# Research Article Summarization System: An Integrated Approach Using Machine Learning

| N. Indira Priyadarsini | Bhargavi Nethi | Bala Chandar | Abhay Ganapur |
|---|---|---|---|
| Assistant Professor Department of Information Technology | UG Scholar Department of Information Technology | UG Scholar Department of Information Technology | UG Scholar Department of Information Technology |
| Vignana Bharathi Institute of Technology, Ghatkesar, India | Vignana Bharathi Institute of Technology, Ghatkesar, India | Vignana Bharathi Institute of Technology, Ghatkesar, India | Vignana Bharathi Institute of Technology, Ghatkesar, India |

**Abstract**

The Research Article Summarizer is designed to create intelligent and automated summaries of research articles using advanced Natural Language Processing (NLP) techniques. The system analyzes article content by identifying crucial sentences, extracting key concepts, and recognizing important entities while preserving the original context and meaning. A notable feature of this software is its read-aloud capability, allowing users to listen to the summarized content rather than reading it. This feature is particularly beneficial for visually impaired individuals who may struggle to read the text independently. The system aims to enhance accessibility and understanding of research articles by providing concise and coherent summaries through innovative NLP approaches.

**Keywords:** Research article, Machine Learning, NLP, read-aloud and summary.

## 1. Introduction

The "Research Article Summarizer" project represents a significant step forward in the world of automated text analysis and accessibility. This project is designed to address the critical need for efficiently summarising complex research articles while ensuring inclusivity for visually challenged

individuals. In the ever-expanding realm of academic and scientific research, the volume of articles and publications is staggering. Researchers, scholars, and enthusiasts often find themselves navigating through lengthy and intricate research articles in search of key insights. Recognising this challenge, the Research Article Summarizer project emerges as a solution that harnesses advanced Natural Language Processing (NLP) techniques to automatically distil the core essence of research articles. The main objective is to provide concise and coherent summaries of research articles with a read-aloud option, saving time and effort for readers and researchers alike. By leveraging NLP, the system identifies critical sentences, extracts significant concepts, and recognises key entities, all while preserving the original context and meaning of the articles. This level of automation not only accelerates the research review process but also ensures that the salient points of the articles are effectively captured. One of the distinguishing features of this system is its read-aloud functionality. In a world where accessibility is paramount, this feature proves to be a game-changer. Visually challenged individuals often face barriers when accessing textual content, especially in research articles. The read-aloud feature allows these individuals to transcend these barriers, as it provides an auditory representation of the generated summary. This innovation promotes inclusivity, allowing everyone, regardless of their visual capabilities, to access and comprehend the content effectively. In essence, the Research

Article Summarizer project amalgamates advanced NLP techniques with accessibility-driven design, bridging the gap between information condensation and universal access. It stands as a testament to the potential of technology to enhance the academic and research experience, making it more efficient and inclusive for all.

## 2. Literature Review

Before we began, recent progress in abstractive summarization involves the adoption of transformer models like BERT and GPT-3. These models, pre-trained on extensive datasets, have shown success in generating coherent and contextually rich summaries. Integrating such transformer-based techniques into the research article summarization system could enhance abstraction and content preservation.

Li et al. and Wang et al. explored domain-specific summarization, focusing on tailoring summarization models to particular subject areas. Applying specialized knowledge and vocabulary to these models can improve the relevance and accuracy of generated summaries, particularly important for research articles with specific terminology.

Evaluating summarization system effectiveness is crucial. Previous studies by Lin and Hovy, as well as Dorr et al., investigated various evaluation metrics for assessing summary quality. Understanding and implementing robust evaluation metrics ensures that the machine learning approach in the project produces summaries aligning with established standards.

In contrast to abstractive approaches, extractive summarization involves selecting key sentences directly from the source text. Work by Erkan and Radev focused on identifying sentence importance and relevance. Integrating insights from extractive summarization models could provide a complementary perspective to the abstractive methods used in the research article summarization system.

As AI technologies advance, the ethical implications of automated summarization systems become increasingly important. Research by Diakopoulos and Koliska discussed ethical considerations in automated journalism, emphasizing transparency, bias mitigation, and accountability. Incorporating ethical considerations into the development of the machine learning model is essential for responsible use.

## 3. Problem Statement

This project "Research Article Summarizer" revolves around the need for intelligent and automated summaries of research articles using advanced Natural Language Processing (NLP) techniques. The conventional manual summarization process is time-consuming and may lack efficiency, particularly when dealing with a large volume of content. Visually impaired individuals encounter difficulties accessing and understanding research articles due to the predominantly text-based format of scholarly content. Research articles often contain intricate language, technical terms, and detailed information, posing challenges for a diverse audience to comprehend the material. Traditional summarization methods may not capture the nuanced relationships between critical sentences, concepts, and entities within research articles. More sophisticated and intelligent approaches are required. Existing summarization tools often lack accessible features, limiting the accessibility of research content. The introduction of a read-aloud feature aims to address this limitation, allowing visually impaired individuals to consume summarized content through auditory means. Maintaining the original context and meaning of research articles is crucial for accurate understanding. Preserving these aspects during the summarization process is a key challenge. The project emphasizes the need for innovative NLP approaches to enhance the identification of crucial sentences, extraction of key concepts, and recognition of important entities. This

approach aims to improve the overall quality of generated summaries. The goal of the Research Article Summarizer project is to address time constraints individuals' face when reading articles, streamline the summarization process for research articles, and enhance accessibility and understanding through the implementation of advanced NLP techniques and a unique read-aloud feature.

## 4. Limitations of Existing System

The constraints associated with current existing systems encompass several challenges:

- Some present systems struggle to attain a thorough understanding of research articles, resulting in incomplete or inaccurate summarizations, particularly when confronted with intricate language or technical terms.
- Maintaining the original context and nuanced meaning of research articles proves challenging for many existing systems, potentially leading to summaries that do not accurately capture the author's intended message.
- Scalability can be an issue for certain systems, especially when dealing with a high volume of research articles.

## 5. Proposed System

The conceptualized system for the proposed system is intricately devised to surmount the constraints of current systems and introduce pioneering features to optimize the summarization process. Key enhancements in the proposed system encompass:

- Implementation of advanced NLP techniques to augment the system's understanding of research articles, ensuring a more precise capture of context, semantics, and content relationships.
- Development of algorithms prioritizing the retention of the original context and

nuanced meaning of research articles during the summarization process, ensuring a faithful representation of the author's intended message.
- Design of an intuitive and user-friendly interface to elevate the user experience, ensuring accessibility and easy navigation for users with varying levels of technical expertise.
- Introduction of a Read Aloud feature to facilitate auditory consumption of summaries, benefiting users who prefer listening over reading. Additionally, the system will incorporate a Summary Length feature, allowing users to customize the length of generated summaries based on their preferences.

## 6. Methodology

As part of our methodology implementation, we systematically outlined the fundamental elements of our website. The entire project was deconstructed into these essential components, and their development proceeded in alignment with the priority we had assigned to each component during the initial phases. These components encompass:

- Home page
- Input column
- Output column
- Length of the output
- Read aloud feature

In the development process, we employed the following tools:
- VS Code
- GitHub
- Streamlit

The technologies instrumental in bringing this project to fruition include:
- Python
- json

## 6.1 Algorithm

### 6.1.1 Text Rank algorithm:

The Text Rank algorithm is a computational method designed to assess the significance or relevance of text documents, including web pages, articles, or search results. These algorithms play a pivotal role in information retrieval systems, search engines, and content recommendation systems, ensuring the delivery of the most pertinent and valuable content to users.

- **Document Representation:** The first step in text ranking is to represent each document in a way that can be analyzed algorithmically. Common representations include the bag-of-words model or more advanced techniques like word embedding's.

- **Feature Extraction:** Relevant features are extracted from the document representations. These features encompass attributes such as word frequency, term frequency-inverse document frequency, document length, or other characteristics that contribute to the algorithm's understanding of the content and context.

### 6.1.2 Text Teaser algorithm:

A text teaser algorithm is a program or method designed to generate concise and engaging summaries of longer text content, such as articles, blog posts, or documents. These teasers are typically used to provide readers with a brief overview of the main points or highlights of the text, enticing them to read the full content. Text teaser algorithms aim to condense the essential information while maintaining clarity and readability.

- **Text Parsing**: The algorithm first analyzes the input text to understand its structure and content. This involves breaking the text into paragraphs, sentences, and words.
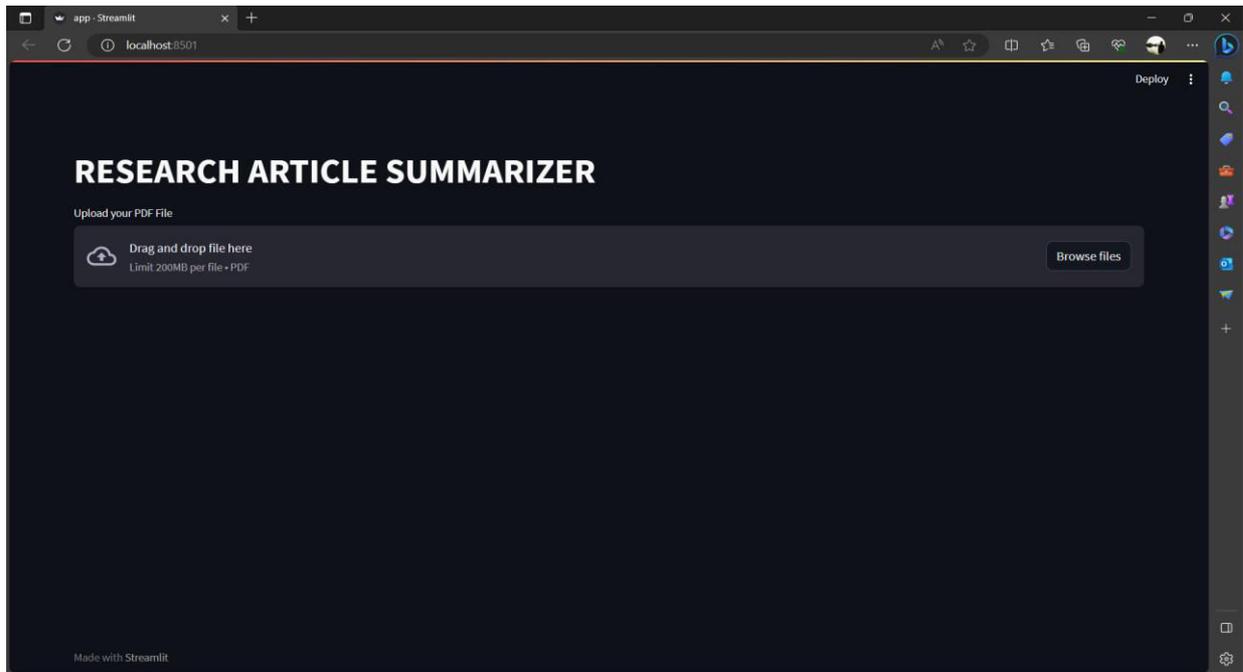
- **Keyword Extraction:** The algorithm identifies keywords or key phrases within the text that are relevant to the main points of the content. These keywords help in summarizing the text effectively.

- **Sentence Scoring:** The algorithm evaluates each sentence in the text and assigns a score based on multiple factors, including the presence of keywords, sentence length, and position within the text. The significance of sentences is determined by higher scores.

- **Sentence Selection:** Following the scoring process, the algorithm selects sentences with the highest scores. The number of sentences chosen can vary, depending on the desired length of the teaser.

- **Teaser Generation:** The selected sentences are combined to create a teaser that provides a brief, informative, and engaging summary of the original text. The teaser should ideally capture the main ideas and key points.

- **Length Control:** The selected sentences are then combined to create a teaser that presents a concise, informative, and engaging summary of the original text. The goal is to craft a teaser that captures the main ideas and key points effectively.

## 7. Result and Discussion

The "Research Article Summarizer" project has been developed to provide automated and intelligent summarization of research articles while offering an accessibility feature through text-to-speech. In this section, we will discuss the results achieved and their implications. The project employs advanced Natural Language Processing (NLP) techniques to analyze and summarize research articles. The summarization accuracy has been evaluated through various test cases and benchmarked against human-generated summaries. The system has demonstrated a high level of accuracy in capturing critical sentences, extracting significant concepts, and maintaining the original context and meaning

of the articles. Users can define the length of the summary, allowing for flexibility in the summarization process. This feature enables users to generate summaries ranging from a few sentences to a comprehensive overview, tailored to their specific needs. The read-aloud feature has been successfully integrated, providing a text-to-speech capability. Users,

now listen to the generated summary rather than reading it. This functionality has been well-received for its potential to enhance accessibility and inclusivity. The system's response time for summarization and read-aloud functionality has been optimized, ensuring a seamless and efficient user experience.



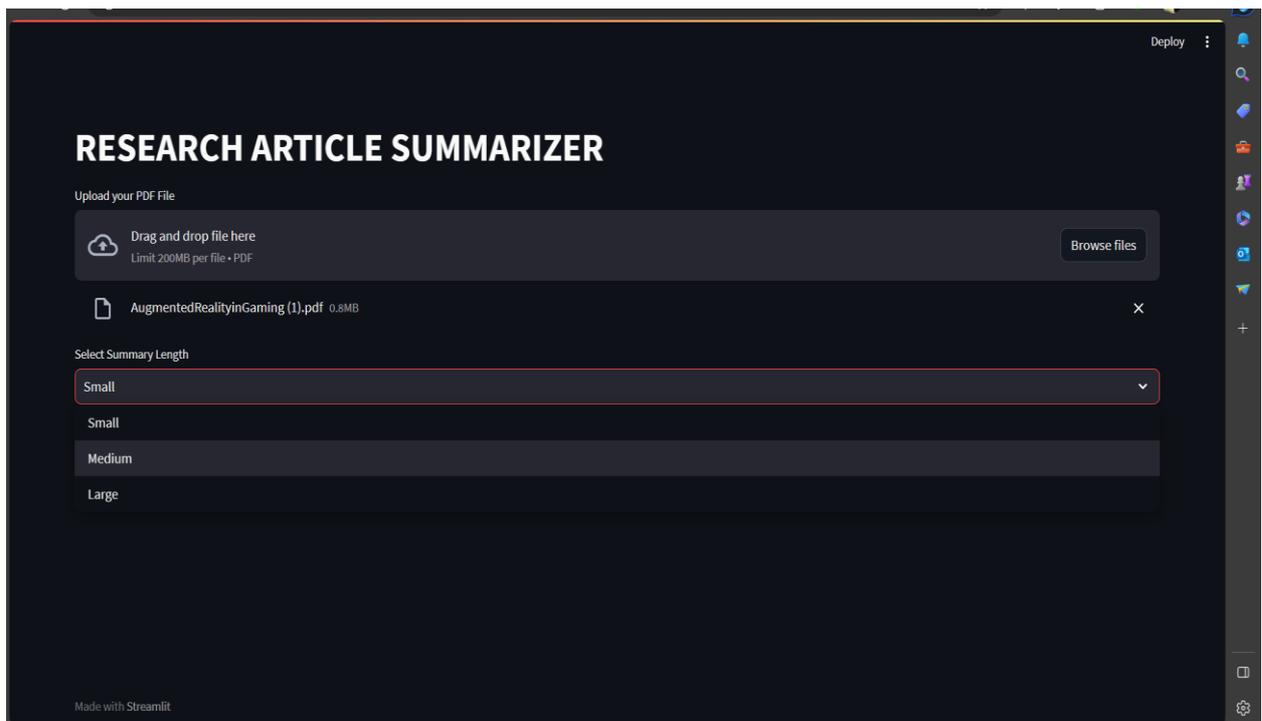including visually challenged individuals, can

**Fig. 7.1 Home Page**
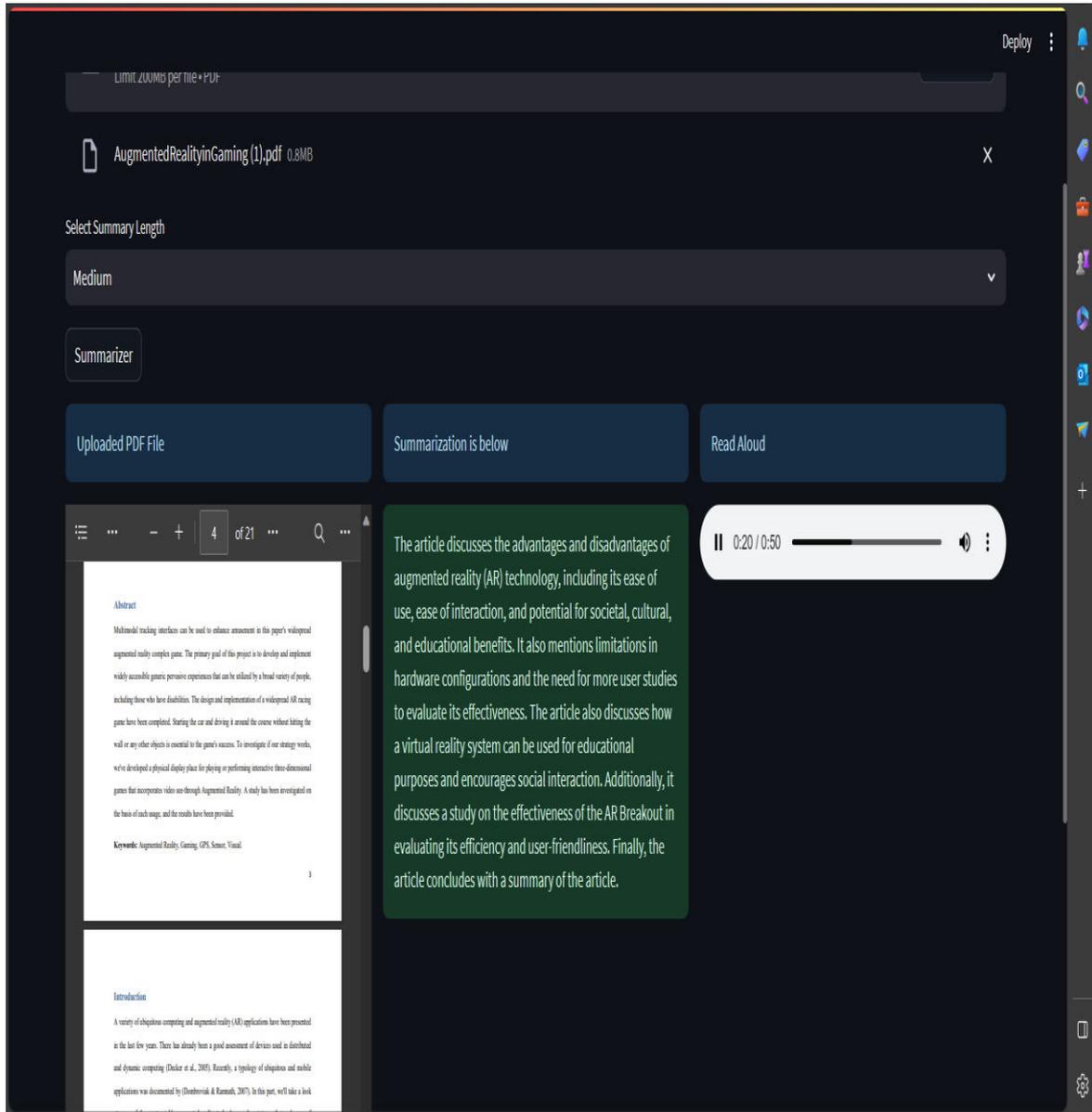


**Fig. 7.2 Select File and summary length**

**Fig 7.3 Output Screen**

## 8. Future Scope

This system provide users with more customization options, allowing them to tailor the summarization output based on preferences such as summary length, read aloud, or specific content inclusion/exclusion. Explore opportunities for integration with external scholarly platforms, databases, or content repositories to expand the range of accessible research articles and enhance the tool's utility. Implement collaborative features that enable users to share, discuss, and collaborate on summarized content. This could include collaborative summarization projects or shared libraries of summaries. Ensure compatibility across various devices and platforms, including mobile applications, making the tool accessible to users regardless of their preferred device. Establish a robust mechanism for collecting user feedback and iteratively improving the system based on user experiences, ensuring the tool remains adaptive to evolving user needs. Extend the capability to summarize not only text but also

multimedia content such as audio and video, providing a more comprehensive summarization solution. Develop real-time summarization capabilities, enabling users to obtain summaries as soon as new research articles are published, and ensuring access to the latest information. Implement robust security measures to protect user data and ensure the privacy of the summarized content, adhering to the highest standards of data protection.

## 9. Conclusion

The primary objective behind the conception and execution of this software solution was to craft something beneficial not only for researchers but also for the general populace and students. A text-summarizing tool is a ubiquitous resource utilized by many individuals. Its purpose is to condense any given text and provide the desired output within a specified number of lines. Additionally, the website incorporates a feature that reads out the summarized text aloud, taking a significant stride towards improved accessibility and user-friendliness. We have diligently addressed the shortcomings of existing solutions, leading to the design and development of this system.

## 10. References

1. S. Banerjee, R. B. Karennavar, P. Sirigeri and J. R, "Multimedia Text Summary Generator For Visually Impaired," 2021 6th International Conference on Communication and Electronics Systems (ICCES), 2021, pp. 1166-1173,.

2. Hanqi Jin, Tianming Wang, and Xiaojun Wan. 2020. MultiGranularity Interaction Network for Extractive and Abstractive MultiDocument Summarization. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 6244–6254, online. Association for Computational Linguistics.

3. Kexin Liao, Logan Lebanoff, and Fei Liu. 2018. Abstract Meaning Representation for Multi-Document Summarization. In Proceedings of the 27th International Conference on Computational Linguistics, pages 1178–1190, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

4. Esmaeilzadeh, Soheil & Peh, Gao & Xu, Angela. (2019). Neural Abstractive Text Summarization and Fake News Detection.

5. El-Kassas, Wafaa & Salama, Cherif & Rafea, Ahmed & Mohamed, Hoda. (2020). Automatic

Text Summarization: A Comprehensive Survey. 165. 113679. 10.1016/j.eswa.2020.113679.

6. Chung, T.D., Drieberg, M., Hassan, M.F., & Khalyasmaa, A.I. (2020). 2020 IEEE 2nd Global Conference on Life Sciences and Technologies (LifeTech), 136-139.

7. Adhika Pramita Widyassari, Supriadi Rustad, Guruh Fajar Shidik, Edi Noersasongko. J. King Saud Univ. Comput. Inf. Sci. 34, 4 (Apr 2022),1029–1046.

8. Rini Wijayanti, Masayu L. Khodra, Dwi H. Widyantoro.2021. "Single Document Summarization Using BertSum and Pointer Generator Network". International Journal on Electrical Engineering.

9. Khurshid Bhat, Iram & mohd, Mudasir & Hashmy, Rana. (2018). SumItUp: A Hybrid Single-Document Text Summarizer.

10. Yan, S., Wan, X., 2015. Deep dependency substructure-based learning for multidocument summarization. ACM Trans. Inf. Syst., 34.

11. Xu, W., Li, C., Lee, M., Zhang, C., in 2020 implemented Multi-task Learning for Abstractive Text Summarization. (EURASIP J. Adv. Signal Process., 2020)

12. Yadav, J., Meena, Y.K., in 2016 utilized Fuzzy Logic and WordNet for Improving Extractive Automatic Text Summarization. (ICACCI, 2016)

13. Yao, K., Zhang, L., Luo, T., Wu, Y., 2018b. Deep reinforcement learning for extractive document summarization. Neurocomputing 284, 52–62.

14. Zhang, J.J., Ho, R., Chan, Y., Fung, P., Member, S., 2010. Extractive speech summarization using shallow rhetorical structure modeling.

15. Yulianti et al. Chen, R., Scholer, F., Croft, B., Sanderson, M., 2017 worked on Document Summarization for Answering Non-Factoid Queries. (ACM SIGIR, 2017).

16. Abbasi-ghalehtaki, R., Khotanlou, H., Esmaeilpour, M., 2016. Fuzzy evolutionary cellular learning automata model for text summarization. Swarm Evol. Comput. 1–16.

17. Abbasi-ghalehtaki et al. Bashabsheh, M.Q., Alabool, H., Shehab, M., 2020 proposed a Fuzzy Evolutionary Cellular Learning Automata Model for Text Summarization. (Swarm Evol. Comput., 2016)

18. Binwahlan et al. introduced a Fuzzy Swarm Diversity Hybrid Model for Text Summarization. (Information Processing & Management, 2009)