# Chatbot Implementation Report

## 1. Overall Approach

The chatbot implementation leverages a combination of advanced machine learning models and efficient text processing techniques to provide accurate and helpful responses to user queries. By using a quantized model, the chatbot can operate on limited computational resources while maintaining high performance. The primary steps involved include loading the model and tokenizer, preparing text embeddings, and creating a conversational retrieval chain to handle user inputs and generate responses.

## 2. Frameworks/Libraries/Tools Used

- Flask-ngrok & Pyngrok: For exposing the local server to the internet during development.

- LangChain: Used for creating conversational retrieval chains and managing different components

of the chatbot.

- Chroma: Employed as the vector store for storing and retrieving text embeddings.

- Sentence-Transformers: Used for generating text embeddings for similarity searches.

- Transformers: Provides the model and tokenizer for text generation.

- Auto-GPTQ: Utilized for loading and running a quantized version of the model.

- PyTorch: Used for model inference, leveraging GPU acceleration if available.

## 3. Problems Faced and Solutions

- High Memory Usage: Initial attempts with standard models led to high memory usage, which was mitigated by using a quantized model (Nous-Hermes-13B-GPTQ).

- Latency Issues: To reduce latency, GPU acceleration was leveraged, and efficient text processing techniques like RecursiveCharacterTextSplitter were used.

- Scalability: Ensuring the chatbot could handle large documents was a challenge. This was addressed by splitting large documents into manageable chunks and using a vector store for efficient similarity searches.

## 4. Future Scope

- Feature Enhancements: Adding support for more languages, improving response generation quality, and incorporating more sophisticated natural language understanding capabilities.

- Performance Optimization: Implementing batch processing, caching mechanisms, and asynchronous processing to further reduce latency and improve responsiveness.

- User Interface: Developing a more user-friendly interface for interacting with the chatbot, potentially integrating it into various platforms like websites, mobile apps, and messaging services.

- Knowledge Base Expansion: Continuously updating and expanding the knowledge base to cover a wider range of topics and provide more comprehensive support.