# Local LLM Interface App

## 1. Objective

The goal of this project is to **install a local LLM** (Large Language Model) and interact with it via a **Streamlit web interface**. Users can input prompts and observe AI-generated responses in real time, while also comparing performance metrics like response time.

## 2. Instructions & Implementation

### Step 1: Install and Setup Local LLM

1. Install a local LLM such as **GPT-2** or **Ollama** (e.g., LLaMA 3.2).
2. Ensure **Python**, **PyTorch**, and **Transformers** library are installed.
3. Verify GPU availability, if possible (torch.cuda.is_available()), otherwise use CPU.

### Step 2: Run Streamlit App

- Save the provided script as streamlit_local_llm_interactive.py.
- Launch the app using:

```bash
streamlit run streamlit_local_llm_interactive.py
```

In the sidebar, select your model (e.g., GPT-2) and verify that the device is detected correctly.

### Step 3: Test Prompt Generation

- Enter a prompt.
- Click **Generate** 🤖.
- The app will display:
  - Generated text
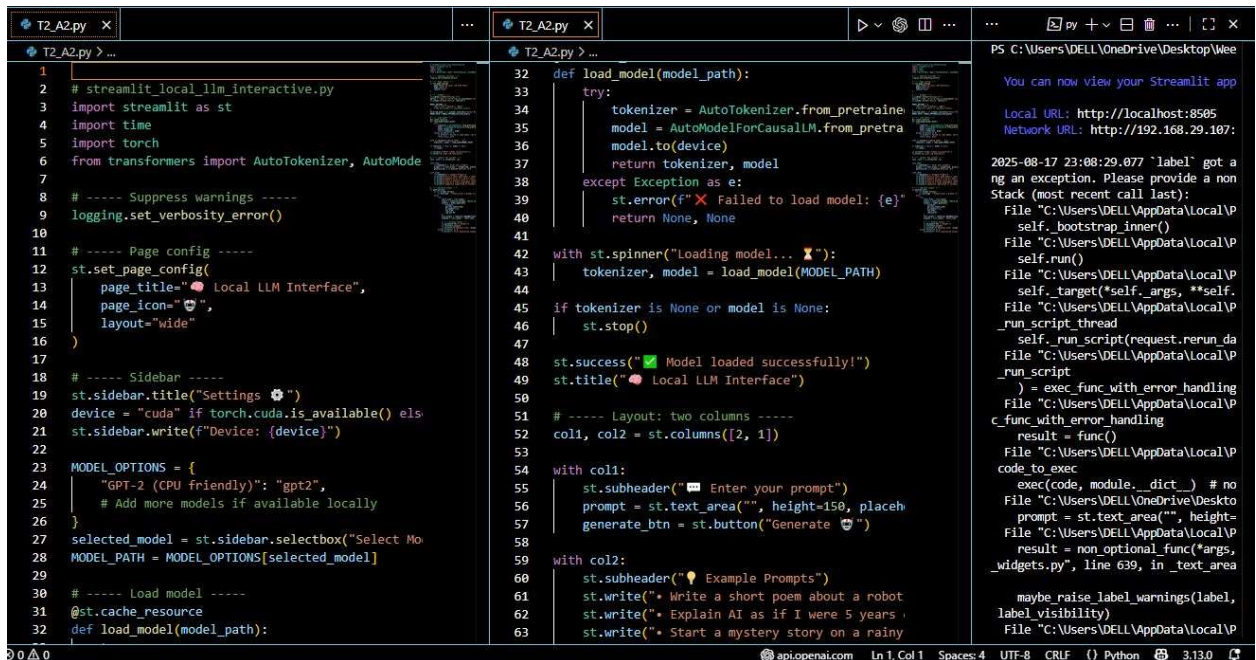  - Response time
- Example of token generation:

```python
inputs = tokenizer(prompt, return_tensors="pt").to(device)
outputs = model.generate(**inputs, max_new_tokens=150, do_sample=True, top_k=50, top_p=0.95)
text_output = tokenizer.decode(outputs[0], skip_special_tokens=True)
```

**Step 4: Troubleshooting**

- If the model fails to load, check:

  - Correct model path
  - PyTorch and Transformers versions
  - GPU availability

- If generation fails:

  - Reduce max_new_tokens
  - Disable sampling parameters (top_k, top_p) for smaller models

## 3. Deliverables

- Streamlit App: Interactive interface for real-time.
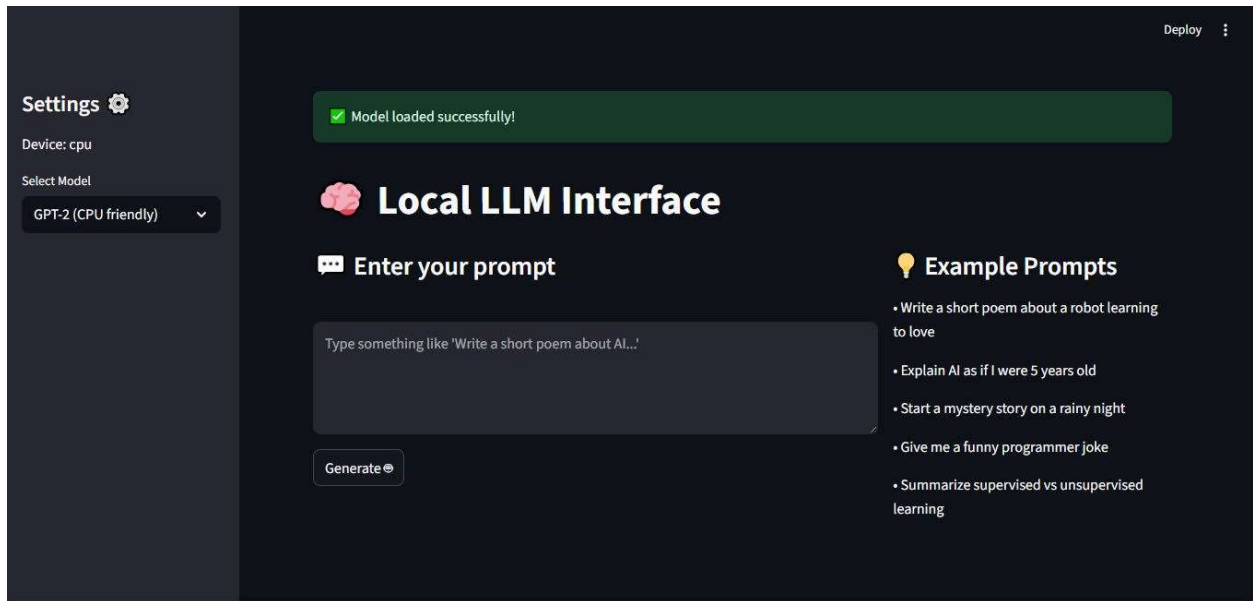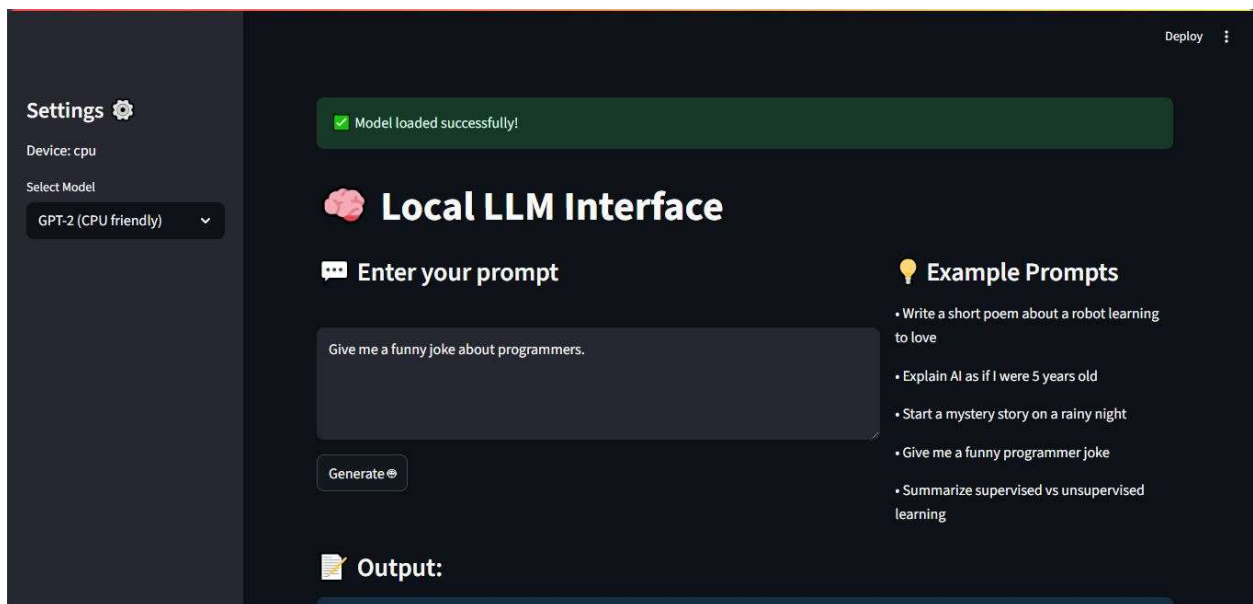


Fig1: Input Snapshot

# 4. Output



Fig2: Output Snapshot1



Fig3: Output Snapshot2

Fig4: Output Snapshot3