# Automated Text Summarization Report
## version 1.0

**Zheng Yuan**

## Abstract

Text summarization is an essential task in natural language processing (NLP) for condensing large volumes of text data while preserving essential information. In this study, I compared the performance of four pre-trained models, namely BART-base, BART-large, T5-base, and FLAN-T5-large, to develop an automated text summarization system. The models were fine-tuned on a given dataset for the task of abstractive summarization. While the T5-base model demonstrated promising results, the other models faced challenges in the fine-tuning process or produced unsatisfactory outcomes.

The T5-base model achieved a training loss reduction from 0.885 for the first epoch to 0.669 at the end, and it outperformed the baselines with ROUGE-1 and ROUGE-2 scores of 0.3147 and 0.1204, respectively. In comparison, the BART-base model demonstrated unstable results, with fluctuating training losses and NaN validation losses in the latter epochs. The fine-tuning process for the FLAN-T5-large model could not be completed due to high computational requirements and long training time. The BART-large model exhibited an increase in training loss after the second epoch.

These results indicate that the choice of pre-trained models and their respective sizes significantly impact the fine-tuning process and the quality of generated summaries, as evidenced by the ROUGE scores. Further investigation is needed to identify optimal configurations and improvements for the models that underperformed. This study emphasizes the importance of careful model selection and fine-tuning in the development of automated text summarization systems. It also highlights the challenges faced in fine-tuning some pre-trained models for abstractive summarization tasks and underscores the need for continued research and optimization to achieve better performance across various models.

## 1 Introduction

In recent years, the rapid growth of digital data has led to an exponential increase in the amount of text data that needs to be processed and analyzed. As a result, there is a growing demand for methods that can efficiently condense large volumes of text data while preserving essential information. Automated text summarization, a field within natural language processing (NLP), addresses this need by generating concise representations of documents or document collections.

Many existing approaches to text summarization have focused on various techniques, including rule-based methods, supervised learning, and unsupervised learning. In this project, I explored a different approach by leveraging pre-trained transformer models, such as BART-base, T5-base, BART-large, and FLAN-T5-large, and fine-tuning them for the task of abstractive text summarization on a news dataset. Although some of these models did not perform as well as expected, this approach still provides valuable insights into the potential of using pre-trained models for text summarization tasks.

In this project, I have fine-tuned the pre-trained models for text summarization and compared their performance against competitive baselines, such as using the first three sentences and the first sentence from the news articles. The results showed that the T5-base fine-tuned model outperformed the lead sentence baseline but did not surpass the first three sentences baseline, indicating the need for further improvements and refinements in this approach. However, other models faced challenges in the fine-tuning process or produced unsatisfactory outcomes. I also analyzed the impact of various training strategies, learning rates, and other hyperparameters on model performance.

Solving the problem of automated text summarization is important because it benefits various stakeholders, including researchers, journalists,

and organizations that require quick access to condensed information for decision-making, analysis, or other purposes. Automating the summarization process can save time and resources, allowing for more effective utilization of human resources in higher-level tasks.

The main contributions of this project are the exploration of pre-trained transformer models for text summarization and the fine-tuning of these models for the specific task. Additionally, the project offers insights into the effectiveness of various training strategies, learning rates, and hyperparameters on model performance, as well as identifying areas where certain models may require further optimization for improved results in the field of text summarization. By building on the insights gained in this project, I aim to contribute to the development of more effective and robust text summarization systems.

## 2 Data

This project employed the CNN/DailyMail dataset, a well-established benchmark dataset in the field of text summarization. The dataset comprised news articles and associated abstractive summaries from CNN and Daily Mail websites. The data is available in its original format from the GitHub repository, which includes the original source code and instructions for data preprocessing and tokenization.

### 2.1 Data Preprocessing

The CNN/DailyMail dataset has three main fields:

- id: the SHA1 hash of the URL where the story was retrieved from

- article: the body of the news article

- highlights: the highlight of the article as written by the author

### 2.2 Data Statistics

The dataset consists of 300,000 instances. As this is not a classification task, there is no class distribution. However, some statistics are presented in Table 4.

The average token count for the articles and the highlights are provided in Table 2:

### 2.3 Annotation

Since the dataset comes with abstractive summaries, no additional annotation is required.

|  | Instance |
|---|---|
| Training data | 287,113 |
| Validation data | 13,368 |
| Test data | 11,49 |

Table 1: The statistics of the CNN/DailyMail dataset

| Feature | Mean Token Count |
|---|---|
| article | 781 |
| highlights | 56 |

Table 2: The average token count for the articles and the highlights

## 3 Related Work

This section provides an overview of the existing literature on text summarization, focusing on six influential papers that have shaped the field. The discussion includes the authors' approaches, differences between their problem formulations, and reasons why my approach may lead to better results.

### 3.1 Abstractive Text Summarization Using Sequence-to-Sequence RNNs and Beyond

Konstantinov et al. (Nallapati et al., 2016) present an abstractive summarization model based on sequence-to-sequence Recurrent Neural Networks (RNNs). The authors design their model to generate new phrases and sentences, moving beyond simple extraction techniques. They evaluate their model on the Gigaword dataset and use ROUGE metrics for comparison with reference summaries. My approach differs in that I explored more recent advancements in neural network architectures, such as transformer models, for improved performance.

### 3.2 Get To The Point: Summarization with Pointer-Generator Networks

See et al. (See et al., 2017) propose a novel attention-based neural network architecture for summarization, combining extractive and abstractive techniques. The pointer-generator network can selectively copy words from the source text, generating a summary that is a mixture of summary-specific words and copied words from the source. The authors use the CNN/Daily Mail dataset for training and evaluation, with ROUGE scores as their evaluation metric. While my approach aimed to build upon this work by further improving the model's ability to generate coherent and informative summaries.

### 3.3 Sentence Simplification with Deep Reinforcement Learning

Zhang et al. (Zhang and Lapata, 2017) present a reinforcement learning-based approach to sentence summarization. The authors train a deep neural network to make a sequence of decisions, selecting a subset of input sentences to form a summary based on rewards from a reward function evaluating summary quality. They use three datasets, including Wikismall, Wikilarge, and Newsela. Evaluation metrics include the Flesch-Kincaid Grade Level index for readability and BLEU for content overlap. The Flesch-Kincaid Grade Level index and BLEU metrics used by Zhang et al. provide valuable insights into summary readability and content overlap, respectively. I considered these metrics in addition to ROUGE scores when evaluating the performance of these fine-tuned models to ensure comprehensive analysis.

### 3.4 BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension

Lewis et al. (Lewis et al., 2019) introduce BART, a denoising autoencoder for pretraining sequence-to-sequence models. BART demonstrates state-of-the-art performance on various NLP tasks, including abstractive summarization. The authors use the CNN/Daily Mail dataset for evaluation and obtain high ROUGE scores. I included BART-base and BART-large as part of the pre-trained transformer models while exploring in my project. Given the success of BART in abstractive summarization tasks, this choice contributed to improved performance on the CNN/Daily Mail dataset.

### 3.5 PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization

Zhang et al. (Zhang et al., 2020) present PEGASUS, a pretraining method specifically designed for abstractive summarization. The authors propose a new pretraining objective, where important sentences are removed from the input, and the model learns to generate them. PEGASUS achieves state-of-the-art results on several summarization benchmarks, including the CNN/Daily Mail dataset. The success of PEGASUS in abstractive summarization tasks serves as motivation for me to explore other innovative techniques and strategies to enhance the

models' performance. By learning from PEGASUS's design and pretraining approach, I refined the models and optimize them for the task of text summarization.

### 3.6 ERNIE-GEN: An Enhanced Multi-Flow Pre-training and Fine-tuning Framework for Natural Language Generation

Xiao et al. (Xiao et al., 2020) present ERNIE-GEN, an enhanced multi-flow pretraining and fine-tuning framework designed for natural language generation tasks, including abstractive summarization. The authors incorporated a multi-flow mechanism that allows the model to learn from both masked language modeling and sequence-to-sequence objectives. They evaluated ERNIE-GEN on several benchmarks, including the CNN/Daily Mail dataset, and reported state-of-the-art results. In my project, I will probably explore the potential of adopting ERNIE-GEN's multi-flow mechanism and pretraining strategies to achieve better text summarization outcomes.

## 4 Methods

This project aims to explore the potential of using pre-trained transformer models for text summarization tasks. The methods employed in this project can be divided into four main steps: data preprocessing, model selection, fine-tuning, evaluation, and libraries used.

### 4.1 Data Preprocessing

Data preprocessing is a crucial step to ensure the quality and consistency of the input data before training and evaluating the models. The preprocessing steps include:

- Tokenization: The text data in the CNN/DailyMail dataset was tokenized using the respective tokenizer for each pre-trained transformer model. Tokenization converted the raw text into a format that can be fed into the models.

- Truncation: To manage the computational resources and maintain consistency across the dataset, the input text was truncated to a fixed length. In this project, a maximum sequence length of 512 tokens was used.

- Padding: The input sequences were padded to ensure that they all have the same length.

This is required for batch processing during training and evaluation.

## 4.2 Model Selection

Four pre-trained transformer models are selected for this project, including BART-base, T5-base, BART-large, and FLAN-T5-large. These models have shown promising results in various natural language processing tasks, such as text classification, translation, and summarization. They are all available in the Hugging Face Model Hub[1].

## 4.3 Fine-tuning

The selected pre-trained transformer models were fine-tuned for the extractive text summarization task using the training data from the CNN/DailyMail dataset. The fine-tuning process involved the following steps:

- Hyperparameter selection: Various hyperparameters, such as learning rate, batch size, and the number of training epochs, are selected and tuned to optimize the model performance.

- Training loop: The models are trained on the training data using the selected hyperparameters. The training process involves passing the input data through the models, computing the loss, and updating the models' parameters using gradient-based optimization algorithms.

- Validation: The models' performance is periodically evaluated on the validation data to monitor the training progress and prevent overfitting.

## 4.4 Libraries Used

In this project, I utilized several established libraries to implement and evaluate my fine-tuned transformer models for text summarization. The following is a list of the primary libraries used, along with their respective URLs:

1. Hugging Face Transformers[2]: The Hugging Face Transformers library provides state-of-the-art natural language processing models, including BART, T5, and FLAN-T5. I used this library to fine-tune my pre-trained transformer models for the extractive text summarization task.

2. Hugging Face Datasets[3]: This library offers an efficient way to access and load various benchmark datasets, including the CNN/Daily Mail dataset used in our project. I used the Hugging Face Datasets library to obtain and preprocess the data for my fine-tuning process.

3. PyTorch[4]: PyTorch is an open-source machine learning framework that enabled me to efficiently build and train deep learning models, such as the transformer models used in our project.

4. Rouge[5]: The Rouge library is a Python implementation of the ROUGE metric for evaluating the performance of text summarization models. I used this library to compute ROUGE-1, ROUGE-2, and ROUGE-L scores for these fine-tuned models, comparing their performance with competitive baselines.

# 5 Evaluation and Results

In this section, I will provide an overview of the evaluation setup for the automated text summarization system, including the methods I compared against and the steps taken to measure the system's success. And the performance for each of the models I have used to fine tune.

## 5.1 Baseline Methods

I compared the method against two baselines:

**Lead sentence baseline:** A simple yet reasonable baseline was using the first sentence of the input articles as the generated summaries. This is a common approach in extractive summarization and often provides a decent summary, as news articles often contain crucial information in the opening sentences.

**First three sentences baseline:** Then I selected first three sentences from the input articles as the generated summaries for the second baselines.

## 5.2 Proposed Method

The proposed method will be compared against the aforementioned baselines. The performance of the Transformer-based abstractive summarization system will be evaluated using the same metrics as the baseline methods.

## 5.3 Evaluation Metrics

To measure the success of my method and compare it against the baselines, I used the following evaluation metric:

---

[1] https://huggingface.co/models
[2] https://github.com/huggingface/transformers
[3] https://github.com/huggingface/datasets

[4] https://pytorch.org/
[5] https://github.com/pltrdy/rouge

**ROUGE:** ROUGE (Recall-Oriented Understudy for Gisting Evaluation) is a widely used set of metrics in text summarization, comparing the generated summaries with the reference summaries. I will report ROUGE-1, ROUGE-2, and ROUGE-L scores.

## 5.4 Results

1. BART-base model: In this study, I employed the BART-base model for training and validation on the CNN/Daily Mail dataset. The training process was conducted over two epochs. After the completion of the first epoch, the training loss was observed to be 0.6991, and the validation loss was found to be 0.5651. Upon completion of the second epoch, the training loss reduced to 0.5072, indicating an improvement in the model's performance on the training data. However, the validation loss remained constant at 0.5651.

The constant validation loss throughout the training process, despite the improvement in training loss, suggests that the model may be overfitting to the training data and not generalizing effectively to the validation data.

2. For the BART-large model: Based on the unsatisfactory performance of the BART-base model, I hypothesized that the model's relatively small size could be the primary reason behind its inability to learn effectively. As a result, I decided to switch to the BART-large model to see if a larger model would perform better.

With the same data preprocessing and training process in place, I adjusted the batch size to 4, the learning rate to 3e-5, and the gradient accumulation steps to 4. I then conducted the training and validation over two epochs.

Upon completing the first epoch, the training loss for the BART-large model was 13.8176, and the validation loss was 16.2573. After the second epoch, the training loss increased to 16.9347. The increase in training loss and relatively high validation loss suggest that the model might not be learning effectively and might require further optimization of its training parameters.

3. For the FLAN-T5-large model: No result (detail in Other Things We Tried section)

4. For T5-Base model:

A table is included to display the results of the proposed method and the baseline methods. The table contained the ROUGE scores for each method.

Our evaluation results are presented in Table 3, which shows the ROUGE-1, ROUGE-2 scores, and ROUGE-L scores for each method:

| Method | ROUGE-1 | ROUGE-2 |
|---|---|---|
| First 3 sentences | 0.3598 | 0.1466 |
| First sentence | 0.2468 | 0.0802 |
| T5-base | 0.3147 | 0.1204 |

Table 3: Evaluation results for lead sentence baseline, first three sentences, and proposed method

The t5-base model with the statistic was fine tuned on the following hyper-parameters:

- Learning Rate: 5e-5 (I think the learning rate is big enough since after warm steps, the loss decreased significantly.)

- Epoch: 3 (After three epoches, my device would be cut off for some unkonw reasons. And It seems like 3 epoches are not enough. Applying to 5 epoches will lead to a better performance)

- Batch Size: 16

- Gradient Accumulation Steps: 4 (In order to prevent out-of-memory issues)

In Table 4, it's showing that in different epoch of fine tuning, the different performance on t5-base model for this task:

| Epoch | ROUGE-1 | ROUGE-2 |
|---|---|---|
| 1 | 0.2433 | 0.1055 |
| 3 | 0.3147 | 0.1204 |

Table 4: Evaluation results of T5-base model for 1 epoch and 3 epoch

## 6 Discussion

In this section, I unpack the results of my evaluation, discussing the performance of the proposed method and the baselines, and interpreting the findings to guide future work.

### 6.1 Overall Performance

My proposed T5-base abstractive summarization system achieved a ROUGE-1 score of 0.3147 and a ROUGE-2 score of 0.1204, outperforming the first sentence baseline but underperforming the first three sentences baseline. While these scores demonstrate the effectiveness of the method in generating informative and coherent summaries, there

is still room for improvement. The performance may be satisfactory for certain end-users, but for more demanding applications or to facilitate scientific advancements, further optimization and enhancements are needed.

## 6.2 Results Analysis

First 3 sentences: This baseline takes the first three sentences of the input text as the summary. It has the highest ROUGE-1 and ROUGE-L scores (0.3598 and 0.3309, respectively), indicating a better overlap of unigrams and longest common subsequences between the generated and reference summaries. The ROUGE-2 score (0.1466) is moderate, suggesting a reasonable overlap of bigrams.

First sentence: This approach takes only the first sentence of the input text as the summary. It has the lowest scores across all ROUGE metrics (0.2468 for ROUGE-1, 0.0802 for ROUGE-2, and 0.2183 for ROUGE-L), indicating that the first sentence alone does not capture the content of the reference summaries well.

T5-base (fine-tuned, 3 epochs): The T5-base model fine-tuned for 3 epochs on the dataset shows a good performance with ROUGE-1 and ROUGE-L scores (0.3147 and 0.3019, respectively) that are close to those of the first 3 sentences approach. The ROUGE-2 score (0.1204) is also higher than that of the first sentence method, but lower than the first 3 sentences method. These results suggest that the fine-tuned T5-base model provides summaries that capture the reference summaries' content reasonably well, but there is still room for improvement.

## 6.3 Comparison with Baselines

The performance of the proposed method was better than the first sentence baseline but worse than the first three sentences baseline. The difference in performance can be attributed to the ability of the T5-base model to generate abstractive summaries, capturing the essence of the input text while rephrasing and condensing the information. The first sentence baseline, being a simpler method, lacks this ability and is thus less effective at generating high-quality summaries. The first three sentences baseline, while having better ROUGE scores, may not always produce coherent and focused summaries since it relies on the input text's structure without any analysis.

## 6.4 Strengths and Weaknesses

My approach performed well in several aspects. The Transformer architecture allowed for the effective capturing of long-range dependencies and complex relationships between words and phrases in the input text. This enabled the generation of more coherent and contextually accurate summaries. Additionally, the use of pretraining and fine-tuning strategies helped the model learn meaningful representations and adapt to the summarization task.

However, the approach also faced some challenges. The performance on long documents was not as strong as I expected, possibly due to the limited capacity of the model to retain and process lengthy sequences of text. Another potential issue was the model's occasional struggle with maintaining factual accuracy or preserving the context of the original text. These weaknesses could be addressed through further research, optimization, and model refinements, such as increasing the number of epochs.

## 7 Conclusion

In this report, I presented the approach to the problem of automated text summarization using a T5-base Transformer-based abstractive summarization system. I compared my proposed method against two baselines: a first three sentences performance baseline and a lead sentence baseline. My evaluation results indicated that the method outperformed the lead sentence baseline but did not surpass the first three sentences baseline, highlighting the need for further improvements and refinements.

Despite the results not meeting the highest expectations, this study has identified certain limitations and areas for improvement. For instance, the method may not be well-suited to summarizing long documents or handling specific types of text data, such as scientific articles or legal documents. Additionally, this model may struggle with retaining factual accuracy or preserving the context of the original text in some cases.

As a part of future work, I plan to address these limitations and explore new avenues for research. This could involve incorporating additional sources of information, experimenting with alternative model architectures or more advanced models, or tackling different types of summarization tasks. I will also continue refining the method to improve its performance, scalability, and applicability to a broader range of text data. By building on the in-

sights gained in this project, I aim to contribute to the development of more effective and robust text summarization systems.

## 8 Other Things I Tried

During the course of our project, I faced numerous challenges and tried various approaches before settling on the final methodology presented in the previous sections.

1. Initially, I spent a considerable amount of time attempting to build and train my own model for text summarization, inspired by the paper "Hierarchical Attention Networks for Document Classification" by Zichao Yang et al. The goal was to develop a hierarchical attention mechanism that could capture both the word-level and sentence-level representations of the input text. I believed that such a model could potentially provide a more nuanced understanding of the document structure, which could in turn lead to improved summarization results.

To implement this approach, I started by building a word-level attention model, where each word in a sentence was assigned a weight based on its importance. Subsequently, I designed a sentence-level attention model that assigned weights to the sentences based on their significance within the document. By combining these two levels of attention, I aimed to create a comprehensive representation of the document that could be used to generate accurate and coherent summaries.

However, I encountered several issues during the implementation of this approach. First, the model proved to be computationally expensive, requiring significant memory and processing power. This made it difficult to scale the model to handle larger datasets or longer documents, limiting its practical applicability. Furthermore, I found that the model struggled to capture the contextual relationships between words and sentences effectively, often leading to the generation of summaries that were semantically disconnected or inaccurate.

Despite spending a considerable amount of time fine-tuning the model's architecture and hyperparameters, I was unable to achieve satisfactory performance with this approach. It became evident that the hierarchical attention mechanism, while conceptually promising, required further development and optimization to be effective in the context of text summarization.

Ultimately, I decided to pivot my focus towards a Transformer-based abstractive summarization sys-

tem. This decision was based on the understanding that the Transformer architecture, with its self-attention mechanism and parallel processing capabilities, could better capture the long-range dependencies and contextual relationships within the input text.

2. And for the FLAN-T5-large model: I utilized the Trainer from the Transformers library to facilitate the training process. However, I encountered challenges due to limited computational resources, particularly GPU resources, which hindered the model's training.

As shown in the log, the training progress was extremely slow, with only 501 out of 35888 steps completed in over 32 minutes. The estimated time to complete one epoch was 38 hours and 19 minutes, which was not feasible given my available resources. The limitations in compute units and GPU resources significantly affected the training process for the FLAN-T5-large model.

## 9 What I Would Have Done Differently and Next

If I was to start some parts of the project over, I would have done the following differently:

1. Model Selection: I experimented with various model architectures, including BART-base, BART-large, and FLAN-T5-large. Given the resource constraints and time limitations, it would have been more prudent to start with smaller models that are easier to train, and gradually scale up to larger models if needed.

2. Compute Resources: I faced challenges with the available computational resources, particularly GPU resources. In hindsight, I would have secured more powerful hardware or cloud resources to ensure a smoother and faster training process.

3. Pre-processing and Data Cleaning: I could have invested more time in cleaning and pre-processing the dataset to ensure that the input to the models was of high quality. This might have helped in improving the overall performance of the models.

4. Evaluation Metrics: I relied on ROUGE scores to evaluate these models, but it would have been helpful to explore other metrics, such as BLEU or METEOR, to gain a more comprehensive understanding of the model's performance.

Ideas I would have wanted to try but didn't have time:

1. Transfer Learning: I could have explored

transfer learning by fine-tuning pre-trained models on domain-specific datasets to improve their performance in specific summarization tasks, such as summarizing scientific articles or legal documents.

2. Multi-modal Summarization: I could have explored multi-modal summarization, which combines textual and visual information to generate more informative and comprehensive summaries.

The results did point towards some potential next steps:

1. Investigate Model Efficiency: My experience with the FLAN-T5-large model showed the importance of efficient models. I could explore model compression techniques, such as pruning or quantization, to create more efficient models that retain their performance while requiring fewer computational resources.

2. Distributed Training: To address the limitations of our available computational resources, I could explore distributed training or model parallelism techniques, which would enable me to leverage multiple GPUs and speed up the training process.

# References

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.

Ramesh Nallapati, Bowen Zhou, Caglar Gulcehre, Bing Xiang, et al. 2016. Abstractive text summarization using sequence-to-sequence rnns and beyond. *arXiv preprint arXiv:1602.06023*.

Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368*.

Dongling Xiao, Han Zhang, Yukun Li, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. 2020. Ernie-gen: An enhanced multi-flow pre-training and fine-tuning framework for natural language generation. *arXiv preprint arXiv:2001.11314*.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *International Conference on Machine Learning*, pages 11328–11339. PMLR.

Xingxing Zhang and Mirella Lapata. 2017. Sentence simplification with deep reinforcement learning. *arXiv preprint arXiv:1703.10931*.