

Term Project

1. Introduction

The project aims to organize a folder named "Unprocessed-Passwords.txt," which contains random passwords, into subfolders based on the first character of each password. While developing the program, I prioritized finding and selecting the most efficient methods to ensure optimal performance and effectiveness.

Key Components:

1. Password Processing:

- **Function:** `processFolder`
- **Task:** Reads passwords from unprocessed files, normalizes them, and writes unique entries to processed files.
- **Normalization:**
 - Converts the first character of each password to lowercase if it's alphabetical.
 - Replaces special characters not allowed in filenames (e.g., \, /, :, *, ?, <, >, |) with underscores (_).

2. Password Indexing:

- **Function:** `index`
- **Task:** Reads processed passwords, computes their MD5, SHA-1, and SHA-256 hashes, and stores them in an indexed folder structure.
- **Folder Structure:** Based on the first character of the password, creating subfolders and storing passwords in files. Ensures that files do not exceed a certain number of lines (500 in this case).

3. Password Searching:

- **Function:** `Search`
- **Task:** Searches for a given password in the indexed files. If not found, computes its hash values, stores it in the appropriate file, and updates the file line count.
- **Output:** Indicates if the password was found and records it if not.

4. Hash Computation:

- **Functions:** `calculate_md5`, `calculate_sha1`, `calculate_sha256`
- **Task:** Computes and returns MD5, SHA-1, and SHA-256 hash values for given password strings.

5. Main Execution:

- **Function:** `main`
- **Task:** Orchestrates the processing, indexing, and searching operations. Measures and reports the time taken for each search operation and calculates the average search time over multiple iterations.

Scope of the Project:

1. Input:

- Unprocessed password files located in a specified directory.

2. Output:

- Processed password files stored in a designated directory.
- Indexed files with hash values stored in a structured directory for efficient searching.

3. Performance Measurement:

- Time taken for each search operation is measured and the average search time is calculated over a set of searches.

2. Tools and Technologies Used:

The project was implemented using the C++ programming language. Visual Studio Code (VS Code) served as the integrated development environment (IDE). The following libraries were utilized to achieve the project's objectives:

Standard Libraries:

- **<iostream>**: For standard input and output operations.
- **<fstream>**: For file input and output operations.
- **<filesystem>**: For directory manipulation and traversal.
- **<unordered_set>**: For efficiently storing unique passwords.
- **<unordered_map>**: For tracking the number of lines in each output file.
- **<iomanip>, <sstream>, <cctype>**: For formatting, string manipulation, and character type checking.

Cryptographic Libraries:

- **<openssl/md5.h>** and **<openssl/sha.h>**: Header files from the OpenSSL library, which provides a comprehensive suite of cryptographic functions.

Cryptographic Functions

OpenSSL is an open-source implementation of SSL and TLS protocols, designed to secure communication over networks. It offers a wide range of cryptographic functions, including those used in this project:

• MD5:

- Produces a 128-bit hash value.
- Initially designed for cryptographic use but later found to have vulnerabilities.

- Still used for database partitioning and file transfer validation due to its checksum capabilities.
- **SHA (Secure Hash Algorithm):**
 - **SHA-1:**
 - Generates a 160-bit hash (20 bytes).
 - Designed for cryptographic applications but has known vulnerabilities, making it unsuitable for secure applications today.
 - **SHA-2:**
 - Includes several variants, with SHA-256 being the most commonly used.
 - Produces a 256-bit hash value (64 hexadecimal digits).
 - Recommended by the National Institute of Standards and Technology (NIST) for its enhanced security compared to MD5 and SHA-1.
 - SHA-256 is about 20-30% slower to compute than MD5 or SHA-1 but offers significantly better security.

By leveraging these tools and technologies, the project ensures efficient and secure processing and organization of password data.

Installing and compiling OpenSSL:

After installing the OpenSSL library, I encountered some difficulties linking the libraries to the program.

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
Code + ⊞ ⌂ ⌄ ⌅ ⌆ ⌇ ⌈ ⌉ ×

⑧ PS C:\Users\ayakh> cd "c:\Users\ayakh\OneDrive\Desktop\Code\" ; if ($?) { g++ pass.cpp -o pass } ; if ($?) { .\pass
}
C:/Strawberry/c/bin/../lib/gcc/x86_64-w64-mingw32/13.1.0/../../../../x86_64-w64-mingw32/bin/ld.exe: C:/Users/ayakh\A
ppData\Local\Temp\ccRnUW5P.o:pass.cpp:(.text+0x5c): undefined reference to `MD5'
C:/Strawberry/c/bin/../lib/gcc/x86_64-w64-mingw32/13.1.0/../../../../x86_64-w64-mingw32/bin/ld.exe: C:/Users/ayakh\A
ppData\Local\Temp\ccRnUW5P.o:pass.cpp:(.text+0x18d): undefined reference to `SHA1'
C:/Strawberry/c/bin/../lib/gcc/x86_64-w64-mingw32/13.1.0/../../../../x86_64-w64-mingw32/bin/ld.exe: C:/Users/ayakh\A
ppData\Local\Temp\ccRnUW5P.o:pass.cpp:(.text+0x2be): undefined reference to `SHA256'
collect2.exe: error: ld returned 1 exit status
○ PS C:\Users\ayakh\OneDrive\Desktop\Code>

```

Through research, I resolved this issue by using the **g++** command, specifying the source file (**pass.cpp**), the output executable name (**-o pass**), and the flags needed to link against the OpenSSL library (**-lssl -lcrypto**).

The command used was:

```
g++ pass.cpp -o pass -lssl -lcrypto
```

This command instructs the compiler to link against the OpenSSL libraries (**libssl** and **libcrypto**) when compiling **pass.cpp**.

Once successfully compiled, the program can be run with the following command:

Once the program is successfully compiled, we can run the executable and it will work just fine 😊 :

```
./pass
```

3. Project Design and Configuration:

The folder structure for the project includes the following directories:

Code: Contains the source code files for the project.

Unprocessed-Passwords: Contains the raw password files to be Processed.

To my understanding “Processing the passwords”: is unifying the passwords and making them suitable to be file or folder name by changing the first character if it is among the characters that are not allowed to be file/folder names. I used the following algorithm to achieve that:

```
char firstChar = line[0];
if (!isalpha(firstChar)) {
    char c = firstChar;
    if (c == '\x5C' || c == '/' || c == ':' || c == '*' ||
c == '?' || c == '"' || c == '<' || c == '>' || c == '|') {
        line[0] = '_';
    }
} else {
    firstChar = tolower(line[0]);
```

Processed: Upon processing, passwords from the Unprocessed-Passwords folder are stored in this directory. Each processed password is saved in a file with the same name as its corresponding raw file.

Index: This folder contains the organized processed passwords, which are stored into subfolders based on the first character of each password. Each subfolder represents an alphabetic character and contains files storing passwords starting with that character.

4.Indexing Process:

To implement the indexing process I used the `index()` function:

This function indexes the passwords stored in the processed directory by computing hash values for each password, organizing them into corresponding subfolders based on their first character, and storing them in output files within those subfolders. It iterates through the files in the processed folder using the ``directory_iterator`` from the ``<filesystem>`` library.

```
for (const auto &entry : fs::directory_iterator(folderPath)) {
```

- **const auto &entry** : declares a loop variable **entry** of type **auto**, which automatically deduces the type of each element returned by the directory iterator
- **fs::directory_iterator(folderPath)** : initializes a directory iterator object for the directory specified by **FolderPath**. It allows us to iterate over the contents of the directory.

```
if (entry.is_regular_file()) {  
    ifstream inputFile(entry.path());
```

- if the file opened is “regular”, it opens an input file using the file’s path.

```
if (!inputFile) {  
    cerr << "Error opening file: " << entry.path() << endl;  
    continue;
```

- If a file cannot be opened, the function displays an error message and proceeds to the next iteration.

Then the file enters a while loop to read and process each line in it, using `getline()`.

- For each line, the function calculates three hash values: MD5, SHA1, and SHA256:

```
    string md5_hash = calculate_md5(line);  
    string sha1_hash = calculate_sha1(line);  
    string sha256_hash = calculate_sha256(line);
```

Subsequently, the function determines the subfolder path based on the first character of the password:

```
string subFolderPath = indexFolderPath + "/" + firstChar;
```

If the subfolder does not exist, it creates one:

```
fs::create_directory(subFolderPath);
```

Next, the function constructs the file path for storing the indexed passwords:

```
string filePath = subFolderPath + "/output_" + firstChar + ".txt";
```

If the file containing passwords surpasses 10,000 lines, the function creates additional files:

```
if (fileLineCounts[filePath] > 10000 - 1) {
    int count = fileLineCounts[filePath] / 10000;
    do {
        count++;
        filePath = subFolderPath + "/output_" + firstChar +
"_" + to_string(count - 1) + ".txt";
    } while (fileLineCounts[filePath] > 10000 - 1);
```

During this phase, an integer variable, **count**, is initialized to count the number of files needed after the first file exceeds 10,000 lines. Within the do-while loop, the program continuously increments the count and updates the file path until it finds or creates a filepath with less than 10,000 lines. This ensures that the passwords are distributed across multiple files, each containing a maximum of 10,000 lines.

Finally, the function appends the password along with its hash values and filename to the designated output file.

The process repeats for each line in each file in the processed folder, ensuring efficient organization and storage of indexed passwords.

5. Search Function and Performance Test:

To implement the searching process I used the search() function:

The `Search` function is designed to locate a specific password within the indexed folder. It takes the password to be searched as input and searches through the indexed folder with the same first character to find a match.

- The function begins by determining the subfolder path based on the first character of the password:

```
string subFolderPath = indexFolderPath + "/" + pass[0];
```

- It then opens that subfolder and iterates through its files:

```
for (const auto &entry : fs::directory_iterator(subFolderPath))
    while (getline(inputFile, line))
```

- For each line, the function extracts the first word (password) and compares it with the password being searched for:

```
if (!line.empty() && line.substr(0, line.find(' | ')) == password) {
    cout << "Found" << entry.path() << endl;
```

If a match is found, the function displays the file path where the password was found along with the corresponding line, and then exits.

If no match is detected in the current file during the search process, the function proceeds to the next file until all files in the subfolder have been examined exhaustively. However, if the password is not located in any of the files, the function takes additional steps to save the searched password so that it will be found in the next search.

Before storing the unfound password, the function processes it. It then determines the appropriate file to store the password based on its first character. Importantly, the function avoids storing the password in a file that has reached its maximum capacity, ensuring efficient management of file sizes and preventing potential data loss or overflow issues.

Performance Test: To measure search performance, the program records the time taken to search for 10 randomly selected passwords. The average search time is calculated based on these measurements. To do this I used `<chrono>` library.

```
for(int i = 1; i <= 10; i++){
    cout << "Enter Password to be searched: " << endl;
    cin >> pass;
    if (pass == "x") break;
    auto start = chrono::high_resolution_clock::now(); // Start timing
    Search(pass);
    auto end = chrono::high_resolution_clock::now(); // End timing
    chrono::duration<double> duration = end - start;
    cout << i << ".Time taken to search: " << duration.count() << "
seconds" << endl;
    total = duration.count() + total;
}
float avg = total/10;
cout << "Average Search time: " << avg << endl;
```

- screenshot of the results of 10 searches:

```

PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    PROFILE

password
Found"D:\\Code\\index/p\\output_p.txt"
1.Time taken to search: 0.0019012 seconds
Enter Password to be searched:
superman
Found"D:\\Code\\index/s\\output_s.txt"
2.Time taken to search: 0.0125746 seconds
Enter Password to be searched:
patrick
Found"D:\\Code\\index/p\\output_p.txt"
3.Time taken to search: 0.0009668 seconds
Enter Password to be searched:
steve
Found"D:\\Code\\index/s\\output_s.txt"
4.Time taken to search: 0.000969 seconds
Enter Password to be searched:
john
Found"D:\\Code\\index/j\\output_j.txt"
5.Time taken to search: 0.0119545 seconds
Enter Password to be searched:
muaddib
Found"D:\\Code\\index/m\\output_m_2.txt"
6.Time taken to search: 0.0040224 seconds
Enter Password to be searched:
popies
Password Not Found Recording....
7.Time taken to search: 0.0203468 seconds
Enter Password to be searched:
abc123
Found"D:\\Code\\index/a\\output_a.txt"
8.Time taken to search: 0.0017515 seconds
Enter Password to be searched:
abcd
Found"D:\\Code\\index/a\\output_a.txt"
9.Time taken to search: 0.0010563 seconds
Enter Password to be searched:
hunter
Found"D:\\Code\\index/h\\output_h.txt"
10.Time taken to search: 0.0124986 seconds
Average Search time: 0.00680417

```

6.Results and Evaluation:

The project effectively sorts and organizes passwords into subfolders based on their initial characters, offering efficient retrieval through the search function. Overall, the project achieves its goals of providing a functional password management system. However, there are areas for improvement, notably in performance optimization and code simplification. The program's execution speed could be enhanced(it took about 6 minutes).

- The index folder after running the program:

Name	Date modified	Type
#	18/05/2024 23:05	File folder
&	18/05/2024 23:05	File folder
_	18/05/2024 23:04	File folder
0	18/05/2024 23:04	File folder
1	18/05/2024 23:05	File folder
2	18/05/2024 23:04	File folder
3	18/05/2024 23:04	File folder
4	18/05/2024 23:04	File folder
5	18/05/2024 23:04	File folder
6	18/05/2024 23:04	File folder
7	18/05/2024 23:04	File folder
8	18/05/2024 23:04	File folder
9	18/05/2024 23:04	File folder
a	18/05/2024 23:05	File folder
b	18/05/2024 23:05	File folder
c	18/05/2024 23:05	File folder

- Inside the 's' folder:

Name	Date modified	Type	Size
Home			
Gallery			
Aya - Personal			
aya khalifa - İstanbul University - Ce			
output_s.txt	18/05/2024 23:11	Text Document	
output_s_1.txt	18/05/2024 23:05	Text Document	
output_s_2.txt	18/05/2024 23:06	Text Document	

- Inside the 's' file:

```

shadow|3bf1114a986ba87ed28fc1b5884fc2f8|ed9d3d832af899035363a69fd53cd3be8f71501|[0]bb9d080600ecc3eb9d7793a6f7859bedde2a2d83899b70bd78e961ed674b32f4|[0]-million-password-list-top-1000.txt
superman|84d961568a65073a3bcf0eb216b2a576|[1]8c28604dd31094a8d69daef6fbcd347f1afc5a|[7]3c1b16c4fb30361ad18a0b29b643a68d4a0075a466dc9e51682f84a847f5|[0]-million-password-list-top-1000.txt
soccer|da43aa0ad979d530df38c1a74e4f80|[5]17fa03ed5fc2475651c1d0ffa70e1bf5e5bd9|[8]f27f432fcbaa405186a1cc7abf81a66a93d43c1bce6f19922d0519d02f4b039|[0]-million-password-list-top-1000.txt
sunshine|0571749e2ac33a7455809c6b0e7af0|[d]de3e4f987851a0599257d3831a1af040886842f|[a]941aa4cd0c01dde6f161bb8e69fb4c1e2b0811c037c3f1835fdff6efc6223|[0]-million-password-list-top-1000.txt
starwars|5badcaf789d3d1d09794d8f021f4f0e|[3]27156a6b287c6aa52c8670e13163fc1bf666a0d4|[7]4ca0325b5fdb3a4b6d40a2581fb5d34417e8d3432952a5abc0929c1246|[0]-million-password-list-top-1000.txt
summer|6b1628b016dff466fa356b48beac9c|[6]420ed4d83b1a436d1e92d25605d18297296374e3|[e]8366425c6963e962b2b0f9fcfaa1b570df5fda69f5444ed37e5260f3ef689|[0]-million-password-list-top-1000.txt
secret|5ebe2294ec0deef08eab7690d2a6ee69|[e]5e9fa1ba31ec1a8e4f75caa447af3a6e3f05f4|[2]bb80d537bd1a3e38bd30361aa855686bde0ead7162f5e6a25fe97b5f27a25b|[0]-million-password-list-top-1000.txt
silver|97f04516561f487e368d015beb3f4|[f]8248e12727710c946f73d1f0e2eb93530ddde|[7]8cde4c3e47f2cbfd9da721f54aa[de337799166837c9de80962898feefac21|[0]-million-password-list-top-1000.txt
scooter|[efbf1c1f03504dfdfc52a9e6lab5e204d|[2]f3cd230e35fbefc59672f75448c446120b|[b]5ad1213079c486741d1d2e45bcbfa70c7f1125e2049c66f658bf40016c4e|[0]-million-password-list-top-1000.txt
samantha|[f]01e0d7992a3b7748538d0291b0bea|[c]5a7c3e21436a8e7617610ce551356f9aa745e|[c]5422c052bf7b0d9764e0467688b62193fc4fa32a1b1af28d1708d5870e|[0]-million-password-list-top-1000.txt
sparky|[a]269d17349c2977a897c7ec4ed0c3d4|[7]5a0a1c981feaf09a01381b13091b6d8e14577fc6|[1]855f8fe0e94e2f8a0b09707b3159a[cc0be5a9f6b1b8fc0c2e75b662398e050|[0]-million-password-list-top-1000.txt
snooky|[e]9646d086a37906e5be4323d3b37c9|[0]9639920909aac2d595b3d2d34e8a5fcab9fe3151|[6]d5ef1fdb8c5ab10b7c90f1796f711153c134f7c5772808bc0b1bf1d3e21|[0]-million-password-list-top-1000.txt

```

References:

- <https://www.freecodecamp.org/news/md5-vs-sha-1-vs-sha-2-which-is-the-most-secure-encryption-hash-and-how-to-check-them/>
- <https://www.openssl.org/>
- https://cplusplus.com/reference/unordered_set/unordered_set/
- <https://www.geeksforgeeks.org/chrono-in-c/>
- <https://www.youtube.com/watch?v=aEgG4pidcKU>

It is also important to mention that I used ChatGPT and other AI tools to refine and format this document.