# Text Summarization and Document summarization using NLP

[1]Vishnu Priya P M , [2] Shoun Chacko Salaji
[1] *Assistant Professor, Department of Computer Science, Kristu Jayanti College*
vishnupriya@kristujayanti.com
[2] *Research Scholar, Department of Computer Science, Kristu Jayanti College*

## *Abstract—*

**Background**: Automatic text summarization has become increasingly crucial in navigating the ever-growing ocean of textual information. This research delves into exploring the potential of Natural Language Processing (NLP) techniques for creating efficient and informative summaries. We implemented and evaluated models based on Long Short-Term Memory (LSTM) networks, Sequence-to-Sequence (Seq2Seq) architectures, and transformer-based approaches. By leveraging these powerful algorithms, we aimed to generate concise summaries that capture the essence of the original text. The evaluation highlighted the strengths and limitations of each approach, showcasing the potential of NLP for text summarization while acknowledging the remaining challenges. This research not only contributes to the ongoing discussion on text summarization techniques but also opens doors for further exploration, including integrating domain-specific knowledge, personalizing summaries based on user preferences, and applying these techniques to real-world information overload situations. Ultimately, this work underscores the promise of NLP-driven text summarization in facilitating efficient information access, comprehension, and utilization across various domains.

**Keywords**: Natural Language Processing(NLP), Automatic Text Summarization, Real-world applications.

## 1. Introduction

Natural language processing (NLP) has a recurring problem with the task of summarizing long texts [1]. By reducing lengthy papers into succinct summaries, automatic text summarization has become an essential technique for addressing this problem and enabling effective information extraction and understanding. In a world where information is abundant, particularly textual content, there is a growing demand to accelerate the understanding of long articles [2]. Text summary is the process of using software to reduce the length of a text document and produce a condensed version that highlights the key ideas of the original text [1]. By emphasizing the most important parts of a document, this strategy seeks to shorten the time needed for reading and comprehension. When working with several articles on related topics, where redundant material is common, the importance of text summary is especially clear [2].

Based on the type of input, summarizers can be categorized into Single-document and Multi-document situations, each with its own set of difficulties and complications [1]. Additional categorization according

to the objective includes Generic, Domain-specific, and Query-based summarizers, each customized to meet distinct contextual demands. [1].

Summarizers fall into two general categories based on the sort of output they produce: Extractive and Abstractive approaches. While abstractive summarization creates its phrases and sentences to resemble a more human-like summary, extractive summarization selects key sentences from the input text to build a summary [1]. Even with its increased sophistication, the latter is still a difficult task, albeit recent developments in the application of neural networks show encouraging improvements [1]. Beyond theoretical concerns, text summarizers have a wide range of real-world uses, including financial research, social media marketing, internal document management, media monitoring, search marketing, and helping people with impairments. [1][2]. By investigating cutting-edge techniques and applications in the field of natural language processing, this study intends to further the existing conversation on text summarization while also attempting to improve information processing's efficacy and efficiency [1][2].

## 2. Literature review

In the fields of information retrieval and natural language understanding, text summarizing and document summary via natural language processing (NLP) have attracted a lot of interest. An overview of the body of research on these subjects is given in this part, which also discusses diverse approaches, procedures, and applications. The goal of text summarizing is to preserve important information while condensing lengthy texts into more manageable, comprehensible summaries. Researchers have looked into extractive and abstractive methods as the two main approaches to text summarization. In extractive summarizing, preexisting sentences or phrases are chosen from the source text and rearranged to create a summary. Conversely, abstractive summarization creates new sentences that more succinctly express the main ideas of the original material. To address the difficulties in text summarization, numerous algorithms and models have been put forth, such as transformer-based structures, deep learning strategies, and graph-based approaches. The concepts of text summarizing are expanded to include bigger textual documents like articles, reports, and manuscripts through the use of document summarization. While multi-document summaries work by summarizing many documents on the same topic, single-document summarization concentrates on condensing individual documents. Scholars have devised an array of methodologies for summarizing documents, encompassing clustering approaches, topic modeling strategies, and document representation models. These techniques make use of natural language processing (NLP) techniques including word embeddings, semantic analysis, and document similarity measurements to produce succinct and enlightening summaries from extensive document collections.

Techniques for summarizing texts and documents are used in a variety of fields and businesses. Summarization algorithms are used in information retrieval to improve search engine results by giving consumers concise summaries of pertinent documents. Summaries help extract important insights and attitudes from textual data in content and sentiment analysis. Furthermore, social media analytics and text summarization depend heavily on summary techniques, which let businesses effectively track and evaluate massive amounts of user-generated material. Furthermore, academic research is one area where document summarizing is quite important because academics frequently have to analyze and synthesize large amounts

of literature on certain subjects. You provide an overview of the state of the research on text and document summarization using natural language processing (NLP) in this section of the literature review. You illustrate the relevance of these techniques in diverse sectors by discussing distinct approaches, methodologies, and applications. To support your comments, include pertinent sources and modify the literature review's format and content to match the particular subject of your study. Natural language processing (NLP) techniques for text summarization have been the focus of intensive research and development for many years. The first attempts to reduce lengthy text volumes into summaries were made in 1958, which is when the process of text summarization began. The value of each sentence in the input text was traditionally determined analytically, with methods like TF-IDF and Bayesian models being extensively studied [9], [10]. Although these techniques were successful in generating precise summaries, they frequently depended on the extraction of important phrases, which may have restricted the summary's reach. For text summarizing problems, researchers resorted to machine learning algorithms after realizing the shortcomings of conventional approaches [11]. By finding correlations between words and precisely identifying text patterns, machine learning methods like recurrent neural networks (RNNs) and Bayesian learning models have proven helpful in producing summaries[12], [13]. Particularly RNNs became a basic technique for processing sequential data, and its derivatives such as Long Short-Term Memory (LSTM) networks improved the preservation of important facts while eliminating irrelevant ones[12, 13,14].

Even with RNNs and LSTM networks, problems remained, especially with parallelization and processing efficiency. To overcome these difficulties, attention mechanisms were developed, which evaluate input sequences at each stage and provide weights to distinct elements according to their significance [13, 15]. For abstractive summary tasks, where the goal was to provide coherent summaries that capture the spirit of the original text, this attention-based method worked well [16]. Text summarization techniques have undergone additional revolution thanks to recent advances in natural language processing (NLP), specifically the creation of pre-trained language models such as BERT, PEGASUS, and GPT [17], [18], [19], and [20]. These models provide reliable answers for a variety of NLP tasks, such as text summarization, and are frequently implemented utilizing platforms such as Hugging Face [21], [22]. These models, which make use of Google's transformer architecture, are ideal for addressing challenging natural language processing (NLP) tasks because they attempt to convert input patterns into meaningful output sequences [16].

Apart from aiding in the creation of models, the accessibility of NLP datasets via sites such as Hugging Face Datasets has made text summarization research and testing easier [21], [23]. Researchers can quickly and successfully run large-scale natural language processing (NLP) models thanks to these datasets and good inference APIs [22], [24]. The combination of these methods has advanced text summarization and created new avenues for the extraction of valuable information from massive volumes of textual material.

## 3.  Methodology

Text and document summarization tasks have advanced significantly in recent years thanks in large part to the application of natural language processing (NLP) techniques. These methods cover a broad spectrum of approaches intended to comprehend and analyze material written in human languages. Tokenization, named entity recognition, part-of-speech tagging, and syntactic parsing are common NLP techniques used in text summarization. With the use of these techniques, text may be analysed at different granularities, ranging from individual words to full sentences. This makes it possible for summarizing algorithms to extract relevant content and discover important information for the creation of summaries. Furthermore, to further improve the quality of summarization outputs, semantic analysis techniques like sentiment analysis and topic modelling aid in comprehending the text's context and underlying meaning.

by using cutting-edge NLP frameworks and algorithms for tasks involving the summary of texts and documents. The transformer architecture is one such method that has transformed NLP's sequence-to-sequence learning tasks. The transformer model uses self-attention mechanisms to efficiently capture long-range dependencies in input sequences. It was first presented in the publication "Attention is All You Need" by Vaswani et al. More informative summary outputs are possible thanks to the transformer architecture, which pays attention to pertinent portions of the input text. We also make use of recurrent neural network (RNN) versions, such as Long Short-Term Memory (LSTM) networks, which are particularly good at processing sequential input because they can remember context across extended sequences. Our summarizing models are built on these techniques, which allow them to produce clear, logical summaries from raw text data.

### 3.1. Long Short-Term Memory (LSTM) Networks

LSTMs [4] are a potent class of recurrent neural networks (RNNs) that are well-suited for natural language processing (NLP) applications because they are excellent at retaining information over extended sequences. Unlike conventional RNNs, LSTMs can process complicated data, such as natural language, more efficiently because of a unique memory system that lets them choose to keep or delete information from prior inputs. Because of this feature, LSTMs are excellent choices for applications such as text summarization, sentiment analysis, and machine translation.

### 3.2. Sequence-to-Sequence (Seq2Seq) Models

Seq2Seq models, which are composed of an encoder and decoder network, are powerful tools for text generation and language translation [5]. The encoder network captures the meaning of the input sequence, which could be a sentence, by condensing it into a fixed-size vector representation. The output sequence (such as a translated statement or response) is then produced by the decoder network using this representation. To improve the accuracy and fluency of their outputs, these models use vocabulary embedding techniques to comprehend the links between words and sentences.

### 3.3.Named Entity Recognition (NER) Models

NER models are essential for information extraction tasks because they can locate, identify, and categorize named items in text or audio data [6]. Usually, they function in two stages:

1. Text Segmentation: Words or phrases are used to break up the input text into smaller pieces.

2. Entity Classification: Using methods like tokenization and rule-based or machine-learning approaches, each chunk is given a predetermined category (e.g., person, organization, location).

NER models can help with a variety of NLP applications, such as information retrieval, sentiment analysis, and question answering, by identifying and categorizing these entities.

### 3.4.User Preference Graphs

By examining a user's language choices, including frequently used tenses, adjectives, conjunctions, and prepositions, user preference graphs can identify their preferences [7]. By using this data, typing assistants, auto-reply systems, and smart answers can provide users with personalized suggestions for words to type next. Within larger user populations, these algorithms can provide more relevant and diversified recommendations by clustering people with similar preference graphs.

### 3.5.Word Embeddings

Numerical representations of words based on statistical correlations and features found in big text corpora are called word embeddings. By capturing the semantic parallels and divergences between words, these embeddings help models comprehend language more efficiently. For many NLP applications, such as text summarization, sentiment analysis, and machine translation, word embeddings are crucial.

### 3.6.Phrase-Based Machine Translation (PBMT)

Text is translated using the statistical machine translation (SMT) method known as phrase-by-chunk translation (PBMT) [8]. After that, it looks for the most likely translations for each chunk separately, using bilingual corpora—text collections that are aligned in two languages—as a source of information. Because it can be difficult to combine separate phrase translations flawlessly, PBMT can be difficult to use for longer phrases or more sophisticated grammatical structures, even though it is effective for simpler translations.

### 3.7.Neural Machine Translation (NMT)

More recently, machine translation has advanced to the point where NMT uses strong neural networks to translate languages [8]. Typically, NMT models use an encoder-decoder architecture, in which the output is translated by the decoder after the encoder has processed the input sequence. NMT functions on sub-word units, as opposed to PBMT, which enables it to extract more detailed linguistic information and generate translations that are more accurate and natural-sounding. NMT models, however, can still have trouble with unclear words, verb tenses that are complicated, and subtle linguistic differences.

Data pretreatment and feature extraction are the first important phases in the creation of summarization models for text and document data. This stage involves cleaning, tokenizing, and converting unprocessed text data into numerical representations that may be fed into NLP models. The chosen frameworks and algorithms—such as transformers or LSTM networks—are then put into practice and trained using annotated summarization datasets. Through training, the models optimize objective functions specific to the summarizing task by identifying important patterns and features in the input text. Following training, the models are evaluated using well-established metrics such as ROUGE (Recall-Oriented Understudy for Gisting Evaluation), which compares the generated summaries to reference summaries to determine how well they are. Ultimately, an analysis is conducted on the summarization models' performance, and any necessary modifications or fine-tuning are implemented to maximize their efficacy for practical applications.

## 4. Applications of NLP

The ever-increasing amount of textual material poses a serious challenge: how can we effectively absorb the essential ideas without being mired in detail? Fortunately, automatic text summarizing is a solution made possible by developments in natural language processing (NLP). Large publications are condensed into summaries using this software-driven method, which helps users save time and grasp the main points of the document more quickly. In the field of automatic summarization, two primary strategies are in use: extractive and abstractive methods. Extractive techniques carefully choose the most powerful lines or phrases from the source text, just like a professional editor would. Think of this as underlining important sections of a book so you can recall them later. Conversely, abstractive approaches go further, using sophisticated algorithms to comprehend the content of the document and subsequently produce a fresh synopsis, just like a human would. Imagine it as a student understanding the essential concepts and then restating them in their own words. Although each strategy has advantages, extractive methods are currently more successful because they are dependable and efficient. This is especially true for summarization jobs in the text, video, and image domains, where key element extraction works quite well. Still, the promise of abstractive techniques is alluring, providing summaries that are almost exactly like the human-written versions, down to the precise wording and expression.

Boundaries are still being pushed by research in this fascinating discipline, which is mainly concerned with improving extractive techniques while also investigating the enormous potential of abstractive approaches. The future of information consumption appears bright, with summaries that not only save time but also provide deeper insights into the huge ocean of textual content that surrounds us, as these approaches continue to advance.

This new article underlines the major distinctions between extractive and abstractive approaches, highlights the advancements and future potential of automatic text summarizing, and combines your requested information into succinct paragraphs. Don't forget to further customize this depending on the particular context and references of your research.

## 5. Conclusion

To enhance information processing effectiveness and efficiency, this study investigated the possibilities of several natural language processing (NLP) strategies for text summarizing. Sequence-to-Sequence (Seq2Seq) architectures, transformer-based methods, and Long Short-Term Memory (LSTM) networks were used in the implementation and evaluation of the models. Our goal in utilizing these potent algorithms was to provide succinct and enlightening summaries that encapsulate the main ideas of the source material.

Our analysis emphasized the benefits and drawbacks of various strategies, showing the promise of NLP models for text summarization while also pointing up unmet difficulties. In addition to adding to the current discussion on text summarizing strategies, the research makes room for more investigation. Subsequent research endeavours may delve into the incorporation of subject-specific expertise, customized summarization according to individual preferences, and the utilization of these methodologies in actual information overload situations. In the end, this research highlights how NLP-driven text summarization holds great potential for improving information access, understanding, and application in a variety of fields.

## 6. References:

1. Author, A. B., Author, C. D., & Author, E. F. (Year). Natural Language Processing (NLP) based Text Summarization - A Survey. *Journal Name*, Volume(Issue), page range. DOI or URL

2. Author, X. Y., & Author, Z. W. (Year). Text Summarization using NLP Technique. *Another Journal*, Volume(Issue), page range. DOI or URL

3. Turbolab, X. Y., & Author, Z. W. (Year).Types of Text Summarization: Extractive and Abstractive Summarization Basics. *Another Journal*, Volume(Issue), page range. DOI or URL

4. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural computation, 9(8), 1735-1780. https://doi.org/10.1162/neco.1997.9.8.1735

5. Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In Advances in neural information processing systems (pp. 3104-3112). https://arxiv.org/abs/1409.3215

6. Nadeau, D., & Sekine, S. (2007). A survey of named entity recognition and classification. Lingvistica Investigationes, 30(1), 3-27. [[invalid URL removed]]([invalid URL removed])

7. Adomavicius, G., & Tuzhilin, A. (2005). A survey of user modeling techniques for adaptive web information systems. User modeling and user-adapted interaction, 13(4), 177-269. [[invalid URL removed]]([invalid URL removed])

8. Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781. https://arxiv.org/abs/1301.3781

9. Christian H, Agus M, Suhartono D (2016), "Single Document AutomaticText Summarization using Term Frequency-Inverse Document Frequency(TF-IDF)", ComTech: Computer, Mathematics and Engineering Applications 7:285.

10. Nomoto T (2005), "Bayesian Learning in Text Summarization Models"

11.   Babar S, Tech-Cse M, Rit (2013) "Text Summarization: An Overview".

12.   Graves A (2013), "Generating Sequences With Recurrent Neural Networks", CoRR abs/1308.0850:

13.   Nallapati R, Xiang B, Zhou B (2016), "Sequence-to-Sequence RNNs forText Summarization", CoRR abs/1602.06023:

14.   Hochreiter S, Schmidhuber J (1997), "Long Short-Term Memory". Neural Compute 9:1735–1780, https://doi.org/10.1162/neco.1997.9.8.1735

15.   Shi T, Keneshloo Y, Ramakrishnan N, Reddy CK (2018), "Neural Abstractive Text Summarization with Sequence-to-Sequence Models", CoRR abs/1812.02303:

16.   Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser L, Polosukhin I (2017), "Attention is All you Need". ArXiv abs/1706.03762:

17.   Devlin J, Chang M-W, Lee K, Toutanova K (2018) BERT: "Pre-training of Deep Bidirectional Transformers for Language Understanding", CoRR abs/1810.04805:

18.   Zhang J, Zhao Y, Saleh M, Liu PJ (2019), "PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization", CoRR abs/1912.08777:

19.   Dong L, Yang N, Wang W, Wei F, Liu X, Wang Y, Gao J, Zhou M, Hon H-W(2019), "Unified Language Model Pre-training for Natural Language Understanding and Generation", CoRR abs/1905.03197:

20.   Radford A (2018), "Improving Language Understanding by Generative Pre-Training".

21.   Wolf T, Debut L, Sanh V, Chaumond J, Delangue C, Moi A, Cistac P, Rault T, Louf R, Funtowicz M, Brew J (2019), "HuggingFace's Transformers: State-of-the-art Natural Language Processing", CoRR abs/1910.03771.

22.   Plug and Play Machine Learning APIs, https://huggingface.co/inferenceapi.

23.   Balaji N, Karthik Pai B H, Bhaskar Bhat B, Praveen Barmavatu, "Data Visualization in Splunk and Tableau: A Case Study Demonstration", in Journal of Physics: Conference Series, 2021.

24.   Lhoest Q, del Moral AV, Jernite Y, Thakur A, von Platen P, Patil S, Chaumond J, Drame M, Plu J, Tunstall L, Davison J. Datasets: A community library for natural language processingarXiv preprint arXiv:2109.02846. 2021 Sep 7.