

# GPT-2 Chatbot Using Streamlit Application

## Project Documentation and Technical Report

### Project Overview

**Objective:** I built a fun, interactive chatbot using some cool tech (GPT-2 and Streamlit). It lets you chat with an AI that generates text in real-time.

### Technical Architecture

- **Framework:** Streamlit •
- **Model:** OpenAI's GPT-2 •
- **Libraries:**
  - Torch
  - Transformers
  - Streamlit

### Key Implementation Components

#### 1. Model Loading Strategy

- I used `@st.cache_resource` decorator for efficient model caching
- I also made a backup plan for a specific technical detail to ensure the AI works smoothly.
- I chose the standard version of the GPT-2 AI model as the foundation for the chatbot.

#### 2. Response Generating Parameters

##### Initial Configuration

```
max_length=150
no_repeat_ngram_size=2
temperature=0.5
top_k=50 top_p=0.95
do_sample=True
```

## Parameter Tuning and Rationale

### Temperature (0.5)

- **Purpose:** Controls randomness in response generation
- **Outcome:** Balanced between creativity and coherence
- **Adjustment Impact:** Moderate predictability with slight randomness

### Top-K Sampling (50)

- **Purpose:** Limits token selection to top 50 most probable tokens
- **Outcome:** Reduced nonsensical generations
- **Adjustment Impact:** Improved response quality and relevance

### Top-P Sampling (0.95)

- **Purpose:** Nucleus sampling to dynamically select token pool
- **Outcome:** More diverse and contextually appropriate responses
- **Adjustment Impact:** Enhanced response variety while maintaining coherence

## 3. User Interface Design

- I designed the interface to look like a typical chat window.
- I made it so the chatbot remembers what you've already talked about.
- I added a sidebar with helpful information about the application.
- I also included some behind-the-scenes safety nets to catch any errors while the AI is generating responses.

## Performance Characteristics

### Strengths

- It generates text super fast, so the conversation feels natural.
- It's easy to get up and running – you don't need a powerful computer.
- It's designed to keep the conversation flowing smoothly.
- It doesn't need a ton of processing power to work.

## Limitations

- The AI's knowledge is limited to what it was trained on, so it might not know about very specific topics.
- Sometimes, the AI might generate responses that don't quite make sense or aren't relevant to the conversation.
- The application lacks the ability to retain conversation history across different sessions.

## Potential Improvements

1. Implement fine-tuning on specific domain datasets
2. Improve its memory
3. Integrate more sophisticated sampling techniques
4. Add safety features

## Technical Challenges Addressed

- Token Generation and Sampling: The process of creating and selecting tokens to form coherent text.
- Session State Management: Maintaining the context and history of the conversation.
- Model Loading and Caching: Efficiently loading and storing the AI model for quick access and reduced computational cost.
- Error Handling in Generative AI: Managing potential errors during the text generation process to ensure application stability.

## Deployment Environment

- Recommended environments:
  - Local development
  - Streamlit Cloud

## Conclusion

This GPT-2 Chatbot Streamlit application effectively demonstrates the use of generative AI, offering an intuitive user interface that highlights the potential of transformer-based models in conversational applications.

---