**Overview**

This survey note provides a comprehensive analysis of a project focused on developing an automatic text summarization system using the BART (Bidirectional and Auto-Regressive Transformers) model, fine-tuned on a subset of the CNN/DailyMail dataset. The project, detailed in an uploaded Jupyter Notebook file named "nlp-project.ipynb," aimed to generate concise summaries from long news articles, with results evaluated against state-of-the-art research. The analysis covers methodology, results, comparisons, challenges, and future directions, ensuring a thorough understanding for researchers and practitioners in NLP.

**Project Context and Objectives**

The project addresses the growing need for automated text summarization in an era of information overload. With the exponential growth of textual data, manual summarization is inefficient, particularly for news articles. The objective was to leverage a pre-trained BART model, known for its effectiveness in text generation tasks, and fine-tune it on the CNN/DailyMail dataset to create a practical summarization tool. This dataset, widely used in NLP research, contains news articles and human-written summaries, making it ideal for training abstractive summarization models.

**Methodology**

The methodology involved several key steps, detailed as follows:

- **Dataset Selection**: The CNN/DailyMail dataset was chosen, with a subset of 10,000 training examples used due to computational constraints. According to Papers with Code - CNN/Daily Mail Dataset, the full dataset includes 286,817 training pairs, 13,368 validation pairs, and 11,487 test pairs, highlighting the scale of data typically used in such tasks.
- **Model Choice**: The pre-trained BART-large-cnn model was selected, as it is fine-tuned for summarization and available through the Hugging Face Transformers library. This model, with 406M parameters, is designed for conditional text generation, making it suitable for summarization.
- **Preprocessing**: Articles were tokenized with a maximum length of 1024 tokens, and summaries with 128 tokens, using truncation and padding to ensure consistency. This preprocessing aligns with standard practices for handling variable-length inputs in transformer models.
- **Training Process**: The model was fine-tuned for one epoch, with a batch size of 2, learning rate of 2e-5, and weight decay of 0.01. Mixed precision training (FP16) was employed to optimize performance on the GPU (NVIDIA Tesla T4), and Weights & Biases (WANDB) logging was disabled for simplicity.

- **Deployment**: A Flask web application, integrated with Ngrok for public access, was developed to allow users to input articles and receive generated summaries, demonstrating practical utility.

**Results**

The results of the training and evaluation are summarized in the following table:

| Metric | Value |
| --- | --- |
| Initial Training Loss | 0.7739 |
| Final Training Loss | 0.5482 |
| Samples per Second | 2.637 |
| Rouge-1 Score | 0.28 (28%) |
| Rouge-2 Score | 0.1667 (16.67%) |
| Rouge-L Score | 0.28 (28%) |

**Training Performance**: The training loss decreased from 0.7739 to 0.5482 over 2,500 steps, indicating effective learning. The model processed approximately 2.637 samples per second, reflecting the efficiency of the GPU setup.

**Summarization Example**: A test article on renewable energy was summarized as: "Renewable energy sources are transforming how the world generates electricity. Solar panels and wind turbines are becoming increasingly efficient and cost-effective. Smart grids are enhancing the reliability of renewable energy." This summary captured the main points, demonstrating coherence.

**Evaluation**: The Rouge metric was used to evaluate the summary against a reference, with scores of 0.28 for Rouge-1, 0.1667 for Rouge-2, and 0.28 for Rouge-L. These scores, when converted to percentages (28%, 16.67%, 28%), indicate moderate overlap with the reference summary, particularly in unigrams and longer sequences.

**Comparison to Other Research**

To contextualize the results, a comparison was made with state-of-the-art performance on the CNN/DailyMail dataset, as reported in recent research. The following table summarizes the comparison:

| Model | Dataset | Rouge-1 | Rouge-2 | Rouge-L |
|---|---|---|---|---|
| Our Model | CNN/DailyMail (10,000 subset) | 28% | 16.67% | 28% |
| BART-large-cnn (SOTA) | CNN/DailyMail (full) | 44.16% | 21.28% | 40.90% |

- **State-of-the-Art Performance**: Research, such as the original BART paper by Lewis et al. (2019), reports Rouge scores of 44.16% for Rouge-1, 21.28% for Rouge-2, and 40.90% for Rouge-L when trained on the full CNN/DailyMail dataset.
- **Our Model's Performance**: Our model's scores are significantly lower, likely due to the limited training data (10,000 examples vs. 286,817) and training for only one epoch. This aligns with findings that larger datasets and longer training durations improve model generalization, as noted in [Papers with Code - CNN / Daily Mail Benchmark (Abstractive Text Summarization)](#).

**Discussion**

The project's outcomes highlight both successes and challenges:

- **Challenges**: The primary challenges were the limited size of the training dataset and the short training duration. Computational constraints, such as GPU availability, restricted the use of the full dataset and longer training epochs. Additionally, the lack of explicit ablation studies or hyperparameter searches, common in research settings, limited the optimization potential.
- **What Worked Well**: Despite these limitations, the model generated coherent and relevant summaries, demonstrating the robustness of the BART architecture. The

deployment as a web application was a practical success, enhancing accessibility for users.

- **Future Directions**: To improve performance, future work should consider using the full CNN/DailyMail dataset, increasing the number of training epochs, and experimenting with hyperparameter tuning. Techniques such as data augmentation, advanced fine-tuning strategies, or comparisons with other models (e.g., PEGASUS, T5) could further enhance results, as suggested in [Text Summarization with BART Model | by Sandeep Sharma | Medium](#).

**Conclusion**

This project successfully demonstrated the fine-tuning of a BART model for text summarization on a subset of the CNN/DailyMail dataset, with deployment as a web application. While the model achieved moderate Rouge scores (28% Rouge-1, 16.67% Rouge-2, 28% Rouge-L), its performance was below state-of-the-art levels (44.16% Rouge-1, 21.28% Rouge-2, 40.90% Rouge-L), primarily due to limited training data and insufficient training time. Future research should focus on leveraging the full dataset, extending training duration, and optimizing hyperparameters to achieve better results, aligning with the evolving standards in NLP research.