# Documentation of Easy-Compare

**Introduction:**

Easy Compare is a document comparison application that helps to detect plagiarism. Every document follows a similar text preprocessing approach, every sentence in the document gets hashed using Murmur3 hashing technique and this hash is stored in PostgreSQL. User can now select two documents and get a plagiarism score which is calculated using Jaccard index.

**The Jaccard index**, also known as Intersection over Union and the Jaccard similarity coefficient  is a statistic used for comparing the similarity and diversity of sample sets. The Jaccard coefficient measures similarity between finite sample sets, and is defined as the size of the intersection divided by the size of the union of the sample sets.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}.$$
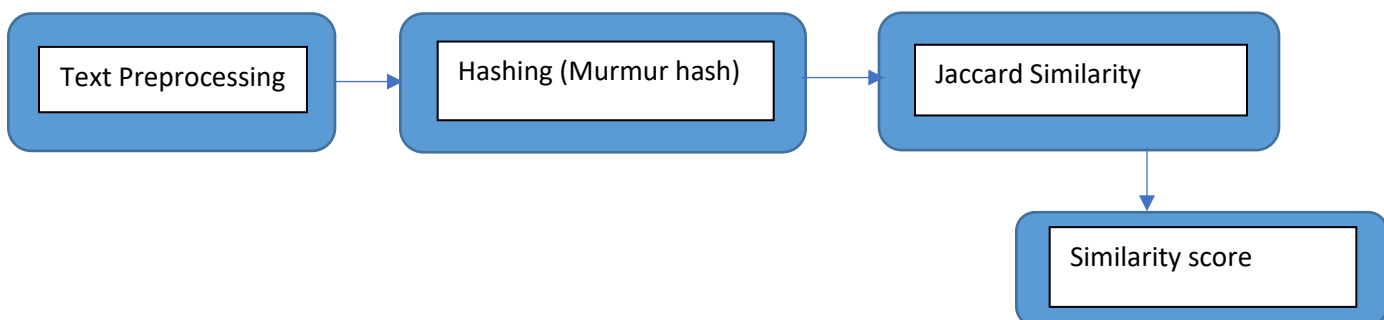
**Murmur Hash 3**:

Murmur Hash is a non-cryptographic hash function suitable for general hash-based lookup. The name comes from two basic operations, multiply (MU) and rotate (R), used in its inner loop. https://en.wikipedia.org/wiki/MurmurHash#MurmurHash3

The current version is MurmurHash3 which yields a 32-bit or 128-bit hash value. When using 128-bits, the x86 and x64 versions do not produce the same values, as the algorithms are optimized for their respective platforms. MurmurHash3 was released alongside SMHasher--a hash function test suite.

**Text Preprocessing**:

Text preprocessing function reads the input file and converts all the text to lowercase and combines every line of the document into a string and removes all the whitespaces, each sentenced is hash and stored in PostgreSQL.

**Block Diagram:**

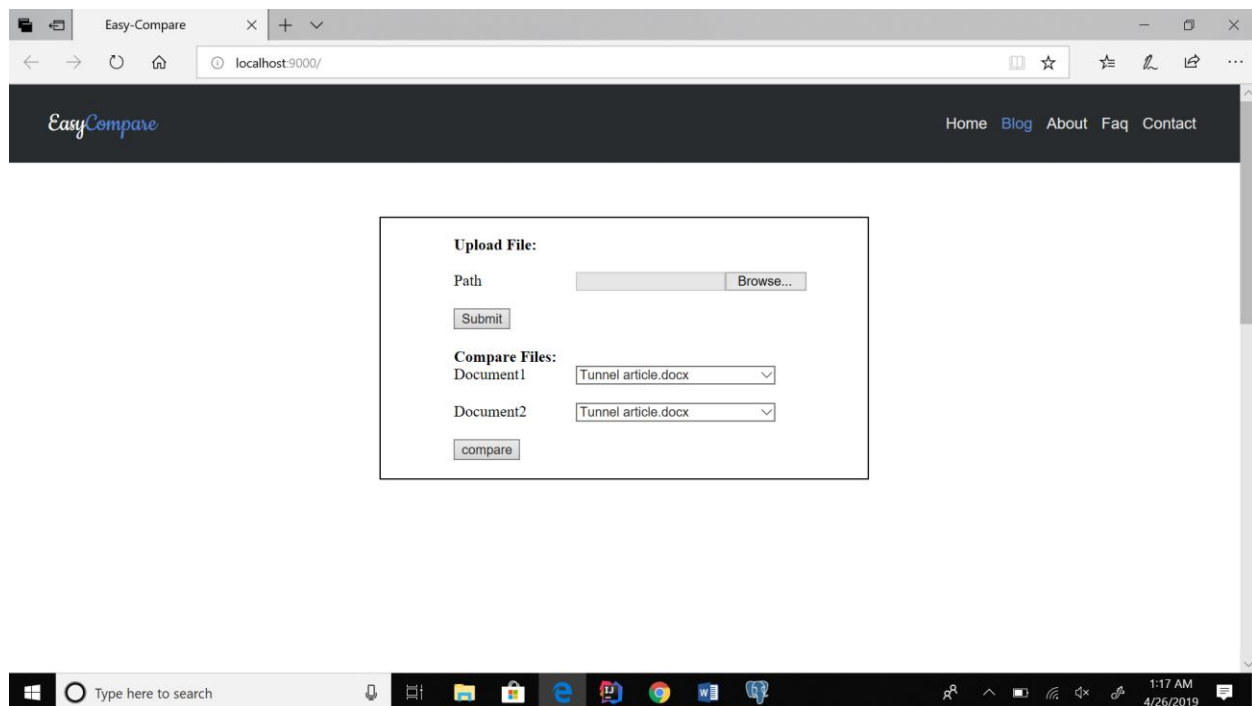Text Preprocessing → Hashing (Murmur hash) → Jaccard Similarity → Similarity score

**How to run and use Easy Compare on your machine:**

1. Install PostgreSQL and pgadmin 4
2. Create a table in PostgreSQL using the script given below

   CREATE TABLE public."tblDocumentInformation"
   (
      "ID" integer NOT NULL,
      "Document_Name" character varying COLLATE pg_catalog."default",
      "Document_Text" character varying COLLATE pg_catalog."default",
      "Sentence_Text" character varying COLLATE pg_catalog."default"
   )
3. Git clone https://github.com/BaraiKaran/Easy-Compare.git
4. Open the play project in your preferred IDE and in the shell type 'sbt run' and start the play application.
5. Once the application is running click on browse and select the file to upload.
6. Once you have uploaded the files, in dropdown below select the two files which you want to compare.
7. When everything is set press compare and obtain plagiarism score.

**Model, View and Controller pattern:**

**Controller –**

| Endpoint | Description |
|----------|-------------|
| / | Initial GET request to load a form. |
| /simpleForm | POST request to store document content in database. |
| /compare | POST request to compare two documents. |

| Services | Description |
|----------|-------------|
| BasicForm | Code for form Mapping. |
| Comparison | Code for comparison between two documents using Jaccard similarity. |
| Preprocess | Code for initial preprocessing of the text. |
| Validation | Code for validation of the user input. |

**Models –**

| Models | Description |
|--------|-------------|
| InteractionWithDb | Contains slick code for interaction with PostgreSQL. |

**Unit Tests –**

Code has 60% coverage i.e. 60% of the code is unit tested currently. There are few integration tests as well which takes coverage to about 65%. Unit tests can be found in "test" directory.
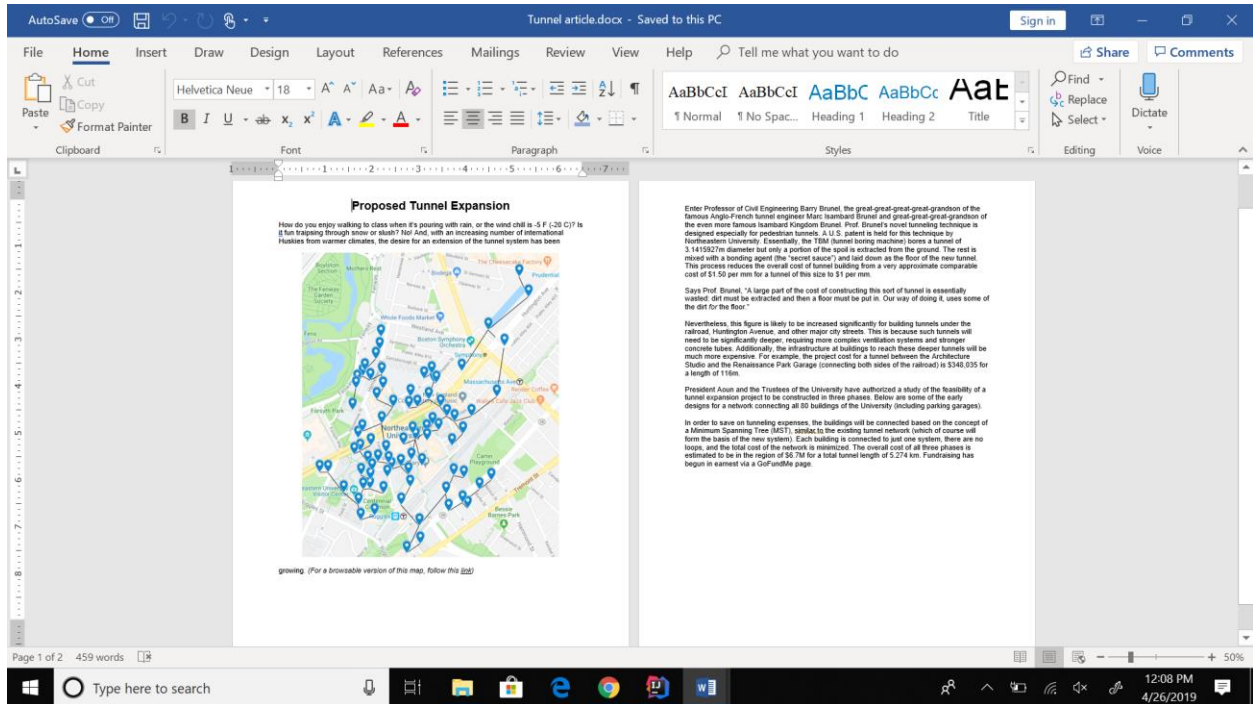
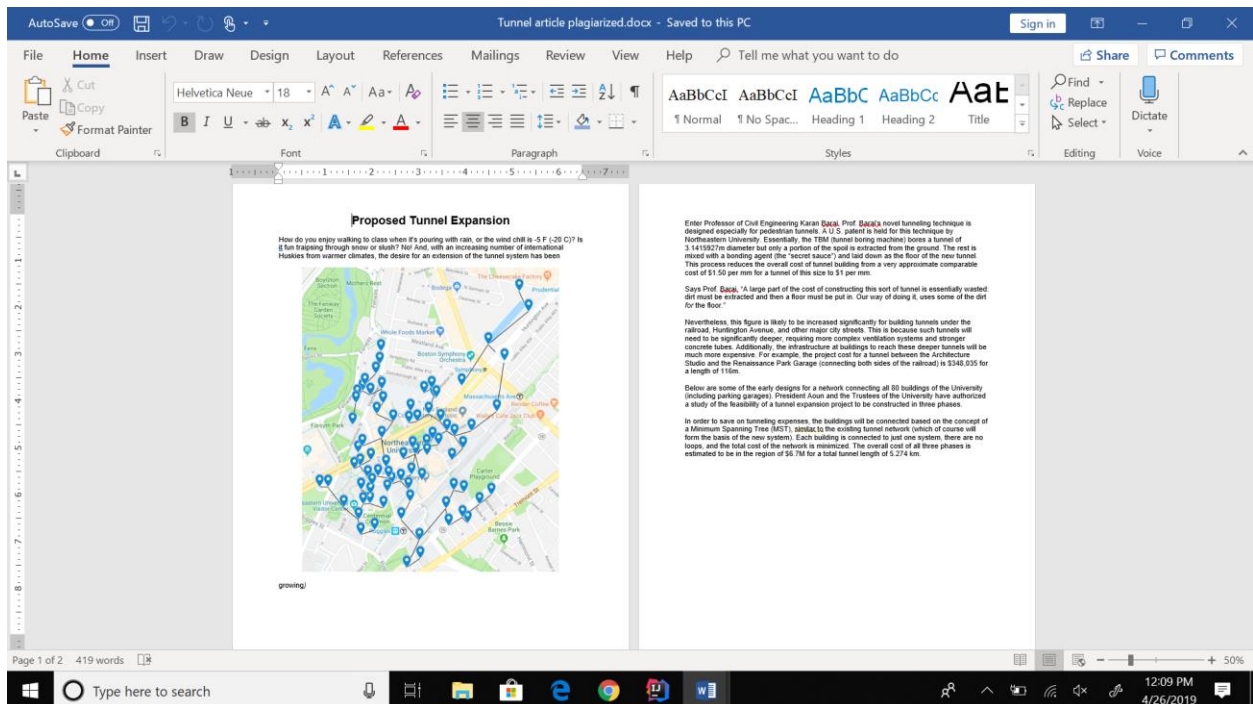| Test File | Unit tests |
|-----------|------------|
| preprocessSpec | Covers unit tests for data preprocessing, Functions tested:<br>1. convertListToString<br>2. removeWhiteSpaces, removeMultipleWhiteSpaces<br>3. splitText<br>4. getFileName<br>5. readTextFile<br>6. apply |
| ComparisonSpec | Compares contents of two documents, functions tested:<br>1. calculateSimilarityScore<br>2. hashContentsOfList, hashContentsOfListNotEqual, hashContentsOfEmptyList<br>3. getListFromOption<br>4. flattenFuture<br>5. getSimilarityScore |
| ValidationSpec | Function tested:<br>1. getFileType |
| InteractionWithDbSpec | Functions tested:<br>1. Sequence, SequenceWhenGivenNull |

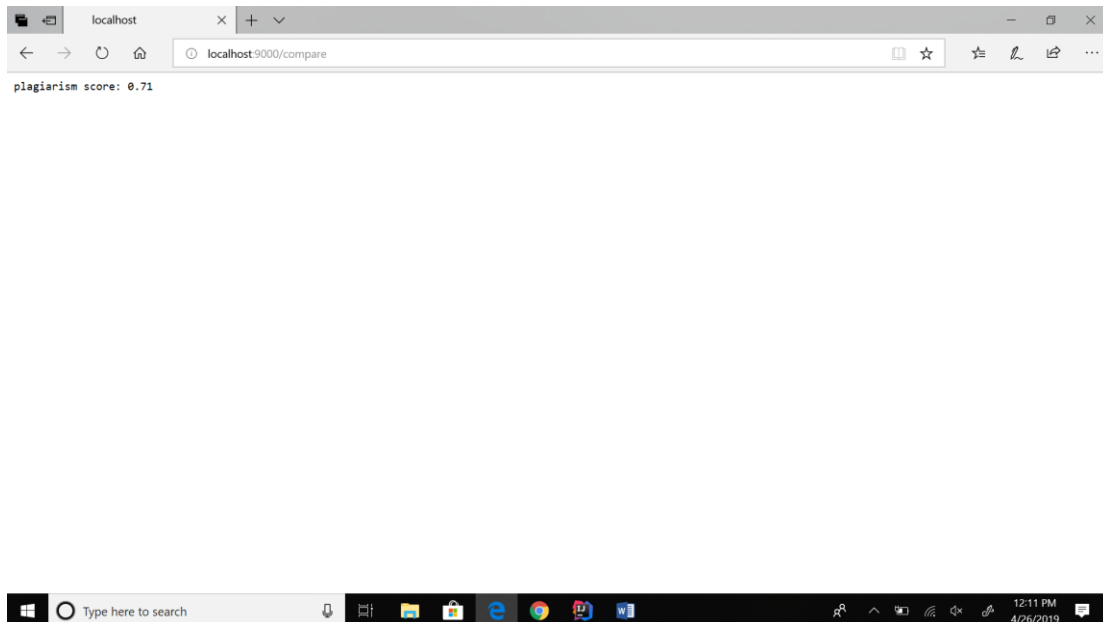| ControllerSpec (Integeration) | This file tests the initial GET controller ("/") which brings up the form in user interface. |
| --- | --- |

## Demo of the application:

### Document 1 -



### Document 2 –

**Plagiarism score:**

The application returns a score of 71%, that means about 71% of the code is plagiarized.



**Continuous Integration:**

To test impact of every commit I have used an integration tool CircleCI that will build the application and run the unit tests to make sure things are working fine after every commit in the git repository.

**Codacy:**

To keep check on the code quality throughout the entire application I have used an automated code review tool Codacy. This tool checks the code quality after every commit or pull request and gives feedback on ways to improve the code. Currently, Easy Compare gets an "A" in terms of code quality.

**Scaladoc** –

Entire code is well documented using scaladoc

**Technical Specs:**

1. Scala 2.12.2
2. Play 2.7
3. PostgreSQL 9.6.5
4. pgAdmin 4
5. Scala Test 3.0.5
6. CircleCI
7. Git and GitHub

**Future Scope:**

1. Include abstract syntax trees for comparison of code and document.
2. Dockerize the application for ease of installation in any machine.

**References:**

1. Jaccard_index - https://en.wikipedia.org/wiki/Jaccard_index
2. Murmur hash - https://en.wikipedia.org/wiki/MurmurHash
3. Rabin Karp - https://en.wikipedia.org/wiki/Rabin%E2%80%93Karp_algorithm

**License:**

1. MIT License