

## LEVERAGING SEMANTIC ANALYSIS FOR MONOLINGUAL PLAGIARISM DETECTION

Sagar Kulkarni<sup>1</sup>, Dr. Sharvari Govilkar<sup>2</sup>, Dhiraj Amin<sup>3</sup>

<sup>1,2,3</sup>Department of Computer Engineering, Pillai College of Engineering, New Panvel, India

<sup>1</sup>[sagark@student.mes.ac.in](mailto:sagark@student.mes.ac.in), <sup>2</sup>[sgovilkar@mes.ac.in](mailto:sgovilkar@mes.ac.in), <sup>3</sup>[amindhiraj@mes.ac.in](mailto:amindhiraj@mes.ac.in)

### **ABSTRACT:**

Monolingual plagiarism detection is a challenging and underexplored area in Natural Language Processing (NLP). Existing systems lack accuracy due to their limited consideration of semantics, Named Entity Recognition (NER), and paraphrases. The proposed research aims to overcome the limitations of existing approaches by creating a new system for MonoLingual Plagiarism Detection (MLPD). The system utilizes Semantic Analysis and advanced methods to accurately check text similarity between Hindi-English language pairs, addressing the challenges of Monolingual plagiarism. Monolingual plagiarism detection addresses the challenges by employing semantic analysis, NER identification, and paraphrase detection. Leveraging semantic algorithms like WordNet, LSA, and BERT, the system captures underlying meaning and detects similarities across languages. NER recognition enhances detection by identifying named entities, while paraphrase detection identifies equivalent expressions. The outcomes of this research contribute to the advancement of plagiarism detection systems, promoting integrity in a monolingual world. leveraging semantic analysis algorithms proved to be out performed in achieving accurate and efficient plagiarism percentage.

**Keywords:** *Monolingual plagiarism detection, Semantic Analysis, LSA, WordNet, BERT Algorithm, Text Similarity, Named Entity Recognition.*

### **1. INTRODUCTION**

Plagiarism is a pervasive issue in academic, professional, and creative fields, undermining the principles of originality, honesty, and intellectual integrity. It involves the unauthorized use or appropriation of someone else's ideas, words, or work without proper attribution. Plagiarism can take various forms, including verbatim copying, paraphrasing without citation, self-plagiarism, and even plagiarism of ideas. It poses a significant threat to the credibility and authenticity of scholarly and creative endeavors. Plagiarism can take various forms, and understanding its nuances is crucial for effective detection.

The types of plagiarism include:

**Verbatim Copying:** Directly copying text from a source without using quotation marks or providing proper citation.

**Paraphrasing Plagiarism:** Rewriting someone else's ideas or work in different words without acknowledging the original source.

**Self-Plagiarism:** Presenting one's own previously published work as new or original without proper citation.

**Mosaic Plagiarism:** Piecing together information from multiple sources without proper attribution, creating an illusion of original work.

**Idea Plagiarism:** Presenting someone else's ideas or concepts as one's own without giving credit.

Plagiarism detection plays a crucial role in safeguarding academic and professional integrity. It involves the use of various techniques and tools to identify instances of plagiarism and provide evidence of originality or identify the sources that have been improperly used. Plagiarism detection serves multiple purposes, including ensuring fair competition, preserving intellectual property rights, maintaining quality standards, and promoting ethical practices. The objectives and importance of plagiarism detection are multifaceted. Firstly, plagiarism detection acts as a deterrent, discouraging individuals from engaging in plagiarism by creating awareness about the consequences and potential repercussions. By fostering a culture of originality and integrity, it upholds the values of academic and professional communities. Secondly, plagiarism detection serves as a proactive mechanism to identify and address instances of plagiarism before they tarnish the reputation of individuals or institutions. By promptly identifying and addressing plagiarism cases, it helps maintain the credibility and trustworthiness of scholarly and creative works. Plagiarism detection is essential for maintaining academic and professional integrity. It acts as a deterrent, identifies instances of plagiarism, and promotes a culture of originality and ethical practices.

**Background:** Plagiarism detection is a field of research that focuses on developing techniques and tools to identify instances of plagiarism, ensuring the integrity and originality of scholarly and creative works. Monolingual plagiarism detection (MPD) involves detecting plagiarism within a specific language. It typically relies on syntactic analysis, where textual patterns and similarities are examined to identify instances of plagiarism. Techniques such as n-gram analysis, fingerprinting, and string-matching algorithms are commonly used in detection. These approaches compare the linguistic structure, vocabulary, and phraseology of documents to identify similarities indicative of plagiarism.

**Problem Statement:** Develop an advanced monolingual plagiarism detection system using semantic analysis and a combination of algorithms, including Jaccard, Cosine, Latent Semantic Analysis (LSA), BERT, and WordNet, to accurately compute the plagiarism percentage. Address the sensitivity and critical issues associated with plagiarism detection to enhance the reliability and effectiveness of the system.

Plagiarism detection is a critical area of research that plays a vital role in maintaining academic integrity, preserving intellectual property rights, and upholding the standards of originality in written content. However, existing plagiarism detection systems often rely on limited algorithms or simplistic approaches, leading to false positives, false negatives, and inadequate detection of plagiarism.

## 2. RELATED WORK

The detection of plagiarism, the unauthorized use or appropriation of someone else's work, is a critical task in ensuring academic and intellectual integrity. While substantial research has

been conducted on plagiarism detection, the majority of existing systems and approaches have focused on monolingual scenarios, primarily for English language.

The authors in [1] introduce an NLP-based approach for detecting plagiarism in short sentences. It emphasizes the importance of considering linguistic patterns in plagiarism detection and proposes a methodology that utilizes natural language processing techniques to identify similarities and patterns in short sentences, enabling effective plagiarism detection. In semantic calculation synonyms of arguments are compared between sentences.

In [2], [3] emphasize on the importance of considering semantic aspects in plagiarism detection and provide methodologies and tools to support the analysis of textual similarities based on semantic content. These approaches offer a promising avenue for improving the accuracy and effectiveness of plagiarism detection systems. The specific algorithms mentioned in the paper include WordNet-based measures, Latent Semantic Analysis (LSA), and Latent Dirichlet Allocation (LDA). These algorithms enable the calculation of semantic relatedness between sentences or documents by considering word semantics, co-occurrence statistics, and probabilistic modeling techniques.

The authors in [4-7], [11] have proposed semantic approaches for various natural language processing tasks, including plagiarism detection, is emphasized. Several algorithms are mentioned in the papers: Devlin et al., in [4], BERT stands for "Bidirectional Encoder Representations from Transformers." It is a pre-training model based on deep bidirectional transformers for language understanding. BERT has been widely used for various NLP tasks, including semantic analysis and text classification. Xie et al., in [5]: XLM is an algorithm for unsupervised cross-lingual representation learning. It focuses on learning meaningful representations of multilingual text data, enabling effective transfer learning across different languages. Vaswani et al.,[6] The Transformer model is a powerful neural network architecture that utilizes self-attention mechanisms to capture long-range dependencies in sequences. It has been widely applied in NLP tasks, including machine translation and language understanding. Semantic Similarity Analysis [8] approach utilizes semantic similarity analysis techniques to compare the textual content of documents and identify potential instances of plagiarism. The specific algorithms employed for semantic analysis are not mentioned in the given study. Artificial Neural Networks with Semantic Analysis [9] approach combines the power of artificial neural networks with semantic analysis techniques to detect plagiarism. The study does not specify the exact neural network architecture used for the task.

Word2Vec-based Semantic Analysis [10], [12] approach utilizes Word2Vec, a popular word embedding model, for semantic analysis. Word2Vec represents words in a high-dimensional space, capturing semantic relationships between words. The model leverages these semantic embeddings to detect instances of plagiarism. The research in [13-17] highlights the importance of combining different similarity measures, such as Jaccard Similarity, Cosine Similarity, and Latent Semantic Analysis, to effectively detect instances of plagiarism. By leveraging semantic relationships and similarity measures, these methods provide a comprehensive approach for plagiarism detection. The study highlights the importance of considering both lexical and semantic aspects of the text to detect plagiarism effectively. LSA, a technique that captures latent semantic relationships between words, is commonly employed to uncover deeper meaning and semantic context in the documents being analyzed. Additionally, measures like

Jaccard Similarity and Cosine Similarity are utilized to assess the overlap and similarity between sets of words or vector representations.

The combination of these techniques allows for a more comprehensive plagiarism detection approach, particularly when applied to various types of documents, including academic papers and programming assignments. By considering both lexical and semantic similarities, these approaches provide a robust framework for identifying instances of plagiarism and promoting academic integrity.

The authors in [18],[19] emphasize the importance of incorporating deep learning and semantic analysis techniques to enhance the effectiveness of plagiarism detection and paraphrase identification in Arabic texts. These approaches provide valuable insights and contribute to the advancement of language-specific plagiarism detection methods, catering to the unique characteristics and challenges of the Arabic language. The authors propose a semantic analysis approach to identify paraphrased sentences or passages in Arabic language documents. By examining the semantic similarities between sentences, the method aims to uncover instances of paraphrasing, which can be helpful in detecting potential cases of plagiarism or content manipulation.

The authors in [20-22] highlight the importance of semantic analysis techniques in detecting plagiarism in Indian languages such as Marathi and Hindi. By considering the specific linguistic features and challenges of these languages, these approaches contribute to the development of effective plagiarism detection methods for Indian language texts. The study highlights the importance of considering language-specific characteristics and challenges in developing plagiarism detection systems for Marathi and Hindi languages. The emphasis is on leveraging semantic analysis to identify similarities, differences, and obfuscated instances of plagiarism in Indian language texts.

### 3. METHOD

The proposed system intends to detect plagiarism by considering semantics of the text document which is cross lingual in nature. It comprises the following modules namely System training module for both text and image contents, Preprocessing, NLP module with different operations such as POS tagging, NER identification etc. At the last semantic similarity identification will be done using few statistical algorithms (such as Jaccard Similarity, Cosine similarity etc.) and semantic analysis algorithms (such as LSA, Wordnet/BoW etc.). The detailed flow of the system working is as shown in figure 1.

#### **Training Module**

In this step, the document files are preprocessed to prepare them for training. The preprocessing steps typically involve: Tokenization, Filtration, Stopword Analysis, Named Entity Recognition (NER) detection, References Detection etc. The preprocessed data from these steps are then stored in the training dataset.

#### **Testing Module**

In the testing phase of your proposed research architecture, the goal is to check the plagiarism of input files, which can be either text documents or image files. Here's an overview of the process:

**Input file format:** The input file can be in various formats such as .txt, .doc, .pdf for text documents. This flexibility allows your system to handle a wide range of input file types.

**Preprocessing:** The input document will go through the preprocessing steps which include tokenization, filtration, stop word analysis etc.

**NLP operations:** Once the preprocessing is done then the document is analyzed for lemmatization, NER detection, and reference resolution. The goal is to normalize and prepare the input data for comparison.

**Comparison with trained documents:** The preprocessed input document is compared with trained documents to determine the plagiarism count. Algorithms such as Jaccard similarity, Cosine similarity, Latent Semantic Analysis (LSA), WordNet, and BERT will be utilized for the comparison process.

**Plagiarism count:** The plagiarism count represents the degree of similarity between the input document and the trained documents. It can be computed based on the comparison results. The threshold for determining what constitutes plagiarism can be defined based on your research requirements or predetermined criteria.

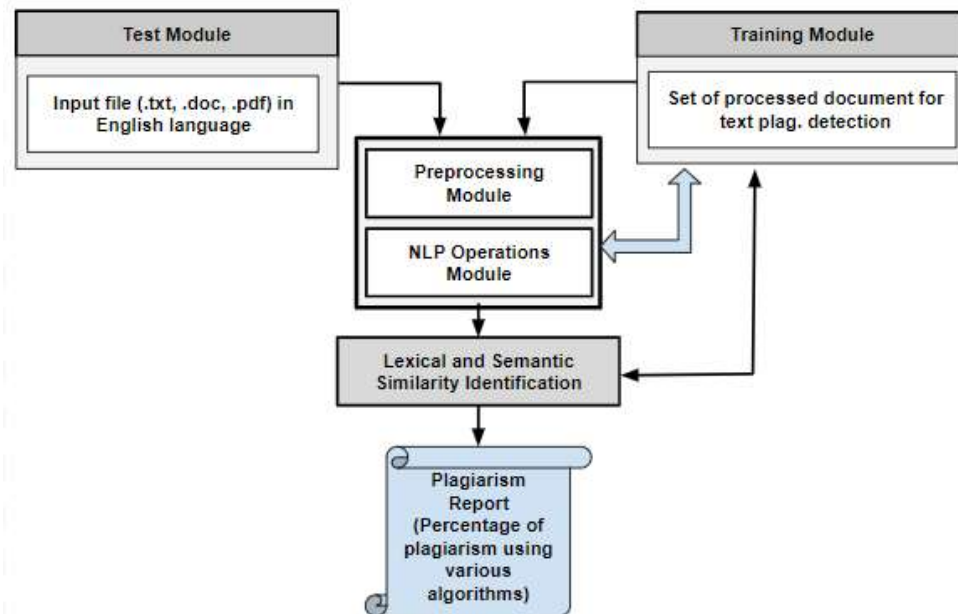


Figure 1. Monolingual Plagiarism Detection- Architecture

#### 4. RESULTS AND DISCUSSION

The obtained results in the plagiarism detection system provide users with a detailed overview of the analysis performed on the document. The system performs plagiarism detection along with handling named entities and reference resolution.

**Monolingual Plagiarism Detection:** The system successfully detects instances of plagiarism within documents written in the same language. By applying algorithms such as Jaccard similarity, cosine similarity, LSA, BERT and WordNet/BoW, the system accurately identifies similarities and determines the percentage of plagiarism within the document.

### Results obtained for Monolingual Plagiarism Detection

The developed system has been already trained with corpus documents which is a collection of .txt , .doc, .pdf files constituting published research work of a few authors. While in the testing phase when input document is given to the system then the plagiarism count will be computed using various algorithms. In the preprocessing phase itself the system performs various NLP operations so as to compute plagiarism detection more efficiently.

The table 1 presents the results of a sample monolingual plagiarism detection using different algorithms for sample input. For each algorithm the system has two outputs: 1. Plagiarism % with named entities included and 2. Plagiarism % with named entities excluded. This is required to have comparative analysis system behavior with and without named entities. The chart in figure 2 associated with the table likely visualizes the data to provide a graphical representation of the plagiarism percentages for each algorithm.

Table 1. Monolingual plagiarism detection results using various methods: Sample input-1

Sr. No.	Algorithm	Plagiarism % when Named Entities Included	Plagiarism % when Named Entities Excluded
1	Jaccard	7.83	7.14
2	Cosine	26.71	26.35
3	LSA	88.48	70.93
4	BERT	50.00	49.30
5	WordNet/BoW	62.93	59.92

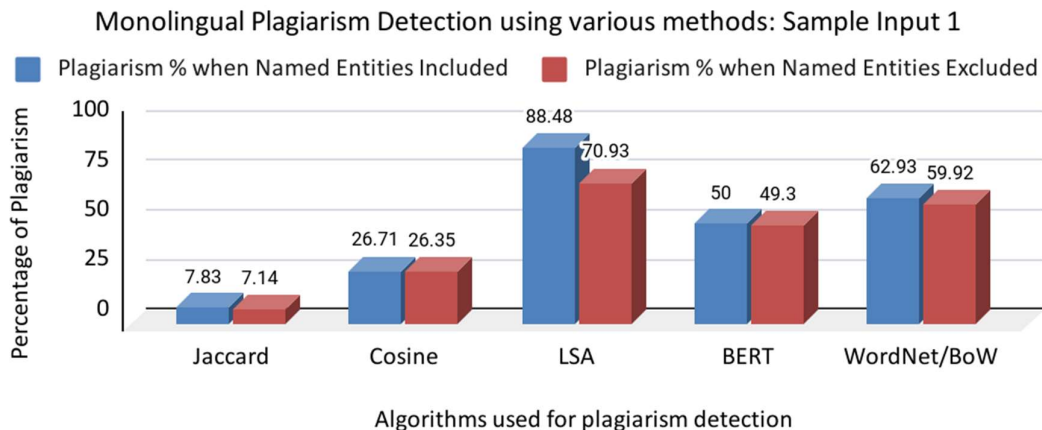


Figure2. Graphical representation of plagiarism % for Sample input-1

The system was tested with various document formats and algorithms, showing satisfactory results with semantic analysis approaches. Table 2 and figure 3 display sample results with and without inclusion of named entities during plagiarism detection.

Table 2. Plagiarism Percentage with Inclusion of Named Entities

Input	Jaccard	Cosine	LSA	BERT	WordNet
Input1	7.83	26.71	88.48	50	62.93
Input2	8.61	8.05	42.08	26.68	23.58
Input3	12.7	10.19	61.97	35.89	26.38
Input4	11.13	9.51	73.3	38.54	26.19
Input5	11.14	9.36	67.23	33.44	28.45
Input6	13.92	10.83	51.37	36.65	28.9
Input7	15.78	11.14	38.58	34.11	29.24

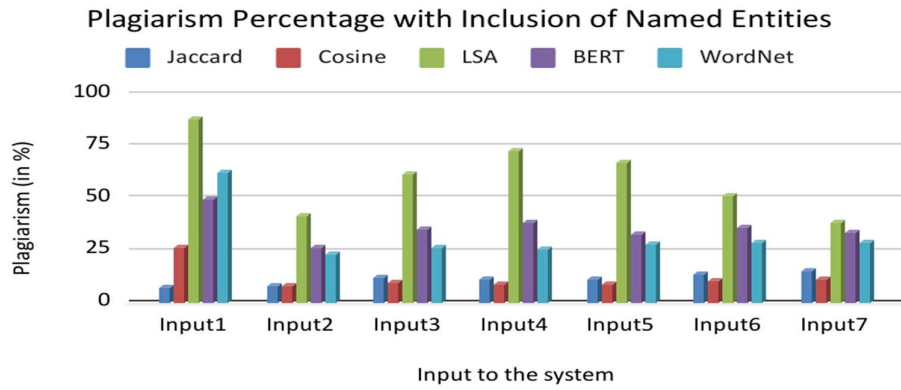


Figure 3. Visualizing Average Plagiarism Percentage with Inclusion of Named Entities

Similarly, table 3 gives sample results when named entities are excluded while in plagiarism detection and the chart representation of the same has been shown in figure 4.

Table 3. Plagiarism Percentage with Exclusion of Named Entities

Input	Jaccard	Cosine	LSA	BERT	WordNet
Input1	7.14	26.35	70.93	49.3	59.92
Input2	7.4	8.51	4.24	30.06	19.01
Input3	11.8	10.66	17.09	36.32	26.74
Input4	9.81	9.53	11.46	37.88	22.27
Input5	9.17	9.6	13.97	34.92	24.88
Input6	11.95	10.74	8.56	36.39	22.82
Input7	14.74	12.1	7.56	37.29	23.4

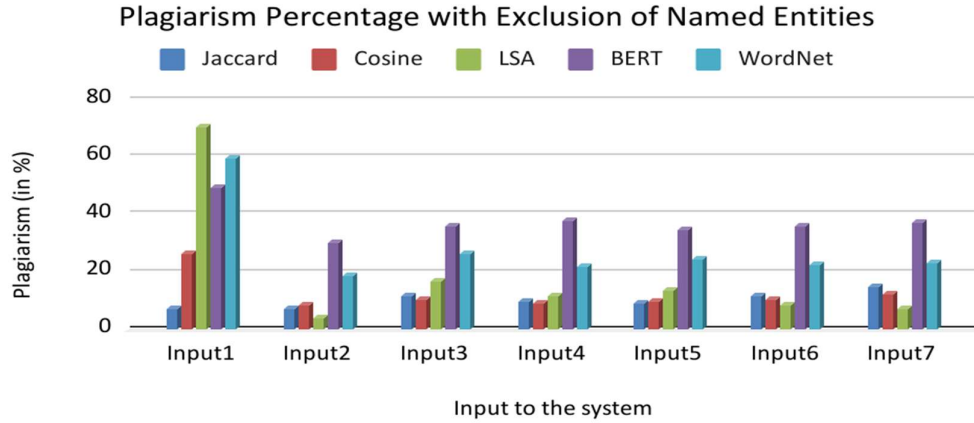


Figure 4. Visualizing Average Plagiarism Percentage with Exclusion of Named Entities

Interestingly, the findings more often suggest that excluding named entities yields favorable outcomes for plagiarism detection. Though there are very few instances that the percentage of plagiarism is a bit higher for some algorithms, still it is recommended to leverage plagiarism detection which excludes named entities because introducing named entities will add false positives and inaccuracies and also results in bias towards surface-level matches. The semantic algorithms such as LSA, BERT, WordNet are shown to be more reliable in plagiarism detection. It has been noticed that results of plagiarism detection using LSA algorithm with and without NER is differentiating more than other algorithms proving the importance of excluding named entities from plagiarism detection.

Including named entities in the analysis enhances the sensitivity to surface-level matches that may not indicate plagiarism. The system generates a comprehensive plagiarism report. The sample snippet of the report is as shown in table 4, comparing the input file with corpus files using different algorithms such as LSA, BERT, and WordNet. The report provides specific percentages for each algorithm and corpus file pair, enabling users to analyze similarities and potential instances of plagiarism. The report focuses on semantic algorithms that have shown satisfactory results when named entities are excluded. The input document is compared against 40 documents from corpus.

Table 4. Plagiarism Percentage of Selected Algorithms with Individual Corpus Files: Comparative Analysis

Sr. No.	Filename	LSA	BERT	WordNet
1	A Comparative Study	3.61%	18.20%	23.14%
2	A Method for Plagiarism Detection	3.58%	18.20%	22.95%
3	A Survey On Plagiarism Detection	2.33%	17.28%	19.87%
4	Academic Plagiarism Detection_A LR	5.30%	18.35%	23.06%



5	An effective approach to candidate retrieval	5.87%	16.90%	23.18%
6	An NLP Based Plagiarism Detection	2.35%	15.77%	21.27%
7	Automatic Plagiarism Detection	1.82%	18.88%	20.25%
8	Comparison of Plagiarism Detection Tools	4.46%	16.79%	23.59%
9	Compiling a text reuse detection	3.07%	16.78%	24.78%
10	Cross language text Alignment	5.61%	17.73%	22.17%
.	.	.	.	.
.	.	.	.	.
25	testfile2	8.75%	19.08%	32.48%
26	testfile3	7.76%	17.08%	31.75%
27	testfile4	5.34%	17.62%	30.50%
28	testfile5	98.28%	100.00%	100.00%
29	testfile6	1.16%	17.25%	18.52%
30	testfile7	11.72%	16.73%	39.68%
.	.	.	.	.
.	.	.	.	.

The combination of semantic analysis and excluding named entities improves monolingual plagiarism detection. This approach consistently outperforms other methods, providing more accurate and reliable results. The overall quality of detection is emphasized over numerical output, showcasing the significance of incorporating semantic understanding and excluding named entities in achieving effective plagiarism detection.

## 5. CONCLUSION

The development of a monolingual plagiarism detection system using semantic analysis and a combination of algorithms, including Jaccard, Cosine, Latent Semantic Analysis (LSA), BERT, and WordNet, holds great promise for accurately computing the plagiarism percentage and addressing the sensitive issue of plagiarism detection. By leveraging semantic analysis techniques, the system goes beyond surface-level matching and focuses on capturing the meaning, context, and semantics of the text to identify instances of potential plagiarism. The overall findings of the research indicate that semantic methods especially LSA and BERT outperform lexical methods in terms of efficiency and accuracy for plagiarism detection. By leveraging semantic analysis and resolving named entities, the system achieved more efficient and accurate detection of plagiarized content across languages. Moreover, the research facilitated the generation of detailed reports or summary reports of plagiarism. Plagiarism remains a substantial challenge across domains, requiring more effective approaches for combating it. The focus is on developing methods that minimize false positives while ensuring accuracy and efficiency. Continuous research and development are crucial to effectively

address the widespread issue of plagiarism. By enhancing plagiarism detection systems, we can promote academic integrity and originality in various fields of study and professional practice.

## REFERENCES

- [1] Pandey, S., & Rawal, A. (2018). An NLP Based Plagiarism Detection Approach for Short Sentences. *International Journal of Recent Technology and Engineering (IJRTE)*, 7(4), November 2018.
- [2] Al-Shamery, E. S., & Ghenni, H. Q. (2016). Plagiarism Detection using Semantic Analysis. *Indian Journal of Science and Technology*, 9(1). DOI: 10.17485/ijst/2016/v9i1/84235
- [3] Rus, V., Lintean, M., Banjade, R., Niraula, N., & Stefanescu, D. (2013). SEMILAR: The Semantic Similarity Toolkit. *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, 163-168.
- [4] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Volume 1 (Long and Short Papers), 4171-4186.
- [5] Xie, R., Lu, Y., Zeng, J., Xie, J., & Huang, D. (2020). Unsupervised cross-lingual representation learning at scale. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4369-4379.
- [6] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 5998-6008.
- [7] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 1873-1883.
- [8] Lashari, A. A., et al. (2020). Plagiarism Detection using Semantic Similarity Analysis. *Proceedings of the 2020 3rd International Conference on Communication, Computing and Digital Systems (C-CODE)*.
- [9] Ahmed, S., et al. (2020). Plagiarism Detection Using Artificial Neural Networks with Semantic Analysis. *Proceedings of the 2020 International Conference on Intelligent Systems Design and Applications (ISDA)*.
- [10] Garg, M., et al. (2021). A Plagiarism Detection Model based on Semantic Analysis using Word2Vec. *Proceedings of the 2021 International Conference on Computational Intelligence and Data Science (ICCIDS)*.
- [11] Liu, X., et al. (2021). A Plagiarism Detection Method based on Text Representation using BERT. *Proceedings of the 2021 International Conference on Machine Learning and Intelligent Systems (ICMLIS)*.
- [12] Gupta, S., & Singhal, S. (2022). Plagiarism Detection Using Cosine Similarity and WordNet. In *Proceedings of the 2nd International Conference on Computational Intelligence in Pattern Recognition (CIPR 2022)* (pp. 112-116). Springer.
- [13] Ali, S. S., Hossain, M. F., & Uddin, M. S. (2021). Plagiarism Detection Using Jaccard Similarity and Latent Semantic Analysis. In *Proceedings of the 11th International Conference on Intelligent Systems and Control (ISCO 2021)* (pp. 74-78). IEEE.

- [14] Gupta, N., & Jain, M. (2022). Plagiarism Detection Using WordNet and Latent Semantic Analysis. *In Proceedings of the 4th International Conference on Computational Intelligence in Data Science (ICCIDS 2022)* (pp. 165-170). IEEE.
- [15] Kumar, S., Kumar, S., & Agarwal, S. (2022). Plagiarism Detection Using a Hybrid Approach of Jaccard Similarity, Cosine Similarity, and Latent Semantic Analysis. *Journal of Information Science*, 48(1), 56-70. DOI: 10.1177/01655515211042668
- [16] Jain, V. P., & Kumar, P. (2021). Plagiarism Detection in Programming Assignments using LSA and Cosine Similarity. *International Journal of Engineering and Advanced Technology*, 10(2), 385-389. DOI: 10.35940/ijeat.B4492.129221
- [17] Islam, M. T., & Hoque, M. A. (2021). Plagiarism Detection Using Hybrid Algorithm of WordNet and Cosine Similarity. *Journal of Computer Science*, 17(1), 30-41. DOI: 10.3844/jcssp.2021.30.41
- [18] Dima Suleiman, Arafat Awajan, Nailah Al-Madi, "Deep Learning Based Technique for Plagiarism Detection in Arabic Texts," *The International Conference on new Trends in Computing Sciences (ICTCS2017)*, DOI: 10.1109/ICTCS.2017.42, 2017
- [19] Adnen Mahmoud, Mounir Zrigui, "Semantic Similarity Analysis for Paraphrase Identification in Arabic Texts," *Proceedings of the 31st Pacific Asia Conference on Language, Information and Computation, The National University (Phillippines)*, pg 274-281, PDF
- [20] Naik Ramesh & Landge, Maheshkumar & C., Namrata. (2017), "Plagiarism Detection in Marathi Language Using Semantic Analysis," *International Journal of Strategic Information Technology and Applications*, 8, 30-39. DOI: 10.4018/IJSITA.2017100103.
- [21] Nilam Shenoy, M. A. Potey, "Semantic Similarity Search Model for Obfuscated Plagiarism Detection in Marathi Language using Fuzzy and Naïve Bayes Approaches," *IOSR Journal of Computer Engineering (IOSR-JCE)*, e-ISSN: 2278-0661, p-ISSN: 2278-8727, Volume 18, Issue 3, Ver. V (May-Jun. 2016), PP 83-88, DOI: 10.9790/0661-1803058388, 2016
- [22] Urvashi Garg and Vishal Goyal, "Maulik: A Plagiarism Detection Tool for Hindi Documents", *Ind. J. Sci. Technology* 9, 12 (2016)