# Text Summarization Techniques - Extractive vs. Abstractive: Exploring extractive and abstractive text summarization techniques for generating concise summaries from longer text documents or articles

*By Dr. Vijay Kumar*

*Professor of Mechanical Engineering and Robotics, University of Pennsylvania (Branch outside normal colleges)*

## Abstract

Text summarization plays a crucial role in condensing large amounts of text into shorter, more manageable summaries. This paper provides an in-depth analysis of extractive and abstractive text summarization techniques, comparing their strengths, weaknesses, and applications. Extractive summarization involves selecting important sentences or phrases from the original text, while abstractive summarization involves generating new sentences to capture the essence of the original text. We discuss various algorithms and models used in both approaches, including TF-IDF, LSA, TextRank, and neural network-based models. Additionally, we examine evaluation metrics and challenges in text summarization, such as maintaining coherence and preserving important information. Finally, we discuss potential future directions in text summarization research, including the integration of machine learning and natural language processing techniques to improve summarization quality and efficiency.

## Keywords

Text Summarization, Extractive Summarization, Abstractive Summarization, TF-IDF, LSA, TextRank, Neural Networks, Evaluation Metrics, Natural Language Processing

## I. Introduction

Text summarization is a vital task in natural language processing (NLP) and information retrieval, aiming to condense large amounts of text into shorter, more concise summaries

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

while retaining the key information. With the exponential growth of digital information, the need for effective text summarization techniques has become increasingly important in various applications, such as news summarization, document clustering, and information retrieval.

There are two primary approaches to text summarization: extractive and abstractive summarization. Extractive summarization involves selecting important sentences or phrases from the original text and arranging them to form a summary. In contrast, abstractive summarization involves generating new sentences that capture the meaning of the original text. Both approaches have their strengths and weaknesses, and their effectiveness often depends on the specific task and the characteristics of the input text.

In this paper, we provide an in-depth analysis of extractive and abstractive text summarization techniques. We discuss the algorithms and models used in each approach, including TF-IDF, Latent Semantic Analysis (LSA), TextRank, and neural network-based models. We also examine the evaluation metrics used to assess the quality of summarization, such as ROUGE (Recall-Oriented Understudy for Gisting Evaluation) and BLEU (Bilingual Evaluation Understudy).

Furthermore, we discuss the applications of text summarization in various domains, such as news summarization for providing concise updates to readers, summarization for information retrieval to quickly find relevant documents, and summarization for document clustering to organize large document collections. We also highlight the challenges and limitations of current text summarization techniques, such as maintaining coherence and preserving important information in the summary.

Finally, we discuss potential future directions in text summarization research, including the integration of machine learning and NLP advancements to improve summarization quality and efficiency. We also discuss ethical considerations in text summarization, such as ensuring fairness and avoiding bias in the summarization process.

## II. Extractive Text Summarization Techniques

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

Extractive text summarization methods aim to identify the most important sentences or phrases from the original text and present them as a summary. These techniques do not generate new sentences but instead select existing content that best represents the key information in the text. Several algorithms and approaches have been developed for extractive summarization, including TF-IDF, LSA, and TextRank.

**TF-IDF (Term Frequency-Inverse Document Frequency)**

TF-IDF is a statistical measure used to evaluate the importance of a word in a document relative to a collection of documents (corpus). The basic idea behind TF-IDF is that words that appear frequently in a document but rarely in the rest of the corpus are considered important and are likely to be included in the summary. TF-IDF assigns a weight to each word based on its frequency in the document (TF) and its rarity in the corpus (IDF).

**Latent Semantic Analysis (LSA)**

LSA is a technique used to analyze the relationships between terms and documents in a corpus. It works by constructing a matrix that represents the relationships between terms and documents based on their co-occurrence in the corpus. LSA then applies singular value decomposition (SVD) to this matrix to identify latent semantic relationships between terms and documents. In extractive summarization, LSA can be used to identify sentences that are most representative of the key concepts in the text.

**TextRank Algorithm**

TextRank is a graph-based algorithm inspired by Google's PageRank algorithm for ranking web pages. In TextRank, sentences are represented as nodes in a graph, and the relationships between sentences are represented as edges. The algorithm then iteratively calculates a score for each sentence based on the scores of its neighboring sentences, similar to how PageRank calculates the importance of a web page based on the links pointing to it. TextRank has been shown to be effective in identifying important sentences for extractive summarization.

Overall, extractive text summarization techniques offer a simple and efficient way to generate summaries by selecting important content from the original text. However, they may struggle with maintaining coherence and readability in the summary, as they rely on selecting and arranging existing content rather than generating new content.

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

### III. Abstractive Text Summarization Techniques

Abstractive text summarization techniques aim to generate a concise summary that captures the key information in the original text, often by paraphrasing and rephrasing the content. Unlike extractive summarization, which selects sentences from the original text, abstractive summarization involves generating new sentences that may not exist in the original text. This approach requires a deeper understanding of the text and the ability to generate coherent and grammatically correct sentences.

### Sequence-to-Sequence Models

Sequence-to-sequence (Seq2Seq) models, based on recurrent neural networks (RNNs) or transformer architectures, have been widely used for abstractive text summarization. These models consist of an encoder that reads the input text and a decoder that generates the summary. The encoder processes the input text and creates a representation (embedding) that captures its meaning, while the decoder uses this representation to generate the summary one word at a time.

### Attention Mechanisms in Summarization

Attention mechanisms have been incorporated into abstractive summarization models to improve the quality of the generated summaries. Attention allows the model to focus on different parts of the input text when generating each word of the summary, enabling it to capture long-range dependencies and important details. This helps in producing more accurate and contextually relevant summaries.

### Transformer-Based Models

Transformer-based models, such as BERT (Bidirectional Encoder Representations from Transformers) and GPT (Generative Pre-trained Transformer), have shown remarkable performance in various NLP tasks, including text summarization. These models leverage self-attention mechanisms to capture the context of each word in the input text, allowing them to generate more coherent and informative summaries.

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

Abstractive text summarization techniques have the advantage of being able to generate summaries that are more concise and coherent compared to extractive techniques. However, they are often more computationally intensive and may struggle with preserving important details from the original text. Nonetheless, recent advancements in deep learning and NLP have significantly improved the performance of abstractive summarization models, making them a promising approach for text summarization.

### IV. Comparison of Extractive and Abstractive Summarization

Both extractive and abstractive summarization techniques have their own strengths and weaknesses, and the choice between them often depends on the specific requirements of the task. In this section, we compare these two approaches based on various criteria, including performance, coherence, and complexity.

### Performance Metrics

Performance metrics such as ROUGE (Recall-Oriented Understudy for Gisting Evaluation) and BLEU (Bilingual Evaluation Understudy) are commonly used to evaluate the quality of summarization output. Extractive summarization techniques often achieve higher ROUGE scores compared to abstractive techniques, as they directly select content from the original text. However, ROUGE scores may not always reflect the overall quality of the summary, as they focus on overlap with reference summaries rather than coherence and readability.

### Coherence and Readability

Abstractive summarization techniques have an advantage in terms of coherence and readability, as they can generate summaries that are more concise and natural-sounding. However, generating coherent and informative summaries remains a challenging task, and abstractive models may sometimes produce summaries that are grammatically incorrect or semantically inconsistent.

### Complexity and Computational Cost

Abstractive summarization techniques are generally more complex and computationally intensive compared to extractive techniques. Models such as transformer-based models

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

require large amounts of training data and computational resources, making them less practical for real-time summarization tasks. Extractive techniques, on the other hand, are more straightforward and computationally efficient, making them suitable for tasks where speed and efficiency are important.

## Preserving Important Information

One of the challenges of extractive summarization is ensuring that important information is preserved in the summary. Extractive techniques tend to prioritize information based on its frequency or importance in the original text, which may lead to the omission of less frequent but crucial details. Abstractive techniques, while better at preserving important information, may still struggle with capturing the nuances and context of the original text.

## V. Applications of Text Summarization

Text summarization techniques have a wide range of applications in various domains, including journalism, academia, and business. In this section, we discuss some of the key applications of text summarization and how different summarization techniques are used in these contexts.

## News Summarization

One of the most common applications of text summarization is in news summarization, where large volumes of news articles are condensed into shorter summaries. Extractive summarization techniques are often used in this context, as they can quickly identify and summarize the main points of a news article. These summaries can then be used to provide readers with concise updates on current events.

## Information Retrieval

Text summarization is also used in information retrieval systems to provide users with concise summaries of documents that match their search queries. Extractive summarization techniques are commonly used in this context to generate summaries that include the most relevant information from the original documents. These summaries can help users quickly assess the relevance of a document to their query.

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

## Document Clustering

Text summarization can be used in document clustering to organize large collections of documents into meaningful clusters. By generating summaries for each document, clustering algorithms can group similar documents together based on the content of their summaries. This can help users navigate and explore large document collections more efficiently.

## Legal and Financial Summarization

In the legal and financial industries, text summarization is used to analyze and summarize large volumes of legal documents, financial reports, and other text-heavy documents. Extractive summarization techniques are often used in this context to identify key information, such as important clauses in a contract or key financial indicators in a report.

## Social Media Summarization

With the proliferation of social media platforms, there is a growing need to summarize and analyze large volumes of social media content. Abstractive summarization techniques are often used in this context to generate summaries that capture the sentiment and key themes of social media posts. These summaries can be used for sentiment analysis, trend detection, and other social media analytics tasks.

Overall, text summarization techniques have a wide range of applications across various domains, helping users quickly extract key information from large volumes of text. The choice of summarization technique depends on the specific application and the desired level of detail and coherence in the summary.

## VI. Evaluation Metrics for Text Summarization

Evaluating the quality of text summarization output is essential for assessing the effectiveness of summarization techniques. Several evaluation metrics have been developed to measure the similarity between a generated summary and one or more reference summaries. In this section, we discuss some of the commonly used evaluation metrics for text summarization.

## ROUGE (Recall-Oriented Understudy for Gisting Evaluation)

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

ROUGE is one of the most widely used evaluation metrics for text summarization. It measures the overlap between the n-grams (sequences of n words) in the generated summary and the reference summaries. ROUGE computes several variants, including ROUGE-N (which considers n-grams of varying lengths), ROUGE-L (which measures the longest common subsequence), and ROUGE-W (which considers weighted overlaps). Higher ROUGE scores indicate a higher level of similarity between the generated summary and the reference summaries.

### BLEU (Bilingual Evaluation Understudy)

BLEU is another popular evaluation metric that measures the n-gram overlap between the generated summary and the reference summaries. However, unlike ROUGE, which focuses on recall, BLEU also considers precision by penalizing the generated summary for producing n-grams that do not appear in the reference summaries. BLEU scores range from 0 to 1, with higher scores indicating better summary quality.

### Other Metrics

In addition to ROUGE and BLEU, other metrics such as METEOR (Metric for Evaluation of Translation with Explicit Ordering) and CIDEr (Consensus-based Image Description Evaluation) have been proposed for text summarization evaluation. METEOR is similar to BLEU but incorporates additional features such as stemming and synonymy matching. CIDEr, on the other hand, is designed to capture the consensus between multiple reference summaries and the generated summary, making it particularly useful for evaluating abstractive summarization techniques.

### Limitations of Evaluation Metrics

While evaluation metrics such as ROUGE and BLEU provide a quantitative measure of summary quality, they have several limitations. These metrics are based on surface-level lexical overlap and may not capture the semantic or contextual similarity between the generated summary and the reference summaries. Additionally, these metrics do not account for factors such as coherence, readability, and informativeness, which are important aspects of a good summary.

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

## VII. Challenges and Future Directions

While text summarization techniques have made significant advancements, several challenges remain in improving the quality and effectiveness of summarization. In this section, we discuss some of the key challenges and potential future directions in text summarization research.

### Coherence and Readability

One of the primary challenges in abstractive summarization is ensuring that the generated summary is coherent and readable. Current models often struggle with producing summaries that flow naturally and maintain the original meaning of the text. Future research may focus on developing techniques to improve coherence and readability, such as better modeling of context and discourse.

### Preserving Important Information

Another challenge in text summarization is preserving important information from the original text. Extractive techniques may overlook less frequent but crucial details, while abstractive techniques may fail to capture the nuances of the original text. Future research may explore techniques to better identify and preserve important information, such as using reinforcement learning to guide the summarization process.

### Integrating Machine Learning and NLP Advancements

Advancements in machine learning and NLP have the potential to significantly improve text summarization techniques. For example, the use of transformer-based models has led to significant improvements in abstractive summarization. Future research may explore the integration of these advancements into summarization models to further enhance their performance and efficiency.

### Ethical Considerations

As text summarization techniques become more sophisticated, there is a growing need to address ethical considerations in summarization. For example, there is a risk of bias in the summarization process, where certain perspectives or information may be prioritized over

others. Future research may focus on developing techniques to ensure fairness and impartiality in summarization.

## Scalability and Real-Time Summarization

Scalability is another challenge in text summarization, particularly for real-time summarization tasks where summaries need to be generated quickly. Current models may be computationally intensive and may not be suitable for real-time applications. Future research may explore techniques to improve the scalability and efficiency of summarization models for real-time use cases.

## VIII. Conclusion

Text summarization plays a crucial role in condensing large amounts of text into concise summaries, enabling users to quickly extract key information. In this paper, we have discussed extractive and abstractive text summarization techniques, comparing their strengths, weaknesses, and applications.

Extractive summarization techniques, such as TF-IDF, LSA, and TextRank, are efficient and often achieve high ROUGE scores. However, they may struggle with coherence and readability. Abstractive summarization techniques, on the other hand, can generate more coherent and readable summaries but are more complex and computationally intensive.

We have also discussed evaluation metrics for text summarization, such as ROUGE and BLEU, which provide a quantitative measure of summary quality. While these metrics are useful, they have limitations and may not fully capture the quality of a summary.

Looking ahead, there are several challenges and opportunities for future research in text summarization. Improving coherence and readability, preserving important information, and integrating machine learning and NLP advancements are key areas for improvement. Addressing ethical considerations and improving scalability for real-time summarization are also important areas for future research.

Overall, text summarization is a vibrant field with many exciting opportunities for advancement. By addressing these challenges and leveraging advancements in machine

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

learning and NLP, we can further improve the quality and effectiveness of text summarization techniques, making them more valuable in a wide range of applications.

**References**

1. Tatineni, Sumanth. "Beyond Accuracy: Understanding Model Performance on SQuAD 2.0 Challenges." *International Journal of Advanced Research in Engineering and Technology (IJARET)* 10.1 (2019): 566-581.

2. Shaik, Mahammad, Srinivasan Venkataramanan, and Ashok Kumar Reddy Sadhu. "Fortifying the Expanding Internet of Things Landscape: A Zero Trust Network Architecture Approach for Enhanced Security and Mitigating Resource Constraints." *Journal of Science & Technology* 1.1 (2020): 170-192.

3. Tatineni, Sumanth. "Cost Optimization Strategies for Navigating the Economics of AWS Cloud Services." *International Journal of Advanced Research in Engineering and Technology (IJARET)* 10.6 (2019): 827-842.

**[Journal of AI in Healthcare and Medicine](#)**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.

**Journal of AI in Healthcare and Medicine**
**Volume 2 Issue 1**
**Semi Annual Edition | Jan - June, 2022**
This work is licensed under CC BY-NC-SA 4.0.