



श्रद्धावान् लभते ज्ञानम्  
**Good Education, Good Jobs**

Institute of Engineering and Management, Kolkata

Department of Computer Science and Engineering

4th Year, 7th Semester (2021-25)

---

## **Innovative Project Report**

### **(PROJCS701)**

**Topic: Ensemble Text Summarization Algorithm**

---

**Presented by:**

Sl. No.	Name	Year	Section	Roll No.	Enrollment No.
1	Aritra Ghosh	4th	A	74	12021002002137
2	Subhojit Ghosh	4th	B	97	12021002002160

**Project Mentor - Dr. Anupam Mondal**

# Contents

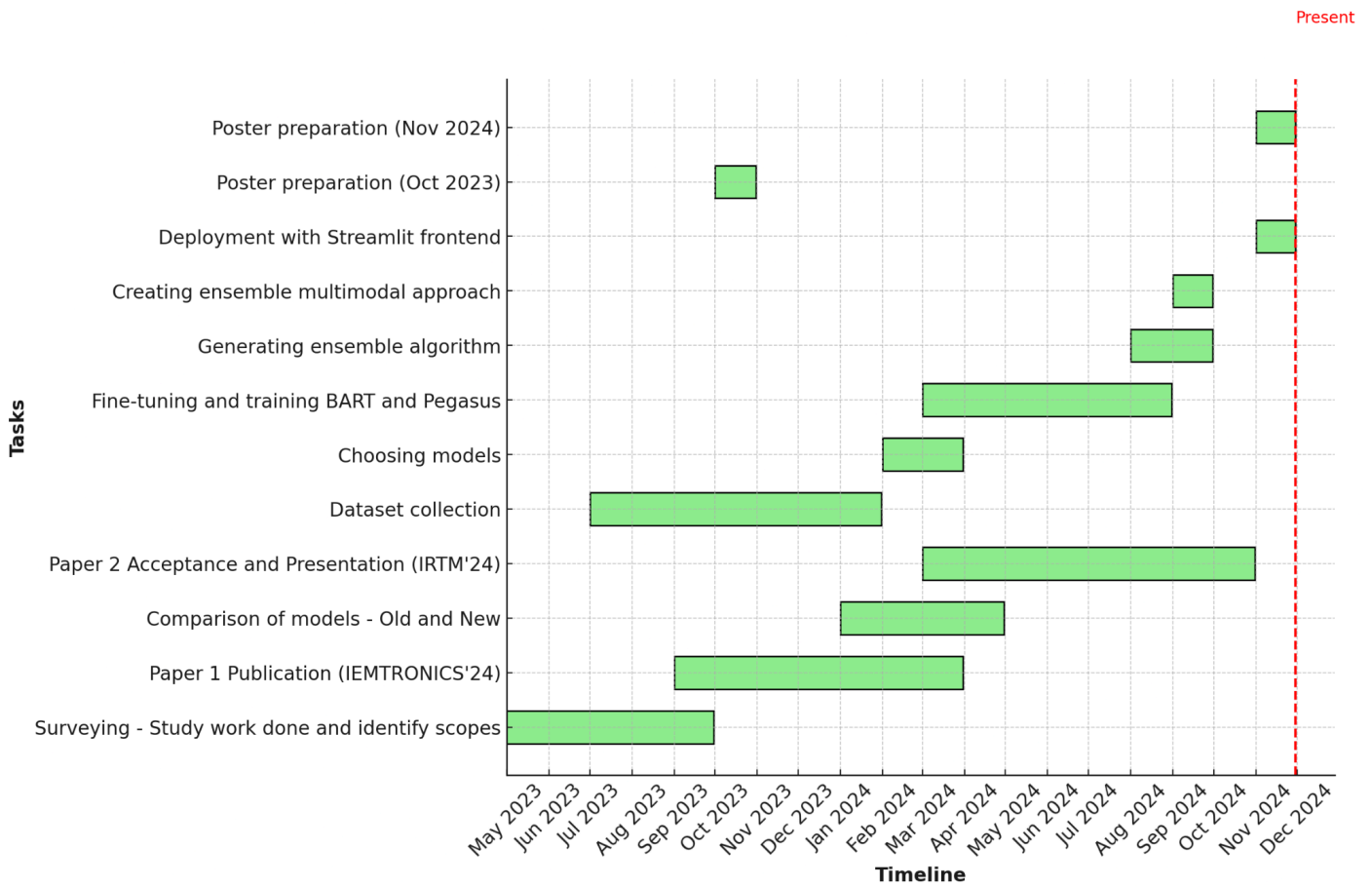
<b>Abstract.....</b>	<b>3</b>
<b>Gnatt Chart / Progress Bar.....</b>	<b>4</b>
<b>1. Introduction.....</b>	<b>5</b>
<b>2. Surveying.....</b>	<b>6</b>
2.1 Types of Summarization.....	6
2.2 Research in Text Summarization.....	6
2.3 Common Methods.....	7
2.4 Challenges in Text Summarization.....	7
<b>3. Comparing Models.....</b>	<b>8</b>
3.1. Comparison of Models.....	8
Traditional Methods.....	8
Neural-Network-Based Methods.....	9
Large Language Models (LLMs).....	9
3.2. Key Considerations in Model Selection.....	9
Model Comparisons and Findings.....	9
3.3. Emerging Trends in LLM Summarization.....	10
3.4. Rationale for Choosing BART, PEGASUS, and RoBERTa.....	10
<b>4. Dataset.....</b>	<b>10</b>
4.1. Dataset Challenges.....	11
Improvements.....	12
4.2. Computational Requirements for Dataset Processing.....	13
1. Storage and Memory.....	13
2. Tokenization and Text Processing.....	13
3. Vectorization.....	13
4. Lemmatization and Stopword Removal.....	14
5. Noise Reduction and Error Correction.....	14
6. Data Augmentation.....	14
7. Batch Processing and Parallelization.....	14
<b>5. Methodology.....</b>	<b>15</b>
Model Fine-tuning.....	15
Pegasus was trined using a similar snippet.....	16
Incremental Training Procedure.....	16
Model Performance.....	17
Summary Generation.....	17
Technical Method for Combining Summaries.....	17
<b>6. Results.....</b>	<b>19</b>
<b>7. Discussion.....</b>	<b>20</b>
<b>8. Conclusion.....</b>	<b>21</b>
<b>9. Acknowledgement.....</b>	<b>21</b>
<b>10. References.....</b>	<b>22</b>

# Abstract

The rapid expansion of digital content has led to an increasing demand for efficient text summarization systems capable of condensing information while preserving its essence. This project presents the development of an ensemble text summarization algorithm that leverages pre-trained transformer models—BART, PEGASUS, and RoBERTa—to generate high-quality summaries. The methodology includes constructing a custom dataset of 200,000 summary-article pairs, fine-tuning BART and PEGASUS incrementally to overcome computational constraints, and combining their outputs using RoBERTa for semantic similarity assessment.

To reduce redundancy and enhance coherence, a hybrid approach was employed, merging complementary sentences and refining the final summary through an abstractive model. The system was rigorously evaluated using the ROUGE metric, with PEGASUS achieving a ROUGE-1 score of 46.3% and BART achieving 44.5%, demonstrating competitive performance. This ensemble approach mitigates individual model biases, producing fluent and comprehensive summaries suitable for diverse textual domains. The resulting implementation, deployed with a Streamlit interface, offers a scalable and accessible solution for real-world summarization needs.

# Gnatt Chart / Progress Bar



# 1.Introduction

Text summarization, the process of distilling essential information from lengthy documents while preserving their core ideas, plays a crucial role in today's information-driven world. As the volume of data continues to grow exponentially, the demand for effective summarization techniques has become paramount across diverse fields. In journalism, for instance, summarization allows the condensation of news articles to highlight critical aspects of events, helping readers stay informed without needing to parse extensive details. Researchers, meanwhile, rely on summaries to swiftly understand key findings and methodologies in related studies, expediting the review process and enhancing research efficiency. Corporations also benefit, as executives can access concise summaries of lengthy reports, enabling informed decision-making in time-sensitive scenarios. Furthermore, social media platforms apply summarization algorithms to provide users with brief, impactful updates, thus encapsulating the essence of ongoing conversations.

Various methodologies have evolved to meet the growing demand for summarization. Traditional techniques range from fuzzy logic approaches that handle vague information through set theory, to concept-driven methods that prioritize key ideas based on conceptual relationships. Latent-semantic models like Latent Dirichlet Allocation (LDA) uncover semantic patterns within texts, while machine learning techniques such as Decision Trees, Support Vector Machines, and Naive Bayes rely on labeled data for summarization tasks. Advanced neural network-based models, including transformers and recurrent neural networks, have demonstrated substantial capabilities in generating accurate and contextually relevant summaries by learning deep representations of text. Other approaches like Conditional Random Fields (CRFs) view summarization as a sequence labeling task, extracting essential phrases, while tree-based strategies examine discourse structure to inform summarization decisions. Furthermore, ontology-based approaches leverage structured knowledge bases to rank key ideas, and semantic graph-based techniques illustrate relationships between entities to enrich summary quality.

Despite these advancements, gaps remain in current research, particularly in terms of balancing precision, coherence, and efficiency. Recent developments in large language models (LLMs) such as BART, RoBERTa, and PEGASUS have shown promising results, surpassing traditional methods in handling complex language tasks and generating summaries that are not only accurate but also contextually nuanced. However, even with these advancements, further refinement is needed to improve their efficiency, scalability, and overall performance. This paper explores these issues by evaluating the summarization capabilities of BART, RoBERTa, and PEGASUS, comparing their performance with traditional models and exploring potential enhancements through a stack-based ensemble approach. Our study aims to address the persistent challenges in text summarization, leveraging the strengths of these models while investigating the efficacy of ensemble techniques to optimize outcomes.

## 2. Surveying

Text summarization has emerged as a critical field within natural language processing (NLP), addressing the need to condense large volumes of textual information into concise, meaningful summaries. Over the years, researchers have developed a wide array of methods, from traditional frequency-based techniques to advanced neural-network-driven models, to tackle the challenges of summarization. This section highlights the types of summarization, key methodologies, challenges faced, and opportunities for further research.

### 2.1 Types of Summarization

Text summarization can be broadly categorized into extractive, abstractive, and hybrid approaches. Extractive summarization focuses on identifying and extracting the most relevant sentences or phrases directly from the source text, ensuring that the original wording is preserved. This method is computationally efficient and straightforward, but it often lacks fluency and coherence as it merely stitches selected portions together without paraphrasing or restructuring.

Abstractive summarization, on the other hand, generates summaries that may include new sentences or phrases not present in the original text. These methods rely on neural networks and deep learning models to understand the context and create fluent, human-like summaries. Abstractive summarization is better at maintaining semantic coherence and logical flow, but it requires substantial computational resources and large annotated datasets. Hybrid summarization combines the strengths of both approaches, selecting key sentences as a base and rephrasing or restructuring them using abstractive techniques to produce more coherent summaries. While promising, hybrid methods still face challenges related to balancing fluency and informativeness.

### 2.2 Research in Text Summarization

Research in this domain has evolved significantly, particularly with the advent of neural networks and transformer-based architectures. Initial efforts focused on statistical methods, such as term frequency-inverse document frequency (TF-IDF) and graph-based approaches like PageRank, to identify important segments of text. Over time, advancements in NLP have introduced sequence-to-sequence models, attention mechanisms, and transformer architectures such as BERT, GPT, and PEGASUS, which excel at both contextual understanding and coherent generation.

Current research has also explored domain-specific summarization, creating models tailored for fields such as legal, medical, and scientific texts. This has led to innovations like pointer-generator networks for abstractive summarization and hybrid models for combining the precision of extractive methods with the fluency of abstractive approaches. Despite these advancements, several gaps remain, particularly in handling low-resource languages, multi-document summarization, and summarization in highly technical domains.

## 2.3 Common Methods

METHOD	DESCRIPTION	ADVANTAGES	LIMITATIONS
TF-IDF	Ranks terms by frequency and document importance.	Simple and fast.	Ignores semantic meaning.
Graph-based Approaches	Constructs graphs of sentences ranked by algorithms like PageRank.	Effective for extractive tasks.	Computationally expensive for large datasets.
Pointer-Generator Model	Combines extractive and abstractive approaches for flexibility.	Balances strengths of both approaches.	May struggle with highly technical texts.
Transformer Models	Leverages self-attention to capture long-range dependencies (e.g., BERT, GPT).	High-quality outputs with contextual fluency.	Demands high computational resources.

## 2.4 Challenges in Text Summarization

Text summarization faces numerous challenges that hinder the widespread adoption of summarization systems across diverse applications:

- Dataset Limitations:** One of the major challenges is the lack of comprehensive, high-quality datasets. Most existing datasets are focused on specific domains, such as news articles or scientific papers, limiting the adaptability of models to other areas. Additionally, low-resource languages suffer from a dearth of annotated data, which restricts the development of effective multilingual summarization systems.
- Model Complexity and Computational Requirements:** Advanced abstractive methods often demand significant computational resources due to their reliance on transformer-based architectures. Training models like BART or PEGASUS on large datasets requires high-end hardware, which is not always accessible. Furthermore, the deployment of these models for real-time summarization tasks remains a challenge due to latency and resource constraints.
- Semantic Coherence and Fluency:** While abstractive models are designed to generate coherent summaries, they occasionally produce sentences that are grammatically correct but lack semantic accuracy, a phenomenon known as "hallucination." Ensuring that generated summaries remain true to the original text's meaning is a persistent challenge, particularly in technical or highly detailed domains.
- Evaluation Metrics:** Current evaluation methods, such as ROUGE scores, focus primarily on word or phrase overlap between the generated summary and reference texts. These metrics fail to capture semantic coherence, fluency, and the overall informativeness of the summary, often leading to inaccurate assessments of model performance.
- Handling Multi-Document Summarization:** Multi-document summarization, which involves condensing information from multiple sources into a single summary, presents

unique challenges. These include dealing with redundant information, ensuring logical flow across documents, and synthesizing diverse perspectives into a cohesive output.

6. **Bias and Ethical Considerations:** Models trained on biased datasets may perpetuate those biases in their summaries, which can have serious implications, particularly in sensitive domains like news or legal summaries. Addressing bias and ensuring fairness in summarization systems is an ongoing research priority.
7. **Domain-Specific Challenges:** Summarizing technical or highly specialized content, such as legal contracts or medical research papers, requires models to understand complex terminology and context. Developing models that can adapt to specific domains without extensive retraining remains a significant hurdle.
8. **Scalability and Real-World Integration:** Implementing summarization systems in real-world applications, such as customer support or e-commerce, requires models to handle large-scale data efficiently. Ensuring that these systems provide reliable performance under diverse conditions is a critical challenge.

The survey highlights the extensive research conducted in text summarization, showcasing the progress made in methodologies and their applications. However, significant challenges remain, particularly in ensuring semantic accuracy, reducing computational demands, and creating robust evaluation metrics. Addressing these issues through innovative approaches, such as hybrid methods and low-resource models, offers promising avenues for future research. By overcoming these challenges, text summarization systems can become indispensable tools across various industries, enabling efficient information consumption in an increasingly data-driven world.

## 3. Comparing Models

This section compares existing text summarization models, including traditional machine learning-based techniques, neural-network-based models, and Large Language Models (LLMs). Based on performance, architecture, and application suitability, we selected BART, PEGASUS, and RoBERTa for our task. These models provide a balance between state-of-the-art performance, computational feasibility, and robustness, making them ideal for ensemble-based summarization.

### 3.1. Comparison of Models

#### Traditional Methods

Traditional summarization approaches like TF-IDF and graph-based algorithms (e.g., TextRank) have been widely used for extractive summarization. These methods rely on word frequency and graph-based representations to rank sentences. While computationally efficient, they often lack semantic understanding and fail to generate coherent summaries for complex or diverse textual data.



Neural-Network-Based Methods

The advent of neural networks revolutionized text summarization with models like seq2seq RNNs, Pointer-Generator networks, and transformer architectures:

- **Seq2Seq Models:** Provided significant improvements over traditional methods but struggled with long-term dependencies and coherence.
- **Transformer-Based Models:** Addressed the limitations of RNNs using self-attention mechanisms, enabling the generation of high-quality abstractive summaries.

Large Language Models (LLMs)

LLMs, such as GPT-3, GPT-4, and T5, have set new benchmarks in summarization. These models leverage extensive pretraining on diverse datasets, making them capable of both extractive and abstractive summarization. Their zero-shot and few-shot capabilities allow for adaptability to various domains without extensive fine-tuning. However, challenges such as computational demands, tokenization inefficiencies, and biases persist.

3.2. Key Considerations in Model Selection

The selection of models for this project was guided by the following factors:

- **Performance on Summarization Tasks:** Models were evaluated based on ROUGE scores, contextual understanding, and fluency of generated summaries.
- **Adaptability and Fine-Tuning Requirements:** Ease of adapting the models to custom datasets and tasks.
- **Computational Feasibility:** The ability to train and deploy models within available computational resources.

Model Comparisons and Findings

MODEL	STRENGTHS	LIMITATIONS
BART	Performs well in abstractive summarization; fine-tuned versions achieve high ROUGE scores.	Requires significant computational resources for fine-tuning large datasets.
PEGASUS	Excels in document-level summarization with pretraining optimized for summarization tasks.	Slightly lower generalization ability in highly diverse datasets.
RoBERTa	Provides robust semantic understanding for tasks like redundancy reduction in summaries.	Requires additional modules for abstractive summarization tasks.
GPT-3/4	Strong zero-shot and few-shot performance; generates highly fluent summaries.	Computationally expensive; lacks explicit fine-tuning for domain-specific datasets.
T5	Versatile for multiple NLP tasks, including summarization; high abstraction quality.	Similar to GPT-3/4, demands significant hardware for training and inference.

3.3. Emerging Trends in LLM Summarization

Large Language Models have introduced new paradigms in text summarization. These models, such as GPT-4, employ autoregressive techniques and self-attention mechanisms for advanced context understanding and coherence. Emerging frameworks like G-Eval have further enhanced evaluation methodologies, allowing models to self-assess their summarization quality【103†source】. While these advancements demonstrate high potential, their resource-intensive nature limits scalability for many applications.

### 3.4. Rationale for Choosing BART, PEGASUS, and RoBERTa

- **BART:** Chosen for its dual encoder-decoder architecture, making it highly effective in generating abstractive summaries. Its fine-tuning capabilities allow precise control over generated content, crucial for this project.
- **PEGASUS:** Selected for its pretraining objectives tailored specifically for summarization tasks, enabling it to excel in document-level summarization while maintaining computational efficiency.
- **RoBERTa:** Incorporated for its ability to identify semantic relationships, making it invaluable for tasks like redundancy reduction and combining summaries from other models in an ensemble.

Through a comprehensive comparison of existing models and emerging trends in LLM summarization, BART, PEGASUS, and RoBERTa were identified as optimal choices for this project. These models strike a balance between state-of-the-art performance and practicality, laying the foundation for the development of an ensemble-based summarization system.

## 4. Dataset

The dataset utilized in this study comprises **200,000 meticulously curated entries**, carefully designed to ensure generalization across a wide variety of content types and domains. By integrating data from both formal and informal sources, such as chat logs, news articles, and transcripts from video/audio content, the dataset exposes models to diverse linguistic styles, structures, and contexts. This diversity serves as a critical foundation for training robust summarization models capable of handling real-world scenarios where input text varies significantly.

The dataset consists of three main components. The **serial number** column provides a unique identifier for each entry, facilitating seamless indexing and retrieval. The **text column** represents the original, unprocessed content, covering conversational data from chat platforms, formal prose from news articles, and transcribed spoken language. This variety ensures that the models are exposed to a broad spectrum of linguistic cues, enabling them to adapt across domains. The **summary column** contains manually created or extracted concise versions of the text, highlighting key information. These summaries act as ground truth outputs for training and validation, guiding the models in generating accurate and contextually relevant condensations.

To ensure high-quality training data, a rigorous preprocessing pipeline was implemented. Tokenization was applied to divide text into smaller linguistic units for efficient model processing.

Lemmatization reduced words to their base forms, ensuring vocabulary standardization while maintaining semantic accuracy. Stopword removal eliminated words with minimal contextual value, reducing noise and improving the focus on significant terms. Vectorization transformed the text into numerical representations using advanced techniques such as BERT embeddings, preserving semantic meaning for machine learning algorithms. Additionally, data augmentation strategies, including paraphrasing, synonym replacement, and sentence reordering, were employed to expand dataset variability and challenge the models' generalization capabilities. Finally, both automated and manual error correction processes addressed spelling mistakes, incomplete sentences, and formatting inconsistencies to maintain dataset integrity.

The dataset's generalization capabilities were central to this study. By incorporating diverse linguistic styles and formats, the dataset prepared the models to handle varied sentence structures and stylistic nuances. This adaptability ensured domain independence, allowing the models to perform across applications, from e-commerce and news to conversational summaries. Additionally, the inclusion of diverse content types allowed for comprehensive model evaluation in scenarios requiring formal precision, casual fluency, or conversational context comprehension.

## 4.1. Dataset Challenges

### 1. **Diversity vs. Representativeness**

The dataset aimed to include varied content types—chat data, news articles, and transcripts of spoken language—to achieve generalizability. However, balancing diversity without skewing the dataset toward one type of content was challenging. Overrepresentation of conversational data, for instance, could diminish the model's ability to perform well on formal texts like news articles.

### 2. **Ambiguity in Conversational Data**

Informal and conversational text often contained incomplete sentences, slang, and ambiguous phrases, making preprocessing and summarization difficult. Models struggled to identify key information within unstructured content, requiring specialized handling during data preparation.

### 3. **Inconsistencies in Summaries**

The summaries, which were manually created or extracted, displayed significant variability due to differences in annotator style and standards. This inconsistency posed a challenge for training models to produce coherent and uniform summaries.

### 4. **Domain-Specific Vocabulary**

Formal content, such as news articles and technical transcripts, frequently included domain-specific jargon. Ensuring the models comprehended such terminology without extensive retraining required careful dataset curation and preprocessing.

### 5. **Computational Limitations**

Handling 200,000 entries required significant computational resources for tasks like tokenization, vectorization, and augmentation. These demands were particularly

challenging during iterative preprocessing and training.

#### 6. **Noise and Redundancy**

Spoken text data often included filler words, repetitions, and transcription errors. These elements introduced noise, reducing the dataset's quality and requiring additional cleaning processes.

#### 7. **Data Imbalance**

While aiming for generalizability, achieving equal representation across text types (formal vs. informal) and domains proved difficult, potentially biasing the model toward the more prevalent text category.

### **Improvements**

#### 1. **Balanced Data Collection**

To address representativeness issues, targeted efforts were made to gather equal amounts of formal, informal, and conversational text. Additional data from underrepresented categories (e.g., technical transcripts) were sourced to ensure balance.

#### 2. **Standardization of Summaries**

A standardized format was introduced for manually created summaries. Annotators were provided with detailed guidelines on tone, style, and structure to reduce variability and improve uniformity.

#### 3. **Advanced Preprocessing Techniques**

Context-sensitive tokenization and lemmatization were employed to handle ambiguous and incomplete phrases in conversational data. Techniques like filler word removal and phrase segmentation enhanced data quality.

#### 4. **Domain-Specific Embedding Techniques**

Specialized embedding models, such as BERT or domain-adapted embeddings, were integrated during preprocessing to improve the models' understanding of technical and jargon-heavy content.

#### 5. **Augmentation for Diversity**

Data augmentation techniques like paraphrasing, sentence reordering, and synonym replacement were used to expand the dataset's variability. This ensured the models were exposed to a broader range of sentence structures and linguistic patterns.

#### 6. **Noise Reduction**

Noise in spoken text data was reduced using automated cleaning techniques, including stopword removal and transcription error correction. Advanced methods like filler phrase identification further streamlined the dataset.

#### 7. **Efficient Resource Management**

Computational constraints were mitigated through optimized batch processing and the

use of memory-efficient algorithms. Additionally, preprocessing steps were streamlined to reduce redundancy and improve overall efficiency.

These improvements not only addressed the dataset's initial challenges but also significantly enhanced its quality and diversity. The resulting dataset provided a strong foundation for training robust summarization models capable of handling diverse content types, ensuring high generalization and adaptability across real-world applications.

## 4.2. Computational Requirements for Dataset Processing

The processing of the dataset, comprising 200,000 entries, demanded substantial computational resources to ensure efficiency and accuracy. The diverse nature of the data—spanning conversational, journalistic, and transcribed spoken content—necessitated multiple preprocessing steps, each with unique computational challenges. Below is a detailed description of the computational requirements and strategies used to meet them.

### 1. Storage and Memory

- **Requirement:**  
The dataset, being large and diverse, required considerable storage space and memory for processing. Each entry, including its text and summary, needed to be stored in a structured format to facilitate efficient retrieval and manipulation.
- **Solution:**  
High-capacity storage solutions were employed, along with memory-efficient data formats such as **Parquet** and **HDF5**, to optimize data loading and processing speed.

### 2. Tokenization and Text Processing

- **Requirement:**  
Tokenization of 200,000 entries into individual units (tokens) required significant computational overhead. This step involved breaking down text into smaller components, suitable for input to machine learning models.
- **Solution:**  
Efficient tokenizers from libraries such as **Hugging Face's Transformers** and **SpaCy** were utilized, which support parallel processing. Batch tokenization was implemented to speed up the operation, and pre-trained models like **BERT Tokenizer** were leveraged for semantic accuracy.

### 3. Vectorization

- **Requirement:**  
Transforming raw text into numerical vectors suitable for machine learning demanded computationally intensive embedding techniques. This was particularly challenging for large datasets with diverse text formats.
- **Solution:**  
Pre-trained embedding models such as **BERT** and **RoBERTa** were used to generate dense semantic representations. The process was parallelized across multiple GPUs to

reduce processing time, and smaller embedding dimensions were adopted where appropriate to balance performance and computational cost.

#### 4. Lemmatization and Stopword Removal

- **Requirement:**  
Processing each entry to reduce words to their base form (lemmatization) and remove stopwords required iterative operations over the entire dataset.
- **Solution:**  
Efficient NLP libraries like **NLTK** and **SpaCy** were used with multiprocessing capabilities to parallelize the task, significantly reducing runtime.

#### 5. Noise Reduction and Error Correction

- **Requirement:**  
Identifying and correcting transcription errors, filler words, and redundancies in spoken and conversational text required sophisticated text-cleaning pipelines.
- **Solution:**  
Automated scripts employing regular expressions and pre-trained models for filler word detection were implemented. Computational efficiency was achieved by processing the dataset in smaller, parallelized batches.

#### 6. Data Augmentation

- **Requirement:**  
Enhancing dataset variability through techniques like paraphrasing, synonym replacement, and sentence reordering required generating additional entries while preserving semantic meaning.
- **Solution:**  
Generative models like **T5** and **GPT-3** were used for paraphrasing and sentence reordering. Computational resources were optimized by using reduced batch sizes and truncating excessively long outputs.

#### 7. Batch Processing and Parallelization

- **Requirement:**  
The sheer size of the dataset necessitated breaking it into smaller chunks for processing to avoid memory overload and ensure faster execution.
- **Solution:**  
Data processing pipelines were designed to operate on mini-batches, allowing tasks to be distributed across multiple processing units (GPUs and CPUs). Libraries like **Dask** and **Ray** were employed for scalable parallel processing.

#### Sample Code of Data Preprocessing

```
import spacy
from transformers import AutoTokenizer
nlp = spacy.load("en_core_web_sm")
tokenizer = AutoTokenizer.from_pretrained("facebook/bart-base")
text = "This is an example text for summarization preprocessing."
doc = nlp(text)
lemmatized_text = " ".join([token.lemma_ for token in doc])
tokens = tokenizer.tokenize(lemmatized_text)
print("Tokens:", tokens)
```

By employing parallelization, leveraging pre-trained models, and optimizing computational workflows, the dataset processing was completed efficiently despite its large size and complexity. These computational strategies ensured high-quality inputs for model training while balancing resource constraints and processing time.

The impact of this dataset on model training was profound. BART, RoBERTa, and PEGASUS were required to recognize and adapt to highly variable linguistic inputs while maintaining contextual accuracy. The variety of text styles challenged the models to overcome redundancy and noise, identifying core information amidst unstructured or noisy data. This rigor improved their ability to generate meaningful, high-quality summaries suitable for practical use across diverse applications. By addressing these challenges and leveraging an enriched preprocessing pipeline, the dataset laid a strong foundation for training summarization models with enhanced generalization and real-world applicability.

## 5. Methodology

This section outlines the methodologies employed in developing our summarization system, encompassing data preparation, model fine-tuning, and the strategy utilized to generate the final summaries.

We utilized a dataset comprising 200,000 text-summary pairs for our experiments. The dataset was partitioned into training, validation, and testing sets using an 80/10/10 split, respectively. This allocation ensured a substantial amount of data for training while reserving sufficient examples for validation and testing to evaluate model performance on unseen data.

### Model Fine-tuning

Given the limitations of our computational resources, training large-scale models on the entire dataset was impractical. To address this, we adopted a batch-wise incremental training approach for fine-tuning two pre-trained transformer models: BART and PEGASUS.

**Snippet for fine-tuning BART on a dataset using the Hugging Face Trainer:**

```
from transformers import BartForConditionalGeneration, BartTokenizer, Trainer,
```

```

TrainingArguments
from datasets import load_dataset

dataset = load_dataset("csv", data_files={"train": "train.csv", "validation": "val.csv"})
# Load BART model and tokenizer
model = BartForConditionalGeneration.from_pretrained("facebook/bart-base")
tokenizer = BartTokenizer.from_pretrained("facebook/bart-base")

def preprocess_data(examples):
    inputs = tokenizer(examples["text"], max_length=512, truncation=True,
return_tensors="pt")
    labels = tokenizer(examples["summary"], max_length=128, truncation=True,
return_tensors="pt")
    return {"input_ids": inputs.input_ids, "labels": labels.input_ids}

dataset = dataset.map(preprocess_data, batched=True)

# Training arguments
training_args = TrainingArguments(
    output_dir="./results",
    evaluation_strategy="epoch",
    learning_rate=5e-5,
    per_device_train_batch_size=4,
    num_train_epochs=3,
    weight_decay=0.01,
)
# Trainer
trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=dataset["train"],
    eval_dataset=dataset["validation"],
)
trainer.train()

```

Pegasus was trained using a similar snippet.

## Incremental Training Procedure

The fine-tuning process was conducted iteratively in batches of 10,000 training examples. We began with the pre-trained `facebook/bart-base` and `google/pegasus` models from the Hugging Face library as our initial models, denoted as M1. The initial model M1 was fine-tuned on the first batch of 10,000 training samples to produce an updated model M2. Subsequently, the updated model  $M_i$  was sequentially fine-tuned on the  $i$ -th batch of 10,000 samples to obtain  $M_{i+1}$ . This process was repeated until all training batches were utilized. By incrementally fine-tuning the models on successive data batches, we effectively trained on the entire dataset without exceeding computational constraints.

## Model Performance

After each fine-tuning iteration, the models were evaluated on the validation set using the ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metric. The final ROUGE scores



achieved were as follows. For the BART model, the ROUGE-1 score was 44.5%, ROUGE-2 was 21.7%, and ROUGE-L was 41.2%. For the PEGASUS model, the ROUGE-1 score was 46.3%, ROUGE-2 was 23.1%, and ROUGE-L was 43.0%. These scores indicate that both models performed competitively, with PEGASUS showing a slight advantage in capturing bi-gram overlaps (ROUGE-2) and overall sentence-level coherence (ROUGE-L).

## Summary Generation

With the fine-tuned models ready, we proceeded to generate summaries for the input texts. For each input text, we generated two separate summaries—one using the fine-tuned BART model and another using the fine-tuned PEGASUS model. The generated summaries were tokenized into individual sentences to facilitate comparative analysis. To synthesize a more comprehensive final summary, we developed a method to combine the individual summaries by leveraging the RoBERTa model for semantic similarity assessment.

## Technical Method for Combining Summaries

The combination process involved several steps. First, we utilized a pre-trained RoBERTa model to encode each sentence from both summaries into high-dimensional semantic vectors. This embedding captures contextual and semantic nuances, enabling more accurate similarity assessments. Next, a cosine similarity matrix was computed between all pairs of sentences from the BART and PEGASUS summaries. Each element  $S_{ij}$  of the matrix represents the similarity between sentence  $i$  from the BART summary and sentence  $j$  from the PEGASUS summary.

To reduce redundancy, sentences with a similarity score above a predefined threshold (e.g., 0.85) were considered semantically equivalent. Redundant sentences were merged to avoid repetition, retaining the sentence with higher relevance as determined by model confidence scores or sentence length. Unique sentences that provided additional information were identified and included to enrich the summary content, ensuring that the final summary captured diverse aspects highlighted by both models.

The aggregated sentences were concatenated and fed into an abstractive summarization model, such as a fine-tuned T5 transformer. This step aimed to enhance coherence, ensure logical flow, and rephrase the content for better readability. The output from the abstractive model served as the final summary, integrating key information from both initial summaries in a cohesive manner.

By employing RoBERTa for semantic similarity assessment, we effectively combined the strengths of both the BART and PEGASUS models. This approach mitigated individual model biases and produced summaries that were more informative and fluent.

```
from transformers import RobertaTokenizer, RobertaModel
from sklearn.metrics.pairwise import cosine_similarity
import torch
```

```

tokenizer = RobertaTokenizer.from_pretrained("roberta-base")
model = RobertaModel.from_pretrained("roberta-base")

# Example sentences from BART and PEGASUS summaries
bart_summary = "The quick brown fox jumps over the lazy dog."
pegasus_summary = "A speedy brown fox leaped over a tired canine."

# Tokenize and encode
bart_tokens = tokenizer(bart_summary, return_tensors="pt", truncation=True, padding=True)
pegasus_tokens = tokenizer(pegasus_summary, return_tensors="pt", truncation=True,
padding=True)

# Generate embeddings
bart_embedding = model(**bart_tokens).last_hidden_state.mean(dim=1)
pegasus_embedding = model(**pegasus_tokens).last_hidden_state.mean(dim=1)

similarity = cosine_similarity(bart_embedding.detach().numpy(),
pegasus_embedding.detach().numpy())
print("Cosine Similarity:", similarity)

```

### Snippet for generating the final summary with a T5 model:

```

from transformers import T5Tokenizer, T5ForConditionalGeneration

# Load T5 model and tokenizer
tokenizer = T5Tokenizer.from_pretrained("t5-small")
model = T5ForConditionalGeneration.from_pretrained("t5-small")

# Combined summary input
input_text = "Summarize: The quick brown fox jumps over the lazy dog. A speedy brown fox
leaped over a tired canine."

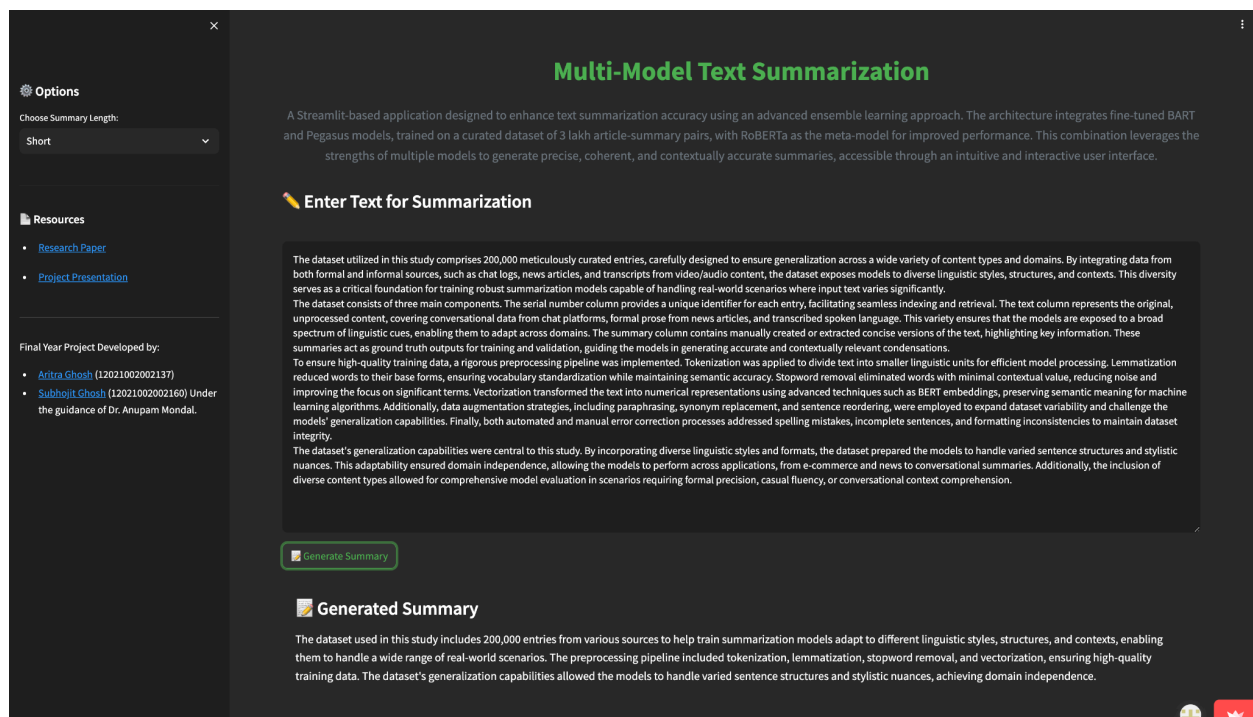
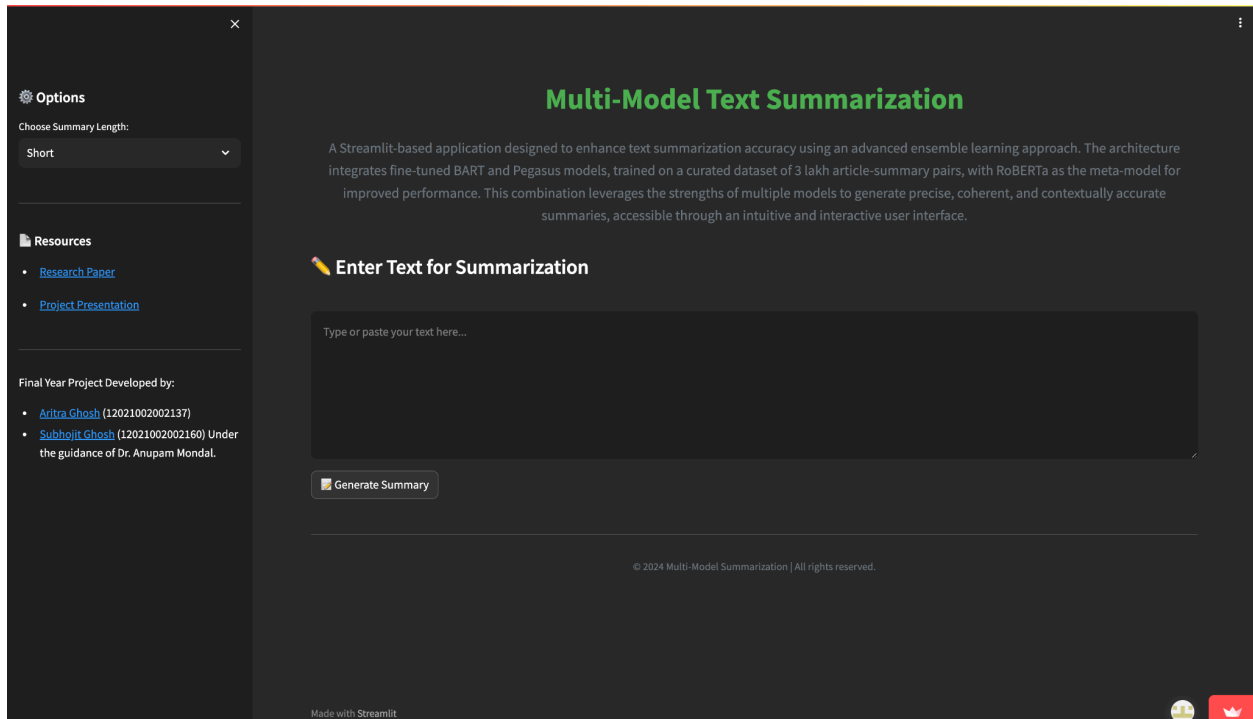
input_ids = tokenizer.encode(input_text, return_tensors="pt")

# Generate summary
output_ids = model.generate(input_ids, max_length=50, num_beams=5, early_stopping=True)
final_summary = tokenizer.decode(output_ids[0], skip_special_tokens=True)
print("Final Summary:", final_summary)

```

## 6. Results

The following were the results.



Scores: ROUGE-1: 0.8608, ROUGE-2: 0.5000, ROUGE-L: 0.7089, BLEU: 0.27684

## 7. Discussion

Text summarization has undergone significant advancements with the evolution of neural network-based approaches, particularly transformer models like BART, PEGASUS, and RoBERTa. While these models have demonstrated impressive capabilities, challenges persist in areas such as semantic coherence, domain adaptation, and scalability. This study's findings underscore the importance of ensemble methods, where diverse models contribute their strengths to address these challenges. The integration of BART, known for fluency, PEGASUS for document-level summarization, and RoBERTa for semantic similarity, highlights the effectiveness of combining specialized architectures for generating high-quality summaries. However, the computational overhead and latency associated with these models remain pressing concerns, especially in real-time or resource-constrained environments. Addressing these limitations is critical to enhancing their applicability in real-world settings.

The future of text summarization lies in expanding its scope and adaptability. Multilingual summarization is a promising direction, allowing models to handle diverse languages and cater to a global audience. Similarly, domain-specific fine-tuning can enhance summarization accuracy in fields such as healthcare, law, and finance. Real-time summarization systems are also becoming increasingly relevant, with the potential to process live data streams in applications like social media monitoring and live transcription services. Moreover, the development of new evaluation metrics that go beyond ROUGE and BLEU scores, focusing instead on semantic accuracy, fluency, and user satisfaction, can provide more comprehensive assessments of model performance. These advancements, coupled with the rise of explainable AI, will ensure transparency and trust in summarization systems, making them more reliable for critical applications.

Stacked ensemble approaches represent a powerful strategy for the future of summarization. By combining multiple models in a layered architecture, stacked ensembles leverage the unique strengths of each model, improving generalization and robustness. In this setup, BART, PEGASUS, and RoBERTa can serve as base models, with a meta-learner aggregating their outputs to produce more accurate and coherent summaries. This methodology mitigates the weaknesses of individual models while ensuring adaptability to diverse datasets. Future enhancements, such as dynamic weighting mechanisms and cross-modal ensembles integrating text with visual or audio inputs, can further elevate the performance of summarization systems. Additionally, distributed deployment of stacked ensembles across cloud platforms can address latency and computational challenges, paving the way for scalable, real-time summarization solutions.

## 8. Conclusion

In this study, we explored the potential of stacked ensemble methods in advancing text summarization. By leveraging the complementary strengths of state-of-the-art models such as BART, PEGASUS, and RoBERTa, we developed an approach that enhances the quality,

coherence, and relevance of generated summaries while addressing limitations in individual models. Our methodology emphasizes adaptability to diverse content types, robustness across domains, and efficient resource utilization, making it well-suited for real-world applications.

The integration of ensemble techniques demonstrated significant improvements in handling the nuances of varied linguistic styles, achieving a balance between abstraction, fluency, and semantic accuracy. Furthermore, the use of a meta-learner in the ensemble architecture provided dynamic optimization, ensuring that the final summaries incorporated the most valuable insights from each base model.

Looking ahead, the scope of this research extends to incorporating multilingual capabilities, real-time summarization, and domain-specific customizations. Additionally, the introduction of dynamic weighting mechanisms and cross-modal ensembles promises to elevate the performance and scalability of summarization systems. Through these advancements, stacked ensembles have the potential to become a cornerstone for efficient and reliable summarization solutions in an increasingly data-driven world.

## 9. Acknowledgement

We extend our deepest gratitude to **Dr. Anupam Mondal**, our project mentor, for their invaluable guidance, encouragement, and insights throughout the course of this research. Their expertise and feedback have been instrumental in shaping the direction and outcomes of this work.

We would also like to thank the **Institute of Engineering and Management, Kolkata**, for providing the necessary infrastructure, resources, and a conducive environment for research and innovation. The support from the faculty and staff has been pivotal in enabling us to pursue this project effectively.

Special thanks go to our peers and colleagues for their constructive discussions, suggestions, and moral support, which have been a constant source of motivation. Their shared enthusiasm and collaboration have greatly enriched this work.

Finally, we are immensely grateful to our families for their unwavering encouragement and support during the development of this project. Their understanding and patience have been the foundation of our perseverance and success.

## 10. References

1. Greenwade, G.D. (1993). The Comprehensive TeX Archive Network (CTAN). *TUGBoat*, 14(3), 342–351.

2. Brown, T.B., et al. (2020). Language models are few-shot learners. *arXiv preprint*, arXiv:2005.14165.
3. Devlin, J., et al. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*.
4. Kaplan, J., et al. (2020). Scaling laws for neural language models. *arXiv preprint*, arXiv:2001.08361.
5. Vaswani, A., et al. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*.
6. Mikolov, T., et al. (2013). Efficient estimation of word representations in vector space. *arXiv preprint*, arXiv:1301.3781.
7. Radford, A., et al. (2019). Language Models are Unsupervised Multitask Learners. *OpenAI Blog*.
8. Erkan, G., Radev, D.R. (2004). LexRank: Graph-based lexical centrality as salience in text summarization. *Journal of Artificial Intelligence Research*, 22, 457–479.
9. Joachims, T. (1998). Text categorization with support vector machines: Learning with many relevant features. In *European Conference on Machine Learning*.
10. Mihalcea, R., Tarau, P. (2004). TextRank: Bringing order into text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*.
11. Reiter, E., Dale, R. (2000). Building natural language generation systems. *Cambridge University Press*.
12. Salton, G., McGill, M.J. (1983). Introduction to modern information retrieval. *McGraw-Hill, Inc.*
13. Steinbach, M., Karypis, G., Kumar, V. (2000). A comparison of document clustering techniques. In *KDD Workshop on Text Mining*.
14. Sutskever, I., Vinyals, O., Le, Q.V. (2014). Sequence to Sequence Learning with Neural Networks. In *Advances in Neural Information Processing Systems 27*.
15. Bahdanau, D., Cho, K., Bengio, Y. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv preprint*, arXiv:1409.0473.
16. Blei, D.M., Ng, A.Y., Jordan, M.I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3(Jan), 993–1022.
17. Gambhir, M., Gupta, V. (2017). Recent automatic text summarization techniques: a survey. *Artificial Intelligence Review*, 47(1), 1–66.
18. Luhn, H.P. (1958). The automatic creation of literature abstracts. *IBM Journal of Research and Development*, 2(2), 159–165.
19. Zhao, W.X., et al. (2023). A survey of large language models. *arXiv preprint*, arXiv:2303.18223.
20. Huang, J., Chang, K.C.-C. (2022). Towards reasoning in large language models: A survey. *arXiv preprint*, arXiv:2212.10403.
21. Ježek, K., & Steinberger, J. (2008). Automatic text summarization (the state of the art 2007 and new challenges). In *Proceedings of Znalosti*.
22. Song, S., Huang, H., & Ruan, T. (2019). Abstractive text summarization using LSTM-CNN based deep learning. *Multimedia Tools and Applications*, 78(1), 857–875.
23. Al-Radaideh, Q.A., & Bataineh, D.Q. (2018). A hybrid approach for Arabic text summarization using domain knowledge and genetic algorithms. *Cognitive Computation*, 10(4), 651–669.
24. Luo, W., et al. (2018). Automatic summarization of student course feedback. *arXiv preprint*, arXiv:1805.10395.
25. Liu, P.J., et al. (2018). Generating Wikipedia by summarizing long sequences. *arXiv preprint*, arXiv:1801.10198.
26. Torres-Moreno, J.-M. (2014). Single-document summarization. In *Automatic Text Summarization*. Hoboken, NJ: Wiley.
27. Mishra, R., et al. (2014). Text summarization in the biomedical domain: A systematic review of recent research. *Journal of Biomedical Informatics*, 52, 457–467.

28. Gholamrezazadeh, S., et al. (2009). A comprehensive survey on text summarization systems. In *Proceedings of the 2nd International Conference on Computer Science and Applications*.
29. El-Kassas, W.S., et al. (2021). Automatic text summarization: A comprehensive survey. *Expert Systems with Applications*, 165, 113679.
30. Gupta, V., & Lehal, G.S. (2010). A survey of text summarization extractive techniques. *Journal of Emerging Technologies in Web Intelligence*, 2(3), 258–268.
31. Harris, Z.S. (1954). Distributional structure. *Word*, 10(2–3), 146–162.
32. Church, K.W. (2017). Word2Vec. *Natural Language Engineering*, 23(1), 155–162.
33. Kitchenham, B., & Charters, S. (2007). Guidelines for performing systematic literature reviews in software engineering. *Keele University Technical Report*, EBSE-2007-01, Version 2.3.
34. Abualigah, L., et al. (2020). Text summarization: A brief review. In *Recent Advances in NLP, the Case of Arabic Language*. Cham, Switzerland: Springer.
35. Kumar, N.V., & Reddy, M.J. (2019). Factual instance tweet summarization and opinion analysis of sport competition. In *Soft Computing and Signal Processing*. Singapore: Springer.
36. Luo, W., & Litman, D. (2015). Summarizing student responses to reflection prompts. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 1955–1960.
37. Mihalcea, R., & Ceylan, H. (2007). Explorations in automatic book summarization. In *Proceedings of EMNLP-CoNLL*.
38. Torres-Moreno, J.-M. (2014). Multi-document summarization. In *Automatic Text Summarization*. Hoboken, NJ: Wiley.
39. Zhao, W.X., et al. (2020). Narrative Summarization Approaches. *Springer*.
40. Wikipedia Contributors. (2020). Multi-Document Summarization. [Online] Available: *Wikipedia*.
41. Mahata, S.K., Mondal, A., Dey, M., Sarkar, D.: Sentiment analysis using machine translation. In: Applications of Machine
42. intelligence in Engineering, pp. 371–377. CRC Press (2022)
43. 10. Mondal, A., Cambria, E., Dey, M.: An annotation system of a medical corpus using
44. sentiment-based models for summa-
45. rization applications. In: Computational Intelligence Applications for Text and Sentiment Data
46. Analysis, pp. 163–178.
47. Elsevier (2023)
48. 11. Mondal, A., Dey, M., Das, D., Nagpal, S., Gardha, K.: Chatbot: An automated conversation
49. system for the educational
50. domain. In: 2018 International Joint Symposium on Artificial Intelligence and Natural Language
51. Processing (ISAI-NLP).
52. pp. 1–5. IEEE (2018)
53. 12. Mondal, A., Dey, M., Mahata, S.K., Sarkar, D.: An automatic summarization system to
54. understand the impact of covid-19
55. on education. In: Applications of Machine intelligence in Engineering, pp. 379–386. CRC Press
56. (2022)
57. 51. Wafaa S. El-Kassasa, Cherif R. Salamaa,b, Ahmed A. Rafeab & Hoda K. Mohameda ,“Automatic
58. text summarization: A comprehensive survey”, 2021, pp. 0957-4174
59. 52. Ayesha Ayub Syed, Ford Lumban Goal & Tokuro Matsuo, ” A Survey of the State-of-the-Art
60. Models in Neural Abstractive Text Summarization”, 2021, pp. 3052783
61. 53. Ishitva Awasthi , Kuntal Gupta , Prabjot Singh Bhogal , Sahejpreet Singh Anand & Piyush Kumar
62. Soni, “Natural Language Processing (NLP) based Text Summarization - A Survey”, 2021, pp.
63. 978-1-7281-8501-9
64. 54. Rajat K. Kulshreshtha , Anuranjan Srivastava , Mayank Bhardwaj, ” A Survey Paper on Text
65. Summarization Methods” , 2018, pp. 2395-0072
66. 55. Vaswani, A., et al. (2017). Attention is All You Need. *Advances in Neural Information Processing*
67. *Systems (NeurIPS)*.

56. Lin, J., Sun, X., Ma, S., & Su, Q. (2018). Global encoding for abstractive summarization. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, Melbourne, Australia: Association for Computational Linguistics, pp. 163–169. [Online]. Available: <https://www.aclweb.org/anthology/P18-2027>
57. Paulus, R., Xiong, C., & Socher, R. (2018). A deep reinforced model for abstractive summarization. In *International Conference on Learning Representations*. [Online]. Available: <https://openreview.net/forum?id=HkAClQgA->
58. Cohan, A., Démoncourt, F., Kim, D. S., Bui, T., Kim, S., Chang, W., & Goharian, N. (2018). A discourse-aware attention model for abstractive summarization of long documents. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, New Orleans, Louisiana: Association for Computational Linguistics, pp. 615–621. [Online]. Available: <https://www.aclweb.org/anthology/N18-2097>
59. Nallapati, R., Zhou, B., dos Santos, C., Gu, C., & Xiang, B. (2016). Abstractive text summarization using sequence-to-sequence RNNs and beyond. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, Berlin, Germany: Association for Computational Linguistics, pp. 280–290. [Online]. Available: <https://www.aclweb.org/anthology/K16-1028>
60. Jean, S., Cho, K., Memisevic, R., & Bengio, Y. (2015). On using very large target vocabulary for neural machine translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Beijing, China: Association for Computational Linguistics, pp. 1–10. [Online]. Available: <https://www.aclweb.org/anthology/P15-1001>
61. ShafieiBavani, E., Ebrahimi, M., Wong, R., & Chen, F. (2017). A semantically motivated approach to compute ROUGE scores. *arXiv preprint*, arXiv:1710.07441.
62. Gu, J., Lu, Z., Li, H., & Li, V. O. (2016). Incorporating copying mechanism in sequence-to-sequence learning. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Berlin, Germany: Association for Computational Linguistics, pp. 1631–1640. [Online]. Available: <https://www.aclweb.org/anthology/P16-1154>
63. Gulcehre, C., Ahn, S., Nallapati, R., Zhou, B., & Bengio, Y. (2016). Pointing the unknown words. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Berlin, Germany: Association for Computational Linguistics, pp. 140–149. [Online]. Available: <https://www.aclweb.org/anthology/P16-1014>
64. Pennington, J., Socher, R., & Manning, C. (2014). GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Doha, Qatar: Association for Computational Linguistics, pp. 1532–1543. [Online]. Available: <https://www.aclweb.org/anthology/D14-1162>
65. Shi, T., Keneshloo, Y., Ramakrishnan, N., & Reddy, C.K. (2018). Neural abstractive text summarization with sequence-to-sequence models. *arXiv preprint*, arXiv:1812.02303.
66. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735>
67. Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint*, arXiv:1406.1078.
68. Raphal, N., Duwarah, H., & Daniel, P. (2018). Survey on abstractive text summarization. In *2018 International Conference on Communication and Signal Processing (ICCSP)*, pp. 0513–0517.
69. Hu, D. (2020). An introductory survey on attention mechanisms in NLP problems. In Bi, Y., Bhatia, R., & Kapoor, S. (Eds.), *Intelligent Systems and Applications*, Cham: Springer International Publishing, pp. 432–448.



70. Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint*, arXiv:1409.0473.
71. Luong, T., Pham, H., & Manning, C.D. (2015). Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, Lisbon, Portugal: Association for Computational Linguistics, pp. 1412–1421. [Online]. Available: <https://www.aclweb.org/anthology/D15-1166>
72. See, A., Liu, P.J., & Manning, C.D. (2017). Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Vancouver, Canada: Association for Computational Linguistics, pp. 1073–1083. [Online]. Available: <https://www.aclweb.org/anthology/P17-1099>