**College Name:** VIT BHOPAL
**Student Name:** NAMRATA BHUTANI

## GEN AI PROJECT PHASE 3 SUBMISSION DOCUMENT

### Phase 3: Final Report and Submission

### 1. Project Title:

Text Generation and Summarization System using Transformers

### 2. Summary of Work Done

#### Phase 1 – Proposal and Idea Submission (10 Marks):

This project aims to improve the limitations of current text generation and summarization approaches by tapping into the potential of pre-trained Transformer models, GPT-2 XL and BART Large CNN. The new system will utilize these models in a user interface developed using Gradio, where the user is able to feed in text and receive high-quality generated or summary output. The main goals of this project include improving text quality, summarization quality, and offering a useful tool for many different NLP applications.

We submitted a detailed proposal including problem definition, objectives, tools required, and expected outcomes.

#### Phase 2 – Execution and Demonstration (15 Marks):

In the second phase, we implemented the proposed idea using Python, HuggingFace Transformers, and Gradio. The following tasks were accomplished:

- Built a web-based interface using Gradio: Developed an interactive interface using Gradio to enable users to input text for both generation and summarization.
- Loaded pre-trained models: Initialized and loaded the pre-trained GPT-2 XL model for text generation and the Facebook BART Large CNN model for text summarization.
- Integrated text generation: Developed a function to take in user input, use the GPT-2 XL model to produce text according to the prompt, and present the output produced in the interface.
- Used text summarization: Implemented a function to accept user input, apply the BART Large CNN model for summarizing text, and show the summarized output in the interface.
- Merged features: Combined text production and summarizing capabilities in one integrated interface, thus facilitating seamless switching between activities.Tested and refined: Conducted thorough testing of both functionalities with various text inputs to ensure performance, relevance, and accuracy.

Sample outputs and the complete code were documented and submitted.

## 3. GitHub Repository Link

You can access the complete codebase, README instructions, and any related resources at the following GitHub link:

🔗 [Text Generation and Summarization using transfomers](#)

(or)

🔗 Visit: https://github.com/nam-bhutani/TextGenerateandSummarize_GenAI

## 4. Testing Phase

### 4.1 Testing Strategy

The system was tested across a variety of use cases to ensure its robustness and accuracy. The testing phase involved both manual testing and automated testing methods to verify the following:

- **Input Handling:** Ensuring the system handles different types of input text for both generation and summarization (e.g., short, long, complex sentences, paragraphs, articles).
- **Contextual Relevance:** Verifying that the generated text is contextually relevant and coherent with the input prompt, and that the generated summaries accurately reflect the main points of the original text.
- **Edge Case Testing:** Testing the models with incomplete sentences, nonsensical input, or extremely long text to observe how the system behaves and handles unexpected scenarios.

### 4.2 Types of Testing Conducted

1. **Unit Testing:**
   o Each function and module (like text generation, summarization, UI components) was tested independently to ensure they work correctly.
2. **Integration Testing:**
   o The integration of the GPT-2 XL and BART Large CNN models with the Gradio interface was tested to ensure smooth interaction between the models and the web interface.
3. **User Testing:**
   o A group of test users interacted with the system to assess its ease of use, interface design, and output relevance for both generation and summarization tasks. Feedback was collected and used for improvements.

4. **Performance Testing:**
   o The system was tested with various input text lengths and complexities to observe any potential delays or slow responses in generating text or summaries.

## 4.3 Results

- **Accuracy:** The system consistently generated contextually relevant text based on the prompts and produced accurate and concise summaries of the input text.
- **Response Time:** The application performed optimally with reasonable response times for both generation and summarization tasks, considering the complexity of the models.
- **Edge Cases:** Incomplete or nonsensical inputs resulted in less coherent or accurate outputs, but the system still attempted to generate text or summaries, demonstrating its ability to handle unexpected inputs to some extent.

## 5. Future Work

While the project successfully implements the **Text Generation and summarization** system, there are several avenues for future enhancement:

1. Model Fine-tuning: Fine-tuning the pre-trained models (GPT-2 XL and BART Large CNN) on specific datasets relevant to target domains can significantly improve their performance and accuracy for those domains.
2. Improved User Interface: The existing Gradio-based interface can be improved with additional features like user login, parameter customization, and integration with other tools and services.
3. Multilingual support: Making the system multi-lingual will expand its user base and usage potential.
4. Real-time Applications: Discussion of how the system could be used to integrate into real-time applications such as chatbots or virtual assistants would allow for text generation and summarization in a dynamic, interactive manner.
5. Diversity of Content and Control: Providing for managing diversity and stylistic features of generated text, i.e., having the option of defining preferred tone or topics, would increase the degree of creative freedom offered to users.
6. Testing in Specific Domains: Implementing a stricter testing of the system's effectiveness in different specific domains like news articles, academic journals, or literature would be very informative of its strengths and weaknesses in different situations.
7. Ethical Aspects: Similar to any system that is based on artificial intelligence, there are ethical concerns to be addressed, such as bias, fairness, and responsible use. Future efforts need to tackle how to counteract these issues and offer ethical use of the system.

## 6. Conclusion

This project effectively proved the capability of using pre-trained Transformer models, i.e., GPT-2 XL and BART Large CNN, in building a strong system that can both generate and summarize text. By combining these models with an interactive interface developed using Gradio, the system allows users to effortlessly provide input text and receive high-quality generated or summarized text.

By intense testing and evaluation, the system proved that it could produce context-relevant text, build accurate summaries, and handle a wide variety of input forms. While there remains room for growth and further research, this project provides a solid foundation for future advancement in the area of generative AI for natural language processing.

The achievement attained in this project enhances the continuous process of automating text-intensive tasks, enriching content production, and supporting effective knowledge attainment. By shattering the confines of existing processes, this project presents new prospects, such as chatbots, content-creation software, and automated report generation. Future development and exploration, as described in the future work section, will further advance the capabilities of the system and address possible shortcomings, thereby enabling more advanced and meaningful applications of generative AI in natural language processing.

## 7. References:

- Hugging Face Transformers: https://huggingface.co/docs
- PyTorch: https://pytorch.org
- OpenAI GPT-2 Research and API: https://openai.com/index/gpt-2-1-5b-release/
- Gradio Documentation: https://www.gradio.app/docs
- Facebook BART Research Paper: https://huggingface.co/facebook/bart-large-cnn
- Research Papers: https://arxiv.org/abs/1706.03762