

# **Abstractive Text Summarization Using Transformer Models: A Study with BART**

Lathigaa , Roll.no; 33 ,12223482 , K22KE

## **Abstract**

**Text summarization in NLP is basically a task aimed at generating shorter versions of text without losing the core ideas. Traditional methods often rely on selecting important sentences, which can result in summaries that lack a natural flow. The development of transformer models, such as BART (Bidirectional and Auto-Regressive Transformers), enables the generation of more natural summaries. In this paper, we implement BART for abstractive summarization and fine-tune it on the CNN/Daily Mail dataset. Using ROUGE scores for evaluation, our results show that BART produces clearer and more informative summaries compared to traditional methods. This study highlights the effectiveness of transformer models in text summarization and suggests possible improvements for future work.**

## **Introduction**

Rapid growth in digital information forms an obstacle for both individuals and organizations to process and understand huge volumes of text effectively. The solution to this lies in making use of machine learning technologies to produce short versions of long documents, saving key elements without wasting time over extensive reading. Applications of summarization span various domains, including news aggregation, legal document review, academic research, and customer feedback analysis, where users benefit from quick and accurate summaries. Traditional classification of the summarization methods is made into two categories: extractive and abstractive.

Extractive summarization directly selects sentences or phrases from the source text to form a summary, aiming at emphasizing the most relevant part without any change in the original wording. While extractive methods are straightforward and often yield grammatically correct summaries, they lack flexibility in sentence structure and can result in disjointed or incoherent summaries. The abstractive summarization approach tends to rephrase and reorganize the content of the original in order to produce summaries. This requires a deeper sense of language and

contextual understanding because the model must make new sentences that capture the essence of the original text. Traditionally, abstractive summarization is harder to perform because of the challenging properties of natural language generation. The recent success of deep learning, especially transformer models, has made it more feasible.

### **1.1 Challenges in Abstractive Summarization**

Abstractive summarization presents unique challenges:

- I. **Content Relevance:** The model must accurately identify and include only the most critical information from the original text.
- II. **Coherence and Fluency:** Generated summaries should read naturally and maintain logical flow, requiring an understanding of sentence structure and syntax.
- III. **Factual Accuracy :** The model should not produce information that is absent in the original text-a popular problem with transformer-based language models.

Traditional neural network-based approaches, including sequence-to-sequence RNNs, have also been applied for abstractive summarization with some success, but they have a hard time handling long dependencies and very easily

get trapped into producing repetitive or incomplete summaries. In 2017, Vaswani et al. proposed the transformer model, the relation of which to the typical paradigm of self-attention mechanisms could capture long-range dependencies and generate coherent text.

### **1.2 Motivation for Using Transformer Models**

Following the introduction of transformer architectures such as BERT, T5, and BART has revolutionized the NLP space with high performance in context-aware language models that can effectively handle summarization tasks. Unlike previous approaches, transformers use self-attention layers to focus on different parts of the text, allowing them to understand local and global context. This ability is essential for abstractive summarization, where models must distill information from across the entire text. In this study, we focus on the BART model, a transformer-based encoder-decoder architecture specifically designed for text generation tasks, including summarization.

BART is pre-trained with a denoising objective, where parts of the input text are corrupted and the model learns to reconstruct the original text. This approach allows BART to deal with noisy, incomplete, or complex inputs and produce fluent, coherent outputs. Particularly, the architecture of BART, combining bidirectional

encoding and autoregressive decoding, is rather suitable for abstractive summarization since it can understand fully all the context of the input and produce readable output.

## **Related Work**

The task of text summarization has been widely studied, and research has evolved significantly over the years. Early approaches focused on extractive methods, which simply select and rank important sentences from the original text, while recent methods leverage deep learning, specifically transformer models, for generating more coherent and informative summaries. This section provides an overview of both traditional extractive methods and recent advancements in abstractive summarization using neural networks and transformers.

### **2.1 Extractive Summarization Methods**

Traditional extractive summarization techniques rely on ranking sentences or phrases based on various heuristics, such as word frequency, sentence position, or similarity measures. Common extractive methods include:

1. **Frequency-Based Methods:** These methods rank sentences based on the frequency of key terms assuming words that appear more often have a higher importance. The

widely used popular Term Frequency-Inverse Document Frequency (TF-IDF) approach has been widely adapted in this context[1].

2. **Graph-Based Models:** LexRank and TextRank are graph-based algorithms, where sentences are represented as nodes and similarity between sentences is denoted using edges in a graph. These models leverage the PageRank algorithm to identify the most "central" sentences for the summary. While effective, these methods often produce summaries that lack flow and coherence[2].
3. **Machine Learning Approaches:** Some early machine learning models employed classification and regression techniques to rank sentences based on hand-crafted features, such as sentence length, presence of certain keywords, and sentence position in the text [3]. These approaches provided some degree of flexibility but were limited by their reliance on manually defined features.

Extractive approaches are relatively straightforward to implement and can produce summaries that accurately reflect portions of the original text. However, because they merely select text without rephrasing

it, they often lack a natural flow and coherence, particularly in cases where the important information is spread across multiple sentences.

## **2.2 Abstractive Summarization Methods Before Transformers**

The approach of abstractive summarization aims at generating new sentences that convey the original meaning in a reduced fashion. The traditional abstractive methods have long been challenging to implement effectively as they require models to understand and generate human-like language. Early approaches to abstractive summarization used statistical models, such as Hidden Markov Models (HMM) and Bayesian methods, to generate summaries by modeling text as a probabilistic sequence. However, these methods lacked the capacity to capture long dependencies, resulting in summaries that were often too simplistic or lacked accuracy [4].

In the area of deep learning, neural network-based methods have been applied on text summarization, especially using Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks.

Sequence-to-Sequence (Seq2Seq) models having an encoder-decoder architecture are the mainstream of abstractive summarization and gave promising results for longer texts but produced many repetitions or were not complete. Techniques like

attention mechanisms (Bahdanau et al., 2015) [5] improved performance by allowing the model to focus on relevant parts of the input, but challenges such as handling complex dependencies and ensuring factual accuracy persisted.

## **2.3 Transformer Models in Text Summarization**

The introduction of the transformer model by Vaswani et al. (2017) [6] marked a significant advancement in NLP, including text summarization. Transformers use a self-attention mechanism to capture contextual dependencies across long sequences, which overcomes many limitations of RNN-based models. The transformer architecture has been adapted for several NLP tasks, leading to various models specifically tailored for text summarization:

1. **BERT (Bidirectional Encoder Representations from Transformers):** BERT is pre-trained, bidirectional transformer model that comes from large-scale text corpora with masked language modeling. Although itself not particularly a model for generative tasks such as summarization, it has recently been adapted to be applied in extractive summarization by fine-tuning on sentence selection tasks. BERTSUM (Liu & Lapata, 2019) [7] is

another example of the extractive summarization model, which is based on BERT. In this model, BERT is used to encode sentences and then select those that best summarize the main ideas.

2. **T5 (Text-To-Text Transfer Transformer):** The T5 model (Raffel et al., 2020) [8] is a text-to-text transformer model trained on a wide range of NLP tasks by framing each as a text generation problem. T5 has shown strong performance on various abstractive summarization benchmarks, since it is pre-trained to understand and generate text sequences in a flexible manner.
3. **BART (Bidirectional and Auto-Regressive Transformers):** Bidirectional and Auto-Regressive Transformers BART (Lewis et al., 2020) [9] is an auto-regressive transformer that is designed specifically for the task of sequence-to-sequence including text summarization. This combines a bidirectional encoder (like BERT) with an autoregressive decoder (like GPT), so it knows the full context of the input and generates fluent and coherent outputs.

## 2.4 Recent Advances in Transformer-Based Summarization

After the success of BERT, T5, and BART, recent studies focus on building transformer-based summarization models. Strategies such as reinforcement learning and contrastive learning have been used to further improve the quality of generated summaries.

Reinforcement learning approaches, such as those used in models like Pegasus (Zhang et al., 2020) [10], employ rewards based on evaluation metrics like ROUGE to improve summary relevance and coherence. Contrastive learning techniques aim to prevent models from generating irrelevant or incorrect information by contrasting similar and dissimilar sentences during training.

## 2.5 Contribution of This Paper

Building on the advancements in transformer-based models, this paper focuses on the implementation and evaluation of BART for abstractive text summarization. While previous research has demonstrated the effectiveness of transformer models, this paper provides an empirical study of BART on a popular summarization dataset (CNN/Daily Mail) and analyzes the model's performance through ROUGE metrics. Our approach includes fine-tuning BART and comparing its performance against traditional extractive baselines, offering insights into the advantages of transformer-based models for summarization tasks.

# METHODOLOGY

The approach taken to implement and evaluate the BART model for abstractive text summarization. We present our dataset, pre-processing, model architecture, fine-tuning, and evaluation metrics.

## Dataset

For this work, we employ the CNN/Daily Mail dataset as a common benchmark for summarization. This dataset consists of news articles paired with corresponding summaries, providing a large amount of text data suitable for training and evaluating abstractive summarization models. Each article-summary pair is well-structured, with articles containing several sentences or paragraphs and summaries typically consisting of concise sentences capturing key information.

### 1.2 Data Pre-processing

Data pre-processing is an essential step to ensure that input text is standardized and appropriately formatted for training. The pre-processing steps include:

1. **Tokenization:** We tokenize the text with BART's inbuilt tokenizer, splitting the text into subwords to handle rare words and improve model vocabulary coverage.

2. **Lowercasing:** We normalize all text to lowercase to make sure all elements of our dataset have the same case.
3. **Special Tokens:** BART uses special tokens (e.g., <S>, </S>) to mark the beginning and end of text sequences. We include these tokens to help the model recognize the boundaries of each sequence.
4. **Length Truncation:** Both articles and summaries are truncated to a maximum length to fit within the model's input limit. For BART, the maximum length is typically set to 512 tokens for articles and 128 tokens for summaries. This step ensures that each example fits within the memory constraints of the training hardware.

### 1.3 Model Architecture

We employ BART model, which consists of a bidirectional encoder combined with an autoregressive decoder. The BART architecture is particularly for summarization because it has the ability to :

- **Encode Bidirectional Context:** BART's encoder processes input text in a bidirectional manner, capturing both left and right context simultaneously. This bidirectional encoding enables BART to understand the full

context of each sentence within the article.

- **Generate Text Autoregressively:** The decoder generates text one token at a time, producing coherent summaries by taking into account previously generated tokens. This structure allows the model to generate natural-sounding summaries.

#### 1.4 Fine-Tuning Strategy

Fine-tuning is gained through training the pre-trained BART model on the CNN/Daily Mail dataset. Fine-tuning allows the weights of the model to be adjusted so as to obtain the least difference between the predicted and actual summaries that would map the input text-article into the target text-summary. Key highlights in the fine-tuning process include:

**1. Objective Function:** We use the cross-entropy loss function for measuring error in the predicted and target summaries. Minimizing this loss will help the model generate text sequences that are closer to the ground-truth summaries.

**2. Batch Size and Learning Rate** **Batch size:** 8 to balance in between the memory usage and train efficiency. A learning rate that decays from  $3e-5$  during training is used in order not to overfit and ensure

convergence.

**3. Early Stopping :** Early stopping is used to prevent the model from overfitting. The training is ended if the performance of the model on the validation set does not improve after a specified number of epochs, usually 3.

**4. Regularization:** Dropout layers are added in both the encoder and decoder to prevent overfitting, thus enhancing the ability of the model to generalize.

Fine-tuning BART on a summarization-specific dataset allows the former to adapt the linguistic and stylistic idiosyncrasies of news summaries, thus producing output that is as accurate as it is contextually apposite.

#### Training Configuration

We train the model on a GPU for faster training. The code follows the documentation in the Hugging Face Transformers library, which provides utilities to load pre-trained models, tokenize data, and train models on custom datasets. The training parameters are thus set as:

- **Epochs:** 5 epochs to allow adequate training without overfitting.
- **Optimizer:** Adam optimizer with weight decay were used to prevent overfitting and ensure it converged.

- **Learning Rate Scheduler:** A linear learning rate scheduler with warm-up steps is used to stabilize training, ensuring that the learning rate starts low and increases gradually.

## 1.5 Evaluation Metrics

The model's performance is evaluated using the ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metrics:

1. **ROUGE-1:** Overlap of unigrams between the predicted and reference summaries to measure the degree of content accuracy.
2. **ROUGE-2:** Evaluates the overlap of bigrams that measures the quality of the phrase structures in the extracted summary.
3. **ROUGE-L:** Estimates the length of the longest common subsequence (LCS) between the prediction and reference summaries, which epitomizes coherence and fluency in the generated text. Summarization quality is evaluated through a comprehensive set of metrics that assesses the closeness between the generated and actual summaries and higher ROUGE scores represent better appropriateness.

## 1.6 Summary of Methodology

In summary, our method fine-tunes the BART model on the CNN/Daily

Mail dataset with a particular pre-processing pipeline, training configurations, and evaluation metrics. Using transformer-based architecture, we would approach generating high-quality abstractive summaries that go beyond the outcomes of more traditional extractive methods.

## Results and Discussion

It presents an analysis of model performance by usage of metrics for evaluation, quality, and limitations of summaries produced. The obtained results are benchmarked against ROUGE scores, while qualitative examples are examined to elaborate on the strengths and areas for further development.

### A. Quantitative Results

The performance of the fine-tuned BART model is evaluated using the ROUGE-1, ROUGE-2, and ROUGE-L metrics, which measure the overlap of unigrams, bigrams, and longest common subsequences between the generated summaries and reference summaries, respectively.

Metric	Score (%)
ROUGE-1	44.5
ROUGE-2	21.3
ROUGE-L	41.7

The results show that the model delivers competitive performance, particularly at capturing key



concepts (ROUGE-1) and coherence (ROUGE-L). High ROUGE-1 and ROUGE-L scores show that the BART model tends to effectively capture relevant information and maintain structural coherence in the summaries. However, the lower ROUGE-2 score proposes that some information could be lost on the phrase level, typical for abstractive summarization due to the model's tendency to rephrase sentences.

**B. Qualitative Analysis** We read through some of the summaries produced by our model as human reference to check for coherence, relevance, and conciseness. Some findings are:

1. **Coherence:** In most summaries generated by BART, coherence and lack of grammar errors are good indicators. Many often denote the main points about the article at hand, suggesting that the model has picked on the core context and can produce natural sentences.
2. **Relevance:** In general, most of the produced summaries contain the main information from the source articles, but peripheral facts may be omitted in some cases. For instance, while summarizing news articles, the model systematically captures significant entities and events but sometimes omits less relevant background information, which is good for

producing short summaries.

3. **Abstractive Quality:** To a great extent, this model would paraphrase sentences more than copying them while reflecting the capability of the actual abstractive summarization model. Rephrased sentences mostly go with human summaries since they are less redundant and do not contain repeated phrases; however, sometimes rephrased sentences may not meet the intended meaning and result in minor inaccuracy.

### **C. Error Analysis and Limitations**

Despite its strong performance, the model exhibits certain limitations, particularly when handling complex sentences or nuanced content. Notable limitations include:

1. **Loss of Specific Details:** The model may omit specific details that contribute to a more nuanced understanding of the article. For example, in some summaries, numerical values or specific names of individuals are occasionally omitted, reducing factual completeness.
2. **Difficulty with Long-Range Dependencies:** BART, while effective for shorter sequences, may struggle with long articles that contain dependencies between distant

sentences. This limitation occasionally results in summaries that lack a comprehensive view of the entire article.

3. **Over-generalization:** In some cases, the model produces summaries that are overly generalized, missing subtle aspects of the text that could add depth to the summary. This could be attributed to the model's pre-training objectives, which may prioritize overall coherence over fine-grained detail.

#### D. Comparison with Extractive Summarization

Comparing the BART model's abstractive summaries to those generated by extractive summarization methods reveals distinct advantages and challenges:

- **Abstractive Advantage:** The BART model's summaries are generally more concise and natural, as they do not rely on direct sentence extraction. This results in a less repetitive and more fluid reading experience.
- **Challenges in Abstractive Summarization:** Unlike extractive methods, which retain exact phrases from the source text, abstractive models like BART are more prone to introducing minor inaccuracies. This trade-off between naturalness and precision highlights the

complexity of abstractive summarization.

#### E. Future Improvements

To further enhance performance, several directions can be explored:

1. **Hybrid Approaches:** Combining extractive and abstractive methods could yield better results, as extractive techniques can highlight essential sentences that the model could then refine.
2. **Data Augmentation:** Introducing additional training data from diverse domains could improve the model's robustness and generalization.
3. **Improved Fine-Tuning Techniques:** Advanced techniques such as reinforcement learning or adversarial training could refine the model's summarization quality by emphasizing more critical elements.

#### F. Summary of Results

Overall, the fine-tuned BART model achieves competitive performance on the CNN/Daily Mail dataset, demonstrating the feasibility of transformer-based models for abstractive summarization tasks. While some challenges remain, the model effectively captures key information and generates coherent summaries, making it a promising approach for real-world summarization applications.

## 5. Conclusion and Future Work

### A. Conclusion

In this research, we implemented and fine-tuned a BART transformer model for the task of abstractive text summarization. The model was evaluated on the CNN/Daily Mail dataset using standard ROUGE metrics, demonstrating competitive performance in capturing the main ideas and providing coherent summaries. The use of transformer-based architectures, such as BART, highlights the potential of these models to generate high-quality, human-like summaries by leveraging their ability to understand and paraphrase content contextually. This approach contributes to the field of natural language processing by offering a powerful method for abstractive summarization that could be applied to various domains, including news, legal documents, and scientific literature.

The findings reveal that transformer-based models can outperform traditional extractive methods in creating concise and meaningful summaries. However, certain limitations persist, especially with handling complex sentence structures and maintaining factual accuracy. These challenges underline the complexity of abstractive summarization tasks and suggest a need for continued research to improve model precision and depth in capturing nuanced information.

### B. Future Work

While the BART model shows promising results, several avenues could enhance its performance and applicability:

1. **Hybrid Summarization Techniques:** Developing hybrid approaches that combine extractive and abstractive techniques may improve the model's ability to retain critical information while maintaining conciseness and coherence.
2. **Enhanced Fine-Tuning Methods:** Techniques such as reinforcement learning could be applied to fine-tune the model's summarization objectives, potentially improving both accuracy and contextual relevance. Additionally, exploring domain-specific fine-tuning, such as on legal or medical texts, could enhance performance in specialized areas.
3. **Incorporating Additional Evaluation Metrics:** Beyond ROUGE scores, future research could include metrics that better capture readability and factual accuracy. This could help evaluate models more holistically, particularly for sensitive applications where information integrity is crucial.
4. **Handling Longer Text Sequences:** Future work could address the limitations of transformer models in

processing long texts by exploring architectures or techniques specifically designed for handling longer sequences, such as Longformer or BigBird, which can better capture long-range dependencies.

5. **Data Augmentation and Diversity:** Expanding training datasets to include diverse text genres and languages may help create a more robust and generalized model, enabling broader applicability across different summarization contexts.

In summary, this research demonstrates the effectiveness of transformer architectures for text summarization while also highlighting areas for improvement. The continued advancement of transformer models, combined with ongoing research in NLP techniques, presents exciting possibilities for future summarization applications and other complex language processing tasks.

## 6. References

- [1] Salton, G., & Buckley, C. "Term-Weighting Approaches in Automatic Text Retrieval." Information Processing & Management, 1988.
- [2] Erkan, G., & Radev, D. R. "LexRank: Graph-based Lexical Centrality as Salience in Text Summarization." Journal of Artificial Intelligence Research, 2004.
- [3] Mihalcea, R., & Tarau, P. "TextRank: Bringing Order into Texts." Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP).
- [4] Knight, K., & Marcu, D. "Summarization beyond Sentence Extraction: A Probabilistic Approach to Sentence Compression." Artificial Intelligence, 2002.
- [5] Bahdanau, D., Cho, K., & Bengio, Y. "Neural Machine Translation by Jointly Learning to Align and Translate." Proceedings of ICLR 2015.
- [6] Vaswani, A., et al. "Attention is All You Need." Proceedings of NeurIPS, 2017.
- [7] Liu, Y., & Lapata, M. "Text Summarization with Pretrained Encoders." Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP).
- [8] Raffel, C., et al. "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer." Proceedings of JMLR, 2020.
- [9] Lewis, M., et al. "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension." Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL), 2020.
- [10] Zhang, X., et al. "PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization." Proceedings of the 37th International Conference on Machine Learning (ICML), 2020.