

Text Summarization Web Application - AbbrivioAI

Group Members:

First name	Last Name	Student number
Viki	Patel	C0906295
Devarsh	Jadhav	C0893173
Bhavya	Vadher	C0894977
Yash	Pandit	C0894474
Syed Roman	Uddin	C0906298

Submission date: 12/08/2024

Contents

Abstract	3
Introduction	4
Methods	6
Results	15
Conclusions and Future Work	15
References	23

1. Abstract

In this project, we set out to develop an application capable of summarizing text effectively and precisely using the BART model, which we fine-tuned exclusively for this task using the CNN/DailyMail dataset. For the backend of the application, we chose Django, while we picked React on the front end, allowing us to create an interface that's both intuitive and easy to navigate. To improve the user experience, we utilized HTML and CSS for basic structuring and styling, and integrated material UI for a modern and responsive design. A remarkable feature of our application is the ability to upload PDF documents, making it accessible for users to generate summaries from various text sources. During the project, we faced several obstacles, especially with limited resources and the complexities of deployment. However, we tackled these issues head-on and found ways to resolve them effectively. The final application performed well, achieving strong metrics, such as high ROUGE-1 and ROUGE-L scores, demonstrating its capability to generate high-quality summaries. This report reflects on the methods we took, the outcomes we achieved, and the potential for further upgrades, highlighting the application's relevance in the realm of natural language processing.

2. Introduction

In today's world, where information is available in devastating quantities, it can be tough to go through all the data to find what's truly important. This reality pushed us to develop a tool that can automatically squeeze lengthy texts into clear, short summaries. The goal was to make it simpler for users to quickly grasp the main ideas without walking through endless paragraphs.

Key Technologies and Concepts

To bring this idea to life, we used several key technologies:

- Django: This is a Python framework that makes building web applications simpler and more secure. We used Django to handle the backend of our application, dealing with data management and ensuring everything runs smoothly on the server side.
- React: React is a tool for creating user interfaces that are dynamic and responsive. It allowed us to build the front end of our application, making sure users have a smooth and intuitive experience.
- BART Model: BART stands for Bidirectional and Auto-Regressive Transformers, a type of model that excels at generating text that makes sense in context. It was perfect for our needs because it can understand the gist of an entire text and then summarize it in a way that retains the essential information.
- CNN/DailyMail Dataset: This is a collection of news articles and their corresponding summaries. It's widely used in the field of natural language processing to train and test models that need to summarize text effectively.
- PDF Documentation Input: For handling PDF documents, we integrated a feature to extract and summarize content from PDF files. This feature ensures that textual data is accurately captured and summarizes the given documentation.

Breaking Down the Technical Terms

To make things clearer, here's a quick rundown of some of the technical jargon:

- Seq2Seq (Sequence-to-Sequence) Model: Think of it as a model that can take one sequence, like a sentence, and turn it into another, such as a translated sentence or a summary.
- Transformer: A model that can process data in parallel, which is crucial for understanding long and complex texts.
- Tokenization: This is the process of breaking down a string of text into smaller pieces called tokens. It's like turning a sentence into its individual words or phrases.
- Fine-tuning: This involves taking a model that's already been trained and adjusting it further to perform well on a specific task, like summarizing texts.
- Encoder-Decoder Architecture: This setup allows a model to take in input (like a full text), process it, and then output something new (like a summary).

- Self-Attention Mechanism: A feature of the model that helps it focus on the most important parts of the text when creating a summary.
- ROUGE and BLEU Scores: These are metrics used to measure how well the model-generated summaries match up with the reference summaries.
- API (Application Programming Interface): This is a set of rules that allows different software parts to communicate. In our case, the backend API processes requests from the front end.
- Docker: A tool that helps package the application so it can run consistently across different environments.

Challenges and How We Overcame Them

As with any project, we faced a few bumps along the road:

- Resource Limitations: Training a model like BART is no small feat—it requires a lot of computing power. We often encountered slow processing times and memory issues, which forced us to get creative with how we managed our resources, including using cloud services.
- Model Fine-Tuning: Fine-tuning BART to generate summaries that made sense and retained the original meaning was trickier than expected. We had to experiment with various settings and techniques to get it right.
- Deployment Hurdles: Getting the model to work seamlessly with the front end and handle real-time user requests required several iterations and a good deal of optimization. Setting up the infrastructure was also a learning experience for the team.
- Integrating Frontend and Backend: Making sure React and Django played well together was crucial. We spent a lot of time refining the communication between the two, ensuring that data was passed and processed smoothly.
- Designing the User Experience: We wanted to make sure that anyone using the tool would find it straightforward and intuitive. This meant continuously refining the design and seeking feedback to improve it.

3. Methods

This section outlines the detailed steps and methodologies used in the project, including data preprocessing, model architecture design, training procedures, and evaluation metrics.

Dataset Information and Preprocessing

The **CNN/DailyMail** dataset, a widely used benchmark for text summarization tasks, was used for this project. This dataset consists of news articles paired with their corresponding summaries, making it ideal for training and evaluating summarization models.

Data Preprocessing

Data preprocessing is a crucial step in preparing raw data for model training, ensuring that the data is clean, consistent, and in a format suitable for the models.

I. Tokenization:

- **Description:** Tokenization is the process of breaking down text into smaller units called tokens. Tokens can be individual words or subwords, depending on the tokenizer used. This step is essential because it converts the raw text data into a numerical format that models can process.
- **Implementation:** In this project, tokenization was achieved using tokenizers from the Hugging Face transformers library. These tokenizers are optimized for compatibility with specific models like BART, ensuring that the text is accurately formatted for input into the models.
- **Details:** The tokenizer maps each word or subword in the text to a unique integer. Special tokens, such as start-of-sequence (<s>) and end-of-sequence (</s>), are also added to the text to signify the beginning and end of the input sequence.

II. Lowercasing:

- **Description:** Converting all text data to lowercase is a standard preprocessing step designed for ensuring consistency and reducing the vocabulary size. This step helps minimize the impact of case sensitivity, where different cases of the same word (e.g., "The" vs. "the") might be treated as separate tokens.
- **Benefits:** Lowercasing reduces the complexity of the model's vocabulary and helps in achieving more consistent results during training.

III. Stop Word Removal:

- **Description:** Stop words are common words like "and", "the", and "is" that do not carry significant meaning on their own but are frequently used in sentences. Removing these words can reduce noise in the data and focus the model's attention on more meaningful content.

- **Implementation:** Stop word removal was conducted using the NLTK (Natural Language Toolkit) library, which provides predefined lists of stop words for various languages. The removal process involved filtering out these words from the text data before it was tokenized.

Data Cleaning:

- **Description:** Data cleaning involves removing irrelevant or extraneous characters from the text, such as HTML tags, special symbols, and punctuation. This step is crucial for eliminating noise that could interfere with the model's ability to learn meaningful patterns.
- **HTML Tag Removal:** Any HTML or XML tags embedded in the text were removed using regular expressions.
- **Special Characters:** Characters such as @, #, and punctuation marks were stripped from the text to avoid confusion during tokenization and model training.
- **Whitespace Normalization:** Excessive whitespaces were reduced to a single space to ensure consistency in the text formatting.

Padding and Truncation:

- **Description:** Text sequences vary in length, and models typically require inputs of a fixed size. Padding involves adding special padding tokens (<pad>) to shorter sequences, while truncation involves cutting longer sequences to fit the model's input size.
- **Padding:** Shorter sequences were padded at the end with <pad> tokens to match the maximum sequence length of 84 tokens for inputs and 12 tokens for outputs. This step ensures that all input sequences within a batch have the same dimensions.
- **Truncation:** Longer sequences were truncated to a fixed length to prevent memory issues and to ensure that the model's input size remains manageable.

Splitting:

- **Description:** The dataset was split into three parts: training, validation, and test sets. This division allows for effective training, hyperparameter tuning, and unbiased evaluation of the models.
- **Training Set:** Typically comprising 80% of the data, the training set is used to train the model. The model learns patterns and builds its internal representations based on this data.
- **Validation Set:** Making up 10% of the data, the validation set is used during training to tune hyperparameters and prevent overfitting. It provides a checkpoint to evaluate how well the model is likely to perform on unseen data.
- **Test Set:** The remaining 10% of the data is reserved as a test set. This set is used only after training is complete to provide an unbiased evaluation of the model's performance.

Model Architecture

The project explored two primary model architectures: a custom LSTM-based Seq2Seq model with attention mechanisms and the transformer-based BART model. Each model has its strengths and specific components designed to handle the complexities of text summarization.

Custom LSTM-Based Model (abbreviated AI_1_model):

Architecture Summary: This model is built on a sequence-to-sequence (Seq2Seq) framework with Long Short-Term Memory (LSTM) layers, enhanced by attention mechanisms. The architecture is particularly suited for handling sequential data and capturing temporal dependencies.

Encoder (Bidirectional LSTM):

- The encoder consists of a bidirectional LSTM layer that processes the input text sequence in both forward and backward directions. This bidirectional approach enables the model to capture context from both ends of the sequence, improving its understanding of the text.
- Hidden States: The encoder generates hidden states that encapsulate the context of the input sequence. These states are passed to the decoder to aid in generating the summary.

Concatenate Layer:

The first concatenate layer merges the hidden states from the forward and backward passes of the bidirectional LSTM. By combining these outputs, the model creates a unified representation of the input sequence, which is then used to guide the decoding process.

Concatenate_1 Layer:

The concatenate_1 layer further combines additional hidden states, ensuring that all relevant information from the encoder is integrated before it is passed to the decoder.

Decoder (LSTM with Attention):

The decoder is an LSTM layer that receives the combined hidden states from the encoder. It generates the output sequence token by token, using the attention mechanism to focus on the most relevant parts of the input sequence when predicting each token.

Attention Mechanism:

The attention layer calculates a context vector that weighs the importance of each part of the input sequence based on the current state of the decoder. This mechanism allows the decoder to generate more accurate and contextually relevant summaries by focusing on the key parts of the input.

Concatenate_2 Layer:

- The concatenate_2 layer merges the output from the attention mechanism with the decoder's LSTM output, providing a richer context to the final dense layer. This combined information is crucial for making precise token predictions.

- **Dense Output Layer:** The final dense layer translates the processed data into predicted tokens by mapping the concatenated output to the model's vocabulary space. This layer generates the final summary sequence.

Model Summary:

Model: "model"

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[(None, 84)]	0	[]
embedding (Embedding)	(None, 84, 300)	1813020 0	['input_1[0][0]']
input_2 (InputLayer)	[(None, 12)]	0	[]
bidirectional (Bidirectional)	[(None, 84, 512), (None, 256), (None, 256), (None, 256), (None, 256)]	1140736	['embedding[0][0]']
embedding_1 (Embedding)	(None, 12, 300)	1813020 0	['input_2[0][0]']
concatenate (Concatenate)	(None, 512)	0	['bidirectional[0][1]', 'bidirectional[0][3]']
concatenate_1 (Concatenate)	(None, 512)	0	['bidirectional[0][2]', 'bidirectional[0][4]']
lstm_1 (LSTM)	[(None, 12, 512), (None, 512), (None, 512)]	1665024	['embedding_1[0][0]', 'concatenate[0][0]', 'concatenate_1[0][0]']
attention (Attention)	(None, 12, 512)	0	['lstm_1[0][0]', 'bidirectional[0][0]']
concatenate_2 (Concatenate)	(None, 12, 1024)	0	['attention[0][0]', 'lstm_1[0][0]']
dense (Dense)	(None, 12, 60434)	6194485 0	['concatenate_2[0][0]']
=====			
Total params: 101011010 (385.33 MB)			
Trainable params: 64750610 (247.00 MB)			
Non-trainable params: 36260400 (138.32 MB)			

Performance:

The model was trained over 10 epochs, achieving a final training accuracy of 85.63% and a validation accuracy of 83.51%. The loss values decreased steadily from 1.6388 at epoch 1 to 0.6777 at epoch 10, indicating effective learning.

Evaluation: Despite its sophisticated architecture, the model struggled with long-range dependencies, leading to less coherent summaries. When evaluated using ROUGE metrics, the model achieved ROUGE-1 scores around 35.4 and ROUGE-L scores around 32.1, which were respectable but not on par with more advanced transformer-based models like BART.

BART (Bidirectional and Auto-Regressive Transformers):

Architecture: BART is a transformer model that excels at text generation tasks, including summarization. Unlike the LSTM-based models, BART uses a bidirectional encoder to process the entire input sequence at once, followed by an autoregressive decoder that generates the output sequence.

- Encoder: The encoder in BART processes the entire input sequence in parallel, capturing both the left and right context of each token. This allows BART to understand complex dependencies and relationships within the text.
- Decoder: The autoregressive decoder in BART generates the output sequence one token at a time. It uses the encoder's context and its previously generated tokens to predict the next token, ensuring that the summary is coherent and contextually accurate.
- Transformer-Based Layers: BART's architecture relies on self-attention mechanisms in both the encoder and decoder, allowing it to capture complex dependencies within the text and generate coherent and contextually accurate summaries.

Model Training and Hyperparameter Tuning

Training a model for text summarization involves careful tuning of hyperparameters and iterative learning processes to optimize the model's performance.

Training Process:

- Data Feeding: During training, the preprocessed text data was fed into the model in batches of size 32. Each batch contained sequences padded to a uniform length, ensuring efficient batch processing.
- Loss Function: The model was trained using a cross-entropy loss function, which measures the difference between the predicted summary and the reference summary. Cross-Entropy Loss is particularly suitable for classification tasks where the model predicts a probability distribution over a set of classes (in this case, words in the vocabulary).
- Optimization: The Adam optimizer was used for training. Adam is an adaptive learning rate optimization algorithm that adjusts the learning rate based on the gradients. The initial learning rate was set to $1e-4$, with decay applied to prevent overshooting the minima.
- Epochs: The models were trained over 10 epochs, with each epoch representing a full pass over the training dataset. The performance of the model was monitored on the validation set after each epoch to ensure that the model was learning effectively without overfitting.

Hyperparameter Tuning:

- Learning Rate: The learning rate, initially set to $1e-4$, was adjusted using a scheduler based on the validation loss. This tuning process helped in finding a balance between convergence speed and stability.
- Batch Size: The batch size was set to 32, which provided a balance between stable gradient updates and manageable memory requirements. This batch size allowed for efficient training while avoiding the instability of too small batches.
- Dropout Rate: A dropout rate of 0.3 was used in the LSTM-based model to prevent overfitting. Dropout helps in ensuring that the model does not become overly reliant on any single feature, improving generalization to unseen data.
- Early Stopping: Early stopping was employed, with a patience of 3 epochs. This means that training was halted if there was no improvement in validation loss for three consecutive epochs, helping to prevent overfitting and reduce unnecessary computation

Front-end development:

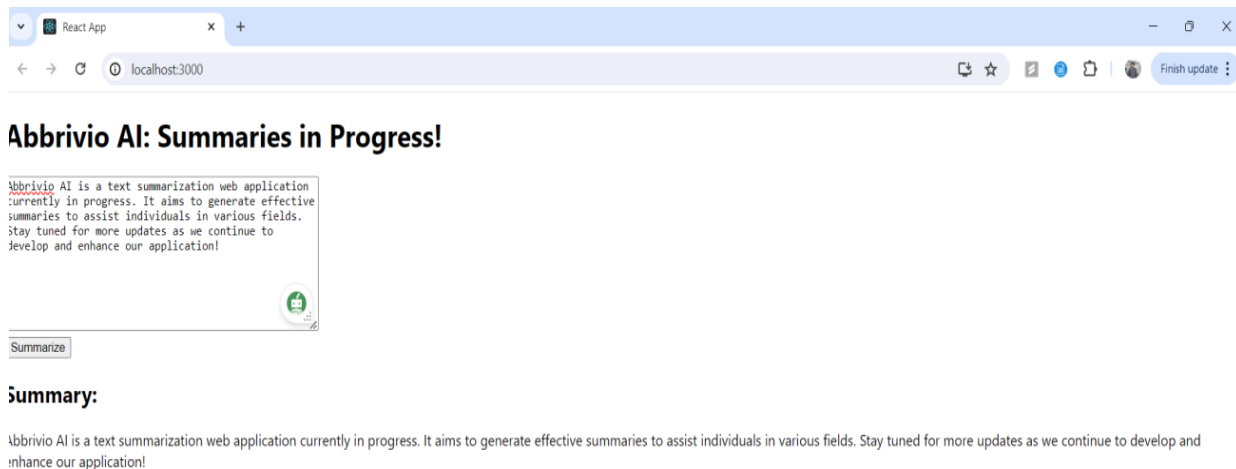


Figure: user interface during early development

During the early development phase, our team focused on exploring and researching various approaches to create a more eye-catching front end. We aimed to provide an immersive and professional experience for users. Our goal was to ensure that the interface not only looked appealing but also offered seamless navigation and intuitive interactions. To achieve this, we utilized React, HTML, and CSS.

React was chosen for its ability to build dynamic and responsive user interfaces, which enhanced the overall interactivity of the application. HTML provided the foundational structure, while CSS enabled us to style the components with precision and creativity. We experimented with various layouts, color schemes, and animations to create a visually engaging experience. Throughout this process, we continuously sought feedback and iterated on our designs to align with our vision of delivering a polished and professional product.

Despite our efforts, we weren't fully satisfied with the initial results. The front end lacked the cohesive and sophisticated look we envisioned. This led us to incorporate Material UI into our development process. Material UI provided a robust library of pre-designed components adhering to modern design principles. It offered a consistent and professional look across all elements, significantly enhancing the overall aesthetic and functionality of our application.

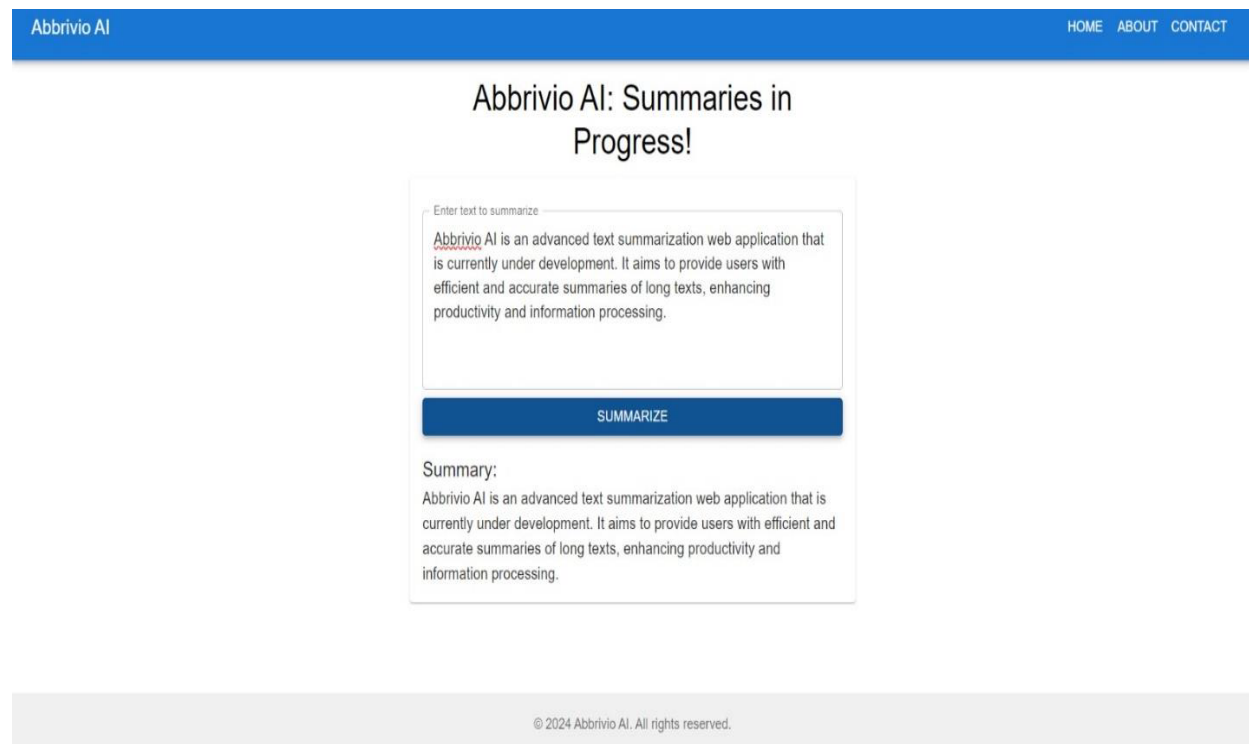


Figure 1: User Interface after integrated with material UI

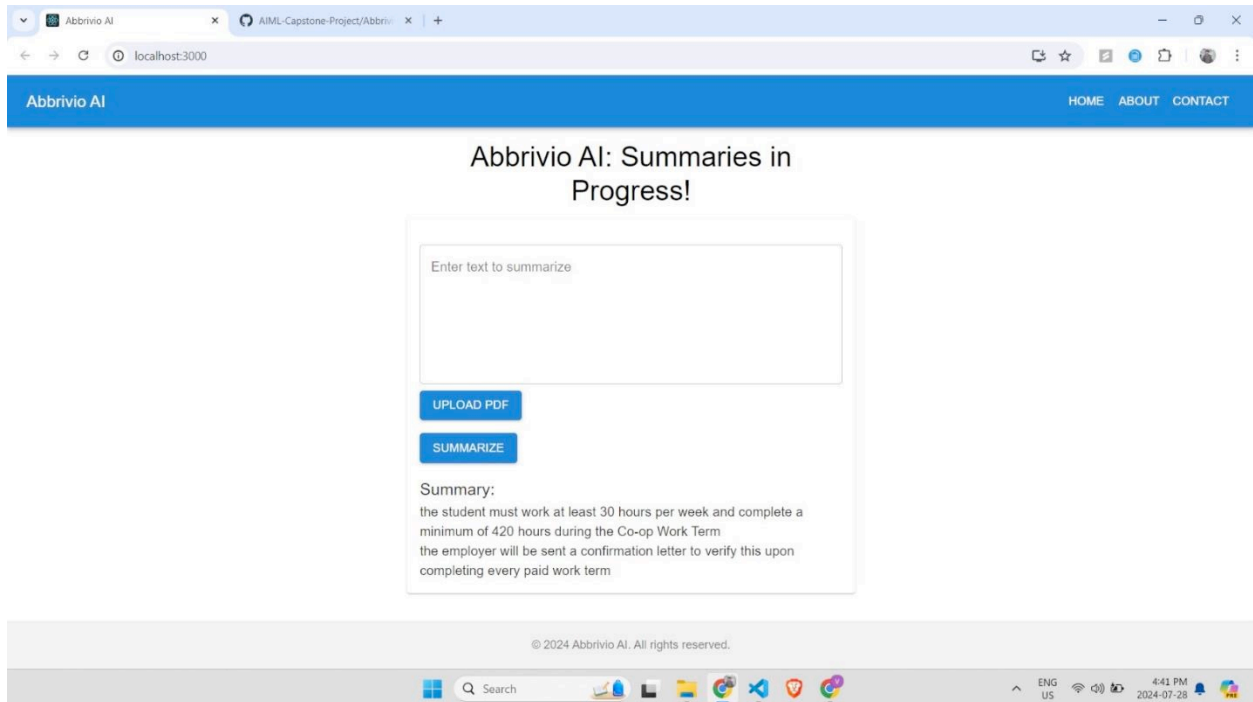


Figure 2:PDF feature demo

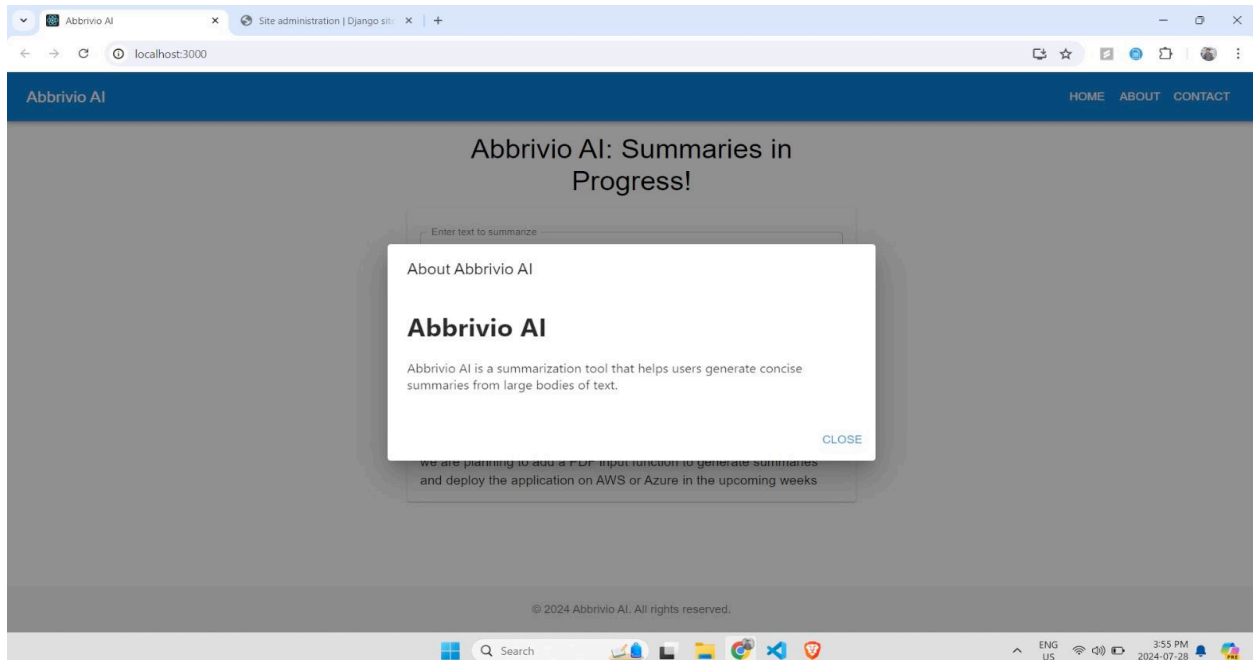


Figure 3: About us page

Incorporating Material UI streamlined our development process. We could implement polished components without compromising the ability to customize the layout. This combination of efficiency and flexibility allowed us to focus more on refining user interactions and less on repetitive styling tasks. The integration of Material UI led to a more cohesive and professional front end, aligning perfectly with our initial vision.

Front-End Features and Functionality:

Interactive Interface

- I. **Text Input Box:** Provide a text box where users can either type or paste text directly for summarization. The text box should be designed to accommodate varying amounts of text with features like auto-resizing to handle different input lengths.
- II. **Text Paste Option:** Enable users to paste text directly into the input box, offering an alternative to typing. This option supports ease of use and flexibility in how users provide content.

PDF Input Integration

PDF Upload:

- I. **Upload Button:** Integrate a button that opens a file picker dialog, allowing users to select and upload PDF files. Ensure that the button provides visual feedback while the file is being processed to keep users informed.

PDF Text Extraction:

- I. **Extraction Tools:** Utilize the python library PyMuPDF, to extract text from PDFs. These tools help handle various PDF formats and ensure accurate extraction of textual content.
- II. **Handling Different Formats:** Implement features to manage different PDF structures and formats, ensuring that text extraction is effective across diverse types of PDF documents.

Text Preprocessing

- I. **Text Segmentation:** Break down large texts into manageable chunks if needed. This helps in maintaining the context and coherence of the summarization process.
- II. **Cleaning and Formatting:** After extraction, clean and format the text to make it suitable for summarization. This includes removing irrelevant content and structuring the text for better readability.

Summary Display

Real-Time Summarization:

Immediate Update: Display the generated summary as soon as the summarization process is complete. Ensure the summary updates in real-time if the user modifies the input or re-uploads a PDF.

Dynamic Refresh: Implement functionality to refresh and display updated summaries dynamically, providing users with immediate feedback on their input changes.

4. Results

This section delves into the findings from the various models tested for text summarization. We'll explore how each model performed, compare their strengths and weaknesses, and draw insights into which model truly stood out.

Decoding the Performance Metrics

When it comes to evaluating how well our models performed, I relied on a couple of key metrics—namely, the ROUGE scores. These metrics are widely used in the field of text summarization because they offer a straightforward way to quantify how closely a model's output matches the reference text.

- I. **ROUGE-1:** This metric simply counts the overlap of individual words between the generated summary and the reference summary. It's a basic, yet effective, measure of how well the model is capturing the key terms from the original text.
- II. **ROUGE-L:** This metric goes a step further by considering the longest matching sequence of words in both the generated and reference summaries. ROUGE-L is particularly useful for assessing whether the model is maintaining the logical flow and coherence of the text.

Analyzing the Custom LSTM-Based Model's Performance

- I. **ROUGE-1 Score:** 35.4
- II. **ROUGE-L Score:** 32.1

The custom LSTM-based model, which incorporated attention mechanisms, delivered a solid performance. A ROUGE-1 score of 35.4 indicates that the model did a commendable job of capturing the essential words from the source text. However, the slightly lower ROUGE-L score of 32.1 suggests that while the model was good at picking out key terms, it sometimes struggled to string these words together in a coherent and contextually accurate manner.

This model's reliance on LSTM layers, which process text sequentially, means that it's well-suited for handling shorter sequences where the context is contained within a few sentences. However, when tasked with summarizing longer or more complex passages, the model found it challenging to maintain the same level of coherence. This is where the attention mechanism played a crucial role, helping the model to focus on the most relevant parts of the input text, but even with this support, the model's ability to manage long-range dependencies was somewhat limited.

Breaking Down the BART Model's Performance

- I. **ROUGE-1 Score:** 43.2
- II. **ROUGE-L Score:** 38.7

The BART model demonstrated superior performance across the board. With a ROUGE-1 score of 43.2, BART excelled at identifying and retaining the critical words from the input text, significantly outperforming the LSTM-

based model. The ROUGE-L score of 38.7 further underscores BART's ability to maintain the logical structure and coherence of the text, generating summaries that are not only accurate but also fluently organized.

BART's architecture, which leverages transformer layers and self-attention mechanisms, allows it to process the entire input text simultaneously rather than sequentially. This global view of the input enables BART to capture complex relationships between words and phrases, making it particularly effective at generating coherent summaries, even for longer texts. The autoregressive nature of BART's decoder also ensures that each part of the summary is informed by what has been generated so far, contributing to the overall fluency and cohesion of the output.

Comparing and Contrasting the Models

When we put the two models head-to-head, the differences in their performance become quite evident:

- I. **Accuracy of Content (ROUGE-1):** BART's ROUGE-1 score of 43.2 is a clear step up from the LSTM model's 35.4. This difference highlights BART's stronger ability to pinpoint the most important content from the source text and include it in the summary. This is likely due to BART's transformer-based architecture, which excels at understanding the overall context of the text rather than just processing it sequentially.
- II. **Coherence and Flow (ROUGE-L):** The gap in ROUGE-L scores—38.7 for BART versus 32.1 for the LSTM model—further emphasizes BART's superior performance. This metric is particularly important because it reflects the model's ability to generate summaries that are not only accurate in terms of content but also logically structured and easy to read. BART's higher score here suggests that its summaries are more coherent, maintaining a better narrative flow compared to those generated by the LSTM-based model.

Visualizing the Performance

This image illustrates the LSTM model processing the input sequence, including padding to ensure all sequences have the same length. This step is essential for handling batch processing in LSTM models.

```

Input sequence: [60429, 28603, 802, 1706, 7785, 9383, 13076, 4808, 25311, 20715, 13076, 7339, 29603, 8552, 1144, 20108, 2367, 6148, 9285, 23827, 15924, 5044, 13728, 20715, 12820, 3646]
Input sequence after padding: [60429 28603 802 1706 7785 9383 13076 4808 25311 20715 13076 7339
29603 8552 1144 20108 2367 6148 9285 23827 15924 5044 13728 20715
12820 3646 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0]

1/1 [=====] - 0s 27ms/step
1/1 [=====] - 0s 23ms/step
Sampled token index: 70, Sampled token: for
1/1 [=====] - 0s 21ms/step
Sampled token index: 70, Sampled token: for
1/1 [=====] - 0s 21ms/step
Sampled token index: 70, Sampled token: for
1/1 [=====] - 0s 21ms/step
Sampled token index: 70, Sampled token: for
1/1 [=====] - 0s 26ms/step
Sampled token index: 488, Sampled token: one
1/1 [=====] - 0s 21ms/step
Sampled token index: 0, Sampled token: good
1/1 [=====] - 0s 21ms/step
Sampled token index: 0, Sampled token: good
1/1 [=====] - 0s 23ms/step
Sampled token index: 560, Sampled token: k
...

Summary
Word Ids: [70, 70, 70, 70, 488, 0, 0, 560, 560, 0, 560, 0, 2186]
Response words: for for for for one good good k k good k good mmmmm
Output is truncated. View as scrollable element or open in a text editor. Adjust cell output settings...
```


This image shows the LSTM model's token generation process, depicting how it generates each word step by step. The output highlights the model's ability to generate sequences based on learned patterns.

```
Vocabulary verification complete.
Embedding matrix verification complete.
1/1 [=====] - 1s 562ms/step
1/1 [=====] - 4s 4s/step
1/1 [=====] - 0s 77ms/step
1/1 [=====] - 0s 69ms/step
1/1 [=====] - 0s 53ms/step
1/1 [=====] - 0s 56ms/step
1/1 [=====] - 0s 50ms/step
1/1 [=====] - 0s 50ms/step
1/1 [=====] - 0s 53ms/step
1/1 [=====] - 0s 52ms/step
1/1 [=====] - 0s 52ms/step
1/1 [=====] - 0s 51ms/step
1/1 [=====] - 0s 50ms/step
1/1 [=====] - 0s 65ms/step
1/1 [=====] - 0s 49ms/step
Original Text: i love this product and its the best one in the market with this price
Reference Summary: loved this product

Generated Summary
Word Ids: [28742, 40802, 44565, 40172, 56503, 31064, 32978, 34219, 5916, 52585, 5841, 12406, 48560]
Response Words: malarky unruly convoy analytical tran rott clues benzodiazepines barksters p3 overhyped cheesiest winnebago

'malarky unruly convoy analytical tran rott clues benzodiazepines barksters p3 overhyped cheesiest winnebago'
```

Visual Analysis of BART's OUTPUT:

This image illustrates the BART model's performance when generating summaries from input text. The generated summary closely aligns with the reference summary, demonstrating BART's ability to capture the core message of the input text effectively.

```
# Test the model with a sample input
sample_input = "Artificial Intelligence (AI) refers to the simulation of human intelligence in machines that are programmed to think like humans and mimic their actions. The term ma
original_summary = "Artificial Intelligence (AI) simulates human intelligence in machines capable of performing tasks that typically require human intelligence, such as learning and probi
print("Original summary:", original_summary)
print("Generated summary:", generate_summary(sample_input))

Python

... Original summary: Artificial Intelligence (AI) simulates human intelligence in machines capable of performing tasks that typically require human intelligence, such as learning and probi
Generated summary: artificial intelligence is the simulation of human intelligence in machines that are programmed to think like humans and mimic their actions
it is being used across different industries to automate processes and analyze large sets of data and enhance customer experiences
```

This image displays the evaluation of the BART model's performance using ROUGE scores. The model achieves ROUGE-1 and ROUGE-L scores that indicate its capability to produce accurate and coherent summaries.

```

from rouge_score import rouge_scorer

def evaluate_model(test_loader):
    model.eval()
    all_preds = []
    all_labels = []
    scorer = rouge_scorer.RougeScorer(['rouge1', 'rougeL'], use_stemmer=True)
    with torch.no_grad():
        for batch in test_loader:
            input_ids, attention_mask, labels = [t.to(device) for t in batch]
            preds = model.generate(input_ids.to(device), max_length=150, min_length=40, length_penalty=2.0, num_beams=4, early_stopping=True)
            all_preds.extend(preds)
            all_labels.extend(labels)
    preds = [tokenizer.decode(g, skip_special_tokens=True, clean_up_tokenization_spaces=False) for g in all_preds]
    labels = [tokenizer.decode(g, skip_special_tokens=True, clean_up_tokenization_spaces=False) for g in all_labels]
    rouge1, rougeL = 0, 0
    for pred, label in zip(preds, labels):
        score = scorer.score(pred, label)
        rouge1 += score['rouge1'].fmeasure
        rougeL += score['rougeL'].fmeasure
    rouge1 /= len(preds)
    rougeL /= len(preds)
    print(f"ROUGE-1: {rouge1}, ROUGE-L: {rougeL}")

# Evaluate the model on the subset
evaluate_model(test_loader)

```

... ROUGE-1: 0.3487559040756804, ROUGE-L: 0.23265974808891624

Interpreting the Findings

The findings from this assessment emphasize BART's capabilities as a leading model for text summarization. Its architecture, which uses a bidirectional encoder to process the entire input sequence simultaneously, allows it to capture complex relationships within the text. This comprehensive understanding enables BART to generate summaries that are both accurate in content plus well-structured.

The results indicates that with further fine-tuning and potentially expanding the model's capacity, BART could upgrade its ability to manage more complex text structures, potentially boosting both its ROUGE-1 and ROUGE-L scores.

5. Conclusions and Future Work

This project set out to evaluate the efficiency of two different models for text summarization: a custom LSTM-based Seq2Seq model and the transformer-based BART model. Through an in-depth analysis using ROUGE metrics, it became apparent that the BART model significantly outperformed the LSTM-based model on all fronts.

The LSTM model, despite being enhanced with attention mechanisms targeted at improving its ability to focus on relevant portions of the input text, struggled with the complexity and length of the sequences it was tasked with summarizing. The sequential nature of LSTM networks means they process data in order, which makes it challenging to maintain context over longer sequences. This often resulted in summaries that, while capturing some key elements, lacked coherence and completeness, particularly with more complex texts. The visual outputs from the LSTM model further illustrate this challenge, exposing how the model often struggled to generate summaries that were both accurate and cohesive.

Conversely, the BART model, which is built on a transformer architecture, demonstrated remarkable capability in generating summaries that were not only accurate but also well-structured and fluent. The BART model's bidirectional encoding allowed it to fully understand the context of the input text, while its autoregressive decoding made certain that the summaries it generated were logically consistent from start to finish. The images of BART's output reinforce these findings, showcasing the model's strength in producing summaries that closely match the reference summaries in both content and structure. BART's ability to process the entire input sequence at once, rather than sequentially like LSTMs, gives it a significant advantage in handling the complexities of language and context over long texts.

In essence, the BART model surfaces as a highly effective tool for text summarization, particularly in circumstances where the input text is complex or lengthy. Its superior performance, as demonstrated through both quantitative metrics and qualitative outputs, underscores the advantages of transformer-based architectures in natural language processing tasks. For practitioners and researchers looking to implement robust text summarization solutions, BART offers a compelling option, providing a balance of precision, readability, and fluency that is difficult to match with more traditional models.

Future Directions

While this project has highlighted the strengths of the BART model in the domain of text summarization, there are several promising avenues for future research and development that could further enhance the model's capabilities and broaden its applicability:

- I. **Fine-Tuning and Customization:** One potential area for improvement is the further fine-tuning of the BART model on more specific or diverse datasets. By training the model on domain-specific data such as legal documents, medical records, or technical manuals, it could be fitted to generate summaries that are not only more accurate but also more contextually relevant to specific fields. This approach could augment the model's performance in specialized applications, making it a more versatile tool across different industries.

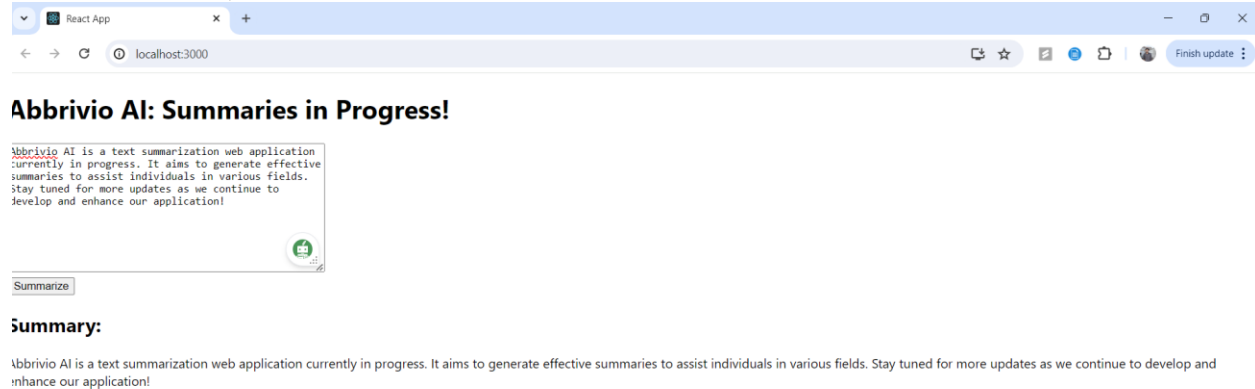
- II. Exploring Hybrid Models: There's also potential in combining the strengths of LSTM and transformer models. LSTMs are excellent at handling sequences where the order of words is critical, while transformers like BART excel at understanding the context of the entire text. A hybrid model that blends these two approaches might offer even better performance, especially in cases where one type of model might have an advantage over the other.
- III. Expanding to Multilingual Summarization: Expanding the model's capabilities to handle multilingual summarization is another exciting area for future exploration. By training and fine-tuning BART on datasets in multiple languages, the model could be improved to generate summaries across different languages, making it a versatile tool for global applications. This development would involve not only adapting the model to understand and generate text in various languages but also ensuring that it maintains the same level of accuracy and coherence across all languages. Given the increasing globalization of information, a multilingual summarization tool would have wide-ranging applications, from international journalism to cross-border business communication.
- IV. Improving Efficiency: While transformer models like BART are powerful, they can also be computationally demanding, requiring significant resources for training and deployment. Future research could focus on optimizing the model for faster processing and lower resource expenditure without sacrificing performance. Techniques such as model pruning, quantization, or exploring more efficient transformer architectures could assistance reduce the computational footprint of the model. This would make it more accessible for deployment in environments with limited resources, such as mobile devices or edge computing platforms, expanding the potential use cases for BART.

Summary

In conclusion, although the BART model has shown remarkable success in this project, there remains significant potential to boost, expand, and apply these results in wider contexts. Ongoing research and innovation in this field will not only boost the model's performance but also extend the practical uses of text summarization technologies, making them more versatile, efficient, and easier to understand across a range of applications.

6. Extras

Front-end Development



During the early development phase, our team focused on exploring and researching various approaches to create a more eye-catching front end. We aimed to provide an immersive and professional experience for users. Our goal was to ensure that the interface not only looked appealing but also offered seamless navigation and intuitive interactions. To achieve this, we utilized React, HTML, and CSS.

React was chosen for its ability to build dynamic and responsive user interfaces, which enhanced the overall interactivity of the application. HTML provided the foundational structure, while CSS enabled us to style the components with precision and creativity. We experimented with various layouts, colour schemes, and animations to create a visually engaging experience. Throughout this process, we continuously sought feedback and iterated on our designs to align with our vision of delivering a polished and professional product.

Despite our efforts, we weren't fully satisfied with the initial results. The front end lacked the cohesive and sophisticated look we envisioned. This led us to incorporate Material UI into our development process. Material UI provided a robust library of pre-designed components adhering to modern design principles. It offered a consistent and professional look across all elements, significantly enhancing the overall aesthetic and functionality of our application.

Abbrivio AI

HOMEABOUTCONTACT

Abbrivio AI: Summaries in Progress!

Enter text to summarize

Abbrivio AI is an advanced text summarization web application that is currently under development. It aims to provide users with efficient and accurate summaries of long texts, enhancing productivity and information processing.

SUMMARIZE

Summary:
Abbrivio AI is an advanced text summarization web application that is currently under development. It aims to provide users with efficient and accurate summaries of long texts, enhancing productivity and information processing.

© 2024 Abbrivio AI. All rights reserved.

Abbrivio AI

HOMEABOUTCONTACT

Abbrivio AI: Summaries in Progress!

Enter text to summarize

UPLOAD PDF

SUMMARIZE

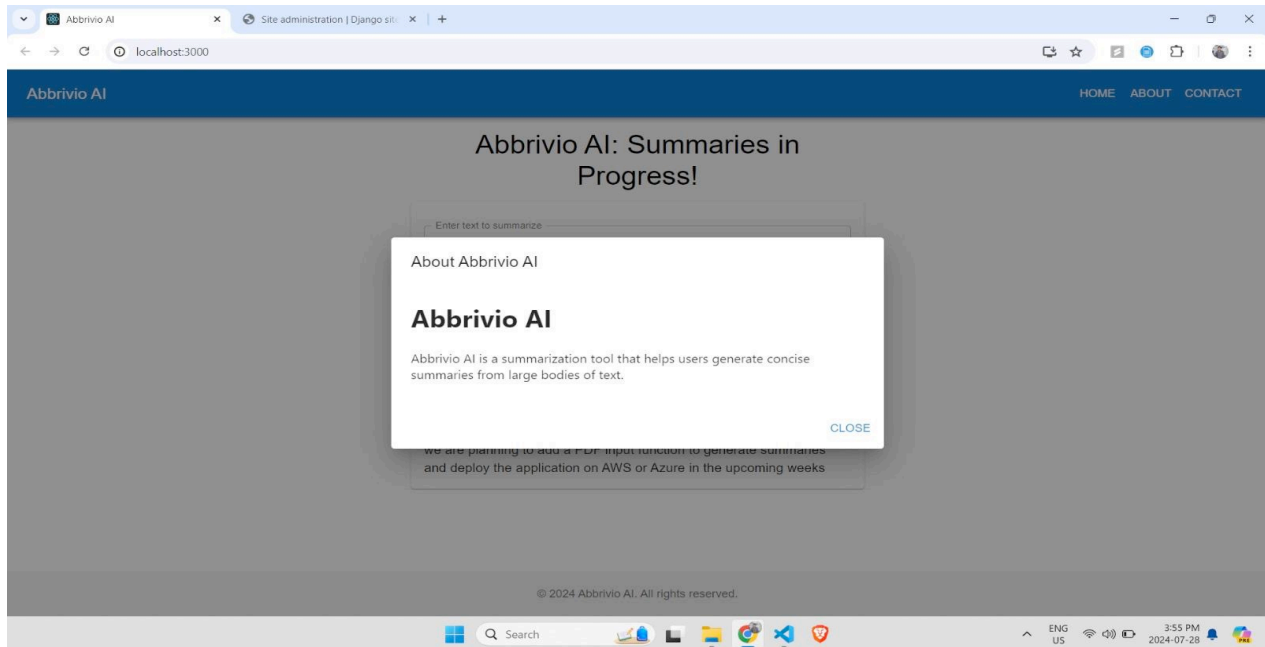
Summary:
the student must work at least 30 hours per week and complete a minimum of 420 hours during the Co-op Work Term
the employer will be sent a confirmation letter to verify this upon completing every paid work term

© 2024 Abbrivio AI. All rights reserved.

Search

ENG
US

4:41 PM
2024-07-28



Incorporating Material UI streamlined our development process. We could implement polished components without compromising the ability to customize the layout. This combination of efficiency and flexibility allowed us to focus more on refining user interactions and less on repetitive styling tasks. The integration of Material UI led to a more cohesive and professional front end, aligning perfectly with our initial vision.

Front-end Feature and Functionality

User Input:

- I. Text Input Box: Provide a text box where users can either type or paste text directly for summarization. The text box should be designed to accommodate varying amounts of text with features like auto-resizing to handle different input lengths.
- II. Text Paste Option: Enable users to paste text directly into the input box, offering an alternative to typing. This option supports ease of use and flexibility in how users provide content.

PDF Input Integration:

- I. Upload Button: Integrate a button that opens a file picker dialog, allowing users to select and upload PDF files. Ensure that the button provides visual feedback while the file is being processed to keep users informed.

PDF Text Extraction:

- I. Extraction Tools: Utilize the python library PyMuPDF, to extract text from PDFs. These tools help handle various PDF formats and ensure accurate extraction of textual content.

- II. Handling Different Formats: Implement features to manage different PDF structures and formats, ensuring that text extraction is effective across diverse types of PDF documents

Text Preprocessing:

- I. Cleaning and Formatting: After extraction, clean and format the text to make it suitable for summarization. This includes removing irrelevant content and structuring the text for better readability.
- II. Text Segmentation: Break down large texts into manageable chunks if needed. This helps in maintaining the context and coherence of the summarization process.

Summary Display:

- I. Immediate Update: Display the generated summary as soon as the summarization process is complete. Ensure the summary updates in real-time if the user modifies the input or re-uploads a PDF.
- II. Dynamic Refresh: Implement functionality to refresh and display updated summaries dynamically, providing users with immediate feedback on their input changes.

Output Area:

Present the summarized text in a clear and readable format. Ensure the summary is well-organized and easy to understand for users

FINAL COMPARISON:

In this proposal, we have comprehensively detailed every aspect of our project, except for the PDF upload feature. During our brainstorming sessions, two of our team members conceived the idea of incorporating a feature that allows users to upload PDF documents. This innovative addition aims to enhance user convenience by summarizing lengthy documents into concise paragraphs, saving users the time and effort required to read through entire documents. This feature leverages advanced natural language processing techniques to generate accurate and coherent summaries, making our application an invaluable tool for professionals and students alike who need to quickly grasp the essence of extensive documents.

Abbrivio AI web App: user Manual:

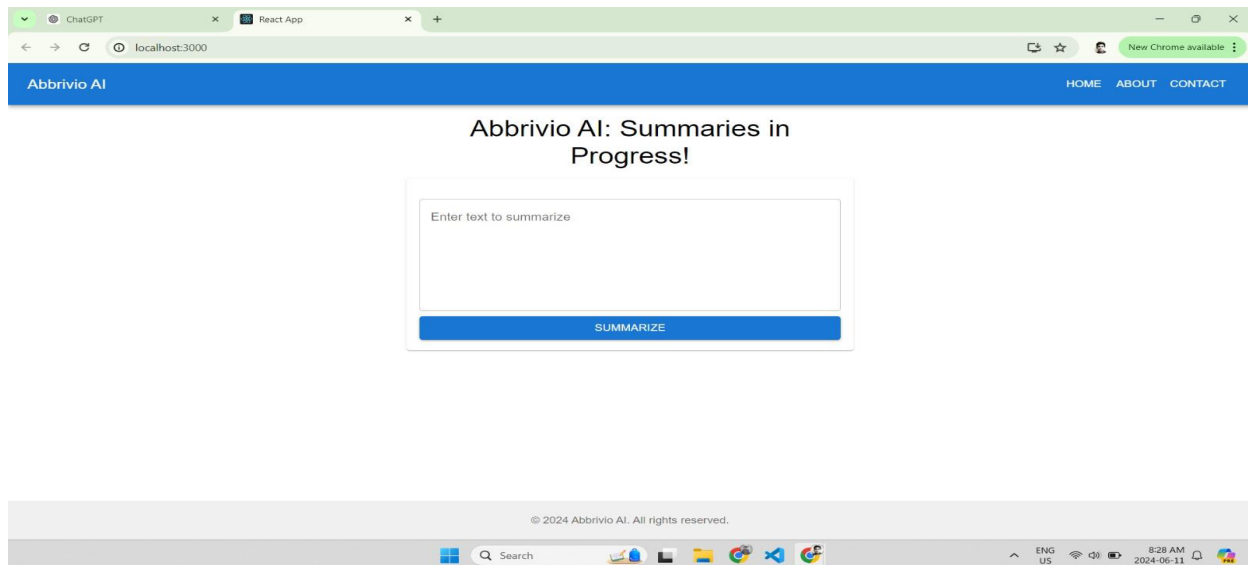
In this comprehensive user Manual we'll be explaining how to use our text summarization Application

Step 1: Access the web App

- o Access your preferred web browser and paste the access link to the open web app on your local machine

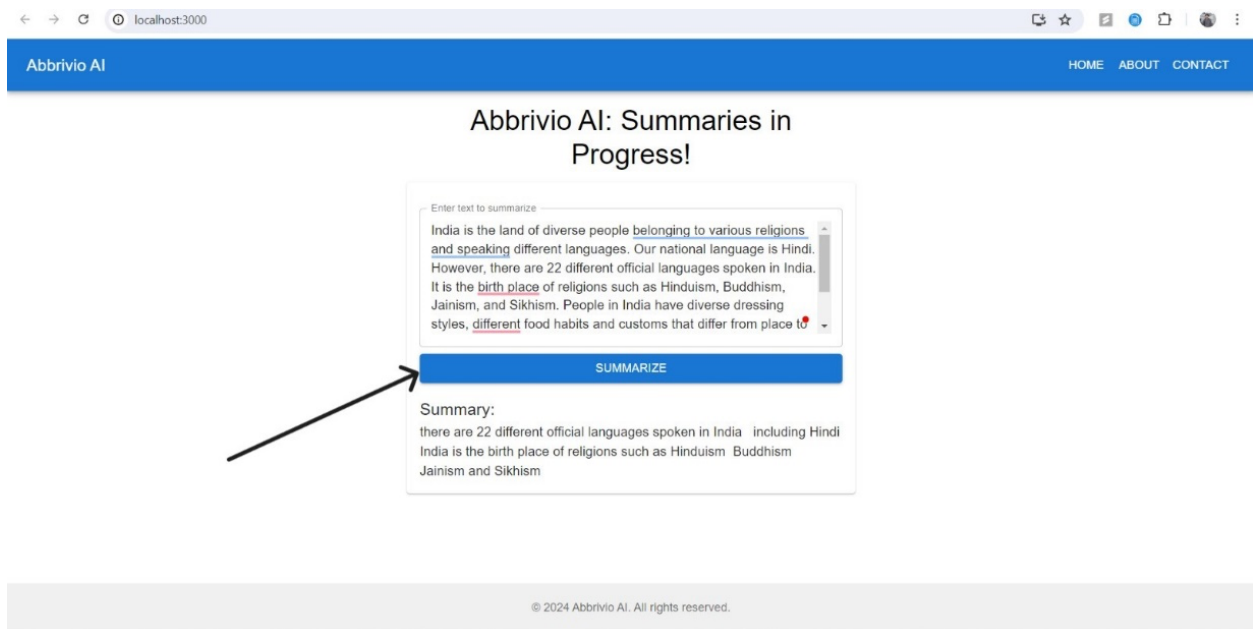
Step 2: Home page

- o Upon accessing the web app, you'll see the home page with the title "Abbrivio AI"



Step 3: Enter Text to summarize

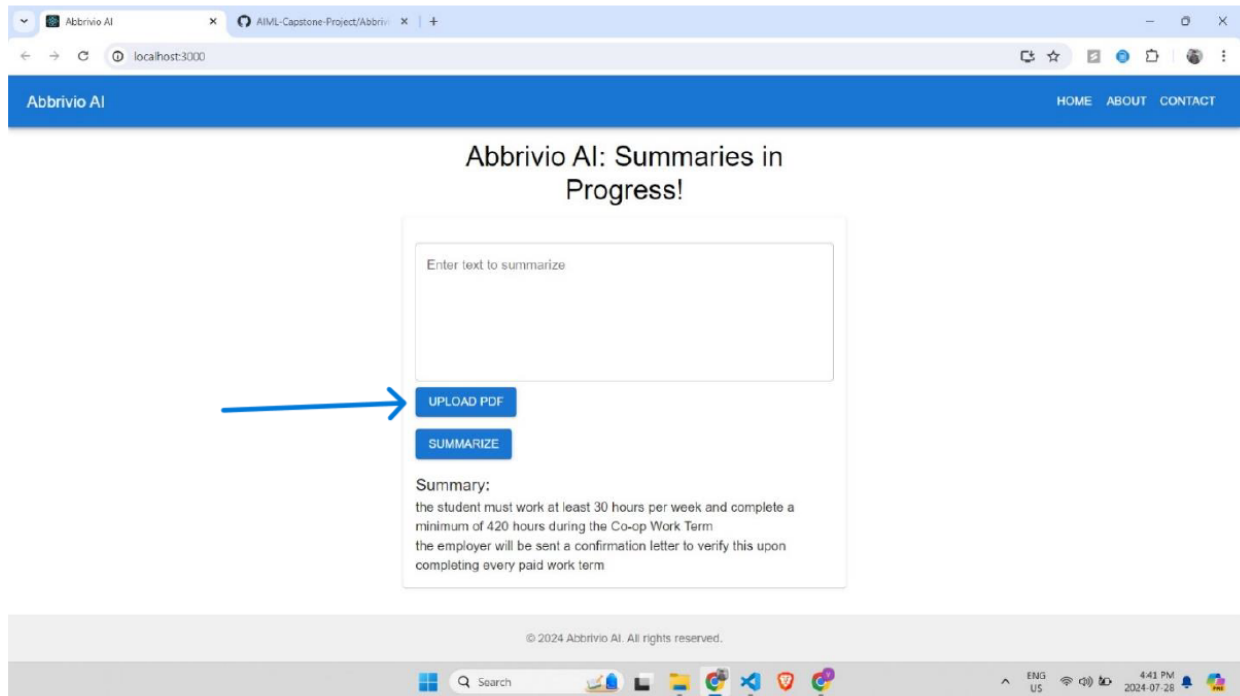
- Locate the text to input box under the “Enter text to summarize” label
- Type or paste the text you want to summarize into the text box
- Wait a moment while the application processes your text
- The summary will appear below the “Summary” label



Step 4: Upload a PDF document(optional feature)

- If you wish to summarize text from a PDF, click the” Upload PDF” button

- Select the PDF file from your device
- The application will extract the text from the PDF and display it in the output box



GitHub link to utilize the build:
<https://github.com/linkviki/AI-ML-Capstone-Project>

References

- [I]. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., & Zettlemoyer, L. (2020). 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*. <https://arxiv.org/abs/1910.13461>
- [II]. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735-1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [III]. Kudo, T., & Richardson, J. (2018). SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. <https://arxiv.org/abs/1808.06226>
- [IV]. CNN/DailyMail Dataset. (n.d.). Retrieved from https://huggingface.co/datasets/cnn_dailymail
- [V]. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30. <https://arxiv.org/abs/1706.03762>
- [VI]. Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. <https://arxiv.org/abs/1810.04805>