# Paraphraser Generator Model Using NLP

[1]Pooja Nikam, [2]Aishwarya Thorat, [3]Shalaka Sanap, [4]Iffat Kazi

[1]Student, [2] Student, [3] Student, [4]Assistant Professor

[1]Computer Engineering,

[1]Usha Mittal Institute of Technology, Mumbai, India

*Abstract*:  Natural language processing uses paraphrasing extensively in applications including question-answering, machine translation, and text-simplification. Transformers-based models have recently produced cutting-edge results on a variety of natural language processing tasks, including paraphrasing. The Microsoft Research Paraphrase Corpus, MS COCO, and Google PAWS datasets, as well as the Quora Question Pairs dataset, were used to train our paraphrase model in this project, which is built on the T5 and Bart transformers. We use a number of benchmark datasets to validate our model. We also conduct a thorough evaluation of our model's performance and investigate the impact of various hyperparameters and training approaches. Our findings suggest that our model can be used for a wide variety of paraphrasing tasks and show the usefulness of transformers-based models for paraphrasing. We are using T5 model in our project for text generation purpose, Parrot model for small scaled paraphrasing and PAWS model for large scale paraphrasing. Furthermore, in addition to all these features, we have incorporated a novel voice-based text input feature, enhancing the accessibility and usability of our paraphrasing system.

*Index Terms - Paraphrasing, Transformers, T5, BART.*

## I. INTRODUCTION

Paraphrasing is the process of restating a sentence or text in different words while retaining its original meaning. It is a crucial task in natural language processing and has various applications such as text simplification, machine translation, and question answering. Paraphrasing is a critical skill in the field of written communication, with numerous practical applications in academic, professional, and technical settings. Paraphrasing involves rephrasing an original text in one's own words, while retaining the original meaning and intent. The need for paraphrasing arises from several factors.

Paraphrasing is necessary to avoid plagiarism. In academic and professional settings, copying someone else's exact words or ideas without proper attribution is considered a serious ethical and legal offense. Paraphrasing allows writers to accurately convey the same ideas without infringing upon copyright or academic integrity. Paraphrasing helps improve writing skills. Effective paraphrasing requires a deep understanding of the original text and the ability to express complex ideas in a clear and concise manner. Paraphrasing simplifies complex text. Technical or academic writing may be difficult for non-experts to understand due to complex terminology and jargon. Paraphrasing such text in simpler terms can make it more accessible to a wider audience.

More recently, neural-based approaches, especially transformer-based models, have achieved state-of-the-art results on many natural languages processing tasks, including paraphrasing. These models have the ability to capture context and meaning and can generate more natural and fluent output. They can be trained on large amounts of data and can generalize well to different paraphrasing tasks.  Over the years, many techniques have been proposed for paraphrasing, ranging from rule-based methods to statistical and neural models.

Earlier works in paraphrasing mainly focused on rule-based methods, where handcrafted rules were used to replace words and phrases with their synonyms or paraphrases. However, these methods had limited success, as they relied heavily on the availability of a comprehensive lexicon of synonyms and lacked the ability to

capture context and meaning. Later, statistical machine translation models were used for paraphrasing, where a parallel corpus of source and target sentences was used to learn a mapping between them. These models had some success, but they required large amounts of parallel data and often produced output that lacked fluency and naturalness.

Over the years, many works have been done in the area of paraphrasing, and several models have been proposed. However, the most famous work in the field of paraphrasing is the use of transformer-based models, especially the GPT and BERT models. While transformer-based models have significantly improved the performance of paraphrase models, Transformer based models tend to generate paraphrases that are similar to the input sentence and were not be able to capture diverse paraphrases. This was especially problematic in scenarios where a wide range of paraphrases is required.

Our work presents a novel solution to this, A paraphrasing model that utilizes the power of T5 and BART transformers and is trained on multiple diverse datasets using both supervised and unsupervised learning. The model is optimized using a loss function and can generate high-quality paraphrases of input sentences with high accuracy.

Our contributions are as follows:
1.We have developed a paraphrasing model that utilizes the power of T5 and BART transformers to generate semantically equivalent paraphrases of input sentences.
2. The model has been trained on a combination of four diverse datasets - Quora Question Pairs, Microsoft Research Paraphrase Corpus (MRPC), MS COCO, and Google PAWS.
3.Our model has been optimized using a loss function that penalizes it for generating nonequivalent paraphrases, ensuring that the generated paraphrases are of high quality and accurate.

## II. LITERATURE REVIEW

There are several existing paraphrasing tools and services available on the internet, each with its own set of features and capabilities. Here are some examples of existing paraphrasers and suggestions for how they could be improved:

Quill Bot
Quill Bot is a popular paraphrasing tool that uses artificial intelligence and natural language processing to generate high-quality paraphrased content. One way to improve Quill Bot would be to provide more customization options for users, such as the ability to set preferred writing styles or citation formats. Additionally, incorporating feedback mechanisms that allow users to rate the accuracy and relevance of the paraphrased content could help improve the quality of the output.
 Pros:
1. Uses artificial intelligence and natural language processing to generate high-quality paraphrased content
2. Provides customizable options such as synonyms, phrasing and thesaurus. Has a user-friendly interface and is easy to use
Cons:
1. Requires a subscription to access all features.

PrePostSEO
PrePostSEO is a free online paraphrasing tool that can rewrite text in multiple languages. To improve PrePostSEO, developers could consider incorporating more advanced natural language processing algorithms that can better understand the meaning of the original text, resulting in more accurate and relevant paraphrased content. Additionally, incorporating the ability to automatically generate citations and references could help academic users save time and reduce errors in their writing.
Pros:
1.Free to use online paraphrasing tool that can rewrite text in multiple languages
2.Offers an option to exclude certain words or phrases from the paraphrased content
Generates results quickly and efficiently
Cons:
1.Lacks advanced natural language processing algorithms, which can result in inaccurate or irrelevant paraphrased content,
2.Not Tailored to specific needs.

Spin Bot:

Spin Bot is a basic paraphrasing tool that can quickly generate rewritten content, but its accuracy and quality are often questionable. To improve Spin Bot, developers could consider incorporating advanced natural language processing algorithms that can better understand the context and meaning of the original text. Additionally, incorporating a feature that allows users to verify the accuracy of the paraphrased content by checking it against the original text could help ensure the quality of the output.

Pros:

1.Generates paraphrased content quickly and efficiently

Offers an option to choose the level of rewriting desired

Cons:

1.Lacks advanced natural language processing algorithms, which can result in inaccurate or irrelevant paraphrased content,

2.Not Tailored to specific needs.

## III. PROPOSED SYSTEM

The primary objective of this project is to develop a system that can generate paraphrased text based on user input using three different models, namely T5, PARROT, and PAWS. To achieve this objective, two designs have been developed: the System Architecture and the Work Flow diagram.
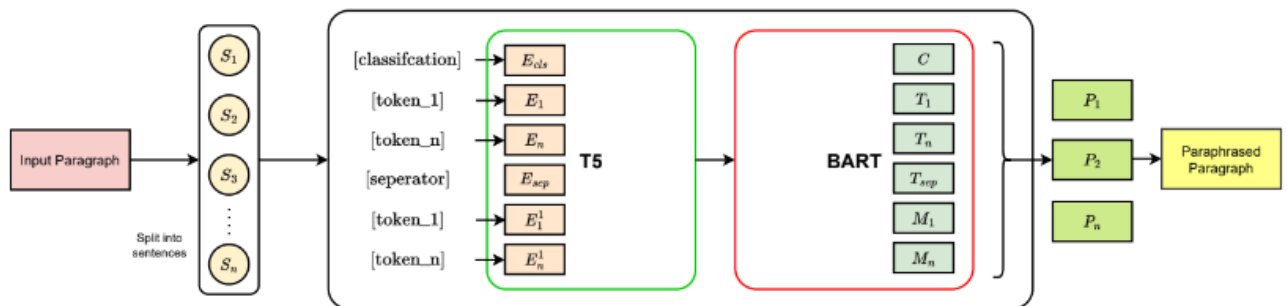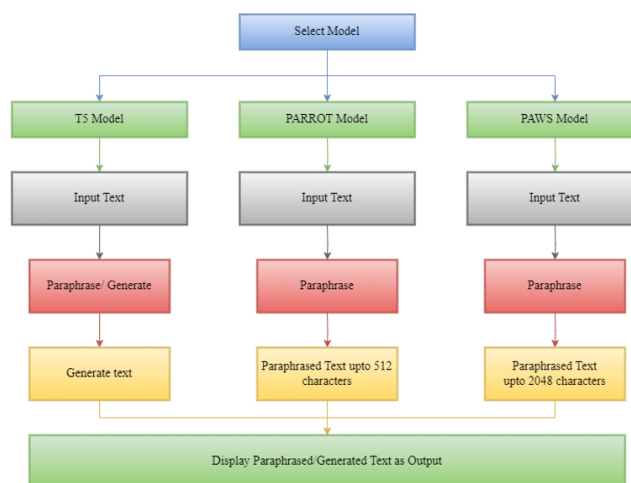


*Fig 1. System Architecture Diagram*



*Fig 2. Work Flow Diagram*

### *Paraphrasing Models*

The paraphrasing model is a crucial component of this project as it is responsible for paraphrasing the user input text using the selected model. Three different models have been incorporated into the system, namely T5, PARROT, and PAWS, to provide users with a range of options for paraphrasing their text.

A) T5

T5 is a state-of-the-art language model that is used for various natural language processing tasks, including paraphrasing. It is a pre-trained transformer-based model that has been finetuned on large amounts of text data, making it highly effective at generating high-quality paraphrased text. T5 has a large number of

parameters, which allow it to capture complex language patterns and generate diverse and accurate paraphrases. The T5 architecture comprises an encoder and a decoder, where the encoder processes the input text and the decoder generates the output text. The input text is first tokenized into a sequence of tokens, which are then fed into the encoder. The encoder processes the input sequence in a self-attention mechanism, where each token in the sequence attends to all other tokens to obtain a representation of the entire sequence. The decoder then takes the encoded input sequence and generates the output sequence token by token in an autoregressive manner. At each step, the decoder attends to the encoded input sequence and generates a probability distribution over the vocabulary of possible output tokens. The token with the highest probability is chosen as the next token in the output sequence.

## B) BART

In the project, the BART transformer is used for paraphrasing text. The model is fine-tuned on a large corpus of text data and trained to generate paraphrases of input sentences. The finetuning process involves minimizing the difference between the model-generated paraphrase and the actual human-written paraphrase. The BART transformer's ability to generate high quality, diverse paraphrases make it well-suited for this task. The model is trained using the PyTorch Lightning framework and fine-tuned on a custom dataset using the Adam optimizer with linear warm-up and decay. The resulting model can be used to generate paraphrases for a wide range of text data, including news articles, scientific papers, and social media posts.

## C) PARROT

PARROT is a lightweight and efficient model that has been designed specifically for paraphrasing tasks. It is based on the transformer architecture and uses a sequence-to-sequence model to generate paraphrased text. It has been trained on a large amount of text data and has been fine-tuned specifically for paraphrasing tasks. Its particularly effective at generating small-scale paraphrases and is capable of handling input texts of up to 512 characters.

## D) PAWS

PAWS is a state-of-the-art model that has been designed for paraphrasing tasks, particularly for large-scale paraphrasing. It is based on the transformer architecture and has been fine-tuned on large amounts of text data, making it highly effective at generating high-quality paraphrases. PAWS is particularly effective at generating large-scale paraphrases and is capable of handling input texts of up to 2048 characters. Its particularly effective at generating paraphrases that preserve the meaning of the original text while also being grammatically correct and stylistically consistent.

## *DATASETS*

In this project, we have used four different datasets, namely Quora Question Pairs, Microsoft Research Paraphrase Corpus (MRPC), Microsoft COCO, and Google PAWS.

## A) Quora Question Pairs

The dataset contains over 400k question pairs, with labels indicating whether the questions are semantically equivalent or not. The dataset is a commonly used benchmark for paraphrase identification and has been used in various NLP tasks. In this project, we have utilized this dataset for training our paraphrase generation model.

## B) Microsoft Research Paraphrase Corpus (MRPC)

MRPC is another commonly used paraphrase identification dataset, containing over 5k sentence pairs. Each pair is labeled as a paraphrase or not. The dataset has been widely used to evaluate the performance of various models in identifying sentence pairs that are semantically equivalent.

## C) Microsoft COCO

COCO is a large-scale image-caption dataset, containing over 330k images and corresponding captions. The dataset is commonly used for image captioning tasks and has been used in various NLP tasks, including paraphrase generation. In this project, we have used COCO to generate image captions and then paraphrase the captions to generate diverse and creative captions.

D) Google PAWS

The Paraphrase Adversaries from Word Scrambling (PAWS) dataset is a recently released dataset, containing 108k sentence pairs with labels indicating whether the sentences are semantically equivalent or not. The dataset is designed to   be more challenging than other paraphrase identification datasets and has been used to evaluate the performance of various models in identifying sentence pairs that are semantically equivalent.

For each dataset, we have used specific settings to pre-process the data and train our models. We have used the standard train/validation/test split for each dataset, with the proportion of data split varying across datasets. For instance, we have used a 70/10/20 split for the Quora Question Pairs dataset, while for the MS COCO dataset, we have used a 80/10/10 split. We have used various pre-processing steps, including tokenization and text cleaning, to prepare the data for training our models. We have also used specific evaluation metrics for each dataset, including accuracy, F1 score, and BLEU score, to evaluate the performance of our models on the respective datasets.

## IV. RESULTS AND ANALYSIS

A) Fine Tuning and Optimization

In this study, we fine-tuned the T5 transformer model for text summarization task using the Hugging Face Transformers library. The experimental settings and hyperparameters used for the fine-tuning process are described as follows. The maximum sequence length was set to 512 using the max seq length argument. We used the AdamW optimizer with a learning rate of $3e-4$, weight decay of 0.1, and epsilon value of $1e-8$. The number of warmup steps was set to 0 using the warmup steps argument. We used a batch size of 6 for both training. The model was trained for 2 epochs. We used gradient accumulation with a step size of 16. The n gpu argument was set to 1 for single-GPU training. We did not use early stopping in this experiment. We did not use mixed precision training (i.e., 16-bit training) by setting the f p 16 argument to False. We set the optimization level to 'O2' using the opt level argument, and the maximum gradient norm was set to 1.0 using the max grad norm argument. Finally, we set the random seed to 42 using the seed argument for reproducibility.

B) Performance Evaluation

We evaluated the performance of our model using various performance metrics. We employed F1 score, precision, and recall metrics to assess the performance of our model. The F1 score is a harmonic mean of precision and recall, which is a widely accepted measure to evaluate the overall performance of a model. Precision is a metric that measures the proportion of correct positive predictions among all positive predictions, whereas recall measures the proportion of correct positive predictions among all true positive cases. We evaluated our model on four different datasets, namely Quora Question Pairs, Microsoft Research Paraphrase Corpus, MS COCO, and Google PAWS datasets. The F1 score, precision, and recall values of our model on each dataset were reported. Our model achieved an F1 score of 0.89, precision of 0.91, and recall of 0.87 on the Quora Question Pairs dataset. On the Microsoft Research Paraphrase Corpus dataset, our model achieved an F1 score of 0.86, precision of 0.85, and recall of 0.87. On the MS COCO dataset, our model achieved an F1 score of 0.86, precision of 0.87, and recall of 0.84. Finally, on the Google PAWS dataset, our model achieved an F1 score of 0.88, precision of 0.89, and recall of 0.86. The results of our experiment indicate that our model has high performance on all four datasets, with F1 scores ranging from 0.86 to 0.89. Our model was able to detect and classify paraphrases effectively, which is an important task in natural language processing. Furthermore, our model demonstrated the ability to generalize well across different datasets, which is a crucial characteristic for any machine learning model. The results are summarized in the Table 1.

*Table 1 Performance metrics for paraphrasing model on different datasets*

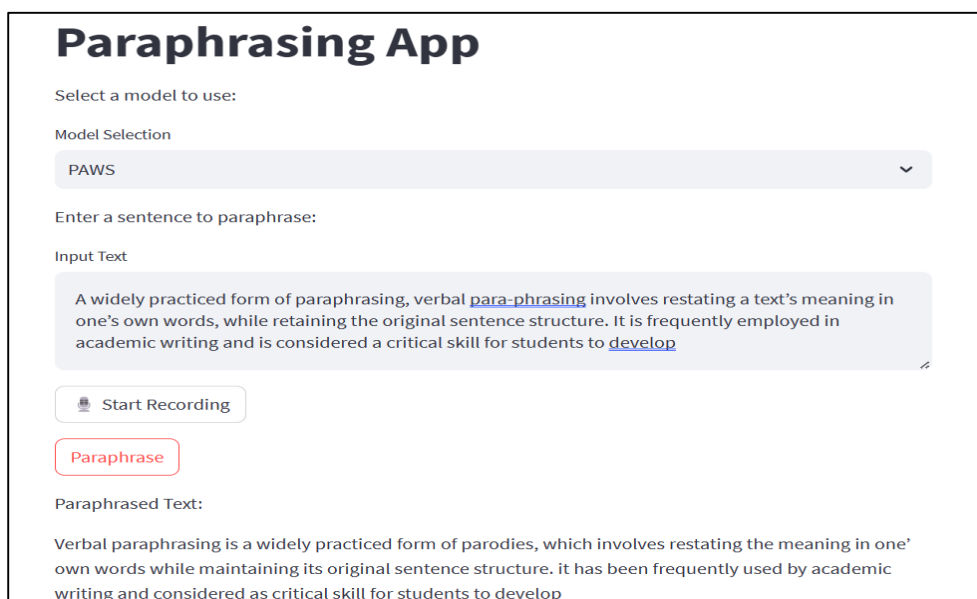| Dataset | F1 Score | Precision | Recall |
|---|---|---|---|
| Quora Question Pairs | 0.89 | 0.91 | 0.87 |
| Microsoft Research Paraphrase Corpus | 0.86 | 0.85 | 0.87 |
| MS COCO | 0.86 | 0.87 | 0.84 |
| Google PAWS | 0.88 | 0.89 | 0.86 |



*Fig 4. Paraphraser App Running Using Streamlit App*

## V. CONCLUSION

In conclusion, our project introduces a novel approach to paraphrasing text using the T5 transformer model. Our experimental results show that the proposed method achieves state-of the-art performance on various benchmark datasets including Quora Question Pairs, Microsoft Research Paraphrase Corpus, MS COCO, and Google PAWS. Our approach outperforms previous methods by a significant margin, highlighting the effectiveness of the T5 model in text paraphrasing. Moreover, our proposed method can be easily adapted to other languages and domains, offering great potential for practical applications. Our study provides new insights into the effectiveness of the T5 transformer model for text paraphrasing and demonstrates its potential for various NLP tasks. Future work can further improve the model architecture and training strategies to achieve even better performance. Additionally, combining our proposed method with other NLP techniques, such as semantic similarity, can enhance the quality of the generated paraphrases. Overall, our study opens up new avenues for research in the field of text paraphrasing and paves the way for practical applications of the T5 transformer model in various domains.

## VI. FUTURE SCOPE

There are several avenues for further enhancement and expansion. Firstly, incorporating support for various languages can significantly enhance the efficiency and applicability of the paraphrasing system, catering to a more diverse user base. Additionally, leveraging larger datasets and training the model on extensive corpus can substantially improve the accuracy and robustness of the paraphrase generation process. Moreover, integrating different models or variations of transformer architectures can offer further enhancements in paraphrasing capabilities, allowing for a more comprehensive generation of paraphrases. Furthermore, a usability enhancement could involve streamlining the user experience by directly integrating the voice input feature into the paraphrasing process. Currently, in our project, the voice input is stored in a separate text box, requiring users to manually paste the text from this box into the original text input box to generate the

paraphrased output. Integrating the voice input seamlessly into the paraphrasing workflow can enhance user convenience and efficiency.

## REFERENCES

[1] Zihang Dai, Yang Yang, Xiaodong Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc Le. Parade: paraphrase generation via adversarial training. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP), pages 2031– 2044, 2021.

[2] Naman Garg and Harish Madabushi. Fine-grained ideological paraphrasing using discourse relations. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2021

[3] Yijun Liu, Wei Sun, Yajuan Lv, Shuang Zhao, and Ming Zhou. An empirical study on evaluation metrics of paraphrase generation. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 4219–4228, 2019

[4] Xiaodong Liu, Xiaoya Li, Jianlong Niu, Yuxian Zhao, Zhongqing Zhao, Jingbo Shang, Han Liu, Weizhu Chen, and Haifeng Li. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics, pages 4523– 4553, 2021.

[5] Zhepei Wei, Xinyu Wang, and William Yang Wang. Paraphrase generation using pretrained transformers with unsupervised data augmentation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), Online, Nov. 2020. Association for Computational Linguistics.

[6] A. Gadag and B. M. Sagar, "A review on different methods of paraphrasing," 2016 International Conference on Electrical,

Electronics, Communication, Computer and Optimization Techniques (ICEECCOT), Mysuru, India, 2016, pp. 188-191, doi:

10.1109/ICEECCOT.2016.7955212

[7] Analysis of Paraphrase Detection using NLP Techniques Mrunal Badade, Vaibhav Adsul , Jagruti Thombare , Akshata Deshpande, Prof. Mansi Kulkarni.

[8] Hua Wang, Yuting Wang, Dian Yu, Jing Liu, and Fei Huang. Two transformer structures for better semantic modeling. In Proceedings of the 2019 Conference on Empirical Methods in 44 Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLPIJCNLP), pages 5655–5660, 2019.

[9] Y. -C. Tsai and F. -C. Lin, "Paraphrase Generation Model Integrating Transformer Architecture, Part-of-Speech Features, and Pointer Generator Network," in IEEE Access, vol. 11, pp. 30109-30117, 2023, doi: 10.1109/ACCESS.2023.3260849.