International Journal of Science and Research (IJSR) ISSN: 2319-7064

Impact Factor 2024: 7.101

# Orgina-Academic Integrity Checker

Rajalekshmi R<sup>1</sup>, Jogimol Joseph<sup>2</sup>

<sup>1</sup>Department of Computer Applications, A P J Abdul Kalam Technological University, Musaliar College of Engineering and Technology, Malayalappuzha, Pathanamthitta, Kerala, India Email: rajalekshmir600[at]gmail.com

<sup>2</sup>Professor, Department of Computer Applications, A P J Abdul Kalam Technological University, Musaliar College of Engineering and Technology, Malayalappuzha, Pathanamthitta, Kerala, India

Abstract: ORGINA is a web-based application designed to streamline the process of assignment submission, plagiarism detection, and performance analysis for students. The system aims to provide an efficient platform for students to register, upload assignments, and check for plagiarism. The application will enable students to receive detailed reports on the level of plagiarism in their work, aiding them in improving the quality and originality of their assignments. For administrators, the system provides functionalities to manage subjects, create assignment topics, and analyze students' submitted assignments. Through advanced plagiarism detection techniques, including tokenization, TF-IDF feature extraction, and cosine similarity calculation, the system can effectively detect and report similarities between documents. Additionally, the system will use the efficient Rabin-Karp algorithm for text matching to ensure high accuracy in plagiarism checking algorithms will be tested, and the best-performing one will be integrated into the final system. With a user-friendly interface and a robust back-end built with Django, this application will help improve the academic integrity and performance tracking of students, providing both students and administrators with valuable insights into assignment submissions.

Keywords: Academic integrity, plagiarism detection, Django, TF-IDF, Cosine similarity, Rabin-Karp algorithm

#### 1. Introduction

A core tenet of education is academic integrity, which ensures students cultivate original thought and uphold ethical research practices. However, the proliferation of digital resources has greatly increased the challenge of plagiarism in academic institutions. To address this, ORGINA – Academic Integrity Checker has been developed as a Student Assignment Analysis System. This system facilitates assignment submission, plagiarism detection, and performance analysis. This web-based tool enables students to verify the originality of their work prior to submission, thus promoting ethical writing and learning. It incorporates advanced plagiarism detection techniques, including the Rabin-Karp stringmatching algorithm, to efficiently identify textual similarities and ensure accurate detection.

The system also provides administrative tools that enable faculty to manage assignments, review plagiarism reports, and assess student performance through detailed analytics. This enhances transparency in academic evaluations and helps to maintain institutional integrity. Developed using Django, ORGINA is designed for scalability, security, and ease of use, providing a seamless experience for both students and administrators. By leveraging data-driven insights, the system supports a fair and accountable academic environment, fostering originality and improving learning outcomes.

## 2. Literature Survey

Existing plagiarism detection systems predominantly rely on document-matching techniques, such as cosine similarity, to identify similarities between texts [8]. These systems compare submitted documents against a database of existing sources, including online resources, previously submitted papers, and academic publications. While effective in detecting direct copying, they often struggle with identifying more sophisticated forms of plagiarism, such as paraphrasing and content re-arrangement. Many of these systems are proprietary software solutions, including Turnitin and Copyscape, which require subscription fees, making them less accessible for institutions with budget constraints. These tools are widely used and have become a standard in many academic settings, but they can be costly for both institutions and students. While effective in detecting copied content, they often lack the flexibility needed to accommodate custom requirements, such as integrating with academic management systems or providing detailed student performance insights. Several open-source plagiarism detection tools are available, offering cost-effective alternatives. These solutions often provide basic text-matching capabilities but may lack the advanced features and accuracy of commercial systems. However, these solutions are often limited in scope, primarily focusing on basic text-matching without the ability to function as a fully integrated educational platform [6,7]. Many of these tools do not include features such as assignment submission tracking, subject management, or comprehensive performance analysis. As a result, educational institutions require a more adaptable and robust plagiarism detection system that not only ensures originality in academic work but also enhances the overall learning experience through detailed analytics and seamless integration with academic workflows. The following published articles have been referred to create a base for the project: -

This project builds upon existing research in plagiarism detection, drawing on a range of techniques and tools. Mirac Suzgun, Stuart M. Shieber, and Dan Jurafsky (2023) introduced string2string, a Python library for efficient string-to-string matching [1], which is fundamental to identifying similarities and differences in text sequences. Weiran Wang, Diamantino Caserio, et al. (2023) explored enhancing contextual biasing in Automatic Speech Recognition (ASR) systems using the Knuth-Morris-Pratt (KMP) algorithm [2], a

# International Journal of Science and Research (IJSR) ISSN: 2319-7064 Impact Factor 2024: 7.101

linear time string-matching algorithm that can be adapted to improve the accuracy of plagiarism detection by efficiently locating specific patterns. Anjali Bohra & N. C. Barwar (2022) presented a deep learning-based approach using BERT (Bidirectional Encoder Representations from Transformers) [3], a model that captures contextual relationships between words for nuanced understanding of semantic similarity and paraphrased content. Ora Amir, Amihod Amir, Aviezri Fraenkel, and David Sarne (2022) examined the efficiency of automata-based approaches, such as the KMP algorithm, in pattern matching, emphasizing their importance for handling large-scale text comparisons [4]. A.A. Adewoyin (2021) discussed the limitations of traditional plagiarism detection tools in identifying paraphrased and translated content and proposed integrating advanced Natural Language Processing (NLP) techniques to address these limitations [5]. Dita Baitu Rahmawati and Muhammad Luthfi Irfani (2020) explored the use of the Rabin-Karp algorithm for plagiarism detection in e-learning systems [6], highlighting its efficiency in detecting exact text matches. Brinardi Leonardo and Seng Hansun (2017) combined the Rabin-Karp algorithm with the Jaro-Winkler distance algorithm to detect both exact and nearduplicate content [7], while Robbi Rahim and Dodi Siregar (2017) investigated the use of K-grams in conjunction with the Rabin-Karp algorithm to enhance plagiarism detection [8]. Sonawane Kiran Shivaji and Prabhudeva S (2015) also focused on improving accuracy by combining the Karp-Rabin algorithm with other string-matching techniques [9]. Vijayarani Mohan and Tamilarasi Angamuthu (2015) introduced an enhanced Knuth-Morris-Pratt (KMP) algorithm designed to improve the accuracy and efficiency of text retrieval from desktop systems [10]. Thailambal Rajagopalan and S. Arumugam (2014) evaluated the performance and efficiency of several string-matching algorithms, including KMP, Rabin-Karp, and Boyer-Moore, in text mining applications [11]. Jiffriya, M. A. C. Akmal Jahan, R. G. Ragel, and S. Deegalla (2014) presented AntiPlag, a plagiarism detection tool that uses tri-gram sequence matching [12], and P. Preethi and S. K. Srivatsa (2012) provided a comparative study of string-matching algorithms for plagiarism detection [13]. Kevin Mote (2012) provided a comprehensive survey of Natural Language Processing (NLP) [14], covering its techniques, challenges, and applications, including plagiarism detection. Bela Gipp and Jöran Beel (2010) presented a detailed review of various plagiarism detection methods [15], categorizing them and discussing their strengths and weaknesses.

# 3. Methodology

Orgina – Academic Integrity Checker adopts a methodical and data-driven approach to ensure accurate plagiarism detection and uphold academic integrity. The process begins with the collection of user-submitted assignments, ensuring real-world authenticity and diversity in data patterns. Unlike traditional systems that depend on publicly available datasets, Orgina operates exclusively on actual academic submissions, enhancing its applicability in educational environments.

The preprocessing stage involves tokenization and the removal of stop words, transforming raw text into analyzable units. Feature extraction is carried out using the Term Frequency–Inverse Document Frequency (TF-IDF) technique, which quantifies the relevance of words in

individual documents relative to the entire collection. This numerical representation enables efficient and meaningful comparison between texts.

To assess similarity, Cosine Similarity is used. This technique calculates the cosine of the angle between vectorized documents, providing a normalized similarity score between 0 and 1. For precise detection of exact matches, the system integrates the Rabin-Karp algorithm, which leverages rolling hash functions to identify duplicate sequences with high efficiency.

Plagiarism reports are automatically generated upon assignment submission. These reports contain a similarity percentage, visual highlights of matching content, and textual explanations, providing valuable feedback for administrators. The modular nature of the system ensures scalability and seamless integration with academic workflows. Continuous updates to detection algorithms and preprocessing techniques keep Orgina responsive to evolving academic practices and plagiarism strategies.

This robust, hybrid methodology—combining traditional string-matching with advanced NLP and vector-based analysis—positions Orgina as a reliable and adaptable tool for plagiarism detection in academic institutions.

## 3.1 Algorithms in Orgina

In Orgina, the **Rabin-Karp algorithm** serves as a key mechanism for identifying copied text by scanning and comparing sequences of characters within documents. It is a hash-based pattern matching technique designed to efficiently locate exact or near-exact matches, making it highly effective for academic plagiarism detection. Its use is essential in scenarios where large volumes of student submissions need to be analyzed rapidly and accurately.

The algorithm simplifies the process of detecting duplicated content by converting both the input and target text fragments into hash values. Rather than directly comparing each character in a string, it compares the hash values—allowing for much faster processing. Only when matching hashes are found does the system proceed to a character-by-character comparison to confirm the match, ensuring accuracy and reducing the chances of false positives.

1) Initial Text Preparation (Preprocessing & **Tokenization**): Before any analytical algorithms are applied, each submitted document undergoes essential preparation to standardize the text and isolate meaningful units. This involves programmatic cleaning, such as converting all text to lowercase for consistency and stripping out punctuation or special characters that could interfere with analysis. The cleaned text is then tokenized-broken down into fundamental components like individual words or specific word sequences (kgrams). A critical part of this stage is the removal of "stop words" (e.g., "and", "the", "of"), which are common words that generally lack significant distinguishing value for similarity comparisons. This rigorous preparation ensures that subsequent analytical steps operate on a clean, consistent, and meaningful representation of the

document's content.

- 2) Content Vectorization (Feature Extraction via TF-IDF): To enable mathematical comparison between documents, the tokenized text must be transformed into a numerical format. Orgina employs the Term Frequency-Inverse Document Frequency (TF-IDF) technique for this conversion. TF-IDF assigns a numerical weight to each token (word) in a document, reflecting its importance not just within that single document but across the entire collection of documents being analyzed.
- 3) **Term Frequency (TF):** This component measures how frequently a specific term appears within the document being analyzed, often normalized to account for document length.

The formula used is:

TF = (Number of times a word appears in the document) / (Total number of words in the document).

- 4) Inverse Document Frequency (IDF): This component assesses the rarity of a term across all documents in the comparison corpus. Terms that appear in many documents (like common jargon) receive lower weights, while rarer, more specific terms receive higher weights. The formula is: IDF = log((Total number of documents) / (Number of documents containing the word)).
- 5) The final **TF-IDF score** for a term is the product of its TF and IDF values. This process generates a high-dimensional vector for each document, where each dimension corresponds to a unique term, and the value in that dimension represents the term's calculated TF-IDF weight. This vector quantitatively captures the document's key lexical content.
- 6) Conceptual Similarity Assessment (Cosine Similarity): To gauge the degree of similarity in content and theme, particularly useful for detecting paraphrased or reworded sections, Orgina utilizes Cosine Similarity. This geometric measure calculates the cosine of the angle between the TF-IDF vectors representing two documents (let's call them A and B). The formula applied is: Cosine Similarity ( $\cos(\theta)$ ) = (A  $\cdot$  B) / (||A|| \* ||B||).

Here,  $A \cdot B$  represents the dot product of the two vectors, while ||A|| and ||B|| are their respective magnitudes (or lengths). A result close to 1 signifies a very small angle between the vectors, indicating strong similarity in their TF-IDF profiles (and thus, likely content overlap), whereas a result close to 0 suggests orthogonality or dissimilarity. This method effectively compares the overall "direction" or thematic focus of documents, irrespective of their absolute lengths.

7) Exact Sequence Identification (Rabin-Karp Algorithm): Complementing the broader similarity measure, the Rabin-Karp algorithm is specifically integrated to efficiently detect instances of direct, verbatim text copying. This algorithm employs a "rolling hash" function - a clever technique that allows for the rapid calculation of hash values for successive segments (substrings) of text without needing to reprocess the entire segment each time the window shifts. It compares the hash value of a segment from the submitted document against the hash values of segments from documents in the comparison database. If the hash values match, indicating a potential exact match, the algorithm then performs a direct, character-by-character comparison of the segments to confirm true equivalence and rule out coincidental hash collisions. This two-phase approach (hash comparison followed by verification) makes Rabin-Karp particularly efficient for scanning large texts for identical passages.

- Linguistic Pattern Analysis (NLP-Based Methods): 8) To enhance the system's ability to detect more nuanced forms of potential plagiarism that might evade lexical or exact-match techniques, Orgina incorporates Natural Language Processing (NLP) methods. While specific models aren't exhaustively detailed, the methodology includes fundamental NLP preprocessing steps like tokenization, stopword removal, and potentially lemmatization (reducing words to their base or dictionary form). Advanced techniques mentioned include Named Entity Recognition (NER), which identifies key entities like names or locations, and the use of Word Embeddings, which represent words as vectors capturing semantic relationships and context. These NLP techniques allow the system to analyze text based on deeper linguistic structure and meaning, improving its capacity to identify sophisticated paraphrasing or conceptual borrowing.
- Integrated Analysis 9) Framework (Synergizing Algorithms): The effectiveness of Orgina's methodology stems from the strategic combination of these diverse algorithms. No single algorithm excels at all forms of plagiarism detection. Rabin-Karp offers speed and precision for exact duplicates. TF-IDF with Cosine Similarity effectively identifies thematic overlap and significant paraphrasing by comparing weighted term distributions. NLP methods provide a layer capable of understanding semantic nuances and complex linguistic alterations. By integrating these techniques, Orgina aims to create a more comprehensive detection net, capturing a wider spectrum of potential academic integrity issues than single-algorithm systems.
- 10) Quantification, Documentation, and Reporting (Similarity Scoring, Storage & Report Generation): Following the multi-faceted analysis, the findings are consolidated and quantified. The system calculates an overall similarity score or percentage, representing the proportion of the submitted document identified as matching content within the comparison database. Specific details about the matches, such as the text segments identified as similar or identical and references to the potential source documents (within the system's accessible database), are recorded and stored. This information is then synthesized into a structured plagiarism report. For administrators, this report provides a detailed breakdown, including the overall similarity score, visual highlighting of the matched sections within the submitted text, and links or references to the corresponding source materials. While students may see their overall similarity rate to encourage self-assessment, the detailed analytical report serves as a crucial tool for educators to evaluate originality and make informed decisions regarding academic integrity

## International Journal of Science and Research (IJSR) ISSN: 2319-7064 Impact Factor 2024: 7.101



Figure 1: Architecture

#### 3.2 Dataset Used

Our project relies on user-submitted text documents rather than publicly available datasets to ensure a more practical and realistic plagiarism detection system. Public datasets often contain pre-labelled or artificially generated data that may not accurately represent real plagiarism cases. By using real documents, the system becomes more effective in real-world scenarios. This approach also allows for better testing and improvement of plagiarism detection algorithms on diverse and unpredictable data. Additionally, using user-submitted documents ensures compliance with data privacy regulations, avoiding legal risks associated with copyrighted content. Unlike many public datasets that primarily focus on simple copy-paste detection, our system is designed to identify more complex plagiarism, including paraphrasing and structural modifications. Since the project is intended for educational institutions, it closely mimics professional plagiarism detection tools like Turnitin, making it highly relevant for academic settings. To enhance detection accuracy, our system integrates text preprocessing techniques, similarity detection methods such as Rabin-Karp and Cosine Similarity, and NLP- based approaches to efficiently compare and highlight similarities in documents.

# 4. Result and Discussion

The following Table 1 showcases the performance of ORGINA – Academic Integrity Checker, a hybrid plagiarism detection system, in comparison with two baseline approaches: traditional keyword matching and a standalone machine learning model using TF-IDF with cosine similarity. The evaluation is carried out using standard performance metrics such as accuracy, precision, recall, and false positive rate (FPR).

Table 1. Comparison of Fragiansin Detection Teeninques			
Metric	ORGINA	Basic Text	Standalone ML
		Matching	Classifier
Accuracy	95.7%	82.5%	89.6%
Precision	94.2%	78.0%	87.1%
Recall	96.3%	80.5%	88.0%
False Positive Rate	4.2%	11.4%	7.2%

Table 1: Comparison of Plagiarism Detection Techniques

The graphical representation based on above model that is given below. Here we can see that the ensemble guard completely outperform.



Figure 2: Graphical representation

The quantitative results clearly demonstrate the superior performance of the Orgina system compared to both traditional keyword matching and a single ML-based classifier approach. With significantly higher scores in Accuracy (95.7%), Precision (94.2%), and Recall (96.3%), Orgina proves highly effective in its detection capabilities. Furthermore, its substantially lower False Positive Rate (4.2%) underscores its reliability and efficiency, minimizing incorrect identifications compared to the baseline methods.

# 5. Conclusion

Orgina represents a significant and practical contribution to addressing the persistent challenge of academic dishonesty in the digital era. This project successfully culminates in a comprehensive, web-based system meticulously designed to streamline the entire lifecycle of assignment management, from submission through to robust plagiarism analysis. It serves the dual needs of students seeking to verify the originality of their work and administrators requiring reliable tools to uphold academic standards. The platform provides an intuitive interface, built with standard web technologies like

## International Journal of Science and Research (IJSR) ISSN: 2319-7064 Impact Factor 2024: 7.101

HTML, CSS, and JavaScript, running on a scalable backend powered by the Django framework, ensuring both usability and maintainability.

The core innovation of Orgina lies in its sophisticated, hybrid approach to plagiarism detection. Rather than relying on a single technique, it strategically integrates multiple algorithms to cast a wider net. The implementation leverages TF-IDF for nuanced feature extraction from text, enabling Cosine Similarity calculations to effectively identify semantic and contextual overlaps indicative of paraphrasing. Simultaneously, the computationally efficient Rabin-Karp algorithm is employed to rapidly pinpoint instances of exact, verbatim copying. This synergistic combination allows Orgina to detect a broader spectrum of potential plagiarism than systems limited to simpler matching methods.

Furthermore, Orgina's design thoughtfully incorporates practical considerations for academic environments. Its reliance on analyzing user-submitted documents creates a dynamic and relevant comparison corpus, mirroring realworld institutional workflows and mitigating concerns associated with static or potentially copyrighted public datasets. Data persistence and management are efficiently handled by an SQLite database backend. For administrators, the system offers detailed reports highlighting similarity percentages and specific matched content, facilitating fair and evidence-based evaluation. While students receive feedback on their submission's similarity rate to encourage selfcorrection, the comprehensive analysis remains a key tool for institutional oversight.

In essence, Orgina provides a functional, scalable, and relevant technological solution aimed directly at reinforcing academic integrity. By automating and enhancing the detection of plagiarism through intelligent algorithm integration, it supports educators in maintaining academic standards and fosters an environment where original work is valued and encouraged. It stands as a testament to the potential of applying targeted data mining and text analysis techniques to solve pressing challenges within the educational domain.

# References

- [1] Suzgun, Mirac; Shieber, Stuart M.; Jurafsky, Dan. string2string: A Modern Python Library for String-to-String Algorithms. 2023.
- [2] Wang, Weiran; Caserio, Diamantino; et al. Contextual Biasing with the Knuth-Morris-Pratt Matching Algorithm. 2023.
- [3] Bohra, Anjali; Barwar, N. C. A Deep Learning Approach for Plagiarism Detection System Using BERT. 2022.
- [4] Amir, Ora; Amir, Amihod; Fraenkel, Aviezri; Sarne, David. On the Practical Power of Automata in Pattern Matching. 2022.
- [5] Adewoyin, A.A. Integrating State-of-the-art NLP Tools into Existing Methods to Improve Plagiarism Detection. 2021.
- [6] Rahmawati, Dita Baitu; Irfani, Muhammad Luthfi. Text Mining to Detect Plagiarism in E-Learning System Using Rabin-Karp Algorithm. 2020.

- [7] Leonardo, Brinardi; Hansun, Seng. Text Documents Plagiarism Detection using Rabin-Karp and Jaro-Winkler Distance Algorithms. 2017.
- [8] Rahim, Robbi; Siregar, Dodi. K-Gram as a Determinant of Plagiarism Level in Rabin-Karp Algorithm. 2017.
- [9] Shivaji, Sonawane Kiran; S, Prabhudeva. Plagiarism Detection by Using Karp-Rabin and String-Matching Algorithm Together. 2015.
- [10] Mohan, Vijayarani; Angamuthu, Tamilarasi. An Efficient Text Pattern Matching Algorithm for Retrieving Information from Desktop. 2015.
- [11] Rajagopalan, Thailambal; Arumugam, S. Performance of Multiple String-Matching Algorithms in Text Mining. 2014.
- [12] Jiffriya; Jahan, M. A. C. Akmal; Ragel, R. G.; Deegalla, S. AntiPlag: Plagiarism Detection on Electronic Submissions of Text-Based Assignments. 2014.
- [13] Preethi, P.; Srivatsa, S. K. A Comparative Study of String-Matching Algorithms for Plagiarism Detection. 2012.
- [14] Mote, Kevin. Natural Language Processing A Survey. 2012.
- [15] Gipp, Bela; Beel, Jöran. A Survey of Plagiarism Detection Methods. 2010.