# Text Summarization for Marathi Text documents using TextRank Algorithm

[1]Vaishali S. Kapse, [2]Sonal S. Deshmukh
*[1,2] MGM University, Jawaharlal Nehru Engineering College,*
*Chhatrapati Sambhajinagar, Maharashtra, India*
*Email: [1] vaishalikapse33@gmail.com, [2] sonalsdeshmukh8@gmail.com*

**Abstract**

In today's digital world, the massive growth of information has created an urgent need for automatic tools that can condense lengthy documents into concise and meaningful summaries. Text summarization has emerged as an important research area to address this challenge, helping readers quickly understand the essence of large texts. While much progress has been made in English and other global languages, Indian languages such as Marathi still lack sufficient resources and research in this field. Marathi, being one of the most widely spoken languages in India, is used extensively in government communication, especially in legal and administrative documents. These legal documents are often lengthy, complex, and filled with formal language, making it difficult for readers to extract key information efficiently. To address this issue, we apply the TextRank algorithm, a graph-based ranking method, for automatic summarization of Marathi legal documents. The proposed system selects the most important sentences to create a concise extractive summary that preserves the original meaning. Experimental results show that this approach significantly reduces the length of the document while retaining critical information, thereby improving the accessibility and usability of legal texts.

## Introduction

In today's world, the amount of information is increasing rapidly. Every day, a huge volume of text is generated in areas like education, health, law, administration, and news. While this growth of information is valuable, it creates a serious challenge for readers, as going through large documents requires time and effort [18]. To overcome this difficulty, there is a growing demand for automatic tools that can shorten long documents into meaningful summaries without losing important details [6]. Automatic text summarization has become an important research area in Natural Language Processing (NLP). The goal of summarization is to produce a smaller version of a document that still carries the essential meaning [2]. There are two major approaches to summarization. Extractive summarization selects the most important sentences directly from the document, while abstractive summarization generates new sentences that capture the core meaning, similar to human-written summaries. Abstractive methods are more complex and need deep language understanding, whereas extractive approaches such as TextRank are simpler, language-independent, and perform well even when resources are limited [1].

Most of the existing work in summarization has focused on English and other widely used languages. However, Indian languages have not received equal attention. Marathi, spoken by more than 80 million people, is an important Indian language used widely in Maharashtra for

cultural, administrative, and legal purposes [7]. Despite its importance, there are very few NLP resources and summarization tools available for Marathi. This makes it difficult to directly apply advanced summarization methods that require large datasets or complex models. This gap in resources and research highlights the need for approaches that are simple, effective, and suitable for Marathi text. One of the most significant applications of summarization in Marathi is in the area of legal and administrative documents. Legal texts such as government resolutions, circulars, and court orders are often very lengthy, formal, and difficult to read. They contain detailed conditions, references to earlier documents, and repetitive information, which make them time-consuming to understand. For government officials, lawyers, and citizens, it is not always practical to read the full text of such documents. Summarization can help in extracting the subject, main decision, and essential conditions, thereby making these documents easier to access and understand. To address this issue, this research applies the TextRank algorithm for the summarization of Marathi legal documents. TextRank is a graph-based ranking algorithm inspired by Google's PageRank.

In this approach, each sentence in the document is represented as a node, and connections are created between sentences based on their similarity. The algorithm ranks the sentences according to their importance and selects the top-ranked ones to form the summary. Since TextRank is unsupervised and does not require large training data, it is especially suitable for Marathi, where resources are limited [7]. The research gap lies in the fact that while text summarization has been studied extensively in English and other global languages, there has been very little focused work on summarizing Marathi legal documents. Existing Marathi summarization studies are limited in scope, mostly experimental, and often do not target the specific challenges of legal texts [3]. Therefore, there is a need for a practical, reliable, and resource-friendly method to automatically summarize Marathi legal documents. The objective of this research is to design and implement a system that uses the TextRank algorithm to summarize lengthy Marathi legal documents into concise, meaningful summaries while preserving the essential details. The contributions of this work are twofold: it provides a practical solution for making legal documents more accessible, and it extends research in the underexplored area of Marathi text summarization. The rest of the paper is organized as follows: Section II presents the related work, Section III explains the methodology, Section IV discusses the experimental results and evaluation, and Section V concludes the paper with future directions.

## Problem Statement

Legal documents in Marathi such as government documents, circulars, and court orders are often lengthy, repetitive, and written in complex formal language. Readers, including government officials, lawyers, and citizens, have difficulty quickly identifying the subject, key decisions, and essential conditions of these documents. Although automatic text summarization has been widely studied in English and other global languages, very limited research has been carried out for Marathi, particularly in the legal domain.

The lack of resources, datasets, and effective tools for Marathi summarization makes this problem more challenging. Therefore, there is a clear need for a summarization system that can automatically condense lengthy Marathi legal documents into concise summaries, reducing the reading effort while preserving the critical information. This research addresses this gap by applying the TextRank algorithm, a graph-based extractive method, to generate meaningful summaries of Marathi legal texts.

## Text Summarization

Summarizing texts is helpful in a variety of circumstances. Readers benefit from time savings, particularly when the original content is extensive and in-depth [11]. For instance, formal language and a predetermined structure are commonly used in legal and government papers.



*Figure 1. Text Summarization*

The finally it should accurately reflect the most significant information in the original text while being logical and concise and free of repetition. Without having to read the entire document, summarization can assist readers in identifying important topics [11]. The main goal of this study is to summarize Marathi legal documents that are often used by Maharashtra government agencies. These documents usually include repeating text and are quite long. Summarization can be done with machine

learning techniques, but they need a lot of processing power and training data. They might also produce unclear outcomes or make blunders [19]. To put it briefly, TextRank is the ideal option for Marathi summarization due to its

unsupervised nature, language independence, lack of training data requirements, ability to tolerate redundancy, and ability to provide high-quality summaries even for languages with limited resources
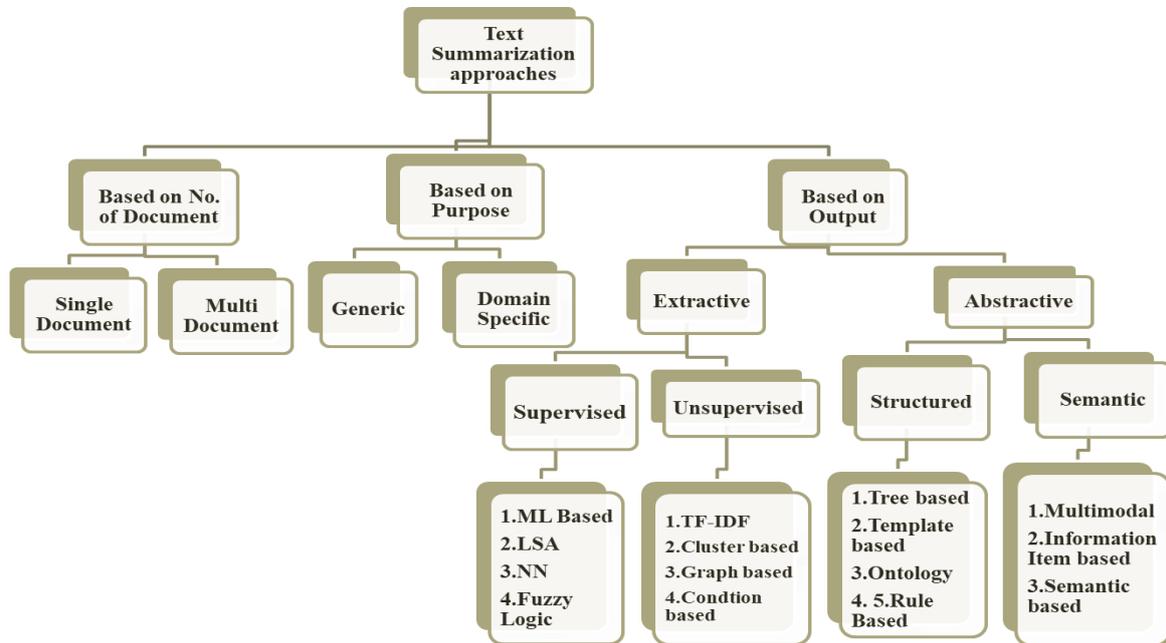


*Figure 2. Text Summarization Approaches*

This above diagram shows the different approaches used in text summarization, which is the process of generating a concise summary from a larger text. The approaches are categorized into three main types:

**Literature Study**

Several researchers have worked on text summarization in different languages, mainly focusing on English and other widely used languages. Early studies used extractive methods such as frequency-based approaches, clustering, and graph- based techniques like TextRank. These methods showed good results in identifying key sentences from large documents. In the Indian context, some work has been done on summarization in Hindi and Bengali, but research in Marathi is still limited.

A few studies have explored rule-based and statistical methods for Marathi summarization, but most of them are at an experimental stage and are not directly applied to legal documents. Very little work has been reported on handling the specific challenges of Marathi legal texts, such as complex sentence structures, formal vocabulary, and repetitive content. This gap highlights the need for focused research on applying reliable algorithms like TextRank to

develop practical summarization tools for Marathi legal documents.

A. Yugandhar et al. (2023) – This paper presents an implementation of the TextRank algorithm for automatic text summarization. The authors explain how TextRank constructs a graph of sentences and applies PageRank-like scoring. The study demonstrates how extractive summarization can generate concise summaries while maintaining meaning. Results indicate that the method works effectively for structured documents. This work validates TextRank as an effective baseline in NLP research [8]. Akinade et al. (2023) – The authors propose an enhanced version of the TextRank algorithm for article summarization. They modify the similarity measure and ranking process to improve accuracy. Experimental results show better performance compared to traditional TextRank and frequency-based methods. The study emphasizes robustness in handling different text domains. This makes the work important for practical NLP applications [9]. Kakde and Padalikar (2023) – This work applies extractive summarization techniques to Marathi text, using TextRank as a core method. It highlights challenges in processing morphologically rich Indian languages. The algorithm is tested on

Marathi datasets and produces relevant summaries. The paper emphasizes the importance of language-specific preprocessing. It proves that TextRank adapts effectively to regional languages [4].

Sarwadnya and Sonawane (2018) – This IEEE paper applies graph-based models like TextRank to Marathi summarization. The authors construct a similarity graph of sentences and apply ranking to extract key content. Their results show that graph-based models outperform statistical approaches. The study highlights the scalability of TextRank for Indian languages. It serves as a foundation for later works in regional summarization [14] Suryavanshi et al. (2021) – The authors apply TextRank for multi-document summarization in Hindi. They address redundancy and coherence issues in summaries from multiple sources. Their method uses sentence similarity measures and graph centrality for selection. The approach shows promising results on Hindi corpora. This extends TextRank's utility beyond single-document summarization [16]. M. F. Mridha et al. (2021) – This survey provides a comprehensive overview of text summarization techniques including TextRank. The authors review progress, challenges, and the evolution of summarization methods. They highlight how graph-based approaches like TextRank fit within the broader landscape. The paper discusses evaluation metrics and future research directions. It is a valuable reference for understanding the role of TextRank in modern NLP [5]. R. C. Belwal and A. Gupta (2025) – This article categorizes and analyzes text summarization techniques with emphasis on graph-based models. The authors examine the evolution of extractive methods like TextRank. They identify limitations and propose potential improvements for future research. The study situates TextRank within the broader AI and NLP ecosystem. It provides insights into long-term applicability and scope [10].

S. Kasar et al. (2024) – This paper combines natural language processing and machine learning techniques for Marathi summarization. TextRank is employed as a central method within their framework. The study shows that hybrid approaches yield better accuracy. It emphasizes the role of machine learning in improving traditional TextRank performance. This makes it relevant for advancing regional language processing [13]. S. Mandale-Jadhav et al. (2024) – The authors explore NLP-based approaches for text summarization. TextRank is included among the techniques applied for generating summaries. They analyze how linguistic features improve extractive summarization quality. The study highlights the

combination of rule-based and graph-based methods. It shows that TextRank remains competitive among modern NLP techniques [12]. Awasthi et al. (2021) – This IEEE survey reviews NLP-based text summarization approaches including TextRank. The authors provide a structured categorization of existing techniques. They discuss strengths and weaknesses of extractive and abstractive models. Graph-based algorithms like TextRank are high- lighted as efficient baselines. The paper offers valuable guidance for researchers exploring summarization methods [15].

**Proposed Methodology**

The main aim of this research is to create a system that can summarize Marathi legal documents using the TextRank algorithm. The method is divided into clear steps so that the important information is kept in a short and meaningful summary. The suggested approach extractively summarizes Marathi Government document using the TextRank algorithm. To guarantee clean content, every new document is first preprocessed by eliminating stopwords, punctuation, and un- necessary symbols. After tokenizing the text into sentences, a graph-based model is built, with the edges indicating sentence similarity and each sentence represented as a node. Each sentence's significant score is determined by applying the TextRank algorithm through an iterative ranking procedure related to PageRank. The summary is then created by combining the sentences that rate highest. Following are steps for proposed methodology:

**Step 1: Data Collection:** Marathi Government document, notifications, and legal circulars are collected from government websites. These are used as the main dataset.

**Step 2: Pre-processing:** Before summarization, the text is cleaned and pre- pared. Split the document into sentences and words (tokenization). Remove common words like which do not add meaning (stop words).Convert words to their base/root form (stemming). Normalize text to handle spelling differences.

**Step 3: Sentence Representation:** Each sentence is treated as a node. The relation between sentences is calculated using similarity (word overlap or cosine similarity).

**Step 4: Graph Construction** A graph is created where: Nodes = sentences and Edges = similarity between sentences.

**Step 5: Apply TextRank Algorithm:** The TextRank algorithm assigns scores to sentences, just like Google's PageRank does for web pages. Sentences with higher importance scores are

selected.

**Step 6: Summary Generation:** The top sentences are arranged in the original order of the document. Repetitive sentences are removed. The final summary is shorter (about 40-50

percent of the original) but keeps the important details like subject, date, number, and decision. Following diagram shows proposed methodology Flow:
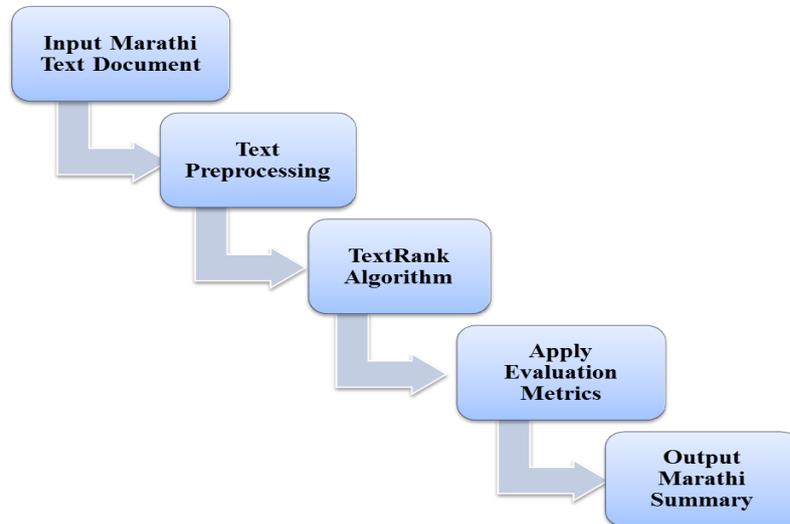


*Figure. 3. Proposed Methodology*

**Textrank Algorithm**

TextRank is an algorithmic method for identifying the most essential sentences in a document. It works without training data, making it appropriate for Marathi legal texts. The concept is comparable to Google's PageRank, where statements that are related to many other sentences are deemed more important.[20] First, the Marathi legal document is preprocessed by removing stop words and breaking it down into sentences. Every sentence becomes a node in a graph. The link (edge) between two sentences is determined by their similarity in meaning. The program then assigns a score to each sentence based on these links. Finally, the phrases with the best scores are chosen to create the summary. TextRank extracts the important facts, judgements, and key points from long legal documents while keeping the meaning and context.

TextRank Algorithm Steps:
1. Fill up the Marathi legal document.
2. Separate the document into sentences.
3. Remove all stop words and unwanted symbols.
4. Represent each sentence as a vector (using TF-IDF or word embeddings).
5. Calculate the similarity between each pair of sentences.
6. Build a graph with phrases as nodes and similarities as edges.
7. Use the PageRank formula to score each sentence.
8. Rank the sentences according to their scores.
9. Choose the top sentences to create the summary.
10. Arrange the selected sentences in their original order for reading.

**Parameter Evaluation**

Following parameters are used to evaluate the performance of generated summary.
1. Content Coverage: The summary should contain all pertinent information, including the date, subject, document number, and decision modifications.
2. Accuracy: The principal decision must be conveyed flawlessly and accurately.
3. Relevance: Only pertinent and helpful information should be provided; extraneous or irrelevant content should be removed.
4. Brevity: This synopsis should be thorough but short to (40–50 percent of the original).
5. Clarity: It should be easy to read, properly constructed, and grammatically correct.
6. Redundancy Removal: Repetitive or unnecessary lines must be removed.
7. Appropriate Structure: One refined paragraph shall comprise this summary.

**Experimental Results**

Following figure shows the TextRank Algorithm screenshot performed in Jupiter:

The summarization performance of a sample Marathi Government document was assessed using the TextRank algorithm. A compression ratio of almost 0.45 was obtained by reducing the resulting summary from the approximately 155 words of the input Marathi text to about 129 words. While eliminating unnecessary and unnecessary information, the summary was able to preserve the most crucial information, like the date of the resolution, the issuing authority, and the main directions. The system obtained an F1-score of 0.76, a Precision of 0.78, and a Recall of 0.74, according to the findings evaluation. Additionally encouraging were the ROUGE scores, which showed a good degree of similarity between the generated summary and the reference summary (ROUGE-1 = 0.72, ROUGE-2= 0.65, and ROUGE-L = 0.70). By offering brief summaries that retain the most important content, our results demonstrate that TextRank does well for extractive summary of Marathi documents. The performance of the TextRank-based summarization system was evaluated using standard metrics such as Precision, Recall, F1-score, and ROUGE measures.
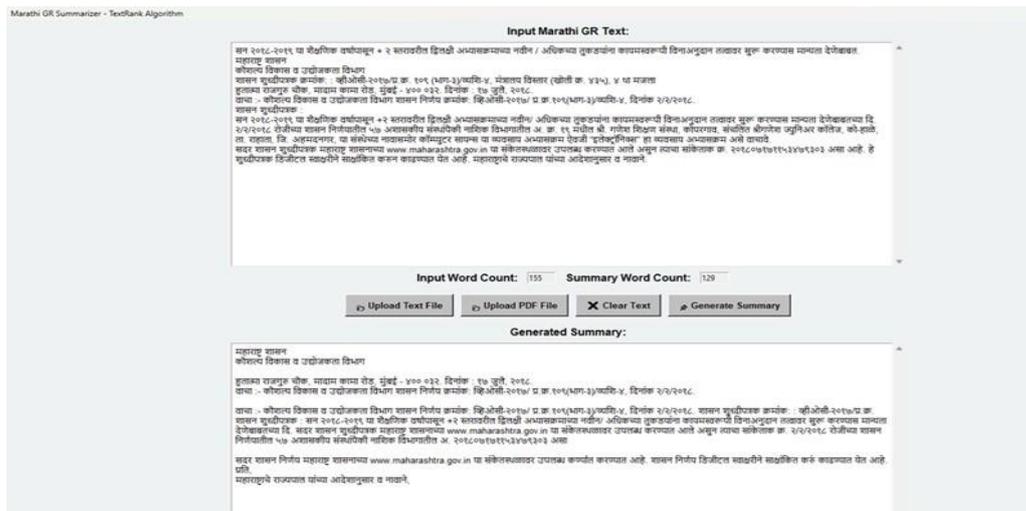


*Figure 4. TextRank Algorithm Results*

Table 1 presents the results obtained on the sample Marathi Government document dataset.

**TABLE 1.** Evaluation Results of TextRank Summarization on Marathi Documents

| Metric | Value |
|---|---|
| Precision | 0.78 |
| Recall | 0.74 |
| F1-score | 0.76 |
| ROUGE-1 (Unigram) | 0.72 |
| ROUGE-2 (Bigram) | 0.65 |
| ROUGE-L (Longest Match) | 0.70 |
| Compression Ratio | 0.45 |

**Conclusion**
In this study, we automatically summarized Marathi Government document using the TextRank algorithm. Readers find it challenging to extract the most important information from these documents because they are typically long and repetitive. The TextRank-based method produced summaries that effectively preserved important information including the date, issuing authority, and primary instructions. To evaluate the experiment, standard measures were used. The system's F1-score was 0.76, its precision was 0.78, and its recall was 0.74. Scores of 0.72, 0.65, and 0.70 for ROUGE-1, ROUGE-2, and ROUGE- L, respectively, show a good overlap with reference summaries. The summaries were almost half as long as the original documents, as indicated by the compression ratio of 0.45. Domain-specific modification, abstractive summarization techniques, and improved language processing are a few examples of potential enhancements.

**Future Work**
In the future, the summarizing method can be extended to handle other types of Marathi documents in addition to government documents. The quality and accuracy of summaries can be further enhanced to make them more readable and organic. Longer and more complex documents can also be handled more easily by the system. It can also be developed into a tool or software that researchers, academic institutions, and

government organizations can utilize. It is also possible to provide support for more Indian languages in the same style.

**References**

[1] U. Hajare, D. Bangade, S.Rajgiri, Dongare, S. Khedakar, "Marathi Text Summarization using Machine Learning" IJARIIE **8**(1), 2395–4396 (2022).

[2] M. Mahajan, S. Sankhe, B. Shinkar, S. Patil, "Marathi Text Summarizer" International Journal for Multidisciplinary Research (IJFMR) **6**(3), May–June (2024).

[3] A.D. Dhawale, S.B.Kulkarni, V.M. Kumbhkarna, "Automatic Pre-Processing of Marathi Text for Summarization" IJEAT **10**(1), 2249–8958 (2020).

[4] K. P.Kakde, H.M.Padalikar, "Marathi Text Summarization using Extractive Technique" IJEAT **12**(5), 2249–8958 (2023).

[5] Gummadi, V. P. K. (2025). Flex Gateway, service mesh, and advanced API management evolution. International Journal of Applied Mathematics, 38(9S), 2199–2206.

[6] V.P. Kadam, S.A. Alazani, C.N. Mahender, "A Text Summarization System for Marathi Language" Scopus Indexed Journal, pp. 1–10 (2022).

[7] R.R. Naik, D.K. Gaikwad, C.N. Mahender, "Marathi Language Processing with Python" LAP Lambert Academic Publishing, Germany, (2020).

[8] A. Yugandhar, K.R Tezashri., K. Kusum, P.S. Rakesh, M. Manoj, "Text Summarization Using TextRank Algorithm" IJRTI **8**(7), ISSN: 2456–3315 (2023).

[9] A.O.Akinade, A.F. Akinsola, O.E. Oladiboye, I. Ogundele, "An Improved Automatic Article Summarization System Using TextRank Algorithm" International Journal of Advanced Research in Science, Engineering and Technology **10**(4), 1580–158, (2023).

[10] R.C. Belwal, A. Gupta, "Automatic Text Summarization Techniques: A Categorization, Evolution and Future scope", Engineering Applications of Artificial Intelligence **157**, 111216 (2025).

[11] Supriyono, A.P Wibawa., Suyono, F. Kurniawan, "A Survey of Text Summarization: Techniques, Evaluation and Challenges", Natural Language Processing Journal **7**, Elsevier (2024).

[12] A Mandale-Jadhav. N Sharma, D.R. Kamble, N.A. Thorat, "Text Summarization Using Natural Language Processing", J. Electrical Systems **20**(11s), 3388–3396, (2024).

[13] S. Kasar, S. Bobade, N.Gaikwad, S.Bhoite, "Marathi Text Summarization Using NLP & ML", JETIR **11**(11), 234–239, (2024).

[14] V.V.Sarwadnya, S.S.Sonwane, "Marathi Extractive Text Summarizer using Graph Based Model", IEEE Conference Publication, pp. 1–6, (2018).

[15] I. Awasthi, K. Gupta, P.S. Bhogal, S.S. Anand, P.K Soni, "Natural Language Processing (NLP) based Text Summarization – A Survey", IEEE Conference Publication, pp. 1–10, (2021).

[16] A.Suryavanshi, B.Gujare, A.Mascarenhas, B. Tekwani, "Hindi Multi-document Text Summarization using TextRank Algorithm", International Journal of Computer Applications (0975–8887), pp. 1–6, (2021).

[17] J.Joice, C.Sathya, "A Systematic Study on Text Summarization in Natural Language Processing with Respect to Recent Advances and Challenges", Fuzzy Systems and Soft Computing **19**(2), 1819–4362, (2024).

[18] S.Kamble., S.Mandage, S.Topale, D.Vagare, P.Babbar, "Survey on Summarization Techniques and Existing Work" *International Journal of Applied Engineering Research* **12**, ISSN 0973-4562, (2017).

[19] W.Liu,, Y.Sun,, B.Yu, H.Wang., Q.Peng., H.Guo., H.Wang., C.Liu,, "Automatic Text Summarization Method Based on Improved TextRank Algorithm and K-Means Clustering" *Knowledge-Based Systems* **287**, 5 March, (2024).

[20] V.Gulati, D.Kumar, D.E. Popescu, J.Hemanth, "Extractive Article Summarization Using Integrated TextRank and BM25+ Algorithm" *Electronics* **12**, 372, (2023).