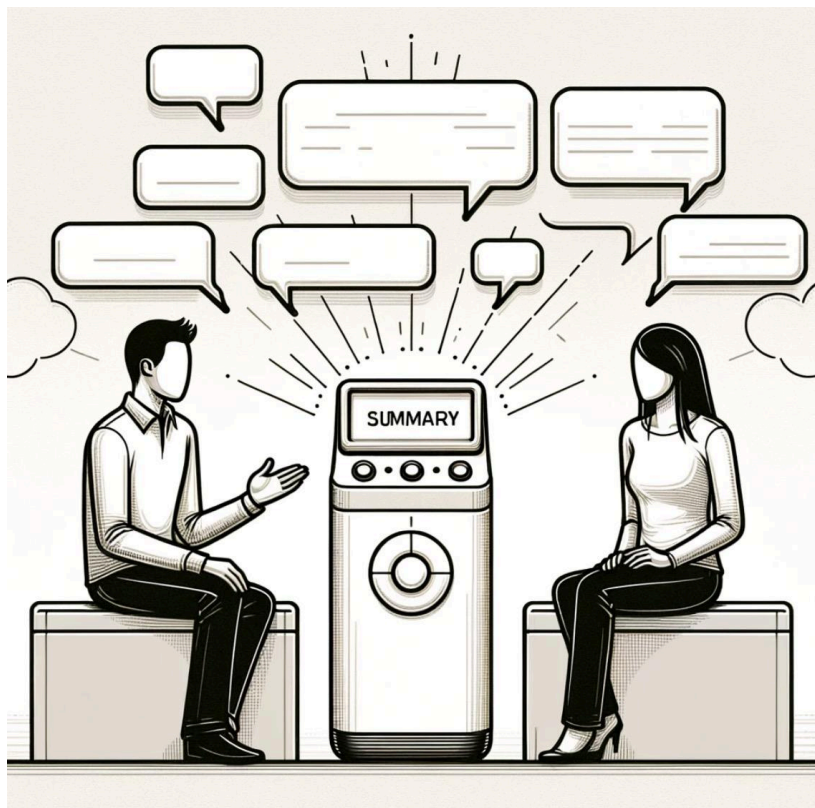# Conversational Summarization Using BART

**INDENG 242B - Machine Learning and Data Analytics II**
**Final Project Report**

Chloe Kang ID: 3039712860
Tunan Li ID:3039669622
Kairan Wang ID: 3034574973
Roxanne Wang ID: 3039679645
Iris Zheng ID: 3039679164

# Introduction

In the realm of natural language processing, the quest for efficient and accurate text summarization has been a longstanding challenge. Leveraging cutting-edge technologies, our project delves into the domain of sentence summarization using the BART (Bidirectional and Auto-Regressive Transformer) architecture. This architecture, known for its prowess in natural language understanding and generation tasks, serves as the cornerstone of our endeavor.

# Motivation

We fine-tune BART models facebook/bart-large-xsum and facebook/bart-large-cnn to distill their verbose text into short summaries. These models are all fine-tuned with a specific dataset to gain the power of summarization. We ensure that no stone is left unturned, from meticulous data preparation, custom training configurations, and rigorous evaluation, both using ROUGE scores. We ensure that the summaries not only give accuracy but also push the text summarization technology in its capabilities to make the results very insightful.

# Data Collection & EDA

The Samsum-train.csv dataset contains 14,732 Training, 818 Validation, and 819 Test conversations with their respective IDs, dialogues, and summaries. There is only one missing value and it was handled  and removed by simply ensuring that the data had conversational elements. In the EDA, different lengths of the dialogue were used (from 31 to 5492 characters) and summary (from 3 to 300 characters) to show different vocabulary patterns and distinct information from dialogue to summary. The dialogues are presented with key points and summaries in a structured manner with the common elements of dialogue. For the training set, these presented many patterns of words, which will be informative for the development of natural language processing-based models. Similar analysis and visualization were performed for the validation and test sets with consistency in data handling and observations across all the three subsets.

# Modeling

Our project employs the BART (Bidirectional and Auto-Regressive Transformer) architecture for the task of sentence summarization, using two different pretrained models: `facebook/bart-large-xsum` and `facebook/bart-large-cnn`. The former is a pre-trained BART model developed by Facebook AI Research fine-tuned on the XSum dataset, which contains single-sentence news articles and their corresponding summaries. The latter was trained on a dataset including articles from the CNN news website and is suitable for various text generation tasks. These models have been selected for their robustness and effectiveness in natural language understanding and generation tasks.

**Model Preparation**

The datasets used include dialogues and respective summaries. We have divided our data into training, testing, and validation sets with 14731, 819, and 818 entries respectively. The datasets were converted from pandas dataframes into Hugging Face's `Dataset` objects to streamline the training and evaluation process.

For both models, we implemented a `preprocess_function` specifically tailored for preparing datasets to be preprocessed by BART. It handles the tokenization of both input and target variables, which is the first step in converting raw text into a form that the model can understand. For the inputs, which are the dialogues in our case, the function extracts each dialogue and encodes it into tokens using the tokenizer1, with a set maximum length of 1024 tokens, applying truncation as needed. The tokenizer is then configured to process the target variables, which are the summaries, encoding them into a tokenized form with a maximum length of 128 tokens, also truncating where necessary. The output of the tokenization includes the `input_ids` for BART, representing the numerical representation of the tokens.The function additionally facilitates labeling, which involves associating each input token with a corresponding label from the target tokenized summary. This is crucial for the model to learn the mapping from inputs to outputs during training.

To complement the preprocessing, a `DataCollatorForSeq2Seq` is used, which is a utility that organizes the tokenized data into batches. It ensures that within each batch, all sequences are padded to the same length, creating uniform input sizes for efficient processing by the sequence-to-sequence model during training. This padding is essential for leveraging the parallel computation capabilities of modern deep learning frameworks and for maintaining consistency across batches during training.

**Model Training and Configuration**
Both models are instantiated with their respective tokenizers and pre-trained weights. The BART model architecture comprises an encoder-decoder structure with multi-layer attention mechanisms, which is ideal for sequence-to-sequence tasks like summarization.

We employed a `Seq2SeqTrainer` from Hugging Face's Transformers library, configured with specific training arguments to optimize the learning process. Throughout the training process, we adjust hyperparameters, such as the learning rate and batch size, and monitor loss metrics to optimize the model's performance. Our arguments specify the number of training epochs, batch sizes, learning rates, and evaluation strategies, among other parameters. For instance, we set a learning rate of 2e-5, a train batch size of 4 per device, and employed gradient accumulation to handle higher effective batch sizes than what single GPU memory could allow. We utilize beam search (k=6) during generation to explore the space of possible summaries and select the one that the model predicts with the highest confidence.

**Evaluation Metrics**
Our custom evaluation function, compute_metrics1, decodes the model's predictions and computes the ROUGE scores to provide a quantitative measure of the summarization quality. The ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metric is a commonly used

metric for evaluating the quality of text summarization or generation tasks. It compares the overlap of n-grams between the generated summaries and the reference summaries. We measure ROUGE-1 (indicates overlap of sole words), ROUGE-2 (measures overlap of two-word sequences), and ROUGE-L (measures the Longest Common Subsequence between the two summaries) scores to assess the quality of our summarizations in terms of unigrams, bigrams, and the longest common subsequence, respectively.

**Fine-tuning and Results**

Through iterative training and validation, we observed how the models adapted to the summarization task.

| Epoch | Training Loss | Validation Loss | Rouge1 | Rouge2 | RougeL | RougeLsum | Gen Len |
|---|---|---|---|---|---|---|---|
| 0 | 1.400500 | 1.474380 | 52.501300 | 27.656500 | 43.656500 | 48.236400 | 29.493300 |
| 2 | 0.842000 | 1.506337 | 52.690900 | 27.688500 | 43.692200 | 48.524100 | 28.652000 |
| 3 | 0.688900 | 1.602949 | 52.639400 | 27.433900 | 43.547400 | 48.267000 | 29.696000 |

Model 1 (facebook/bart-large-xsum) Results

| Epoch | Training Loss | Validation Loss | Rouge1 | Rouge2 | RougeL | RougeLsum | Gen Len |
|---|---|---|---|---|---|---|---|
| 0 | 1.371200 | 1.476597 | 40.101500 | 19.775400 | 30.595000 | 37.298800 | 60.138000 |
| 2 | 0.823200 | 1.471878 | 41.121000 | 20.387700 | 31.620100 | 38.093000 | 59.837600 |
| 3 | 0.674600 | 1.598011 | 41.125700 | 20.499300 | 31.584600 | 37.994100 | 60.084200 |

Model 2 (facebook/bart-large-cnn) Results

The models showed different levels of efficacy, with Model 1 `facebook/bart-large-xsum` generally performing better across the metrics, likely due to its training on a summarization-specific dataset (XSum). As we set `load_best_model_at_end = True` in the trainer's arguments, the training process utilizes a trainer that automatically selects the final model based on the smallest validation loss recorded during training. Therefore, for Model 1, the best-performing epoch according to this criterion would be epoch 0, as it has the lowest validation loss. Similarly, for Model 2, the trainer would select the model from epoch 2, where the validation loss is at its minimum.

Model 1's performance surpasses standard baseline metrics, especially notable in ROUGE-1 scores, consistently above 50. This means that 52.639% of the unigrams in the reference summary are also present in the generated summary, a benchmark indicative of strong unigram matching. Similarly, Model 1's ROUGE-2 (-L) scores remain superior to Model 2's and above 27 (43), denoting an adeptness at capturing multigram relationships. The Gen Len metric (length of the generated summary) further reveals Model 1's capacity to generate succinct summaries, a desirable trait that underscores the model's ability to condense text into its most informative components. See the appendix for an example of input conversation and output summary generated by model 1.

4

# Discussion

**Implications**

The BART model could change the way educational content is being presented by providing brief synopsis that help the learning person understand tough material easily and hence provide more personalized ways of learning. Their effectiveness, however, is dependent more on accurate understanding of various academic disciplines and especially new vocabularies and concepts involved through constant learning, more so in fields that have highly technical jargon be it in legal studies or science. This, therefore, implies the necessity of an ongoing model fine-tuning process that would assure not only its relevance but effectiveness in light of actuality.

**Limitations**

Meanwhile, the heavy computation resources needed for BART models will be not only limiting their application to resource-poor environments but also symbolizing a much larger problem for the broad accessibility of advanced AI technologies. This could disadvantage institutions and human beings in poorer areas, who are getting pushed farther into a corner in order to take advantage of the great opportunities for educational or professional development that these powerful tools offer, and thereby increase the digital divide. This difference indicates the need to create more resource-efficient models that can still be fully functional but are available to a larger audience.

Further, the biases that are embedded in the training data for these models carry considerable risk. If the training data used to train these summarization models is not curated with a fine and careful eye to balance and comprehensiveness, then the produced summaries will do more than just reflect back existing prejudices; they may also introduce biases or distort and even drop important information. This can have dire consequences, especially in those circumstances where precise and objective information is paramount, such as in legal, medical, or educational fields. Correcting for these biases demands effort both in the design of the dataset and in actively monitoring and updating the model to produce fair and truthful output.

# Conclusion

Ultimately, these BART models represent a major advancement beyond the state-of-the-art in NLP, yet their real-world effectiveness and ethical deployment will rely on continued progress in model training, ethical considerations, and computational optimization. Should further advance, it is recommended that the systems minimize biases and computational demands for wider accessibility and fair treatment during automated text summarization.

# Appendix

Below is the text summarization simulation we tested:
**Original Dialogue:**

Marta: Hi, I'm at the supermarket now to make some shopping for todays dinner.
Do you have any wishes?
Nick: Hm I don't know. I haven't eat spaghetti in a while
Marta: Oh no, I've got spaghetti yesterday by Patric and the day before too.
Nick: Okay maybe some fish?
Marta: Yeah fish is great, I'll go and search for something
Nick: Text me what do you find.
Marta: Actually there is one small fish left and I don't think we will be full
from it.
Nick: Let's make a lasagne
Marta: Do you know how much work that is?
Nick: No i don't know
Marta: It's a lot...
Nick: Maybe I could help you?
Marta: That's not a bad idea
Marta: Did you cook something yet?
Nick: No, but I can learn really fast :D
Marta: Okay, I'll buy the meat and sauce an let's do it
Nick: Should I look for a recipe in the Internet?
Marta: No need I did lasagne many times before
Nick: I can't wait ntil you teach me how to cook
Marta: I hope it will be eatable :D


**Reference Summary:**

Marta is grocery shopping for dinner. She and Nick will make lasagne.


**Model-generated Summary:**

Marta is at the supermarket. She will buy fish, meat and sauce for today's
dinner with Nick. They will make lasagna.

Appendix on heatmaps:
Heatmaps of dialogues and summary for unigram, bigram, and trigram in train set

Unigrams Heatmap - Train Summary

Bigrams Heatmap - Train Dialogue

Bigrams Heatmap - Train Summary

Trigrams Heatmap - Train Dialogue

Trigrams Heatmap - Train Summary

Unigrams Heatmap - Validation Summary

Bigrams Heatmap - Validation Dialogue

Bigrams Heatmap - Validation Summary

Trigrams Heatmap - Validation Summary

Trigrams Heatmap - Validation Dialogue



Distribution of Dialogue Lengths in Validation Set



Distribution of Summary Lengths in Validation Set

Distribution of Dialogue Lengths in Text Set


Distribution of Summary Lengths in Text Set