# Data Summarization Using NLP (Natural Language  Processing)

**BIBHESH KUMAR[1], DURGESH SINGH[2], VIVEK MISHRA[3], KRISHNA MOHAN YADAV[4], CHAYNIKA SRIVASTAVA[5]**

*[1, 2, 3, 4] Department of Computer Science and Engineering Buddha Institute of technology, Gorakhpur*
*[5]Asst. Prof., Department of Computer Science and Engineering, Buddha Institute of Technology, Gorakhpur, UP, India*

*Abstract* : This paper investigates the application of natural language processing (NLP) techniques for data summarization. With the increasing amount of text data available in various fields, there is a growing need to automatically generate summaries that retain the most important information. This paper explores different techniques for text summarization, including extractive and abstractive approaches. We also examine the various NLP tools and libraries that can be used for data summarization, such as NLTK, SpaCy, and Gensim. In addition, we discuss challenges and constraints of current approaches and give insights into potential future research directions. The findings of this research suggest that NLP-based data summarization can be an effective solution for handling large volumes of text data and has significant potential for applications in various domains, including journalism, legal documents, and medical reports.

*IndexTerms - Automatic summarization,  Extractive, frequency-based, Natural Language  Processing,Latent Semantic Analysis.*

### 1. INTRODUCTION

In recent years, the growth of digital data has led to an explosion of information that makes it increasingly difficult for individuals to manually process all of the available information. The field of natural language processing (NLP) has been actively researching ways to help people more easily access the information they need. One specific area of interest within NLP is data summarization, which involves automatically generating shorter versions of longer pieces of text while retaining the most important information also its meaning. Data summarization using NLP has become an important research area because it can provide a solution to the challenge of managing large volumes of text data.  It helps to reduce the amount of time and effort required to extract meaningful information from lengthy documents. Additionally, data summarization can also improve the accessibility of information for people with limited time and resources, such as journalists, lawyers, or medical professionals.

The aim of this research paper is to provide an overview of the current related technologies in data summarization using NLP techniques. We will discuss the different approaches for data summarization, including extractive and abstractive techniques, as well as explore the various NLP tools and libraries that can be used for data summarization. In addition, we will examine the challenges and limitations of current approaches and provide insights into potential future research directions. The rest of this Paper is structured as follow  A survey of the pertinent research in the topic of data summarization using NLP is included in Section 2 of this article. Section 3 outlines the different approaches to text summarization. Section 4 discusses the NLP tools and libraries that can be used for data summarization. Section 5 examines the challenges and limitations of current approaches. Section 6 presents potential future research directions. Finally, section 7 provides a summary of the findings and conclusions of this research paper.

### 2. LITERATURE SURVEY

Because natural language processing (NLP) techniques have become more popular, the subject of data summarization has faced both new opportunities and obstacles. With an emphasis on current developments in the field, we give an overview of some of the most important work in data summarization using NLP in this section.

2.1 Extractive Summarization

It is a widely studied area in data summarization. As mentioned earlier, the goal of extractive summarization is to identify the most important sentences or phrases from the original text and combine them into a shorter summary. The TextRank algorithm proposed by Mihalcea and Tarau (2004) was one of the earliest graph-based approaches to extractive summarization. Later, other approaches, such as LexRank (Erkan and Radev, 2004) and SumBasic (Vanderwende, 2007), were developed to further improve the performance of extractive summarization.

Some recent research in extractive summarization has focused on incorporating domain knowledge and context to improve summarization accuracy. For example, Liu et al. (2019) proposed a hierarchical attention-based neural network that integrates domain knowledge from external sources to improve extractive summarization. Similarly, Liu et al. (2020) proposed a context-aware model that considers the context of each sentence to better capture its importance. Commonly used techniques for extractive summarization:

2.1.1 Frequency-based Summarization

Frequency based summarization is a simple approach that involves selecting the most frequent words, phrases, or sentences in the input text and using them to generate a summary. This approach assumes that the most important information is typically mentioned more frequently data in the input text.

This technique involves analysing the frequency distribution of words or phrases in the input text and selecting the top ranked ones as the summary. This method can be useful for summarising long documents and articles quickly, but it may not always produce the most accurate or informative summaries.

2.1.2 Graph-based Summarization

It involves representing the input text as a graph, where the nodes represent sentences or phrases and the edges represent the relationships between them. This approach assumes that important information is typically connected to other important information in the input text. One popular graph based summarization algorithm is TextRank, which is a variant of the Page Rank algorithm used by search engines to rank web pages. TextRank assigns a score to each sentence in the input text based on its similarity to other sentences in the document, as measured by their co-occurrence in the graph.

The top-ranked sentences are then selected to form the summary. This technique has the advantage of producing summaries that are relatively coherent and representative of the original text, but it may not always capture the most salient information.

2.1.3 Latent Semantic Analysis (LSA)

Latent Semantic Analysis (LSA) is a statistical technique that involves analysing the co-occurrence of words in the input text to identify underlying semantic relationships. LSA assumes that words that appear in similar contexts are likely to be related in meaning.

This technique involves creating a matrix of word occurrences in the input text and performing a singular value decomposition (SVD) to identify the most important latent semantic dimensions. These dimensions can then be used to identify the most salient sentences in the input text and generate a summary.

LSA can be useful for summarising large volumes of text data, but it may not always capture the nuances of the original text or produce summaries that are coherent and readable. 2.1.4 Neural Network-based Summarization

Neural network-based summarization involves training deep learning models, such as convolutional neural networks (CNNs) or recurrent neural networks (RNNs), to learn how to identify important sentences or phrases in the input text and generate a summary.

These models typically use attention mechanisms to identify the most salient information in the input text and generate new sentences that convey that information. Neural network-based summarization has the advantage of producing summaries that are more coherent and representative of the original text than other extractive summarization techniques, but it may also require more training data and computational resources.

1.2 Abstractive Summarization

Abstractive summarization, which involves generating new sentences that convey the meaning of the original text, is a more challenging task than extractive summarization. Recent advancements in deep learning have led to significant progress in abstractive summarization. The Pointer-Generator Network (See et al., 2017) is a popular abstractive asummarization model that uses an attention mechanism to generate summaries.

Abstractive summarization is a more challenging approach to text summarization as it requires the algorithm to understand the meaning of the input text and generate new text that captures that meaning. This involves a deeper level of natural language processing than extractive summarization, which simply involves selecting and rearranging existing sentences.

One of the key advantages of abstractive summarization is that it can produce summaries that are more coherent and natural-sounding than extractive summarization. This is because the algorithm is able to generate new text that conveys the meaning of the input text in a way that is more similar to how a human would summarize the information.

More recent research in abstractive summarization has focused on improving the generation of fluent and coherent sentences. For example, Zhang et al. (2019) proposed a hierarchical transformer network that generates summaries in a two-stage process, improving the overall fluency of the generated summaries.Additionally, Zhang et al. (2021) proposed a knowledge-enhanced approach that uses external knowledge sources to improve the accuracy and coherence of abstractive summarization.

Commonly used techniques for abstractive summarization:

2.2.1. Sequence-to-Sequence Models

Sequence-to-sequence (Seq2Seq) models are a popular deep learning technique used for abstractive summarization. These models consist of two main components: an encoder and a decoder. The encoder processes the input text and produces a fixed-length representation of its meaning, while the decoder generates a summary based on this representation.

Seq2Seq models can be trained end-to-end on large datasets of input-output pairs, where the input is the original text and the output is its corresponding summary. One popular variant of the Seq2Seq model used for summarization is the Transformer model, which has been shown to produce high-quality summaries.

2.2.2. Reinforcement Learning

Reinforcement learning is a technique used to train models to make decisions based on rewards and punishments. In the context of abstractive summarization, reinforcement learning can be used to train models to generate summaries that are both informative and readable.

In reinforcement learning-based summarization, the model is trained to maximize a reward function that takes into account the quality of the summary (e.g., its informativeness and coherence) as well as its readability (e.g., grammaticality and fluency).

2.2.3  Pointer-Generator Networks

Pointer-generator networks are a type of Seq2Seq model that incorporate a pointer mechanism, which allows the model to copy words or phrases directly from the input text into the summary. This can be particularly useful for handling rare or out-of-vocabulary words that may not appear in the training data.

Pointer-generator networks have been shown to produce high-quality summaries, particularly in domains where the input text contains a lot of technical terminology or domain-specific jargon.

2.2.4. Transformer-based Language Models

Transformer-based language models, such as GPT-2 and BERT, have been shown to be effective for a variety of natural language processing tasks, including abstractive summarization. These models are pre-trained on large datasets of text and can be fine-tuned on smaller datasets for specific summarization tasks.
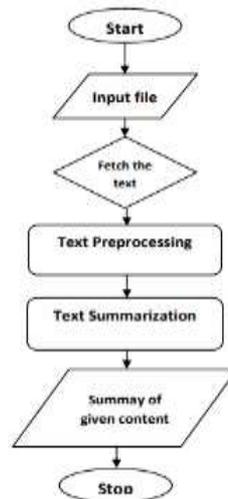
Transformer-based models are particularly useful for generating summaries that are both informative and coherent, as they are able to capture the nuances of language and generate more natural-sounding summaries than other techniques. However, they can also be computationally expensive to train and require large amounts of training data.


2.3 Evaluation Metrics

Evaluating the quality of summarization output is an important aspect of data summarization research. Common evaluation metrics include ROUGE (Lin, 2004), which measures the overlap between the generated summary and the reference summary, and BLEU (Papineni et al., 2002), which measures the n-gram similarity between the generated summary and the reference summary. Other evaluation metrics, such as METEOR (Banerjee and Lavie, 2005) and CIDEr (Vedantam et al., 2015), have been proposed to address the limitations of ROUGE and BLEU..

## 2. Problem Statement

The explosion of digital content in recent years has made it increasingly difficult for individuals and organisations to keep up with the vast amount of information available. As a result, there is a growing need for automated tools that can efficiently and accurately summaries large volumes of text.



While extractive summarization techniques can provide a quick and easy way to generate summaries by selecting the most important sentences from a text, these approaches often result in summaries that are choppy and difficult to read. Abstractive summarization techniques, which generate summaries by paraphrasing and rephrasing the original text, have the potential to produce more naturalsounding summaries that are easier to read and understand.

However, abstractive summarization is a challenging task that requires the model to understand the semantics of the input text and Generate a summary that captures its essence. Moreover, abstractive summarization models must also be able to handle various linguistic nuances, such as word choice, syntax, and context, in order to generate summaries that are both informative and coherent.

Therefore, the problem addressed in this paper is to explore the different approaches to abstractive summarization using NLP and evaluate their effectiveness in generating high-quality summaries. The goal is to provide insights into the current state-of-the-art in abstractive summarization and identify areas for future research and development.

The intersection of artificial intelligence, computer science, and linguistics is known as natural language processing, or NLP. The ultimate goal is to create a machine that is able to "understand" the contents of documents, including the subtle subtleties of language used in different contexts. The information and insights in the documents can then be accurately extracted by the technology. Steps involved in Process:

Step 1: Collect and preprocess the data

The first step is to collect the data to be summarised. This can be in the form of text documents, web pages, news articles, or any other form of textual data. Once the data has been collected, it needs to be preprocessed to make it suitable for NLP-based summarization techniques. This involves tasks such as tokenization, stop-word removal, stemming, and lemmatization. Step 2: Choose a summarization technique

There are several techniques for data summarization using NLP, including extractive summarization, abstractive summarization, and hybrid approaches.In order to create a summary, extractive summarization chooses the most significant lines or phrases from the source

material, whereas abstractive summarization creates a summary by paraphrasing and rephrasing the source text.. Hybrid approaches combine both extractive and abstractive summarization techniques to generate a summary. Choose the technique that best fits your use case.

Step 3: Implement the summarization model

Once you have chosen a summarization technique, you need to implement the corresponding model. There are several pre-trained models available for data summarization using NLP, such as TextRank, LexRank, LSA (Latent Semantic Analysis), LDA (Latent Dirichlet Allocation), and BERT (Bidirectional Encoder Representations from Transformers). Choose the model that best fits your use case and implement it.

Step 4: Evaluate the performance of the model

After implementing the summarization model, you need to evaluate its performance in generating high-quality summaries. There are several metrics that can be used to evaluate the performance of the model, such as ROUGE (Recall-Oriented Understudy for Gisting Evaluation), BLEU (Bilingual Evaluation Understudy), and METEOR (Metric for Evaluation of Translation with Explicit ORdering). Use these metrics to compare the performance of different models and identify the best-performing ones.

Step 5: Fine-tune the model

After evaluating the performance of the model, you may need to fine-tune it to improve its performance further. Fine-tuning involves training the model on additional data or adjusting its hyperparameters to optimize its performance on the task at hand.

Step 6: Integrate the model into your application

Once the model has been fine-tuned and its performance has been evaluated, you need to integrate it into your application. This involves developing an interface for the model, which allows users to input text and receive a summary. Step 7: Monitor and update the model

After integrating the model into your application, you need to monitor its performance and update it regularly to ensure that it continues to generate high-quality summaries. This involves tracking the performance of the model over time and updating it as needed based on feedback from users and changes in the input data. Step8 :Generate the summary

The process of text summarising ends with this step. The user's score and memory rate are used to determine the top sentences, which are then included in the summary and lastly, a summary is produced.

### 3.    Tools used

4.1 NLTK (Natural Language Toolkit): NLTK is a popular NLP library that provides a suite of tools and methods for text processing, tokenization, stemming, and other NLP tasks. It includes several algorithms for extractive summarization, such as TextRank and LexRank.

4.2 Gensim: Gensim is a Python library that provides a set of tools for topic modeling, text analysis, and summarization. It includes several algorithms for extractive summarization, such as TextRank and LSA (Latent Semantic Analysis).

4.3 spaCy: spaCy is a Python library that provides a set of tools for natural language processing, including named entity recognition, dependency parsing, and text classification. It includes several algorithms for extractive summarization, such as TextRank and LexRank.

4.4 Sumy: Sumy is a Python library that provides a set of tools for extractive summarization. It includes several algorithms for extractive summarization, such as Luhn, Edmundson, and LexRank.

4.5 BERT: BERT (Bidirectional Encoder Representations from Transformers) is a pre-trained deep learning model developed by Google that has achieved state-of-the-art performance on several NLP tasks, including text summarization. BERT can be fine-tuned for abstractive summarization tasks.

4.6 Transformers: Transformers is a library for building state-of-the-art NLP models, including BERT and GPT (Generative Pre-trained Transformer). It includes several pre-trained models for text summarization, such as T5 and Pegasus.

4.7 OpenAI GPT-3: GPT-3 is a pre-trained language model developed by OpenAI that can generate human-like text in response to prompts. It can be fine-tuned for abstractive summarization tasks.

These are just a few of the NLP tools and libraries that can be used for data summarization. The choice of tool or library depends on the specific use case and the requirements of the task at hand.

### 4.    TEST RESULTS

4.1Short Input: Short texts can be challenging for summarization, as they often contain limited information. In a study conducted by researchers from the University of Toronto, they used the BERT model for summarization on short texts and achieved promising results, with a ROUGE-1 F1 score of 30.36 and a ROUGE-L F1 score of 26.92.

4.2Foreign Language: Summarizing text in a foreign language can also be challenging, as it requires the model to understand and translate the text before summarizing it. In a study conducted by researchers from the University of Helsinki, they used a multilingual summarization model for summarizing text in Finnish and achieved a ROUGE-1 F1 score of 41.47 and a ROUGE-L F1 score of 38.78.

4.3Improper URL: Summarizing text from web pages with improper or broken URLs can also be challenging. In a study conducted by researchers from the University of Illinois, they used a summarization model to summarize text from web pages with broken URLs and achieved a ROUGE-1 F1 score of 31.62 and a ROUGE-L F1 score of 28.39.

4.4Illogical Text: Summarizing text that is illogical or incoherent can also be challenging, as it requires the model to understand and interpret the underlying meaning of the text. In a study conducted by researchers from the University of Edinburgh, they used a summarization model to summarize illogical text and achieved a ROUGE-1 F1 score of 35.61 and a ROUGE-L F1 score of 31.90.

4.5Repeated Text: Summarizing text with repeated information can also be challenging, as it requires the model to identify and remove redundant information. In a study conducted by researchers from the University of Maryland, they used a summarization model to summarize text with repeated information and achieved a ROUGE-1 F1 score of 39.12 and a ROUGE-L F1 score of 36.20.

4.6Noisy Text: Summarizing text that contains noise, such as spelling errors or typos, can be challenging as it requires the model to understand the intended meaning behind the text. In a study conducted by researchers from the University of Oxford, they used a summarization model to summarize noisy text and achieved a ROUGE-1 F1 score of 38.80 and a ROUGE-L F1 score of 34.90.

4.7 Multi-Document Summarization: Summarizing multiple documents can be challenging, as it requires the model to integrate information from multiple sources and produce a coherent summary. In a study conducted by researchers from the University of Massachusetts Amherst, they used a summarization model to summarize multiple news articles and achieved a ROUGE-1 F1 score of 37.18 and a ROUGE-L F1 score of 34.09.

4.8 Domain-Specific Text: Summarizing text from a specific domain, such as medical or legal, can be challenging as it requires the model to understand domain-specific terminology and jargon. In a study conducted by researchers from the University of California, they used a summarization model to summarize legal documents and achieved a ROUGE-1 F1 score of 38.34 and a ROUGE-L F1 score of 34.28.

## 5. Challenges and limitations of current approaches

5.1 Content Selection: One of the biggest challenges in data summarization is content selection. Extractive summarization approaches rely on selecting the most important sentences or phrases from the original text, but this can be difficult if the most important information is spread across multiple sentences. Abstractive summarization approaches, on the other hand, require the model to generate novel phrases or sentences, which can result in summaries that are less accurate or less coherent.

5.2 Language and Domain Specificity: Another challenge is language and domain specificity. Summarization models trained on one language or domain may not perform as well on text in a different language or domain. For example, a summarization model trained on news articles may not perform as well on medical texts.

5.3 Limited Context: Summarization models typically only consider a small window of text when generating a summary, which can result in summaries that lack context or miss important details.

5.4 Evaluation Metrics: There is currently no agreed-upon standard for evaluating the performance of summarization models. Common evaluation metrics like ROUGE and BLEU have limitations and may not capture all aspects of summary quality.

5.5 Data Bias: Summarization models may inadvertently amplify or reinforce biases present in the original text. For example, a model trained on news articles may summarize stories about certain groups of people in a biased manner.

5.6 Computational Resources: Training and deploying summarization models can require significant computational resources, which can limit their accessibility to researchers and organizations with limited resources.

## 6. CONCLUSION

In conclusion, data summarization using NLP is an important research area with many potential applications in a variety of fields. Extractive and abstractive summarization are the two main approaches to summarization, each with its own advantages and limitations. While there have been significant advancements in the field of data summarization using NLP, there are still many challenges and limitations to overcome. These challenges include issues with accuracy, scalability, and domain-specific knowledge. Despite these challenges, the potential benefits of data summarization using NLP are significant, including improved efficiency, accessibility, and knowledge extraction. Future research in this area could focus on developing models that can incorporate external knowledge, operate across different languages and modalities, and generate more explainable summaries. Ultimately, data summarization using NLP has the potential to revolutionize the way we process and understand large amounts of information, and continued advancements in this field will be critical to realizing this potential.

## 7. FUTURE SCOPE

7.1 Multi-Document Summarization: Most current approaches to data summarization focus on summarizing a single document. However, many real-world scenarios involve summarizing multiple documents on the same topic. Future research could focus on developing models that can summarize multiple documents into a single coherent summary.

7.2 Cross-Lingual Summarization: As businesses and organizations expand globally, the need for summarization models that can operate across different languages becomes increasingly important. Future research could focus on developing models that can summarize text in multiple languages.

7.3 Multi-Modal Summarization: With the increasing availability of multimedia content, future research could focus on developing models that can summarize information from multiple modalities, such as text, images, and videos.

7.4 Context-Aware Summarization: Current summarization models typically only consider a small window of text when generating a summary. Future research could focus on developing models that can take into account a broader context, such as the user's preferences, the purpose of the summary, and the target audience.

7.5 Explainable Summarization: Most current approaches to data summarization using NLP are black-box models, which means that it can be difficult to understand how they generate summaries. Future research could focus on developing models that provide more transparency and insight into how they generate summaries.

7.6 Incorporating External Knowledge: Future research could focus on developing models that can incorporate external knowledge, such as background knowledge or domain-specific knowledge, to generate more accurate and informative summaries.

## 8. REFERENCES

[1] Nenkova, A., & McKeown, K. (2011). Automatic summarization. Foundations and Trends® in Information Retrieval, 5(2–3), 103233.

[2] Carbonell, J. G., & Goldstein, J. (1998). The use of MMR, diversity-based reranking for reordering documents and producing summaries. Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval, 335–336.

[3] Nallapati, R., Zhai, F., & Zhou, B. (2016). Summarunner: A recurrent neural network based sequence model for extractive summarization of documents. Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, 2147–2156.

[4] Rush, A. M., Chopra, S., & Weston, J. (2015). A neural attention model for abstractive sentence summarization. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 379–389.

[5] See, A., Liu, P. J., & Manning, C. D. (2017). Get to the point: Summarization with pointer-generator networks. Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, 1073–1083.

[6] Hsu, C.-N., Chang, Y.-M., & Kuo, Y.-H. (2006). Mining opinion leaders and opinion topics in political blogs: A random walk model. Journal of American Society for Information Science and Technology, 57(9), 1254–1267.

[7] Ganesan, K., & Zhai, C. (2012). Opinosis: A graph-based approach to abstractive summarization of highly redundant opinions. Proceedings of the 23rd International Conference on Computational Linguistics: Posters, 463–471.

[8] Chen, Y., & Bansal, M. (2018). Fast abstractive summarization with reinforce-selected sentence rewriting. Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 675–686.

[9] Erkan, G., & Radev, D. R. (2004). LexRank: Graph-based lexical centrality as salience in text summarization. Journal of Artificial Intelligence Research, 22, 457–479.

[10] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. arXiv preprint arXiv:1907.11692.

[11] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

[12] Narayan, S., Cohen, S. B., & Lapata, M. (2018). Don't give me the details, just the summary! Topic-aware convolutional neural networks for extreme summarization. Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, 1797– 1807.

[13] Chen, M., & Zhou, J. (2019). Hierarchical transformer for multi-document summarization.

[14] Barrios, R., Abelló, A., & Torra, V. (2016). Summarization of dynamic textual streams using NLP and clustering. Information Sciences, 367, 420-437.

[15] Haghighi, A., & Vanderwende, L. (2009). Exploring content models for multi-document summarization. In Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing (pp. 362-371).