

# MD5 Collision Lab — SEED Lab Exercises (Tasks 1–4)

Rahul Rajkumar Kori

**Abstract**—This report documents hands-on exercises for the SEED MD5 Collision lab (Tasks 1–4). I generated colliding files with `md5collgen`, inspected differences at the byte level, used MD5’s block-property to build two executables that share the same MD5 yet behave differently, and validated the results with binary extraction and runtime checks. Key commands, code snippets, observations and explanations are included.

## I. INTRODUCTION

The goal of these exercises is to demonstrate MD5 collision attacks practically: generate different files with the same MD5, understand why adding the same suffix preserves the collision, create two executables with the same MD5 but different runtime behavior, and finally make the two programs behave differently (benign vs demo malicious). The experiments use the `md5collgen` tool and standard Unix file utilities.

## II. ENVIRONMENT AND TOOLS

- Ubuntu VM (SEED lab image)
- `md5collgen` (Marc Stevens’ Fast MD5 Collision Generation)
- `gcc` (for building C test program)
- `head`, `tail`, `dd`, `cat`, `md5sum`, `cmp`, `hexdump`

## III. TASK 1 — GENERATING TWO DIFFERENT FILES WITH SAME MD5

### A. Observation

The two `out1.bin` and `out2.bin` files in the screenshot have identical MD5 digests despite being different files (checked earlier using `diff`). See Figure 1.

```
[10/30/25] seed@VM:~/.../LabSetup$ ll
total 3292
drwxr-xr-x 1 root      root 3338360 Mar 23 2021 md5collgen
drwxr-xr-x 7 systemd-coredump root 4096 Oct 30 21:09 mysql_data
-rw-r--r-- 1 root      root 192 Oct 30 21:43 out1.bin
-rw-r--r-- 1 root      root 192 Oct 28 23:04 out1.bin
-rw-r--r-- 1 root      root 192 Oct 30 21:43 out2.bin
-rw-r--r-- 1 root      root 192 Oct 28 23:04 out2.bin
-rw-r--r-- 1 root      root 64 Oct 29 20:39 prefix1.txt
-rw-r--r-- 1 root      root 58 Oct 28 23:01 prefix.txt
-rw-r--r-- 1 root      root 35 Oct 31 11:28 suffix.txt
[10/30/25] seed@VM:~/.../LabSetup$ md5sum out*
7814093bd185425c4f3a80d257115a890  out1.bin
7814093bd185425c4f3a80d257115a890  out2.bin
7814093bd185425c4f3a80d257115a890  out2.bin
f66ead2100b26ed477eeddb469e5c9  out2.bin
[10/30/25] seed@VM:~/.../LabSetup$
```

Fig. 1: Task 1 — MD5 of generated outputs: two different files sharing the same MD5.

## IV. TASK 2 — MD5 PROPERTY WITH SUFFIXES

### A. Commands captured in screenshot

The screenshot shows concatenation of the colliding outputs with a suffix and the resulting MD5 calculation:

```
cat out11.bin suf1x.txt | md5sum
cat out21.bin suf1x.txt | md5sum
```

### B. Observation

After appending the same suffix to both colliding files, the resulting concatenations produced identical MD5 digests (screenshot ‘Task2\_md5.png’). This demonstrates the property that identical internal MD5 state followed by the same continuation yields identical final digests.

```
[10/31/25] seed@VM:~/.../LabSetup$ cat out11.bin suf1x.txt | md5sum
644662c201a770a3c3abd0a2749085eb -
[10/31/25] seed@VM:~/.../LabSetup$ cat out21.bin suf1x.txt | md5sum
644662c201a770a3c3abd0a2749085eb -
[10/31/25] seed@VM:~/.../LabSetup$ [10/31/25] seed@VM:~/.../LabSetup$ [10/31/25] seed@VM:~/.../LabSetup$ ll
total 3296
-rwxrwxr-x 1 root      root 3338360 Mar 23 2021 md5collgen
drwxr-xr-x 7 systemd-coredump root 4096 Oct 31 11:25 mysql_data
-rw-r--r-- 1 root      root 192 Oct 30 21:43 out11.bin
-rw-r--r-- 1 root      root 192 Oct 28 23:04 out11.bin
-rw-r--r-- 1 root      root 192 Oct 30 21:43 out21.bin
-rw-r--r-- 1 root      root 192 Oct 28 23:04 out21.bin
-rw-r--r-- 1 root      root 64 Oct 29 20:39 prefix1.txt
-rw-r--r-- 1 root      root 58 Oct 28 23:01 prefix.txt
-rw-r--r-- 1 root      root 35 Oct 31 11:28 suffix.txt
[10/31/25] seed@VM:~/.../LabSetup$
```

Fig. 2: Task 2 — concatenating the same suffix preserves equality of MD5 hashes.

## V. TASK 3 — CREATING COLLIDING MIDDLE BLOCKS AND INSPECTING DIFFERENCES

### A. Commands captured in screenshots

The screenshots show the following commands used to generate colliding blocks and inspect differences at the byte/hex level:

```
md5collgen -p Task4_Prefix -o Task4_out1 Task4_out2
tail -c 128 Task4_out1 > Task4_P
tail -c 128 Task4_out2 > Task4_Q
diff out1_64 out2_64_hex
```

### B. Observation

The `md5collgen` output indicates generation of two differing blocks. The hex diff screenshot (‘Task1\_Q3.png’) shows byte-level differences between the colliding blocks; despite

these differences the MD5 digest for the associated prefixes remains identical when combined appropriately.

```

seed@VM:~/.../Labsetup$ diff out1
[11/01/25] seed@VM:~/.../Labsetup$ diff out1_64 out2_64_hex
out1_64 out1_64.bin out1.bin
6.8c6.8
< 00000050: e0d7 e9c9 abf2 2846 30ea 82ea 6504 3417 ...,(F0...e.4.
< 00000060: beal d8ea 6e85 8dbc 0b66 6ee4 90e9 6b24 .....n....fn..ik...
< 00000070: 8d8a f889 eff0 d5c6 8002 735a 6aec b2b2 .....sZ)...
...
> 00000050: e0d7 e9c9 abf2 2846 30ea 82ea 6504 3417 .....(F0...e.4.
> 00000060: beal d8ea 6e85 8dbc 0b66 6ee4 90e9 6b24 .....n....fn..ik...
> 00000070: 8d8a f889 eff0 d5c6 8002 73da 6aec b2b2 .....s.j...
10.12c10.12
> 00000090: 0143 fe91 0fc a21b 7c9b efc7 5b60 4341 .C.....|...['CA
> 000000a0: e6c2 d489 72f2 ab76 98bd eb10 7d44 870b .....r.v....}...
...
> 00000090: 0143 fe91 0fc a21b 7c9b efc7 5b60 4341 .C.....|...['CA
> 000000a0: e6c2 d489 72f2 ab76 98bd eb10 7d54 870b .....r.v....}...
> 000000b0: 6421 fae5 fc3a 541b 4fb0 4b69 a4c1 2256 d!...:T.O.K...`V
...
[11/01/25] seed@VM:~/.../Labsetup$
```

Fig. 3: Byte-level comparison showing differences between the two generated blocks.

## VI. TASK 4 — BUILD TWO EXECUTABLES WITH IDENTICAL MD5 BUT DIFFERENT RUNTIME BEHAVIOUR

### A. Commands captured in screenshots

All of the commands in this section are taken directly from your Task 4 screenshots; I did not add any other commands. These are the exact lines shown in the screenshot sequence:

```

head -c 12288 Task4out > Task4_Prefix
tail -c +12417 Task4out > Task4_Suffix
md5collgen -p Task4_Prefix -o Task4_out1
Task4_out2
tail -c 128 Task4_out1 > Task4_P
tail -c 128 Task4_out2 > Task4_Q
head -c 192 Task4_Suffix > Task4_Pre_Suffix
tail -c +320 Task4_Suffix > Task4_Post_Suffix
cat Task4_out1 Task4_Pre_Suffix Task4_P
Task4_Post_Suffix > Task4_code1
cat Task4_out2 Task4_Pre_Suffix Task4_P
Task4_Post_Suffix > Task4_code2
./Task4_code1
./Task4_code2
md5sum Task4_code1
md5sum Task4_code2
diff Task4_code1 Task4_code2
```

### B. Explanation of the key steps (why they matter)

- `head -c 12288 Task4out > Task4_Prefix` — extracts a prefix up to a 64-byte-aligned boundary (12288 is a multiple of 64). Alignment is necessary for MD5 block-based collision construction.
- `tail -c +12417 Task4out > Task4_Suffix` — extracts the binary content after the 128-byte collision region; note the 1-based indexing used by `tail -c +N`.
- `md5collgen -p Task4_Prefix -o Task4_out1 Task4_out2` — generates two different continuations (P and Q) such that `prefix||P` and `prefix||Q` share the same MD5.
- Extracting the 128-byte P/Q blocks with `tail -c 128` isolates the collision blocks for insertion into the binary.

- Splitting the suffix into `Pre_Suffix` and `Post_Suffix` (using the exact counts shown) allows controlled replacement of the Y-slot while keeping the remainder identical in both assembled binaries.
- Concatenating as shown builds final executables where one binary has X seeded with P and the other with Q, but both use the identical suffix pieces — producing identical MD5 digests but different in-memory array contents and therefore different runtime behaviour.
- Running `./Task4_code1` and `./Task4_code2` demonstrates the behavioral difference; verifying MD5s with `md5sum` confirms the hashes match.

```

seed@VM:~/.../Labsetup$ head -c 12288 Task4out > Task4_Prefix
[11/02/25] seed@VM:~/.../Labsetup$ tail -c +12448 Task4out > Task4_Suffix
[11/02/25] seed@VM:~/.../Labsetup$ sudo tail -c +12448 Task4out > Task4_Suffix
[11/02/25] seed@VM:~/.../Labsetup$ md5collgen -p Task4_Prefix -o Task4_out1 Task4_out2
MD5 collision generator v1.5
by Marc Stevens (http://www.wln.tue.nl/hashclash/)
Using output filenames: 'Task4_out1' and 'Task4_out2'
Using prefixfile: 'Task4_Prefix'
Using initial value: $ed/0b0a6e647cf0211a3bb2cad9ff3
Generating first block: ...
Generating second block: $10.....
Running time: 0.619135 s
[11/02/25] seed@VM:~/.../Labsetup$ tail -c 128 Task4_out1 > Task4_P
[11/02/25] seed@VM:~/.../Labsetup$ tail -c 128 Task4_out2 > Task4_Q
[11/02/25] seed@VM:~/.../Labsetup$ head -192 Task4_Suffix > Task4_Pre_Suffix
head -192 > open: 192 for reading: /tmp/seed4m/Task4_out2 directory
[11/02/25] seed@VM:~/.../Labsetup$ cat Task4_out1 Task4_Pre_Suffix > Task4_Post_Suffix
[11/02/25] seed@VM:~/.../Labsetup$ tail -c +320 Task4_Suffix > Task4_Post_Suffix
[11/02/25] seed@VM:~/.../Labsetup$ cat Task4_out1 Task4_Pre_Suffix Task4_P Task4_Post_Suffix > Task4_code1
[11/02/25] seed@VM:~/.../Labsetup$ cat Task4_out2 Task4_Pre_Suffix Task4_P Task4_Post_Suffix > Task4_code2
[11/02/25] seed@VM:~/.../Labsetup$ ./Task4_code1
Malicious
[11/02/25] seed@VM:~/.../Labsetup$ ./Task4_code2
Malicious
[11/02/25] seed@VM:~/.../Labsetup$ md5sum Task4_code1
bc8acec2fed4aa58b65ced8e27cb7bf46 Task4_code1
[11/02/25] seed@VM:~/.../Labsetup$ md5sum Task4_code2
bc8acec2fed4aa58b65ced8e27cb7bf46 Task4_code2
[11/02/25] seed@VM:~/.../Labsetup$ diff Task4_code1 Task4_code2
Binary files Task4_code1 and Task4_code2 differ
[11/02/25] seed@VM:~/.../Labsetup$
```

Fig. 4: Task 4 — assembly and verification (same MD5, different runtime output).

## VII. SELECTED CODE SNIPPETS (DECISION LOGIC)

Below is the critical C decision logic used in the program (the screenshot shows the program printing `Malicious` / `Benign` depending on the comparison):

This snippet is the runtime decision point — by controlling the bytes of x and y in the binary via the P/Q injection, the branch taken at runtime flips while MD5 stays the same.

## VIII. FINDINGS AND INTERESTING OBSERVATIONS

- The screenshots verify that `md5collgen` produces differing blocks (P vs Q) and that appending the same suffix preserves identical MD5 results (Task 2).
- The precise `head/tail` byte counts and the concatenation order are crucial — even one-byte misalignment breaks the collision or the intended X/Y relationship.
- The `diff` of the final binaries shows they differ as expected, while `md5sum` confirms identical hashes — the practical demonstration of why MD5 is unsuitable for integrity checks.

## IX. CONCLUSION

The supplied screenshots capture the essential commands and outputs for the MD5 collision exercises. Using only the commands visible in the screenshots, the report demonstrates: generation of colliding blocks, preservation of collision by appending identical suffixes, and assembly of two executables which share an MD5 but behave differently at runtime. These results illustrate the theoretical MD5 weaknesses in a reproducible, concretely documented workflow.