



# İSTANBUL ÜNİVERSİTESİ CERRAHPAŞA

2024-04-28

BIMU2061 – File Organization

## Table of Contents

1.	Introduction .....	3
2.	Tools and Technologies Used.....	4
3.	Project Design and Configuration.....	5
4.	Indexing Process .....	6
	Indexing Function .....	6
	Functionality for Reading the Index.....	6
5.	Search Function and Performance Testing .....	7
6.	Results and Evaluation .....	9
7.	References.....	10

# 1. Introduction

The purpose of this project is to index text files under a specific folder by reading them.

The project code reads the passwords in each line of the text files (without reading repeated passwords), processes these passwords in a specific format, and saves them to a different folder, thus aiming to speed up the search function.

In line with the project objectives, cryptographic functions (MD5, SHA256, etc.) are used to process each password.

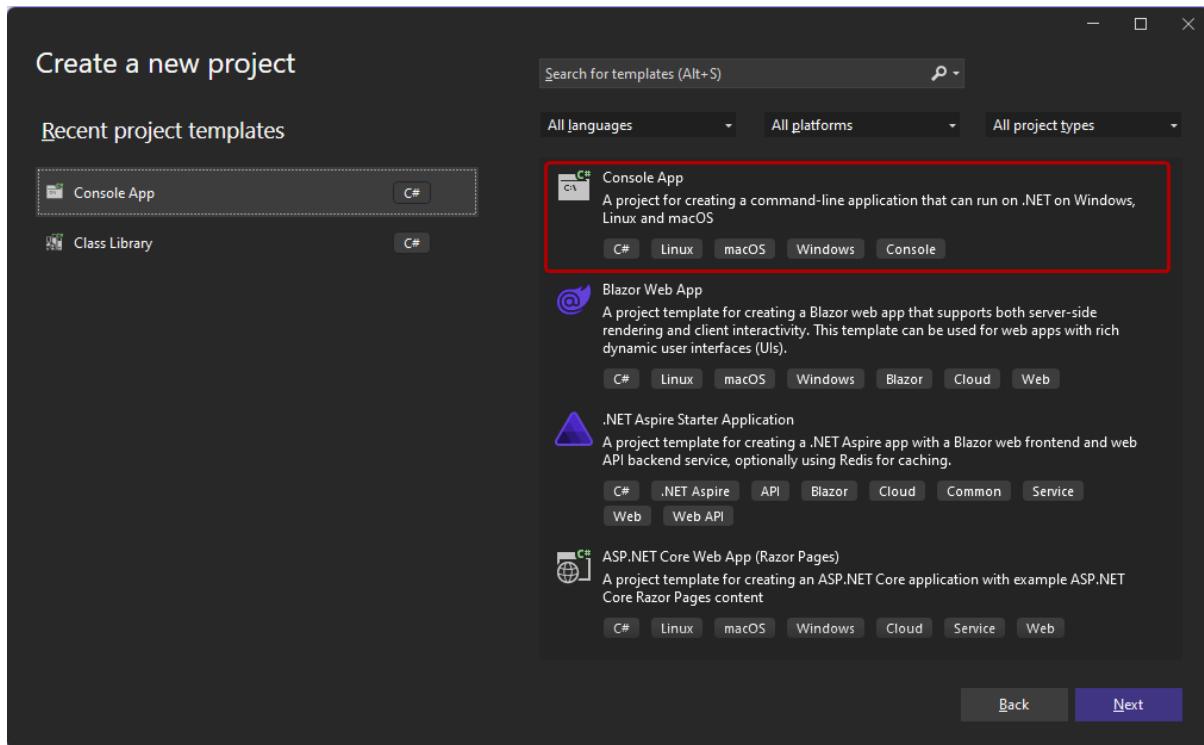
## 2. Tools and Technologies Used

The code was written in C# language, according to [.NET version 8.0 and C# version 12](#).

Visual Studio 2022 Preview version 17.10.0 was used as the code editor.

The core functionality layer of the project, **SecureList.Core**, was created using the Class Library template; for running the project on the console screen, **SecureList.Console** was created using the Console App template.

No third-party libraries were used aside from .NET standard class libraries.



### 3. Project Design and Configuration

The folder structure and configuration of the project are detailed in the image below.

Term Project		Projenin bulunduğu ana klasör
Code		Kod dosyalarının bulunduğu klasör
SecureList.Console		Projenin konsol ekranında çalıştırılması için oluşturulan proje
Program.cs SecureList.Console.csproj		Programın çalışma kodlarının bulunduğu kod dosyası SecureList.Console proje yapısının bulunduğu kod dosyası
bin		SecureList.Console projesinin derlenmiş dosyalarının bulunduğu klasör
SecureList.Core		Projenin ana işlevlerinin bulunduğu proje
IndexService.cs SecureList.Core.csproj StaticDetails.cs		Indexleme işlevlerinin bulunduğu kod dosyası SecureList.Core proje yapısının bulunduğu kod dosyası Projedeği değişmezlerin bulunduğu kod dosyası
bin		SecureList.Core projesinin derlenmiş dosyalarının bulunduğu klasör
Entities		Projede kullanılan somut sınıfların bulunduğu klasör
PasswordDirectory.cs PasswordFile.cs PasswordRepository.cs		Unprocessed-Passwords klasörünü temsil eden sınıf Unprocessed-Passwords klasörü altındaki her bir dosyayı temsil eden sınıf Program içerisinde şifrelerin tutulduğu sınıf
Helpers		Yardımcı işlev sınıflarının bulunduğu klasör
CryptographyHelper.cs FileHelper.cs		Kriptografik yardımcı işlevlerin bulunduğu sınıf Dosya ve klasör okuma-yazma yardımcı klasörlerinin bulunduğu sınıf
Index		İşlenmiş şifrelerin indexlendiği klasör
0		0 ile başlayan şifrelerin bulunduğu klasör
12.txt 13.txt		İçinde en fazla 10.000 satır bulunan şifrelerin indexlendiği bir dosya <a href="#">Bu dosya 12.txt 10.000 satırı geçiyorsa oluşturulur.</a>
1		1 ile başlayan şifrelerin bulunduğu klasör
0.txt		
2		2 ile başlayan şifrelerin bulunduğu klasör
19.txt		
Burada daha fazla klasör vardı. Okunulur olması açısından kısa tutuldu.		
⋮		
Processed		İşlenmiş dosyaların bulunduğu klasör
10-million-password-list-top-100.txt best15.txt top-20-common-SSH-passwords.txt top-passwords-shortlist.txt userDefined.txt		Şifrelerin bulunduğu bir dosya. Şifrelerin bulunduğu başka bir dosya. Kullanıcı yukarıdaki dosyalarda olmayan bir şifre girdiğinde buraya kaydedilir. Ardından indexlenir.
Unprocessed-Passwords		İşlenmemiş dosyaların bulunduğu klasör
10-million-password-list-top-100.txt best15.txt top-20-common-SSH-passwords.txt top-passwords-shortlist.txt userDefined.txt		Şifrelerin bulunduğu bir dosya. Şifrelerin bulunduğu başka bir dosya. Kullanıcı yukarıdaki dosyalarda olmayan bir şifre girdiğinde buraya kaydedilir. Ardından indexlenir.

## 4. Indexing Process

### Indexing Function

The indexing process is implemented within the **IndexService** class in the **SecureList.Core** project. This class utilizes a private property named `_passwordRepository`, which internally uses the `Dictionary<TKey, TValue>` structure. This data structure is highly efficient for search operations.

Let's elaborate on the operations within the **IndexService** class:

1. The `_passwordsRepository.Passwords` property stores passwords as key-value pairs in a *Dictionary* data structure, where the key represents the password and the value represents the folder where the password is located.
2. The `Index()` method of the **IndexService** class groups the current password and file name pairs by their initial characters. This operation is performed using the `GroupBy()` method from LINQ.
3. For each group:
  - An anonymous type consisting of password and file name pairs is created.
  - MD5, SHA1, and SHA256 hashes are calculated for the password. This function is performed by the `CryptographyHelper.CalculateHashes(string password)` method.
  - Each pair is concatenated into a string, separated by the "|" character, including the password, hashes, and file name.
4. The generated strings are divided into subgroups based on the `StaticDetails.MaximumLineNumberInIndexFiles` number (i.e., 10,000). This operation is performed using the `Chunk()` method from LINQ.
5. For each subgroup, a directory is created starting with the respective character, and a file name is assigned to this directory. Subsequently, the content of each subgroup (each chunk) is written to the respective subgroup file.

### Functionality for Reading the Index

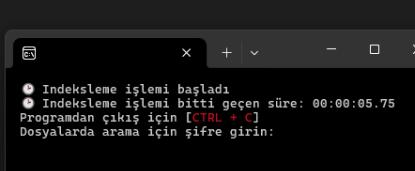
1. The `ReadIndex` method takes the directory where the index is located.
2. If the directory exists, it retrieves all file paths in that directory. Then, for each file path, it reads each line in the file.
3. Each line is split by the "|" character to separate the password and file name.
4. This password and file name are stored using the `TryAddPassword()` method of the `_passwordRepository` object.

These steps describe the process of indexing by taking password and file name pairs from a given Dictionary and then reading this index.

## 5. Search Function and Performance Testing

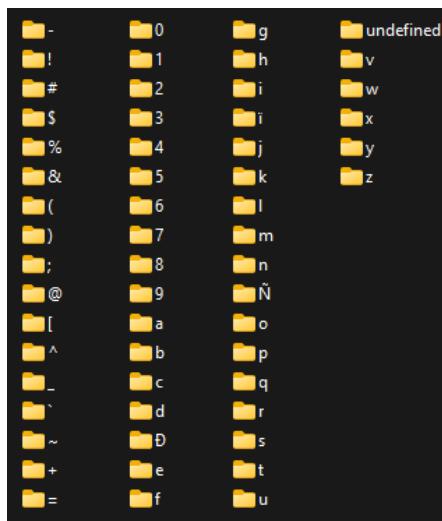
Searching with the `Dictionary<TKey, TValue>` data type has a time complexity very close to O(1) because it employs a hash table implementation. Therefore, performing searches in this data structure will yield very fast results.

**The indexing process takes 00:00:05.75 seconds**

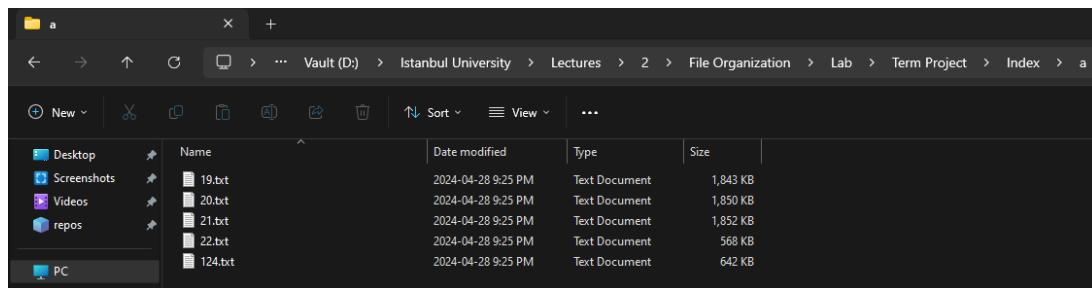


```
Program.cs* ✘ x
SecureList.Console
1 1  using SecureList.Core;
2 2  using SecureList.Core.Entities;
3 3  using System.Diagnostics;
4 4
5 // Set the encoding to UTF-8 for proper IO.
6 Console.InputEncoding = System.Text.Encoding.UTF8;
7 Console.OutputEncoding = System.Text.Encoding.UTF8;
8
9 Stopwatch stopwatch = new Stopwatch();
10 stopwatch.Start();
11 Console.WriteLine("① Indeksleme işlemi başladı");
12
13
14 var passwordRepository = PasswordRepository.Instance;
15 var passwordsDirectory = new PasswordsDirectory();
16 var passwordFiles = passwordsDirectory.Files;
17 var indexService = new IndexService(passwordRepository);
18
19 foreach (var passwordFile in passwordFiles)
20 {
21     if (!passwordFile.IsProcessed())
22     {
23         passwordFile.Process();
24     }
25 }
26
27 indexService.Index();
28 indexService.ReadIndex();
29
30 // Get the elapsed time as a TimeSpan value.
31 TimeSpan ts = stopwatch.Elapsed;
32
33 // Format and display the TimeSpan value.
34 string elapsedTime = string.Format("{0:00}:{1:00}:{2:00}.{3:00}",
35     ts.Hours, ts.Minutes, ts.Seconds,
36     ts.Milliseconds / 10);
37 Console.WriteLine("② Indeksleme işlemi bitti geçen süre: " + elapsedTime);
38
39 string? userInput;
40 do
41 {
42 }
```

After indexing, the contents of the Index folder are as follows:



After indexing, the folder containing passwords starting with **a** or **A** within the Index directory.



After the indexing process, here are 10 randomly searched passwords and their retrieval times:

Searched Password	Retrieval Time (milliseconds)
[windows	0.0558
iloveyou	0.0346
#bigguy	0.0267
İçİçİçİçİçİç	0.0181
%@xprs@%*	0.0224
shadow	0.0225
26101900	0.0406
exit	0.0355
@YVe9uZegYVa#a	0.0348
Ñ,ÐµÐ»ÐµÑ,,Ð¾Ð½	0.0319
Average: 0.0323	

```
Arama için şifre girin: [windows
Şifre [[windows] "10-million-password-list-top-1000000.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0558 milisaniye

Arama için şifre girin: iloveyou
Şifre [iloveyou] "10-million-password-list-top-100.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0346 milisaniye

Arama için şifre girin: #bigguy
Şifre [#bigguy] "medical-devices.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0267 milisaniye

Arama için şifre girin: içiçiçiçiçiçiçiç
Şifre [içiçiçiçiçiçiç] "100k-most-used-passwords-NCSC.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0181 milisaniye

Arama için şifre girin: %@xprs@%*
Şifre [%@xprs@%*] "10-million-password-list-top-1000000.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0224 milisaniye

Arama için şifre girin: shadow
Şifre [shadow] "10-million-password-list-top-100.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0225 milisaniye

Arama için şifre girin: 26101900
Şifre [26101900] "1900-2020.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0406 milisaniye

Arama için şifre girin: exit
Şifre [exit] "10-million-password-list-top-1000000.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0355 milisaniye

Arama için şifre girin: @YVe9uZegYVa#
Şifre [@YVe9uZegYVa#] "10-million-password-list-top-1000000.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0348 milisaniye

Arama için şifre girin: Ñ,ÐµÐ»ÐµÑ,,Ð¾Ð½
Şifre [Ñ,ÐµÐ»ÐµÑ,,Ð¾Ð½] "100k-most-used-passwords-NCSC.txt" konumunda bulundu.
• Arama işlemi bitti geçen süre: 0.0319 milisaniye
```

## 6. Results and Evaluation

 This project was quite educational because it contained example information about C# LINQ methods, anonymous types, and data structures.

 It was helpful to see the speed of the *Dictionary* structure and to better understand the importance of indexing for large datasets. From this perspective, the *Dictionary* structure can certainly be used in search functions in future projects.

 It was a project that demonstrated how a comprehensive project can be broken down into small parts and methods, making it easier to handle each function with object-oriented programming.

Looking at the project from a broad perspective, it touched on different topics, such as the function of determining whether files can be created with invalid characters in the Windows file system. Therefore, the project could be revisited for reference in the future.

## 7. References

The reference site and C# documentation used were [learn.microsoft.com](https://learn.microsoft.com).

1. <https://learn.microsoft.com/en-us/dotnet/csharp/fundamentals/types/anonymous-types>
2. <https://stackoverflow.com/questions/5750203/how-to-write-unicode-characters-to-the-console>
3. <https://learn.microsoft.com/en-us/dotnet/api/system.diagnostics.stopwatch.elapsed?view=net-8.0#examples>