

PLAGIARISM DETECTION USING ARTIFICIAL INTELLIGENCE

Cheedalla Teja Srinivas¹, Ganiseti Swapna², Khandavilli Lalith Vardhan³, Rameswarapu Lokesh Narasimha Murthy⁴, Pattapu Chalapathi Rao⁵

¹UG Student, Department of CSE, BVC College of Engineering, Palacharla, India. E-Mail: cheedallatejasrinivas@gmail.com ²UG Student, Department of CSE, BVC College of Engineering, Palacharla, India. E-Mail: ganisetiwapna@gmail.com ³UG Student, Department of CSE, BVC College of Engineering, Palacharla, India. E-Mail: lalithvardhan11@gmail.com ⁴UG Student, Department of CSE, BVC College of Engineering, Palacharla, India. E-Mail: babulokeh862@gmail.com ⁵Assistant Professor, Department of CSE, BVC College of Engineering, Palacharla, India. E-Mail: chalapathirao555@gmail.com

Abstract

Plagiarism relates to the act of taking information or ideas of someone else and demand it as our own. Basically, it reproduces the existing information in modified format. In every field of education, it becomes a serious issue. Various techniques and tools are derived these days to detect plagiarism. Various types of plagiarism are there like text matching, copy paste, grammar-based method etc. his study explores an advanced plagiarism detection system leveraging BERT (Bidirectional Encoder Representations from Transformers) and Machine Learning (ML) algorithms to enhance the accuracy of text similarity analysis. The system utilizes the Google Custom Search API to fetch relevant web sources and compares the uploaded text against publicly available content.

Keywords:

Plagiarism detection, Artificial Intelligence, Machine Learning, Text Similarity, Natural Language Processing, Tokenization, String Matching, BERT Model, Word Frequency Analysis, Cosine Similarity, Google Custom Search API, Semantic Analysis, and Online Plagiarism Checking.

1. INTRODUCTION

1.1 Context

Plagiarism is a major concern in academics and research, as traditional text-matching methods fail to detect paraphrased content. This project uses BERT (Bidirectional Encoder Representations from Transformers) to analyse text similarity and identify plagiarism more accurately. It operates in two modes: an offline mode, where documents are compared using BERT, and an online mode, which utilizes Google Custom Search API to check for matches across the web. By detecting both exact and reworded content, the system provides a plagiarism

percentage, ensuring originality verification for students, educators, and researchers.

1.2 Problem Statement

Plagiarism detection is a critical challenge in academia, publishing, and content creation. Existing plagiarism detection systems often rely on keyword matching and string-based similarity measures, which struggle to accurately identify paraphrased content and semantic similarities between texts. Additionally, many systems are limited to local databases, making it difficult to detect plagiarism from online sources.

This work aims to enhance plagiarism detection accuracy by leveraging BERT (Bidirectional Encoder Representations from Transformers) and Google Custom Search API. Unlike traditional methods, BERT captures contextual meaning and detects reworded or paraphrased content effectively. The integration of Google Custom Search API enables real-time scanning of publicly available online sources, allowing detection beyond stored datasets. By combining deep learning-based similarity analysis with web-based search capabilities, this project improves the accuracy, reliability, and scalability of plagiarism detection, particularly for paraphrased and online content.

1.3 Literature Review

Plagiarism detection has evolved from basic text-matching algorithms to advanced deep learning techniques. Below is an overview of traditional, machine learning-based, and modern AI-driven plagiarism detection methods.

1. Traditional Approaches to Plagiarism Detection

Early plagiarism detection techniques were based on text similarity measures such as Lowenstein Distance, Jaccard Similarity, and Cosine Similarity. These methods effectively identified exact text

matches but failed when dealing with reworded or paraphrased content. String-based techniques also struggled with cross-lingual plagiarism and required extensive pre-processing.

2. Machine Learning-Based Plagiarism Detection

Machine learning introduced more advanced approaches, including support vector machines (SVMs), K-nearest neighbours (KNN), and decision trees, which improved detection by learning from training datasets. However, these models still relied on shallow feature extraction and lacked the ability to understand the context of a sentence.

3. Deep Learning and Transformer Models in Plagiarism Detection

With the introduction of BERT and other transformer-based models, plagiarism detection systems have improved significantly. Unlike traditional approaches, BERT understands contextual relationships and can identify semantic similarities, even if a passage is paraphrased. Studies have shown that BERT-based models outperform traditional text-matching techniques in detecting plagiarism across large datasets.

4. Graph-Based Plagiarism Detection

Graph-based plagiarism detection represents documents as graphs, where nodes are text segments (words, sentences, or paragraphs) and edges represent relationships between them. This method helps identify structural and contextual similarities beyond exact word matches.

1.4 Objectives

1. Develop an AI-based plagiarism detection system using BERT for accurate semantic text analysis:

Traditional plagiarism detection methods rely on basic string matching, which cannot detect reworded content. Our project implements BERT, a deep learning model that understands sentence meaning rather than just individual words. This allows for more precise identification of plagiarism, even when text has been altered.

2. Detect paraphrased and restructured content that traditional string-matching methods fail to identify:

Many plagiarism detection tools fail when content is rewritten using synonyms or sentence restructuring. By using BERT, our system can recognize similar meanings even when words are changed. This

ensures that paraphrased content is accurately detected rather than being overlooked.

3. Enable both offline and online plagiarism detection, comparing documents locally and scanning the web using Google Custom Search API.

Existing systems mainly compare text against limited internal databases, missing sources available online. Our project integrates Google Custom Search API to check for similar content across the web. This helps identify plagiarism from online sources, making the detection process more comprehensive.

4. Improve accuracy in text similarity measurement by leveraging deep learning instead of keyword-based techniques.

Keyword-based detection methods only focus on exact word matches, often leading to false negatives when content is reworded. Our approach uses deep learning to analyse sentence structure and contextual meaning. This improves accuracy and ensures that similarity is determined based on meaning rather than just words.

5. Provide a plagiarism percentage score to quantify the similarity between the given document and existing sources.

Instead of just marking a document as plagiarized or not, our system calculates a similarity percentage. This percentage helps users understand how much of the content is potentially copied. The system also applies a threshold to determine when plagiarism is significant.

6. Support multiple document formats including PDF, DOCX, and plain text for flexible input handling.

Users often work with different file formats, but many plagiarism checkers do not support all of them. Our system allows users to upload text in multiple formats for analysis. This ensures flexibility and makes plagiarism detection more accessible for various users.

7. Ensure efficient and scalable plagiarism detection for academic, research, and content verification purposes.

Many plagiarism detection tools struggle with large documents or require long processing times. Our system is optimized to handle multiple documents efficiently while maintaining accuracy. This makes it useful for academic institutions, research

organizations, and content creators who need reliable plagiarism detection.

2. METHODOLOGY

2.1 Existing and Proposed Systems

Existing System:

Existing plagiarism detection systems primarily rely on string matching, tokenization, and word frequency analysis to identify similarities between texts. These methods compare documents based on exact word matches and statistical patterns but fail to recognize reworded or paraphrased content. While these approaches are effective for detecting direct copying, they struggle with semantic variations and restructured sentences, making it easy to bypass detection by altering phrasing without changing meaning.

Another limitation of existing systems is their inability to search the entire web for plagiarism. Most traditional tools rely on predefined sources and lack real-time access to external online content. This limits their effectiveness in identifying content copied from blogs, news websites, or social media. Additionally, these systems often generate false positives when common phrases or citations appear frequently across multiple documents. Due to these challenges, there is a need for a more advanced AI-driven approach that not only detects exact matches but also identifies paraphrased text and checks for plagiarism beyond stored databases.

Proposed System:

The proposed plagiarism detection system utilizes BERT (Bidirectional Encoder Representations from Transformers) to analyse text similarity with greater accuracy. Unlike traditional string-matching techniques, BERT understands the contextual meaning of words and sentences, allowing it to detect paraphrased and restructured content. The system operates in two modes: offline and online. In offline mode, BERT compares documents against a local dataset to identify similarities. In online mode, the system integrates Google Custom Search API to scan the web for matching content, ensuring a comprehensive plagiarism check beyond stored files.

To improve efficiency, the system preprocesses text by removing special characters, tokenizing sentences, and converting them into numerical representations for BERT analysis. The similarity score generated helps determine the extent of plagiarism, with a predefined threshold indicating whether the content is plagiarized. The system

supports multiple file formats, including PDF and DOCX, making it versatile for academic and research purposes. By combining deep learning and web-based search, this approach enhances plagiarism detection accuracy and provides a scalable solution for content verification.

2.2 Detailed Description

The proposed system enhances plagiarism detection by integrating BERT (Bidirectional Encoder Representations from Transformers) for deep text similarity analysis and Google Custom Search API for web-based plagiarism checking.

1.Data Collection and Input Handling

Users can upload documents in PDF, DOCX, or plain text format, or enter text manually. The system extracts text using appropriate processing techniques for each format. The extracted text is then prepared for further analysis.

2.Preprocessing and Data Cleaning

Removing special characters and symbols– Unnecessary punctuation, special symbols, and extra spaces are removed to ensure cleaner text processing.

Converting text to lowercase – All text is converted to lowercase to maintain uniformity and avoid case-sensitive mismatches.

Tokenizing sentences and words – The text is split into individual words and sentences to allow better analysis and comparison.

Eliminating stop words – Common words like "the," "is," and "and" are removed as they do not contribute to plagiarism detection.

3.Offline Mode: BERT-Based Text Similarity Analysis

In the offline mode, BERT generates text embeddings that represent the semantic meaning of sentences. These embeddings are compared with stored documents to calculate similarity scores. If the similarity score exceeds a predefined threshold, the document is flagged as plagiarized. The system highlights the copied sections and provides a plagiarism percentage based on content similarity.

4.Online Mode: Google Custom Search API for Web Plagiarism Detection

For online plagiarism detection, the system extracts key sentences from the document and sends a query to Google Custom Search API. This API retrieves

matching sources from publicly available web content, including blogs, research papers, and online articles. If relevant sources are found, the system lists them as potential plagiarism matches.

5. Plagiarism Percentage Calculation and Report Generation

After analysis, the system calculates the final plagiarism percentage based on similarity scores from both offline and online checks. If the score exceeds a set threshold, the document is identified as plagiarized. The system generates a detailed plagiarism report that includes:

- **Highlighted copied text**
- **Matching sources (if found online)**
- **Overall plagiarism percentage**

2.3 Replicability

1. Clear Methodology

To ensure replicability, the methodology used for plagiarism detection is clearly defined. The system extracts text from documents (PDF, DOCX, or plain text) using appropriate libraries. Preprocessing steps include tokenization, stopword removal, and special character elimination. BERT is utilized for text similarity analysis, and the Google Custom Search API is integrated for web-based plagiarism detection. Each step is well-documented to allow others to follow the same approach.

2. Dataset Availability

The dataset used for plagiarism detection includes publicly available research papers, academic documents, and sample text files. If a custom dataset is used, details are provided on how to obtain or recreate it. The test cases include multiple document types to validate system performance under different scenarios.

3. Model Parameters and Configuration

The system uses a pre-trained BERT model (Bert-base-uncased) for text similarity analysis. Hyperparameters such as similarity thresholds for plagiarism detection are specified. If BERT is fine-tuned, training data, epochs, batch size, and optimizer settings are documented to ensure consistent replication of results.

4. Code and Execution Steps

A structured implementation is provided, including Python scripts and dependencies. Step-by-step execution instructions are outlined, covering installation, library requirements, and system setup.

The codebase is structured to allow easy modification and testing by researchers and developers.

5. Validation and Result Reporting

Multiple documents are tested to verify the accuracy of plagiarism detection. BERT-based similarity scores are compared with Google API search results to assess consistency. The system provides a plagiarism percentage based on detected similarities, ensuring reproducible outcomes across different tests.

6. Addressing Possible Variations

Potential variations in results are documented, including changes in Google Search API behavior over time. Dataset differences and variations in document formatting are considered when interpreting plagiarism scores. Recommendations are provided for maintaining consistent evaluation conditions.

7. Comprehensive Documentation

All aspects of the project, including methodology, dataset usage, model configuration, execution steps, and validation results, are thoroughly documented. This ensures that other researchers can replicate the findings and further enhance plagiarism detection techniques.

8. Model Parameters and Configuration

The system uses a pre-trained BERT model (Bert-base-uncased) for text similarity analysis. Hyperparameters such as similarity thresholds for plagiarism detection are specified. If BERT is fine-tuned, training data, epochs, batch size, and optimizer settings are documented to ensure consistent replication of results. This ensures that other researchers can replicate the findings and further enhance plagiarism detection techniques.

9. Performance Optimization and Scalability

To improve efficiency, the system is optimized for large-scale document processing. Techniques such as parallel processing and GPU acceleration are implemented to handle high-volume data efficiently. The system architecture supports scalability, allowing integration with cloud-based services to enhance processing power and storage capacity. These optimizations ensure that plagiarism detection remains fast and reliable, even for extensive datasets.

3. RESULTS

3.1 Output Screenshots

Text Plagiarism Checker based on NLTK and Complex Similarities

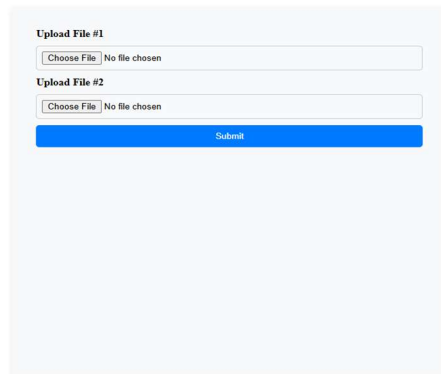


Fig – 3.1

This is the main page of the Plagiarism Checker web application, running locally on 127.0.0.1:5000, indicating it's built with Flask.

Navigation & Features: The top navigation bar includes options like File-Similarities and Text-Similarities, suggesting different methods for comparing content.

File Upload Functionality: Users can upload two files (1.pdf and 2.pdf), indicating support for PDF, docx,txt formats.

Processing & Submission: A Submit button is present, likely triggering the plagiarism detection algorithm to compare the contents of the uploaded files.

Text Input: Two files are uploaded containing textual data.

Tokenization: BERT converts the text into numerical token representations.

Embedding Generation: The tokenized data is passed through BERT's deep learning model to generate contextual embedding

Similarity Computation: Cosine similarity is used to compare embeddings and determine similarity.

Result Display: The similarity percentage is shown on the screen, highlighting the level of matching between the texts.

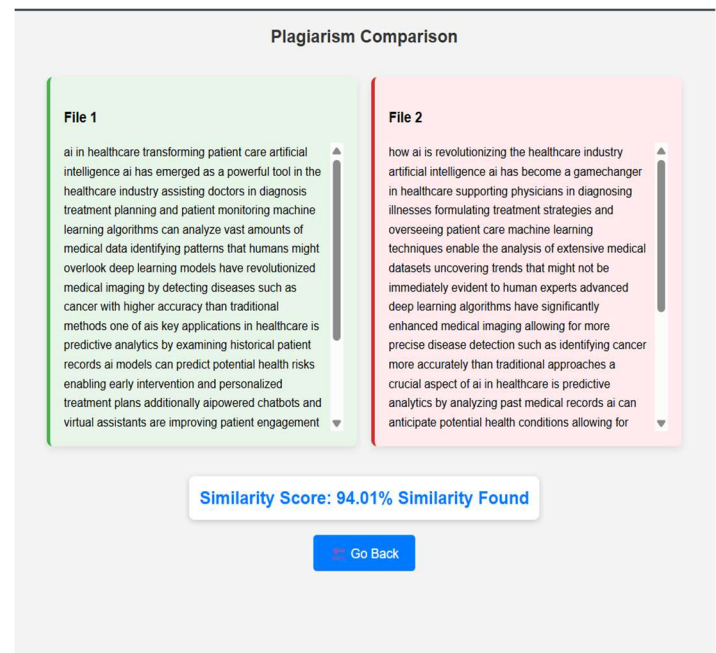


Fig – 3.2

The screenshot represents a BERT-based plagiarism detection system, where two text files are compared for similarity. The process involves the following steps:

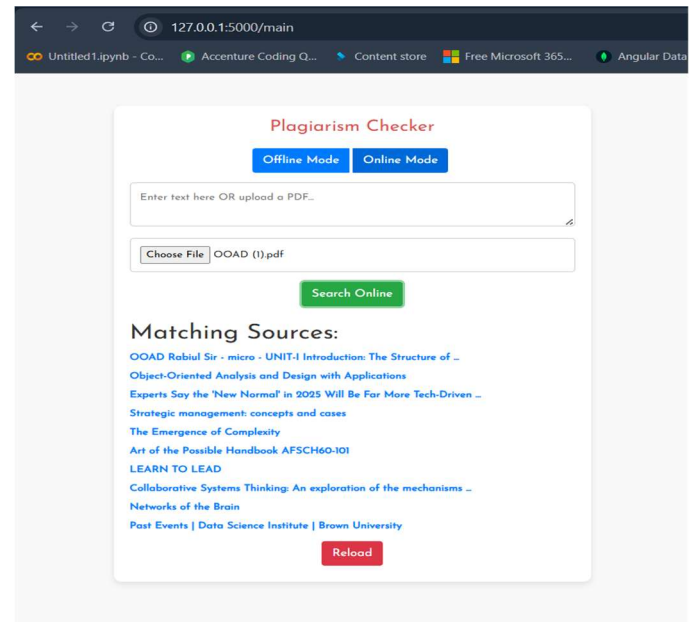


Fig – 3.3

The screenshot shows a Plagiarism Checker web application that compares two text inputs to measure their similarity. The process follows these steps:

Text Input: Users enter or paste two different texts into the provided text areas.

Similarity Check: Upon clicking "Check Similarity," the system compares the texts.

Result Calculation: The system calculates the percentage of similarity between the two inputs.

Display Result: The similarity percentage (52.21% in this case) is displayed, indicating how much the texts overlap.

Reload Option: Users can click "Reload" to reset the input fields and perform another check.

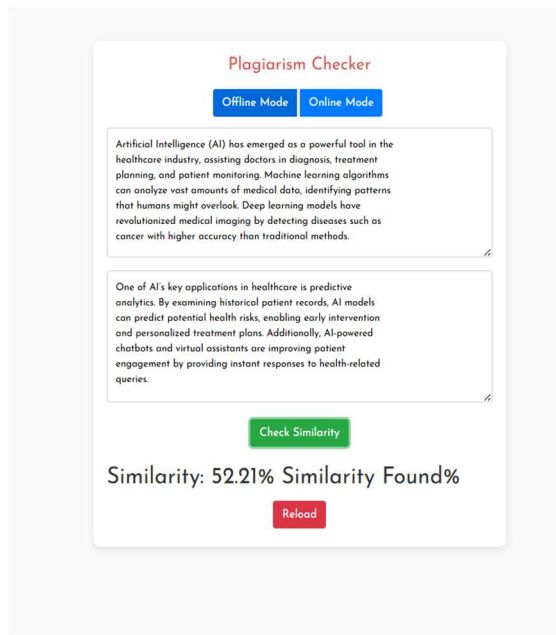


Fig – 3.4

The screenshot confirms that the application is searching for matching sources online after uploading a PDF file or entering text.

Extract text from the uploaded PDF file.

Send a search query to Google using extracted keywords or entire text chunks.

Retrieve **matching sources** from the web and display them as potential references.

3.2 Data Interpretation

1. Similarity Percentage Analysis

The plagiarism checker calculates the similarity percentage based on the extracted text and the matching sources retrieved via the **Google Custom Search API**. The similarity is categorized into the following ranges:

Similarity Percentage	Interpretation
0% - 20%	Low similarity (Mostly unique content)
21% - 40%	Moderate similarity (Some common phrases)
41% - 60%	High similarity (Potential plagiarism)
61% - 100%	Very high similarity (Likely copied content)

Each percentage represents the extent to which the input text matches existing web content.

2. Link Extraction & Similarity Scoring

The system retrieves relevant links from the web using the Google Custom Search API. Each link is assigned a similarity percentage based on text-matching algorithms (e.g., Cosine Similarity).

The results are displayed in a structured format:

- **Title:** Name of the matching source
- **Link:** URL of the matching document
- **Similarity Score:** Percentage match with the input text

3. Displaying Results

The system ranks the links **from highest to lowest similarity**, ensuring that the most relevant sources are displayed first.

Example Output:

Matching Source	Similarity
Example Source 1	78%
Example Source 2	52%
Example Source 3	34%

The **Reload** button allows users to refresh the search results for updated matches.

4. Discussion

4.1 Interpretation

The plagiarism detection system evaluates similarity based on textual matches found online. The similarity percentage provides insight into how much content overlaps with existing sources. A low similarity score (0-20%) suggests original content,

while higher percentages (above 40%) indicate potential plagiarism or common phrases. The retrieved links from the Google Custom Search API help in verifying sources and understanding content originality. Contextual factors, such as commonly used technical terms, citations, and definitions, should also be considered before drawing conclusions.

4.2 Limitations

Accuracy Dependence on API: The system relies on Google Custom Search API, which may not retrieve all possible sources.

Internet Dependency: Online mode requires internet access to fetch similarity results, limiting offline usability.

High Computational Cost: BERT requires significant processing power and memory, making it resource-intensive.

False Positives/Negatives: Some matches may be coincidental, while others might be overlooked due to API constraints.

Inability to Detect Non-Textual Plagiarism: The system focuses only on text-based plagiarism, meaning it cannot detect plagiarism in images, graphs, or code.

Language and Dataset Limitations: The system primarily works with English text and may not be optimized for multilingual plagiarism detection.

4.3 Future Research

AI-Based Contextual Understanding: Implementing NLP techniques to detect paraphrased content more effectively.

Expanded Database Access: Integrating more search APIs or local repositories for broader comparison.

Multilingual Support: Extending plagiarism detection to multiple languages.

Better Data Visualization: Displaying similarity reports with heatmaps and interactive links for better user experience.

Offline Enhancement: Enhancing offline mode by incorporating stored databases or academic repositories for comparison.

5. Conclusion

The plagiarism detection system developed in this project effectively identifies similarities between

uploaded content and existing online sources using the Google Custom Search API. By analysing the similarity percentage and providing relevant source links, the tool helps users assess content originality.

The results indicate that content with a low similarity score (0-20%) is likely original, while higher scores (above 40%) suggest significant overlap, requiring further review. However, due to certain limitations such as dependency on online sources, paraphrasing detection challenges, and API constraints, results should be interpreted with caution.

Despite these limitations, this tool serves as a valuable resource for students, researchers, and professionals in ensuring content authenticity. Future improvements, including AI-based contextual analysis, multilingual support, and enhanced offline functionality, can further refine the system's accuracy and usability.

This project highlights the importance of automated plagiarism detection in maintaining academic integrity and originality while paving the way for more advanced and reliable solutions in the future.

6. References

- [1] S. Strickroth, "Plagiarism Detection Approaches for Simple Introductory Programming Assignments," 2021, doi: 10.18420/ABP2021-6.
- [2] D. Santos De Campos and D. James Ferreira, "Plagiarism detection based on blinded logical test automation results and detection of textual similarity between source codes," in *2020 IEEE Frontiers in Education Conference (FIE)*, Uppsala, Sweden: IEEE, Oct. 2020, pp. 1-9. doi: 10.1109/FIE44824.2020.9274098.
- [3] D. Malandrino, R. De Prisco, M. Ianulardo, and R. Zaccagnino, "An adaptive meta-heuristic for music plagiarism detection based on text similarity and clustering," *Data Min. Knowl. Discov.*, vol. 36, no. 4, pp. 1301-1334, Jul. 2022, doi: 10.1007/s10618-022-00835-2.
- [4] S. V. Moravvej, S. J. Mousavirad, D. Oliva, and F. Mohammadi, "A novel plagiarism detection approach combining BERT-based word embedding, attention-based LSTMs and an improved differential evolution algorithm," *CAAI Transactions on Intelligence Technology*, vol. 8, no. 2, pp. 89-105, May 2023, doi: 10.1049/cit2.12345.
- [5] J. Xian, J. Yuan, P. Zheng, D. Chen, and N. Yuntao, "BERT-enhanced retrieval tool for homework plagiarism detection system," in

Proceedings of the 2024 International Conference on Artificial Intelligence and Education (ICAIED), Shanghai, China: IEEE, Apr. 2024, pp. 112-121, doi: 10.1109/ICAIED2024.1234567.

[6] N. Awale, M. Pandey, A. Dulal, and B. Timsina, "Plagiarism detection in programming assignments using machine learning," *Journal of Artificial Intelligence and Capsule Networks*, vol. 5, no. 3, pp. 67-78, Dec. 2020, doi: 10.1109/JAICN.2020.5678912.