**V.V. Voitsekhovskyi**
**I.V.Bogach**
**A.Y. Baranovska**

# THE SYSTEM OF INDEXING DOCUMENTS AND SEARCHING FOR THEIR INDEXES STORED IN DATABASE WRITTEN IN JAVA USING GRADLE BUILD AUTOMATION TOOL

Вінницький національний технічний університет

*Анотація:*

*У даній доповіді розглянуто систему індексування документів та їх пошуку по збереженим в базу даних індексам написану на мові Java використовуючи інструмент автоматичного збирання проектів Gradle та особливості її використання.*

*Ключові слова: програмування, програмне забезпечення, індексація, бази даних, автоматична збірка проєктів, Java, Gradle.*

*Abstract:*

*This article examines the system of indexing documents and searching for their indexes stored in database written in Java language using Gradle build automation tool and peculiarity of its use.*

*Keywords: programming, software engineering, indexing, databases, build automation tool, Java, Gradle.*

## Introduction

Indexing is the process of viewing content on your computer and classifying information about them, such as the words and metadata they contain.

This is helps to search for files fast; in addition, it takes up little space – less than 10 percent of the size of indexed files. This process used on many various operating systems and platforms in order to optimize a process of search [1].

## Methods of Indexing

There are two types of indexes – forward and inverted. Forward – consists in comparing the document with a list of words encountered in it. Inverted (named to contrast a forward type) compares the word is matched with a list of documents in which it is. It is logical to assume that an inverted index is best for fast searches.

Typically, search engines rank a list of documents containing queries after using an inverted index to list documents from the query. An inverted index is the most popular data structure used in information retrieval.

There are also two variants of the inverted index:
1.    Index that contains only a list of documents for each word;
2.    Index, additionally includes the position of the word in each document [6].

## MySQL and formats of files

MySQL is highly effective DBMS that is fast and scalable. It has a number of advantages such as multithreading, support for multiple simultaneous requests; optimization of connections with joining many data in one pass; fixed and variable length records; flexible support for number formats, variable length strings and timestamps. All of these aspects is very important to create a quality indexing system so that is why it's used here [3].

Doc, docx and pdf formats are used because it considered as a documents standard (using widely) and it is easy to extract and read data from this types of files.

## Java Language and Gradle Tool

A Java programming language was used to create this system of indexing. Java has many advantages over other programming languages, which allows you to solve almost any problem with it, for example:

Java is an object-oriented language. This allows you to create modular programs whose source code that can be used repeatedly;

Java is easy to learn comparing to other OOP languages;

One of the main advantages of the Java language is the ability to transfer programs from one system to another. This language is platform independent

The disadvantage is that compared to other languages, Java is rather slow in execution [4].

Gradle is an open-source tool that helps us to create software with mechanization. This tool is widely used for the creation of different kinds of software due to its high performance. It works on Java and a Groovy-based Domain-Specific Language (DSL) for developing the project structure. Gradle supports the creation of mobile and web applications with testing and deploying on various platforms.

Main benefits of Gradle [2]:

This tool is highly customizable as it supports a variety of IDE's. It avoids compilation.

Performance is better than in Maven and other build tools.

More flexible than other tools. You even can write some code for configuring Gradle in a script. A lot of different plugins and platforms which support Gradle

Usability in different version control systems such as GitHub

Open-source tool, so it is free.

Some disadvantages:

More difficult and complex than other build automation tool

Has own script language, not like XML using in Maven. So it can be difficult to understand it without learning documentation

Maven has more examples and more dependencies than Gradle [6].

Graddle script listing:

```
plugins {
    id 'java'
    id 'maven-publish'
}

repositories {
    mavenLocal()
    maven {
        url = uri('https://repo.maven.apache.org/maven2/')
    }
}

dependencies {
    implementation 'org.apache.poi:poi-ooxml:3.16'
    implementation 'org.apache.poi:poi-scratchpad:3.16'
    implementation 'org.apache.pdfbox:pdfbox:2.0.19'
    implementation 'mysql:mysql-connector-java:8.0.18'
}

group = 'org.example'
version = '1.0-SNAPSHOT'
description = 'TestTask'
java.sourceCompatibility = JavaVersion.VERSION_1_8
```

```
publishing {
    publications {
        maven(MavenPublication) {
            from(components.java)
        }
    }
}

tasks.withType(JavaCompile) {
    options.encoding = 'UTF-8'
}
```

File reader method for reading pdf files example code Java:

```java
public static String readPdfFile(String filePath){
        try {
            PDDocument document = PDDocument.load(new File(filePath));
            if (!document.isEncrypted()) {
                PDFTextStripper stripper = new PDFTextStripper();
                return stripper.getText(document);
            }
            document.close();
        } catch (IOException e) {
            e.printStackTrace();
        }
        return "";
    }
```

Example method of searching file by searched word Java:

```java
public DefaultListModel<String> getFileListIncludesSearchingWord(String
wordForSearch ) {
        DefaultListModel<String> fileList = new DefaultListModel<>();

        try {
            Connection connection = DriverManager.getConnection(connectionUrl,
userName, password);
            Statement statement = connection.createStatement();

            String query = "select file_name, file_type from Files where
file_content like " + "\"%" + wordForSearch + "%\"";
            ResultSet rs = statement.executeQuery(query);

            while (rs.next())
                fileList.addElement(rs.getString("file_name"));

            connection.close();

            return fileList;
        }
        catch (SQLException throwables) {
            throwables.printStackTrace();
            return null;
        }
    }
```
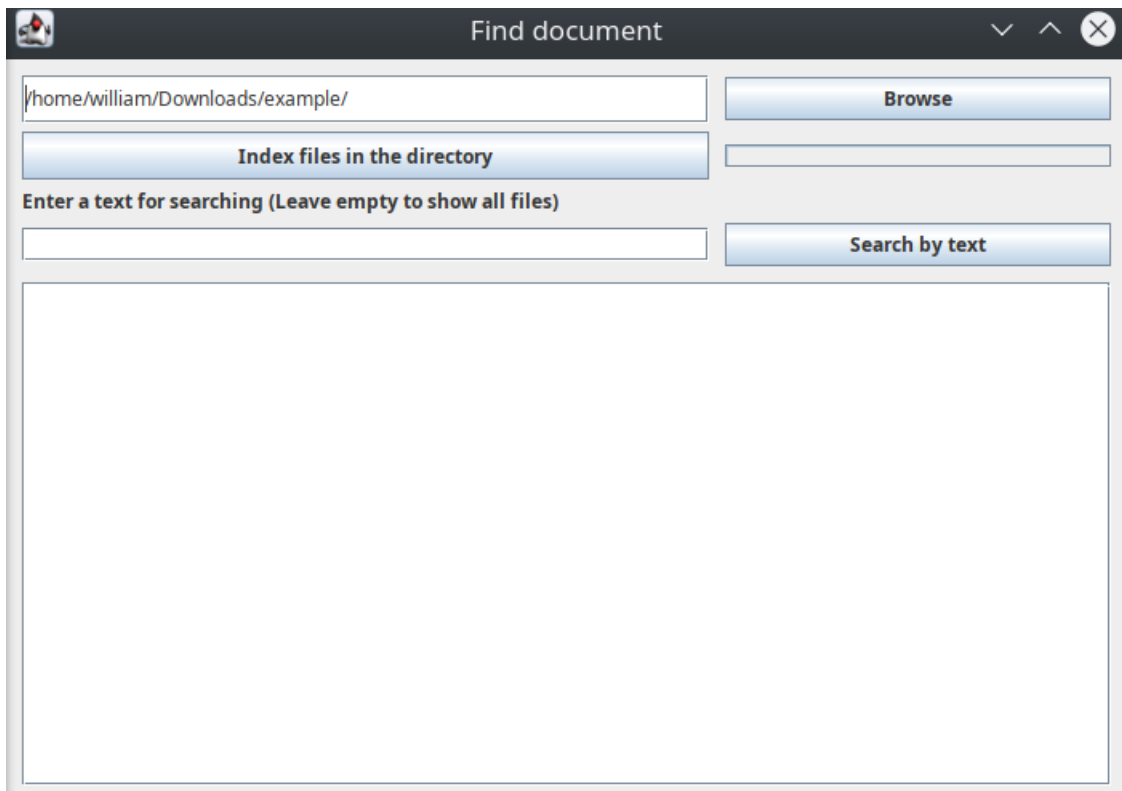
**Using the system of indexing**



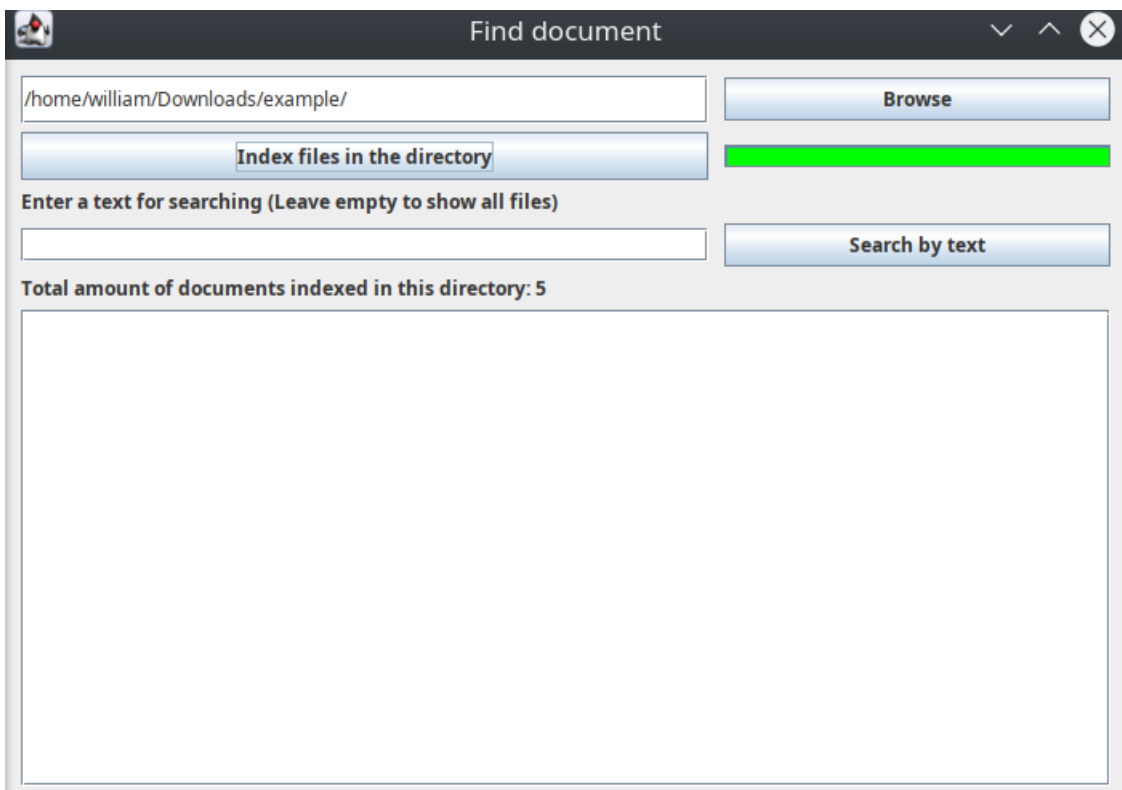Figure 1 – General view of the program



Figure 2 – A specific folder is selected and a certain number
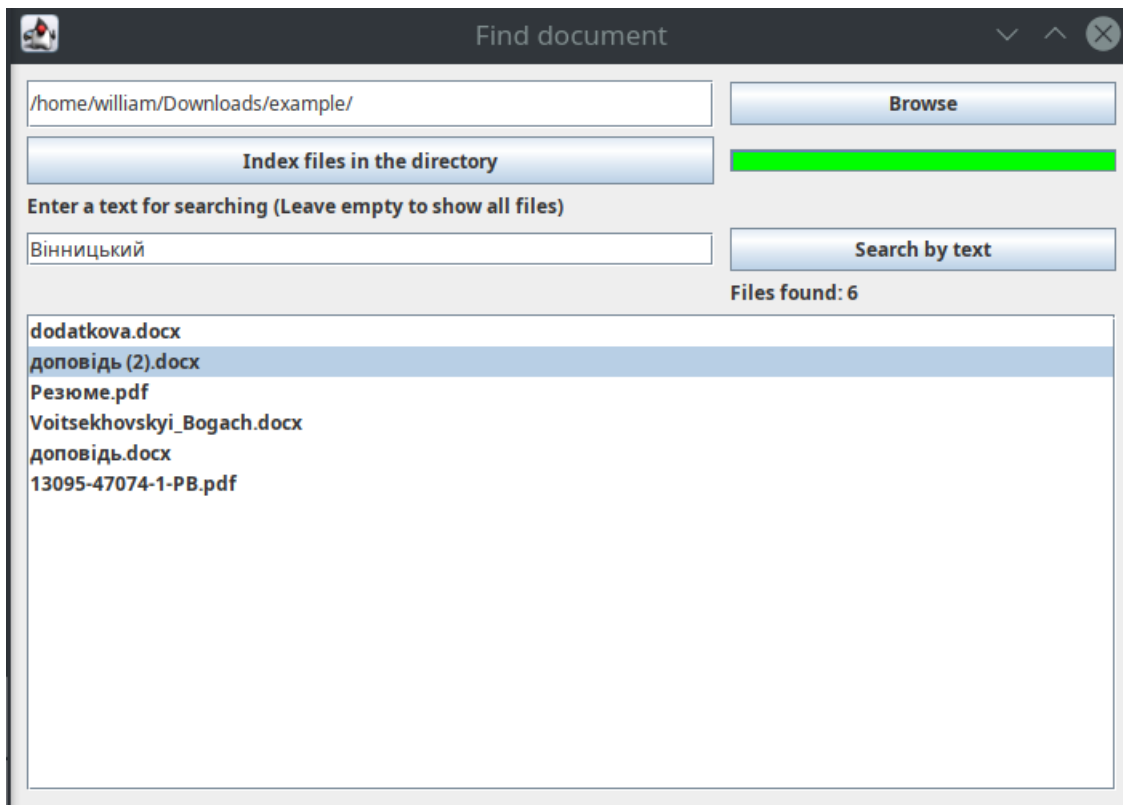of documents that fit the file type are indexed

Figure 3 – Showing files that match the search query.
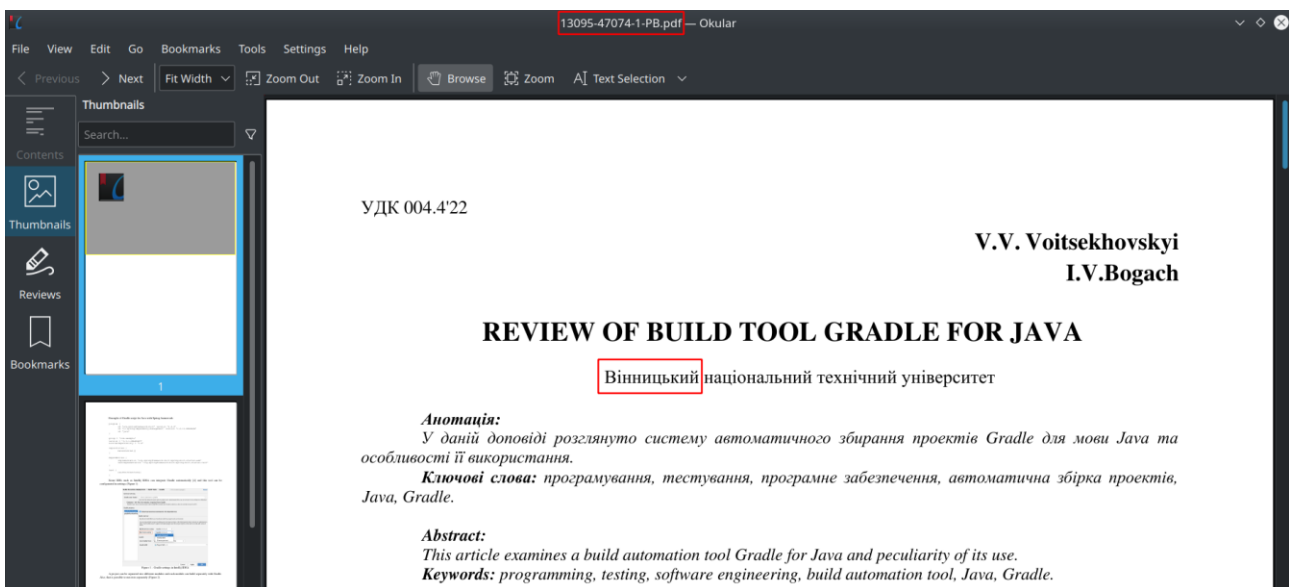Open the file selected in the screenshot



Figure 4 – The specified file contains the word from the search entry

Figure 4 – As we can see the program ignored files with the incorrect data type

**Conclusion**

In this article, we reviewed and described the system of indexing documents and searching for their indexes stored in database written in Java language using Gradle build automation tool. Also summed up pros and cons of every technology. We added some code examples, configurations, and use. In conclusion, we can say that we used most proper tools and algorithms to create a fast and accurate system for both indexing and search.

## СПИСОК ВИКОРИСТАНОЇ ЛІТЕРАТУРИ

1. Search indexing in Windows [Електронний ресурс] : [Веб-сайт]. – Режим доступу: https://support.microsoft.com/en-us/windows/search-indexing-in-windows-10-faq-da061c83-af6b-095c-0f7a-4dfecda4d15a . – Назва з екрана.
2. Gradle. Wikipedia [Електронний ресурс] : [Веб-сайт]. – Режим доступу: https://en.wikipedia.org/wiki/Gradle. – Назва з екрана.
3. 8 Advantages using MySQL [Електронний ресурс] : [Веб-сайт]. – https://devops.com/8-advantages-using-mysql/. – Назва з екрана.
4. Advantages of Java [Електронний ресурс] : [Веб-сайт]. – Режим доступу: https://www.ibm.com/docs/en/aix/7.1?topic=monitoring-advantages-java .– Назва з екрана.
5. Gradle User Manual [Електронний ресурс] : [Веб-сайт]. – Режим доступу: https://docs.gradle.org/current/userguide/userguide.html. – Назва з екрана.
6. Inverted index, Wikipedia [Електронний ресурс] : [Веб-сайт]. – Режим доступу: https://en.wikipedia.org/wiki/Inverted_index . – Назва з екрана.

*Войцеховський Вільям Вільямович* – *студент групи 1АКІТ-18Б, кафедра автоматизації та інтелектуальних інформаційних технологій, Факультет комп'ютерних систем і автоматики, Вінницький національний технічний університет, м.Вінниця, e-mail: fkca.1akit18.VVV@gmail.com*

*Богач Ілона Віталіївна* – *к.т.н., доцент кафедри Автоматизації та інтелектуальних інформаційних технологій, Вінницький національний технічний університет, м.Вінниця, e-mail: ilona.bogach@gmail.com*

*Барановська Анастасія Юріївна* – *студентка групи 1ІСТ-19Б, кафедра автоматизації та інтелектуальних інформаційних технологій, Факультет комп'ютерних систем і автоматики, Вінницький національний технічний університет, м.Вінниця, e-mail: 01-19-051.stud@vntu.edu.ua*

*Voitsekhovskyi Viliam Viliyamovych* – *student of 1AKIT-18B group, Department of Automatization and Intellectual Informational Technologies, Faculty of Computer Systems and Automatics, Vinnytsia National Technical University, Vinnytsia, e-mail: fkca.1akit18.VVV@gmail.com*

*Bogach Ilona Vitaliivna* - *Associate Professor of Automation and Intelligent Information Technologies, Vinnytsia National Technical University, Vinnytsia, e-mail: ilona.bogach@gmail.com*

*Baranovska Anastasiia Yuriivna* – *student of 1IST-19B group, Department of Automatization and Intellectual Informational Technologies, Faculty of Computer Systems and Automatics, Vinnytsia National Technical University, Vinnytsia, e-mail: 01-19-051.stud@vntu.edu.ua*