

# MD5 Collision Attack Writeup

Liam Ryan

## Contents

## Contents

|  |    |
|--|----|
| Contents .....   | 1  |
| Introduction .....   | 1  |
| The objective of this lab is to gain insight into the importance of the collision-resistance property of a hash. For a one-way hash function to be considered secure there are several properties it must satisfy. The first is the one-way property; consider $hash(M) = h$ ..... | 1  |
| Environment .....  | 1  |
| Performed inside an Ubuntu VM (20.04) using the md5collgen program written by Marc Stevens. ....   | 1  |
| Task 1: Generating Two Different Files with the Same MD5 Hash .....  | 2  |
| Task 2: Further Understanding MD5's Property .....   | 6  |
| Task 3: Generating Two Executable Files with the Same MD5 Hash .....   | 7  |
| Task 4: Making the Two Programs Behave Differently .....   | 11 |
| Citations .....  | 16 |

## Introduction

The objective of this lab is to gain insight into the importance of the collision-resistance property of a hash. For a one-way hash function to be considered secure there are several properties it must satisfy. The first is the one-way property; consider  $hash(M) = h$ . Given  $h$  it must be infeasible to find the input  $M$ .

The second one and the focus of this lab is the collision-resistance property. It must be computationally infeasible (cost more to do than the information is worth) to find the hash of two different inputs equal. A collision occurs when this property is violated. i.e.,

$$Hash(M_1) == Hash(M_2)$$

In this Lab we will take advantage of this collision property (or lack thereof) to create two distinct programs that produce the same MD5 checksum (hash).

## Environment

Performed inside an Ubuntu VM (20.04) using the md5collgen program written by Marc Stevens.

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

To start this lab the md5collgen program and a sample text file are required. We want to ensure the program works and can generate the same hash for two different text files before we proceed with the actual attack.

```
liam@ubuntu:~/Lab1$ touch lab1md5.txt
liam@ubuntu:~/Lab1$ ls
lab1md5.txt  md5collgen
liam@ubuntu:~/Lab1$ vi lab1md5.txt
liam@ubuntu:~/Lab1$
```

```
i am the contents of the file.
```

I then ran the program using the previously created text file.

```
liam@ubuntu:~/Lab1$ ./md5collgen -p lab1md5.txt -o out1.bin out2.bin
MD5 collision generator v1.5
by Marc Stevens (http://www.win.tue.nl/hashclash/)

Using output filenames: 'out1.bin' and 'out2.bin'
Using prefixfile: 'lab1md5.txt'
Using initial value: fb7b32aa3ca2612bb54b3900568725ec

Generating first block: .....
Generating second block: W.....
Running time: 4.67827 s
liam@ubuntu:~/Lab1$
```

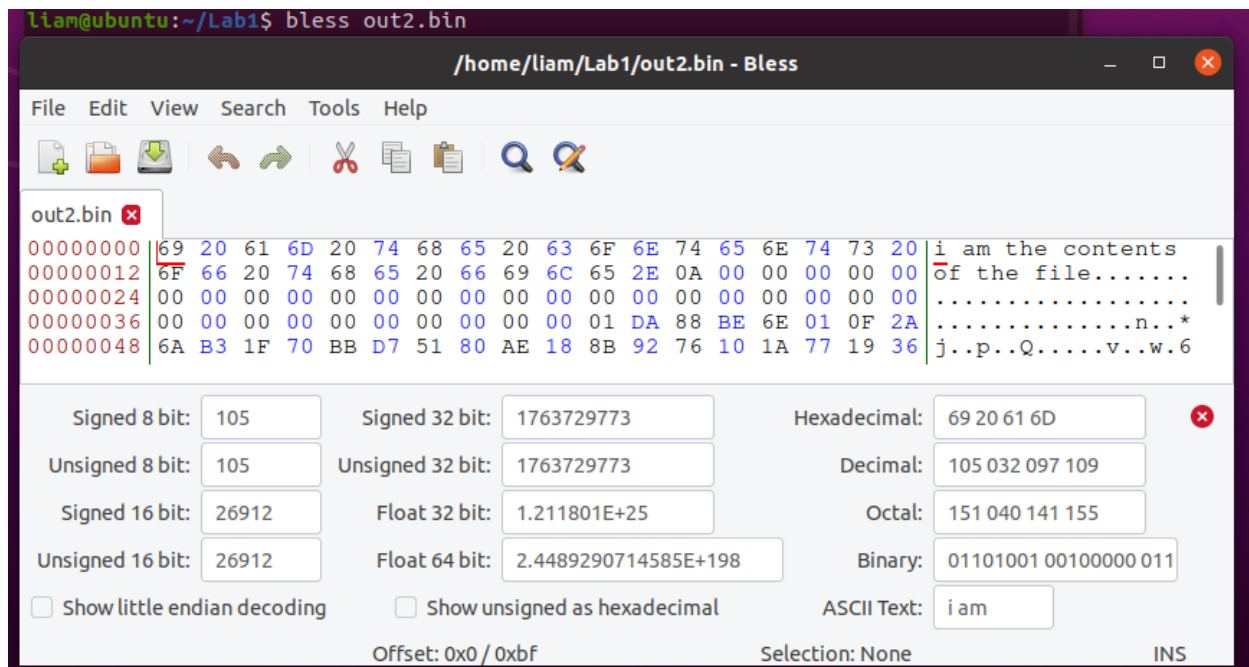
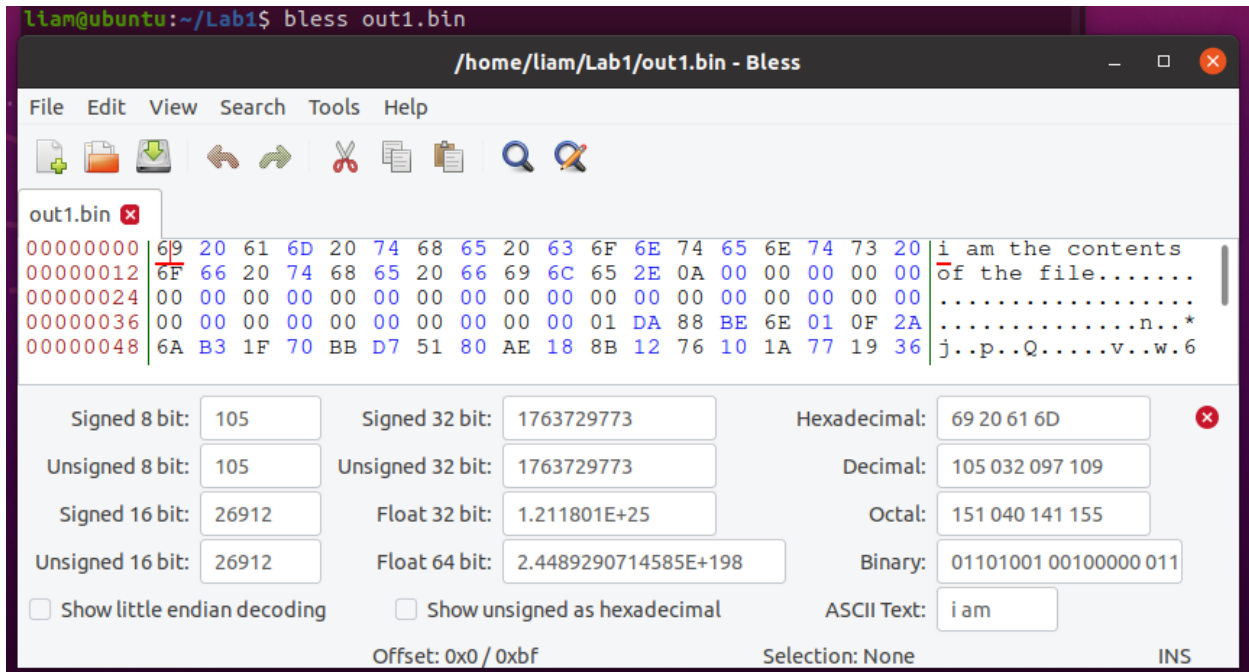
Using *diff* on the two output files we can determine if they are different, it determined that they are.

```
liam@ubuntu:~/Lab1$ diff out1.bin out2.bin
Binary files out1.bin and out2.bin differ
```

When performing a checksum on the two output files, despite them having an alternate 128-byte ending, they produce the exact same MD5 hash, meaning a collision has occurred.

```
liam@ubuntu:~/Lab1$ md5sum out1.bin
cbaacb8298be9dd80e128e5be0ede11a  out1.bin
liam@ubuntu:~/Lab1$ md5sum out2.bin
cbaacb8298be9dd80e128e5be0ede11a  out2.bin
liam@ubuntu:~/Lab1$
```

Using *bless* (a hex editor), I can dump the hex for both output files. The prefix appears on the right in plain text followed by 0 padding, and then what I assume is the appended P and some other data. The results of both dumps appear to be identical at first glance. But if you read through the hex, you can see the difference.



## If the length of your prefix file is not multiple of 64, what is going to happen?

The remaining space gets padded with zeros; this is because MD5 encrypts 64-byte blocks.

## Create a prefix file with exactly 64 bytes, and run the collision tool again, and see what happens.

Recreate the file with a prefix size of 64, and then ensure its size is 64.

```
liam@ubuntu:~/Lab1$ ls -l lab1md5.txt
-rw-rw-r-- 1 liam liam 64 Jan 12 12:57 lab1md5.txt
```

Re-run the MD5collgen program.

```
liam@ubuntu:~/Lab1$ ./md5collgen -p lab1md5.txt -o out1.bin out2.bin
MD5 collision generator v1.5
by Marc Stevens (http://www.win.tue.nl/hashclash/)

Using output filenames: 'out1.bin' and 'out2.bin'
Using prefixfile: 'lab1md5.txt'
Using initial value: b2291b48842948151ea655bc673edbf0

Generating first block: .....
Generating second block: S11...
Running time: 14.6277 s
```

I dump the hex again and can see that there is no zero padding.

| out1.bin  | out2.bin             |
|---|----------------------|
| 00000000 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41    | AAAAAAAAAAAAAAAAAAAA |
| 00000012 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41    | AAAAAAAAAAAAAAAAAAAA |
| 00000024 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41    | AAAAAAAAAAAAAAAAAAAA |
| 00000036 41 41 41 41 41 41 41 0A 0A EF 03 79 C5 90 49 25 41 | AAAAAAAA...y..I%A    |

| out1.bin   | out2.bin             |
|--|----------------------|
| 00000000 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41       | AAAAAAAAAAAAAAAAAAAA |
| 00000012 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41       | AAAAAAAAAAAAAAAAAAAA |
| 00000024 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41       | AAAAAAAAAAAAAAAAAAAA |
| 00000036 41 41 41 41 41 41 41 0A 0A EF 03 79 C5 90 49 25 41    | AAAAAAAA...y..I%A    |
| 00000048 42 C7 C4 AB 68 B3 C1 68 23 47 C0 58 BA 52 04 5F 2D 46 | B...h...h#G.X.R._-F  |

**Are the data (128 bytes) generated by md5collgen completely different for the two output files? Please identify all the bytes that are different.**

A small portion of the bytes are different and scattered throughout.

Lines affected in p.txt are denoted by <

Lines affected in q.txt are denoted by >

The diff command outputs a set of instructions for how to change one file to make it identical to the second file [1]. In this case we have 6,8c6,8 and 10,12c10,12 where c means change.

Thus, this output can be interpreted as lines 6-8 in the first file (p.txt) need to be changed to lines 6-8 in the second file (q.txt) and lines 10-12 in p.txt need to be changed to lines 10-12 in q.txt.

```
lab1md5.txt md5collgen out1.bin out2.bin
liam@ubuntu:~/Lab1$ xxd out1.bin > p.txt
liam@ubuntu:~/Lab1$ xxd out2.bin > q.txt
liam@ubuntu:~/Lab1$ diff p.txt q.txt
6,8c6,8
< 00000050: 2347 c0d8 ba52 045f 2d46 4136 b396 1ea3 #G...R._-FA6....
< 00000060: 62ee cb05 3c6e 3b7d 6f0c e9f9 d151 e243 b...<n;}o....Q.C
< 00000070: b330 3ee6 30ea 7c45 d79f 354c a16a c5f6 .0>.0.|E..5L.j..
---
> 00000050: 2347 c058 ba52 045f 2d46 4136 b396 1ea3 #G.X.R._-FA6....
> 00000060: 62ee cb05 3c6e 3b7d 6f0c e9f9 d1d1 e243 b...<n;}o.....C
> 00000070: b330 3ee6 30ea 7c45 d79f 35cc a16a c5f6 .0>.0.|E..5..j..
10,12c10,12
< 00000090: 2903 cb1a ca0e 19ef 74e7 cbcf 4f4d 43c1 ).....t...OMC.
< 000000a0: ea46 c5ea 7465 9357 9d94 6fd9 d329 d9e6 .F..te.W..o..)..
< 000000b0: af87 12f7 1282 a4e8 c04e 5b4f 9528 7c9a .....N[0.(|.
---
> 00000090: 2903 cb9a ca0e 19ef 74e7 cbcf 4f4d 43c1 ).....t...OMC.
> 000000a0: ea46 c5ea 7465 9357 9d94 6fd9 d3a9 d8e6 .F..te.W..o.....
> 000000b0: af87 12f7 1282 a4e8 c04e 5bcf 9528 7c9a .....N[...(|.

```

## Task 2: Further Understanding MD5's Property

If two inputs to the MD5 algorithm create a collision, then adding the same suffix (T) to them will result in a collision as well. To demonstrate this, I created two new files f1.txt and f2.txt with the same contents inside each file. To ensure a collision occurred I ran a checksum on both files and sure enough the MD5 hashes for each file are the same. This satisfies the initial part of the task in that  $MD5(M_1) = MD5(M_2)$ .

```
liam@ubuntu:~/Lab1$ md5sum f1.txt
f30834534bedae23808b9a4d2b1464e9  f1.txt
liam@ubuntu:~/Lab1$ md5sum f2.txt
f30834534bedae23808b9a4d2b1464e9  f2.txt
```

The next part of the task is to satisfy the following condition: *for any input T,  $MD5(M_1 || T) = MD5(M_2 || T)$ , where  $||$  means concatenate*. To do this I created a new text file and concatenate it to both f1.txt and f2.txt. Finally, I re-run the md5 check sum on the concatenated files and it again returns two identical MD5 Hashes.

```
liam@ubuntu:~/Lab1$ echo "suffix" >> suffix.txt

liam@ubuntu:~/Lab1$ cat f1.txt suffix.txt > f3.txt
liam@ubuntu:~/Lab1$ cat f2.txt suffix.txt > f4.txt
liam@ubuntu:~/Lab1$ ls
f1.txt  f3.txt  lab1md5.txt  out1.bin  p.txt  suffix.txt
f2.txt  f4.txt  md5collgen  out2.bin  q.txt
liam@ubuntu:~/Lab1$ md5sum f3.txt
1b798f4c1d8d556a6dc2cfc7b7187a63  f3.txt
liam@ubuntu:~/Lab1$ md5sum f4.txt
1b798f4c1d8d556a6dc2cfc7b7187a63  f4.txt
liam@ubuntu:~/Lab1$
```

Thus, we know that if two objects have the same hash, the concatenation of data onto both objects will result in them having the same hash as well.

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

Below is a sample executable C program. The array contents are what will vary in this program.

```

/* c
#include <stdio.h>
unsigned char xyz[200] = {
/* The contents of this array will vary */
};
int main()
{
    int i;
    for (i=0; i<200; i++){
        printf("%x", xyz[i]);
    }
    printf("\n");
}
*/

```

I had previously typed out 64 As in task 1 and then double checked the length with the `/s -/` command and option. However, I was not about to type out 200 A's and since making files of some size with arbitrary content seemed to be a reoccurring theme, I searched for a method and found one that would streamline this process for me [2]. These 200 A's will be inserted into the array so when I dump the hex I can easily find where it starts.

```
liam@ubuntu:~/Lab1$ echo -e "print('\0x41','*199)" | python3 > fileofA
```

Then copying all these contents into the array in the provided c program, I compiled with gcc and created an output file xyz.o.

```
liam@ubuntu:~/Lab1$ gcc xyz.c -o xyz.o
```

Skimming the hex, it was easy to find where the array started, and the offset.

xyz.o

00002fe2

00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00

.....

00002f4f

00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00

.....

00003006

00 00 08 40 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00

...@.....

00003018

00 00 00 00 00 00 00 00 41 41 41 41 41 41 41 41 41 41 41 41

.....AAAAAAAA

0000302a

41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41 41

AAAAAAAAAAAAAAAA

|  |       |   |                  |                  |                       |
|--|-------|---|------------------|------------------|-----------------------|
| Signed 8 bit:  | 65    | Signed 32 bit:  | 1094795585       | Hexadecimal:     | 41 41 41 41           |
| Unsigned 8 bit:                                      | 65    | Unsigned 32 bit:                                      | 1094795585       | Decimal:         | 065 065 065 065       |
| Signed 16 bit:                                       | 16705 | Float 32 bit:   | 12.07843         | Octal:           | 101 101 101 101       |
| Unsigned 16 bit:                                     | 16705 | Float 64 bit:   | 2261634.50980392 | Binary:          | 01000001 01000001 010 |
| <input type="checkbox"/> Show little endian decoding |       | <input type="checkbox"/> Show unsigned as hexadecimal |                  | ASCII Text: AAAA |                       |
| Offset: 0x3020 / 0x4257                              |       |   |                  | Selection: None  |                       |
| INS  |       |   |                  |                  |                       |

The offset is given in hex, converting to decimal it is 12,320. Dividing by 64 tells us if it is a multiple of 64, and it is not so I must jump to the next multiple which is at 0x3040 or decimal 12352.

| Programmer |                     |
|------------|---------------------|
| 3020       |                     |
| HEX        | 3020                |
| DEC        | 12,320              |
| OCT        | 30 040              |
| BIN        | 0011 0000 0010 0000 |

The prefix required will contain everything before the array and a portion of the array, to ensure our prefix remains a multiple of 64 I just jumped to the next offset that was a multiple of 64.

The prefix is followed by a 128-bytegap (where a p or q value goes) and then the suffix. So, our prefix +

128 + 1 is the offset for suffix.  $12352 + 128 + 1 = 12,481$ .

Now I want to use these offsets to create the prefixes and suffixes.

```
liam@ubuntu:~/Lab1$ head -c 12352 xyz.o > prefix
liam@ubuntu:~/Lab1$ tail -c +12481 xyz.o > suffix
liam@ubuntu:~/Lab1$
```

Generate collisions with the prefix.

```
liam@ubuntu:~/Lab1$ ./md5collgen -p prefix -o p.bin q.bin
MD5 collision generator v1.5
by Marc Stevens (http://www.win.tue.nl/hashclash/)

Using output filenames: 'p.bin' and 'q.bin'
Using prefixfile: 'prefix'
Using initial value: b9a6f7bbceb320b0dc6128ad2226f1c2

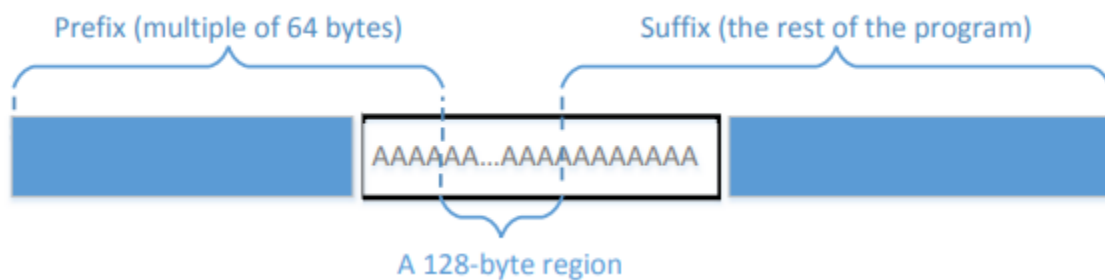
Generating first block: .....
Generating second block: S00.....
Running time: 22.3524 s
liam@ubuntu:~/Lab1$
```





To elaborate on the above results what we have done is taken a program, dissected its hex to find the offset where the array starts and then found some offset within the array that was a multiple of 64. It must be a multiple of 64 so that no padding is added. Finding this offset leaves us with three keys pieces:

- **Prefix:** some multiple of 64 that goes from the beginning of the program to a point within the array.
- **128-byte region:** A portion within the array after the prefix and before the suffix that contains either a P value or a Q value. These P and Q values are like the magic values that will allow us to have two different programs with the same checksum.
- **Suffix:** Everything in the program after the 128-byte region.



[3]

## Task 4: Making the Two Programs Behave Differently

The idea here is that the certificate check is only concerned with the hash value of the program and not the code itself, so in code I can have malicious and benign routes of execution. For this lab this requires us to create a collision. The execution route is determined by whether two lists are identical, if they are identical, execute benign code, if they aren't do something malicious. However, in general this would cause a change in the program data that would change the hash value of it, nullifying the certificate. If we can change the contents of the list without changing the hash, we can circumvent this. This is where the collision comes in. By going from the start of the program to the start of the first array (or where inside it is the first multiple of 64) we can use MD5collgen to create a prefix that can be used to generate a collision, that is, the start of the program with a differing 128-byte ending (p & q). Then we append a suffix (the rest of the program after our prefix). If we copy the 128-byte P value and paste it into the same part in the second array as the first, we have a program that will run the benign code and produce a hash value of x. This is what we would send in to get certified. Then we can substitute the 128-byte p value in one of the arrays with the 128-byte q value. The contents will be different, but the hash will still be x, and now the execution will be malicious. I started by creating a program that would do something benign when the contents of two arrays are the same and then do something else when they are not the same, then compiled and ran it to make sure it worked.

[illegible]

```

19 unsigned char y[200] = {
20     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
21     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
22     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
23     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
24     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
25     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
26     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
27     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
28     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
29     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
30     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
31     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41,
32     0x41,0x41,0x41,0x41,0x41,0x41,0x41,0x41
33 };
34
35 int main()
36 {
37     int i;
38
39     for (i = 0; i<200; i++){
40         if (x[i] != y[i])l
41             printf("malicious\n");
42         break;
43     }else if(i==199){
44         printf("benign\n");
45     }
46 }
47 return 1;
48 }

```



```
liam@ubuntu:~/Lab1/lab1t4$ cat prefixP.bin suffix > p_injection
liam@ubuntu:~/Lab1/lab1t4$ cat prefixQ.bin suffix > q_injection
liam@ubuntu:~/Lab1/lab1t4$ diff prefixP.bin prefixQ.bin
Binary files prefixP.bin and prefixQ.bin differ
liam@ubuntu:~/Lab1/lab1t4$ md5sum prefixP.bin
bd93f80081f77a523c3ae24b609d033a  prefixP.bin
liam@ubuntu:~/Lab1/lab1t4$ md5sum prefixQ.bin
bd93f80081f77a523c3ae24b609d033a  prefixQ.bin
liam@ubuntu:~/Lab1/lab1t4$
```

[illegible][illegible]



Interesting bumps and hurdles I ran into:

- Apparently, bless doesn't have a tmp file to save to by default in Ubuntu when modifying hex directly, so I had to go into the .config/preferences.xml of bless and add a file path for it to save to so that I could do task 4 the way I did, where I copy bytes from one array and paste them into the other and then save it.
  - o It looks something like this `<pref name="ByteBuffer.TempDir"></pref>`
  - o between the two red arrows I added a file and it fixed it.

## Citations

<https://www.computerhope.com/unix/udiff.htm> [1]

<https://stackoverflow.com/questions/2043453/executing-multi-line-statements-in-the-one-line-command-line> [2]

[SEED Project \(seedsecuritylabs.org\)](https://seedsecuritylabs.org) [3]