 8a87 5caf d118 0bed ..F..0....\.....
000000b0: fe51 a77e 005b b614 e68c fc38 1614 59c4 .Q.~.[....8..Y.
(END)
```

Taking the two binary files sourced from `prefix64.txt`, I was able to find the differing bytes between them. The command and its output are displayed below.

(Figure 3)

```
(base) ariblack@owls Labsetup % cmp -l out1.bin out2.bin
 84  41  241
110 113 313
124 151 351
148 102 302
174  30 230
175  13  12
188  70 270
(base) ariblack@owls Labsetup %
```

Discussion (Q1): If the length of your prefix file is not multiple of 64, what is going to happen?

If the length of the prefix file is not divisible by 64, its last lines are padded in order to fill the remaining bytes. This is represented by the sequence of 0's in the hex in Figure 1, and with the sequence of ‘.‘ characters in the ASCII.

Discussion (Q2): Create a prefix file with exactly 64 bytes, and run the collision tool again, and see what happens.

If the length is already 64 bytes, then there will be no padding appended to the end of the prefix file. The hex dump thus looks more random than with the previous prefix.

Discussion (Q3): Are the (128 bytes) generated by `md5collgen` completely different for the two output files? Please identify all the bytes that are different.

The data is mostly the same, but a handful of bytes are different between the two files. From the output in Figure 3, we can see that the byte at offset 84 (byte 20) has value 41 in the first binary, and it is 241 in the second. So, at offsets 84, 110, 124, 148, 174, 175, and 188, we see different values for bytes 20, 46, 60, 84, 110, 111, and 124, respectively.

Task 2: Understanding MD5's Property

Overview: Students will demonstrate that, for MD5, the following property holds: if inputs M and N have the same hash, then adding the same suffix T to them will result in two outputs that contain the same hash value.

Methods: I first ran the `md5collgen` script on an arbitrary prefix text file and verified that the hash for both output files were the same.

(Figure 4)

```
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ ./md5collgen -p prefix.txt -o output1 output2
MD5 collision generator v1.5
by Marc Stevens (http://www.win.tue.nl/hashclash/)

Using output filenames: 'output1' and 'output2'
Using prefixfile: 'prefix.txt'
Using initial value: 126007d7ce6e252762b3494daf706d3f

Generating first block: .....
Generating second block: S01..
Running time: 9.26406 s
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ md5sum output1; md5sum output2
801e4b28441950c68c4d71f6a1249722  output1
801e4b28441950c68c4d71f6a1249722  output2
ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ ]
```

I then created an arbitrary suffix text file and concatenated it to each of the binary output files. I then verified that, again, the MD5 hash values for the files were the same.

(Figure 5)

```
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ cat output1 suffix1.txt > output1
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ cat output2 suffix1.txt > output2
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ md5sum output1; md5sum output2
db25193c92fa2efff207cc75aaea3ad0  output1
db25193c92fa2efff207cc75aaea3ad0  output2
ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ ]
```

Discussion: From this exercise, I was able to better understand the property of MD5 previously presented in the lab document. I developed a deeper understanding of how the MD5 algorithm works through practical examples.

Task 3: Generating Two Executable Files with the Same MD5 Hash

Overview: Students will generate two different executable files which produce the same MD5 hash. They will take a provided program and create two versions of it, such that the contents of their arrays are different, but the hash values of their executables are the same.

Methods: Utilizing the skeleton code provided in the lab document, I created a C program which initialized an unsigned char array of 0x41 (the hex value for ASCII character 'A') and printed its values.

(Figure 6)

```
Users > ariblack > Documents > grad > 5850 > project > Labsetup > program.c > ...
1  #include <stdio.h>
2
3
4
5
6  unsigned char xyz[200] = {
7  /* The actual contents of this array are up to you */
8  | [0 ... 199] = 0x41
9  };
10
11 int main() {
12
13     int i;
14     for (i=0; i<200; i++) {
15
16         xyz[i] = 0x41;
17
18         // print the 200-character array
19         printf("%x", xyz[i]);
20     }
21
22     printf("\n");
23 }
24
25 // compile cmd: gcc program.c -o program.out
26 |
```

I compiled the program to get its binary executable, then stored its hex dump in a text file.

(Figure 6)

```
(base) ariblack@owls Labsetup % gcc program.c -o program.out
(base) ariblack@owls Labsetup % hexdump -C -v program.out > hex
(base) ariblack@owls Labsetup %
```


(Figure 8)

```
|(base) ariblack@owls Labsetup % head -c 32832 program.out > prefix
|(base) ariblack@owls Labsetup % hexdump -C -v prefix
00000000 cf fa ed fe 07 00 00 01 03 00 00 00 02 00 00 00 |................|
00000010 11 00 00 00 a8 04 00 00 85 00 20 00 00 00 00 00 |................|
00000020 19 00 00 00 48 00 00 00 5f 5f 50 41 47 45 5a 45 |...H...PAGEZE|
00000030 52 4f 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |R0.....|
00000040 00 00 00 00 01 00 00 00 00 00 00 00 00 00 00 00 |................|
00000050 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |................|
00000060 00 00 00 00 00 00 00 00 19 00 00 00 88 01 00 00 |................|
00000070 5f 5f 54 45 58 54 00 00 00 00 00 00 00 00 00 00 |_TEXT.....|
00000080 00 00 00 00 01 00 00 00 00 40 00 00 00 00 00 00 |.....@....|
00000090 00 00 00 00 00 00 00 00 00 40 00 00 00 00 00 00 |.....@....|
000000a0 05 00 00 00 05 00 00 00 04 00 00 00 00 00 00 00 |................|
000000b0 5f 5f 74 65 78 74 00 00 00 00 00 00 00 00 00 00 |_text.....|
000000c0 5f 5f 54 45 58 54 00 00 00 00 00 00 00 00 00 00 |_TEXT.....|
000000d0 30 3f 00 00 01 00 00 00 74 00 00 00 00 00 00 00 |?.....t....|
000000e0 30 3f 00 00 04 00 00 00 00 00 00 00 00 00 00 00 |?.....|
000000f0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |................|
00000100 5f 5f 73 74 75 62 73 00 00 00 00 00 00 00 00 00 |_stubs.....|
00000110 5f 5f 54 45 58 54 00 00 00 00 00 00 00 00 00 00 |_TEXT.....|
00000120 a4 3f 00 00 01 00 00 00 06 00 00 00 00 00 00 00 |?.....|
00000130 a4 3f 00 00 01 00 00 00 00 00 00 00 00 00 00 00 |?.....|
00000140 08 04 00 80 00 00 00 00 06 00 00 00 00 00 00 00 |.....|
00000150 5f 5f 63 73 74 72 69 6e 67 00 00 00 00 00 00 00 |_cstring.....|
00000160 5f 5f 54 45 58 54 00 00 00 00 00 00 00 00 00 00 |_TEXT.....|
00000170 aa 3f 00 00 01 00 00 00 05 00 00 00 00 00 00 00 |?.....|
00000180 aa 3f 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |?.....|
00000190 02 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |.....|
000001a0 5f 5f 75 6e 77 69 6e 64 5f 69 6e 66 6f 00 00 00 |_unwind_info...|
000001b0 5f 5f 54 45 58 54 00 00 00 00 00 00 00 00 00 00 |_TEXT.....|
000001c0 b0 3f 00 00 01 00 00 00 48 00 00 00 00 00 00 00 |?.....H....|
000001d0 b0 3f 00 00 02 00 00 00 00 00 00 00 00 00 00 00 |?.....|
000001e0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |.....|
000001f0 19 00 00 00 98 00 00 00 5f 5f 44 41 54 41 5f 43 |....._DATA_C|
00000200 4f 4e 53 54 00 00 00 00 00 40 00 00 01 00 00 00 |ONST...@....|
00000210 00 40 00 00 00 00 00 00 00 40 00 00 00 00 00 00 |@.....@....|
```

The length of my suffix would be the length of the remainder of the executable, after the array (so 16472 bytes). With this information, I generated my suffix (which should also be a multiple of 64).

(Figure 9)

```
|(base) ariblack@owls Labsetup % tail -c 16472 program.out > suffix
|(base) ariblack@owls Labsetup % cat suffix
AAAAAAA PT@_printf_ _mh_execute_headermain*xyz/?~??~0??! __mh_execute_header_main_xyx_printf
|(base) ariblack@owls Labsetup %
```

To get my P and Q values, I ran the `md5collgen` script on the prefix and extracted 128 bits (`p`, `q`) from each of the resulting binaries (`outp`, `outq`).

(Figure 10)

```
(base) ariblack@owls Labsetup % tail -c 128 outp > p
(base) ariblack@owls Labsetup % tail -c 128 outq > q
(base) ariblack@owls Labsetup % cat p; cat q
,?r|x?Z4v?E?/?zBbCI?????????I???s97?K%?|??Z??6??~5?K(t???'@??O`[?詳0_*`?x(???
?:???
H???)j??},?r|x?Z4v?E?/_?zBbCI?????????I???s97tL%?|?????6??~5?K(t???'@??O`[?詳0_*`?x(???
?:???
H???)j??})
```


Task 4: Making the Two Programs Behave Differently

Overview: Students will generate two different programs which produce the same MD5 hash. Execute an attack demonstrating that it is possible to obtain a certificate (hash value) which applies to two different sets of instructions.

Methods: I used the provided skeleton pseudo code as a base to develop my main C program, shown below. My program initializes two unsigned char arrays, both filled with the value 0x41 (ASCII character 'A'). It then includes a check: if the two arrays are the same, then it will execute a set of benign instructions; otherwise, it will execute its malicious instructions.

(Figure 14)

```
Users > ariblack > Documents > grad > 5850 > project > Labsetup > benign.c > ...
1  #include <stdio.h>
2
3
4  unsigned char x[200] = {
5  | [0 ... 199] = 0x41
6  };
7  unsigned char y[200] = {
8  | [0 ... 199] = 0x41
9  };
10
11 int main() {
12
13     for (unsigned int i = 0; i < 200; i++) {
14         // check if arrays are same
15
16         if (x[i] != y[i]) {
17             // execute malicious code
18             printf("hehe <3\n");
19             break;
20         }
21         if (i == 199) {
22             // execute benign code
23             printf("not malicious!\n");
24         }
25     }
26
27
28     return 0;
29 }
30 // compile cmd: gcc benign.c -o benign.out
```

I compiled the code to get its binary executable, then stored its hex dump in another text file, shown below. From there, I found the locations of the two arrays. The first began and ended at the same locations as the array from the previous task (0x8000/32768 and 0x80c0/32960, respectively). The second began and ended at 0x80d0/32976 and 0x8190/33168.

To generate the p and q values, I called the `md5collgen` script on the prefix, verified that their hashes were the same, and extracted 128 bytes from both p and q.

(Figure 17)

```
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ ./md5collgen -p prefix -o p q
MD5 collision generator v1.5
by Marc Stevens (http://www.win.tue.nl/hashclash/)

Using output filenames: 'p' and 'q'
Using prefixfile: 'prefix'
Using initial value: c0702d127728014e0fae9ef70e6774eb

Generating first block: ..
Generating second block: $00.....
Running time: 4.32342 s
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ md5sum p q
cd474e29219591d451322bfadd508e65 p
cd474e29219591d451322bfadd508e65 q
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ tail -c 128 p > p2
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ tail -c 128 q > q2
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ cat p2; cat q2
?
???o?;cHn??\??;C?Q??F
????*??o?D??????`Y&Ó????j?ikH1??M0??3??7{?I?
?"_??Qr}2?_??wÈ???C,??c7t??v??#LC?
???o?;cHn??\?L;C?Q??F
????*??o?D?_?%
?2`Y&Ó?????ikH1??M0??3??7{?I??"_??Qr}2?_??wÈ???1C,??c7t??v #LC
ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labset
up$ ]
```

I was then able to concatenate the parts to create my executables. My benign program had the structure `prefix || p || middle || p || suffix`, so both arrays within the program were the same, leading its control flow to execute the benign instructions. The malicious program had a slightly different structure, `prefix || q || middle || p || suffix`, where the one array is modified, leading the control flow to execute the malicious instructions. After adding execution permissions for both programs, I executed both as shown below (where `b` is the name of the benign program, and `m` is the malicious one) and verified that their hashes remained the same.

(Figure 18)

```
((base) ariblack@owls Labsetup % cat prefix p2 tmp2 p2 suffix > b
((base) ariblack@owls Labsetup % ./b
not malicious!
((base) ariblack@owls Labsetup % ./m
hehe <3
((base) ariblack@owls Labsetup % ]
```

(Figure 19)

```
[ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ md5sum m b
8cf42cb8cfe55be49bd736d378cc5578  m
8cf42cb8cfe55be49bd736d378cc5578  b
ariblack@lima-default:/Users/ariblack/Documents/grad/5850/project/Labsetup$ ]
```

The lengths of the files I used in my program concatenation were (in bytes) 32832, 128, 80, 128, and 16288 for the prefix, p (labeled p2 below), middle (labeled tmp2 below), q (labeled q2 below), and suffix splice sections, respectively.

(Figure 20)

```
|(base) ariblack@owls Labsetup % ll
total 7224
-rwxr-xr-x  1 ariblack  staff   49456 Nov  9 01:01 b
-rw-r--r--  1 ariblack  staff    518 Nov  8 15:47 benign.c
-rwxr-xr-x  1 ariblack  staff   49456 Nov  8 14:59 benign.out
-rwxr-xr-x  1 ariblack  staff   49456 Nov  9 01:01 m
-rwxr-xr-x  1 ariblack  staff   49456 Nov  8 15:04 malicious.out
-rwxrwxr-x@ 1 ariblack  staff  3338360 Mar 23  2021 md5collgen
-rw-r--r--  1 ariblack  staff     80 Nov  9 00:44 middle
drwxr-xr-x  5 ariblack  staff    160 Nov  9 01:04 misc
-rw-r--r--  1 ariblack  staff    6822 Nov  9 01:02 notes.txt
-rw-r--r--  1 ariblack  staff   32960 Nov  8 15:14 p
-rw-r--r--  1 ariblack  staff    128 Nov  8 15:15 p2
-rw-r--r--  1 ariblack  staff   32832 Nov  9 01:01 prefix
-rw-r--r--  1 ariblack  staff   32960 Nov  8 15:14 q
-rw-r--r--  1 ariblack  staff    128 Nov  8 15:15 q2
-rw-r--r--  1 ariblack  staff   16288 Nov  9 01:01 suffix
-rw-r--r--  1 ariblack  staff   16496 Nov  9 01:01 tmp
-rw-r--r--  1 ariblack  staff     80 Nov  9 01:01 tmp2
|(base) ariblack@owls Labsetup % ]
```

Discussion: This exercise was the most challenging task within the lab for me as, although I had a good general idea of how I planned to splice the different pieces together (very similarly to in Task 3), I struggled with finding the precise number of bytes that each piece should hold. I spent a good portion of time experimenting with different values for the `head` and `tail` commands, saving the resulting hex dump to a file, and observing the architecture of the output in order to see precisely why I was receiving the errors I found; eventually, I realized that my suffix and middle splices were off by 1 byte (when fixed, my scripts worked as expected).

Discussion

Overall, I believe that this lab was successful in its goal to teach me about one-way hash function vulnerabilities and about collision attacks, specifically using MD5. While I previously had a good general understanding of how collisions work and their impact on security, stepping through each of the four tasks above allowed me to understand the security risks at a more sophisticated level.

Upon completion of the lab, I can conclude that the MD5 hash function's security is severely compromised, and that it generally should no longer be used in favor of other, more secure hash functions.