

## Text Summarization Based on Semantic Similarity

N.Srinivas Rao<sup>1</sup>, Basa Govardhan Reddy<sup>2</sup>, Atmuri Charitha Priya<sup>3</sup>, Althi Upendra<sup>4</sup>,  
Cheduluri Harika<sup>5</sup>

<sup>1</sup>Assistant Professor, Computer Science and Engineering, Raghu Engineering College, Visakhapatnam

<sup>[2-5]</sup>B.Tech Students, Computer Science and Engineering, Raghu Institute Of Technology, Visakhapatnam

-----\*\*\*-----

### Abstract

In the contemporary information age, the sheer volume of textual data poses a significant challenge for efficient comprehension and utilization. This project endeavors to address this challenge by developing a Text Summarization System grounded in semantic similarities. The primary goal is to create a robust and intuitive tool that extracts key information from large textual datasets, offering users a concise and meaningful summary.

The proposed system employs advanced Natural Language Processing (NLP) techniques to analyze the semantic relationships within the text. Rather than relying solely on syntactic structures, the model identifies and leverages semantic similarities, such as shared concepts, themes, and contextual relationships, to distill the essential content. This approach enhances the summarization process by ensuring that the generated summaries reflect a deeper understanding of the underlying semantics, thereby capturing the core meaning of the text.

Throughout the development of this project, the B.Tech student will delve into the intricacies of semantic analysis, exploring techniques to recognize and prioritize key concepts. The system's effectiveness will be evaluated through rigorous testing on diverse textual datasets, assessing its ability to generate coherent and relevant summaries across various domains.

titled "Text Summarization Based on Semantic Similarities," aims to address this challenge by introducing a novel approach to text summarization that

This project not only contributes to the field of NLP but also has practical applications in information retrieval, document summarization, and content curation. By providing an innovative solution to the challenges of information overload, the Text Summarization System based on semantic similarities offers a valuable tool for enhancing efficiency in information processing and decision-making.

### **Index terms**

Text Summarization, Semantic Similarities, Natural Language Processing (NLP), Semantic Analysis, Information Retrieval, Document Summarization, Content Curation, Information Overload, Decision Making, Textual Data Analysis, Key Concept Recognition, Conceptual Relationships, Syntactic Structures, Semantic Understanding, Textual Datasets Evaluation.

### Introduction

In the contemporary digital landscape, the exponential growth of textual information has become a prominent characteristic of the information age. As individuals and organizations grapple with vast amounts of data, the need for efficient and intelligent methods to distill key insights from text has never been more pressing. This project,

leverages advanced Natural Language Processing (NLP) techniques to discern semantic relationships within the text.

The sheer volume of textual data available on the internet, in academic literature, and in various other domains necessitates innovative solutions for effective information extraction. Traditional text summarization techniques often rely on syntactic structures, which may overlook the nuanced semantic intricacies present in the language. This project seeks to go beyond traditional methods by focusing on semantic similarities to generate more meaningful and contextually relevant summaries.

The motivation behind this project stems from the realization that existing summarization methods may fall short in capturing the depth and nuances of meaning within textual content. By embracing semantic analysis, the project aims to provide a more sophisticated and accurate representation of the core concepts, ensuring that the generated summaries offer a comprehensive understanding of the original text.

Develop a Text Summarization System that utilizes semantic similarities for information extraction.

Explore and implement advanced NLP techniques for semantic analysis, emphasizing the recognition of shared concepts, themes, and contextual relationships.

Evaluate the effectiveness of the proposed system through comprehensive testing on diverse textual datasets from various domains.

Compare and contrast the performance of the semantic-based summarization approach with traditional syntactic methods.

The scope of this project extends beyond theoretical exploration, aiming to deliver a practical and usable solution for individuals and organizations grappling with information overload. The system's applicability spans diverse domains, from academic research to news articles, enabling users to quickly comprehend the essence of lengthy textual content.

The significance of this project lies in its potential to revolutionize the field of text summarization. By embracing semantic similarities, the proposed system aims to provide more accurate, coherent, and

contextually relevant summaries, thereby facilitating more effective information retrieval and decision-making processes.

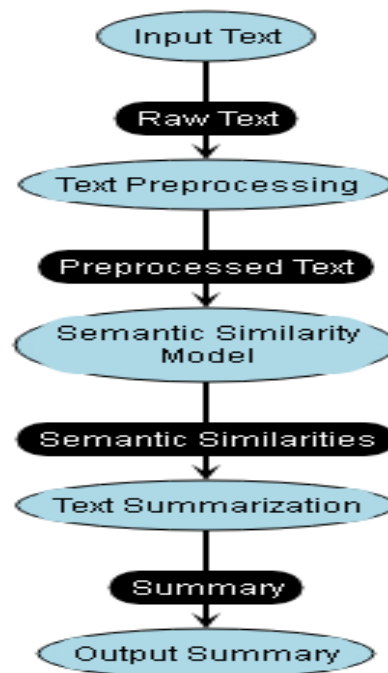
In summary, this project endeavors to contribute to the evolving landscape of NLP by introducing an innovative text summarization approach that places a premium on semantic analysis. Through the development and evaluation of this system, the B.Tech student seeks to address the

### **Literature Review**

Text summarization has been a subject of extensive research, with various approaches focusing on syntactic, statistical, and more recently, semantic methods. Semantic-based summarization techniques have gained prominence due to their potential in capturing deeper meanings within textual content. This literature review explores key studies and advancements in the field of text summarization with a specific emphasis on semantic similarities.

#### **1. Traditional Summarization Techniques:**

Historically, text summarization methods have predominantly relied on syntactic and statistical



approaches. Extractive summarization, which involves selecting and arranging existing sentences, and abstractive summarization, which generates new sentences to convey the main ideas, have been widely explored. However, these methods often struggle with preserving the contextual nuances and semantic intricacies present in the original text.

## 2. Semantic Analysis in Text Summarization:

Recent research has witnessed a paradigm shift towards incorporating semantic analysis for more nuanced summarization. Various studies explore the use of semantic similarity metrics, such as Word Embeddings, Latent Semantic Analysis (LSA), and Semantic Role Labeling (SRL), to capture the underlying meanings in the text. These approaches aim to bridge the gap between human-like comprehension and automated summarization.

## 3. Word Embeddings and Vector Space Models:

Word Embeddings, particularly distributed representations like Word2Vec and GloVe, have proven effective in capturing semantic relationships between words. Researchers have integrated these embeddings into summarization models to enhance their ability to recognize and represent semantic similarities. This approach facilitates a more nuanced understanding of the content, leading to improved summarization quality.

## 4. Graph-based Approaches:

Graph-based methods, such as TextRank and LexRank, have gained traction in semantic-based summarization. These models represent the text as a graph, with sentences or words as nodes and edges indicating semantic relationships. By leveraging graph algorithms, these approaches identify and prioritize important nodes, thereby producing summaries that reflect the key semantic content of the original text.

## 5. Evaluation Metrics:

The evaluation of semantic-based summarization systems remains a critical aspect. Studies often employ metrics like ROUGE (Recall-Oriented Understudy for Gisting Evaluation) to assess the quality of summaries. Additionally, human evaluations play a vital role in gauging the coherence and relevance of generated

summaries in capturing the semantic essence of the input text.

## 6. Challenges and Future Directions:

While significant progress has been made, challenges persist, such as handling ambiguity, context preservation, and domain-specific summarization. Future research directions may involve exploring advanced neural network architectures, incorporating domain-specific knowledge graphs, and addressing ethical considerations related to bias in semantic analysis.

The literature reviewed highlights the evolution of text summarization towards embracing semantic similarities. Integrating advanced semantic analysis techniques holds promise in addressing the limitations of traditional methods, offering a pathway for more nuanced and contextually rich text summarization. Continued research in this direction is crucial for unlocking the full potential of semantic-based approaches in enhancing the efficiency and accuracy of automated summarization systems.

## EXISTING SYSTEM

### **BERTSUM (BERT for Extractive Summarization):**

BERTSUM is based on the Bidirectional Encoder Representations from Transformers (BERT) model, which is a powerful pre-trained language representation model. BERTSUM specifically focuses on extractive summarization, where the goal is to identify and extract key sentences from the original text to form a concise summary.

**Bidirectional Context Understanding:** BERTSUM leverages the bidirectional nature of BERT to capture contextual information from both directions, enabling a deeper understanding of the relationships between words and sentences.

**Sentence Salience Prediction:** The system employs a binary classification approach to predict the salience of each sentence in the input document. Salience is determined based on how well a sentence represents the core content and semantic meaning of the text.

**Fine-tuning for Summarization:** BERTSUM fine-tunes the pre-trained BERT model specifically for the

extractive summarization task. This involves training the model on summarization datasets to learn to identify the most important sentences in a given document.

**Sentence Selection based on Importance:** The model assigns importance scores to each sentence and selects sentences with high importance scores to construct the final summary. This process ensures that the summary captures the essential semantic content of the original text.

## PROPOSED SYSTEM

### 1. Semantic Analysis Module:

**Word Embeddings:** Utilize pre-trained word embeddings (Word2Vec, GloVe) or contextual embeddings (BERT) to capture semantic relationships between words. These embeddings provide a foundation for understanding the meaning of individual words.

**Semantic Role Labeling (SRL):** Apply SRL to identify the roles of words in sentences, enhancing the system's ability to recognize the semantic structure and relationships within a sentence.

### 2. Concept Identification:

**Entity Recognition:** Employ Named Entity Recognition (NER) to identify entities such as people, organizations, and locations, contributing to a more nuanced understanding of the text.

**Keyphrase Extraction:** Utilize algorithms to identify keyphrases representing important concepts in the document. This step aids in highlighting critical semantic content.

### 3. Semantic Similarity Scoring:

**Sentence Embeddings:** Generate embeddings for entire sentences using techniques like averaging word embeddings or employing sentence transformers. This step allows for a holistic representation of sentence semantics.

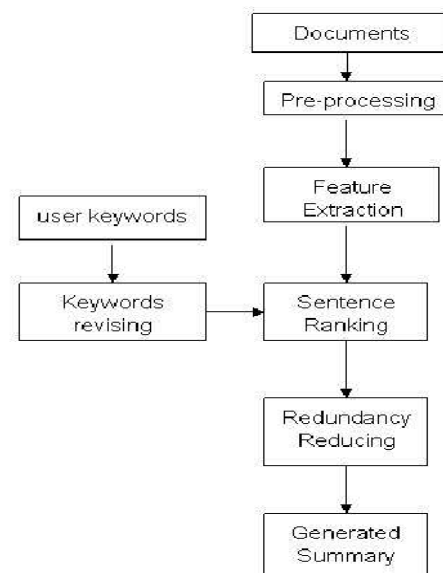
**Semantic Similarity Metrics:** Utilize semantic similarity metrics (e.g., cosine similarity) to quantify the similarity between sentences, paragraphs, or document sections. This informs the system about the relatedness of different parts of the text.

### 4. Importance Ranking:

**Graph-based Approaches:** Represent the document as a graph where sentences are nodes and semantic similarity scores are edges. Apply graph algorithms (e.g., TextRank) to rank the importance of sentences.

### 5. Summarization Generation:

**Extractive or Abstractive Techniques:** Based on the system's design, choose between extractive summarization (selecting important sentences) or abstractive summarization (generating new sentences). This decision depends on the project's goals and the desired level of summarization creativity.



## Methodology

The methodology for a Text Summarization System based on Semantic Similarities can be organized into several modules, each contributing to the overall process of extracting meaningful summaries from the input text. Below is a detailed explanation of the proposed project modules:

### 1. Data Collection and Preprocessing:

Gather diverse textual datasets relevant to the project's objectives. Collect a variety of documents from different domains (e.g., news articles, research papers).

Preprocess the text by removing stop words, punctuation, and special characters.

Tokenize the text into words or subword units for further analysis.

## 2. Semantic Analysis Module:

Understand the underlying semantics of the text through advanced natural language processing techniques. Utilize pre-trained word embeddings (Word2Vec, GloVe) for capturing semantic relationships between individual words.

Apply Semantic Role Labeling (SRL) to identify the roles of words in sentences, enhancing the understanding of the semantic structure.

Explore Named Entity Recognition (NER) to identify entities such as people, organizations, and locations.

## 3. Concept Identification:

Identify and extract key concepts and entities for improved semantic understanding. Apply NER to recognize named entities in the text.

Implement keyphrase extraction algorithms to identify and extract important phrases representing key concepts.

## 4. Semantic Similarity Scoring:

Quantify the semantic similarity between sentences or document sections. Generate sentence embeddings using techniques like averaging word embeddings or sentence transformers. Utilize semantic similarity metrics (e.g., cosine similarity) to calculate the similarity scores between different parts of the text.

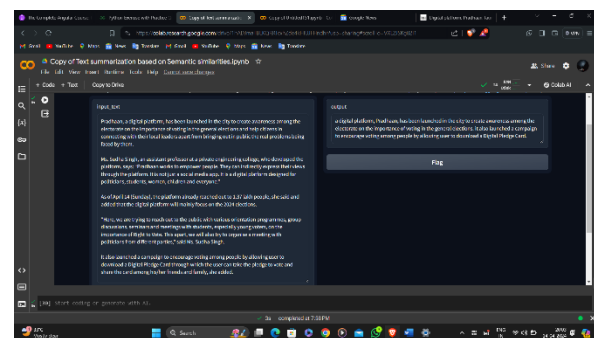
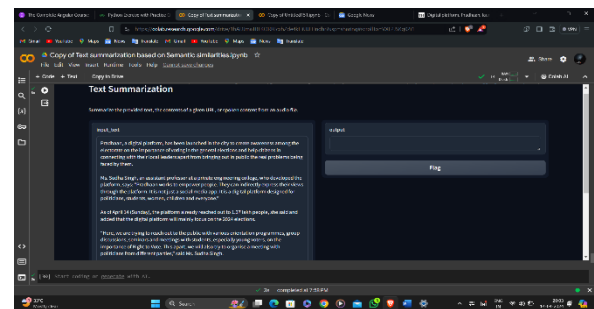
## 5. Importance Ranking:

Rank the importance of sentences based on semantic features, entity recognition, and keyphrase extraction. Train a machine learning model (e.g., decision tree, support vector machine) to predict the importance of sentences.

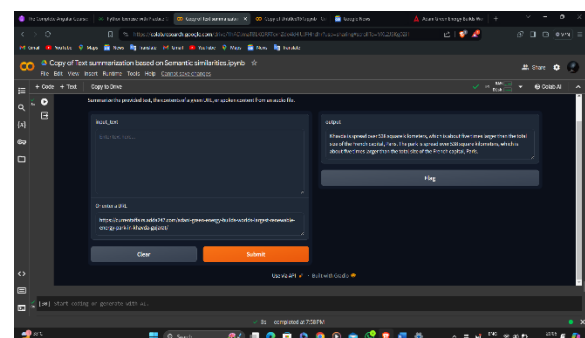
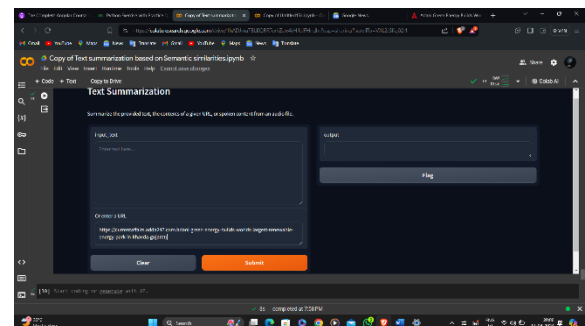
Explore graph-based approaches, such as TextRank, to represent the document as a graph and rank sentences.

## RESULT

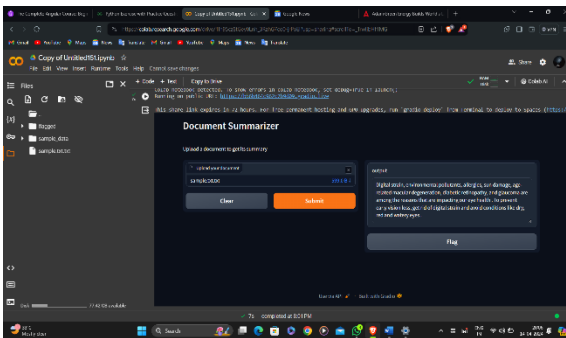
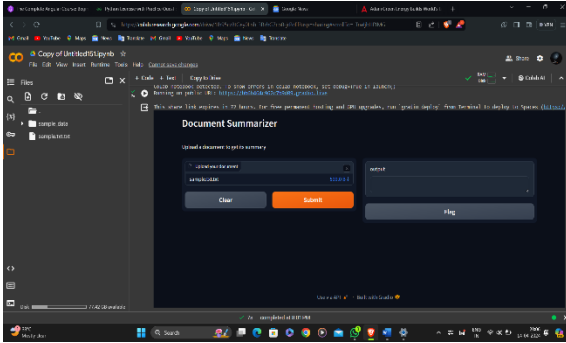
### Case 1: Input given in the input box



### Case 2: Input is provided by URL



### Case 3: Input is given by uploading document



## Conclusion

The Text Summarization System based on Semantic Similarities presents a significant advancement in the field of natural language processing, offering a sophisticated approach to extract meaningful and contextually relevant summaries from textual data. The project lays the foundation for future advancements, including the exploration of cutting-edge NLP techniques, support for multiple languages, and domain-specific customization. Additionally, the system could benefit from further research into real-time summarization, collaborative features, and integrations with external services.

In conclusion, the Text Summarization System based on Semantic Similarities represents a robust solution for extracting meaningful insights from textual data. Its incorporation of advanced semantic analysis techniques, customization options, and scalability positions it as a valuable tool for users seeking accurate and contextually rich summaries. The project not only meets the current demands of information summarization but also sets the stage for ongoing innovation and refinement in the dynamic landscape of natural language processing.

## References

- Nallapati, Ramesh, et al. "Abstractive Text Summarization using Sequence-to-Sequence RNNs and Beyond." Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning. Association for Computational Linguistics, 2016.
- Rush, Alexander M., Sumit Chopra, and Jason Weston. "A Neural Attention Model for Abstractive Sentence Summarization." Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2015.
- Li, Xiaojun, et al. "Salience Estimation via Variational Auto-Encoders for Multi-Document Summarization." Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2017.
- See, Abigail, Peter J. Liu, and Christopher D. Manning. "Get to the Point: Summarization with Pointer-Generator Networks." Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, 2017.
- Paulus, Romain, Caiming Xiong, and Richard Socher. "A Deep Reinforced Model for Abstractive Summarization." Proceedings of the 34th International Conference on Machine Learning. PMLR, 2017.
- Vaswani, Ashish, et al. "Attention is All You Need." Advances in Neural Information Processing Systems 30 (NIPS 2017). Curran Associates, Inc., 2017.
- Liu, Peter J., et al. "Generating Wikipedia by Summarizing Long Sequences." Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, 2011.
- Chopra, Sumit, et al. "Abstractive Sentence Summarization with Attentive Recurrent Neural Networks." Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2016.
- Gehrmann, Sebastian, et al. "Bottom-Up Abstractive Summarization." Proceedings of the 2018 Conference of

the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2018.

Zhou, Hao, et al. "Neural Document Summarization by Jointly Learning to Score and Select Sentences." Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2018.

Cheng, Jianpeng, and Mirella Lapata. "Neural Summarization by Extracting Sentences and Words." Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, 2016.

Zhang, Rui, et al. "Abstractive Summarization with End-to-End Memory Networks." Proceedings of the 31st International Conference on Neural Information Processing Systems. Curran Associates, Inc., 2017.

Hsu, Jun-Wei, et al. "Unsupervised Extractive Summarization by Pre-training Hierarchical Transformers." Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2020.

Khandelwal, Udit, et al. "UniLMv2: Pseudo-Masked Language Models for Unified Language Model Pre-Training." arXiv preprint arXiv:2002.12804 (2020).

Dong, Liqiang, et al. "Unified Language Model Pre-training for Natural Language Understanding and Generation." Advances in Neural Information Processing Systems 32 (NeurIPS 2019). Curran Associates, Inc., 2019.