

# Project Title: Text Summarization System

## Authors Information

1. Tanuboddi, Manish Reddy, tanuboddi.m@northeastern.edu
2. Nara, Kusuma, nara.k@northeastern.edu
3. Vuyyuru, Harish Reddy, vuyyuru.h@northeastern.edu

## Abstract

This study presents an automated text summarization system developed using a sequence-to-sequence (seq2seq) model and a T5 (Text-to-Text Transfer Transformer) model. T5 models are renowned for their adeptness in comprehending intricate textual relationships, rendering them particularly suited for summarization tasks. Our system undertakes the analysis of lengthy documents and articles, adeptly extracting key points and synthesizing them into concise and informative summaries. Trained on an extensive data set comprising more than 40,000 paired documents and summaries, the system showcases its proficiency in generating summaries. This technology holds promise in empowering users across diverse disciplines to effectively get valuable insights from vast expanses of text.

## Introduction

The digital age's information explosion, from research papers to articles, necessitates efficient navigation. Text summarization tackles this challenge by automatically condensing complex documents into concise summaries that capture the main points. By transforming lengthy materials into review sheets or core news summaries, this technique empowers users to grasp essential information and navigate the ever-growing sea of text.

In the age of information overload, staying on top of relevant content can be overwhelming. Our project aims to tackle this challenge by developing a text summarization system. This system will leverage the latest advancements in Natural Language Processing, specifically pre-trained Transformer models like T5. These models are exceptional at understanding complex relationships within text, making them ideal for tasks like summarization.

Our system will automatically analyze lengthy documents and articles, extracting the key points and condensing them into concise and informative summaries. This will empower users across various fields, from journalists and academics to business professionals, to navigate the information landscape efficiently and gain valuable insights from vast amounts of text.

## Motivation:

The project is driven by the increasing demand for efficient information extraction tools across various fields. Here's a breakdown of our key motivations:

- Information Overload: The digital age has led to an explosion of textual data, making it difficult to keep up with and extract crucial themes.

- **Improved Decision-Making:** Distilling key points from lengthy documents allows users to gain a swift understanding of essential concepts, facilitating informed decisions.
- **Enhanced Efficiency:** Text summarization automates the process of summarizing information, saving users valuable time and effort.

### **Approach:**

- *Data Collection and Preparation:*  
We collected news articles for our summarization dataset and cleaned them up to make sure the model could understand them better.
- *Using seq2seq LSTM and T5 Model:*  
We're trying out two models, one called seq2seq LSTM, which helps the model understand sequences of words, and another called T5, which is good at summarizing text by understanding long bits of information.
- *Fine-tuning and Making it Better:*  
We're not just stopping with the default T5 model. We're making it better for our task by training it more on our dataset and tweaking it to understand what makes a good summary.
- *Optimization for Readability and Accuracy:*  
We fine-tuned the pre-trained T5 model on our specific summarization dataset to further improve its performance and tailor it to our needs. We also incorporated additional loss functions during training to optimize generated summaries.

### **Corpus:**

The corpus used for our project consisted primarily of the news dataset, which contains news articles paired with corresponding summaries. We also explored using TensorFlow's News Dataset and potentially incorporating additional datasets to improve the model's ability to handle various text types. This multi-dataset approach aimed to create a more robust and adaptable text summarization system.

### **Background:**

Text summarization techniques have come a long way, evolving from simply picking statistically important sentences to using machine learning and neural networks for a deeper understanding of text. This allows for summaries that capture the essence of the original content, but challenges remain in ensuring the summaries are grammatically sound. Research continues to tackle these issues and push the boundaries of what summarization can achieve.

### **Relevant Prior Work:**

1. *Early Approaches to Text Summarization: Luhn, H. P. (1958). The automatic creation of literature abstracts. IBM Journal of Research and Development, 2(2), 159-165.*  
Luhn's pioneering work introduced statistical methods for generating abstracts, laying the foundation for text summarization. However, these methods had limitations in capturing semantic nuances and producing summaries.
2. *Introduction of Neural Networks for Summarization: Rush, A. M., & Chopra, S. (2015). A neural attention model for abstractive sentence summarization. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP) (pp. 379-389).*

Rush and Chopra's neural attention model marked a significant shift towards advanced deep learning techniques. Leveraging attention mechanisms, this model improved performance by focusing on relevant parts of the input text.

3. *Seq2seq Models for Text Summarization: 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) (Vol. 1, pp. 1073-1083).* See et al. introduced pointer-generator networks, a seq2seq variation, addressing vocabulary mismatches and generating summaries resembling human-written ones. This represented a notable advancement in handling diverse vocabularies.
4. *Advancements in Abstractive Summarization: Paulus, R., Xiong, C., & Socher, R. (2017). A deep reinforced model for abstractive summarization. arXiv preprint arXiv:1705.04304.* Paulus et al. proposed a deep reinforced model, combining reinforcement learning with seq2seq architectures. This approach enables the model to learn directly from reward signals, resulting in more fluent summaries compared to traditional methods.

## Approach

### Overview:

Our project aims to automate the summarization of text documentation, enabling a quick understanding of content through concise summaries. The approach integrates data collection, preprocessing, training, and evaluation. Multi-News dataset is used from tensor flow, ensuring a diverse range of text for summarization. Preprocessing includes text cleaning, stop word removal, and contraction mapping to prepare the data for model input. T5 and LSTM models are then employed to generate concise summaries, capturing the essence of the original documents.

### Data Preprocessing:

To ensure the optimal performance of our text summarization model, a series of preprocessing steps were undertaken:

- *Text Lowercasing:* All text data was converted to lowercase to maintain uniformity and eliminate any potential inconsistencies caused by differing case styles.
- *Remove Punctuation:* Using the translate method punctuations are removed to maintain data consistency and eliminate extraneous symbols.
- *Replacing Special Characters:* Using a regular expression, the special characters are removed and replaced with spaces.
- *Removing Multiple Spaces:* General multiple spaces and spaces that replaced special characters are replaced by single spaces.
- *Contraction Mapping:* Expands contractions using a predefined dictionary mapping, to improve processing and ease the understanding of the text.
- *Remove Stop words:* Filters out common words that typically do not contribute much information. which helps to focus on the significant words for summarization tasks.

```
def text_strip(s):
    """
    This function is used to remove the punctuations, special characters, and stopwords from the text.
    """
    # Convert the text to lowercase
    s = s.lower()

    # Removing punctuations using the translate method
    s = s.translate(str.maketrans('', '', string.punctuation))

    # Replacing special characters with spaces using regular expressions
    s = re.sub(r"<(>|&@#%&quot;'\\";~*!]", ' ', str(s))

    # Replacing multiple spaces with a single space
    s = re.sub("\\.\\s+", ' ', str(s))
    s = re.sub("\\-\\s+", ' ', str(s))
    s = re.sub("\\:\\s+", ' ', str(s))

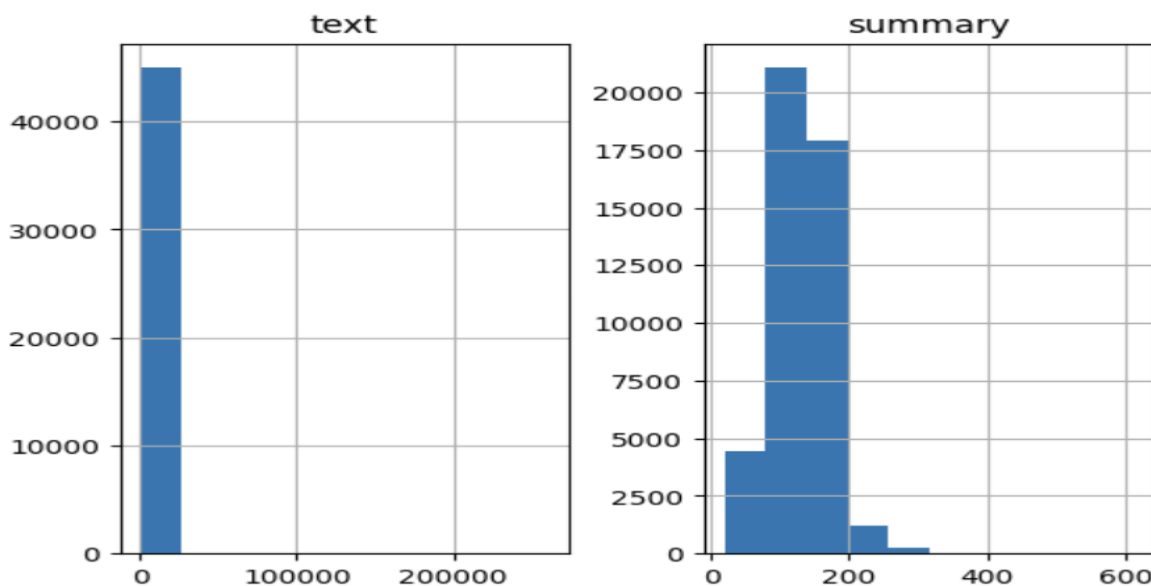
    # Expand contractions using a predefined mapping
    s = ' '.join([contraction_mapping[t] if t in contraction_mapping else t for t in s.split(" ")])

    # Removing punctuations again after expanding contractions
    s = s.translate(str.maketrans('', '', string.punctuation))

    # Removing stopwords
    stop_words = set(stopwords.words('english'))
    return ' '.join(word for word in s.split() if word not in stop_words)
```

## Document and summary length distribution:

Analyzing the length distribution of both documents and their corresponding summaries. Histograms were generated to visualize this distribution, revealing that most documents contain between 0 to 1024 words and summaries typically range from 0 to 128 words. This analysis is critical for determining the maximum length parameters for our text summarization model. By setting appropriate lengths, we can optimize our model to handle the common lengths effectively, thus improving the model performance and computational efficiency.



## Tokenization and vocabulary building:

Initially tokenized the corpus, transforming word sequences into integer sequences suitable for the model's input. To safeguard against any potential information leak from the test data only train data is tokenized. To optimize the model's performance and reduce the chance of overfitting. We employ padding strategies to standardize the length of input sequences, guided by the data distribution analysis visualized through previously plotted histograms. These steps in preparing the text ensure the model is fed data of uniform length, allowing for efficient computation. The vocabulary size is a deliberate choice, balancing the model's ability to discern nuances within the text and the practical considerations of data volume and computational demand.

## Seq2Seq using LSTM Model Architecture:

This model is designed using a sequences-to-sequence framework with LSTM layers. The model comprises two main components, an encoder, and a decoder, each with multiple LSTM layers that are used in handling long-range dependencies in text documentation:

1. **Encoder:** The encoder in our seq2seq architecture forms the foundation for processing input text. It begins with an input layer that takes seq of variable length, the subsequent embedding layer is tasked with converting these integer sequences into dense vectors of size 200. Three stacked LSTM layers follow, each building on the previous one to understand the text better which helps in retaining the important part of the sentences and forgetting less important details. These LSTMs are configured with dropout techniques with dropout rate and recurrent dropout of 0.4, To avoid the model getting overfitted by randomly omitting nodes during the training phase. The final LSTM layer outputs a contextualized state, concluding the encoding process and providing the initial context for the decoder to generate the summary. This encoder is crucial for capturing the complex parts of language within the input text.
2. **Decoder:** After the encoder captures the context of the input text, the decoder begins its process of generating a summary. It starts with its embedding layer, which is like the encoder's embedding layer but is specific to the summary data. The embedded summary data is then passed to another LSTM layer. This layer uses the final hidden states from the encoder as its initial state with a dropout rate =0.4 and a recurrent dropout rate of 0.2.

The output of this LSTM is then run through a dense layer, where each neuron in this layer is connected to every neuron in the previous layer. The dense layer has a SoftMax activation function, which calculates probability distribution over all possible summary words and decides which word should come next in the sequence.

Finally, the model is defined by specifying the inputs from the encoder and the outputs from the decoder, tying together the full sequence-to-sequence model.

This sequence of operations culminates in the formation of a summary that closely aligns with the semantic content of the input text, embodying the goal of our text summarization model.

```

latent_dim = 300 # Dimensionality of the latent space
embedding_dim = 200 # Dimensionality of the word embeddings

# Encoder
encoder_inputs = Input(shape=(max_text_len, ))

# Embedding layer for input sequences
enc_emb = Embedding(x_voc, embedding_dim, trainable=True)(encoder_inputs)

# First LSTM layer of the encoder
encoder_lstm1 = LSTM(latent_dim, return_sequences=True,
                     return_state=True, dropout=0.4,
                     recurrent_dropout=0.4)
(encoder_output1, state_h1, state_c1) = encoder_lstm1(enc_emb)

# Second LSTM layer of the encoder
encoder_lstm2 = LSTM(latent_dim, return_sequences=True,
                     return_state=True, dropout=0.4,
                     recurrent_dropout=0.4)
(encoder_output2, state_h2, state_c2) = encoder_lstm2(encoder_output1)

# Third LSTM layer of the encoder
encoder_lstm3 = LSTM(latent_dim, return_state=True,
                     return_sequences=True, dropout=0.4,
                     recurrent_dropout=0.4)
(encoder_outputs, state_h, state_c) = encoder_lstm3(encoder_output2)

# Setting up the decoder, using encoder_states as the initial state
decoder_inputs = Input(shape=(None, ))

# Embedding layer for decoder input sequences
dec_emb_layer = Embedding(y_voc, embedding_dim, trainable=True)
dec_emb = dec_emb_layer(decoder_inputs)

# Decoder LSTM
decoder_lstm = LSTM(latent_dim, return_sequences=True,
                    return_state=True, dropout=0.4,
                    recurrent_dropout=0.2)
(decoder_outputs, decoder_fwd_state, decoder_back_state) = \
    decoder_lstm(dec_emb, initial_state=[state_h, state_c])

# Dense layer to output probabilities over the target vocabulary
decoder_dense = TimeDistributed(Dense(y_voc, activation='softmax'))
decoder_outputs = decoder_dense(decoder_outputs)

# Defining the model
model = Model([encoder_inputs, decoder_inputs], decoder_outputs)
model.summary()

```

## T5 Model Architecture:

The foundation of our model's architecture is derived from AutoModelForSeq2SeqLM, utilizing the pre-trained t5-small model known for its efficiency, which is specifically tailored for sequence-to-sequence tasks such as document summarization. This approach capitalizes on a model pre-trained on a diverse corpus, providing a robust starting point for fine-tuning our summarization dataset.

DataCollatorForSeq2Seq helps streamline the input data, ensuring that batches are correctly preprocessed, which is essential for maintaining training loop stability and achieving optimal model performance.

In the subsequent setup, we configure the training parameters through Seq2SeqTrainingArguments. This configuration is designed, considering the optimal learning rate, batch size, and the total number of

training epochs that balance computational efficiency and accuracy. Using half-precision computation (fp16=True) helps in substantially reducing the memory consumption and expedites the training process on compatible GPUs. The fine-tuning of these hyperparameters is critical to the training's success, as it influences both the training duration and the quality of the resulting summarization model.

The culmination of this setup is the initialization of Seq2SeqTrainer, which integrates the model, tokenizer, data collator, and training configuration into a unified framework for training. Leveraging the Hugging Face Transformers library, we invoke `trainer.train()` to begin the fine-tuning process. During this phase, the model iteratively learns from the training data, adjusting its internal parameters to minimize the summarization task's loss function. This fine-tuning stage enhances the model's ability to generate precise summaries and ensures that the nuances and specificities of the domain are captured, which is important for producing summaries that are both informative and reflective of the source material's intent.

```
# Loading Model
# Data collator for sequence-to-sequence (seq2seq) tasks, specifically designed for T5 model
data_collator = DataCollatorForSeq2Seq(tokenizer=tokenizer, model='t5-small')

# Loading the pre-trained T5-small model for sequence-to-sequence language modeling
# from the Hugging Face model hub
model = AutoModelForSeq2SeqLM.from_pretrained("t5-small")
```

```
# Hyperparameters for T5
training_args = Seq2SeqTrainingArguments(
    output_dir="./results",
    evaluation_strategy="epoch",
    learning_rate=2e-5,
    per_device_train_batch_size=10,
    per_device_eval_batch_size=10,
    weight_decay=0.01,
    save_total_limit=3,
    num_train_epochs=10,
    fp16=True,
)
```

lizing Trainer

```
# Initializing the Seq2SeqTrainer for training and evaluation

trainer = Seq2SeqTrainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_multi_news["train"],
    eval_dataset=tokenized_multi_news["test"],
    tokenizer=tokenizer,
    data_collator=data_collator,
)
```

## Results

### Corpus:

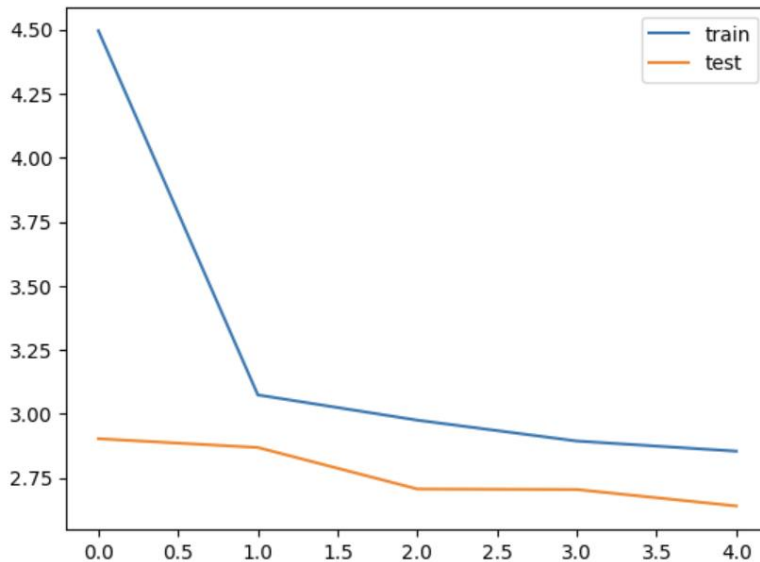
The dataset used in this project is the multi-news dataset which consists of news articles and their corresponding summaries. We leverage this paired structure (document and summary) for training the text summarization model. The dataset contains a substantial size for training, containing approximately 16,357 training examples and separate validation and test sets.

### Experiments and Performance Evaluation:

In our pursuit to explore various neural network architectures for the text summarization task, we conducted experiments with a sequence-to-sequence (Seq2Seq) model leveraging Long Short-Term Memory (LSTM) networks and a T5 (Text-to-Text Transfer Transformer) model.

**LSTM model:** The Seq2Seq model is composed of LSTM layers both in the encoder and the decoder. These LSTMs employ a sophisticated architecture with memory cells, which enables them to capture and utilize long-range dependencies in text data. This setup is essential for capturing the temporal dependencies within the text data. Due to computational constraints, we proceeded with a subset of the complete dataset. Training the model was a calculated endeavor, with epochs limited by computational power. Each epoch was a complete pass over the sample dataset, with the model learning to map input sequences to their summaries. Although the training was constrained, it was invaluable for obtaining insights into the LSTM's performance dynamics and setting a baseline for comparison with more complex models such as the Transformer.

### LSTM Loss Function Metrics:



The LSTM graph shows a loss function that starts at a higher value and displays a steep decline as training progresses. This shows that the LSTM model is learning and improving its performance on the summarization task as it goes through more epochs. The loss for the training set is declining at a steady rate, which suggests that the model fits well to the training data. However, the test validation loss starts to plateau early, which could imply that further training might not lead to significant improvements on unseen data, or it could indicate a need for additional training data or epochs to achieve further gains.



**T5 model:** Our experiments with the T5 model were conducted under controlled computational settings, focusing on fine-tuning the model on our chosen dataset. Due to the model's pre-trained nature, we hypothesized that even a limited number of training epochs would yield a significant understanding of the summarization context. This approach allowed us to conduct iterative training and evaluation cycles, progressively refining the model's summarization capabilities without necessitating extensive computational resources.

### **Transformer Loss Function Metrics:**



The Transformer model graph shows a more dramatic decrease in the training loss, showing that the model is rapidly learning from the data. The validation loss also decreases steadily, but the convergence with the training loss is much tighter than with the LSTM model. This could indicate that the Transformer model is generalizing better and is more capable of applying what it has learned from the training data to new, unseen data. The gap between training and validation loss is smaller, which often signifies a good fit and less overfitting.

### **Rouge Values:**

The evaluation exposes a clear advantage for the Transformer model. While the LSTM's ROUGE score of 0.02916 indicates minimal overlap with reference summaries, the Transformer's score of 0.1687 suggests a significantly better capture of essential information. This aligns with the Transformer's strength in handling long-range dependencies and context, leading to more comprehensive summaries

```
# Seq2seq LSTM - rouge 1
df_results['rouge1'].mean()
```

0.029162156663601575

```
# T5 - Transformer - rouge 1
val_data['rouge1'].mean()
```

0.16871276145837286

Both models have room for improvement, but the Transformer's superior performance here reflects current research trends favoring Transformers for text summarization tasks.

Finally, using the transformer model, we were able to predict the summary of the news articles. The table below displays the output of the text summarization model across different documents, showcasing each document, its predicted summary, and the corresponding ROUGE score. The predicted summary is displayed in the table. The ROUGE score evaluates the quality of summaries by measuring their similarity to a set of reference summaries, with higher scores indicating greater similarity.

	document	summary	pred_summary	rouge1	
0	Whether a sign of a good read; or a comment on...	– The Da Vinci Code has sold so many copies—th...	<pad><extra_id_0>,<extra_id_1>,<extra_id_2>,<e...	0.207921	
1	The deaths of three American soldiers in Afgha...	– A major snafu has hit benefit payments to st...	<pad><extra_id_0> is a national disgrace, one ...	0.179372	
2	DUBAI Al Qaeda in Yemen has claimed responsibi...	– Yemen-based al-Qaeda in the Arabian Peninsul...	<pad><extra_id_0> has claimed responsibility f...	0.161491	
3	Cambridge Analytica, a data firm that worked f...	– Cambridge Analytica is calling it quits. The...	<pad><extra_id_0>,<extra_id_1>,<extra_id_2>, a...	0.169231	
4	The N.S.A.'s Evolution: The National Security ...	– A lengthy report in the New York Times, base...	<pad><extra_id_0>,<extra_id_1>,<extra_id_2>,<e...	0.140704	
5	× \n \n remaining of \n \n Thank you for Readi...	– Don Juan de Oñate sought a city of gold when...	<pad><extra_id_0> the name Etzanao, a long-los...	0.147287	
6	Rep. Anthony Weiner rejected Democratic leader...	– Another bad day for Anthony Weiner: Nancy Pe...	<pad><extra_id_0>,<extra_id_1>,<extra_id_2>,<e...	0.187845	
7	Augmented reality startup Magic Leap is being ...	– An augmented reality startup is being sued f...	<pad><extra_id_0> the<extra_id_1>,<extra_id_2>...	0.136126	
8	Research suggests that the ratio of the length...	– The length of a man's index and ring fingers...	<pad><extra_id_0>, 100 matched healthy males, ...	0.193548	
9	HOUSTON – Almost two years ago, a pair of Unit...	– A 71-year-old lawyer is suing United Airline...	<pad><extra_id_0> a 71-year-old man is suing U...	0.252336	

## Conclusion

The project aimed to develop an efficient text summarization system using advanced neural network models, specifically the T5 and based sequence-to-sequence models. Our system, trained on a substantial corpus of news articles paired with summaries, demonstrated its capacity to effectively distill complex documents into concise summaries. The T5 model, leveraging its pre-trained capabilities and fine-tuned on a targeted dataset, showcased superior performance in learning and generalizing from the data, as reflected in the closely aligned training and validation loss curves. The LSTM model, while also effective, displayed potential limitations in scalability and generalization compared to the T5 model.

The key takeaways from this project are the superior efficacy of the T5 model, whose architecture and pre-training on diverse datasets effectively managed the complexities of text summarization. In contrast, the LSTM model, despite its ability to capture long-range dependencies, exhibited signs of possible overfitting and a plateau in validation loss, indicating a more limited adaptability. Key to the project's success was the iterative training and evaluation process. Effective preprocessing, such as tokenization and the removal of extraneous textual elements was crucial, showcasing the importance of clean and structured input data for optimal NLP performance.

## References

- Alla, S. Introduction to Encoder-Decoder Sequence-to-Sequence Models (Seq2Seq). Paperspace Blog. Online Available: <https://blog.paperspace.com/introduction-to-seq2seq-models/>
- B. N, D. Kumari, B. N, M. N, S. K. P and S. R. A, "Text Summarization using NLP Technique," 2022 International Conference on Distributed Computing, VLSI, Electrical Circuits and Robotics (DISCOVER), Shivamogga, India, 2022, pp. 30-35, doi: 10.1109/DISCOVER55800.2022.9974823. keywords: {Training;Computational modeling;Machine learning;Manuals;Very large scale integration;Transformers;Data mining;Summarization;Natural Language Processing;Transformers;Deep Learning},
- J. Ranganathan and G. Abuka, "Text Summarization using Transformer Model," 2022 Ninth International Conference on Social Networks Analysis, Management and Security (SNAMS), Milan, Italy, 2022, pp. 1-5, doi: 10.1109/SNAMS58071.2022.10062698. keywords: {Drugs;Measurement;Analytical models;Social networking (online);Transformers;Natural language processing;Security;Deep Learning;Natural Language Processing;Text Summarization;T5 Model},
- Alla, S. Implementing Seq2Seq Models for Text Summarization with Keras. Paperspace Blog. Online Available: <https://blog.paperspace.com/implement-seq2seq-for-text-summarization-keras/>

## CONTRIBUTIONS

MRT, NK, HRV contributed equally in the whole project.