**Research Article**

# Enhancing Regional Plagiarism Detection Using a Backtrack Matching Model: A Precision, Recall, and F1 Score-Based Evaluation

Mr. Prashanth Kumar HM[1] & Prof. Subrahmanya Bhat[2]

[1]*Student, College of Computer Science, Srinivas University, Mangalore, India*

[2] *Professor, College of Computer Science, Srinivas University, Mangalore, India*

*E-mail: prashanth.hm02@gmail.com, itsbhat@gmail.com*

| ARTICLE INFO | ABSTRACT |
|---|---|
| | In the era of digital content creation and academic publishing, ensuring the originality of written work has become crucial. Plagiarism detection tools have evolved to combat this issue, but their effectiveness varies based on regional linguistic nuances, such as local idioms, translations, and cultural expressions. This paper presents an evaluation framework for a regional plagiarism checker (Using Hindi, Marathi, Tamil, Gujarati, etc....), assessing its performance using key metrics: precision, recall, and F1 score. Precision measures the accuracy of identified plagiarism instances, while recall evaluates the tool's ability to detect all actual instances of plagiarism. The F1 score provides a harmonic meaning, balancing both precision and recall. The evaluation highlights the challenges of detecting plagiarism in regional contexts and comparison, where language variations can lead to detect or undetected content from similarity. Our results show that by tailoring plagiarism detection algorithms to regional linguistic characteristics, it is possible to significantly improve detection accuracy and reduce false positives, thus providing a more effective and reliable tool for regional academic integrity.<br><br>**Keywords:** Pression, Recall, F1 Score, Accuracy, Plagiarism, Integrity. |

## 1. INTRODUCTION

Plagiarism, the act of using someone else's work or ideas without proper attribution, has become a major concern in academic, professional, and creative writing. With the proliferation of digital content, access to vast amounts of information has made it easier for individuals to commit plagiarism, whether intentionally or unintentionally. As a result, the development of plagiarism detection tools has gained increasing importance in educational institutions, publishing, and research organizations. These tools help ensure the integrity of content, maintain academic honesty, and foster creativity by safeguarding intellectual property. However, many existing plagiarism detection systems are designed with global or generic language models that do not account for regional language variations. This often leads to reduced detection accuracy in specific cultural or linguistic contexts. Therefore, the need for regionally tailored plagiarism checkers has become evident.

DrillBit regional plagiarism checker is a tool designed to specifically address the linguistic, cultural, and regional nuances that are often missed by general plagiarism detection systems. These nuances may include local expressions, idiomatic phrases, cultural references, and differences in language structure that are unique to specific geographic areas. By developing a DrillBit anti-plagiarism detection system that is sensitive to these regional variations, the goal is to enhance the tool's ability to detect plagiarism more accurately and efficiently in localized contexts.

The effectiveness of plagiarism detection tools is generally measured through three key performance metrics: precision, recall, and F1 score. Precision measures the accuracy of the tool in identifying true positives (i.e., actual cases of plagiarism) versus false positives (i.e., cases where the tool incorrectly flags content as plagiarized). Recall, on the other hand, assesses the tool's ability to detect all actual instances of plagiarism, capturing the rate at which true positives are correctly identified. The F1 score is the harmonic means of precision and recall, providing a balanced measure of the system's overall performance. Together, these metrics allow researchers and developers to

**Research Article**

evaluate the strengths and weaknesses of plagiarism detection tools in different contexts, particularly in regional scenarios where language complexities may affect performance.

Academic institutions under the UGC implemented strict policies to combat plagiarism, often incorporating the use of plagiarism detection tools. These tools have become an essential part of maintaining academic integrity, enabling educators to verify the originality of students' work and researchers to ensure that their publications are free from unethical content duplication. In response to the growing problem of plagiarism, software companies have developed a wide range of plagiarism detection systems, each with varying levels of effectiveness.

## 2. INSIGHTS

Until December 2024, we relied on a generic plagiarism detection tool to check for regional plagiarism. However, in early 2025, we implemented a new "Backtrack Matching Model" specifically designed to handle regional plagiarism, with a focus on improving accuracy. This model is optimized for Indian literature and supports a variety of regional languages, addressing the linguistic nuances often overlooked by generic tools. The Backtrack Matching Model has performed as expected, providing exact matching capabilities and accurately identifying the source locations of plagiarized content.

Compared to the previous generic detection tool, the new model has resulted in a 30-35% improvement in accuracy. This significant enhancement was measured and analysed using key performance metrics using Precision, Recall, and F1 Score to ensure comprehensive evaluation. Precision allowed us to assess the model's accuracy in detecting true positives, while Recall measured its ability to identify all instances of plagiarism. The F1 Score provided a balanced view of both metrics, reflecting the overall effectiveness of the new model. This improvement demonstrates the effectiveness of tailoring plagiarism detection tools to regional linguistic and cultural contexts, resulting in a more reliable and accurate plagiarism detection system.

## 3. THE IMPORTANCE OF PRECISION, RECALL, AND F1 SCORE

**Precision:**

Precision measures the accuracy of the plagiarism detection tool by calculating the proportion of correctly identified plagiarism instances (true positives) out of all instances flagged as plagiarism (true positives + false positives). High precision means that the tool is accurate in identifying plagiarism without mistakenly flagging original content as plagiarized. In regional contexts, precision is crucial because the tool must be able to distinguish between actual plagiarism and legitimate use of regional idiomatic expressions or cultural references.

**Recall:**

Recall, also known as sensitivity, assesses the tool's ability to detect all instances of plagiarism. It is calculated as the proportion of true positives (correctly identified plagiarism instances) out of all actual cases of plagiarism (true positives + false negatives). A high recall score indicates that the tool is effective in capturing the full range of plagiarism, even if some cases are more subtle or hidden within regional language variations. In regional plagiarism detection, recall is important to ensure that the system does not miss instances of plagiarism due to local linguistic differences.

**F1 Score:**

The F1 score is the harmonic means of precision and recall, providing a balanced measure of the plagiarism detection system's overall performance. A high F1 score indicates that the tool performs well in both accuracy and completeness, ensuring that it can identify plagiarism while minimizing errors. In regional plagiarism detection, the F1 score is particularly useful for evaluating the tool's ability to handle the complexities of regional languages and dialects.

## 4. EVALUATING REGIONAL PLAGIARISM CHECKERS

The evaluation of a regional plagiarism checker involves a comprehensive analysis of how well the tool performs in detecting plagiarism in localized linguistic contexts. This process typically begins with the development of a region-specific language model that incorporates local idioms, expressions, and cultural references. The model is then

**Research Article**

integrated into the plagiarism detection system and tested on a dataset of texts that include both original content and plagiarized content from the target region.

By applying precision, recall, and F1 score metrics to the evaluation process, researchers can determine the effectiveness of the regional plagiarism checker in comparison to generic plagiarism detection systems. The goal is to achieve a balance between precision and recall, ensuring that the tool accurately identifies instances of plagiarism while minimizing false positives and false negatives. The rise of digital content and the increased accessibility of information have amplified the need for effective plagiarism detection systems. While generic plagiarism detection tools have proven useful in many contexts, they often fall short when applied to regional linguistic variations. By developing regional plagiarism checkers that account for local language and cultural differences, and evaluating their performance using precision, recall, and F1 score, it is possible to improve the accuracy and reliability of plagiarism detection in specific regions. This approach ensures that regional nuances are considered, leading to more effective tools that uphold academic and creative integrity on a global scale.

Limitations of Generic Plagiarism Detection Systems Plagiarism checkers often struggle to detect paraphrasing, especially when content is reworded using local language variations or cultural expressions. Unlike direct copying, paraphrased text may retain the original meaning while altering sentence structure, vocabulary, or idiomatic usage. Generic plagiarism detection tools, and even some regional ones, may fail to recognize these subtle shifts, leading to undetected plagiarism. This is particularly challenging in regions with unique linguistic patterns, where paraphrased content can appear significantly different from its source, making it difficult for the checker to accurately flag such instances without advanced paraphrase detection algorithms.

The checkers often face challenges in detecting plagiarism in multilingual content, particularly in regions where multiple languages are spoken. Plagiarized text can be translated from one language to another, making it difficult for detection systems to identify the original source. Without robust cross-language detection capabilities, these tools may fail to flag instances where content is lifted and translated. This limitation is especially pronounced in multilingual regions, where language diversity complicates the detection of copied material. As a result, plagiarism in translated texts often goes unnoticed, highlighting the need for advanced algorithms to handle multilingual plagiarism effectively.

For example, plagiarism detection systems that are developed primarily for English-speaking regions may struggle to accurately detect plagiarism in languages that use different sentence structures, such as those found in East Asian or Middle Eastern languages. Similarly, content written in regional dialects or vernaculars may be misinterpreted by generic tools, leading to inaccurate results. In light of these challenges, there is a growing need for plagiarism detection systems that are specifically tailored to regional linguistic and cultural contexts.

## 5. OBJECTIVES

**1. Accuracy:** Accuracy in regional plagiarism checkers is crucial to ensure the integrity of academic and professional work. These tools must account for local languages, dialects, and regional publications to detect plagiarism effectively. An accurate checker should recognize nuanced phrasing, regional expressions, and diverse citation styles specific to a region. It must also access local databases, universities' research papers, and smaller publications to provide a thorough analysis. Without such precision, regional variations may go undetected, leading to incomplete or skewed results. Therefore, a well-tuned plagiarism checker significantly enhances the reliability of detecting copied content in regional contexts.
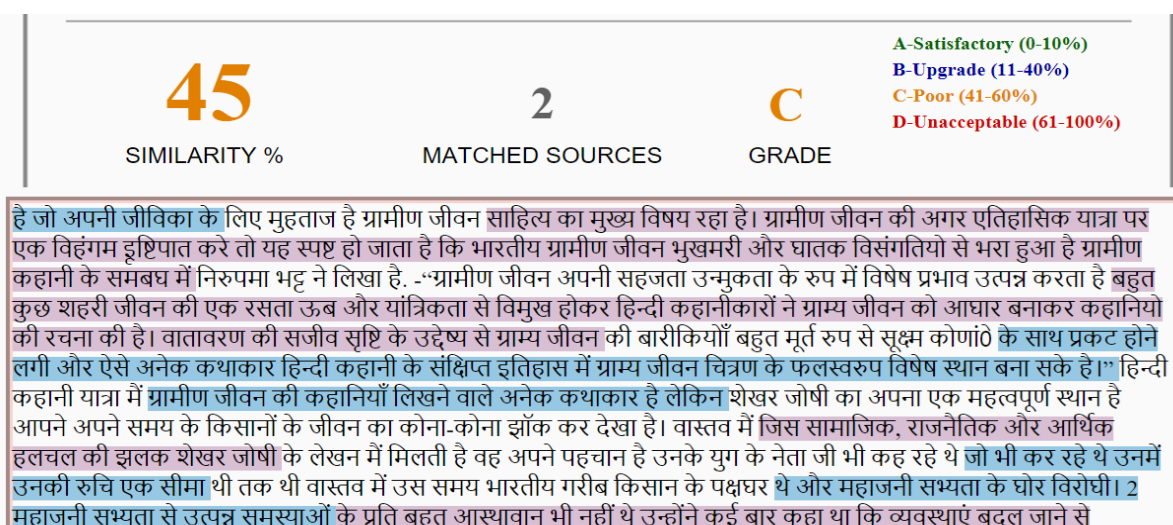
**2. Speed:** The speed at which a regional plagiarism checker generates results is a key factor in its efficiency. Users, especially in academic and professional settings, rely on timely feedback to revise and improve their work. A fast checker saves time and enables quick decision-making. However, speed should not compromise accuracy. A well-optimized plagiarism checker balances swift processing with in-depth analysis, using advanced algorithms and access to local databases for quick comparisons. Regional nuances, like language variations, must be processed efficiently to maintain both speed and accuracy. Fast result generation is essential for meeting deadlines and ensuring workflow productivity.

**Research Article**

**3. Repository System:** A robust repository system is essential for a regional plagiarism checker to deliver accurate and comprehensive results. This system stores a vast collection of local academic papers, research articles, publications, and other region-specific content, allowing the checker to compare submissions against relevant sources. By having access to regional universities' databases, archives, and local publishers, the checker can identify matches and similarities more effectively. A well-maintained repository system ensures continuous updates, expanding its database with new materials over time. This enables the plagiarism checker to stay current and provide precise, regionally tailored plagiarism detection, improving overall accuracy and reliability.
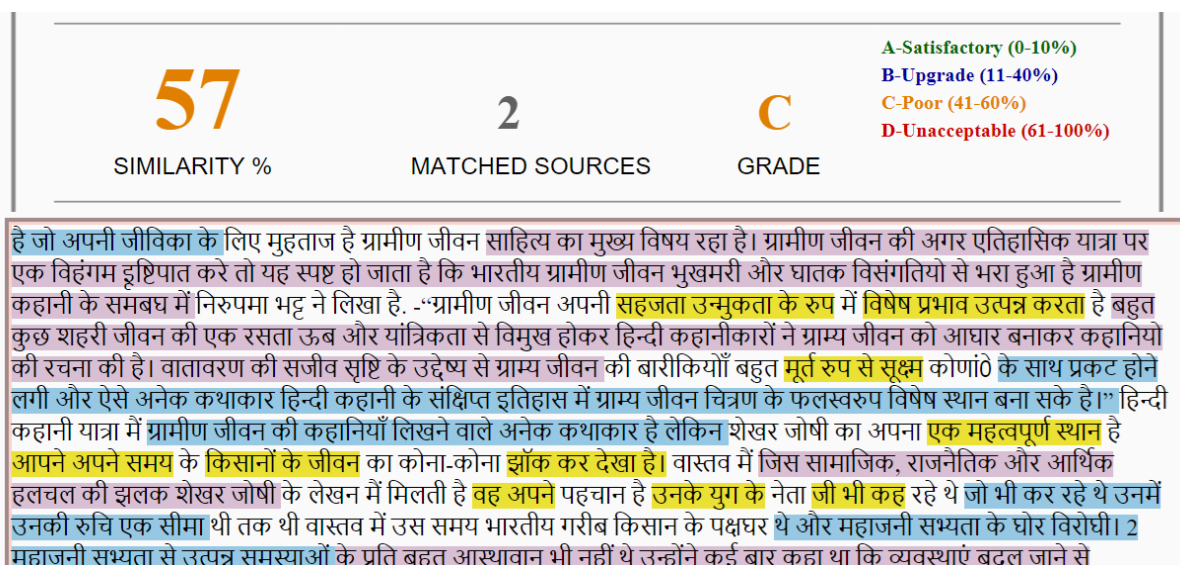
**4. Integrity:** Academic integrity is a fundamental principle that regional plagiarism checkers uphold by ensuring originality in scholarly work. These checkers help maintain ethical standards by detecting unauthorized copying, ensuring students and researchers produce their own work. In regions with specific academic traditions, citation styles, and languages, a regional plagiarism checker is crucial for upholding these values. It identifies even subtle forms of plagiarism, like improper paraphrasing or translation issues. By fostering a culture of honesty and accountability, regional plagiarism checkers play a vital role in reinforcing academic integrity, promoting fairness, and safeguarding the credibility of local educational institutions.

## 6. COMPARISON

As comparing the result after implemented backtracking method then the two below images show plagiarism report summaries of the same text from the same report from the regional language, but with different similarity percentages. In the first image, the similarity index is 45% (before implementing backtrack matching method), while in the second image, it has increased to 57% (after implementing backtrack matching method). Both have the same number of matched sources (2) and are graded C, which corresponds to a Poor (41-60%) rating.



In the below image highlighted sections within the text indicate the portions that have been flagged for similarity. In the above image, fewer segments are highlighted, and the similarity detection is primarily concentrated in blue and purple colours. In contrast, the second image includes more yellow highlights which is detected by backtracking model algorithm, and increased matching with the same source materials, hence the higher similarity percentage.

**Research Article**



The core content in both images appears largely unchanged, but the second image reflects a higher detection of duplicated content, possibly due to new backtracking algorithm detection system. The key difference lies in the increase in similarity percentage from 45% to 57% due to more extensive content being flagged as matching, which may indicate the new implemented algorithm is creating more accuracy in the second version.

## 6. RESULT

**Precision:** The table below forms by comparing previous and present accuracy status, here the time gap is min 6 months and max 2 year, so randomly selected different types of sources. In the below chart we categorized five level of testing called Journal data, live internet data, Student paper, Repository data and UGC Shoudhganga data. The content we tested minimum 100 to 200 pages with respect to words.

| | | | Previous % | | Present % | |
|---|---|---|---|---|---|---|
| Type | Pages | Words | TP | FP | TP | FP |
| Journal | 100 | 25,000 | 65 | 35 | 84 | 21 |
| Internet | 200 | 50,000 | 68 | 33 | 90 | 23 |
| Student paper | 100 | 25,000 | 90 | 08 | 95 | 01 |
| Repository Data | 100 | 25,000 | 92 | 07 | 97 | 01 |
| Shodh ganga | 100 | 25,000 | 90 | 15 | 95 | 08 |
| Total Avg. | | | 81 | 19.6 | 92.2 | 10.8 |

To calculate precision, use the formula: Precision = True Positives / (True Positives + False Positives). So, if we calculate present status Precision=92.2/(92.2+10.8) = 0.895, similarly the value of previous precession value is 0.81. According to the comparison the final precession average value in Previous TP should be lesser than or equal to Present TP values, in the same way Previous FP value is higher or equal to Present FP average value, so 0.81 <= 0.895.

**Recall:** In the comparing of previous and present accuracy status, here the time gap is min 6 months and max 2 year, so randomly selected different types of sources. In the below chart we categorized five level of testing called Journal data, live internet data, Student paper, Repository data and UGC Shoudhganga data. The content we tested minimum 100 to 200 pages with respect to words.

| | | | Previous % | | Present % | |
|---|---|---|---|---|---|---|
| Type | Pages | Words | TP | FN | TP | FN |
| Journal | 100 | 25,000 | 65 | 05 | 84 | 01 |
| Internet | 200 | 50,000 | 68 | 03 | 90 | 02 |
| Student paper | 100 | 25,000 | 90 | 05 | 95 | 01 |

**Research Article**

| Repository Data | 100 | 25,000 | 92 | 02 | 97 | 01 |
|---|---|---|---|---|---|---|
| Shodh ganga | 100 | 25,000 | 90 | 02 | 95 | 01 |
| **Total Avg.** | | | **81** | **3.4** | **92.2** | **1.2** |

To calculate recall, use the formula: Recall = True Positives / (True Positives + False Negative). So, if we calculate present status Precision=92.2/(92.2+1.2) = 0.987, similarly the value of previous precession value is 0.96. According to the comparison the final precession average value in Previous TP should be lesser than or equal to Present TP values, in the same way Previous FN value is higher or equal to Present FN average value, so 0.96 <= 0.987.

**F1 Score:** In this research methodology, the F1 score is a metric used to evaluate the performance of classification models, particularly useful when dealing with imbalanced datasets. It combines precision and recall into a single score using their harmonic mean. The formula for F1 score is: F1 = 2 * (Precision * Recall) / (Precision + Recall).

1. F1 value for Previous data: F1 = 2 * ( 0.81 * 0.96 ) / ( 0.81 + 0.96 ), = 2 * 0.77 / 1.77 = 0.87
2. F1 value for Present data: F1 = 2 * ( 0.895 * 0.987 ) / ( 0.895 + 0.987 ), = 2 * 0.88 / 1.88 = 0.94

The F1 score effectively measures the balance between precision and recall, making it a reliable metric for evaluating classification performance, especially on imbalanced datasets. In this study, the F1 score improved from 0.87 in the previous dataset to 0.94 in the present dataset. This increase indicates a significant enhancement in the model's ability to correctly classify instances. Overall, the results demonstrate improved precision and recall, reflecting a more accurate and robust classification model in the current implementation.

## CONCLUSION

The findings of this study emphasize the importance of developing plagiarism detection tools that are sensitive to regional linguistic and cultural nuances. Generic systems often fail to capture subtle variations in language, idioms, and structure, which can lead to undetected plagiarism or false positives. By incorporating a Backtrack Matching Model tailored to regional languages like Hindi, Marathi, Tamil, and Gujarati, detection accuracy was significantly improved. Evaluation metrics are precision, recall, and F1 score has confirmed this improvement, showing a marked increase in reliability and efficiency. Specifically, the F1 score rose from 0.87 to 0.94 means overall average accuracy improved from 5% - 15%. The implementation of a regionally optimized repository and enhanced algorithmic speed further bolstered the tool's effectiveness. Challenges such as paraphrasing, translation, and multilingual content remain, but targeted algorithms show promise in addressing them. The results underscore the value of localized solutions in upholding academic integrity. This study advocates for broader adoption of regional models to ensure fair and accurate plagiarism detection across diverse linguistic settings. Continued development and refinement will help bridge the gaps left by traditional detection systems and support ethical integrity worldwide.

## REFERENCES:

[1] Xian, J., Yuan, J., Zheng, P., Chen, D., & Yuntao, N. (2024). *BERT-Enhanced Retrieval Tool for Homework Plagiarism Detection System*. arXiv preprint arXiv:2404.01582.

[2] Wahle, J. P., Ruas, T., Foltýnek, T., Meuschke, N., & Gipp, B. (2021). *Identifying Machine-Paraphrased Plagiarism*. arXiv preprint arXiv:2103.11909.

[3] Wahle, J. P., Ruas, T., Kirstein, F., & Gipp, B. (2022). *How Large Language Models are Transforming Machine-Paraphrased Plagiarism*. arXiv preprint arXiv:2210.03568.

[4] Humayoun, M., Hashmi, M. A., & Khan, A. H. (2022). *Measuring Plagiarism in Introductory Programming Course Assignments*. arXiv preprint arXiv:2205.08520.

[5] *International Journal of Advanced Research in Computer and Communication Engineering. (2024). Plagiarism Detection Based on Machine Learning*. IJARCCE, 13(4), 107.

[6] Foltýnek, Tomáš, Norman Meuschke, and Bela Gipp. "*Academic Plagiarism Detection: A Systematic Literature Review*." ACM Computing Surveys, vol. 52, no. 6, 2019, pp. 1–39. doi:10.1145/3345317.

[7] Eisa, Taiseer A. E., Naomie Salim, and Salha Alzahrani. "*Existing Plagiarism Detection Techniques: A Systematic Mapping of the Scholarly Literature*." Journal of Theoretical and Applied Information Technology, vol. 78, no. 2, 2015, pp. 221–229.

[8] Meuschke, Norman. "*Analyzing Non-Textual Content Elements to Detect Academic Plagiarism.*" arXiv preprint arXiv:2106.05764, 2021.

[9] Ihle, Cornelius, et al. "*A First Step Towards Content Protecting Plagiarism Detection.*" arXiv preprint arXiv:2005.11504, 2020.

[10] Wahle, Jan Philip, et al. "*Identifying Machine-Paraphrased Plagiarism.*" arXiv preprint arXiv:2103.11909, 2021.

[11] Kamat, Omraj, et al. "*Plagiarism Detection Using Machine Learning.*" arXiv preprint arXiv:2412.06241, 2024.

[12] Kumar, Abhishek, and Suraiya Jabin. "*Plagiarism Detection in Text Document Using Sentence Based Semantic Similarity.*" Procedia Computer Science, vol. 89, 2016, pp. 36–41. doi:10.1016/j.procs.2016.06.006.

[13] Zargari, Sima, et al. "*A New Method for Plagiarism Detection Using N-gram and Semantic Similarity.*" International Journal of Information and Education Technology, vol. 5, no. 2, 2015, pp. 89–92.

[14] Alzahrani, Salha M., Naomie Salim, and Ajith Abraham. "*Understanding Plagiarism Linguistic Patterns, Textual Features, and Detection Methods.*" IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews), vol. 42, no. 2, 2011, pp. 133–149.

[15] Shakeri, Mohadese, and Nasser Ghadiri. "*Paraphrase Plagiarism Detection in Text Documents Using Sentence Similarity Based on Semantic Role Labeling.*" Expert Systems with Applications, vol. 168, 2021, 114210.

[16] Groza, Adrian, and Andrei Marcu. "*Using Ontologies and Argumentation for Detecting Plagiarism in Scientific Publications.*" Studies in Logic, Grammar and Rhetoric, vol. 57, no. 1, 2019, pp. 7–24.

[17] Sharma, Chetan, and Jitender Kumar Chhabra. "*Source Code Plagiarism Detection Techniques: A Review.*" Computer Applications: An International Journal (CAIJ), vol. 2, no. 3, 2015, pp. 55–65.

[18] Gupta, Gaurav, and Ashok Kumar. "*Plagiarism Detection Using Machine Learning Techniques.*" International Journal of Advanced Research in Computer Science, vol. 8, no. 5, 2017, pp. 167–170.

[19] Singh, Uday Pratap, and Sudipta Roy. "*Techniques for Detecting Plagiarism: A Review.*" International Journal of Computer Applications, vol. 125, no. 11, 2015, pp. 8–12.

[20] Mozgovoy, Max. "*A Fast String Matching Algorithm with Application to Plagiarism Detection.*" Journal of Educational Computing Research, vol. 43, no. 2, 2010, pp. 211–228.

[21] Clough, Paul. "*Plagiarism in Natural and Programming Languages: An Overview of Current Tools and Technologies.*" Research Memorandum, Department of Computer Science, University of Sheffield, 2000.

[22] Maurer, Hermann A., Frank Kappe, and Bilal Zaka. "*Plagiarism–A Survey.*" Journal of Universal Computer Science, vol. 12, no. 8, 2006, pp. 1050–1084.

[23] Faidhi, Jalal A., and Susan K. Robinson. "*An Empirical Approach for Detecting Program Similarity within a University Programming Environment.*" Computers & Education, vol. 11, no. 1, 1987, pp. 11–19.

[24] Lancaster, Thomas, and Frank M. Lancaster. "*Computer-Based Assessment for Computing: A Review of Plagiarism Literature.*" Learning and Teaching in Computing and Engineering, IEEE, 2006, pp. 56–60.