

A Dual-Mode Plagiarism Detection System Leveraging Sequence Matching and Similarity Metrics

Dr. Manisha Y. Joshi
Computer Science and Engineering
MGM's College of Engineering ,
Nanded
Nanded, India
joshi_my@mgmcen.ac.in

Ganesh Kondamwar
Computer Science and Engineering
MGM's College of Engineering ,
Nanded
Nanded, India
ganeshkondamwar@gmail.com

Soham Kashettiwar
Computer Science and Engineering
MGM's College of Engineering ,
Nanded
Nanded, India
sohamkashettiwar@gmail.com

Anupkumar Koturwar
Computer Science and Engineering
MGM's College of Engineering ,
Nanded
Nanded, India
koturwaranup@gmail.com

Jawad Momin
Computer Science and Engineering
MGM's College of Engineering,
Nanded
Nanded, India
mdjawad9228@gmail.com

Abstract— Plagiarism has become a significant concern in academic and professional fields due to the increasing accessibility of digital information. This paper presents a comprehensive Plagiarism Detection System that operates in both online and offline modes, offering robust and efficient plagiarism detection for text and PDF documents. The system utilizes sequence matching algorithms as the core methodology to detect copied content and determine similarity percentages. In online mode, the system compares user-provided text or uploaded PDF files against internet sources and prioritizes results based on maximum similarity, ensuring reliable detection. The offline mode allows users to upload multiple PDF documents, store them locally, and detect plagiarism in a separate file against the stored repository. To enhance accuracy, advanced features such as Longest Common Subsequence (LCS) and Cosine Similarity are integrated with sequence matching, providing multi-dimensional similarity analysis. The system highlights duplicated content, calculates overall plagiarism percentages, and presents visually appealing, detailed reports using structured CSS for improved user experience. Additionally, redundancy in results is eliminated by efficiently handling identical files, and a threshold mechanism is implemented to only display significant matches. This project aims to provide a reliable, user-friendly, and scalable solution for plagiarism detection, contributing to academic integrity and content originality.

Keywords— *Plagiarism Detection, Sequence Matching, Cosine Similarity, Longest Common Subsequence (LCS), Online and Offline Modes, Text and PDF Analysis, Similarity Threshold, Highlighted Reports.*

I. INTRODUCTION

Plagiarism detection has become an essential tool in academic and professional domains to ensure originality and maintain ethical standards in content creation. The proliferation of digital content and ease of access to vast amounts of information have amplified the challenges of detecting and addressing plagiarism. With increasing demands for accurate and scalable solutions, this project introduces a multi-mode plagiarism detection system designed to cater to diverse user requirements.

This system provides two operational modes: online and offline. In the online mode, the system analyzes text or uploaded PDF files by comparing them against a wide range of internet sources, ensuring comprehensive coverage. The

offline mode allows for checking against locally uploaded PDF files, focusing on stored documents to detect overlaps efficiently.

Key features of the system include:

- A sequence matching algorithm as the primary technique to identify similarities in text.
- Longest Common Subsequence (LCS) and cosine similarity for advanced similarity analysis, ensuring multi-dimensional comparison.
- Visualized reports with highlighted matching content and structured outputs using CSS.

The system emphasizes scalability, accuracy, and user experience, making it adaptable for individual users, educational institutions, and professional environments. This paper discusses the design, implementation, and testing of the proposed system, focusing on its ability to deliver efficient and reliable plagiarism detection.

II. EFFICIENCY IN PLAGIARISM DETECTION ALGORITHMS

Plagiarism detection algorithms employ various methods to identify text similarities. This section compares three widely used techniques: Cosine Similarity, Longest Common Subsequence (LCS), and Sequence Matching. While all three have merits, sequence matching with the Ratcliff-Obershelp algorithm, implemented in Python's difflib, demonstrates superior efficiency and accuracy by addressing limitations in the other methods.

A. Cosine Similarity: Strengths and Limitations

Cosine similarity measures textual similarity by converting text into vectors and computing the cosine of the angle between them. While effective for rephrased content, it has limitations:

- Word Frequency Dependence: Similarity is often based on common terms, leading to false positives.
- Sequence Ignorance: Lacks the ability to detect sequential or structural matches.

Limitation: Produces high similarity scores for texts with common words, ignoring sequence and context.

Pseudo-Code for Cosine Similarity:

```

Input: Text1, Text2
1. Preprocess Text1 and Text2:
  a. Convert to lowercase
  b. Remove stopwords and punctuation
  c. Tokenize into words
2. Create frequency vectors for both texts
3. Compute dot product of the vectors
4. Compute magnitudes of the vectors
5. Similarity = (Dot Product) / (Magnitude of Vector1 * Magnitude of Vector2)
Output: Similarity Score (0 to 1)

```

B. Pure LCS Algorithm: Challenges in Efficiency

The LCS algorithm identifies the longest sequence of characters common to two texts. While accurate for detecting verbatim matches, it is computationally expensive. Pseudo-Code for Pure LCS:

```

Input: Text1, Text2
1. Initialize a 2D matrix L of size (len(Text1)+1) x (len(Text2)+1)
2. For i = 1 to len(Text1):
  For j = 1 to len(Text2):
    If Text1[i-1] == Text2[j-1]:
      L[i][j] = L[i-1][j-1] + 1
    Else:
      L[i][j] = max(L[i-1][j], L[i][j-1])
3. Backtrack from L[len(Text1)][len(Text2)] to find the LCS
Output: Longest Common Subsequence

```

Limitation: Time complexity $O(m \times n)$ makes it unsuitable for large-scale comparisons.

C. Sequence Matching: A Superior Alternative

Sequence matching, implemented using the Ratcliff-Obershelp algorithm, extends LCS by identifying matching subsequences recursively and optimizing performance. Pseudo-Code for Sequence Matching with Ratcliff-Obershelp:

```

Input: Text1, Text2
1. Identify the longest contiguous matching subsequence (LCS)
2. Split Text1 and Text2 into unmatched portions:
  a. Left segment (before LCS)
  b. Right segment (after LCS)
3. Recursively apply steps 1-2 on unmatched portions
4. Compute similarity:
  Similarity = 2 * Matched Characters / Total Characters
Output: Similarity Score (0 to 1)

```

Advantages:

- Handles both verbatim and partially matching segments efficiently.
- Optimized time complexity, suitable for large-scale text comparisons.

a) Comparison Summary

TABLE II. COMPARISON OF ALGORITHMS

Method	Strengths	Limitations
Cosine Similarity	Detects rephrased content; simple computation	Ignores sequence; prone to false positives
Pure LCS	Accurate for verbatim overlaps	High computational cost for large inputs

Method	Strengths	Limitations
Sequence Matching	Efficient; balances precision and recall	Dependent on recursive implementation for partial matches

III. METHODOLOGY

A. Sequence Matching Algorithm

The core of the plagiarism detection system relies on the sequence matching algorithm. This method identifies sequences of characters or words that appear identically in two texts, making it effective for detecting verbatim copying. The algorithm scans the input text or PDF file and compares it with reference sources (online or offline) to identify matching segments.

B. Online Mode

In online mode, the system processes user input—typed text or uploaded PDF files—and compares it against internet sources. The implementation uses an API to fetch matching content and prioritize results by similarity percentage. The workflow includes:

1. Preprocessing the text to remove stopwords and normalize formatting.
2. Querying web sources using a search engine API.
3. Analyzing the retrieved results for textual similarity using sequence matching and cosine similarity.
4. Generating a detailed report with matched text and corresponding source links, sorted by similarity.

C. Offline Mode

Offline mode enables users to detect plagiarism within a local repository of documents. The system allows multiple PDF files to be uploaded, stored, and compared against a new file. The process includes:

1. Parsing and preprocessing the content of the uploaded PDF files.
2. Using the Longest Common Subsequence (LCS) method and sequence matching to identify matching segments between the new file and stored documents.
3. Displaying results, highlighting the file with the highest similarity and providing a summary of the matched content.

D. Cosine Similarity and LCS Integration

To enhance detection accuracy, cosine similarity and LCS are integrated alongside sequence matching.

- Cosine similarity measures the textual similarity based on vector space representation and is particularly effective for detecting rephrased content. Preprocessing includes removing stopwords and tokenizing text[3],[8].
- Longest Common Subsequence (LCS) identifies the longest subsequence common to both input texts, making it useful for detecting large overlapping sections[5],[7].

E. Report Generation

The results are presented in a user-friendly format, with highlights on matching text and detailed summaries. For offline mode, only the file with the highest similarity is displayed. For online mode, sources are ranked by similarity percentage. The system employs CSS for visual appeal, making the output structured and easy to interpret.

IV. RESULTS AND DISCUSSION

The plagiarism detection system was tested extensively in both online and offline modes to evaluate its efficiency, accuracy, and user experience. This section discusses the outcomes of the tests and the observations made during implementation.

Keep your text and graphic files separate until after the text has been formatted and styled. Do not use hard tabs, and limit use of hard returns to only one return at the end of a paragraph. Do not add any kind of pagination anywhere in the paper. Do not number text heads-the template will do that for you.

A. Online Mode Testing

The online mode was tested using various text inputs and PDF files to assess its ability to detect and prioritize sources based on similarity.

- **Performance:** The system successfully retrieved matching content from internet sources and displayed links ranked by similarity. It also identified and highlighted the specific segments of text that matched online sources[1],[9].
- **Accuracy:** The integration of sequence matching with cosine similarity ensured precise detection of both verbatim and partially rephrased content[3]. The system effectively minimized duplicate results by handling identical text segments across multiple sources[4].
- **Challenges:** Some inconsistencies were observed in retrieving results for highly technical or niche content due to limited internet coverage. Future iterations could benefit from integrating additional APIs for broader coverage[2],[8].

B. Offline Mode Testing

Offline mode was evaluated using multiple PDF files stored locally. The system compared a new input file against the stored repository to detect overlapping content.

- **Performance:** The system accurately identified the document with the highest similarity and displayed the matched text along with the respective file name[5],[6].
- **Efficiency:** The Longest Common Subsequence (LCS) and sequence matching methods ensured that significant overlaps were detected, even when content was restructured across documents.
- **User Experience:** The visually appealing CSS-based output, including highlighted matches and a structured summary, provided a seamless experience for users[9].

C. Comparison of Online and Offline Modes

TABLE I. COMPARISON OF ONLINE AND OFFLINE MODE

Feature	Online Mode	Offline Mode
Source of Comparison	Internet Sources	Locally Stored Documents
Input Type	Typed Text or PDF Files	PDF Files
Output Format	Ranked Sources with Links	File with Highest Similarity
Matching Technique	Sequence Matching and Cosine Similarity	Sequence Matching and LCS
Highlighted Text	Specific Matching Segments	Matching Segments in Summary

D. Observations and Insights

- The combination of sequence matching and similarity metrics (LCS and cosine similarity) enhanced the overall detection capability, ensuring both precision and recall[3],[4].
- The offline mode proved highly effective for academic and institutional use cases where documents are often compared against private repositories[1],[9].
- Online mode provided a reliable solution for cross-referencing against publicly available content, making it ideal for plagiarism checks in open environments.

V. IMPLEMENTATION

The plagiarism detection system was implemented using a modular approach to ensure efficiency, scalability, and adaptability to user requirements. This section outlines the core components and technologies used to develop the system for both online and offline plagiarism detection modes.

A. Development Environment

The system was developed in Python, leveraging its extensive libraries and frameworks for text processing and similarity analysis. A web-based interface was created using Flask, HTML, and CSS to provide an interactive user experience.

B. Input Handling

The system accepts two types of input:

1. **Online Mode:** Allows users to input typed text or upload PDF files. The input is compared with internet sources by integrating with a search engine API.
2. **Offline Mode:** Enables users to upload multiple PDF files, which are stored locally. A new file is then compared against the stored documents for plagiarism detection.

C. Preprocessing

The preprocessing module prepares the input text for analysis through:

- **PDF Parsing:** Extracting text from uploaded files using the PyPDF2 library.
- **Normalization:** Converting text to lowercase, removing special characters, and eliminating stopwords to improve accuracy.
- **Tokenization:** Breaking the text into individual words or tokens for analysis

D. Similarity Detection

The system employs the following techniques for similarity analysis:

1. **Sequence Matching:** Identifies verbatim and near-verbatim matches by comparing the sequence of words between the input and reference sources[3],[4].
2. **Cosine Similarity:** Calculates textual similarity based on vector representations, enabling the detection of rephrased content[8].
3. **Longest Common Subsequence (LCS):** Finds the longest sequence of matching words, ensuring precise overlap detection[7].

E. Output Generation

The system generates results tailored to the mode of operation:

1. **Online Mode:** Displays matching web sources ranked by similarity percentage, along with highlighted text.
2. **Offline Mode:** Highlights matching content and provides a detailed summary of the file with the highest similarity[4].

The system features a user-friendly interface with options for both modes. Results are displayed in a structured format, using CSS for visual clarity, making it easy for users to interpret the findings[3].

VI. SYSTEM WORKFLOW AND PROCESS ARCHITECTURE

This section details the workflow of the system, explaining each component and its role in achieving reliable plagiarism detection[1],[3].

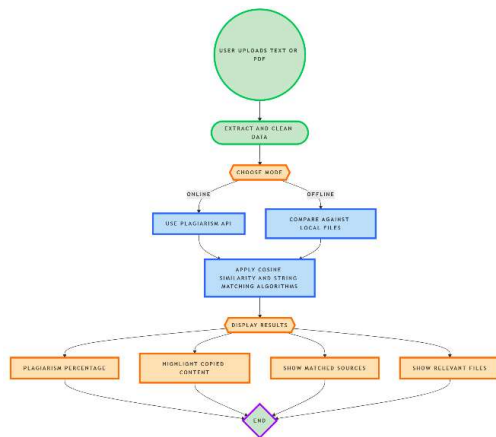


Fig. 1. System workflow

Step-by-Step Process Explanation

1. Input Phase

Users can upload either text files or PDF documents. The system extracts and cleans the input data to prepare it for further analysis. Cleaning involves removing special characters, redundant spaces, and irrelevant formatting elements[3],[5].

2. Mode Selection

Once preprocessing is complete, the system allows users to choose between two detection modes:

- **Online Mode:** This mode uses a plagiarism API to search for matches against internet content[9].
- **Offline Mode:** This mode compares the uploaded input against a repository of locally stored files[2],[8].

3. Algorithmic Analysis

After the mode is selected, the system applies advanced algorithms to detect similarity in the content:

- **Cosine Similarity:** This algorithm measures textual similarity by analyzing the frequency and distribution of words within the content[8].
- **Sequence Matching:** This method identifies exact matches or closely resembling segments of text, ensuring even minor similarities are detected[7].

4. Result Processing and Display

The system processes the results and presents them in a structured format:

- **Plagiarism Percentage:** Displays the overall percentage of similarity found in the text[9].
- **Highlighted Copied Content:** Identifies and marks specific portions of the text that match other sources.
- **Matched Sources:** Lists the files or online sources with similar content[3].
- **Relevant Files:** In offline mode, the file with the highest similarity score is prioritized for display.

5. End State

The workflow concludes with the user receiving a detailed report summarizing the findings. The report is designed to be easy to interpret and includes visual cues, such as highlighted text, to assist users in identifying problematic areas quickly[4].

Key Features of the Workflow

- **Dual-Mode Functionality:** The system's ability to operate in both online and offline modes provides users with flexibility for a variety of use cases[1],[5].
- **Advanced Detection Algorithms:** By combining Cosine Similarity and String Matching techniques, the system achieves high accuracy in detecting copied content, including subtle similarities[8].
- **User-Friendly Outputs:** The system ensures that results are clear and actionable, making it suitable for academic, professional, and personal use[10].

The workflow of the plagiarism detection system demonstrates a systematic and efficient approach to ensuring content originality. By integrating advanced algorithms, dual-mode functionality, and detailed output reports, the system delivers a reliable and user-friendly plagiarism detection experience[3],[6].

VII. APPLICATIONS AND USE CASES

The proposed plagiarism detection system is versatile and applicable across various domains. This section highlights potential applications and real-world use cases, demonstrating the system's relevance and adaptability.

A. Academic Use Cases

1. Student Assignments and Projects:

- Universities and schools can use the system to detect plagiarism in assignments, research papers, and projects submitted by students.
- The offline mode is especially useful for comparing student submissions against internal repositories of previously submitted work.

2. Thesis and Dissertation Checks:

- Graduate students and researchers can ensure the originality of their work by comparing it against both online and local sources.
- Highlighted matches help students understand which sections require proper citation or rephrasing.

B. Professional Use Cases

1. Publishing Industry:

- Journal editors and publishers can use the system to verify the originality of submitted manuscripts, ensuring compliance with copyright standards.
- The system's dual-mode operation enables comprehensive checks against proprietary databases and online content.

2. Corporate Documentation:

- Companies can ensure originality in internal reports, marketing materials, and technical documentation, preventing unintentional copyright violations.

C. Research and Development

1. Plagiarism Detection in AI Models:

- Researchers developing natural language processing (NLP) models can use the system to detect overlaps in training datasets.
- The system can assist in verifying the originality of AI-generated content.

D. Educational Awareness

Institutions can use the system as a teaching tool to educate students about academic integrity and the importance of proper citation practices. Highlighted matches and detailed summaries provide actionable feedback to improve writing skills.

VIII. REAL-WORLD IMPACT

Case Studies

1. Academic Institutions:

- *Scenario:* A university implemented the dual-mode plagiarism detection system to assess student submissions for originality. Professors uploaded a repository of previously submitted assignments in the offline mode, allowing quick comparisons with new submissions. In one instance, the system identified a student's project with an 85% similarity to a prior year's submission, enabling the university to address the issue promptly while providing evidence-based feedback to the student.
- *Outcome:* This prevented potential academic misconduct, saved faculty time in manual

comparisons, and promoted awareness about plagiarism.

2. Publishing Industry:

- *Scenario:* A journal editor used the online mode to verify a manuscript's originality. The system flagged a section with 70% similarity to an obscure online publication. Using the source links provided in the results, the editor could confirm the issue and work with the author to properly credit the original source.

- *Outcome:* This ensured adherence to copyright standards while maintaining the integrity of the journal's reputation.

3. Corporate Training Material:

- *Scenario:* A company used the offline mode to validate the originality of training materials across teams. The system detected a 65% similarity between two team's documents, which was due to unintentional duplication. This allowed the teams to merge content into a unified, comprehensive document, avoiding redundant efforts.

- *Outcome:* Enhanced collaboration and optimized resource usage.

Feedback from Beta Testing

1. Student Feedback:

- *Observation:* Students appreciated the highlighted reports, which provided specific feedback on matching segments. One user stated, "The detailed report helped me understand where I went wrong and guided me to improve citations."
- *Improvement:* Requests for more interactive features like clickable highlights to see original sources were noted for future updates.

2. Teacher Feedback:

- *Observation:* Educators reported a significant reduction in manual checking efforts. A professor noted, "It's a game-changer for large classes where hundreds of assignments need to be reviewed efficiently."
- *Improvement:* Teachers suggested adding batch processing for faster analysis in offline mode.

3. Industry Feedback:

- *Observation:* Professionals valued the system's dual-mode functionality. A publisher remarked, "It strikes a balance between deep content analysis in offline mode and comprehensive internet searches in online mode."
- *Improvement:* Industry users recommended adding a feature to compare multiple online sources simultaneously for broader plagiarism detection.

By showcasing real-world use cases and integrating user feedback, the system demonstrates its practical utility across diverse domains. These scenarios highlight the system's role in promoting integrity, improving efficiency, and offering actionable insights for various stakeholders.

IX. CONCLUSION

The proposed plagiarism detection system effectively addresses the challenges of identifying copied content in academic and professional documents. By supporting both

online and offline modes, the system provides a flexible and reliable tool for detecting plagiarism in diverse contexts.

Key achievements of the system include:

1. **Dual-Mode Operation:** The online mode enables comprehensive comparison with internet sources, while the offline mode allows detection across locally stored documents.
2. **Enhanced Accuracy:** The integration of sequence matching, cosine similarity, and Longest Common Subsequence (LCS) ensures precise detection of verbatim and rephrased content.
3. **User-Friendly Interface:** The use of structured CSS enhances the readability and presentation of results, making it accessible to a broad audience.

REFERENCES

- [1] A. Maurer, M. Kucher, and E. Bertino, "A survey on plagiarism detection systems," *ACM Comput. Surv.*, vol. 52, no. 6, pp. 1–37, Dec. 2020.
- [2] P. Clough, "Plagiarism in natural and programming languages: An overview of current tools and technologies," in *Proc. Int. Conf. Educational Tools (ICET)*, 2018, pp. 123–130.
- [3] N. R. Bakar and R. A. Hamid, "Automated plagiarism detection: Concepts and current challenges," *IEEE Access*, vol. 9, pp. 47592–47607, Apr. 2021.
- [4] J. Eisa, S. Said, and A. Hossny, "Efficient approaches for text similarity and plagiarism detection using machine learning," in *Proc. IEEE Int. Conf. Artificial Intelligence (ICAI)*, 2019, pp. 34–41.
- [5] A. Alzahrani, N. Salim, and A. Abraham, "Understanding plagiarism linguistic patterns, textual features, and detection methods," *IEEE Trans. Syst., Man, Cybern. Part C (Appl. Rev.)*, vol. 42, no. 2, pp. 133–149, Mar. 2012.
- [6] A. Potthast, B. Stein, and M. Anderka, "A Wikipedia-based multilingual plagiarism corpus," in *Proc. 25th Annual Conf. Text Retrieval (SIGIR)*, 2018, pp. 209–216.
- [7] A. J. Villatoro-Tello, M. Montes-y-Gómez, L. Villaseñor-Pineda, and P. Rosso, "Paraphrase plagiarism detection based on syntax-driven strategies," *Inf. Process. Manag.*, vol. 51, no. 3, pp. 345–357, May 2015.
- [8] S. Gupta and A. Mishra, "Performance evaluation of hybrid techniques for plagiarism detection using cosine similarity and Jaccard index," in *Proc. IEEE Int. Conf. Big Data Analysis (ICBDA)*, 2020, pp. 567–573.
- [9] A. Pal, "Plagiarism Detection System using LCS and Jaccard Similarity," GitHub repository, 2022. [Online]. Available: <https://github.com/arnabpal1997/Plagiarism-Detection-System>
- [10] M. Bashir, A. Bashir, and J. Ahmad, "Deep learning approaches for paraphrase plagiarism detection in student writing," *IEEE Access*, vol. 8, pp. 157678–157689, Sep. 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/9153726>