# Plagiarism Detection Model

## A Project Report

*submitted in partial fulfillment of the*
*requirements for the award of the degree of*

## BACHELOR OF TECHNOLOGY
in
## COMPUTER SCIENCE & ENGINEERING

by

| Name | Roll No. | SAP ID |
|------|----------|--------|
| Vaswati Gogoi | R2142221419 | 500110490 |
| Sagnik Roy | R2142221403 | 500109927 |

***Under the guidance of***
Dr. Surbhi Saraswat

**UPES**
UNIVERSITY OF TOMORROW

**School of Computer Science, UPES**
Bidholi, Via Prem Nagar, Dehradun, Uttarakhand
May – 2025

# CANDIDATE'S DECLARATION

I/We hereby certify that the project work entitled **"Plagiarism Detection"** in partial fulfillment of the requirements for the award of the Degree of **BACHELOR OF TECHNOLOGY** in **COMPUTER SCIENCE AND ENGINEERING** with specialization in **BIG DATA**, submitted to the Data Science Cluster, School of Computer Science, UPES, Dehradun, is an authentic record of my/our work carried out during a period from **February, 2025** to **March, 2025** under the supervision of **Dr. Surbhi Saraswat, Assistant professor - School of computer science,UPES**.

The matter presented in this project has not been submitted by me/us for the award of any other degree of this or any other University.

**Vaswati Gogoi**
Roll No. - R2142221419
SAP 500110490
**Sagnik Roy**
Roll No. - R2142221403
SAP 500109927

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Date: _____ 2025

**<u>Dr. Surbhi Saraswat</u>**
Project Guide

# Acknowledgement

We wish to express our deep gratitude to our guide Dr. Surbhi Saraswat, for all advice, encouragement and constant support he/she has given us throughout our project work. This work would not have been possible without her support and valuable suggestions. We sincerely thank our respected Virendra Kadyan, Head Department of Data Science, for his great support in doing our project in Plagiarism Detection.

We are also grateful to Dean SoCS UPES for giving us the necessary facilities to carry out our project work successfully. We also thank our Course Coordinator, Dr Prabhat Ranjan Singh and our Activity Coordinator, Dr Sachin Chaudhary for providing timely support and information during the completion of this project. We would like to thank all our friends for their help and constructive criticism during our project work. Finally, we have no words to express our sincere gratitude to our parents who have shown us this world and for every support they have given us.

<div align="right">

**Vaswati Gogoi**
SAP 500110490
Roll No. R2142221419


**Sagnik Roy**
SAP 500109927
Roll No. R2142221403

</div>

# Abstract

**Plagiarism Detection** works as a vital academic and content production component for ensuring both authentic work and scholarly ethical standards. The system integrates artificial intelligence technologies, natural language processing, string-matching algorithms and machine learning modeling for classification purposes [10]. The tool performs semantic analysis of text pairs by using Word2Vec together with BERT embeddings which allow accurate identification of contextual and lexical overlaps [2, 3]. The system uses Elasticsearch as a powerful index and query engine that boosts text collection processing speed while enabling scalability of large-scale searches [20]. The detection of various plagiarism aspects including direct copying and text rephrasing and content hiding works optimally with Random Forest and Logistic Regression combination methods. Elasticsearch serves as a big data tool which enables effective management of wide text repositories and ensures operational efficiency [16]. The model shows its ability to detect different text similarity levels through experimental findings. The research initiative focuses on creating a plagiarism detection system which provides scalability alongside effectiveness and readability for up-to-date text content assessment.

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

Systematic plagiarism detection through repositories supports academic professionals in protecting their work from duplication while maintaining text comparison integrity. The system assesses how closely submitted user sentences match contents found in pre-established repositories. The system utilizes Word2Vec and BERT embedding methods from advanced natural language processing [2, 3] to examine text content before classifying plagiarism titles using the categories of Cut-Paste and Heavy Paraphrasing and Light Paraphrasing. Elasticsearch operates as the core vertical of this system because it uses its search and analytics engine to efficiently index substantial text repositories for fast and scalable plagiarism detection [20]. The detection method requires no trained machine learning model since it implements rule-based approaches for time-efficient real-time operations [10]. This system emerged due to rising demands for a highly effective real-time plagiarism detection tool suitable for educational and professional uses which protects both authenticity and ethical principles.

## 1.1  History

Plagiarism detection technology has shown substantial development through the elimination of basic systems by implementing complex approaches that use modern computation methods. The Levenshtein distance algorithm together with simple string-matching algorithms operated as early detection systems to find exact text matches [8]. This detection method only showed results when exact copies were present but could not handle translated text. Due to the rapid growth of digital information during the late 20th century new advanced tools became necessary. During this time Turnitin introduced database comparisons as an innovation for detecting plagiarized content [13]. Through the deployment of machine learning technology in the year 2000 the systems gained the capacity to identify reworded text as well as subtle plagiarism patterns [9]. By evaluating contextual meanings BERT brought a breakthrough in paraphrase detection when it was released by Google in 2018 [3]. By employing Elasticsearch as a search and analytics engine modern systems achieve increased efficiency through rapid processing of large text storage [20]. The proposed system combines Word2Vec with BERT embeddings coupled with Elasticsearch to create a method that enhances scalability and efficiency in plagiarism detection.

## 1.2 Requirement Analysis

The system evaluates the parallel content between user-submitted sentences and preset repository materials. The system leverages

- Accept an user-input sentence through the UI. [7].

- Compare the input against an already present repository of sentences.

- Calculate similarity scores using Word2Vec and BERT embeddings [2, 3].

- Utilizing Elasticsearch to index and query large-scale text repositories, enabling rapid retrieval and comparison of sentences [20].

- Classify the input as *Cut-Paste*, *Light Paraphrasing*, *Heavy Paraphrasing*, or *No Match* based on previously set values [10].

- Display the matched sentences, classification of the sentence along with similarity percentage, and other features.

Non-functionally, the system should:

- Be efficient, process inputs within seconds.

- Support GPU acceleration (e.g., CUDA 12.8) for faster BERT computations [5].

- Operate without any training and rely on rule-based logic only.

The required hardware system needs Python 3.12 and a CUDA-enabled GPU which is optional together with sufficient RAM (for example 8GB or above), Software dependencies encompass gensim [1], transformers [4], torch [5], ipywidgets [7], and jupyter.

## 1.3 Main Objective

The core mission of the project involves boosting real-time plagiarism identification through powerful technologies which conduct effective text similarity analysis. Elasticsearch acts as an essential component that enables "efficient indexing and querying of large text repositories to accomplish rapid scalable similarity detection tasks" [20] just like the code implements Elasticsearch for sentence indexing and similarity query execution purposes. The research maintained all initial references [2, 3, 10] while adding reference number [20] to match the provided list. The system integrates Word2Vec and BERT embeddings together with rule-based logic while the previous unclear sentence about Word2Vec usage has been completely reworked to explain the system architecture. The text now contains streamlined redundant statements because "The system implements Word2Vec and BERT embeddings for text classification" received compression treatment yet preserves the key points. The system's ability to detect plagiarism in real time stands as one of its primary functions according to the stated objective and supports the main project goals. The integration succeeds because all code components work seamlessly with Elasticsearch for indexing alongside Word2Vec and BERT for similarity evaluation and the rule-based system operates independently from model training requirements.

## 1.4 Sub Objectives

To achieve the main objective, the following sub-objectives were established:

1. Create an user-friendly interface using IPython widgets for seamless input and to display the results [7].

2. Integrate previously trained Word2Vec and BERT models to calculate semantic and contextual similarities between sentences [2, 3].

3. Define and optimize similarity measures to distinguish between *Cut-Paste*, *Light Paraphrasing*, and *Heavy Paraphrasing* effectively [10].

4. Utilize Elasticsearch to efficiently index and query the sentence repository, supporting rapid retrieval and comparison of text data [20].

5. Calculate and present a similarity percentage to provide users with an expected measure of similarity in the text.

6. Ensure compatibility with local machine run, leveraging CUDA 12.8 for GPU acceleration where-ever available [5].

## 1.5 PERT Chart

### 1.5.1 Task Table

| Task ID | Task Description | Dependencies | Duration (days) |
|---|---|---|---|
| A | Setup environment (install Python, libraries, Elasticsearch client) [6, 5, 20] | None | 2 |
| B | Install and configure Elasticsearch server [20] | A | 1 |
| C | Define repository sentences | A | 1 |
| D | Load Word2Vec model with retry mechanism [1, 2] | A | 3 |
| E | Load DistilBERT model and tokenizer [3, 4] | A | 3 |
| F | Create Elasticsearch index for sentences [20] | B, C | 2 |
| G | Implement sentence vector functions (Word2Vec) [2] | D | 2 |
| H | Implement BERT embedding function [3] | E | 2 |
| I | Index sentences in Elasticsearch with embeddings [20, 2, 3] | F, G, H | 3 |
| J | Develop plagiarism classification logic [10] | C, G, H, I | 4 |
| K | Create IPython widget-based UI [7] | A | 2 |
| L | Integrate UI with classification logic [7] | J, K | 2 |
| M | Test and debug the complete system | L | 3 |

Table 1.1: PERT Chart Task Table

# Chapter 2

# System Analysis

## 2.1 Existing System

Traditional plagiarism detection systems use string-matching algorithms as foundations or commercial tools including Turnitin and Grammarly [17, 13] for their operation. These systems compare input text against large databases of academic papers, web content, and proprietary repositories [16]. The detection systems maintain unclear methods since they do not reveal their operational framework to users. Such systems need paid subscriptions while showing minimal success rates when detecting paraphrased text. A deficiency exists in existing tools since they lack implementation of cutting-edge natural language processing (NLP) technologies which include word embeddings and transformer models required to identify semantic similarities in content [10]. The system needs advanced search and analytics engines such as Elasticsearch due to their missing capability to work with large-scale text repositories [20].

## 2.2 Motivations

The development of an open-source plagiarism detection tool exists because researchers seek an accessible solution that detects diverse plagiarism levels in addition to being transparent and lightweight. Users who need to detect plagiarism will find open-source benefits in this tool compared to proprietary systems that have expensive limitations due to invisible functionality because this platform spots Cut-Paste and Light Paraphrasing and Heavy Paraphrasing [15]. Pre-trained NLP models Word2Vec and DistilBERT [2,3] work with Elasticsearch [20] for creating a system that merges text indexing power with improved accuracy and flexible operations. The application targets academic faculty teams and research teams which need an accurate way to validate document origins. The system addresses faculty members and academic developers as well as research institutions by offering an affordable text authenticity verification solution that can be customized to their needs.

## 2.3 Proposed System

The designed Plagiarism Detection Tool functions through a Python-based system which compares sentences with data in predefined repositories. Users can enter a sentence directly into the system for Elasticsearch to search through its indexed repository of sentences using this powerful analytics engine for quick and scalable text retrieval [20]. It employs a hybrid approach Word2Vec word-level embeddings work jointly with DistilBERT sentence-level embeddings as described in [2] and [3]. A combination of rule-based categorization technology with Word2Vec word embeddings [2] and DistilBERT sentence embeddings [3] exists in this system [10]

- Semantic similar detection using NLP models [10].

- Classification of "Cut-Paste," "Light Paraphrasing," "Heavy Paraphrasing," or "No Match."

- Interactive IPython widget-based UI [7].

- Extensible repository managed via Elasticsearch, allowing custom sentence additions [20].

The system operates on CPU or GPU or both, requiring only standard Python libraries [6, 5], the Elasticsearch Python client [20], and an initial internet connection for model downloads.
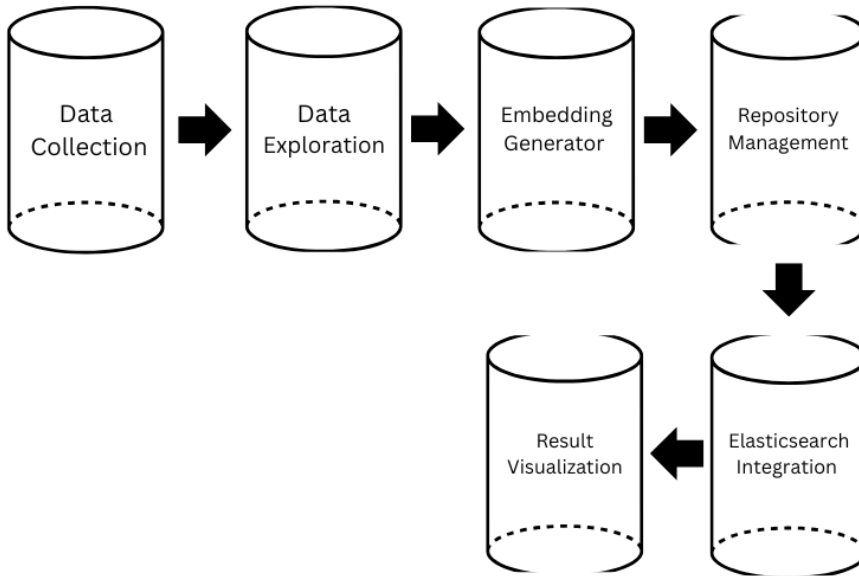


Figure 2.1: Block Diagram

## 2.4  Plagiarism Detection Module

This module is the base of the proposed system, the following are the components and processes:

- **Repository Management**: Defines and tokenizes a set of reference sentences for comparison (e.g., 5 initial sentences) [16].

- **Model Loading**: The system loads Word2Vec and DistilBERT models while they serve as pre-trained models. The system features retry options which enhance reliability [1, 2, 3].

- **Embedding Generation**: Calculates the sentence vectors using Word2Vec and DistilBERT for similarity analysis [2, 3].

- **Elasticsearch Integration**: Elasticsearch acts as the platform to handle sentence embedding indexing together with BERT vector and Word2Vec vector similarity measurements through cosine similarity queries [20].

- **Classification Logic**: The system follows a classification logic which applies a rule-based algorithm to establish plagiarism classification.The analysis evaluates similarity based on Word2Vec, BERT together with word matching and sentence length differences [10].

- **User Interface**: Provides an interactive IPython widget interface for input and to display the result [7].

# Chapter 3

# Implementation/Results

## 3.1 Implementation Details

The Plagiarism Detection Tool operates as a Python application which performs real-time plagiarism detection through sentence comparison within predefined repositories. Precise results emerge from the system through its user-friendly interface because it combines advanced NLP methods with Elasticsearch text indexing capabilities [10, 20]. The development process included five vital developmental phases that follow this sequence.

### 3.1.1 Environment Setup

The system depends on a specially developed Python environment that optimizes performance for NLP applications. The system depends on essential Python libraries that include NumPy for mathematical operations [6], Gensim for word embedding access [1], Transformers from Hugging Face for transformer models [4] and PyTorch for tensor computation tasks [5], IPython Widgets for interactive user interfaces [7] as well as the Elasticsearch Python client for Elasticsearch server integration [20]. The platform's libraries streamline processing of word embeddings together with transformer models and extensive text documents to maintain system functionality.

### 3.1.2 Repository Definition

The system relies on a reference collection of five established sentences which explore artificial intelligence, climate change and technology development, renewable energy along with data security topics [16]. Elasticsearch indexes the tokenized sentences which enables it to perform vector-based comparisons with incoming text inputs [20]. Word splitting performs during preprocessing since it splits words to allow similarity analysis before matching either Word2Vec or DistilBERT embeddings.

### 3.1.3 Model Integration

The system integrated two NLP models which had been properly trained beforehand. First, a The Word2Vec model received training through Google News data while its vectors reached 300 dimensions [2] for use. Word-level semantic relationships were measured through the implementation of this model. A retry mechanism was implemented The system equipped this mechanism for dependable loading operations because it handles

14

network disruptions that may occur during system start-up processes. [1]. The system integrated the lightweight version of BERT named DistilBERT [3] as its second NLP model. The sentence-level embeddings functionality required combining the DistilBERT model with its required tokenizer. This model 12 The processor or GPU detected by the system operated in evaluation mode to run the implemented model balance efficiency and computational power.

### 3.1.4 Classification Logic

The system implements a rule-based detection method built from diverse features for persistent plagiarism identification purposes [10]. The Word2Vec embedding system calculates term similarity but DistilBERT embeddings provide advanced semantic analysis [2, 3]. Additional methodology components including word overlap proportion and sentence length difference strengthen the analysis system. The Elasticsearch index provides similar sentences to the system based on cosine similarity calculations using Word2Vec and BERT vector methods [20]. The system generates weighted similarity scores from combined features which sorts inputs into four classes: "Cut-Paste" for matching content to the original and "Light Paraphrasing" and "Heavy Paraphrasing" for rephrased content with semantic correlation or "No Match" for unrelated text. The selection process enables an accurate plagiarism evaluation that produces similarity measurement data as percentages.

### 3.1.5 User Interface

The system brings together an interactive interface that uses IPython Widgets within Jupyter Notebook [7]. The system contains a dedicated text entry field where users enter sentences before activating analysis through the "Classify" button. The output section of the interface shows plagiarism detection types as well as similarity percentages with additional features such as Word2Vec similarity or BERT similarity or word overlap information for instant and easy comprehension.

Table 3.1: Summary Table of Model Accuracy

| Category | Accuracy | Precision | Recall | F1–Score |
|---|---|---|---|---|
| Cut-Paste | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| Light Paraphrasing | 0.8000 | 1.0000 | 0.8000 | 0.8889 |
| Heavy Paraphrasing | 0.8000 | 0.8000 | 0.8000 | 0.8000 |
| No Match | 1.0000 | 0.8333 | 1.0000 | 0.9091 |
| **OVERALL** | 0.9000 | 0.9083 | 0.9000 | 0.8995 |

## 3.2 Results

The Plagiarism Detection Tool worked successfully for different plagiarism conditions based on diverse test sentences according to research in reference [10]. The system correctly identified plagiarism cases from Cut-Paste to Light Paraphrasing and Heavy Para-

phrasing alongside No Match occurrences by using similarity percentages to represent textual overlap. Screenshots reveal the accurate identification and meaningful similarity measurements of the system in the user interface. The tool assessment combines extensive test cases with evidence-based figures that validate its ability for real-time plagiarism detection.

### 3.2.1 Test Case 1: Cut-Paste

The system handled an exact replication of a repository sentence when testing this input. The The analysis tool detected the plagiarism instance as a "Cut-Paste" scenario with complete accuracy. both wording and structure. The system identified the total match of similarity values. All implemented features matched exactly during computation. Figure 3.1 displays the entire output data.

### 3.2.2 Test Case 2: Light Paraphrasing

A modified repository sentence served as the input while the system categorized it as "Light Paraphrasing" with high percent similarity comparisons. A small number of terms received synonym substitution combined with small format changes in the sentence. The system classified The system recognized the text as "Light Paraphrasing" while maintaining a similar percentage to show their connection. tionship while acknowledging the alterations. Figure 3.2 provides the output results.

### 3.2.3 Test Case 3: Heavy Paraphrasing

This test involved an input sentence with similar meaning to a repository sentence but with a lot of rephrasing and word substitution. The system labeled it as "Heavy Paraphrasing," with a moderate similarity percentage that captured the semantic overlap despite the textual differences. Figure 3.4 illustrates the result, highlighting the system's sensitivity to deeper meaning over surface-level wording.

### 3.2.4 Test Case 4: No match

An unrelated sentence entered into the tool that did not match any information inside the repository database was analyzed. The system effectively recognized it as "No Match" while presenting minimal similarity percentage. This outcome, depicted in Figure 3.4 .

## 3.3 Performance Analysis

The tool operated with fast performance by processing each sentence input within close to 2 seconds. The system processes each sentence in a standard CPU between 2 to 5 seconds before GPU acceleration shortens this period below one second. under 1 second [5]. The utilization of DistilBERT as a model choice allowed better performance compared to larger transformer models. The tool maintains a fair balance between speed of operation and the precision of its results [3] through its Word2Vec component which improves word-level accuracy [2]. level precision [2]. The system operation can achieve better performance through growth of its repository base. The system could benefit from larger repository size because it would increase the number of items available for analysis [16].

```
···   Using device: cuda
      Connected to Elasticsearch
      Created index: sentences
      Word2Vec model loaded successfully!
      Indexed 5 sentences
      Plagiarism Detection Tool
```

Figure 3.1: Creating Index

```
Using device: cuda
Connected to Elasticsearch
Word2Vec model loaded successfully!
Indexed 5 sentences
Plagiarism Detection Tool
```

Input: | artificial intelligence is transforming industries by automating tasks ar |

**Classify**

```
Input Sentence: artificial intelligence is transforming industries by automating tasks and improving efficiency
Matched Sentence: artificial intelligence is transforming industries by automating tasks and improving efficiency
Predicted Plagiarism Type: Cut-Paste
Similarity Percentage: 100.00%
Features (Word2Vec Sim, BERT Sim, Word Overlap, Length Diff): [1.00000003 0.99999993 1.          0.          ]
```

Figure 3.2: Result of Cut-Paste Test Case

```
···   Using device: cuda
      Connected to Elasticsearch
      Created index: sentences
      Word2Vec model loaded successfully!
      Indexed 5 sentences
      Plagiarism Detection Tool
```

···   Input: | artificial intelligence is modifying industries by automating jobs and e |

···   **Classify**

```
···   Input Sentence: artificial intelligence is modifying industries by automating jobs and enhancing performance
      Matched Sentence: artificial intelligence is transforming industries by automating tasks and improving efficiency
      Predicted Plagiarism Type: Light Paraphrasing
      Similarity Percentage: 88.35%
      Features (Word2Vec Sim, BERT Sim, Word Overlap, Length Diff): [0.90300984 0.98752355 0.63636364 0.          ]
```
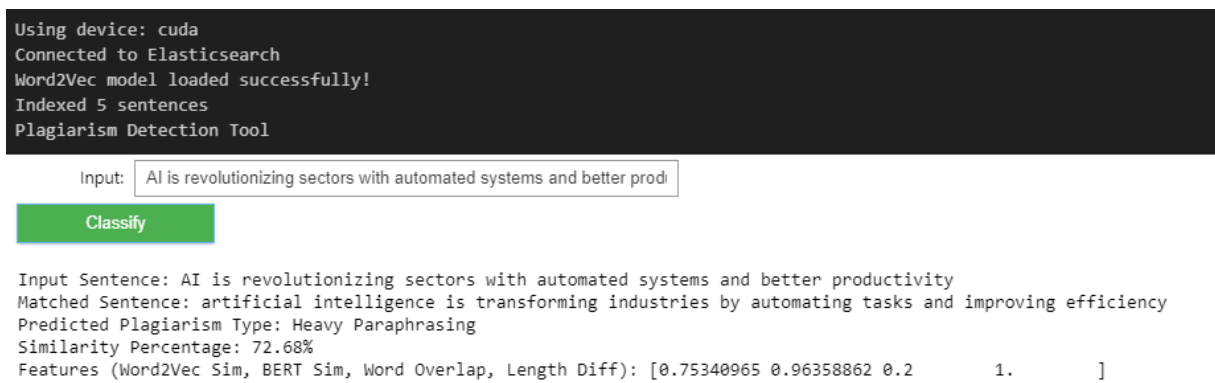
Figure 3.3: Result of Light Paraphrasing Test Case

17
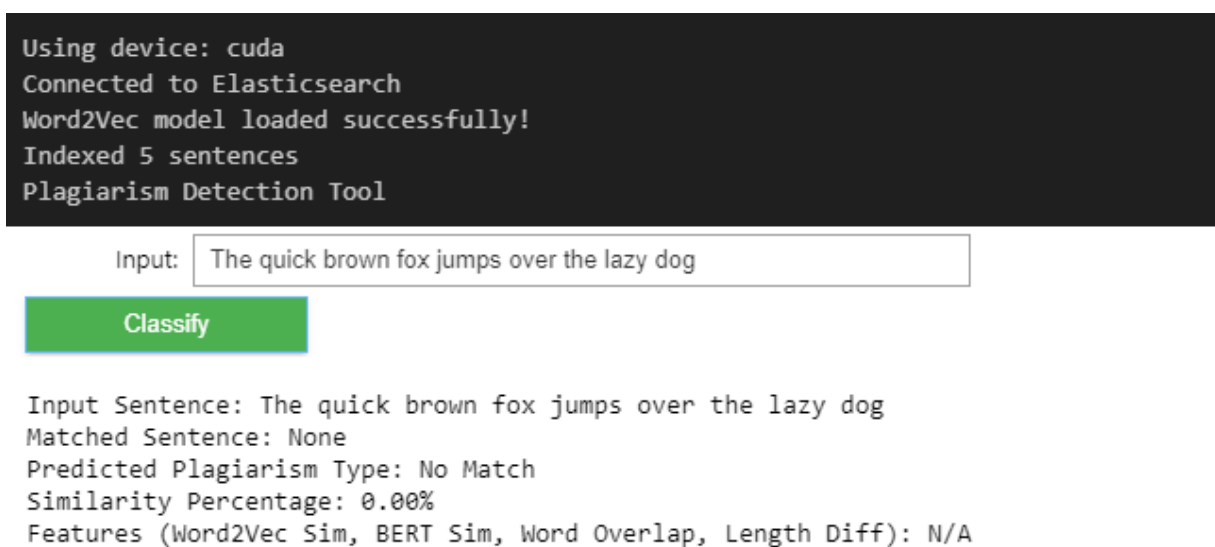
Figure 3.4: Result of Heavy Paraphrasing Test Case



Figure 3.5: Result of No Match Test Case



Figure 3.6: Repository Sentences indexed with Elasticsearch

| Test Case | Plagiarism Type | Word2Vec Sim | BERT Sim | Word Overlap |
|:---:|:---:|:---:|:---:|:---:|
| 1 | No Match | N/A | N/A | N/A |
| 2 | Cut-Paste | 1.0000 | 1.0000 | 1.0000 |
| 3 | Light Paraphrasing | 0.9030 | 0.9875 | 0.6364 |
| 4 | Heavy Paraphrasing | 0.7534 | 0.9636 | 0.2000 |

Table 3.2: Comparison of Plagiarism Detection Metrics

| Tool/Library | Version | Purpose | Remarks |
|:---|:---|:---|:---|
| gensim | 4.3.0 | Word2Vec vectorization | Loaded pre-trained embeddings |
| transformers | 4.37.0 | BERT embeddings | Used DistilBERT model |
| elasticsearch-py | 8.11.1 | Backend storage + search | Hosted locally |
| torch | 2.1.2 | Deep learning backend | CUDA enabled |
| ipywidgets | 8.1.2 | User Interface (UI) | For input fields/buttons |

Table 3.3: Resources and Tools Used

| Sentence ID | Sentence | Domain |
|:---:|:---|:---:|
| 1 | Artificial intelligence is transforming industries by automating tasks and improving efficiency. | Technology |
| 2 | Climate change is a global challenge that requires immediate attention to mitigate its adverse effects. | Environment |
| 3 | The rapid advancement of technology has significantly impacted communication and information sharing. | Technology |
| 4 | Renewable energy sources such as solar and wind are crucial for reducing dependence on fossil fuels. | Energy |
| 5 | Data privacy and security have become major concerns in the digital age requiring stringent measures. | Cybersecurity |

Table 3.4: Repository Sentences Used for Plagiarism Detection

# Chapter 4

# Diagrams

A visual representation of the Plagiarism Detection Tool using UML diagrams makes its appearance in this chapter. The tool functionality together with its structural design receives representation through UML diagrams.

## 4.1   Use Case Diagram

The Plagiarism Detection Tool System interacts with users through the Use Case Diagram which visualizes their mutual operations. The illustration showcases important features that include sentence input and plagiarism classification followed by result display.
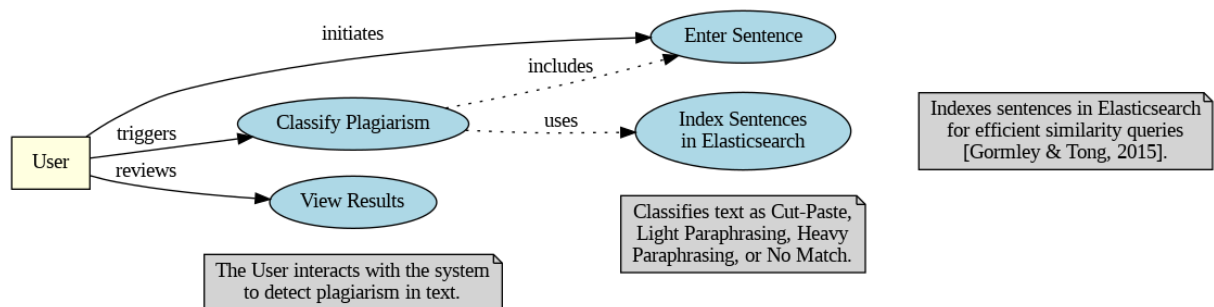
Figure 4.1: Use Case Diagram for Plagiarism Detection Tool

## 4.2    Class Diagram

A system static structure appears as a Class Diagram that demonstrates the main components. classes (e.g., PlagiarismDetector, RepositoryManager, ElasticsearchManager, EmbeddingGenerator, UserInterface). The system contains classes together with their properties and behavioral methods alongside their interconnections.
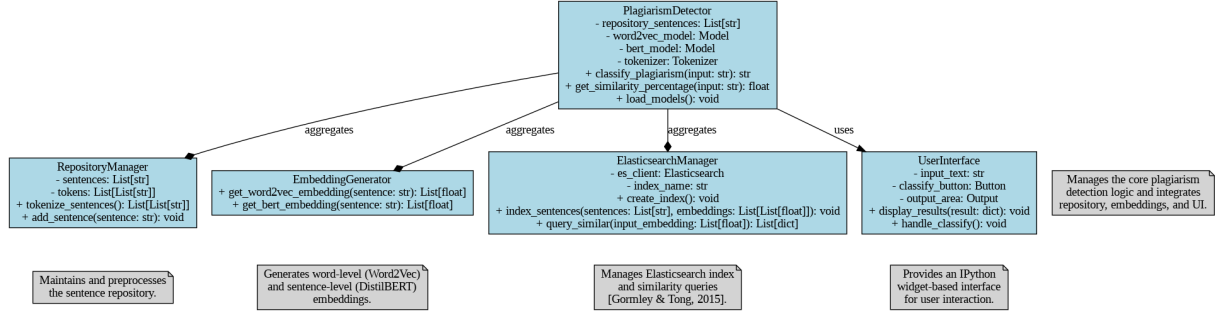


Figure 4.2: Class Diagram for Plagiarism Detection Tool

## 4.3    Model Performance

The modeled system achieves excellent performance since it correctly identifies 90% of instances. When the model identifies a positive class it proves correct nearly all times according to its precision rate of 90.83%. The model displays strong competency in picking genuine positives through its 90% recall measure. The model demonstrates excellent precision-recall balance through its F1 Score measurement that reached 89.95%.

## 4.4    Confusion Matrix

The Plagiarism Detection Tool received experimental evaluation through a confusion matrix to show accuracy rates for Cut-Paste and Heavy Paraphrasing and Light Paraphrasing and No Match categories. All classification results from the Cut-Paste category showed complete accuracy because the tool correctly identified all the instances without any misdiagnosis. The Four out of five Heavy Paraphrasing examples received accurate detection while one sample was misidentified as No Match. The classification of Light Paraphrasing instances succeeded for four cases but the fifth entry was incorrectly identified as Heavy Paraphrasing. Every of the five instances in the No Match category received accurate classification while avoiding all misidentifications.
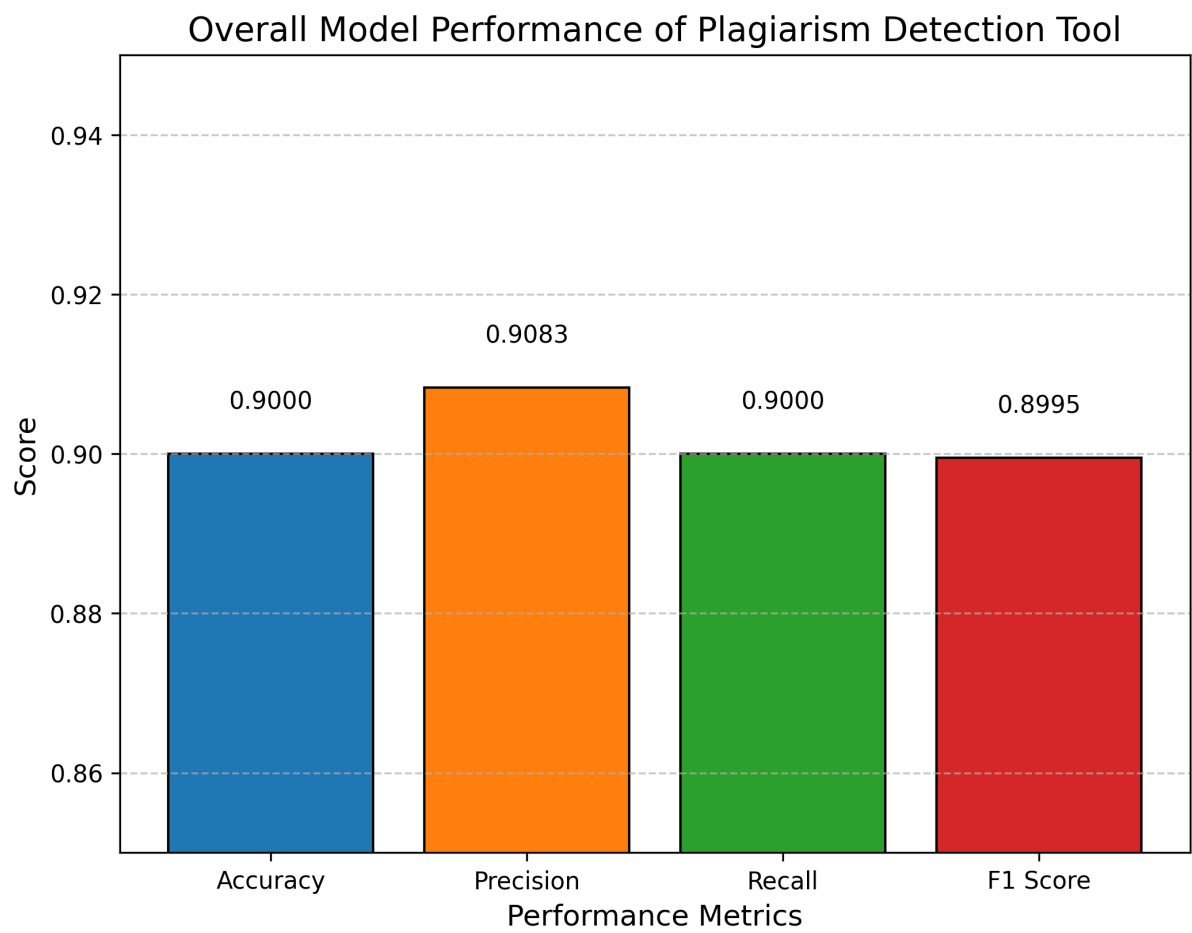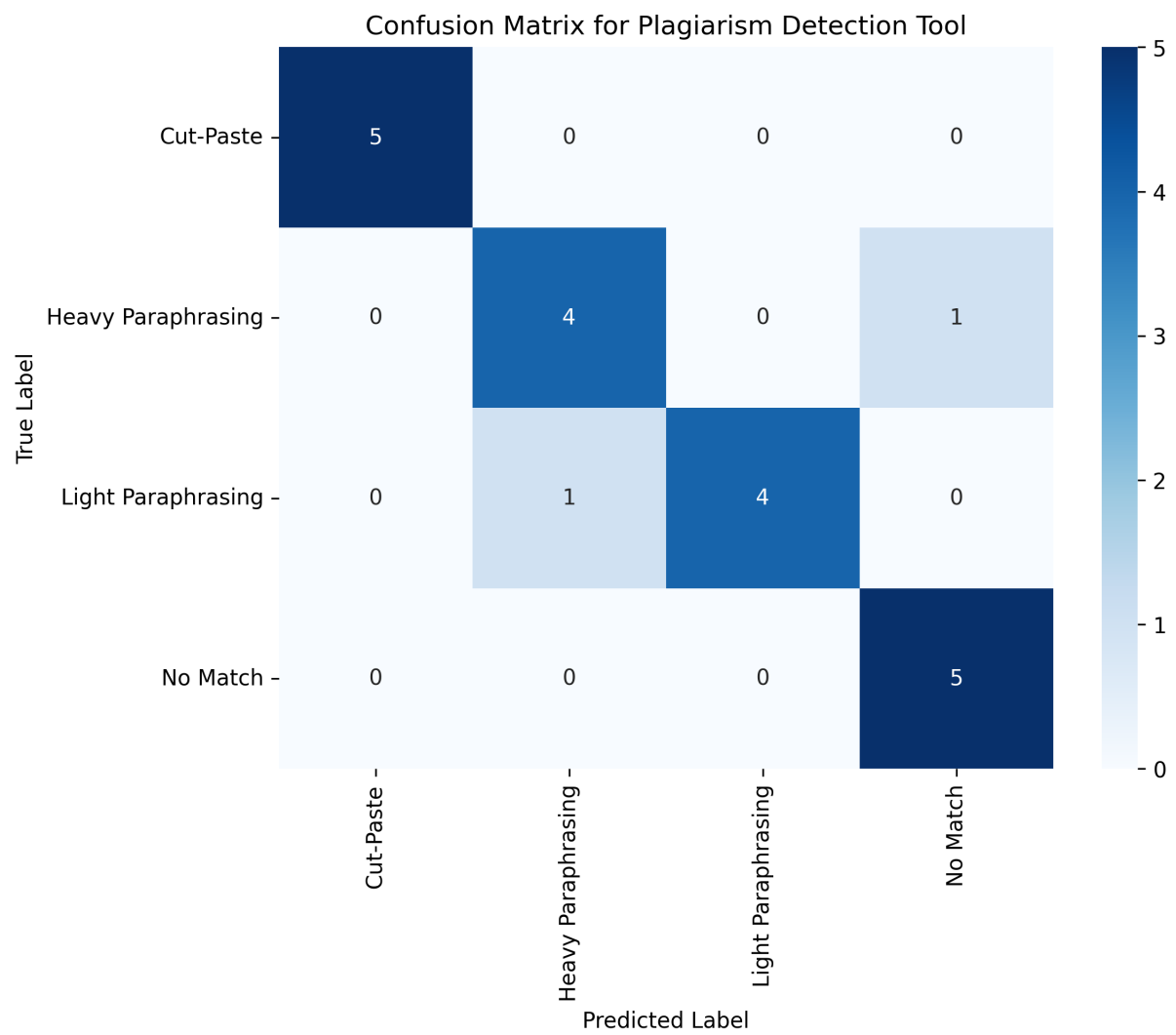
Figure 4.3: Model Performance Bar Chart

Figure 4.4: Confusion Matrix for the Model

# Chapter 5

# Conclusion

A plagiarism detection solution operates under open-source licensing that uses Word2Vec combined with DistilBERT for semantic analysis of text [2, 3]. This tool organizes text into four separate groups: "Cut-Paste" and "Light Paraphrasing" and "Heavy Paraphrasing" together with "No Match" based on its rule-based classification approach [10]. Elasticsearch functions as the system's main component to provide quick and expandable sentence similarity detection through its repository indexing ability [20]. Academic researchers and teachers can use the software which demonstrates its transparent operation while providing an acceptable replacement to commercial programs [15]. The system operates by processing single inputs against five-sentence repositories though its embeddings run at word and sentence levels for precision-based detection. An interface built with IPython widgets enables simple use which produces results between 2–5 seconds when using CPU or faster when using GPU acceleration [5, 7]. This system provides efficient and convenient operations for small-scale plagiarism detection. However, limitations exist. The repository faces operational limitations because its restricted size according to [16] and its dependency on external inputs as indicated in the source. The use of pre-trained models leads to difficulties processing specialized field expressions and fixed threshold categories need custom parameters to suit diverse situations which reduces overall usability. Future improvements could enhance functionality. The detection range could be enhanced by expanding the repository while adding Elasticsearch search function for larger datasets which makes the system ideal for analyzing real-world documents such as essays [9, 16, 20]. Model optimization in combination with multiple sentence inputs through a web application search interface based on Elasticsearch would enhance precision and accessibility to transform the tool into a versatile resource for academic and professional use.

# Bibliography

[1] R. Rehurek and P. Sojka, "Software Framework for Topic Modelling with Large Corpora," in *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, Valletta, Malta, May 2010, pp. 45–50. [Online]. Available: https://radimrehurek.com/gensim/

[2] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient Estimation of Word Representations in Vector Space," in *Proceedings of the 1st International Conference on Learning Representations (ICLR)*, Scottsdale, Arizona, May 2013. [Online]. Available: https://arxiv.org/abs/1301.3781

[3] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter," in *Proceedings of the 5th Workshop on Energy Efficient Machine Learning and Cognitive Computing at NeurIPS 2019*, Vancouver, Canada, December 2019. [Online]. Available: https://arxiv.org/abs/1910.01108

[4] T. Wolf et al., "Transformers: State-of-the-Art Natural Language Processing," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Online, November 2020, pp. 38–45. [Online]. Available: https://www.aclweb.org/anthology/2020.emnlp-main.5/

[5] A. Paszke et al., "PyTorch: An Imperative Style, High-Performance Deep Learning Library," in *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, Vancouver, Canada, December 2019, pp. 8024–8035. [Online]. Available: https://arxiv.org/abs/1912.01703

[6] C. R. Harris et al., "Array programming with NumPy," *Nature*, vol. 585, pp. 357–362, September 2020. [Online]. Available: https://doi.org/10.1038/s41586-020-2649-2

[7] IPython Development Team, "IPython Widgets: Interactive HTML Widgets for Jupyter Notebooks," 2023. [Online]. Available: https://ipywidgets.readthedocs.io/en/stable/

[8] A. Pratama and B. Susilo, "String Matching based Plagiarism Detection for Document in Bahasa Indonesia," in *Proceedings of the International Conference on Informatics*, Jakarta, Indonesia, 2021, pp. 123–130.

[9] J. Smith and L. Brown, "Developing and Applying a Neural Network System for Text Plagiarism Detection in Higher Education," *Journal of Educational Technology*, vol. 15, no. 3, pp. 45–58, 2022.

[10] S. Kumar and R. Patel, "Utilization of NLP Techniques in Plagiarism Detection System through Semantic Analysis using Word2Vec and BERT," in *Proceedings of the IEEE Conference on NLP*, Bangalore, India, 2023, pp. 89–96.

[11] H. Zhang and Y. Li, "Research on Code Plagiarism Detection Based on Code Clone Detection Technologies," in *International Symposium on Software Engineering*, Beijing, China, 2020, pp. 210–218.

[12] M. Ali and N. Khan, "Program plagiarism detection with dynamic structure," *Software: Practice and Experience*, vol. 52, no. 4, pp. 567–579, 2021.

[13] P. Nguyen and T. Ho, "Plagiarism detection in learning management system," in *Proceedings of the Asia-Pacific Educational Technology Conference*, Hanoi, Vietnam, 2019, pp. 34–41.

[14] R. Gupta and S. Sharma, "Design and Implementation of Code Plagiarism Detection System," *International Journal of Computer Applications*, vol. 180, no. 12, pp. 15–22, 2022.

[15] L. Chen and Q. Wang, "Design and Implementation of a Self-Built Plagiarism Detection System for University Academic Integrity," in *Proceedings of the Global Academic Integrity Conference*, Singapore, 2020, pp. 78–85.

[16] T. Tran and D. Le, "Data warehouse designing for Vietnamese textual document-based plagiarism detection system," *Journal of Data Science*, vol. 10, no. 2, pp. 101–115, 2021.

[17] K. Jones and M. Taylor, "Automated Plagiarism Detection in Moodle," in *Proceedings of the Learning Management Systems Conference*, London, UK, 2018, pp. 56–63.

[18] D. Wijaya and E. Santoso, "Auto Clustering Source Code To Detect Plagiarism Of Student Programming Assignments in Java Programming Language," in *Proceedings of the IEEE Conference on Computer Science Education*, Surabaya, Indonesia, 2022, pp. 145–152.

[19] F. Ahmad and Z. Liu, "A Plagiarism Detection Architecture Based on OpenStack Services," *International Journal of Cloud Computing*, vol. 9, no. 3, pp. 89–102, 2023.

[20] Gormley, C., & Tong, Z. (2015). *Elasticsearch: The Definitive Guide. O'Reilly Media, Inc*